



# Empirical Risk Minimization with Relative Entropy Regularization

Samir M. Perlaza, Gaetan Bisson, Iñaki Esnaola, Alain Jean-Marie, Stefano Rini

## ► To cite this version:

Samir M. Perlaza, Gaetan Bisson, Iñaki Esnaola, Alain Jean-Marie, Stefano Rini. Empirical Risk Minimization with Relative Entropy Regularization. [Research Report] RR-9454, Inria. 2022. hal-03560072v6

**HAL Id: hal-03560072**

**<https://hal.science/hal-03560072v6>**

Submitted on 21 Nov 2023 (v6), last revised 27 Feb 2024 (v7)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License



# Empirical Risk Minimization with Relative Entropy Regularization

Samir M. Perlaza, Gaetan Bisson, Iñaki Esnaola, Alain Jean-Marie,  
and Stefano Rini

**RESEARCH  
REPORT**

**N° 9454**

February 2022

Project-Team NEO

ISRN INRIA/RR--9454--FR+ENG

ISSN 0249-6399





# Empirical Risk Minimization with Relative Entropy Regularization

Samir M. Perlaza, Gaetan Bisson, Iñaki Esnaola,  
Alain Jean-Marie, and Stefano Rini

Project-Team NEO

Research Report n° 9454 — version 6 — initial version February 2022 —  
revised version November 2023 — 100 pages

**Abstract:** The empirical risk minimization (ERM) problem with relative entropy regularization (ERM-RER) is investigated under the assumption that the reference measure is a  $\sigma$ -finite measure, and not necessarily a probability measure. Under this assumption, which leads to a generalization of the ERM-RER problem allowing a larger degree of flexibility for incorporating prior knowledge, numerous relevant properties are stated. Among these properties, the solution to this problem, if it exists, is shown to be a unique probability measure, often mutually absolutely continuous with the reference measure. Such a solution exhibits a probably-approximately-correct guarantee for the ERM problem independently of whether the latter possesses a solution. For a fixed dataset, the empirical risk is shown to be a sub-Gaussian random variable when the models are sampled from the solution to the ERM-RER problem. The generalization capabilities of the solution to the ERM-RER problem (the Gibbs algorithm) are studied via the sensitivity of the expected empirical risk to deviations from such a solution towards alternative probability measures. Finally, an interesting connection between sensitivity, generalization error, and mutual information is established.

**Key-words:** Supervised Learning, PAC-Learning, Regularization, Relative Entropy, Empirical Risk Minimization, Gibbs Measure, Gibbs Algorithm, Generalization, and Sensitivity.

Samir M. Perlaza and Alain Jean-Marie are with INRIA, Centre Inria d'Université Côte d'Azur, Sophia Antipolis 06902, France. Gaetan Bisson and Samir M. Perlaza are with the GAATI Laboratory at the Université de la Polynésie Française, Faaa 98702, French Polynesia. Iñaki Esnaola is with the ACSE Dept. at The University of Sheffield, Sheffield S1 3JD, UK. Stefano Rini is with the ECE Dept. at the National Yang Ming Chiao Tung University (NYCU), Hsinchu, Taiwan 30010, ROC. Samir M. Perlaza and Iñaki Esnaola are also with the ECE Dept. at Princeton University, Princeton N.J. 08544, USA. This work has been presented in part at the IEEE International Symposium on Information Theory (ISIT) in [1].

RESEARCH CENTRE  
SOPHIA ANTIPOLIS – MÉDITERRANÉE

2004 route des Lucioles - BP 93  
06902 Sophia Antipolis Cedex

## Minimisation du Risque Empirique avec Régularisation par l'Entropie Relative

**Résumé :** Le problème de minimisation du risque empirique (ERM) avec régularisation d'entropie relative (ERM-RER) est étudié sous l'hypothèse que la mesure de référence est une mesure  $\sigma$ -finie, et pas nécessairement une mesure de probabilité. Sous cette hypothèse, qui conduit à une généralisation du problème ERM-RER permettant une plus grande flexibilité pour l'incorporation des connaissances antérieures, de nombreuses propriétés pertinentes sont énoncées. Parmi ces propriétés, la solution à ce problème, si elle existe, se révèle être une mesure de probabilité unique, souvent mutuellement absolument continue avec la mesure de référence. Une telle solution présente une garantie probablement-approximativement-correcte pour le problème ERM indépendamment du fait que ce dernier possède ou non une solution. Pour un ensemble de données fixe, le risque empirique s'avère être une variable aléatoire sous-gaussienne lorsque les modèles sont échantillonnés à partir de la solution au problème ERM-RER. Les capacités de généralisation de la solution au problème ERM-RER (l'algorithme de Gibbs) sont étudiées via la sensibilité de la valeur espérée du risque empirique aux déviations d'une telle solution vers des mesures de probabilité alternatives. Enfin, un lien intéressant entre la sensibilité, l'erreur de généralisation et l'information lautum est établi.

**Mots-clés :** Apprentissage Supervisé, Apprentissage PAC, Régularisation, Entropie Relative, Minimisation du Risque Empirique, Mesure de Gibbs, Algorithme de Gibbs, Généralisation, et Sensitivité.

## Contents

<b>1</b>	<b>Introduction</b>	<b>5</b>
<b>2</b>	<b>Empirical Risk Minimization (ERM)</b>	<b>7</b>
2.1	Notation and Main Assumptions . . . . .	8
2.2	Relative Entropy Extended to $\sigma$ -Finite Measures . . . . .	8
2.3	ERM with Relative Entropy Regularization . . . . .	10
2.4	Type-I and Type-II Relative Entropy Regularization . . . . .	10
<b>3</b>	<b>The Solution to the ERM-RER Problem</b>	<b>11</b>
3.1	Examples . . . . .	13
3.1.1	ERM with Discrete Entropy Regularization . . . . .	13
3.1.2	ERM with Differential Entropy Regularization . . . . .	14
3.1.3	Information-Risk Minimization . . . . .	14
3.2	Bounds on the Radon-Nikodym Derivative . . . . .	14
3.3	Asymptotes of the Radon-Nikodym Derivative . . . . .	15
<b>4</b>	<b>Reference Measures</b>	<b>18</b>
4.1	Coherent and Consistent Reference Measures . . . . .	18
4.2	Gibbs Reference Measures . . . . .	21
<b>5</b>	<b>The Log-Partition Function</b>	<b>25</b>
5.1	Separable Empirical Risk Functions . . . . .	25
5.2	Properties of the Log-Partition Function . . . . .	26
<b>6</b>	<b>Expectation of the Empirical Risk</b>	<b>28</b>
<b>7</b>	<b>Variance of the Empirical Risk</b>	<b>31</b>
<b>8</b>	<b>Cumulant Generating Function of the Empirical Risk</b>	<b>35</b>
<b>9</b>	<b>Concentration of Probability</b>	<b>44</b>
9.1	The Limit Set . . . . .	44
9.2	The Nonnegligible Limit Set . . . . .	45
<b>10</b>	<b><math>(\delta, \epsilon)</math>-Optimality</b>	<b>47</b>
<b>11</b>	<b>Sensitivity and Generalization</b>	<b>49</b>
11.1	Sensitivity . . . . .	49
11.2	Generalization Error . . . . .	50
<b>12</b>	<b>Conclusions and Final Remarks</b>	<b>54</b>
	<b>Appendices</b>	<b>55</b>
<b>A</b>	<b>Proof of Theorem 2.2</b>	<b>55</b>

<b>B</b>	<b>Proof of Lemma 3.1</b>	<b>57</b>
<b>C</b>	<b>Proof of Theorem 3.1</b>	<b>58</b>
<b>D</b>	<b>Proof of Lemma 3.2</b>	<b>62</b>
<b>E</b>	<b>Proof of Lemma 3.3</b>	<b>63</b>
<b>F</b>	<b>Proof of Lemma 3.4</b>	<b>63</b>
<b>G</b>	<b>Proof of Lemma 3.6</b>	<b>64</b>
<b>H</b>	<b>Proof of Lemma 3.7</b>	<b>67</b>
<b>I</b>	<b>Proof of Lemma 4.1</b>	<b>68</b>
<b>J</b>	<b>Proof of Theorem 4.1</b>	<b>69</b>
<b>K</b>	<b>Proof of Lemma 5.2</b>	<b>72</b>
<b>L</b>	<b>Proof of Lemma 5.3</b>	<b>74</b>
<b>M</b>	<b>Proof of Lemma 5.4</b>	<b>75</b>
<b>N</b>	<b>Proof of Theorem 6.1</b>	<b>77</b>
	N.1 Preliminaries . . . . .	78
	N.2 The proof . . . . .	80
<b>O</b>	<b>Proof of Lemma 6.1</b>	<b>82</b>
<b>P</b>	<b>Proof of Theorem 6.2</b>	<b>83</b>
<b>Q</b>	<b>Proof of Theorem 9.1</b>	<b>83</b>
<b>R</b>	<b>Proof of Theorem 9.2</b>	<b>85</b>
<b>S</b>	<b>Proof of Lemma 9.2</b>	<b>90</b>
<b>T</b>	<b>Proof of Theorem 11.2</b>	<b>90</b>
<b>U</b>	<b>Proof of Theorem 11.4</b>	<b>93</b>

## 1 Introduction

In statistical machine learning, the problem of empirical risk minimization (ERM) with relative entropy regularization (ERM-RER) has been the workhorse for building probability measures on the set of models, without any additional assumption on the statistical description of the datasets. See for instance [2–4] and [5]. Instead of additional statistical assumptions on the datasets, which are typical in Bayesian methods [6], relative entropy regularization requires a reference probability measure on the set of models, which is external to the ERM problem. Often, such a reference measure represents prior knowledge or side information and is chosen for guiding the search of models towards those inducing low empirical risks with high probability over seen and unseen datasets. From this perspective, the reference measure can be seen as an additional degree of freedom to improve the generalization capabilities of machine learning algorithms based on ERM-RER, e.g. Gibbs algorithms [4, 7–14] and [15]. This new degree of freedom is one of the main motivations for regularizing the ERM problem using relative entropy, or more generally, any  $f$ -divergence regularization, as discussed in [16, 17] and [18]. Beyond probability measures, as shown in this paper, the reference measure can be any  $\sigma$ -finite measure with arbitrary support. The flexibility introduced by this generalization becomes particularly relevant for the case in which priors are available in the form of probability distributions that can be evaluated up to some normalizing factor, cf. [19], or cannot be represented by probability distributions, e.g., equal preferences among elements of infinite countable sets. For some specific choices of  $\sigma$ -finite reference measures, the ERM-RER boils down to particular cases of special interest: (i) the information-risk minimization problem presented in [20]; (ii) the ERM with differential entropy regularization (ERM-DiffER); and (iii) the ERM with discrete entropy regularization (ERM-DisER). See for instance [21] and references therein. From this perspective, the proposed ERM-RER formulation yields a unified mathematical framework that comprises a large class of problems.

When the reference measure is a probability measure, the solution to the ERM-RER problem is known to be unique and correspond to a Gibbs probability measure. Such a Gibbs probability measure has been studied using measure theoretic and information theoretic notions in [8, 20, 22–29]; statistical physics in [2]; PAC (Probably Approximately Correct)-Bayesian learning theory in [30–33]; and proved to be of particular interest in classification problems in [4, 11, 17, 34–36] and [37]. In the general case in which the reference is a  $\sigma$ -finite measure, a solution to the ERM-RER problem does not always exist. Nonetheless, if it exists, it is shown to be a unique Gibbs probability despite the fact that its partition function is defined with respect to a  $\sigma$ -finite measure. The condition for the existence is mild and is always satisfied when the reference measure is a probability measure, as highlighted above. Interestingly, such a solution is mutually absolutely continuous with the reference measure in most practical cases. Interestingly, most of the properties known for the classical ERM-RER problem are shown to hold in the most general case. For instance,



the empirical risk observed when models are sampled from the ERM-RER-optimal probability measure is a sub-Gaussian random variable that exhibits a PAC guarantee for the ERM problem without regularization.

When the solution to the ERM-RER problem is used to sample models to label unseen patterns, the process is known as the Gibbs algorithm. One of the traditional performance metrics to evaluate the generalization capabilities of the Gibbs algorithm is the generalization error. When the reference measure is a probability measure, a closed-form expression for the generalization error of the Gibbs algorithm is presented in [8], while upper bounds have been derived in [15, 20, 27–33, 38–51], and references therein. In this work, a new performance metric coined *sensitivity*, which quantifies the variations of the expected empirical risk due to deviations from the solution of the ERM-RER problem is introduced. The sensitivity is defined as the difference between two quantities: (a) The expectation of the empirical risk with respect to the solution to the ERM-RER problem; and (b) the expectation of the empirical risk with respect to an alternative measure. The absolute value of the sensitivity is shown to be upper bounded by a term that is proportional to the squared-root of the relative entropy of the alternative measure with respect to the ERM-RER-optimal measure. Such bound allows providing lower and upper bounds on the expected empirical risk after a deviation from the ERM-RER-optimal measure towards an alternative probability measure. More interestingly, the expectation (with respect to the probability distribution of the datasets) of the sensitivity to deviations to a specific measure is shown to be equal to the generalization error of the Gibbs algorithm. Using this result, the closed-form expression for the generalization error of the Gibbs algorithm presented in [8] is shown to hold even in the case in which the reference measure is a  $\sigma$ -finite measure. Moreover, the generalization error is shown to be upper bounded by a term that is proportional to the squared-root of the lautum information between the models and the datasets, cf. [52]. This bound is reminiscent of the result in [29, Theorem 1] in which a similar bound is presented using the mutual information instead of the lautum information. While [29, Theorem 1] follows immediately from the variational representation of relative entropy, c.f., [53, Lemma 4.18 (Transportation Lemma)], the new result follows from the fact that the empirical risk when models are sampled from the ERM-RER-optimal probability measure is a sub-Gaussian random variable. Interestingly, the new upper-bound does not require any of the conditions in [29, Theorem 1].

The remainder of this work is organized as follows. Section 2 introduces two optimization problems: the ERM and the ERM-RER. The asymmetry of the relative entropy is analyzed in the context of the ERM-RER and two variants, coined Type-I and Type-II, are distinguished. The former considers the case in which the regularization is the relative entropy of the optimization measure with respect to the reference measure. The latter considers a regularization by the relative entropy of the reference measure with respect to the optimization measure. Section 3 presents the solution to the ERM-RER problem in the general case and introduces its main properties. Section 4 introduces two new classes

of reference measures and the solution of the ERM-RER problem is shown to exhibit different properties for each class. This section ends by studying the ERM-RER problem in the special case in which the reference measure is a Gibbs probability measure. This special case exhibits a solution that is identical to the solution to an ERM-RER problem whose reference measure is the same used to build the above mentioned Gibbs measure. Section 5 studies the properties of the log-partition function of the ERM-RER-optimal probability measure. The first, second, and third cumulants of the empirical risk when the models are sampled from the ERM-RER-optimal measure and the reference measure are respectively characterized. Section 6 and Section 7 study the properties of the expectation and variance of the empirical risk when the models are sampled from the ERM-RER-optimal probability measure. These mean and variance are compared with the mean and variance of the empirical risk when models are sampled from the reference measure. Section 8 introduces several explicit expressions for the cumulant generating function of the empirical risk when the models are sampled from the ERM-RER-optimal measure. Using these equivalent expressions, it is shown that empirical risk is a sub-Gaussian random variable when models are sampled from the ERM-RER-optimal measure. Section 9 describes the monotonic concentration of the ERM-RER-optimal probability measure when the regularization factor tends to zero. Section 10 show that the empirical risk when the models are sampled from the ERM-RER-optimal probability measure exhibits a PAC-type guarantee with respect to the ERM problem without regularization. Finally, Section 11 studies the sensitivity of the expected empirical risk with respect to deviations from the ERM-RER-optimal measure to alternative measures and shows connections with the generalization error and the lautum information. Section 12 ends this work with conclusions and a discussion on the results.

## 2 Empirical Risk Minimization (ERM)

Let  $\mathcal{M}$ ,  $\mathcal{X}$  and  $\mathcal{Y}$ , with  $\mathcal{M} \subseteq \mathbb{R}^d$  and  $d \in \mathbb{N}$ , be sets of *models*, *patterns*, and *labels*, respectively. A pair  $(x, y) \in \mathcal{X} \times \mathcal{Y}$  is referred to as a *labeled pattern* or as a *data point*. Given  $n$  data points, with  $n \in \mathbb{N}$ , denoted by  $(x_1, y_1)$ ,  $(x_2, y_2)$ ,  $\dots$ ,  $(x_n, y_n)$ , the corresponding dataset is represented by the tuple

$$\mathbf{z} = ((x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)) \in (\mathcal{X} \times \mathcal{Y})^n. \quad (1)$$

Let the function  $f : \mathcal{M} \times \mathcal{X} \rightarrow \mathcal{Y}$  be such that the label assigned to the pattern  $x$  according to the model  $\boldsymbol{\theta} \in \mathcal{M}$  is  $f(\boldsymbol{\theta}, x)$ . Let also the function

$$\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow [0, +\infty] \quad (2)$$

be such that given a data point  $(x, y) \in \mathcal{X} \times \mathcal{Y}$ , the risk induced by a model  $\boldsymbol{\theta} \in \mathcal{M}$  is  $\ell(f(\boldsymbol{\theta}, x), y)$ . In the following, the risk function  $\ell$  is assumed to be nonnegative and for all  $y \in \mathcal{Y}$ ,  $\ell(y, y) = 0$ .

The *empirical risk* induced by the model  $\theta$ , with respect to the dataset  $\mathbf{z}$  in (1) is determined by the function  $\mathbf{L}_{\mathbf{z}} : \mathcal{M} \rightarrow [0, +\infty]$ , which satisfies

$$\mathbf{L}_{\mathbf{z}}(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(f(\theta, x_i), y_i). \quad (3)$$

Using this notation, the ERM consists of the following optimization problem:

$$\min_{\theta \in \mathcal{M}} \mathbf{L}_{\mathbf{z}}(\theta). \quad (4)$$

Let the set of solutions to the ERM problem in (4) be denoted by

$$\mathcal{T}(\mathbf{z}) \triangleq \arg \min_{\theta \in \mathcal{M}} \mathbf{L}_{\mathbf{z}}(\theta). \quad (5)$$

Note that if the set  $\mathcal{M}$  is finite, the ERM problem in (4) always possesses a solution, and thus,  $|\mathcal{T}(\mathbf{z})| > 0$ . Nonetheless, in general, the ERM problem might not necessarily possess a solution, i.e.,  $|\mathcal{T}(\mathbf{z})| = 0$ .

## 2.1 Notation and Main Assumptions

In the following, given a measurable space  $(\Omega, \mathcal{F})$ , the notation  $\Delta(\Omega, \mathcal{F})$  is used to represent the set of  $\sigma$ -finite measures that can be defined over  $(\Omega, \mathcal{F})$ . Given a measure  $Q \in \Delta(\Omega, \mathcal{F})$ , the subset  $\Delta_Q(\Omega, \mathcal{F})$  of  $\Delta(\Omega, \mathcal{F})$  contains all  $\sigma$ -finite measures that are absolutely continuous with respect to the measure  $Q$ . Alternatively, the subset  $\nabla_Q(\Omega, \mathcal{F})$  of  $\Delta(\Omega, \mathcal{F})$  contains all probability measures  $P$  such that  $Q$  is absolutely continuous with respect to  $P$ . Given a set  $\mathcal{A} \subset \mathbb{R}^d$ , the Borel  $\sigma$ -field over  $\mathcal{A}$  is denoted by  $\mathcal{B}(\mathcal{A})$ .

The main assumption adopted in this work is that the function  $\mathbf{L}_{\mathbf{z}}$  in (3) is measurable with respect to the Borel measurable spaces  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$  and  $([0, +\infty], \mathcal{B}([0, +\infty]))$ .

## 2.2 Relative Entropy Extended to $\sigma$ -Finite Measures

In this work, the *relative entropy*, which is usually defined for probability measures, is extended to  $\sigma$ -finite measures.

**Definition 2.1** (Generalized Relative Entropy). *Given two  $\sigma$ -finite measures  $P$  and  $Q$  on the same measurable space, such that  $P$  is absolutely continuous with respect to  $Q$ , the relative entropy of  $P$  with respect to  $Q$  is*

$$D(P\|Q) = \int \frac{dP}{dQ}(x) \log \left( \frac{dP}{dQ}(x) \right) dQ(x), \quad (6)$$

where the function  $\frac{dP}{dQ}$  is the Radon-Nikodym derivative of  $P$  with respect to  $Q$ .

The relative entropy exhibits a property often referred to as the *information inequality* [54, Theorem 2.6.3] in the case of probability measures on  $(\Omega, \mathcal{F})$ , with  $\Omega$  a countable set. The following theorem explores this property in a more general scenario.

**Theorem 2.1.** *If  $P$  and  $Q$  are both probability measures on a general measurable space  $(\Omega, \mathcal{F})$ , then,*

$$D(P\|Q) \geq 0, \quad (7)$$

*with equality if and only if  $P$  and  $Q$  are identical.*

*Proof:* Consider the function  $f : [0, \infty) \rightarrow \mathbb{R}$  such that for all  $x \in (0, +\infty)$ ,  $f(x) = x \log(x)$  and  $f(0) = 0$ . Note that  $f$  is strictly convex. If  $P$  and  $Q$  are both probability measures on the measurable space  $(\Omega, \mathcal{F})$ , the following holds:

$$D(P\|Q) = \int \frac{dP}{dQ}(x) \log \left( \frac{dP}{dQ}(x) \right) dQ(x) \quad (8)$$

$$= \int f \left( \frac{dP}{dQ}(x) \right) dQ(x) \quad (9)$$

$$\geq f \left( \int \frac{dP}{dQ}(x) dQ(x) \right) \quad (10)$$

$$= f(1) \quad (11)$$

$$= 0, \quad (12)$$

where the inequality (11) follows from Jensen's inequality [55, Section 6.3.5]. Equality in (11) holds if and only if for all  $x \in \text{supp } Q$ ,  $\frac{dP}{dQ}(x) = 1$ , which implies that both  $P$  and  $Q$  are identical. This completes proof. ■

If  $Q$  is not a probability measure, then it might be observed that  $D(P\|Q) < 0$ . Consider for instance the case in which  $P$  is a zero-mean Gaussian probability measure with variance  $\sigma^2$  and  $Q$  is the Lebesgue measure on  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ . Hence, the Radon-Nikodym derivative  $\frac{dP}{dQ}$  is the Gaussian probability density function such that for all  $x \in \mathbb{R}$ ,

$$\frac{dP}{dQ}(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left( -\frac{x^2}{2\sigma^2} \right). \quad (13)$$

Under this assumption, the relative entropy of  $P$  with respect to  $Q$  is the negative of the differential entropy of  $P$ . That is,

$$D(P\|Q) = -\frac{1}{2} \log(2\pi\epsilon\sigma^2), \quad (14)$$

with  $\epsilon$  being Néper's constant. See for instance [54, Example 8.1.2]. Hence,  $D(P\|Q)$  is negative for all  $\sigma^2 \in (\frac{1}{2\pi\epsilon}, +\infty)$  and nonnegative for all  $\sigma^2 \in (0, \frac{1}{2\pi\epsilon}]$ . Finally, note also that

$$\lim_{\sigma^2 \rightarrow 0} D(P\|Q) = +\infty, \text{ and} \quad (15)$$

$$\lim_{\sigma^2 \rightarrow +\infty} D(P\|Q) = -\infty. \quad (16)$$

A central observation from (14) is that the equality  $D(P\|Q) = 0$  does not necessarily imply that  $P$  and  $Q$  are identical measures. For instance, when

$\sigma^2 = \frac{1}{2\pi\epsilon}$  in (15), it holds that  $D(P\|Q) = 0$ , while  $P$  is a Gaussian probability measure and  $Q$  is the Lebesgue measure.

The following property, known for the case of probability measures as the *joint-convexity of the relative entropy*, is extended by the following theorem.

**Theorem 2.2.** *Let  $P_1$  and  $P_2$  be two probability measures and  $Q_1$  and  $Q_2$  be two  $\sigma$ -finite measures, all on the same measurable space. For all  $i \in \{1, 2\}$ , let  $P_i$  be absolutely continuous with respect to  $Q_i$ . Then, for all  $\lambda \in [0, 1]$ ,*

$$\begin{aligned} D(\lambda P_1 + (1 - \lambda)P_2 \| \lambda Q_1 + (1 - \lambda)Q_2) \\ \leq \lambda D(P_1 \| Q_1) + (1 - \lambda)D(P_2 \| Q_2). \end{aligned} \quad (17)$$

Equality in (17) holds if and only if  $P_1 = P_2$  and  $Q_1 = Q_2$ .

*Proof:* The proof is presented in Appendix A. ■

### 2.3 ERM with Relative Entropy Regularization

Given a dataset, the *expected empirical risk* induced by a probability measure  $P \in \Delta(\mathcal{M}, \mathcal{B}(\mathcal{M}))$  is defined as follows.

**Definition 2.2** (Expected Empirical Risk). *Let  $P$  be a probability measure in  $\Delta(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ . The expected empirical risk with respect to the dataset  $\mathbf{z}$  in (1) induced by the measure  $P$  is*

$$R_{\mathbf{z}}(P) = \int L_{\mathbf{z}}(\boldsymbol{\theta}) dP(\boldsymbol{\theta}), \quad (18)$$

where the function  $L_{\mathbf{z}}$  is in (3).

The ERM-RER problem is parametrized by a  $\sigma$ -finite measure in  $\Delta(\mathcal{M}, \mathcal{B}(\mathcal{M}))$  and a positive real, which are referred to as the *reference measure* and the *regularization factor*, respectively. Let  $Q \in \Delta(\mathcal{M}, \mathcal{B}(\mathcal{M}))$  be a  $\sigma$ -finite measure and let  $\lambda$  be a positive real. The ERM-RER problem, with parameters  $Q$  and  $\lambda$ , consists of the following optimization problem:

$$\min_{P \in \Delta_Q(\mathcal{M}, \mathcal{B}(\mathcal{M}))} R_{\mathbf{z}}(P) + \lambda D(P\|Q), \quad (19a)$$

$$\text{s. t.} \quad \int dP(\boldsymbol{\theta}) = 1, \quad (19b)$$

where the dataset  $\mathbf{z}$  is in (1), and the functional  $R_{\mathbf{z}}$  is defined in (18).

### 2.4 Type-I and Type-II Relative Entropy Regularization

The optimization problem in (19) is coined Type-I ERM-RER in [56] in the aim of distinguishing it from the optimization problem

$$\min_{P \in \nabla_Q(\mathcal{M}, \mathcal{B}(\mathcal{M}))} R_{\mathbf{z}}(P) + \lambda D(Q\|P), \quad (20a)$$

$$\text{s. t.} \quad \int dP(\boldsymbol{\theta}) = 1, \quad (20b)$$

which is coined Type-II ERM-RER.

The Type-II ERM-RER problem in (20), when  $Q$  is a probability measure, exhibits a solution that is identical to the solution to the following Type-I ERM-RER problem [56, Theorem 1]:

$$\min_{P \in \Delta_Q(\mathcal{M}, \mathcal{B}(\mathcal{M}))} \int \log(\beta + \mathbf{L}_{\mathbf{z}}(\boldsymbol{\nu})) dP(\boldsymbol{\nu}) + D(P \| Q), \quad (21a)$$

$$\text{s. t.} \quad \int dP(\boldsymbol{\theta}) = 1, \quad (21b)$$

where  $\beta$  is a constant chosen to satisfy

$$\int \frac{\lambda}{\beta + \mathbf{L}_{\mathbf{z}}(\boldsymbol{\nu})} dQ(\boldsymbol{\nu}) = 1. \quad (21c)$$

Essentially, by appropriately transforming the objective function, an equivalence can be established between Type-I and Type-II ERM-RER problems. Hence, without loss of generality, the remainder of this work focuses exclusively on Type-I ERM-RER, which is simply referred to as ERM-RER.

### 3 The Solution to the ERM-RER Problem

The solution to the ERM-RER problem in (19) is presented in terms of two objects. First, the function  $K_{Q,\mathbf{z}} : \mathbb{R} \rightarrow \mathbb{R} \cup \{+\infty\}$  such that for all  $t \in \mathbb{R}$ ,

$$K_{Q,\mathbf{z}}(t) = \log \left( \int \exp(t \mathbf{L}_{\mathbf{z}}(\boldsymbol{\theta})) dQ(\boldsymbol{\theta}) \right), \quad (22)$$

with  $\mathbf{L}_{\mathbf{z}}$  in (3). Second, the set  $\mathcal{K}_{Q,\mathbf{z}} \subset (0, +\infty)$ , which is defined by

$$\mathcal{K}_{Q,\mathbf{z}} \triangleq \left\{ s \in (0, +\infty) : K_{Q,\mathbf{z}} \left( -\frac{1}{s} \right) < +\infty \right\}. \quad (23)$$

The notation for the function  $K_{Q,\mathbf{z}}$  and the set  $\mathcal{K}_{Q,\mathbf{z}}$  are chosen such that their parametrization by (or dependence on) the dataset  $\mathbf{z}$  in (1) and the  $\sigma$ -finite measure  $Q$  in (19) are highlighted.

The following lemma describes the set  $\mathcal{K}_{Q,\mathbf{z}}$ .

**Lemma 3.1.** *The set  $\mathcal{K}_{Q,\mathbf{z}}$  in (23) is a convex subset of  $\mathbb{R}$ . If the measure  $Q$  in (19) is a probability measure, then, the set  $\mathcal{K}_{Q,\mathbf{z}}$  in (23) satisfies*

$$\mathcal{K}_{Q,\mathbf{z}} = (0, +\infty). \quad (24)$$

*Proof:* The proof is presented in Appendix B. ■

Using this notation, the solution to the ERM-RER problem in (19) is presented by the following theorem.

**Theorem 3.1.** *If  $\lambda \in \mathcal{K}_{Q,z}$ , with  $\mathcal{K}_{Q,z}$  in (23), the solution to the optimization problem in (19) is a unique probability measure, denoted by  $P_{\Theta|Z=z}^{(Q,\lambda)}$ , which satisfies for all  $\theta \in \text{supp } Q$ ,*

$$\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) = \exp\left(-K_{Q,z}\left(-\frac{1}{\lambda}\right) - \frac{1}{\lambda}L_z(\theta)\right), \quad (25)$$

where the function  $L_z$  is defined in (3) and the function  $K_{Q,z}$  is defined in (22).

*Proof:* The proof is presented in Appendix C. ■

Contrary to the ERM problem in (4), which does not necessarily possess a solution, the ERM-RER problem in (19) always possess a solution when  $Q$  is a probability measure. This is essentially because the set  $\mathcal{K}_{Q,z}$  is the set of all positive reals (Lemma 3.1), and thus, the condition  $\lambda \in \mathcal{K}_{Q,z}$  is always verified. On the contrary, when  $Q$  is a  $\sigma$ -finite measure, the solution to the ERM-RER problem in (19) depends on whether  $\lambda \in \mathcal{K}_{Q,z}$ . If the solution exists, it is  $P_{\Theta|Z=z}^{(Q,\lambda)}$  in (25), which is a unique probability measure referred to as the Gibbs measure [57]. The function  $K_{Q,z}$  is often referred to as the *log-partition function*, see for instance, [58, Section 7.3.1].

The following lemma shows that the Radon-Nikodym derivative in (25) is both nonnegative and finite.

**Lemma 3.2.** *The Radon-Nikodym derivative  $\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}$  in (25) satisfies for all  $\theta \in \text{supp } Q$  that*

$$0 \leq \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) < +\infty, \quad (26)$$

where the equality  $\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) = 0$  holds if and only if  $L_z(\theta) = +\infty$ .

*Proof:* The proof is presented in Appendix D. ■

An immediate consequence of Lemma 3.2 is the equality

$$P_{\Theta|Z=z}^{(Q,\lambda)}(\{\theta \in \mathcal{M} : L_z(\theta) = +\infty\}) = 0.$$

Theorem 3.1 shows that the probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  is absolutely continuous with respect to the measure  $Q$ . The following lemma shows that the converse is also true if and only if the set of models that lead to an infinite empirical risk exhibit zero measure with respect to the reference measure  $Q$ .

**Lemma 3.3.** *The  $\sigma$ -finite measure  $Q$  and the probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  in (25) are mutually absolutely continuous if and only if*

$$Q(\{\theta \in \mathcal{M} : L_z(\theta) = +\infty\}) = 0. \quad (27)$$

*Proof:* The proof is presented in Appendix E. ■

The relevance of Lemma 3.3 is that it shows that if  $\lambda \in \mathcal{K}_{Q,z}$ , the collection of negligible sets with respect to the measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  in (25) is identical to the collection of negligible sets with respect to the measure  $Q$  in (19), under the assumption in (27). Such an assumption is trivially true when the function  $\ell$  in (2) is bounded.

The following lemma shows that the negligible sets with respect to the measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  in (25) are invariant with respect to  $\lambda$ .

**Lemma 3.4.** *For all  $(\alpha, \beta) \in \mathcal{K}_{Q,z} \times \mathcal{K}_{Q,z}$ , with  $\mathcal{K}_{Q,z}$  in (23), assume that the probability measures  $P_{\Theta|Z=z}^{(Q,\alpha)}$  and  $P_{\Theta|Z=z}^{(Q,\beta)}$  satisfy (25) with  $\lambda = \alpha$  and  $\lambda = \beta$ , respectively. Then,  $P_{\Theta|Z=z}^{(Q,\alpha)}$  and  $P_{\Theta|Z=z}^{(Q,\beta)}$  are mutually absolutely continuous.*

*Proof:* The proof is presented in Appendix F. ■

Particular assumptions on the set  $\mathcal{M}$  and the reference measure  $Q$  lead to well-known instances of the ERM-RER problem in (19), as discussed hereunder.

### 3.1 Examples

Three examples are of particular interest: (a) The set  $\mathcal{M} \subset \mathbb{R}^d$  is countable and the measure  $Q$  is the counting measure in  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ , which leads to the ERM-DisER problem; (b) The set  $\mathcal{M}$  is an uncountable subset of  $\mathbb{R}^d$ , and  $Q$  is the Lebesgue measure on  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ , which leads to the ERM-Differ problem; and (c) The set  $\mathcal{M}$  and the measure  $Q$  form a Borel probability measure space  $(\mathcal{M}, \mathcal{B}(\mathcal{M}), Q)$ , which leads to the information-risk minimization problem.

#### 3.1.1 ERM with Discrete Entropy Regularization

When the set  $\mathcal{M} \subset \mathbb{R}^d$  is countable and the  $\sigma$ -finite measure  $Q$  in (19) is the counting measure in  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ , given a probability measure  $P$  on the same measurable space, the Radon-Nikodym derivative  $\frac{dP}{dQ}$  is a probability mass function, denoted by  $p$ . Thus, the relative entropy  $D(P||Q)$  is equivalent to the negative of the discrete entropy induced by  $p$  [54, Chapter 2], denoted by  $H(p)$ . In this case, the ERM-RER in (19) can be re-written as the following ERM-DisER problem:

$$\min_p \sum_{\theta \in \mathcal{M}} \mathbf{L}_z(\theta) p(\theta) - \lambda H(p), \quad (28)$$

where the optimization domain in (28) is the set of probability mass functions that can be defined over the measure space  $\Delta(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ . In this special case, the probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  in (25) whose probability mass function is the



solution to the ERM-DisER problem in (28) satisfies

$$\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) = \frac{\exp\left(-\frac{\mathbf{L}_z(\theta)}{\lambda}\right)}{\sum_{\nu \in \mathcal{M}} \exp\left(-\frac{\mathbf{L}_z(\nu)}{\lambda}\right)}, \quad (29)$$

which describes the discrete Gibbs probability measure on  $\Delta(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ , with temperature parameter  $\lambda$ , and energy function  $\mathbf{L}_z$  in (3).

### 3.1.2 ERM with Differential Entropy Regularization

When  $\mathcal{M} \subseteq \mathbb{R}^d$  is uncountable and the  $\sigma$ -finite measure  $Q$  in (19) is the Lebesgue measure in  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ , for all probability measures  $P \in \Delta_Q(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ , the Radon-Nikodym derivative  $\frac{dP}{dQ}$  is a probability density function, denoted by  $g$ . Thus, the relative entropy  $D(P||Q)$  is equivalent to the negative of the differential entropy induced by  $g$  [54, Chapter 8], denoted by  $h(g)$ . In this special case, the ERM-RER in (19) can be re-written as the following ERM-Differ problem:

$$\min_g \int_{\mathcal{M}} \mathbf{L}_z(\theta) g(\theta) d\theta - \lambda h(g), \quad (30)$$

where the optimization domain in (30) is the set of probability density functions that can be defined over the measure space  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ . The probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  in (25) whose probability density function is the solution to the ERM-RER problem in (30) satisfies

$$\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) = \frac{\exp\left(-\frac{\mathbf{L}_z(\theta)}{\lambda}\right)}{\int_{\mathcal{M}} \exp\left(-\frac{\mathbf{L}_z(\nu)}{\lambda}\right) d\nu}, \quad (31)$$

which describes the absolutely continuous Gibbs probability measure with temperature parameter  $\lambda$  and energy function  $\mathbf{L}_z$  in (3).

Both, the ERM-Differ and ERM-DisER problems are closely related to those typically arising while using Jayne's maximum entropy principle [59, 60] for classification problems such as those in [34–36], and [61].

### 3.1.3 Information-Risk Minimization

When  $Q$  is a probability measure, the ERM-RER in (19) is equivalent to the *information-risk minimization* (IRM) problem in [20]. The IRM problem in (19) is known to possess a unique solution equal to the Gibbs probability measure in (25), as independently shown in [20, 29, 57, 62, 63] and [64].

## 3.2 Bounds on the Radon-Nikodym Derivative

The Radon-Nikodym derivative  $\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}$  in (25) is bigger for models inducing smaller empirical risks, as shown by the following corollary of Theorem 3.1.

**Corollary 3.1.** *The Radon-Nikodym derivative  $\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}$  in (25) satisfies for all  $(\theta_1, \theta_2) \in \text{supp } Q \times \text{supp } Q$ , with  $L_z(\theta_2) \leq L_z(\theta_1)$ , that*

$$\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta_1) \leq \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta_2), \quad (32)$$

with equality if and only if  $L_z(\theta_1) = L_z(\theta_2)$ .

The intuition that follows from corollary 3.1 is that under the assumption that the ERM problem in (4) possesses a solution in the support of the reference measure, i.e.,  $\mathcal{T}(z) \cap \text{supp } Q$  is not empty, with  $\mathcal{T}(z)$  in (5), the maximum of the function  $\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}$  in (25) is achieved by the models in  $\mathcal{T}(z) \cap \text{supp } Q$ . When the Radon-Nikodym derivative  $\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}$  in (25) is either the probability mass function in (29) or the probability density function in (31), Corollary 3.1 shows that the elements of the set  $\mathcal{T}(z) \cap \text{supp } Q$  are the *modes* of the corresponding probability density function or probability mass function.

### 3.3 Asymptotes of the Radon-Nikodym Derivative

The following lemma describes the asymptotic behavior of the Radon-Nikodym derivative  $\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}$  in (25) when the regularisation factor increases, i.e.,  $\lambda \rightarrow +\infty$  and the reference measure  $Q$  is a probability measure.

**Lemma 3.5.** *Let the measure  $Q$  in (19) be a probability measure. Then, for all  $\theta \in \text{supp } Q$ , the Radon-Nikodym derivative  $\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}$  in (25) satisfies*

$$\lim_{\lambda \rightarrow +\infty} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) = 1. \quad (33)$$

*Proof:* From Theorem 3.1, it follows that for all  $\theta \in \text{supp } Q$ ,

$$\lim_{\lambda \rightarrow +\infty} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) = \lim_{\lambda \rightarrow +\infty} \frac{\exp\left(-\frac{L_z(\theta)}{\lambda}\right)}{\int \exp\left(-\frac{L_z(\nu)}{\lambda}\right) dQ(\nu)} \quad (34)$$

$$= \frac{1}{\int dQ(\nu)} \quad (35)$$

$$= 1, \quad (36)$$

where the function  $L_z$  is defined in (3). This completes the proof. ■

Lemma 3.5 unveils the fact that, when  $Q$  is a probability measure, in the limit when  $\lambda \rightarrow +\infty$ , both probability measures  $P_{\Theta|Z=z}^{(Q,\lambda)}$  and  $Q$  are identical. This is consistent with the fact that when  $\lambda$  tends to infinity, the optimization problem

in (19) boils down to exclusively minimizing the relative entropy. Such minimum is zero and is observed when both probability measures  $P_{\Theta|Z=z}^{(Q,\lambda)}$  and  $Q$  are identical (Theorem 2.1). Such intuition breaks when the reference measure is a  $\sigma$ -finite measure, but not a probability measure. In such a case, the relative entropy term in (19) might be negative and a minimum might not exist. See for instance, the case of the relative entropy between a Gaussian measure and the Lebesgue measure in (14), which satisfies (16).

The limit of the Radon-Nikodym derivative  $\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}$  in (25), when  $\lambda$  tends to zero from the right, can be studied using the following set

$$\mathcal{L}_z(\delta) \triangleq \{\theta \in \mathcal{M} : \mathbb{L}_z(\theta) \leq \delta\}, \quad (37)$$

where the function  $\mathbb{L}_z$  is defined in (3) and  $\delta \in [0, +\infty)$ . In particular consider the nonnegative real

$$\delta_{Q,z}^* \triangleq \inf \{\delta \in [0, +\infty) : Q(\mathcal{L}_z(\delta)) > 0\}. \quad (38)$$

Let also  $\mathcal{L}_{Q,z}^*$  be the following level set of the empirical risk function  $\mathbb{L}_z$  in (3):

$$\mathcal{L}_{Q,z}^* \triangleq \{\theta \in \text{supp } Q : \mathbb{L}_z(\theta) = \delta_{Q,z}^*\}. \quad (39)$$

Using this notation, the limit of the Radon-Nikodym derivative  $\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}$  in (25), when  $\lambda$  tends to zero from the right, is described by the following lemma.

**Lemma 3.6.** *If  $Q(\mathcal{L}_{Q,z}^*) > 0$ , with the set  $\mathcal{L}_{Q,z}^*$  in (39) and  $Q$  the  $\sigma$ -finite measure in (19), then for all  $\theta \in \text{supp } Q$ , the Radon-Nikodym derivative  $\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}$  in (25) satisfies*

$$\lim_{\lambda \rightarrow 0^+} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) = \frac{1}{Q(\mathcal{L}_{Q,z}^*)} \mathbb{1}_{\{\theta \in \mathcal{L}_{Q,z}^*\}}. \quad (40)$$

*Alternatively, if  $Q(\mathcal{L}_{Q,z}^*) = 0$ . Then, for all  $\theta \in \text{supp } Q$ ,*

$$\lim_{\lambda \rightarrow 0^+} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) = \begin{cases} +\infty & \text{if } \theta \in \mathcal{L}_{Q,z}^* \\ 0 & \text{otherwise.} \end{cases} \quad (41)$$

*Proof:* The proof is presented in Appendix G. ■

Consider that  $Q(\mathcal{L}_{Q,z}^*) > 0$ , with  $\mathcal{L}_{Q,z}^*$  in (39). Under this assumption, from Lemma 3.6, it holds that the probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  asymptotically concentrates on the set  $\mathcal{L}_{Q,z}^*$  when  $\lambda$  tends to zero from the right. More specifically,

note that for all measurable sets  $\mathcal{A} \subseteq \mathcal{L}_{Q,z}^* \cap \text{supp } Q$ , it holds that

$$\lim_{\lambda \rightarrow 0^+} P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{A}) = \lim_{\lambda \rightarrow 0^+} \int_{\mathcal{A}} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) dQ(\theta) \quad (42)$$

$$= \int \lim_{\lambda \rightarrow 0^+} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) \mathbb{1}_{\{\theta \in \mathcal{A}\}} dQ(\theta) \quad (43)$$

$$= \int \frac{1}{Q(\mathcal{L}_{Q,z}^*)} \mathbb{1}_{\{\theta \in \mathcal{L}_{Q,z}^*\}} \mathbb{1}_{\{\theta \in \mathcal{A}\}} dQ(\theta) \quad (44)$$

$$= \frac{1}{Q(\mathcal{L}_{Q,z}^*)} \int \mathbb{1}_{\{\theta \in \mathcal{A}\}} dQ(\theta) \quad (45)$$

$$= \frac{Q(\mathcal{A})}{Q(\mathcal{L}_{Q,z}^*)}, \quad (46)$$

where the equality in (43) follows from Lemma 3.2 and the dominated convergence theorem [55, Theorem 2.6.9]. The equality in (44) follows from Lemma 3.6. In the particular case in which  $\mathcal{A} = \mathcal{L}_{Q,z}^*$  in (46), it holds that  $\lim_{\lambda \rightarrow 0^+} P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_{Q,z}^*) = 1$ , which verifies the asymptotic concentration of the probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  on the set  $\mathcal{L}_{Q,z}^*$ .

Another interesting observation is that the Radon-Nikodym derivative  $\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}$  in (25) is a constant among the elements of the set  $\mathcal{L}_{Q,z}^*$ . This can be assimilated to a uniform distribution of the probability among the elements of the set  $\mathcal{L}_{Q,z}^*$  in the limit when  $\lambda$  tends to zero from the right, as previously highlighted in [22–24] and [25]. This becomes more evident in the case in which the set  $\mathcal{M}$  is finite and  $Q$  is the counting measure. In such a case, the asymptotic probability of each of the elements in  $\mathcal{L}_{Q,z}^*$  when  $\lambda$  tends to zero from the right is  $\frac{1}{|\mathcal{L}_{Q,z}^*|}$ .

Consider now that  $Q(\mathcal{L}_{Q,z}^*) = 0$ , with  $\mathcal{L}_{Q,z}^*$  in (39). Under this assumption, in the asymptotic regime when  $\lambda \rightarrow 0$ , the measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  is not a probability measure but either the trivial measure or the infinite measure. This is typically the case in which  $\mathcal{M} = \mathbb{R}^d$ , the measure  $Q$  is absolutely continuous with respect to the Lebesgue measure, and the solution to the ERM problem in (4) has a unique solution on the support of  $Q$ , i.e.,  $\mathcal{L}_{Q,z}^* = \mathcal{T}(z)$  and  $|\mathcal{T}(z)| = 1$ , which implies  $Q(\mathcal{L}_{Q,z}^*) = 0$ .

An interesting question, which is left out of the scope of this paper, is the rate at which  $P_{\Theta|Z=z}^{(Q,\lambda)}$  converges to such limiting measure. The interested reader is referred to [22, 25], and references therein.

The following lemma shows that independently of whether the set  $\mathcal{L}_{Q,z}^*$  is negligible with respect to the measure  $Q$ , the limit when  $\lambda$  tends to zero from the right of  $P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_{Q,z}^*)$  is equal to one.

**Lemma 3.7.** *The measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  in (25) and the set  $\mathcal{L}_{Q,z}^*$  in (39) satisfy,*

$$\lim_{\lambda \rightarrow 0^+} P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_{Q,z}^*) = 1. \quad (47)$$

*Proof:* The proof is presented in Appendix H. ■

Note that if the ERM problem in (4) possesses at least one solution and such solution is within the support of the measure  $Q$ , i.e.,  $\mathcal{T}(z) \cap \text{supp } Q \neq \emptyset$ , then, when  $\lambda$  tends to zero from the right, the probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  asymptotically concentrates on the solution (or the set of solutions within the support of  $Q$ ) to the ERM problem in (4). Alternatively, in the case in which  $\mathcal{L}_{Q,z}^* \cap \mathcal{T}(z) = \emptyset$ , when  $\lambda$  tends to zero from the right, the probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  asymptotically concentrates on a set that does not contain the set of solutions to the ERM problem in (4). This observation leads to the introduction to two new classes of reference measures, namely, *coherent* and *consistent* measures, in the following section.

## 4 Reference Measures

This section introduces two classes of reference measures, namely *coherent* and *consistent* measures, and discusses the special case of Gibbs reference measures.

### 4.1 Coherent and Consistent Reference Measures

A class of reference measures of particular importance to establish connections between the set of solutions to the ERM problem in (4) and the solution to the ERM-RER problem in (19) is that of *coherent* measures. Let  $\rho^* \geq 0$  be the infimum of the empirical risk  $\mathcal{L}_z$  in (3). That is,

$$\rho^* \triangleq \inf\{\mathcal{L}_z(\theta) : \theta \in \mathcal{M}\}. \quad (48)$$

Using this notation, coherent measures are defined as follows.

**Definition 4.1** (Coherent Measures). *The  $\sigma$ -finite measure  $Q$  in (19) is said to be coherent if, for all  $\delta \in (\rho^*, +\infty)$ , with  $\rho^*$  in (48), it holds that*

$$Q(\mathcal{L}_z(\delta)) > 0, \quad (49)$$

where the set  $\mathcal{L}_z(\delta)$  is defined in (37).

When the reference measure  $Q$  in the EMR-RER problem in (19) is a coherent measure, it holds that for all  $\delta > \rho^*$ , the set  $\mathcal{L}_z(\delta)$  in (37) exhibits positive probability with respect to the probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  in (25). The following lemma highlights this observation.

**Lemma 4.1.** *The probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  in (25) satisfies for all  $\delta \in (\rho^*, +\infty)$ , with  $\rho^*$  in (48), that*

$$P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_z(\delta)) > 0, \quad (50)$$

with  $\mathcal{L}_z(\delta)$  in (37), if and only if the  $\sigma$ -finite measure  $Q$  in (19) is coherent.

*Proof:* The proof is presented in Appendix I. ■

Under the assumption that the ERM problem in (4) possesses a solution, it holds that

$$\min_{\theta \in \mathcal{M}} \mathbb{L}_z(\theta) = \inf\{\mathbb{L}_z(\theta) : \theta \in \mathcal{M}\}. \quad (51)$$

Hence, when the  $\sigma$ -finite measure  $Q$  in (19) is coherent, then

$$\delta_{Q,z}^* = \rho^*, \quad (52)$$

with  $\delta_{Q,z}^*$  in (38) and  $\rho^*$  in (48), which implies that

$$\mathcal{L}_{Q,z}^* \subseteq \mathcal{T}(z), \quad (53)$$

with  $\mathcal{T}(z)$  in (5) and  $\mathcal{L}_{Q,z}^*$  in (39). This observation, together with Lemma 3.7, leads to the following result.

**Lemma 4.2.** *Assume that the ERM problem in (4) possesses a solution. Then, the probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  in (25) and the sets  $\mathcal{T}(z)$  in (5) and  $\mathcal{L}_{Q,z}^*$  in (39) satisfy*

$$\lim_{\lambda \rightarrow 0^+} P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_{Q,z}^* \cap \mathcal{T}(z)) = 1, \quad (54)$$

if and only if the  $\sigma$ -finite measure  $Q$  in (19) is coherent.

*Proof:* The proof follows by observing that if  $Q$  is a coherent measure and the ERM problem in (4) possesses a solution, the inclusion in (53) holds. Thus, from Lemma 3.7, the equality in (54) holds. Alternatively, when the measure  $Q$  in (19) is noncoherent, then  $\delta_{Q,z}^* > \rho^*$ , which implies that  $\mathcal{L}_{Q,z}^* \cap \mathcal{T}(z) = \emptyset$ . Hence, from Lemma 3.7, it follows that

$$\lim_{\lambda \rightarrow 0^+} P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_{Q,z}^* \cap \mathcal{T}(z)) = 0, \quad (55)$$

and completes the proof. ■

The relevance of coherent measures in ERM-RER problems is well highlighted by Lemma 4.2. Essentially, when the ERM problem in (4) possesses at least one solution, the concentration of the probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  in (25) on the set (or a subset) of solutions to the ERM problem in (4) occurs asymptotically when  $\lambda$  tends to zero from the right, if only if the reference measure  $Q$  in (19) is coherent. Nonetheless, such asymptotic concentration is not a guarantee that for strictly positive values of  $\lambda$  in (19), the set  $\mathcal{T}(z)$  in (5) and the measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  in (25) satisfy  $P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{T}(z)) > 0$ . In order to ensure this, another class of reference measures, known as *consistent measures*, is introduced.

**Definition 4.2** (Consistent Measure). *The  $\sigma$ -finite measure  $Q$  in (19) is said to be consistent if  $Q(\mathcal{L}_{Q,z}^*) > 0$ , with  $\mathcal{L}_{Q,z}^*$  in (39).*

Note that every consistent measure is not necessarily coherent. For instance, if  $Q$  is consistent but  $\delta_{Q,z}^* > \rho^*$ , with  $\rho^*$  in (48) and  $\delta_{Q,z}^*$  in (38), then, for all  $\delta \in (\rho^*, \delta_{Q,z}^*)$ , it follows that  $Q(\mathcal{L}_z(\delta)) = 0$ , and thus,  $Q$  is not coherent. Alternatively, every coherent measure is not necessarily consistent. For instance, if  $|\mathcal{L}_{Q,z}^*| = 1$  and  $Q$  is coherent and absolutely continuous with respect to the Lebesgue measure, it follows that  $Q(\mathcal{L}_{Q,z}^*) = 0$ , and thus,  $Q$  is not consistent.

The relevance of consistent measures is highlighted by the following lemma.

**Lemma 4.3.** *The probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  in (25) and the set  $\mathcal{L}_{Q,z}^*$  in (39) satisfy*

$$P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_{Q,z}^*) > 0, \quad (56)$$

*if and only if the  $\sigma$ -finite measure  $Q$  in (19) is consistent.*

*Proof:* When  $Q$  is nonconsistent, it holds that  $Q(\mathcal{L}_{Q,z}^*) = 0$  and thus, from the fact that the measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  in (25) is absolutely continuous with respect to  $Q$ , it holds that  $P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_{Q,z}^*) = 0$ . When  $Q$  is consistent, it holds that  $Q(\mathcal{L}_{Q,z}^*) > 0$ . Moreover, for all  $\theta \in \mathcal{L}_{Q,z}^*$ , it holds that  $L_z(\theta) < +\infty$  and thus, from Lemma 3.2, it follows that  $\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) > 0$ . Hence,

$$P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_{Q,z}^*) = \int_{\mathcal{L}_{Q,z}^*} dP_{\Theta|Z=z}^{(Q,\lambda)}(\theta) \quad (57)$$

$$= \int_{\mathcal{L}_{Q,z}^*} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) dQ(\theta) > 0, \quad (58)$$

which completes the proof. ■

The following lemma highlights a central property of consistent measures when the ERM problem in (4) possesses a solution.

**Lemma 4.4.** *Assume that the ERM problem in (4) possesses a solution in the support of  $Q$ . The probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  in (25) and the sets  $\mathcal{T}(z)$  in (5) and  $\mathcal{L}_{Q,z}^*$  in (39) satisfy*

$$P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_{Q,z}^* \cap \mathcal{T}(z)) > 0, \quad (59)$$

*if and only if the  $\sigma$ -finite measure  $Q$  in (19) is consistent.*

*Proof:* The proof follows from Lemma 4.3 by noticing that when the ERM problem in (4) possesses a solution in the support of  $Q$ , the inclusion in (53) holds. ■

The distinction between coherent and consistent measures becomes more evident under certain conditions. Consider the case in which  $\mathcal{M}$  is finite. In this case, if the solution to the ERM problem in (4) is in the support of the  $\sigma$ -finite measure  $Q$ , then  $Q$  is both coherent and consistent. This is essentially because all measurable singletons (models) in  $\text{supp } Q$  exhibit positive measure with respect to  $Q$ . Alternatively, if the solution to the ERM problem in (4) is not in the support of  $Q$ , then  $Q$  is consistent but not coherent. Consider the case in which  $\mathcal{M}$  is the set  $\mathbb{R}^d$ ; the loss function  $\ell$  in (2) is continuous; and the ERM problem in (4) admits a unique solution. In this case, any probability measure  $Q$  absolutely continuous with respect to the Lebesgue measure is a coherent measure, but it is not a consistent measure. Alternatively, if the set of solutions to the ERM problem in (4) exhibits positive Lebesgue measure, then, the measure  $Q$  is both coherent and consistent.

## 4.2 Gibbs Reference Measures

In model selection, a natural idea is to proceed by successive approximations in the seek of lower computation complexity. From this perspective, one might wonder whether the solution to a current instance of an ERM-RER problem might serve as reference measure for the next instance. In this section, it is shown that this yields no benefit. Composing two successive ERM-REM problems boils down to a unique ERM-RER problem with the initial reference measure and a particular regularization factor. Under the assumption that  $\lambda \in \mathcal{K}_{Q,z}$ , with  $\mathcal{K}_{Q,z}$  in (23), the problem of interest is:

$$\min_{P \in \Delta_Q(\mathcal{M}, \mathcal{B}(\mathcal{M}))} R_z(P) + \alpha D(P \| P_{\Theta|Z=z}^{(Q,\lambda)}), \quad (60a)$$

$$\text{s. t.} \quad \int dP(\theta) = 1, \quad (60b)$$

where  $\alpha > 0$ ; the reference measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$ , which satisfies (25), is the solution of the ERM-RER problem in (19); and the functional  $R_z$  is defined in (18). From Theorem 3.1, the solution to the ERM-RER problem in (60), which is denoted by  $P_{\Theta|Z=z}^{(P_{\Theta|Z=z}^{(Q,\lambda)}, \alpha)}$ , satisfies for all  $\theta \in \text{supp } Q$  that

$$\frac{dP_{\Theta|Z=z}^{(P_{\Theta|Z=z}^{(Q,\lambda)}, \alpha)}}{dP_{\Theta|Z=z}^{(Q,\lambda)}}(\theta) = \exp\left(-K_{P_{\Theta|Z=z}^{(Q,\lambda)}, z}\left(-\frac{1}{\alpha}\right) - \frac{1}{\alpha} L_z(\theta)\right). \quad (61)$$

The log-partition functions  $K_{Q,z}$  in (22) and  $K_{P_{\Theta|Z=z}^{(Q,\lambda)}, z}$  in (61) are strongly related, as shown by the following lemma.

**Lemma 4.5.** *The functions  $K_{Q,z}$  in (22) and  $K_{P_{\Theta|Z=z}^{(Q,\lambda)}, z}$  in (61) satisfy for all  $t \in \mathbb{R}$ ,*

$$K_{P_{\Theta|Z=z}^{(Q,\lambda)}, z}(t) = K_{Q,z}\left(t - \frac{1}{\lambda}\right) - K_{Q,z}\left(-\frac{1}{\lambda}\right). \quad (62)$$



Moreover, for all  $t \leq 0$ ,

$$K_{P_{\Theta|Z=z}^{(Q,\lambda)}}(t) \leq 0. \quad (63)$$

*Proof:* The proof of (62) relies on the fact that for all

$$t \in \left\{ \nu \in \mathbb{R} : K_{P_{\Theta|Z=z}^{(Q,\lambda)}}(\nu) < \infty \right\},$$

the function  $K_{P_{\Theta|Z=z}^{(Q,\lambda)}}$  in (61) satisfies

$$K_{P_{\Theta|Z=z}^{(Q,\lambda)}}(t) = \log \left( \int \exp(t \mathbb{L}_z(\theta)) dP_{\Theta|Z=z}^{(Q,\lambda)}(\theta) \right) \quad (64)$$

$$= \log \left( \int \exp(t \mathbb{L}_z(\theta)) \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) dQ(\theta) \right) \quad (65)$$

$$= \log \left( \int \exp \left( \left( t - \frac{1}{\lambda} \right) \mathbb{L}_z(\theta) - K_{Q,z} \left( -\frac{1}{\lambda} \right) \right) dQ(\theta) \right) \quad (66)$$

$$= \log \left( \int \exp \left( \left( t - \frac{1}{\lambda} \right) \mathbb{L}_z(\theta) \right) dQ(\theta) \right) - K_{Q,z} \left( -\frac{1}{\lambda} \right) \quad (67)$$

$$= K_{Q,z} \left( t - \frac{1}{\lambda} \right) - K_{Q,z} \left( -\frac{1}{\lambda} \right), \quad (68)$$

where the equality in (66) follows from (25). Moreover, from Lemma 5.2, it follows that the function  $K_{P_{\Theta|Z=z}^{(Q,\lambda)}}$  is continuous and nondecreasing. Let  $s^* \in \mathbb{R} \cup \{+\infty\}$  be defined by

$$s^* \triangleq \sup \left\{ \nu \in \mathbb{R} : K_{P_{\Theta|Z=z}^{(Q,\lambda)}}(\nu) < \infty \right\}. \quad (69)$$

If  $s^* = +\infty$ , then for all  $t \in \mathbb{R}$ ,  $K_{P_{\Theta|Z=z}^{(Q,\lambda)}}(t) < +\infty$ , and the proof of (62) is completed.

Alternatively, if  $s^* < +\infty$ , it follows that for all  $t > s^*$ ,  $K_{P_{\Theta|Z=z}^{(Q,\lambda)}}(t) = +\infty$ , which implies that  $K_{Q,z} \left( t - \frac{1}{\lambda} \right) = +\infty$ , as the function  $K_{Q,z}$  is also continuous (Lemma 5.2) and  $K_{Q,z} \left( -\frac{1}{\lambda} \right) < \infty$  (due to the choice of  $\lambda$ ). Hence, in this case, the equality in (62) is of the form  $+\infty = +\infty$ . This completes the proof of (62).

The proof of (63) follows by noticing that for all  $t \leq 0$  and for all  $\theta \in \text{supp } Q$ , it holds that  $\exp(t \mathbb{L}_z(\theta)) \leq 1$ . Hence,

$$K_{P_{\Theta|Z=z}^{(Q,\lambda)}}(t) = \log \left( \int \exp(t \mathbb{L}_z(\theta)) dP_{\Theta|Z=z}^{(Q,\lambda)}(\theta) \right) \quad (70)$$

$$\leq \log \left( \int dP_{\Theta|Z=z}^{(Q,\lambda)}(\theta) \right) \quad (71)$$

$$= 0, \quad (72)$$

which completes the proof.  $\blacksquare$

The following lemma establishes that the solution to the ERM-RER problem in (60) is identical to the solution to another ERM-RER problem of the form

$$\min_{P \in \Delta_Q(\mathcal{M}, \mathcal{B}(\mathcal{M}))} R_z(P) + \left( \frac{1}{\frac{1}{\alpha} + \frac{1}{\lambda}} \right) D(P \| Q), \quad (73a)$$

$$\text{s. t.} \quad \int dP(\boldsymbol{\theta}) = 1, \quad (73b)$$

with  $\lambda \in \mathcal{K}_{Q,z}$ , with  $\mathcal{K}_{Q,z}$  in (23), and whose solution, denoted by  $P_{\boldsymbol{\Theta}|Z=z}^{(Q, \frac{1}{\frac{1}{\alpha} + \frac{1}{\lambda}})}$ , satisfies for all  $\boldsymbol{\theta} \in \text{supp } Q$ ,

$$\frac{dP_{\boldsymbol{\Theta}|Z=z}^{(Q, \frac{1}{\frac{1}{\alpha} + \frac{1}{\lambda}})}}{dQ}(\boldsymbol{\theta}) = \exp \left( -K_{Q,z} \left( -\frac{1}{\lambda} - \frac{1}{\alpha} \right) - \left( \frac{1}{\lambda} + \frac{1}{\alpha} \right) L_z(\boldsymbol{\theta}) \right). \quad (74)$$

The formal statement is as follows.

**Lemma 4.6.** *Let  $\alpha \in (0, +\infty)$  and  $\lambda \in \mathcal{K}_{Q,z}$ , with  $\mathcal{K}_{Q,z}$  in (23). Then, the probability measures  $P_{\boldsymbol{\Theta}|Z=z}^{(P_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}, \alpha)}$  in (61) and  $P_{\boldsymbol{\Theta}|Z=z}^{(Q, \frac{1}{\frac{1}{\alpha} + \frac{1}{\lambda}})}$  in (74) are identical.*

*Proof:* For all  $\boldsymbol{\theta} \in \text{supp } Q$ ,

$$\frac{dP_{\boldsymbol{\Theta}|Z=z}^{(P_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}, \alpha)}}{dQ}(\boldsymbol{\theta}) = \frac{dP_{\boldsymbol{\Theta}|Z=z}^{(P_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}, \alpha)}}{dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}}(\boldsymbol{\theta}) \frac{dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}}{dQ}(\boldsymbol{\theta}) \quad (75)$$

$$= \exp \left( -K_{P_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}, z} \left( -\frac{1}{\alpha} \right) - K_{Q,z} \left( -\frac{1}{\lambda} \right) - \left( \frac{1}{\alpha} + \frac{1}{\lambda} \right) L_z(\boldsymbol{\theta}) \right) \quad (76)$$

$$= \exp \left( -K_{Q,z} \left( -\frac{1}{\alpha} - \frac{1}{\lambda} \right) - \left( \frac{1}{\alpha} + \frac{1}{\lambda} \right) L_z(\boldsymbol{\theta}) \right) \quad (77)$$

$$= \frac{dP_{\boldsymbol{\Theta}|Z=z}^{(Q, \frac{1}{\frac{1}{\alpha} + \frac{1}{\lambda}})}}{dQ}(\boldsymbol{\theta}), \quad (78)$$

where the equality in (75) follows from the fact that the measure  $P_{\boldsymbol{\Theta}|Z=z}^{(P_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}, \alpha)}$  is absolutely continuous with respect to  $P_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}$  and  $P_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}$  is absolutely continuous with respect to the measure  $Q$ ; the equality in (76) follows from Lemma 4.5; and the equality in (78) follows from Theorem 3.1.

For all measurable subsets  $\mathcal{A}$  of  $\mathcal{M}$ , the following holds:

$$P_{\Theta|Z=z}^{(P^{(Q,\lambda)}, \alpha)}(\mathcal{A}) = \int_{\mathcal{A}} \frac{dP_{\Theta|Z=z}^{(P^{(Q,\lambda)}, \alpha)}}{dQ}(\theta) dQ(\theta) \quad (79)$$

$$= \int_{\mathcal{A}} \frac{dP_{\Theta|Z=z}^{(Q, \frac{1}{\lambda} + \frac{1}{\alpha})}}{dQ} dQ(\theta) \quad (80)$$

$$= \int_{\mathcal{A}} dP_{\Theta|Z=z}^{(Q, \frac{1}{\lambda} + \frac{1}{\alpha})}(\theta) \quad (81)$$

$$= P_{\Theta|Z=z}^{(Q, \frac{1}{\lambda} + \frac{1}{\alpha})}(\mathcal{A}), \quad (82)$$

where the equality in (80) follows from (78). This completes the proof. ■

The following theorem establishes a relation between the solutions to the following optimization problems

$$\min_{P \in \Delta_Q(\mathcal{M}, \mathcal{B}(\mathcal{M}))} R_z(P), \quad (83a)$$

$$\text{s. t.} \quad D(P \| P_{\Theta|Z=z}^{(Q,\lambda)}) \leq c, \text{ and} \quad (83b)$$

$$\int dP(\theta) = 1, \quad (83c)$$

and

$$\min_{P \in \Delta_Q(\mathcal{M}, \mathcal{B}(\mathcal{M}))} R_z(P) + \omega D(P \| Q), \quad (84a)$$

$$\text{s. t.} \quad \int dP(\theta) = 1, \quad (84b)$$

with  $c > 0$  and  $\omega \in \mathcal{K}_{Q,z}$ , with  $\mathcal{K}_{Q,z}$  in (23), two constants;  $P_{\Theta|Z=z}^{(Q,\lambda)}$  the probability measure in (25); and  $R_z$  the functional in (18).

From Theorem 3.1, the solution to the ERM-RER problem in (84), which is denoted by  $P_{\Theta|Z=z}^{(Q,\omega)}$ , satisfies for all  $\theta \in \text{supp } Q$  that

$$\frac{dP_{\Theta|Z=z}^{(Q,\omega)}}{dQ}(\theta) = \exp\left(-K_{Q,z}\left(-\frac{1}{\omega}\right) - \frac{1}{\omega} L_z(\theta)\right), \quad (85)$$

where the function  $K_{Q,z}$  is in (22).

The following theorem formalizes the relation between both optimization problems.

**Theorem 4.1.** *Assume that  $c$  and  $\omega$  in (83) and (84) satisfy*

$$D(P_{\Theta|Z=z}^{(Q,\omega)} \| P_{\Theta|Z=z}^{(Q,\lambda)}) = c, \quad (86)$$

with  $P_{\Theta|Z=z}^{(Q,\lambda)}$  and  $P_{\Theta|Z=z}^{(Q,\omega)}$  being the probability measures in (25) and (85), respectively. Then, the solution to the optimization problem in (83) is the probability measure  $P_{\Theta|Z=z}^{(Q,\omega)}$ .

*Proof:* The proof is presented in Appendix J. ■

## 5 The Log-Partition Function

This section introduces some properties of the log-partition function  $K_{Q,z}$  in (22) using the notion of *separable* empirical risk functions.

### 5.1 Separable Empirical Risk Functions

Separable empirical risk functions are defined with respect to a measure  $P \in \Delta(\mathcal{M})$ .

**Definition 5.1** (Separable Empirical Risk Function). *The empirical risk function  $L_z$  in (3) is said to be separable with respect to a  $\sigma$ -finite measure  $P \in \Delta(\mathcal{M})$ , if there exist a positive real  $c > 0$  and two subsets  $\mathcal{A}$  and  $\mathcal{B}$  of  $\mathcal{M}$  that are nonnegligible with respect to  $P$ , and for all  $(\theta_1, \theta_2) \in \mathcal{A} \times \mathcal{B}$ ,*

$$L_z(\theta_1) < c < L_z(\theta_2) < +\infty. \quad (87)$$

In a nutshell, a nonseparable empirical risk function with respect to the measure  $Q$  is a constant almost surely. More specifically, there exists a real  $a \geq 0$ , such that

$$Q(\{\theta \in \mathcal{M} : L_z(\theta) = a\}) = 1. \quad (88)$$

From this perspective, nonseparable empirical risk functions exhibit little practical interest for model selection.

The definition of separability in Definition 5.1 and Lemma 3.3 lead to the following lemma.

**Lemma 5.1.** *The empirical risk function  $L_z$  in (3) is separable with respect to the  $\sigma$ -finite measure  $Q$  in (19) if and only if it is separable with respect to the probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  in (25).*

*Proof:* Consider first that the function  $L_z$  is separable with respect to the  $\sigma$ -finite measure  $Q$ . Hence, there exist a positive real  $c > 0$  and two subsets  $\mathcal{A}$  and  $\mathcal{B}$  of  $\mathcal{M}$  that are nonnegligible with respect to  $Q$ , such that for all  $(\theta_1, \theta_2) \in \mathcal{A} \times \mathcal{B}$  the inequality in (87) holds. Hence, from (87) the following inequalities hold:

$$-\frac{1}{\lambda}L_z(\theta_1) > -\frac{c}{\lambda} > -\frac{1}{\lambda}L_z(\theta_2) > -\infty, \text{ and} \quad (89)$$

$$\exp\left(-\frac{1}{\lambda}L_z(\theta_1)\right) > \exp\left(-\frac{c}{\lambda}\right) > \exp\left(-\frac{1}{\lambda}L_z(\theta_2)\right) > 0. \quad (90)$$

This implies that

$$\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta_1) > \exp\left(-K_{Q,z}\left(-\frac{1}{\lambda}\right) - \frac{c}{\lambda}\right) \quad (91)$$

$$> \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta_2) \quad (92)$$

$$> 0. \quad (93)$$

Using the inequality in (91) and the facts that  $Q(\mathcal{A}) > 0$  and  $Q(\mathcal{B}) > 0$ , the following holds

$$P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{A}) = \int_{\mathcal{A}} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) dQ(\theta) > 0, \quad (94)$$

and

$$P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{B}) = \int_{\mathcal{B}} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) dQ(\theta) > 0. \quad (95)$$

which implies that the function  $L_z$  is separable with respect to the probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$ .

Consider now that the function  $L_z$  is separable with respect to the probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$ . Hence, there exist a positive real  $c > 0$  and two subsets  $\mathcal{A}$  and  $\mathcal{B}$  of  $\mathcal{M}$  that are nonnegligible with respect to  $P_{\Theta|Z=z}^{(Q,\lambda)}$ , such that for all  $(\theta_1, \theta_2) \in \mathcal{A} \times \mathcal{B}$  the inequality in (87) holds. More specifically,  $P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{A}) > 0$  and  $P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{B}) > 0$ . From Lemma 3.2 and the inequality in (87), it follows that for all pairs  $(\theta_1, \theta_2) \in \mathcal{A} \times \mathcal{B}$ ,  $\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta_1) > 0$  and  $\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta_2) > 0$ . Hence, from the fact that  $P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{A}) > 0$  and  $P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{B}) > 0$ , it follows that  $Q(\mathcal{A}) > 0$  and  $Q(\mathcal{B}) > 0$ , which implies that the function  $L_z$  is separable with respect to the  $\sigma$ -finite measure  $Q$ . This completes the proof. ■

Lemma 5.1 shows that separable empirical risk functions, and only these functions, lead to ERM-RER-optimal probability measures from which models are sampled with different probabilities. For the case of nonseparable empirical risk functions, all models are sampled from the ERM-RER-optimal probability measure with the same probability.

## 5.2 Properties of the Log-Partition Function

The log-partition function  $K_{Q,z}$  in (22) is a nondecreasing continuous convex function as shown by the following lemmas.

**Lemma 5.2.** *The function  $K_{Q,z}$  in (22) is nondecreasing and differentiable infinitely many times in the interior of  $\{t \in \mathbb{R} : K_{Q,z}(t) < +\infty\}$ .*

*Proof:* The proof is presented in Appendix K. ■

**Lemma 5.3.** *The function  $K_{Q,z}$  in (22) is convex in  $\{t \in \mathbb{R} : K_{Q,z}(t) < +\infty\}$ . Moreover, it is strictly convex if and only if the empirical risk function  $\mathbb{L}_z$  in (3) is separable with respect to the  $\sigma$ -finite measure  $Q$  in (19).*

*Proof:* The proof is presented in Appendix L. ■

In Lemma 5.2, it has been established that the log-partition function  $K_{Q,z}$  in (22) is differentiable infinitely many times in the interval

$$\{t \in \mathbb{R} : K_{Q,z}(t) < +\infty\}.$$

Let the  $m$ -th derivative of the function  $K_{Q,z}$  in (22) be denoted by  $K_{Q,z}^{(m)} : \mathbb{R} \rightarrow \mathbb{R}$ , with  $m \in \mathbb{N}$ . Hence, for all  $s \in \mathcal{K}_{Q,z}$ ,

$$K_{Q,z}^{(m)}\left(-\frac{1}{s}\right) \triangleq \frac{d^m}{dt^m} K_{Q,z}(t) \Big|_{t=-\frac{1}{s}}. \quad (96)$$

The following lemma provides explicit expressions for the first, second and third derivatives of the function  $K_{Q,z}$  in (22).

**Lemma 5.4.** *The first, second and third derivatives of the function  $K_{Q,z}$  in (22), denoted respectively by  $K_{Q,z}^{(1)}$ ,  $K_{Q,z}^{(2)}$ , and  $K_{Q,z}^{(3)}$ , satisfy for all  $\lambda \in \text{int}\mathcal{K}_{Q,z}$ , with  $\mathcal{K}_{Q,z}$  in (23),*

$$K_{Q,z}^{(1)}\left(-\frac{1}{\lambda}\right) = \int \mathbb{L}_z(\boldsymbol{\theta}) dP_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}(\boldsymbol{\theta}), \quad (97)$$

$$K_{Q,z}^{(2)}\left(-\frac{1}{\lambda}\right) = \int \left(\mathbb{L}_z(\boldsymbol{\theta}) - K_{Q,z}^{(1)}\left(-\frac{1}{\lambda}\right)\right)^2 dP_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}(\boldsymbol{\theta}), \quad (98)$$

$$K_{Q,z}^{(3)}\left(-\frac{1}{\lambda}\right) = \int \left(\mathbb{L}_z(\boldsymbol{\theta}) - K_{Q,z}^{(1)}\left(-\frac{1}{\lambda}\right)\right)^3 dP_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}(\boldsymbol{\theta}), \quad (99)$$

where the function  $\mathbb{L}_z$  is defined in (3) and the measure  $P_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}$  satisfies (25).

*Proof:* The proof is presented in Appendix M. ■

From Lemma 5.4, it follows that if  $\boldsymbol{\Theta} \sim P_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}$ , with  $P_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}$  in (25), the random variable

$$W \triangleq \mathbb{L}_z(\boldsymbol{\Theta}), \quad (100)$$

with the function  $\mathbb{L}_z$  in (3), possesses a mean, variance, and third cumulant that are equivalent to  $K_{Q,z}^{(1)}(-\frac{1}{\lambda})$  in (97),  $K_{Q,z}^{(2)}(-\frac{1}{\lambda})$  in (98), and  $K_{Q,z}^{(3)}(-\frac{1}{\lambda})$  in (99), respectively.

Note that if there exists a  $\delta > 0$  such that the log-partition function  $K_{Q,z}$  is differentiable within the open interval  $(-\delta, \delta)$  and  $Q$  in (19) is a probability measure, the function  $K_{Q,z}$  is the cumulant generating function of the random variable

$$V \triangleq \mathbb{L}_z(\boldsymbol{\Theta}), \text{ with } \boldsymbol{\Theta} \sim Q. \quad (101)$$

The following lemma leverages this observation.

**Lemma 5.5.** *Assume that  $Q$  in (19) is a probability measure and that there exists real  $\delta > 0$  such that the log-partition function  $K_{Q,z}$  in (22) is differentiable within  $(-\delta, \delta)$ . Then, the first, second and third derivatives of  $K_{Q,z}$ , denoted respectively by  $K_{Q,z}^{(1)}$ ,  $K_{Q,z}^{(2)}$ , and  $K_{Q,z}^{(3)}$ , satisfy*

$$K_{Q,z}^{(1)}(0) = \int \mathbb{L}_z(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}), \quad (102)$$

$$K_{Q,z}^{(2)}(0) = \int \left( \mathbb{L}_z(\boldsymbol{\theta}) - K_{Q,z}^{(1)} \left( -\frac{1}{\lambda} \right) \right)^2 dQ(\boldsymbol{\theta}), \quad (103)$$

$$K_{Q,z}^{(3)}(0) = \int \left( \mathbb{L}_z(\boldsymbol{\theta}) - K_{Q,z}^{(1)} \left( -\frac{1}{\lambda} \right) \right)^3 dQ(\boldsymbol{\theta}), \quad (104)$$

where the function  $\mathbb{L}_z$  is defined in (3).

*Proof:* The proof follows along the same arguments of the proof of Lemma 5.4. ■

The mean, variance, and third cumulant of the random variable  $V$  in (101) are  $K_{Q,z}^{(1)}(0)$  in (102),  $K_{Q,z}^{(2)}(0)$  in (103), and  $K_{Q,z}^{(3)}(0)$  in (104), respectively.

## 6 Expectation of the Empirical Risk

The mean of the random variable  $W$  in (100) is equivalent to the expectation of the empirical risk function  $\mathbb{L}_z$  with respect to the probability measure  $P_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}$  in (25), which is equal to  $R_z \left( P_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)} \right)$ , with the functional  $R_z$  in (18). Often,  $R_z \left( P_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)} \right)$  is referred to as the ERM-RER-optimal expected empirical risk to emphasize that this is the expected value of the empirical risk when models are sampled from the solution of the ERM-RER problem in (19). The following corollary of Lemma 5.4 formalizes this observation.

**Corollary 6.1.** *The probability measure  $P_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}$  in (25) verifies that*

$$R_z \left( P_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)} \right) = K_{Q,z}^{(1)} \left( -\frac{1}{\lambda} \right), \quad (105)$$

where the functional  $R_z$  and the function  $K_{Q,z}^{(1)}$  are defined in (18) and (97), respectively.

The expected empirical risk  $R_z \left( P_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)} \right)$  in (105) exhibits the following property.

**Theorem 6.1.** *The expected empirical risk  $R_z \left( P_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)} \right)$  in (105) is nondecreasing with  $\lambda \in K_{Q,z}$ , with  $K_{Q,z}$  in (23). Moreover,  $R_z \left( P_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)} \right)$  is strictly increasing with  $\lambda \in K_{Q,z}$  if and only if the function  $\mathbb{L}_z$  in (3) is separable with respect to the measure  $Q$ .*

*Proof:* The proof is presented in Appendix N. ■

The expected empirical risk  $R_z \left( P_{\Theta|Z=z}^{(Q,\lambda)} \right)$  in (105) has been shown to be non-decreasing with  $\lambda$  in [8, Appendix E.4] for the special case in which  $Q$  is a probability measure.

A question that arises from Theorem 6.1 is whether the value  $R_z \left( P_{\Theta|Z=z}^{(Q,\lambda)} \right)$  in (105) can be made arbitrarily close to  $\delta_{Q,z}^*$ , with  $\delta_{Q,z}^*$  in (38), by making  $\lambda$  arbitrarily small. The following lemma shows that the value  $R_z \left( P_{\Theta|Z=z}^{(Q,\lambda)} \right)$  is often bounded away from  $\delta_{Q,z}^*$ , even for arbitrarily small values of  $\lambda$ .

**Lemma 6.1.** *The expected empirical risk  $R_z \left( P_{\Theta|Z=z}^{(Q,\lambda)} \right)$  in (105) satisfies,*

$$R_z \left( P_{\Theta|Z=z}^{(Q,\lambda)} \right) \geq \delta_{Q,z}^*, \quad (106)$$

where  $\delta_{Q,z}^*$  is defined in (38). Moreover, the inequality in (106) is strict if and only if the function  $L_z$  in (3) is separable with respect to the measure  $Q$  in (19).

*Proof:* The proof is presented in Appendix O. ■

In the asymptotic regime when  $\lambda$  tends to zero, the expected empirical risk  $R_z \left( P_{\Theta|Z=z}^{(Q,\lambda)} \right)$  in (105) is equal to  $\delta_{Q,z}^*$ , as shown by the following lemma.

**Theorem 6.2.** *The expected empirical risk  $R_z \left( P_{\Theta|Z=z}^{(Q,\lambda)} \right)$  in (105) satisfies,*

$$\lim_{\lambda \rightarrow 0^+} R_z \left( P_{\Theta|Z=z}^{(Q,\lambda)} \right) = \delta_{Q,z}^*, \quad (107)$$

where  $\delta_{Q,z}^*$  is defined in (38).

*Proof:* The proof is presented in Appendix P. ■

The following lemma determines the value of the objective function of the ERM-RER problem in (19) when it is evaluated at its solution. This result appeared first in [10, Lemma 3].

**Lemma 6.2** (Lemma 3 in [10]). *The probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  in (25) and the  $\sigma$ -finite measure  $Q$  in (19) satisfy*

$$R_z \left( P_{\Theta|Z=z}^{(Q,\lambda)} \right) + \lambda D \left( P_{\Theta|Z=z}^{(Q,\lambda)} \| Q \right) = -\lambda K_{Q,z} \left( -\frac{1}{\lambda} \right). \quad (108)$$

Moreover, if the condition in (27) holds, then,

$$R_z(Q) - \lambda D \left( Q \| P_{\Theta|Z=z}^{(Q,\lambda)} \right) = -\lambda K_{Q,z} \left( -\frac{1}{\lambda} \right), \quad (109)$$

where the functional  $R_z$  is defined in (18); and the function  $K_{Q,z}$  is defined in (22).



*Proof:* From Theorem 3.1, it follows that for all  $\theta \in \text{supp } Q$ ,

$$\log \left( \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) \right) = -K_{Q,z} \left( -\frac{1}{\lambda} \right) - \frac{1}{\lambda} \mathbf{L}_z(\theta), \quad (110)$$

where the function  $\mathbf{L}_z$  is defined in (3). Thus,

$$D(P_{\Theta|Z=z}^{(Q,\lambda)} \| Q) = \int \log \left( \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) \right) dP_{\Theta|Z=z}^{(Q,\lambda)}(\theta) \quad (111)$$

$$= -K_{Q,z} \left( -\frac{1}{\lambda} \right) - \frac{1}{\lambda} \int \mathbf{L}_z(\theta) dP_{\Theta|Z=z}^{(Q,\lambda)}(\theta) \quad (112)$$

$$= -K_{Q,z} \left( -\frac{1}{\lambda} \right) - \frac{1}{\lambda} \mathbf{R}_z(P_{\Theta|Z=z}^{(Q,\lambda)}), \quad (113)$$

where the functional  $\mathbf{R}_z$  is defined in (18). This completes the proof of (108).

From Lemma 3.3 and (110), it follows that

$$D(Q \| P_{\Theta|Z=z}^{(Q,\lambda)}) = - \int \log \left( \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) \right) dQ(\theta) \quad (114)$$

$$= K_{Q,z} \left( -\frac{1}{\lambda} \right) + \frac{1}{\lambda} \int \mathbf{L}_z(\theta) dQ(\theta) \quad (115)$$

$$= K_{Q,z} \left( -\frac{1}{\lambda} \right) + \frac{1}{\lambda} \mathbf{R}_z(Q), \quad (116)$$

which completes the proof of (109).  $\blacksquare$

The following corollary of Lemma 6.2 characterizes the difference between the expected values of the random variables  $W$  and  $V$  in (100) and (101), respectively.

**Corollary 6.2.** *If measures  $Q$  and  $P_{\Theta|Z=z}^{(Q,\lambda)}$  in (25) are both probability measures, then,*

$$\mathbf{R}_z(Q) - \mathbf{R}_z(P_{\Theta|Z=z}^{(Q,\lambda)}) = \lambda \left( D(Q \| P_{\Theta|Z=z}^{(Q,\lambda)}) + D(P_{\Theta|Z=z}^{(Q,\lambda)} \| Q) \right). \quad (117)$$

The right-hand side of (117) is a symmetrized Kullback-Liebler divergence, also known as Jeffrey's divergence [65], between the measures  $Q$  and  $P_{\Theta|Z=z}^{(Q,\lambda)}$ . More importantly, when  $Q$  is a probability measure, it follows that  $D(P_{\Theta|Z=z}^{(Q,\lambda)} \| Q) \geq 0$  and  $D(Q \| P_{\Theta|Z=z}^{(Q,\lambda)}) \geq 0$ , which leads to the following corollary from Lemma 6.2.

**Corollary 6.3.** *If the  $\sigma$ -finite measure  $Q$  in (19) is a probability measure, then, the probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  in (25) satisfies*

$$\mathbf{R}_z(P_{\Theta|Z=z}^{(Q,\lambda)}) \leq \mathbf{R}_z(Q), \quad (118)$$

where, the functional  $\mathbf{R}_z$  is defined in (18).

## 7 Variance of the Empirical Risk

In Lemma 5.2, it has been established that if there exists a  $\delta > 0$  such that the log-partition function  $K_{Q,z}$  in (22) is finite within the open interval  $(-\delta, \delta)$  the log-partition function  $K_{Q,z}$  is differentiable infinitely many times within the interval  $(-\infty, \delta)$ . This together with the *mean value theorem* [66, Theorem 5.10] lead to the following characterization of the differences of the values  $K_{Q,z}^{(2)}(-\frac{1}{t})$  and  $K_{Q,z}^{(2)}(0)$ , with  $t > 0$ .

**Lemma 7.1.** *If the measure  $Q$  in (19) is a probability measure and there exists a  $\delta > 0$  such that the function  $K_{Q,z}$  in (22) is differentiable within the open interval  $(-\delta, \delta)$ , then for all  $t > 0$ ,*

$$K_{Q,z}^{(2)}\left(-\frac{1}{t}\right) - K_{Q,z}^{(2)}(0) = -\frac{1}{t} K_{Q,z}^{(3)}\left(-\frac{1}{\beta}\right) < +\infty, \quad (119)$$

for some  $\beta \in (t, +\infty)$ , where the functions  $K_{Q,z}^{(2)}$  and  $K_{Q,z}^{(3)}$  are defined in (96).

*Proof:* The proof is an immediate consequence of Lemma 5.2 and the mean value theorem [66, Theorem 5.10]. ■

The relevance of Lemma 7.1 lies on the fact that  $K_{Q,z}^{(2)}(-\frac{1}{\lambda})$  and  $K_{Q,z}^{(2)}(0)$  are the variances of the random variables  $W$  in (100) and  $V$  in (101). See Lemma 5.4 and Lemma 5.5. Under the assumptions of Lemma 7.1, it follows that the function  $K_{Q,z}^{(3)}$  is continuous in  $(-\infty, \delta)$ , where  $\delta > 0$ . Hence, for all  $t > 0$ , the function  $K_{Q,z}^{(3)}$  achieves a maximum and a minimum within the interval  $[-\frac{1}{t}, 0]$ . Such extrema allow providing lower and upper bounds on the variance  $K_{Q,z}^{(2)}(-\frac{1}{\lambda})$  of the random variable  $W$  in terms of the variance  $K_{Q,z}^{(2)}(0)$  of the random variable  $V$ , as shown hereunder.

**Corollary 7.1.** *If the measure  $Q$  in (19) is a probability measure and there exists a  $\delta > 0$  such that the function  $K_{Q,z}$  in (22) is differentiable within the open interval  $(-\delta, \delta)$ , then for all  $t > 0$ ,*

$$K_{Q,z}^{(2)}(0) - \frac{1}{t} c_2 \leq K_{Q,z}^{(2)}\left(-\frac{1}{t}\right) \leq K_{Q,z}^{(2)}(0) - \frac{1}{t} c_1, \quad (120)$$

where,

$$c_1 = \min_{s \in [-\frac{1}{t}, 0]} K_{Q,z}^{(3)}(s) \quad \text{and} \quad (121)$$

$$c_2 = \max_{s \in [-\frac{1}{t}, 0]} K_{Q,z}^{(3)}(s), \quad (122)$$

and the functions  $K_{Q,z}^{(2)}$  and  $K_{Q,z}^{(3)}$  are defined in (96).

The inequality in (120) reveals that under the assumptions of Corollary 7.1, in the asymptotic regime when  $t \rightarrow +\infty$ , the variances of the random variables  $W$

in (100) and  $V$  in (101) are identical. Additionally, unlike the means  $K_{Q,z}^{(1)}(-\frac{1}{\lambda})$  and  $K_{Q,z}^{(1)}(0)$  of the random variables  $W$  and  $V$ , which satisfy  $K_{Q,z}^{(1)}(-\frac{1}{\lambda}) \leq K_{Q,z}^{(1)}(0)$  (Corollary 6.3), their variances  $K_{Q,z}^{(2)}(-\frac{1}{\lambda})$  and  $K_{Q,z}^{(2)}(0)$  might satisfy  $K_{Q,z}^{(2)}(-\frac{1}{\lambda}) < K_{Q,z}^{(2)}(0)$  or  $K_{Q,z}^{(2)}(-\frac{1}{\lambda}) \geq K_{Q,z}^{(2)}(0)$  depending on whether the function  $K_{Q,z}^{(3)}$  is positive or negative within the interval  $[-\frac{1}{\lambda}, 0]$ . Using this observation the values  $K_{Q,z}^{(2)}(-\frac{1}{t})$ , with  $t > 0$ , and  $K_{Q,z}^{(2)}(0)$  can be compared as follows.

**Lemma 7.2.** *Assume that the measure  $Q$  in (19) is a probability measure and there exists a  $\delta > 0$  such that the function  $K_{Q,z}$  in (22) is differentiable within the open interval  $(-\delta, \delta)$ . Hence, the following holds for all  $t > 0$ :*

- If for all  $s > t$ ,  $K_{Q,z}^{(3)}(-\frac{1}{s}) < 0$ , then

$$K_{Q,z}^{(2)}(0) < K_{Q,z}^{(2)}\left(-\frac{1}{t}\right) < +\infty; \quad (123)$$

- If for all  $s > t$ ,  $K_{Q,z}^{(3)}(-\frac{1}{s}) > 0$ , then

$$K_{Q,z}^{(2)}\left(-\frac{1}{t}\right) < K_{Q,z}^{(2)}(0) < +\infty; \quad (124)$$

- If for some  $s > t$ ,  $K_{Q,z}^{(3)}(-\frac{1}{s}) = 0$ , then there exists two positive reals  $c_1$  and  $c_2$  such that

$$c_1 \leq \min\{K_{Q,z}^{(2)}\left(-\frac{1}{t}\right), K_{Q,z}^{(2)}(0)\} \quad (125)$$

$$\leq \max\{K_{Q,z}^{(2)}\left(-\frac{1}{t}\right), K_{Q,z}^{(2)}(0)\} \quad (126)$$

$$\leq c_2. \quad (127)$$

*Proof:* The proofs of the inequalities in (123) and (124) are immediate consequences of Lemma 7.1. The inequalities in (125) and (127) follow from the fact that the function  $K_{Q,z}^{(2)}$ , which is continuous, exhibits critical points at  $-\frac{1}{s}$ , with  $s$  satisfying  $K_{Q,z}^{(3)}(-\frac{1}{s}) = 0$ . Some of such critical points might be local extrema of the function  $K_{Q,z}^{(2)}$ , either local minima or local maxima. Hence, the inequalities (125) and (127) follow by choosing  $c_1$  as the smallest minimum of the function  $K_{Q,z}^{(2)}$  within the interval  $[-\frac{1}{t}, 0]$ ; and  $c_2$  as the biggest maximum of the function  $K_{Q,z}^{(2)}$  within the interval  $[-\frac{1}{t}, 0]$ . If none of such critical points is a local extremum, then, (125) and (127) hold with equality. ■

Lemma 7.1 and Lemma 7.2 show that the monotonicity of the expectation of the random variable  $W$  in (100), stated by Theorem 6.1, is not a property exhibited by the variance nor the third cumulant. The following example highlights this observation.

**Example 7.1.** Consider the ERM-RER problem in (19), under the assumption that  $Q$  is a probability measure and the empirical risk function  $L_z$  in (3) is such that for all  $\theta \in \mathcal{M}$ ,

$$L_z(\theta) = \begin{cases} 0 & \text{if } \theta \in \mathcal{A} \\ 1 & \text{if } \theta \in \mathcal{M} \setminus \mathcal{A}, \end{cases} \quad (128)$$

where the sets  $\mathcal{A} \subset \mathcal{M}$  and  $\mathcal{M} \setminus \mathcal{A}$  are nonnegligible with respect to the reference probability measure  $Q$ . In this case, the function  $K_{Q,z}$  in (22) satisfies for all  $\lambda > 0$ ,

$$K_{Q,z}\left(-\frac{1}{\lambda}\right) = \log\left(Q(\mathcal{A}) + \exp\left(-\frac{1}{\lambda}\right)(1 - Q(\mathcal{A}))\right). \quad (129)$$

The derivatives  $K_{Q,z}^{(1)}$ ,  $K_{Q,z}^{(2)}$ , and  $K_{Q,z}^{(3)}$  in (96) of the function  $K_{Q,z}$  in (129) satisfy for all  $\lambda > 0$ ,

$$K_{Q,z}^{(1)}\left(-\frac{1}{\lambda}\right) = \frac{\exp\left(-\frac{1}{\lambda}\right)(1 - Q(\mathcal{A}))}{Q(\mathcal{A}) + \exp\left(-\frac{1}{\lambda}\right)(1 - Q(\mathcal{A}))}; \quad (130)$$

$$K_{Q,z}^{(2)}\left(-\frac{1}{\lambda}\right) = \frac{Q(\mathcal{A})(1 - Q(\mathcal{A}))\exp\left(-\frac{1}{\lambda}\right)}{(Q(\mathcal{A}) + \exp\left(-\frac{1}{\lambda}\right)(1 - Q(\mathcal{A})))^2}; \text{ and} \quad (131)$$

$$K_{Q,z}^{(3)}\left(-\frac{1}{\lambda}\right) = K_{Q,z}^{(2)}\left(-\frac{1}{\lambda}\right) \left( \frac{Q(\mathcal{A}) - (1 - Q(\mathcal{A}))\exp\left(-\frac{1}{\lambda}\right)}{Q(\mathcal{A}) + \exp\left(-\frac{1}{\lambda}\right)(1 - Q(\mathcal{A}))} \right). \quad (132)$$

Note that  $K_{Q,z}^{(3)}\left(-\frac{1}{\lambda}\right) > 0$  if and only if

$$Q(\mathcal{A}) - (1 - Q(\mathcal{A}))\exp\left(-\frac{1}{\lambda}\right) > 0. \quad (133)$$

Assume that  $Q(\mathcal{A}) \geq \frac{1}{2}$ . Thus, it holds that for all  $\lambda > 0$ , the inequality in (133) is always satisfied. This follows from observing that for all  $\lambda > 0$ ,

$$\exp\left(-\frac{1}{\lambda}\right) < 1 \leq \frac{Q(\mathcal{A})}{1 - Q(\mathcal{A})}. \quad (134)$$

Hence, if  $Q(\mathcal{A}) \geq \frac{1}{2}$ , for all decreasing sequences of positive reals  $\lambda_1 > \lambda_2 > \dots > 0$ , it holds that

$$\frac{1}{4} \geq K_{Q,z}^{(2)}\left(-\frac{1}{\lambda_1}\right) > K_{Q,z}^{(2)}\left(-\frac{1}{\lambda_2}\right) > \dots > 0. \quad (135)$$

Alternatively, assume that  $Q(\mathcal{A}) < \frac{1}{2}$ . In this case, the inequality in (133) is satisfied if and only if

$$\lambda < \left( \log\left(\frac{1 - Q(\mathcal{A})}{Q(\mathcal{A})}\right) \right)^{-1}. \quad (136)$$

Hence, if  $Q(\mathcal{A}) < \frac{1}{2}$ , then for all decreasing sequences of positive reals

$$\left( \log \left( \frac{1 - Q(\mathcal{A})}{Q(\mathcal{A})} \right) \right)^{-1} > \lambda_1 > \lambda_2 > \dots > 0,$$

it holds that

$$\frac{1}{4} > K_{Q,z}^{(2)} \left( -\frac{1}{\lambda_1} \right) > K_{Q,z}^{(2)} \left( -\frac{1}{\lambda_2} \right) > \dots > 0. \quad (137)$$

Moreover, for all decreasing sequences of positive reals

$$\lambda_1 > \lambda_2 > \dots > \left( \log \left( \frac{1 - Q(\mathcal{A})}{Q(\mathcal{A})} \right) \right)^{-1},$$

it holds that

$$K_{Q,z}^{(2)} \left( -\frac{1}{\lambda_1} \right) < K_{Q,z}^{(2)} \left( -\frac{1}{\lambda_2} \right) < \dots < \frac{1}{4}. \quad (138)$$

The upperbound by  $\frac{1}{4}$  in (135), (137) and (138) follows by noticing that the value  $K_{Q,z}^{(2)} \left( -\frac{1}{\lambda} \right)$  is maximized when  $\lambda = \left( \log \left( \frac{1 - Q(\mathcal{A})}{Q(\mathcal{A})} \right) \right)^{-1}$  and  $K_{Q,z}^{(2)} \left( -\frac{1}{\lambda} \right) = \frac{1}{4}$ .

Example 7.1 provides important insights on the choice of the reference measure  $Q$ . Note for instance that when the reference measure assigns a probability to the set of models  $\mathcal{T}(z)$  in (5) that is greater than or equal to the probability of suboptimal models  $\mathcal{M} \setminus \mathcal{T}(z)$ , i.e.,  $Q(\mathcal{T}(z)) \geq \frac{1}{2}$ , the variance is strictly decreasing to zero when  $\lambda$  decreases. See for instance, Figure 1 and Figure 2. That is, when the reference measure assigns higher probability to the set of solutions to the ERM problem in (4), the variance is monotone with respect to the parameter  $\lambda$ .

Alternatively, when the reference measure assigns a probability to the set  $\mathcal{T}(z)$  that is smaller than the probability of the set  $\mathcal{M} \setminus \mathcal{T}(z)$ , i.e.,  $Q(\mathcal{T}(z)) < \frac{1}{2}$ , there exists a critical point for  $\lambda$  at  $\left( \log \left( \frac{1 - Q(\mathcal{A})}{Q(\mathcal{A})} \right) \right)^{-1}$ . See for instance, Figure 3. More importantly, such a critical point can be arbitrarily close to zero depending on the value  $Q(\mathcal{A})$ . The variance strictly decreases when  $\lambda$  decreases beyond the value  $\left( \log \left( \frac{1 - Q(\mathcal{A})}{Q(\mathcal{A})} \right) \right)^{-1}$ . Otherwise, reducing  $\lambda$  above the value  $\left( \log \left( \frac{1 - Q(\mathcal{A})}{Q(\mathcal{A})} \right) \right)^{-1}$  increases the variance.

In general, these observations suggest that reference measures  $Q$  that allocate small measures to the sets containing the set  $\mathcal{T}(z)$  might require reducing the value  $\lambda$  beyond a small threshold in order to observe small values of  $K_{Q,z}^{(2)} \left( -\frac{1}{\lambda} \right)$ , which is the variance of the random variable  $W$ , in (100). These observations are central to understanding the concentration of probability that occurs when  $\lambda$  decreases to zero, as discussed in Section 9.

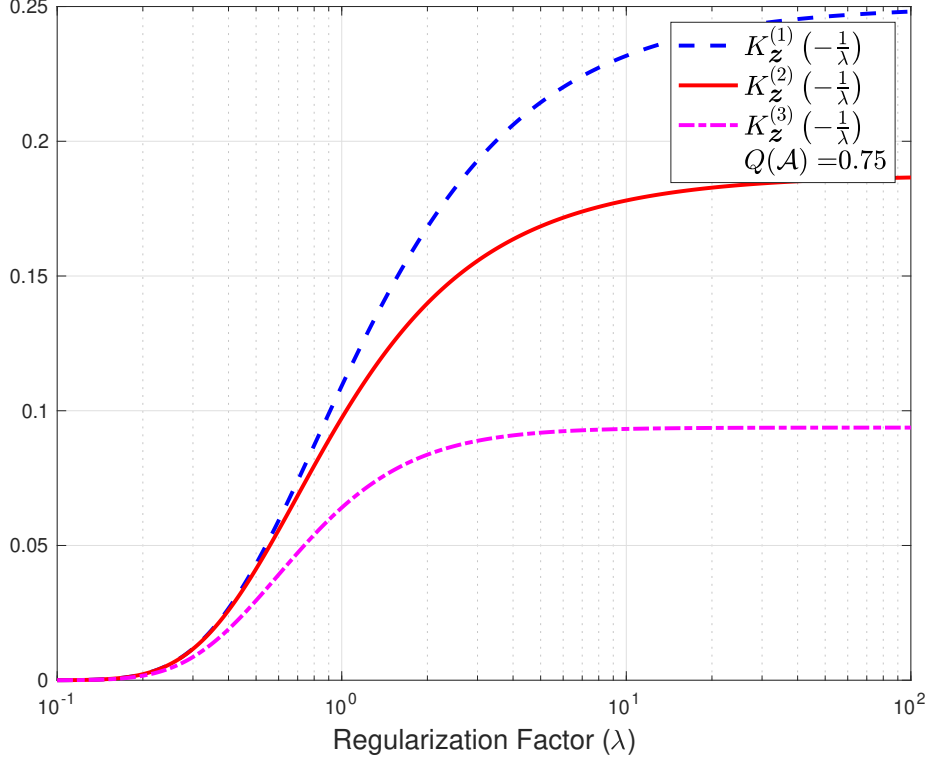


Figure 1: Mean  $K_{Q,z}^{(1)}(-\frac{1}{\lambda})$ , variance  $K_{Q,z}^{(2)}(-\frac{1}{\lambda})$ , and third central moment  $K_{Q,z}^{(3)}(-\frac{1}{\lambda})$  of the empirical risk in Example 7.1, with  $Q(\mathcal{A}) = \frac{3}{4}$

## 8 Cumulant Generating Function of the Empirical Risk

Consider the transport of the measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  in (25) from  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$  to  $([0, +\infty], \mathcal{B}([0, +\infty]))$  through the function  $L_z$  in (3). Denote the resulting probability measure in  $([0, +\infty], \mathcal{B}([0, +\infty]))$  by  $P_{W|Z=z}^{(Q,\lambda)}$ . That is, for all  $\mathcal{A} \in \mathcal{B}([0, +\infty])$ ,

$$P_{W|Z=z}^{(Q,\lambda)}(\mathcal{A}) = P_{\Theta|Z=z}^{(Q,\lambda)}(L_z^{-1}(\mathcal{A})), \quad (139)$$

where the term  $L_z^{-1}(\mathcal{A})$  represents the set

$$L_z^{-1}(\mathcal{A}) \triangleq \{\nu \in \mathcal{M} : L_z(\nu) \in \mathcal{A}\}. \quad (140)$$

Note that the random variable  $W$  in (100) induces the probability measure  $P_{W|Z=z}^{(Q,\lambda)}$  in  $([0, +\infty], \mathcal{B}([0, +\infty]))$ . The objective of this section is to study the properties of the cumulant generating function of the probability measure  $P_{W|Z=z}^{(Q,\lambda)}$ , denoted by  $J_{z,Q,\lambda} : \mathbb{R} \rightarrow \mathbb{R} \cup \{+\infty\}$ , which satisfies for all  $t \in$

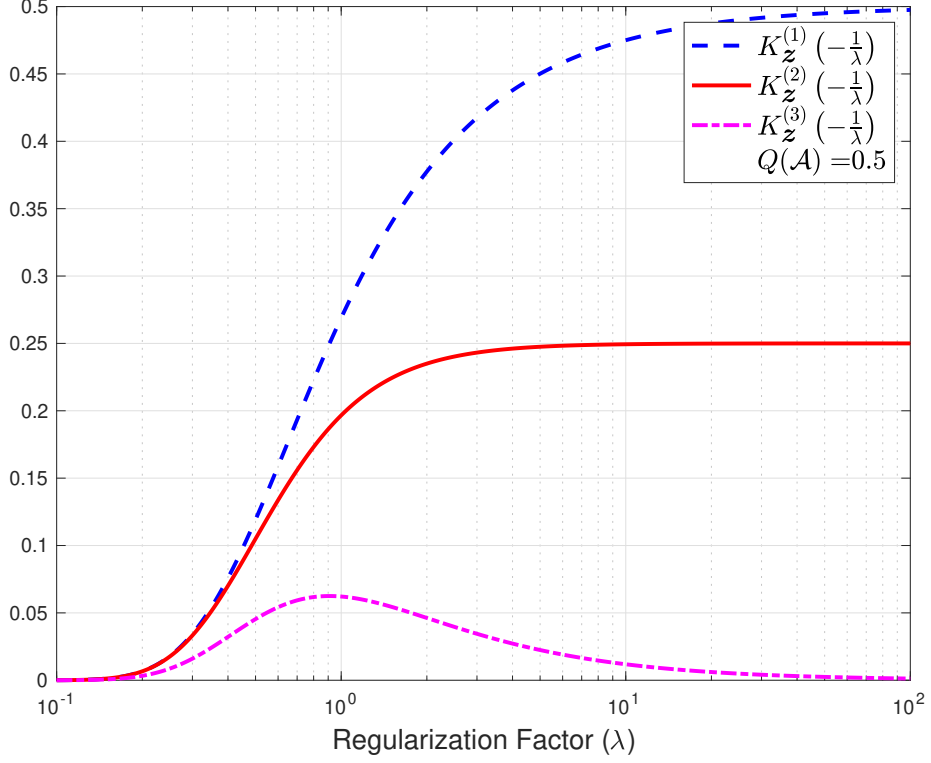


Figure 2: Mean  $K_{Q,z}^{(1)}(-\frac{1}{\lambda})$ , variance  $K_{Q,z}^{(2)}(-\frac{1}{\lambda})$ , and third central moment  $K_{Q,z}^{(3)}(-\frac{1}{\lambda})$  of the empirical risk in Example 7.1, with  $Q(\mathcal{A}) = \frac{1}{2}$

$\mathbb{R}$ ,

$$J_{z,Q,\lambda}(t) = \log \left( \int \exp(tw) dP_{W|Z=z}^{(Q,\lambda)}(w) \right) \quad (141)$$

$$= \log \left( \int \exp(tL_z(\theta)) dP_{\Theta|Z=z}^{(Q,\lambda)}(\theta) \right), \quad (142)$$

where the equality in (142) follows from [55, Theorem 1.6.12].

The following lemma provides an expression for  $J_{z,Q,\lambda}$  in terms of the log-partition function  $K_{Q,z}$  in (22).

**Lemma 8.1.** *If  $\lambda \in \mathcal{K}_{Q,z}$ , with  $\mathcal{K}_{Q,z}$  in (23), then, the function  $J_{z,Q,\lambda}$  in (141),*

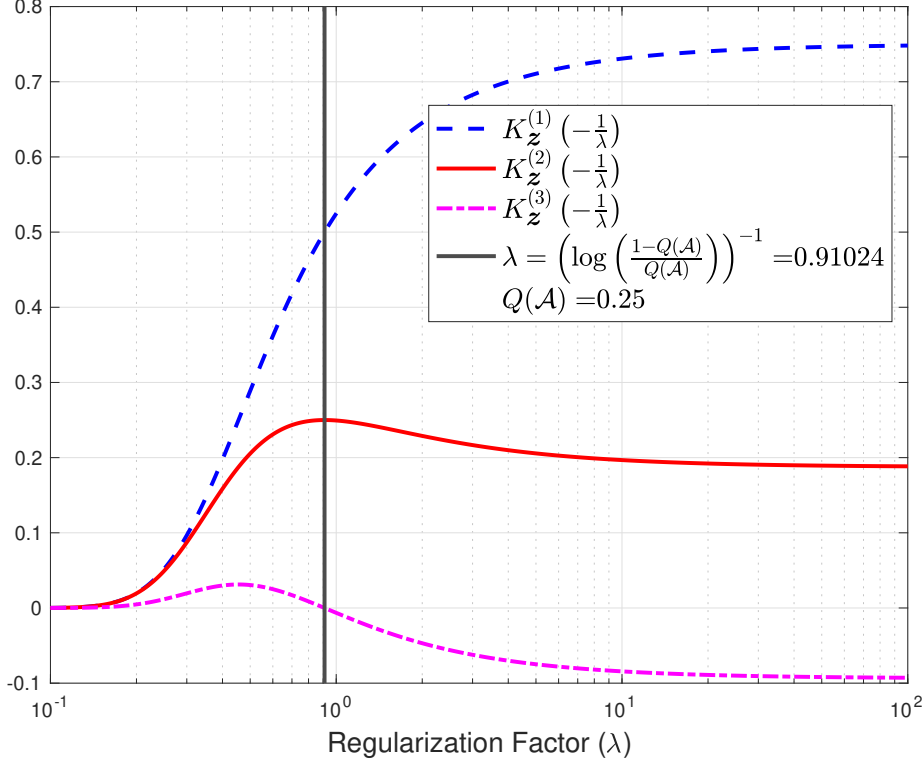


Figure 3: Mean  $K_{Q,z}^{(1)}(-\frac{1}{\lambda})$ , variance  $K_{Q,z}^{(2)}(-\frac{1}{\lambda})$ , and third central moment  $K_{Q,z}^{(3)}(-\frac{1}{\lambda})$  of the empirical risk in Example 7.1, with  $Q(\mathcal{A}) = \frac{1}{4}$

verifies for all  $t \in \mathbb{R}$ ,

$$J_{z,Q,\lambda}(t) = K_{P_{\Theta|Z=z,z}^{(Q,\lambda)}}(t) \quad (143)$$

$$= K_{Q,z}\left(t - \frac{1}{\lambda}\right) - K_{Q,z}\left(-\frac{1}{\lambda}\right) \quad (144)$$

$$= \sum_{m=1}^{+\infty} \frac{t^m}{m!} K_{Q,z}^{(m)}\left(-\frac{1}{\lambda}\right), \quad (145)$$

with the function  $K_{Q,z}$  in (22) and the function  $K_{Q,z}^{(m)}$  in (96).

*Proof:* The proof of (143) follows immediately from (22) and (142). The proof of (144) follows from Lemma 4.5. Finally, the proof of (145) follows by observing that a Taylor expansion of the function  $K_{Q,z}$  in (22) at the point  $-\frac{1}{\lambda}$ , yields for all  $t \in \{\nu \in \mathbb{R} : K_{Q,z}(\nu) < +\infty\}$ ,

$$K_{Q,z}(t) = K_{Q,z}\left(-\frac{1}{\lambda}\right) + \sum_{s=1}^{+\infty} \frac{K_{Q,z}^{(s)}\left(-\frac{1}{\lambda}\right)}{s!} \left(t + \frac{1}{\lambda}\right)^s. \quad (146)$$



Choosing  $\alpha \in \{\nu \in \mathbb{R} : K_{Q,z}(\nu - \frac{1}{\lambda}) < +\infty\}$  such that  $t = \alpha - \frac{1}{\lambda}$  in (146) yields

$$K_{Q,z}\left(\alpha - \frac{1}{\lambda}\right) = K_{Q,z}\left(-\frac{1}{\lambda}\right) + \sum_{s=1}^{+\infty} \frac{\alpha^s}{s!} K_{Q,z}^{(s)}\left(-\frac{1}{\lambda}\right), \quad (147)$$

which implies that for all  $t \in \{\nu \in \mathbb{R} : K_{Q,z}(\nu - \frac{1}{\lambda}) < +\infty\}$ ,

$$K_{Q,z}\left(t - \frac{1}{\lambda}\right) - K_{Q,z}\left(-\frac{1}{\lambda}\right) = \sum_{s=1}^{+\infty} \frac{t^s}{s!} K_{Q,z}^{(s)}\left(-\frac{1}{\lambda}\right). \quad (148)$$

Let  $s^* \in \mathbb{R} \cup \{+\infty\}$  be defined by

$$s^* \triangleq \sup \left\{ \nu \in \mathbb{R} : K_{Q,z}\left(\nu - \frac{1}{\lambda}\right) < \infty \right\}. \quad (149)$$

If  $s^* = +\infty$ , then for all  $t \in \mathbb{R}$ ,  $K_{Q,z}\left(t - \frac{1}{\lambda}\right) - K_{Q,z}\left(-\frac{1}{\lambda}\right) < +\infty$ , and thus,

$$+\infty > J_{z,Q,\lambda}(t) = K_{P_{\Theta|Z=z}^{(Q,\lambda)},z}(t) \quad (150)$$

$$= K_{Q,z}\left(t - \frac{1}{\lambda}\right) - K_{Q,z}\left(-\frac{1}{\lambda}\right) \quad (151)$$

$$= \sum_{m=1}^{+\infty} \frac{t^m}{m!} K_{Q,z}^{(m)}\left(-\frac{1}{\lambda}\right). \quad (152)$$

Alternatively, if  $s^* < +\infty$ , it follows that for all  $t > s^*$ ,  $K_{Q,z}\left(t - \frac{1}{\lambda}\right) = +\infty$ . From the fact that the function  $K_{Q,z}$  is continuous (Lemma 5.2) and  $K_{Q,z}\left(-\frac{1}{\lambda}\right) < \infty$  (due to the fact that  $\lambda \in \mathcal{K}_{Q,z}$  in (23)), it follows that

$$+\infty = J_{z,Q,\lambda}(t) = K_{Q,z}\left(t - \frac{1}{\lambda}\right) - K_{Q,z}\left(-\frac{1}{\lambda}\right) \quad (153)$$

$$= \sum_{m=1}^{+\infty} \frac{t^m}{m!} K_{Q,z}^{(m)}\left(-\frac{1}{\lambda}\right), \quad (154)$$

which implies that  $\sum_{m=1}^{+\infty} \frac{t^m}{m!} K_{Q,z}^{(m)}\left(-\frac{1}{\lambda}\right) = +\infty$ . Hence, in this case, the equality in (145) is of the form  $+\infty = +\infty$ . This completes the proof. ■

Alternative expressions for  $J_{z,Q,\lambda}$  in (141) are provided hereunder.

**Lemma 8.2.** *If  $\lambda \in \mathcal{K}_{Q,z}$ , with  $\mathcal{K}_{Q,z}$  in (23), then, the function  $J_{z,Q,\lambda}$  in (141), verifies for all  $t \in (0, +\infty)$ ,*

$$J_{z,Q,\lambda}\left(-\frac{1}{t}\right) = -\frac{1}{t} R_z \left( P_{\Theta|Z=z}^{(Q, \frac{1}{\lambda+t})} \right) - D \left( P_{\Theta|Z=z}^{(Q, \frac{1}{\lambda+t})} \| P_{\Theta|Z=z}^{(Q,\lambda)} \right) \quad (155)$$

$$= -\frac{1}{t} R_z \left( P_{\Theta|Z=z}^{(Q,\lambda)} \right) + D \left( P_{\Theta|Z=z}^{(Q,\lambda)} \| P_{\Theta|Z=z}^{(Q, \frac{1}{\lambda+t})} \right) \quad (156)$$

$$\leq 0, \quad (157)$$

where the functional  $R_z$  is in (18); the function  $K_{P_{\Theta|Z=z}, z}^{(Q, \lambda)}$  is in (61); and the

probability measures  $P_{\Theta|Z=z}^{(Q, \lambda)}$  and  $P_{\Theta|Z=z}^{(Q, \frac{1}{\lambda} + \frac{1}{t})}$  are respectively in (25) and (74).

*Proof:* The proof of (155) follows from (108) in Lemma 6.2 by observing that for all  $t \in (0, +\infty)$ ,

$$-tK_{P_{\Theta|Z=z}, z}^{(Q, \lambda)}\left(-\frac{1}{t}\right) = R_z\left(P_{\Theta|Z=z}^{(P_{\Theta|Z=z}^{(Q, \lambda)}, t)}\right) + tD\left(P_{\Theta|Z=z}^{(P_{\Theta|Z=z}^{(Q, \lambda)}, t)} \| P_{\Theta|Z=z}^{(Q, \lambda)}\right) \quad (158)$$

$$= R_z\left(P_{\Theta|Z=z}^{(Q, \frac{1}{\lambda} + \frac{1}{t})}\right) + tD\left(P_{\Theta|Z=z}^{(Q, \frac{1}{\lambda} + \frac{1}{t})} \| P_{\Theta|Z=z}^{(Q, \lambda)}\right), \quad (159)$$

where the equality in (159) follows from Lemma 4.6. The proof of (156) follows from (109) in Lemma 6.2 by observing that for all  $t \in (0, +\infty)$ ,

$$-tK_{P_{\Theta|Z=z}, z}^{(Q, \lambda)}\left(-\frac{1}{t}\right) = R_z\left(P_{\Theta|Z=z}^{(Q, \lambda)}\right) - tD\left(P_{\Theta|Z=z}^{(Q, \lambda)} \| P_{\Theta|Z=z}^{(P_{\Theta|Z=z}^{(Q, \lambda)}, t)}\right) \quad (160)$$

$$= R_z\left(P_{\Theta|Z=z}^{(Q, \lambda)}\right) - tD\left(P_{\Theta|Z=z}^{(Q, \lambda)} \| P_{\Theta|Z=z}^{(Q, \frac{1}{\lambda} + \frac{1}{t})}\right), \quad (161)$$

where the equality in (161) follows from Lemma 4.6, which completes the proof.  $\blacksquare$

From Lemma 5.2 and Lemma 8.1, it follows that the function  $J_{z, Q, \lambda}$  in (141) is increasing and differentiable infinitely many times in the interior of

$$\left\{t \in \mathbb{R} : K_{Q, z}\left(t - \frac{1}{\lambda}\right) < +\infty\right\}.$$

Moreover, note that  $(-\infty, \frac{1}{\lambda}] \subset \{t \in \mathbb{R} : K_{Q, z}\left(t - \frac{1}{\lambda}\right) < +\infty\}$ . Denote by  $J_{z, Q, \lambda}^{(m)} : \mathbb{R} \rightarrow \mathbb{R} \cup \{+\infty\}$ , with  $m \in \mathbb{N}$ , the  $m$ -th derivative of the function  $J_{z, Q, \lambda}$  in (141). That is, for all  $s \in \mathbb{R}$ ,

$$J_{z, Q, \lambda}^{(m)}(s) = \frac{d^m}{dt^m} J_{z, Q, \lambda}(t) \Big|_{t=s}. \quad (162)$$

From Lemma 8.1, it follows that for all  $m \in \mathbb{N}$ , and for all  $\alpha \in \mathbb{R}$ , the following holds,

$$J_{z, Q, \lambda}^{(m)}(\alpha) = K_{Q, z}^{(m)}\left(\alpha - \frac{1}{\lambda}\right), \quad (163)$$

where the function  $K_{Q, z}^{(m)}$  denotes the  $m$ -th derivative of the function  $K_{Q, z}$  in (22). See for instance, Lemma 5.4. The equality in (163) establishes a relation between the cumulant generating function  $J_{z, Q, \lambda}$  and the function  $K_{Q, z}$ . This observation becomes an alternative proof to Lemma 5.4.

The following theorem presents the relation between the cumulant generating function  $J_{z, Q, \lambda}$  and the functions  $K_{Q, z}^{(1)}$  and  $K_{Q, z}^{(2)}$  in (97) and (98).

**Theorem 8.1.** *For all  $\alpha \in \mathbb{R}$ , the function  $J_{\mathbf{z},Q,\lambda}$  in (141) verifies the following equality*

$$J_{\mathbf{z},Q,\lambda}(\alpha) = \alpha K_{Q,\mathbf{z}}^{(1)}\left(-\frac{1}{\lambda}\right) + \frac{1}{2}\alpha^2 K_{Q,\mathbf{z}}^{(2)}(\xi) \quad (164)$$

with

$$\xi \in \left( \min \left\{ -\frac{1}{\lambda}, \alpha - \frac{1}{\lambda} \right\}, \max \left\{ -\frac{1}{\lambda}, \alpha - \frac{1}{\lambda} \right\} \right), \quad (165)$$

where the functions  $K_{Q,\mathbf{z}}^{(1)}$  and  $K_{Q,\mathbf{z}}^{(2)}$  are defined in (97) and (98), respectively.

*Proof:* From Lemma 5.2, it follows that the function  $K_{Q,\mathbf{z}}$  is differentiable infinitely many times in the interior of  $\{t \in \mathbb{R} : K_{Q,\mathbf{z}}(t) < +\infty\}$ . Then, a Taylor expansion of the function  $K_{Q,\mathbf{z}}$  in (22) at the point  $-\frac{1}{\lambda}$  yields for all  $t \in \{\nu \in \mathbb{R} : K_{Q,\mathbf{z}}(\nu) < +\infty\}$ ,

$$K_{Q,\mathbf{z}}(t) = K_{Q,\mathbf{z}}\left(-\frac{1}{\lambda}\right) + \sum_{s=1}^{+\infty} \frac{1}{s!} \left(t + \frac{1}{\lambda}\right)^s K_{Q,\mathbf{z}}^{(s)}\left(-\frac{1}{\lambda}\right). \quad (166)$$

Choosing  $t = \alpha - \frac{1}{\lambda}$ , with  $\alpha \in \{\nu \in \mathbb{R} : K_{Q,\mathbf{z}}(\nu - \frac{1}{\lambda}) < +\infty\}$  in (166), it holds from the Taylor-Lagrange theorem [67, Theorem 2.5.4] that

$$K_{Q,\mathbf{z}}\left(\alpha - \frac{1}{\lambda}\right) = K_{Q,\mathbf{z}}\left(-\frac{1}{\lambda}\right) + \alpha K_{Q,\mathbf{z}}^{(1)}\left(-\frac{1}{\lambda}\right) + \frac{1}{2}\alpha^2 K_{Q,\mathbf{z}}^{(2)}(\xi), \quad (167)$$

where  $\xi \in (\min\{-\frac{1}{\lambda}, \alpha - \frac{1}{\lambda}\}, \max\{-\frac{1}{\lambda}, \alpha - \frac{1}{\lambda}\})$ .

Let  $s^* \in \mathbb{R} \cup \{+\infty\}$  be defined by

$$s^* \triangleq \sup \left\{ \nu \in \mathbb{R} : K_{Q,\mathbf{z}}\left(\nu - \frac{1}{\lambda}\right) < \infty \right\}. \quad (168)$$

If  $s^* = +\infty$ , then for all  $\alpha \in \mathbb{R}$ ,  $K_{Q,\mathbf{z}}(\alpha - \frac{1}{\lambda}) - K_{Q,\mathbf{z}}(-\frac{1}{\lambda}) < +\infty$ , and thus, the proof is completed by noticing that from Lemma 8.1, it holds that  $J_{\mathbf{z},Q,\lambda}(\alpha) = K_{Q,\mathbf{z}}(\alpha - \frac{1}{\lambda}) - K_{Q,\mathbf{z}}(-\frac{1}{\lambda})$ .

Alternatively, if  $s^* < +\infty$ , it follows that for all  $\alpha > s^*$ ,  $K_{Q,\mathbf{z}}(\alpha - \frac{1}{\lambda}) = +\infty$ . From Lemma 8.1, it holds that  $J_{\mathbf{z},Q,\lambda}(\alpha) = +\infty$ , which implies that  $+\infty = \alpha K_{Q,\mathbf{z}}^{(1)}(-\frac{1}{\lambda}) + \frac{\alpha^2}{2} K_{Q,\mathbf{z}}^{(2)}(\xi)$ , and thus,  $K_{Q,\mathbf{z}}^{(2)}(\xi)$  is infinite. Hence, in this case, the equality in (164) is of the form  $+\infty \leq +\infty$ . This completes the proof. ■

In (164), the parameter  $\xi$  depends on  $\alpha$ , as shown in (165). To highlight this dependence, in the following, the parameter  $\xi$  is denoted by  $\xi_\alpha$ . Using this notation, the focus is now on the term  $K_{Q,\mathbf{z}}^{(2)}(\xi_\alpha)$ , when  $\alpha \in \{t \in \mathbb{R} : J_{\mathbf{z},Q,\lambda}(t) < +\infty\}$ .

**Theorem 8.2.** *The function  $J_{\mathbf{z},Q,\lambda}$  in (141) verifies the following inequality, for all  $\alpha \in \{t \in \mathbb{R} : J_{\mathbf{z},Q,\lambda}(t) < +\infty\}$ ,*

$$J_{\mathbf{z},Q,\lambda}(\alpha) \leq \alpha K_{Q,\mathbf{z}}^{(1)}\left(-\frac{1}{\lambda}\right) + \frac{1}{2}\alpha^2 \beta_{Q,\mathbf{z}}^2 \quad (169)$$

where  $\beta_{Q,\mathbf{z}}$  is finite, and satisfies

$$\beta_{Q,\mathbf{z}} = \sup \left\{ \sqrt{K_{Q,\mathbf{z}}^{(2)}(\alpha)} : \alpha \in \left( -\infty, b - \frac{1}{\lambda} \right) \right\}, \quad (170)$$

with

$$b \triangleq \sup \{ t \in \mathbb{R} : J_{\mathbf{z},Q,\lambda}(t) < +\infty \}, \quad (171)$$

and the functions  $K_{Q,\mathbf{z}}^{(1)}$  and  $K_{Q,\mathbf{z}}^{(2)}$  defined in (97) and (98), respectively.

*Proof:* The proof of the inequality in (169) is trivial from Theorem 8.1 and the choice of  $\beta_{Q,\mathbf{z}}$  in (170). Hence, the remainder of the proof focuses on proving that  $\beta_{Q,\mathbf{z}} < +\infty$ . From Lemma 5.2 and Lemma 8.1, it holds that

$$\{ t \in \mathbb{R} : J_{\mathbf{z},Q,\lambda}(t) < +\infty \} = \left\{ t \in \mathbb{R} : K_{Q,\mathbf{z}} \left( t - \frac{1}{\lambda} \right) < +\infty \right\},$$

which implies that the set  $\{ t \in \mathbb{R} : J_{\mathbf{z},Q,\lambda}(t) < +\infty \}$  is an interval of the form  $(-\infty, b)$ , with  $b$  in (171). This follows from the fact that the function  $K_{Q,\mathbf{z}}$  is continuous and nondecreasing (Lemma 5.2) and the fact that

$$\lim_{t \rightarrow b} J_{\mathbf{z},Q,\lambda}(t) = +\infty. \quad (172)$$

For all  $\alpha \in (-\infty, b)$ , the function  $K_{Q,\mathbf{z}}^{(2)}$  is continuous (Lemma 5.2). Hence, for all  $t \in (\min \{-\frac{1}{\lambda}, \alpha - \frac{1}{\lambda}\}, \max \{-\frac{1}{\lambda}, \alpha - \frac{1}{\lambda}\}) \subset (-\infty, b)$ , the value  $K_{Q,\mathbf{z}}^{(2)}(t)$  is finite. Moreover, the values  $K_{Q,\mathbf{z}}^{(2)}(\min \{-\frac{1}{\lambda}, \alpha - \frac{1}{\lambda}\})$  and  $K_{Q,\mathbf{z}}^{(2)}(\max \{-\frac{1}{\lambda}, \alpha - \frac{1}{\lambda}\})$  are both finite. This implies that the function  $K_{Q,\mathbf{z}}^{(2)}$  achieves a minimum and maximum within the closed interval  $[\min \{-\frac{1}{\lambda}, \alpha - \frac{1}{\lambda}\}, \max \{-\frac{1}{\lambda}, \alpha - \frac{1}{\lambda}\}]$ . Thus, the corresponding term  $K_{Q,\mathbf{z}}^{(2)}(\xi_\alpha)$  is finite.

In the asymptotic regime, when  $\alpha \rightarrow -\infty$ , the following holds:

$$\lim_{\alpha \rightarrow -\infty} \xi_\alpha \in \left( -\infty, -\frac{1}{\lambda} \right). \quad (173)$$

The function  $K_{Q,\mathbf{z}}^{(2)}$  is continuous in  $(-\infty, -\frac{1}{\lambda})$ , as a consequence of the inclusion  $(-\infty, -\frac{1}{\lambda}) \subset (-\infty, b)$ , and thus, for all  $t \in (-\infty, -\frac{1}{\lambda})$ ,  $K_{Q,\mathbf{z}}^{(2)}(t) < +\infty$ . Moreover, from the assumption that  $\lambda \in \mathcal{K}_{Q,\mathbf{z}}$ , with  $\mathcal{K}_{Q,\mathbf{z}}$  in (23), it holds that

$$\lim_{t \rightarrow -\frac{1}{\lambda}} K_{Q,\mathbf{z}}^{(2)}(t) = K_{Q,\mathbf{z}}^{(2)} \left( -\frac{1}{\lambda} \right) < +\infty. \quad (174)$$

Alternatively,

$$\lim_{t \rightarrow -\infty} K_{Q,z}^{(2)}(t) = \lim_{t \rightarrow 0^+} K_{Q,z}^{(2)}\left(-\frac{1}{t}\right) \quad (175)$$

$$= \lim_{t \rightarrow 0^+} \int \left( \mathbb{L}_z(\boldsymbol{\theta}) - K_{Q,z}^{(1)}\left(-\frac{1}{t}\right) \right)^2 dP_{\boldsymbol{\Theta}|Z=z}^{(Q,t)}(\boldsymbol{\theta}), \quad (176)$$

$$= \lim_{t \rightarrow 0^+} \int (\mathbb{L}_z(\boldsymbol{\theta}))^2 dP_{\boldsymbol{\Theta}|Z=z}^{(Q,t)}(\boldsymbol{\theta}) - \lim_{t \rightarrow 0^+} \left( K_{Q,z}^{(1)}\left(-\frac{1}{t}\right) \right)^2, \quad (177)$$

$$= \lim_{t \rightarrow 0^+} \int (\mathbb{L}_z(\boldsymbol{\theta}))^2 dP_{\boldsymbol{\Theta}|Z=z}^{(Q,t)}(\boldsymbol{\theta}) - (\delta_{Q,z}^*)^2 \quad (178)$$

$$= \lim_{t \rightarrow 0^+} \int (\mathbb{L}_z(\boldsymbol{\theta}))^2 \frac{dP_{\boldsymbol{\Theta}|Z=z}^{(Q,t)}}{dQ}(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) - (\delta_{Q,z}^*)^2 \quad (179)$$

$$= \int (\mathbb{L}_z(\boldsymbol{\theta}))^2 \left( \lim_{t \rightarrow 0^+} \frac{dP_{\boldsymbol{\Theta}|Z=z}^{(Q,t)}}{dQ}(\boldsymbol{\theta}) \right) dQ(\boldsymbol{\theta}) - (\delta_{Q,z}^*)^2 \quad (180)$$

$$= \int (\mathbb{L}_z(\boldsymbol{\theta}))^2 \left( \frac{1}{Q(\mathcal{L}_{Q,z}^*)} \mathbb{1}_{\{\boldsymbol{\theta} \in \mathcal{L}_{Q,z}^*\}} \right) dQ(\boldsymbol{\theta}) - (\delta_{Q,z}^*)^2 \quad (181)$$

$$= \frac{1}{Q(\mathcal{L}_{Q,z}^*)} \int_{\mathcal{L}_{Q,z}^*} (\mathbb{L}_z(\boldsymbol{\theta}))^2 dQ(\boldsymbol{\theta}) - (\delta_{Q,z}^*)^2 \quad (182)$$

$$= (\delta_{Q,z}^*)^2 - (\delta_{Q,z}^*)^2 = 0, \quad (183)$$

where the equality in (176) follows from Lemma 5.4; the equality in (178) follows from Theorem 6.2, with  $\delta_{Q,z}^*$  in (38); the equality in (180) follows from the dominated convergence theorem [55, Theorem 1.6.9]; the equality in (181) follows from Lemma 3.6; and the equality in (183) follows from the definition of the set  $\mathcal{L}_{Q,z}^*$  in (39).

Hence, from (173), (174), and (183), it follows that

$$\lim_{\alpha \rightarrow -\infty} K_{Q,z}^{(2)}(\xi_\alpha) \in \left[ 0, \max_{c \in (-\infty, -\frac{1}{\lambda}]} K_{Q,z}^{(2)}(c) \right], \quad (184)$$

where the maximum exists and is finite.

On the other hand, in the asymptotic regime, when  $\alpha \rightarrow b^-$ , two cases are considered: (i)  $b \geq 0$ ; and (ii)  $b < 0$ . In the first case, the following holds from (165):

$$\lim_{\alpha \rightarrow b^-} \xi_\alpha \in \left[ -\frac{1}{\lambda}, b - \frac{1}{\lambda} \right). \quad (185)$$

The function  $K_{Q,z}^{(2)}$  is continuous in  $(-\frac{1}{\lambda}, b - \frac{1}{\lambda})$ , as a consequence of the inclusion  $(-\frac{1}{\lambda}, b - \frac{1}{\lambda}) \subset (-\infty, b)$ , and thus, for all  $t \in (-\frac{1}{\lambda}, b - \frac{1}{\lambda})$ ,  $K_{Q,z}^{(2)}(t) < +\infty$ .

Moreover,

$$K_{Q,\mathbf{z}}^{(2)}\left(-\frac{1}{\lambda}\right) < +\infty, \text{ and} \quad (186)$$

$$K_{Q,\mathbf{z}}^{(2)}\left(b - \frac{1}{\lambda}\right) < +\infty. \quad (187)$$

This implies that when  $b \geq 0$ ,

$$\lim_{\alpha \rightarrow b^-} K_{Q,\mathbf{z}}^{(2)}(\xi_\alpha) \in \left[0, \max_{c \in \left[-\frac{1}{\lambda}, b - \frac{1}{\lambda}\right]} K_{Q,\mathbf{z}}^{(2)}(c)\right], \quad (188)$$

where the maximum exists and is finite. Finally, In the second case, the following holds from (165):

$$\lim_{\alpha \rightarrow b^-} \xi_\alpha \in \left(b - \frac{1}{\lambda}, -\frac{1}{\lambda}\right). \quad (189)$$

The function  $K_{Q,\mathbf{z}}^{(2)}$  is continuous in  $\left(b - \frac{1}{\lambda}, -\frac{1}{\lambda}\right)$ , as a consequence of the inclusion  $\left(b - \frac{1}{\lambda}, -\frac{1}{\lambda}\right) \subset (-\infty, b)$ , and thus, for all  $t \in \left(b - \frac{1}{\lambda}, -\frac{1}{\lambda}\right)$ ,  $K_{Q,\mathbf{z}}^{(2)}(t) < +\infty$ . Moreover,

$$K_{Q,\mathbf{z}}^{(2)}\left(b - \frac{1}{\lambda}\right) < +\infty, \text{ and} \quad (190)$$

$$K_{Q,\mathbf{z}}^{(2)}\left(-\frac{1}{\lambda}\right) < +\infty. \quad (191)$$

This implies that when  $b < 0$ ,

$$\lim_{\alpha \rightarrow b^-} K_{Q,\mathbf{z}}^{(2)}(\xi_\alpha) \in \left[0, \max_{c \in \left(b - \frac{1}{\lambda}, -\frac{1}{\lambda}\right]} K_{Q,\mathbf{z}}^{(2)}(c)\right], \quad (192)$$

where the maximum exists and is finite. From all the above, it holds that for all  $\alpha \in \{t \in \mathbb{R} : J_{\mathbf{z},Q,\lambda}(t) < +\infty\}$ , the value  $K_{Q,\mathbf{z}}^{(2)}(\xi_\alpha)$  is finite, and this completes the proof.  $\blacksquare$

The main implication of Theorem 8.2 is that the random variable  $W$  in (100) is a sub-Gaussian random variable with sub-Gaussianity parameter  $\beta_{Q,\mathbf{z}}$  in (170) [53, Section 2.3]. This follows by noticing that the function  $J_{\mathbf{z},Q,\lambda}$  in (142) is the cumulant generating function of the random variable  $W$ . Hence, whenever it is finite, it is upper bounded as shown in Theorem 8.2. The following corollary of Theorem 8.2 highlights this observation.

**Corollary 8.1.** *The random variable  $W$  in (100) is a sub-Gaussian random variable with sub-Gaussianity parameter  $\beta_{Q,\mathbf{z}}$  in (170).*

The relevance of Corollary 8.1 is that it highlights the fact that when the models are sampled from the ERM-RER optimal measure  $P_{\Theta|Z=\mathbf{z}}^{(Q,\lambda)}$  in (25), the empirical risk with respect to the dataset  $\mathbf{z}$  is a sub-Gaussian random variable with sub-Gaussianity parameter  $\beta_{Q,\mathbf{z}}$  in (170).

## 9 Concentration of Probability

Consider the following set,

$$\mathcal{N}_{Q,z}(\lambda) \triangleq \left\{ \theta \in \mathcal{M} : \mathcal{L}_z(\theta) \leq \mathcal{R}_z \left( P_{\Theta|Z=z}^{(Q,\lambda)} \right) \right\}, \quad (193)$$

where the function  $\mathcal{L}_z$  is defined by (3); the functional  $\mathcal{R}_z$  is defined by (18); and the probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  is in (25). This section introduces two results. First, in Theorem 9.1, it is shown that when  $\lambda$  tends to zero, the set  $\mathcal{N}_{Q,z}(\lambda)$  forms an indexed family of sets that is monotonic and decreases to the set

$$\mathcal{N}_{Q,z}^* \triangleq \mathcal{L}_z(\delta_{Q,z}^*), \quad (194)$$

where  $\delta_{Q,z}^*$  is defined in (38); and the set  $\mathcal{L}_z(\delta_{Q,z}^*)$  is defined in (37). Second, in Theorem 9.2, it is shown that the probability  $P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{N}_{Q,z}(\lambda))$  strictly increases when  $\lambda$  tends to zero. More importantly, in Theorem 9.3, it is shown that the limit of the probability  $P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{N}_{Q,z}(\lambda))$ , when  $\lambda \rightarrow 0$ , is equal to one. These observations justify referring to the set  $\mathcal{N}_{Q,z}^*$  as the *limit set*. These observations are complementary to those stated in Section 3.2 and Section 3.3. This section ends by showing that the probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  concentrates on a specific subset  $\mathcal{L}_{Q,z}^*$  in (39) of the set  $\mathcal{N}_{Q,z}^*$ . At the light of this observation, the set  $\mathcal{L}_{Q,z}^*$  is referred to as the *nonnegligible limit set*. Finally, it is shown that when the  $\sigma$ -finite measure  $Q$  in (19) is coherent, the sets  $\mathcal{N}_{Q,z}^*$  and  $\mathcal{L}_{Q,z}^*$  are identical.

### 9.1 The Limit Set

The set  $\mathcal{N}_{Q,z}(\lambda)$  in (193), with  $\lambda \in \mathcal{K}_{Q,z}$  and  $\mathcal{K}_{Q,z}$  in (23), contains all the models that induce an empirical risk that is smaller than or equal to  $\mathcal{R}_z(P_{\Theta|Z=z}^{(Q,\lambda)})$ , i.e., the ERM-RER-optimal expected empirical risk in (105). This observation unveils the existence of a relation between the set  $\mathcal{N}_{Q,z}^*$  in (194) and the set  $\mathcal{T}(z)$  in (5), as shown by the following lemma.

**Lemma 9.1.** *The set  $\mathcal{N}_{Q,z}^*$  in (194) satisfies*

$$\mathcal{T}(z) \subseteq \mathcal{N}_{Q,z}^*, \quad (195)$$

where the set  $\mathcal{T}(z)$  is in (5). Moreover,

$$\mathcal{T}(z) = \mathcal{N}_{Q,z}^*, \quad (196)$$

if and only if (a) the ERM problem in (4) possesses a solution; and (b) the reference measure  $Q$  in (19) is coherent.

*Proof:* If the set  $\mathcal{T}(z)$  in (5) is empty, the inclusion in (195) is trivially true. Assume that  $|\mathcal{T}(z)| > 0$ . Hence, the proof of the inclusion in (195) follows from observing that for all  $\theta \in \mathcal{T}(z)$ , it holds that  $\mathcal{L}_z(\theta) = \rho^* \leq \delta_{Q,z}^*$ , with  $\delta_{Q,z}^*$  in (38) and  $\rho^*$  in (48). Hence,  $\theta \in \mathcal{N}_{Q,z}^*$ . This completes the proof of the inclusion in (195).

The proof of the equality in (196) is presented in two parts. In the first part, it is proved that if (196) holds, then the ERM problem in (4) possesses a solution and the measure  $Q$  is coherent. The second part proves the converse. The proof of the first part is as follows. Under the assumption that  $\mathcal{T}(z) = \mathcal{N}_{Q,z}^*$  holds, it follows that  $\delta_{Q,z}^* = \rho^*$ , with  $\rho^*$  in (48), which implies that the ERM problem in (4) possesses a solution. Moreover, for all  $\delta \in (\rho^*, +\infty)$ , it holds that  $Q(\mathcal{L}_z(\delta)) > 0$ , which verifies that the measure  $Q$  is coherent and completes the proof of the first part. The proof of the second part is as follows. Under the assumption that the ERM problem in (4) possesses a solution and the measure  $Q$  is coherent, it follows that  $\delta_{Q,z}^* = \rho^*$ . Hence,  $\mathcal{T}(z) = \mathcal{N}_{Q,z}^*$ , which completes the proof of the second part. ■

The following theorem highlights that the set  $\mathcal{N}_{Q,z}(\lambda)$  is decreasing with  $\lambda$ .

**Theorem 9.1.** *For all  $(\lambda_1, \lambda_2) \in \mathcal{K}_{Q,z} \times \mathcal{K}_{Q,z}$ , with  $\mathcal{K}_{Q,z}$  in (23) and  $\lambda_1 > \lambda_2$ , the sets  $\mathcal{N}_{Q,z}(\lambda_1)$  and  $\mathcal{N}_{Q,z}(\lambda_2)$  in (193) satisfy*

$$\mathcal{M} \supseteq \mathcal{N}_{Q,z}(\lambda_1) \supseteq \mathcal{N}_{Q,z}(\lambda_2) \supseteq \mathcal{N}_{Q,z}^*, \quad (197)$$

with  $\mathcal{N}_{Q,z}^*$  being the set defined in (194). Moreover, if the empirical risk function  $\mathbb{L}_z$  in (3) is continuous on  $\mathcal{M}$  and separable with respect to the measure  $Q$  in (19), then,

$$\mathcal{M} \supset \mathcal{N}_{Q,z}(\lambda_1) \supset \mathcal{N}_{Q,z}(\lambda_2) \supset \mathcal{N}_{Q,z}^*. \quad (198)$$

*Proof:* The proof is presented in Appendix Q. ■

An interesting observation is that for all  $\lambda \in \mathcal{K}_{Q,z}$ , with  $\mathcal{K}_{Q,z}$  in (23), only a subset of  $\mathcal{N}_{Q,z}(\lambda)$  might exhibit nonzero probability with respect to the measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  in (25). Consider for instance that the measure  $Q$  in (19) is noncoherent (Definition 4.1). That is,  $\delta_{Q,z}^* > \rho^*$ , with  $\delta_{Q,z}^*$  in (38) and  $\rho^*$  in (48). Thus, for all  $\gamma \in (\rho^*, \delta_{Q,z}^*)$ , it holds that  $Q(\mathcal{L}_z(\gamma)) = 0$ , with the set  $\mathcal{L}_z(\cdot)$  in (37). From Lemma 3.3, this implies that for all  $\gamma \in (\rho^*, \delta_{Q,z}^*)$ , the measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  in (25) satisfies  $P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_z(\gamma)) = 0$ , while verifying that  $\mathcal{L}_z(\gamma) \subseteq \mathcal{N}_{Q,z}(\lambda)$ . These observations lead to the analysis of the asymptotic concentration of probability in the following section.

## 9.2 The Nonnegligible Limit Set

The first step in the analysis of the asymptotic concentration of the probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  in (25) is to show that the probability  $P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{N}_{Q,z}(\lambda))$  increases when  $\lambda$  tends to zero, as shown by the following theorem.

**Theorem 9.2.** *For all  $(\lambda_1, \lambda_2) \in \mathcal{K}_{Q,z} \times \mathcal{K}_{Q,z}$ , with  $\mathcal{K}_{Q,z}$  in (23) and  $\lambda_1 > \lambda_2$ , assume that the measures  $P_{\Theta|Z=z}^{(Q,\lambda_1)}$  and  $P_{\Theta|Z=z}^{(Q,\lambda_2)}$  satisfy (25) with  $\lambda = \lambda_1$  and  $\lambda = \lambda_2$ , respectively. Then, the set  $\mathcal{N}_{Q,z}(\lambda_2)$  in (193) satisfies*

$$0 < P_{\Theta|Z=z}^{(Q,\lambda_1)}(\mathcal{N}_{Q,z}(\lambda_2)) \leq P_{\Theta|Z=z}^{(Q,\lambda_2)}(\mathcal{N}_{Q,z}(\lambda_2)), \quad (199)$$



where strict inequality holds if and only if the function  $L_z$  is separable with respect to the  $\sigma$ -finite measure  $Q$ .

*Proof:* The proof is presented in Appendix R. ■

The following lemma highlights a case in which a stronger concentration of probability is observed.

**Lemma 9.2.** *Let the function  $L_z$  in (3) be separable with respect to the  $\sigma$ -finite measure  $Q$  in (19). Let also  $(\lambda_1, \lambda_2) \in \mathcal{K}_{Q,z} \times \mathcal{K}_{Q,z}$ , with  $\mathcal{K}_{Q,z}$  in (23), be two positive reals such that  $\lambda_1 > \lambda_2$  and*

$$Q\left(\mathcal{N}_{Q,z}(\lambda_1) \cap (\mathcal{N}_{Q,z}(\lambda_2))^c\right) = 0, \quad (200)$$

with the complement with respect to the set of models  $\mathcal{M}$ . Then, two measures  $P_{\Theta|Z=z}^{(Q,\lambda_1)}$  and  $P_{\Theta|Z=z}^{(Q,\lambda_2)}$  that respectively satisfy (25) with  $\lambda = \lambda_1$  and  $\lambda = \lambda_2$  verify that

$$P_{\Theta|Z=z}^{(Q,\lambda_1)}(\mathcal{N}_{Q,z}(\lambda_1)) < P_{\Theta|Z=z}^{(Q,\lambda_2)}(\mathcal{N}_{Q,z}(\lambda_2)), \quad (201)$$

where, the set  $\mathcal{N}_{Q,z}(\cdot)$  is defined in (193).

*Proof:* The proof is presented in Appendix S. ■

The following example shows the relevance of Lemma 9.2 in the case in which the empirical risk function  $L_z$  in (3) is a simple function and separable with respect to the  $\sigma$ -finite measure  $Q$  in (19).

**Example 9.1.** *Consider Example 7.1. Note that, for all  $\lambda > 0$ ,*

$$0 < R_z\left(P_{\Theta|Z=z}^{(Q,\lambda)}\right) < 1, \quad (202)$$

where  $R_z\left(P_{\Theta|Z=z}^{(Q,\lambda)}\right)$  is the ERM-RER-optimal expected empirical risk in (105). The equality in (202) implies that given two reals  $\lambda_1$  and  $\lambda_2$  such that  $\lambda_1 > \lambda_2 > 0$ , it holds that,

$$\mathcal{N}_{Q,z}(\lambda_1) \cap (\mathcal{N}_{Q,z}(\lambda_2))^c = \left\{ \nu \in \mathcal{M} : R_z\left(P_{\Theta|Z=z}^{(Q,\lambda_2)}\right) < L_z(\nu) \leq R_z\left(P_{\Theta|Z=z}^{(Q,\lambda_1)}\right) \right\} \quad (203)$$

$$= \emptyset, \quad (204)$$

and moreover,  $\mathcal{N}_{Q,z}(\lambda_1) = \mathcal{N}_{Q,z}(\lambda_2)$ . Finally, from Lemma 9.2,

$$P_{\Theta|Z=z}^{(Q,\lambda_1)}(\mathcal{N}_{Q,z}(\lambda_1)) < P_{\Theta|Z=z}^{(Q,\lambda_2)}(\mathcal{N}_{Q,z}(\lambda_2)). \quad (205)$$

The main result of this section is presented by the following theorem.

**Theorem 9.3.** *The probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  in (25) satisfies*

$$\lim_{\lambda \rightarrow 0^+} P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{N}_{Q,z}(\lambda)) = 1, \quad (206)$$

where, the set  $\mathcal{N}_{Q,z}(\lambda)$  is defined in (193).

*Proof:* The proof follows immediately from Lemma 3.7 and by noticing that for all  $\lambda \in \mathcal{K}_{Q,z}$ , with  $\mathcal{K}_{Q,z}$  in (23), the sets  $\mathcal{L}_{Q,z}^*$  in (39) and  $\mathcal{N}_{Q,z}(\lambda)$  in (193) satisfies  $\mathcal{L}_{Q,z}^* \subseteq \mathcal{N}_{Q,z}(\lambda)$ . ■

Note that Theorem 9.3 and Lemma 3.7 lead to the following conclusion

$$\lim_{\lambda \rightarrow 0^+} P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{N}_{Q,z}(\lambda) \setminus \mathcal{L}_{Q,z}^*) = 0, \quad (207)$$

which follows from the fact that  $\mathcal{L}_{Q,z}^* \subset \mathcal{N}_{Q,z}(\lambda)$ , with  $\mathcal{L}_{Q,z}^*$  in (39). This justifies referring to the set  $\mathcal{L}_{Q,z}^*$  as the nonnegligible limit set.

## 10 $(\delta, \epsilon)$ -Optimality

This section introduces a PAC guarantee of optimality for the models that are sampled from the probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  in (25) with respect to the ERM problem in (4). Such guarantee is defined as follows.

**Definition 10.1**  $((\delta, \epsilon)$ -Optimality). *Given a pair of positive reals  $(\delta, \epsilon)$ , with  $\epsilon < 1$ , the probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  in (25) is said to be  $(\delta, \epsilon)$ -optimal, if the set  $\mathcal{L}_z(\delta)$  in (37) satisfies*

$$P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_z(\delta)) > 1 - \epsilon. \quad (208)$$

If the probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  in (25) is  $(\delta, \epsilon)$ -optimal, then it assigns a probability that is always greater than  $1 - \epsilon$  to a set that contains models that induce an empirical risk that is smaller than  $\delta$ . From this perspective, particular interest is given to the smallest  $\delta$  and  $\epsilon$  for which  $P_{\Theta|Z=z}^{(Q,\lambda)}$  is  $(\delta, \epsilon)$ -optimal.

The main result of this section is presented by the following theorem.

**Theorem 10.1.** *For all  $(\delta, \epsilon) \in (\delta_{Q,z}^*, +\infty) \times (0, 1)$ , with  $\delta_{Q,z}^*$  in (38), there exists a real  $\lambda \in \mathcal{K}_{Q,z}$ , with  $\mathcal{K}_{Q,z}$  in (23), such that the probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  is  $(\delta, \epsilon)$ -optimal.*

*Proof:* Let  $\delta$  be a real in  $(\delta_{Q,z}^*, +\infty)$ , with  $\delta_{Q,z}^*$  in (38). Let also  $\lambda \in \mathcal{K}_{Q,z}$  satisfy the following equality:

$$K_{Q,z}^{(1)}\left(-\frac{1}{\lambda}\right) \leq \delta. \quad (209)$$

Note that from Lemma 5.2, it follows that the function  $K_{Q,z}^{(1)}$  is continuous. Moreover, from Theorem 6.2, it follows that such a  $\lambda$  in (209) always exists. From (37) and (193), it holds that

$$\mathcal{N}_{Q,z}(\lambda) \subseteq \mathcal{L}_z(\delta), \quad (210)$$

and thus,

$$P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_z(\delta)) \geq P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{N}_{Q,z}(\lambda)). \quad (211)$$

Let  $\gamma$  be a positive real such that  $\gamma \leq \lambda$  and

$$P_{\Theta|Z=z}^{(Q,\gamma)}(\mathcal{N}_{Q,z}(\gamma)) > 1 - \epsilon. \quad (212)$$

The existence of such a positive real  $\gamma$  follows from Theorem 9.3. Hence, from (212), it holds that,

$$1 - \epsilon < P_{\Theta|Z=z}^{(Q,\gamma)}(\mathcal{N}_{Q,z}(\gamma)) \quad (213)$$

$$\leq P_{\Theta|Z=z}^{(Q,\gamma)}(\mathcal{L}_z(\delta)), \quad (214)$$

where the inequality in (214) follows from the fact that  $\mathcal{N}_{Q,z}(\gamma) \subseteq \mathcal{N}_{Q,z}(\lambda) \subseteq \mathcal{L}_z(\delta)$ . Finally, the inequality in (214) implies that the probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  is  $(\delta, \epsilon)$ -optimal (Definition 10.1). This completes the proof. ■

A stronger optimality claim can be stated when the reference measure is coherent.

**Theorem 10.2.** *For all  $(\delta, \epsilon) \in (\rho^*, +\infty) \times (0, 1)$ , with  $\rho^*$  in (48), there always exists a  $\lambda \in \mathcal{K}_{Q,z}$ , with  $\mathcal{K}_{Q,z}$  in (23), such that the probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  is  $(\delta, \epsilon)$ -optimal if and only if the reference measure  $Q$  is coherent.*

*Proof:* The proof is divided into two parts. The first part shows that if for all  $(\delta, \epsilon) \in (\rho^*, +\infty) \times (0, 1)$ , there always exists a  $\lambda \in \mathcal{K}_{Q,z}$ , with  $\mathcal{K}_{Q,z}$  in (23), such that the probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  in (25) is  $(\delta, \epsilon)$ -optimal, then, the measure  $Q$  is coherent. The second part deals with the converse.

The first part is as follows. Let  $\gamma \in \mathcal{K}_{Q,z}$  be such that

$$P_{\Theta|Z=z}^{(Q,\gamma)}(\mathcal{L}_z(\delta)) > 1 - \epsilon, \quad (215)$$

then, for all measurable subsets  $\mathcal{A}$  of  $\mathcal{L}_z(\delta)$ , it holds that

$$\begin{aligned} 1 - \epsilon &< P_{\Theta|Z=z}^{(Q,\gamma)}(\mathcal{L}_z(\delta)) \\ &= \int_{\mathcal{A}} \frac{dP_{\Theta|Z=z}^{(Q,\gamma)}(\nu)}{dQ} dQ(\nu) + \int_{\mathcal{L}_z(\delta) \setminus \mathcal{A}} \frac{dP_{\Theta|Z=z}^{(Q,\gamma)}(\nu)}{dQ} dQ(\nu), \end{aligned} \quad (216)$$

which, together with Lemma 3.2, implies that there exists at least one measurable subset  $\mathcal{A}$  for which  $Q(\mathcal{A}) > 0$ , and thus,

$$Q(\mathcal{L}_z(\delta)) > Q(\mathcal{A}) > 0, \quad (217)$$

which implies that the measure  $Q$  is coherent. This completes the first part of the proof.

The second part of the proof is as follows. Under the assumption that the measure  $Q$  is coherent, it follows that  $\delta_{Q,z}^* = \rho^*$ . Then, from Theorem 10.1, it

follows that for all  $(\delta, \epsilon) \in (\delta_{Q,z}^*, +\infty) \times (0, 1)$ , there always exists a  $\lambda \in \mathcal{K}_{Q,z}$ , with  $\mathcal{K}_{Q,z}$  in (23), such that the probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  is  $(\delta, \epsilon)$ -optimal. This completes the second part of the proof. ■

## 11 Sensitivity and Generalization

This section introduces the notion of sensitivity and establishes its connections with the notion of generalization error of the Gibbs algorithm, cf. [8].

### 11.1 Sensitivity

The sensitivity of the expected empirical risk  $R_z$  in (18) to deviations from the probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  in (25) towards an alternative probability measure  $P \in \Delta(\mathcal{M}, \mathcal{B}(\mathcal{M}))$  is introduced as a novel metric to evaluate the generalization capabilities of the ERM-RER-optimal measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$ . Deviations from the probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  towards an alternative probability measure  $P$  would allow comparing the ERM-RER-optimal measure with alternative measures (or algorithms). For instance, if new datasets become available, a new ERM-RER problem can be formulated using a larger dataset obtained by aggregating the old and the new datasets, cf. [10] and [68]. Intuitively, the ERM-RER-optimal measure obtained after the aggregation of datasets might exhibit better generalization capabilities, see for instance [10]. This analysis is the motivation of the sensitivity, which is defined as follows.

**Definition 11.1** (Sensitivity). *Given the  $\sigma$ -finite measure  $Q$  and the positive real  $\lambda > 0$  in (19), let  $S_{Q,\lambda} : (\mathcal{X} \times \mathcal{Y})^n \times \Delta_Q(\mathcal{M}, \mathcal{B}(\mathcal{M})) \rightarrow (-\infty, +\infty]$  be a functional such that*

$$S_{Q,\lambda}(z, P) = \begin{cases} R_z(P) - R_z(P_{\Theta|Z=z}^{(Q,\lambda)}) & \text{if } \lambda \in \mathcal{K}_{Q,z} \\ +\infty & \text{otherwise,} \end{cases} \quad (218)$$

where the functional  $R_z$  is defined in (18) and the probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  is in (25). The sensitivity of the expected empirical risk  $R_z$  due to a deviation from  $P_{\Theta|Z=z}^{(Q,\lambda)}$  to  $P$  is  $S_{Q,\lambda}(z, P)$ .

Recently, the following exact expression for the sensitivity  $S_{Q,\lambda}(z, P)$  in (218) was introduced in [10].

**Theorem 11.1** (Theorem 1 in [10]). *The sensitivity  $S_{Q,\lambda}(z, P)$  in (218) satisfies*

$$S_{Q,\lambda}(z, P) = \lambda \left( D(P_{\Theta|Z=z}^{(Q,\lambda)} \| Q) + D(P \| P_{\Theta|Z=z}^{(Q,\lambda)}) - D(P \| Q) \right), \quad (219)$$

where the probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  is in (25).

The following theorem introduces an upper bound on the absolute value of the sensitivity  $S_{Q,\lambda}(z, P)$  in (218), which requires the calculation of only one of the relative entropies in Theorem 11.1.

**Theorem 11.2.** *For all  $P \in \Delta_Q(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ , the sensitivity  $S_{Q,\lambda}(z, P)$  in (218) satisfies*

$$|S_{Q,\lambda}(z, P)| \leq \sqrt{2\beta_{Q,z}^2 D(P \| P_{\Theta|Z=z}^{(Q,\lambda)})}, \quad (220)$$

where the constant  $\beta_{Q,z}$  is defined in (170).

*Proof:* The proof is presented in Appendix T. ■

Note that equality holds in (220) in the trivial case in which the empirical risk function is not separable with respect to  $Q$  (Definition 5.1). In such case, for all  $P \in \Delta_Q(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ , it holds that  $S_{Q,\lambda}(z, P) = 0$  and  $\beta_{Q,z} = 0$ .

Theorem 11.2 establishes an upper and a lower bound on the increase and decrease of the expected empirical risk that can be obtained by deviating from the optimal solution of the ERM-RER problem in (19). More specifically, note that for all probability measures  $P \in \Delta_Q(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ , it holds that,

$$R_z(P) \geq R_z(P_{\Theta|Z=z}^{(Q,\lambda)}) - \sqrt{2\beta_{Q,z}^2 D(P \| P_{\Theta|Z=z}^{(Q,\lambda)})} \quad \text{and} \quad (221)$$

$$R_z(P) \leq R_z(P_{\Theta|Z=z}^{(Q,\lambda)}) + \sqrt{2\beta_{Q,z}^2 D(P \| P_{\Theta|Z=z}^{(Q,\lambda)})}. \quad (222)$$

## 11.2 Generalization Error

This section unveils the interesting connection between the notion of sensitivity and the notion of generalization error of the Gibbs algorithm, cf. [8]. The Generalization error is defined under the assumption that datasets are sampled from a probability measure

$$P_Z \in \Delta((\mathcal{X} \times \mathcal{Y})^n, \mathcal{F}), \quad (223)$$

where  $\mathcal{F}$  denotes a given  $\sigma$ -field on the set  $(\mathcal{X} \times \mathcal{Y})^n$ . For such a probability measure  $P_Z$  in (223), let the set  $\mathcal{K}_{Q,P_Z} \subset \mathbb{R}$  be

$$\mathcal{K}_{Q,P_Z} = \bigcap_{z \in \text{supp } P_Z} \mathcal{K}_{Q,z}, \quad (224)$$

where the  $\sigma$ -finite measure  $Q$  is in (19). The set  $\mathcal{K}_{Q,P_Z}$  in (224) can be empty for some choices of the  $\sigma$ -finite measure  $Q$ . Nonetheless, from Lemma 3.1, it follows that if  $Q$  is a probability measure, then,

$$\mathcal{K}_{Q,P_Z} = (0, +\infty). \quad (225)$$

Under the assumption that datasets are sampled from  $P_Z$  in (223), the generalization error of the Gibbs algorithm with parameters  $Q$  and  $\lambda$ , is defined as the

expectation with respect to the product measure  $P_{\Theta|Z}^{(Q,\lambda)} \cdot P_Z$ , with  $P_{\Theta|Z}^{(Q,\lambda)}$  in (25), of difference between: (a) the *population risk* due to a model  $\theta \in \mathcal{M}$ ,

$$\int \mathbf{L}_z(\theta) dP_Z(z) \quad (226)$$

with the function  $\mathbf{L}_z$  defined in (3); and (b) The empirical risk induced by the model  $\theta$  with respect to a training dataset  $z$ , that is,  $\mathbf{L}_z(\theta)$ . More specifically, the generalization error of the Gibbs algorithm with parameters  $Q$  and  $\lambda$  is

$$\begin{aligned} & \int \int \left( \int \mathbf{L}_z(\theta) dP_Z(z) - \mathbf{L}_\nu(\theta) \right) dP_{\Theta|Z=\nu}^{(Q,\lambda)}(\theta) dP_Z(\nu) \\ &= \int \left( \int \mathbf{L}_z(\theta) dP_Z(z) \right) dP_{\Theta}^{(Q,\lambda)}(\theta) \\ & \quad - \int \mathbf{L}_\nu(\theta) dP_{\Theta|Z=\nu}^{(Q,\lambda)}(\theta) dP_Z(\nu) \end{aligned} \quad (227)$$

$$\begin{aligned} &= \int \left( \int \mathbf{L}_z(\theta) dP_{\Theta}^{(Q,\lambda)}(\theta) \right) dP_Z(z) \\ & \quad - \int \mathbf{L}_\nu(\theta) dP_{\Theta|Z=\nu}^{(Q,\lambda)}(\theta) dP_Z(\nu) \end{aligned} \quad (228)$$

$$= \int \left( \mathbf{R}_\nu(P_{\Theta}^{(Q,\lambda)}) - \mathbf{R}_\nu(P_{\Theta|Z=\nu}^{(Q,\lambda)}) \right) dP_Z(\nu), \quad (229)$$

where the probability measure  $P_{\Theta}^{(Q,\lambda)}$  satisfies for all sets  $\mathcal{A} \in \mathcal{B}(\mathcal{M})$ ,

$$P_{\Theta}^{(Q,\lambda)}(\mathcal{A}) = \int P_{\Theta|Z=\nu}^{(Q,\lambda)}(\mathcal{A}) dP_Z(\nu), \quad (230)$$

and the functional  $\mathbf{R}_\nu$  is defined in (18).

The following theorem establishes a connection between sensitivity and generalization error in the particular case in which  $Q$  in (19) is a probability measure.

**Theorem 11.3.** *Under the assumption that datasets are sampled from  $P_Z$  in (223), the generalization error of the Gibbs algorithm with parameters  $Q$  (a probability measure) and  $\lambda > 0$ , is*

$$\int S_{Q,\lambda}(\nu, P_{\Theta}^{(Q,\lambda)}) dP_Z(\nu), \quad (231)$$

where the functional  $S_{Q,\lambda}$  is in (218); and the probability measure  $P_{\Theta}^{(Q,\lambda)}$  is in (230).

*Proof:* The proof uses the fact that under the assumption that  $Q$  is a probability measure, for all  $\nu \in \text{supp } P_Z$ , it follows from Lemma 3.1 that  $\mathcal{K}_{Q,\nu} = (0, +\infty)$ . This implies that for all  $z \in \text{supp } P_Z$  and for all  $\lambda > 0$ , the ERM-RER problem in (19), always possesses as solution the measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  in (25). Thus,

the measure  $P_{\Theta}^{(Q,\lambda)}$  in (230) is well defined. Moreover,  $S_{Q,\lambda}(z, P_{\Theta}^{(Q,\lambda)}) = R_z(P_{\Theta}^{(Q,\lambda)}) - R_z(P_{\Theta|Z=z}^{(Q,\lambda)})$  and the integral in (229) is also well defined, which completes the proof. ■

Theorem 11.3 provides an interesting viewpoint of the generalization error. For instance, the probability measure  $P_{\Theta}^{(Q,\lambda)}$  in (223) can be understood as the barycenter of a subset of  $\Delta(\mathcal{M}, \mathcal{B}(\mathcal{M}))$  containing the solutions to ERM-RER problems of the form in (19), with  $z \in \text{supp } P_Z$  in (223). Hence, the generalization error of the Gibbs algorithm is the expectation (with respect to  $P_Z$ ) of the sensitivity of the expected empirical risks  $R_z$  in (18) to variations from the ERM-RER-optimal measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  towards the barycenter, i.e., the measure  $P_{\Theta}^{(Q,\lambda)}$ .

The following definition extends the notion of generalization error to Gibbs algorithms obtained by assuming that the reference measure  $Q$  in (19) is a  $\sigma$ -finite measure. This definition also exploits the relation between the notions of sensitivity and generalization error introduced by Theorem 11.3.

**Definition 11.2** (Generalization Error of the Gibbs Algorithm). *Given a  $\sigma$ -finite measure  $Q \in \Delta(\mathcal{M}, \mathcal{B}(\mathcal{M}))$  and a real  $\lambda > 0$ , let the functional  $G_{Q,\lambda} : \Delta((\mathcal{X} \times \mathcal{Y})^n, \mathcal{F}) \rightarrow (-\infty, +\infty]$  be such that*

$$G_{Q,\lambda}(P_Z) = \begin{cases} \int S_{Q,\lambda}(\nu, P_{\Theta}^{(Q,\lambda)}) dP_Z(\nu) & \text{if } \lambda \in \mathcal{K}_{Q,P_Z} \\ +\infty & \text{otherwise,} \end{cases} \quad (232)$$

where the functional  $S_{Q,\lambda}$  is in (218); the set  $\mathcal{K}_{Q,P_Z}$  is in (224); and the probability measure  $P_{\Theta}^{(Q,\lambda)}$  is in (230). The generalization error induced by the Gibbs algorithm with parameters  $Q$  and  $\lambda$  under the assumption that datasets are sampled from the probability measure  $P_Z$ , is  $G_{Q,\lambda}(P_Z)$ .

The main difficulty for extending the notion of generalization error to Gibbs algorithms obtained under the assumption that the reference measure is not a probability measure, but a  $\sigma$ -finite measure, is that the integrals in (229) and (230) might not be well defined. This is essentially due to the fact that, while the ERM-RER problem in (19) always possesses a solution when  $Q$  is a probability measure, the existence of a solution when  $Q$  is not a probability measure is subject to the condition that for all  $z \in \text{supp } P_Z$ ,  $\lambda \in \mathcal{K}_{Q,z}$ , with  $\mathcal{K}_{Q,z}$  in (23). This leads to the condition that  $\lambda \in \mathcal{K}_{Q,P_Z}$ , with the set  $\mathcal{K}_{Q,P_Z}$  in (224). When such a condition is not met, the definition of sensitivity is void.

The following theorem provides a closed-form expression for the generalization error of the Gibbs algorithm in the general case in which the reference measure  $Q$  in (19) is a  $\sigma$ -finite measure.

**Theorem 11.4.** *If  $\lambda \in \mathcal{K}_{Q,P_Z}$ , with  $\mathcal{K}_{Q,P_Z}$  in (224), the generalization error*

$G_{Q,\lambda}(P_Z)$  in (232) satisfies

$$\begin{aligned} & G_{Q,\lambda}(P_Z) \\ &= \lambda \left( \int D(P_{\Theta|Z=\nu}^{(Q,\lambda)} \| P_{\Theta}^{(Q,\lambda)}) dP_Z(\nu) + \int D(P_{\Theta}^{(Q,\lambda)} \| P_{\Theta|Z=\nu}^{(Q,\lambda)}) dP_Z(\nu) \right), \end{aligned} \quad (233)$$

where for all  $z \in \text{supp } P_Z$ , the probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  is in (25); and the probability measure  $P_{\Theta}^{(Q,\lambda)}$  is defined in (230).

*Proof:* The proof is presented in Appendix U. ■

The terms  $\int D(P_{\Theta|Z=\nu}^{(Q,\lambda)} \| P_{\Theta}^{(Q,\lambda)}) dP_Z(\nu)$  and  $\int D(P_{\Theta}^{(Q,\lambda)} \| P_{\Theta|Z=\nu}^{(Q,\lambda)}) dP_Z(\nu)$  in the right-hand side of (233) are respectively the mutual and the lautum information [52] induced by a joint probability measure  $P_{\Theta,Z}$  whose marginals are  $P_Z$  in (223) and  $P_{\Theta}^{(Q,\lambda)}$  in (230). When the reference measure  $Q$  in (19) is a probability measure, Theorem 11.4 reduces to [8, Theorem 1]. Interestingly, independently of whether the reference measure  $Q$  in (19) is a probability measure, or whether the  $n$  data points in the datasets are independent and identically distributed, the generalization error  $G_{Q,\lambda}(P_Z)$  in (232) is always a factor of the sum of the mutual and lautum information induced by the joint probability measure  $P_{\Theta,Z}$  mentioned above.

Theorem 11.4 also provides an alternative interpretation of the generalization error  $G_{Q,\lambda}(P_Z)$  in (232). Note that by writing one of the factors in the right-hand side of (233) as

$$\int \left( D(P_{\Theta|Z=\nu}^{(Q,\lambda)} \| P_{\Theta}^{(Q,\lambda)}) + D(P_{\Theta}^{(Q,\lambda)} \| P_{\Theta|Z=\nu}^{(Q,\lambda)}) \right) dP_Z(\nu),$$

it becomes clear that  $G_{Q,\lambda}(P_Z)$  is the expectation with respect to  $P_Z$  of the symmetrized Kullback-Leibler divergence, also known as Jeffrey's divergence [65], of the probability measures  $P_{\Theta|Z=z}^{(Q,\lambda)}$  and  $P_{\Theta}^{(Q,\lambda)}$ . That is, the solution to the ERM-RER problem in (19) and the barycenter induced by  $P_Z$ .

The following theorem provides an upper-bound on the generalization error of the Gibbs algorithm only in terms of the lautum information induced by such a joint probability measure  $P_{\Theta,Z}$ .

**Theorem 11.5.** *The generalization error  $G_{Q,\lambda}(P_Z)$  in (232) satisfies for all  $\lambda \in \mathcal{K}_{Q,P_Z}$ ,*

$$0 \leq G_{Q,\lambda}(P_Z) \leq \sqrt{2\sigma_Q^2 \int D(P_{\Theta}^{(Q,\lambda)} \| P_{\Theta|Z=\nu}^{(Q,\lambda)}) dP_Z(\nu)}, \quad (234)$$

where for all  $z \in \text{supp } P_Z$ , the probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  is in (25); the probability measure  $P_{\Theta}^{(Q,\lambda)}$  is defined in (230); and

$$\sigma_Q = \sup \{ \beta_{Q,z} : z \in (\mathcal{X} \times \mathcal{Y})^n \}, \quad (235)$$

with  $\beta_{Q,z}$  in (165).



*Proof:* The proof of the inequality  $G_{Q,\lambda}(P_Z) \geq 0$  follows from observing that for all  $\nu \in (\mathcal{X} \times \mathcal{Y})^n$ , the terms  $D(P_{\Theta|Z=\nu}^{(Q,\lambda)} \| P_{\Theta}^{(Q,\lambda)})$  and  $D(P_{\Theta}^{(Q,\lambda)} \| P_{\Theta|Z=\nu}^{(Q,\lambda)})$  in (233) are nonnegative (Theorem 2.1). The proof of the remaining inequality follows from (232) and the following inequalities:

$$G_{Q,\lambda}(P_Z) = \left| \int S_{Q,\lambda}(\nu, P_{\Theta}^{(Q,\lambda)}) dP_Z(\nu) \right| \quad (236)$$

$$\leq \int |S_{Q,\lambda}(\nu, P_{\Theta}^{(Q,\lambda)})| dP_Z(\nu) \quad (237)$$

$$\leq \int \sqrt{2\beta_{Q,\nu} D(P_{\Theta}^{(Q,\lambda)} \| P_{\Theta|Z=\nu}^{(Q,\lambda)})} dP_Z(\nu), \quad (238)$$

$$\leq \int \sqrt{2\sigma_Q^2 D(P_{\Theta}^{(Q,\lambda)} \| P_{\Theta|Z=\nu}^{(Q,\lambda)})} dP_Z(\nu) \quad (239)$$

$$\leq \sqrt{2\sigma_Q^2 \int D(P_{\Theta}^{(Q,\lambda)} \| P_{\Theta|Z=\nu}^{(Q,\lambda)}) dP_Z(\nu)}, \quad (240)$$

where the equality in (236) follows from (232); the inequality in (237) follows from [55, Theorem 1.5.9(c)]; the inequality in (238) follows from Theorem 11.2; the inequality in (239) follows from (235); and the inequality in (240) follows from Jensen's inequality [55, Section 6.3.5]. This completes the proof. ■

In a nutshell, the generalization error  $G_{Q,\lambda}(P_Z)$  in (232) is upper bounded up to a constant factor by the square root of the lautum information induced by the joint probability measure  $P_{\Theta,Z}$  mentioned above. Theorem 11.5 is reminiscent of [29, Theorem 1], which provides a similar upper-bound on  $G_{Q,\lambda}(P_Z)$  using the mutual information instead of the lautum information induced by the joint probability measure  $P_{\Theta,Z}$ . The interest in Theorem 11.5 for the specific case of the Gibbs algorithm, lies on the fact that it holds under milder conditions than those in [29, Theorem 1]. For instance, no additional conditions on the loss function  $\ell$  in (2) concerning sub-Gaussianity are assumed. Moreover, the probability measure  $P_Z$  from which datasets are sampled is not necessarily a product measure.

## 12 Conclusions and Final Remarks

The classical ERM-RER problem in (19) has been studied under the assumption that the reference measure  $Q$  is a  $\sigma$ -finite measure, instead of a probability measure, which leads to a more general problem that includes the ERM problem with (discrete or differential) entropy regularization and the information-risk minimization problem. While in the case in which the reference measure is a probability measure the solution to the ERM-RER problem always exists, in this general case, the existence of a solution is subject to a condition that depends on the loss function, the reference measure, the regularization factor, and the training dataset. When a solution exists, it has been proved that it is unique. Additionally, if it exists, such a solution and the reference measure are mutually

absolutely continuous in most of the practical cases of interest. Interestingly, the empirical risk observed when models are sampled from the ERM-RER-optimal probability measure is a sub-Gaussian random variable that exhibits a PAC guarantee for the ERM problem. That is, for some positive  $\delta$  and  $\epsilon$ , it is shown that there always exist some parameters for the ERM-RER problem such that the set of models that induce an empirical risk smaller than  $\delta$  exhibits a probability that is not smaller than  $1 - \epsilon$ . Interestingly, none of these results relies on statistical assumptions on the datasets.

The sensitivity of the expected empirical risk to deviations from the ERM-RER-optimal measure to alternative measures is introduced as a new performance metric to evaluate the generalization capabilities of the Gibbs algorithm. In particular, an upper bound on the absolute value of the sensitivity, which depends on the training dataset, is presented. This bound is formed by a constant factor and the square root of the relative entropy of the alternative measure (the deviation) with respect to the ERM-RER solution. Finally, it is shown that the expectation of the sensitivity (with respect to the datasets) to deviations towards a particular measure is equivalent to the generalization error of the Gibbs algorithm. Equipped with this observation, the generalization error is shown to be in the most general case, up to a constant factor, the sum of the mutual and lautum information between the models and the datasets, which was a result known exclusively for the case in which the reference is a probability measure, cf. [8]. From this perspective, it is argued that the study of the generalization capabilities of the Gibbs algorithm based on generalization error is a significantly narrow view. This is essentially because it is looking at an expectation of the sensitivity to deviations to a particular measure, i.e., the barycenter of the set of ERM-RER solutions induced by a prior on the datasets. A broader view is offered by the study of the sensitivity to deviations towards other measures, i.e., ERM-RER-optimal measures obtained with different training data sets. This approach has lead already to a few initial results in [10] that highlight the connections to sensitivity, training error, and test error. Nonetheless, the study of the sensitivity in the aim of describing the generalization capabilities of learning algorithms remains by now as an open problem.

## Appendices

### A Proof of Theorem 2.2

Consider the function  $f : [0, +\infty) \rightarrow \mathbb{R}$  such that

$$f(x) = \begin{cases} x \log(x) & \text{if } x > 0 \\ 0 & \text{if } x = 0, \end{cases} \quad (241)$$

and note that it is strictly convex. From the assumption that for all  $i \in \{1, 2\}$ ,  $P_i$  and  $Q_i$  are both measures on the same measurable space  $(\Omega, \mathcal{F})$ , with  $P_i$  absolutely continuous with respect to  $Q_i$ , let  $g : \Omega \rightarrow [0, \infty)$  be the function

$$g(x) = \frac{d(\lambda P_1 + (1-\lambda)P_2)}{d(\lambda Q_1 + (1-\lambda)Q_2)}(x), \quad (242)$$

where  $\frac{d(\lambda P_1 + (1-\lambda)P_2)}{d(\lambda Q_1 + (1-\lambda)Q_2)}$  is the Radon-Nikodym derivative of the measure  $\lambda P_1 + (1-\lambda)P_2$  with respect to  $\lambda Q_1 + (1-\lambda)Q_2$ . Using this notation, for all  $\lambda \in (0, 1)$ ,

$$\begin{aligned} & D(\lambda P_1 + (1-\lambda)P_2 \| \lambda Q_1 + (1-\lambda)Q_2) \\ & - \lambda D(P_1 \| Q_1) + (1-\lambda) D(P_2 \| Q_2) \end{aligned} \quad (243)$$

$$\begin{aligned} &= \int \log(g(x)) d(\lambda P_1 + (1-\lambda)P_2)(x) \\ &= \lambda \int \log\left(\frac{dP_1}{dQ_1}(x)\right) dP_1(x) - (1-\lambda) \int \log\left(\frac{dP_2}{dQ_2}(x)\right) dP_2(x) \\ &= \lambda \int \log(g(x)) dP_1(x) + (1-\lambda) \int \log(g(x)) dP_2(x) \\ &= \lambda \int \log\left(\frac{dP_1}{dQ_1}(x)\right) dP_1(x) - (1-\lambda) \int \log\left(\frac{dP_2}{dQ_2}(x)\right) dP_2(x) \\ &= \lambda \int \log\left(\left(\frac{dP_1}{dQ_1}(x)\right)^{-1} g(x)\right) dP_1(x) \\ &\quad + (1-\lambda) \int \log\left(\left(\frac{dP_2}{dQ_2}(x)\right)^{-1} g(x)\right) dP_2(x) \\ &= \lambda \int \frac{dP_1}{dQ_1}(x) \log\left(\left(\frac{dP_1}{dQ_1}(x)\right)^{-1} g(x)\right) dQ_1(x) \\ &\quad + (1-\lambda) \int \frac{dP_2}{dQ_2}(x) \log\left(\left(\frac{dP_2}{dQ_2}(x)\right)^{-1} g(x)\right) dQ_2(x) \\ &= \lambda \int \frac{g(x) \frac{dP_1}{dQ_1}(x)}{g(x)} \log\left(\left(\frac{dP_1}{dQ_1}(x)\right)^{-1} g(x)\right) dQ_1(x) \\ &\quad + (1-\lambda) \int \frac{g(x) \frac{dP_2}{dQ_2}(x)}{g(x)} \log\left(\left(\frac{dP_2}{dQ_2}(x)\right)^{-1} g(x)\right) dQ_2(x) \\ &= -\lambda \int g(x) f\left(\frac{dP_1}{dQ_1}(x)(g(x))^{-1}\right) dQ_1(x) - (1-\lambda) \int g(x) f\left(\frac{dP_2}{dQ_2}(x)(g(x))^{-1}\right) dQ_2(x), \end{aligned} \quad (244)$$

where the function  $f$  is defined in (241). Let  $\beta_1$  and  $\beta_2$  be the following constants:

$$\beta_1 \triangleq \int g(\nu) dQ_1(\nu) \text{ and } \beta_2 \triangleq \int g(\nu) dQ_2(\nu). \quad (245)$$

From (244) and (245), it follows that for all  $\lambda \in (0, 1)$ ,

$$\begin{aligned} & D(\lambda P_1 + (1 - \lambda)P_2 \| \lambda Q_1 + (1 - \lambda)Q_2) - \lambda D(P_1 \| Q_1) + (1 - \lambda)D(P_2 \| Q_2) \\ &= -\lambda \beta_1 \int \frac{g(x)}{\beta_1} f\left(\frac{dP_1}{dQ_1}(x)(g(x))^{-1}\right) dQ_1(x) \\ &\quad - (1 - \lambda) \beta_2 \int \frac{g(x)}{\beta_2} f\left(\frac{dP_2}{dQ_2}(x)(g(x))^{-1}\right) dQ_2(x) \end{aligned} \quad (246)$$

$$\leq -\lambda \beta_1 f\left(\int \frac{g(x)}{\beta_1} \frac{dP_1}{dQ_1}(x)(g(x))^{-1} dQ_1(x)\right) \quad (247)$$

$$\begin{aligned} & - (1 - \lambda) \beta_2 f\left(\int \frac{g(x)}{\beta_2} \frac{dP_2}{dQ_2}(x)(g(x))^{-1} dQ_2(x)\right) \\ &= -\lambda \beta_1 f\left(\frac{1}{\beta_1} \int dP_1(x)\right) - (1 - \lambda) \beta_2 f\left(\frac{1}{\beta_2} \int dP_2(x)\right) \end{aligned} \quad (248)$$

$$= -\lambda \beta_1 f\left(\frac{1}{\beta_1}\right) - (1 - \lambda) \beta_2 f\left(\frac{1}{\beta_2}\right) \quad (249)$$

$$\leq -f\left(\lambda \beta_1 \frac{1}{\beta_1} + (1 - \lambda) \beta_2 \frac{1}{\beta_2}\right) \quad (250)$$

$$= -f(1) \quad (251)$$

$$= 0, \quad (252)$$

where the inequalities in (247) and (250) follow from Jensen's inequality [55, Section 6.3.5] and the fact that the function  $f$  in (244) is strictly concave. Note that from (245), in (247), for all  $i \in \{1, 2\}$ ,  $\int \frac{g(x)}{\beta_i} dQ_i(x) = 1$ ; while in (250),

$$\lambda \beta_1 + (1 - \lambda) \beta_2 = \int g(\nu) d(\lambda Q_1 + (1 - \lambda)Q_2)(\nu) \quad (253)$$

$$= \int d(\lambda P_1 + (1 - \lambda)P_2)(\nu) \quad (254)$$

$$= \lambda \int dP_1(\nu) + (1 - \lambda) \int dP_2(\nu) \quad (255)$$

$$= 1. \quad (256)$$

Given the strict convexity of the function  $f$  in (241), equality in (247) and (250) hold if and only if  $P_1 = P_2$  and  $Q_1 = Q_2$ . This completes the proof.

## B Proof of Lemma 3.1

The proof is divided into two parts. The first part is as follows. Under the assumption that the set  $\mathcal{K}_{Q,z}$  in (23) is empty, there is nothing to prove. Alternatively, under the assumption that the set  $\mathcal{K}_{Q,z}$  is not empty, there always exists a real  $b \in \mathcal{K}_{Q,z}$ , such that  $K_{Q,z}\left(-\frac{1}{b}\right) < +\infty$ . Note that for all  $\theta \in \mathcal{M}$ ,

$$\frac{d}{dt} \exp\left(-\frac{1}{t} \mathbf{L}_z(\theta)\right) = \frac{1}{t^2} \mathbf{L}_z(\theta) \exp\left(-\frac{1}{t} \mathbf{L}_z(\theta)\right) \geq 0, \quad (257)$$

with  $L_z$  in (3). Thus, from (22), it follows that  $K_{Q,z}(-\frac{1}{b})$  is nondecreasing with  $b$ . This implies that  $(0, b] \subseteq \mathcal{K}_{Q,z}$ .

Let  $b^* \in (0, +\infty]$  be

$$b^* = \sup \mathcal{K}_{Q,z}. \quad (258)$$

Hence, if  $b^* = +\infty$ , it follows from (23) that

$$\mathcal{K}_{Q,z} = (0, +\infty). \quad (259)$$

Alternatively, if  $b^* < +\infty$ , it holds that

$$(0, b^*) \subseteq \mathcal{K}_{Q,z} \subseteq (0, b^*]. \quad (260)$$

In either case, it follows that  $\mathcal{K}_{Q,z}$  is a convex set. This completes the first part of the proof.

The second part of the proof is under the assumption that  $Q$  is a probability measure. Under this assumption, for all  $\theta \in \mathcal{M}$  and for all  $t > 0$ , it follows that

$$\exp\left(-\frac{1}{t} L_z(\theta)\right) \leq 1, \quad (261)$$

with  $L_z$  in (3). Thus,

$$K_{Q,z}\left(-\frac{1}{t}\right) = \log\left(\int \exp\left(-\frac{1}{t} L_z(\theta)\right) dQ(\theta)\right) \quad (262)$$

$$\leq \log\left(\int dQ(\theta)\right) \quad (263)$$

$$= 0, \quad (264)$$

which implies that  $(0, +\infty) \subseteq \mathcal{K}_{Q,z}$ . Thus, if  $Q$  is a probability measure, from (23), it holds that  $\mathcal{K}_{Q,z} = (0, +\infty)$ , which completes the proof.

## C Proof of Theorem 3.1

The optimization problem in (19) can be re-written in terms of the Radon-Nikodym derivative of the optimization measure  $P$  with respect to the measure  $Q$ , denoted by  $\frac{dP}{dQ} : \mathcal{M} \rightarrow [0, \infty)$ , which yields:

$$\min_{P \in \Delta_Q(\mathcal{M}, \mathcal{B}(\mathcal{M}))} \int L_z(\theta) \frac{dP}{dQ}(\theta) dQ(\theta) + \lambda \int \frac{dP}{dQ}(\theta) \log\left(\frac{dP}{dQ}(\theta)\right) dQ(\theta) \quad (265a)$$

$$\text{s. t.} \quad \int \frac{dP}{dQ}(\theta) dQ(\theta) = 1. \quad (265b)$$

The remainder of the proof focuses on the problem in which the optimization is over the function  $\frac{dP}{dQ}$  instead of the measure  $P$ . This is due to the fact that for all  $P \in \Delta_Q(\mathcal{M})$ , the Radon-Nikodym derivate  $\frac{dP}{dQ}$  is unique up to sets of zero measure with respect to the measure  $Q$ . Let  $\mathcal{M}$  be the set of measurable

functions  $\mathcal{M} \rightarrow \mathbb{R}$  with respect to the measurable spaces  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$  and  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$  that are absolutely integrable with respect to  $Q$ . That is, for all  $\hat{g} \in \mathcal{M}$ , it holds that

$$\int |\hat{g}(\boldsymbol{\theta})| dQ(\boldsymbol{\theta}) < \infty. \quad (266)$$

Hence, the optimization problem of interest is:

$$\min_{g \in \mathcal{M}} \int \mathbf{L}_z(\boldsymbol{\theta}) g(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) + \lambda \int g(\boldsymbol{\theta}) \log(g(\boldsymbol{\theta})) dQ(\boldsymbol{\theta}) \quad (267a)$$

$$\text{s. t. } \int g(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) = 1. \quad (267b)$$

Let the Lagrangian of the optimization problem in (267) be the functional  $L : \mathcal{M} \times \mathbb{R} \rightarrow \mathbb{R}$  such that

$$\begin{aligned} L(g, \beta) = & \int \mathbf{L}_z(\boldsymbol{\nu}) g(\boldsymbol{\nu}) dQ(\boldsymbol{\nu}) + \lambda \int g(\boldsymbol{\nu}) \log(g(\boldsymbol{\nu})) dQ(\boldsymbol{\nu}) \\ & + \beta \left( \int g(\boldsymbol{\nu}) dQ(\boldsymbol{\nu}) - 1 \right), \end{aligned} \quad (268)$$

where  $\beta$  is a real that acts as a Lagrangian multiplier due to the constraint (267b). Let  $\hat{g} : \mathcal{M} \rightarrow \mathbb{R}$  be a function in  $\mathcal{M}$ . The Gateaux differential of the functional  $L$  in (268) at  $(g, \beta) \in \mathcal{M} \times \mathbb{R}$  in the direction of  $\hat{g}$ , if it exists, is

$$\partial L(g, \beta; \hat{g}) \triangleq \left. \frac{d}{d\gamma} L(g + \gamma \hat{g}, \beta) \right|_{\gamma=0}. \quad (269)$$

The proof continues under the assumption that the functions  $g$  and  $\hat{g}$  are such that the Gateaux differential in (269) exists. Under such an assumption, let the function  $r : \mathbb{R} \rightarrow \mathbb{R}$  satisfy for all  $\alpha \in (-\epsilon, \epsilon)$ , with  $\epsilon$  arbitrarily small, that

$$\begin{aligned} r(\alpha) = & \int \mathbf{L}_z(\boldsymbol{\nu}) (g(\boldsymbol{\nu}) + \alpha \hat{g}(\boldsymbol{\nu})) dQ(\boldsymbol{\nu}) \\ & + \beta \left( \int (g(\boldsymbol{\nu}) + \alpha \hat{g}(\boldsymbol{\nu})) dQ(\boldsymbol{\nu}) - 1 \right) \\ & + \lambda \int (\hat{g}(\boldsymbol{\nu}) + \alpha \hat{g}(\boldsymbol{\nu})) \log(g(\boldsymbol{\nu}) + \alpha \hat{g}(\boldsymbol{\nu})) dQ(\boldsymbol{\nu}) \end{aligned} \quad (270)$$

$$\begin{aligned} = & \int g(\boldsymbol{\nu}) (\mathbf{L}_z(\boldsymbol{\nu}) + \beta) dQ(\boldsymbol{\nu}) - \beta \\ & + \alpha \left( \int \hat{g}(\boldsymbol{\nu}) (\mathbf{L}_z(\boldsymbol{\nu}) + \beta) dQ(\boldsymbol{\nu}) \right) \\ & + \lambda \int (g(\boldsymbol{\nu}) + \alpha \hat{g}(\boldsymbol{\nu})) \log(g(\boldsymbol{\nu}) + \alpha \hat{g}(\boldsymbol{\nu})) dQ(\boldsymbol{\nu}), \end{aligned} \quad (271)$$

where the last equality is simply an algebraic re-arrangement of terms. From the assumption that the functions  $g$  and  $\hat{g}$  are such that the Gateaux differential

in (269) exists, it follows that the function  $r$  in (271) is differentiable at zero. Note that the first two terms in (271) are independent of  $\alpha$ ; the third term is linear with  $\alpha$ ; and the fourth term can be written using the function  $\hat{r} : \mathbb{R} \rightarrow \mathbb{R}$  such that for all  $\alpha \in (-\epsilon, \epsilon)$ , with  $\epsilon$  arbitrarily small, satisfies

$$\hat{r}(\alpha) = \lambda \int (g(\boldsymbol{\nu}) + \alpha \hat{g}(\boldsymbol{\nu})) \log(g(\boldsymbol{\nu}) + \alpha \hat{g}(\boldsymbol{\nu})) dQ(\boldsymbol{\nu}) \quad (272)$$

$$= \lambda \int f(g(\boldsymbol{\nu}) + \alpha \hat{g}(\boldsymbol{\nu})) dQ(\boldsymbol{\nu}), \quad (273)$$

where  $f : (0, +\infty) \rightarrow \mathbb{R}$  is such that  $f(t) = t \log(t)$ . Under the same assumption, it follows that the function  $\hat{r}$  in (272) is differentiable at zero. That is, the limit

$$\lim_{\delta \rightarrow 0} \frac{1}{\delta} (\hat{r}(\gamma + \delta) - \hat{r}(\gamma)) \quad (274)$$

exists for all  $\gamma \in (-\epsilon, \epsilon)$ , with  $\epsilon$  arbitrarily small. Note that the function  $f$  in (273) is continuous and differentiable (with finite derivate) in  $(0, +\infty)$ . Thus, the function  $f$  is also Lipschitz continuous. This implies that for all  $\boldsymbol{\theta} \in \text{supp } Q$ , and for all  $\gamma \in (-\epsilon, \epsilon)$ , with  $\epsilon > 0$  arbitrarily small, it holds that

$$|f(g(\boldsymbol{\theta}) + (\gamma + \delta)\hat{g}(\boldsymbol{\theta})) - f(g(\boldsymbol{\theta}) + \gamma\hat{g}(\boldsymbol{\theta}))| \leq c|\hat{g}(\boldsymbol{\theta})||\delta|, \quad (275)$$

with  $\delta > 0$ , for some constant  $c$  positive and finite. This implies that

$$\left| \frac{f(g(\boldsymbol{\theta}) + (\gamma + \delta)\hat{g}(\boldsymbol{\theta})) - f(g(\boldsymbol{\theta}) + \gamma\hat{g}(\boldsymbol{\theta}))}{\delta} \right| \leq c|\hat{g}(\boldsymbol{\theta})|. \quad (276)$$

Using these arguments, the limit in (274) satisfies for all  $\gamma \in (-\epsilon, \epsilon)$ , with  $\epsilon > 0$  arbitrarily small, that

$$\begin{aligned} & \lim_{\delta \rightarrow 0} \frac{1}{\delta} (\hat{r}(\gamma + \delta) - \hat{r}(\gamma)) \\ &= \lambda \lim_{\delta \rightarrow 0} \int \frac{f(g(\boldsymbol{\theta}) + (\gamma + \delta)\hat{g}(\boldsymbol{\theta})) - f(g(\boldsymbol{\theta}) + \gamma\hat{g}(\boldsymbol{\theta}))}{\delta} dQ(\boldsymbol{\theta}) \\ &= \lambda \int \dot{f}(g(\boldsymbol{\theta}) + \gamma\hat{g}(\boldsymbol{\theta})) \hat{g}(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) \end{aligned} \quad (277)$$

$$< \infty, \quad (278)$$

where the function  $\dot{f} : (0, +\infty) \rightarrow \mathbb{R}$  is the derivative of  $f$ . That is,  $\dot{f}(t) = 1 + \log(t)$ . The equality in (277) and the inequality in (278) follow from noticing that the conditions for the dominated convergence theorem hold [55, Theorem 1.6.9], namely:

- For all  $\gamma \in (-\epsilon, \epsilon)$ , with  $\epsilon > 0$ , the inequality in (276) holds;
- The function  $\hat{g}$  in (276) satisfies the inequality in (266); and

- For all  $\boldsymbol{\theta} \in \text{supp } Q$  and for all  $\gamma \in (-\epsilon, \epsilon)$ , with  $\epsilon > 0$  arbitrarily small, it holds that

$$\lim_{\delta \rightarrow 0} \frac{f(g(\boldsymbol{\theta}) + (\gamma + \delta)\hat{g}(\boldsymbol{\theta})) - f(g(\boldsymbol{\theta}) + \gamma\hat{g}(\boldsymbol{\theta}))}{\delta} = \frac{d}{d\gamma} f(g(\boldsymbol{\theta}) + \gamma\hat{g}(\boldsymbol{\theta})) \quad (279)$$

$$= \dot{f}(g(\boldsymbol{\theta}) + \gamma\hat{g}(\boldsymbol{\theta}))\hat{g}(\boldsymbol{\theta}). \quad (280)$$

Hence, the derivative of the real function  $r$  in (271) is

$$\begin{aligned} \frac{d}{d\alpha} r(\alpha) &= \int \mathbf{L}_{\mathbf{z}}(\boldsymbol{\nu}) \hat{g}(\boldsymbol{\nu}) dQ(\boldsymbol{\nu}) + \beta \int \hat{g}(\boldsymbol{\nu}) dQ(\boldsymbol{\nu}) \\ &\quad + \lambda \int \hat{g}(\boldsymbol{\nu}) (1 + \log(g(\boldsymbol{\nu}) + \alpha\hat{g}(\boldsymbol{\nu}))) dQ(\boldsymbol{\nu}). \end{aligned} \quad (281)$$

From (269) and (281), it follows that

$$\partial L(g, \beta; \hat{g}) = \int \hat{g}(\boldsymbol{\nu}) (\mathbf{L}_{\mathbf{z}}(\boldsymbol{\nu}) + \lambda(1 + \log(g(\boldsymbol{\nu}))) + \beta) dQ(\boldsymbol{\nu}). \quad (282)$$

The relevance of the Gateaux differential in (282) stems from [69, Theorem 1, page 178], which unveils the fact that a necessary condition for the functional  $L$  in (268) to have a minimum at  $\left(\frac{dP_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}}{dQ}, \beta\right) \in \mathcal{M} \times \mathbb{R}$  is that for all functions  $\hat{g} \in \mathcal{M}$ ,

$$\partial L\left(\frac{dP_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}}{dQ}, \beta; \hat{g}\right) = 0. \quad (283)$$

From (283), it follows that  $\frac{dP_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}}{dQ}$  must satisfy for all functions  $\hat{g}$  in  $\mathcal{M}$  that

$$0 = \int \hat{g}(\boldsymbol{\nu}) \left( \mathbf{L}_{\mathbf{z}}(\boldsymbol{\nu}) + \lambda \left( 1 + \log \left( \frac{dP_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}}{dQ}(\boldsymbol{\nu}) \right) \right) + \beta \right) dQ(\boldsymbol{\nu}),$$

which implies that for all  $\boldsymbol{\nu} \in \text{supp } Q$ ,

$$\mathbf{L}_{\mathbf{z}}(\boldsymbol{\nu}) + \lambda \left( 1 + \log \left( \frac{dP_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}}{dQ}(\boldsymbol{\nu}) \right) \right) + \beta = 0, \quad (284)$$

and thus,

$$\frac{dP_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}}{dQ}(\boldsymbol{\nu}) = \exp\left(-\frac{\beta + \lambda}{\lambda}\right) \exp\left(-\frac{\mathbf{L}_{\mathbf{z}}(\boldsymbol{\nu})}{\lambda}\right), \quad (285)$$

with  $\beta$  chosen to satisfy (265b). That is,

$$\frac{dP_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}}{dQ}(\boldsymbol{\nu}) = \frac{\exp\left(-\frac{\mathbf{L}_{\mathbf{z}}(\boldsymbol{\nu})}{\lambda}\right)}{\int \exp\left(-\frac{\mathbf{L}_{\mathbf{z}}(\boldsymbol{\theta})}{\lambda}\right) dQ(\boldsymbol{\theta})} \quad (286)$$

$$= \exp\left(-K_{Q,\mathbf{z}}\left(-\frac{1}{\lambda}\right) - \frac{1}{\lambda}\mathbf{L}_{\mathbf{z}}(\boldsymbol{\nu})\right). \quad (287)$$



The proof continues by verifying that the measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  that satisfies (286) is the unique solution to the ERM-RER problem in (19). Such verification is done by showing that the objective function in (19) is strictly convex with the optimization variable. Let  $P_1$  and  $P_2$  be two different probability measures in  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$  and let  $\alpha$  be in  $(0, 1)$ . Hence,

$$\begin{aligned} & R_z(\alpha P_1 + (1 - \alpha)P_2) + \lambda D(\alpha P_1 + (1 - \alpha)P_2 \| Q) \\ &= \alpha R_z(P_1) + (1 - \alpha)R_z(P_2) + \lambda D(\alpha P_1 + (1 - \alpha)P_2 \| Q) \\ &> \alpha(R_z(P_1) + \lambda D(P_1 \| Q)) + (1 - \alpha)(R_z(P_2) + \lambda D(P_2 \| Q)) \end{aligned}$$

where the functional  $R_z$  is defined in (18). The equality above follows from the properties of the Lebesgue integral, while the inequality follows from Theorem 2.2. This proves that the solution is unique due to the strict concavity of the objective function, which completes the proof.

## D Proof of Lemma 3.2

From Theorem 3.1, it follows that for all  $\theta \in \text{supp } Q$ ,

$$\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) = \exp\left(-K_{Q,z}\left(-\frac{1}{\lambda}\right) - \frac{1}{\lambda}L_z(\theta)\right) \quad (288)$$

$$\leq \exp\left(-K_{Q,z}\left(-\frac{1}{\lambda}\right)\right) \quad (289)$$

$$< +\infty, \quad (290)$$

where the inequality in (289) follows from the fact that the function  $L_z$  is non-negative; and the equality in (290) follows from the fact that  $\lambda \in \mathcal{K}_{Q,z}$ . This completes the proof of finiteness.

The proof of positivity follows from observing that  $\lambda \in \mathcal{K}_{Q,z}$  and thus,  $K_{Q,z}\left(-\frac{1}{\lambda}\right) < +\infty$ , and thus,  $\exp\left(-K_{Q,z}\left(-\frac{1}{\lambda}\right)\right) > 0$ . Moreover, for all  $\theta \in \text{supp } Q$ , it holds that  $L_z(\theta) \leq +\infty$ , which implies that  $-\frac{1}{\lambda}L_z(\theta) \geq -\infty$ , and thus,  $\exp\left(-\frac{1}{\lambda}L_z(\theta)\right) \geq 0$ , with equality if and only if  $L_z(\theta) = +\infty$ . These two observations put together yield

$$\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) = \exp\left(-K_{Q,z}\left(-\frac{1}{\lambda}\right) - \frac{1}{\lambda}L_z(\theta)\right) \quad (291)$$

$$= \exp\left(-K_{Q,z}\left(-\frac{1}{\lambda}\right)\right) \exp\left(-\frac{1}{\lambda}L_z(\theta)\right) \quad (292)$$

$$\geq 0, \quad (293)$$

with equality if and only if  $L_z(\theta) = +\infty$ . This completes the proof.

## E Proof of Lemma 3.3

The probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  in (25) satisfies for all  $\mathcal{C} \in \mathcal{B}(\mathcal{M})$ ,

$$P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{C}) = \int_{\mathcal{C}} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}), \quad (294)$$

and thus, if  $Q(\mathcal{C}) = 0$ , then

$$P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{C}) = 0, \quad (295)$$

which implies the absolute continuity of  $P_{\Theta|Z=z}^{(Q,\lambda)}$  with respect to  $Q$ .

Alternatively, given a set  $\mathcal{C} \in \mathcal{B}(\mathcal{M})$ , assume now that  $P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{C}) = 0$ . Hence, it follows that

$$0 = P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{C}) \quad (296)$$

$$= \int_{\mathcal{C}} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}). \quad (297)$$

From Lemma 3.2, and the assumption  $Q(\{\boldsymbol{\theta} \in \mathcal{M} : \mathcal{L}_z(\boldsymbol{\theta}) = +\infty\}) = 0$ , it holds that for all  $\boldsymbol{\theta} \in \text{supp } Q$ ,

$$\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\boldsymbol{\theta}) > 0, \quad (298)$$

which implies that

$$\int_{\mathcal{C}} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) = 0, \quad (299)$$

if and only if  $Q(\mathcal{C}) = 0$ . This verifies the absolute continuity of  $Q$  with respect to  $P_{\Theta|Z=z}^{(Q,\lambda)}$ , and completes the proof.

## F Proof of Lemma 3.4

Consider the function  $g : \mathcal{M} \rightarrow [0, +\infty)$ ,

$$g(\boldsymbol{\theta}) = \frac{dP_{\Theta|Z=z}^{(Q,\alpha)}}{dQ}(\boldsymbol{\theta}) \left( \frac{dP_{\Theta|Z=z}^{(Q,\beta)}}{dQ}(\boldsymbol{\theta}) \right)^{-1}, \quad (300)$$

and note that for all  $\boldsymbol{\theta} \in \text{supp } Q \setminus \{\boldsymbol{\nu} \in \mathcal{M} : \mathcal{L}_z(\boldsymbol{\nu}) = +\infty\}$ ,  $g(\boldsymbol{\theta}) > 0$ . Alternatively, for all  $\boldsymbol{\theta} \in \{\boldsymbol{\nu} \in \mathcal{M} : \mathcal{L}_z(\boldsymbol{\nu}) = +\infty\}$ ,  $g(\boldsymbol{\theta}) = 0$ , which follows from the assumption  $0 \cdot \frac{1}{0} = 0$ .

Consider a measure  $P$  on  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ , such that for all sets  $\mathcal{A} \in \mathcal{B}(\mathcal{M})$ ,

$$P(\mathcal{A}) = \int_{\mathcal{A}} g(\boldsymbol{\theta}) dP_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\beta)}(\boldsymbol{\theta}), \quad (301)$$

and note that if  $P_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\beta)}(\mathcal{A}) = 0$ , then  $P(\mathcal{A}) = 0$ . This implies that  $P$  is absolutely continuous with respect to  $P_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\beta)}(\mathcal{A})$ . Moreover, from (301), it follows that

$$P(\mathcal{A}) = \int_{\mathcal{A}} \frac{dP_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\alpha)}(\boldsymbol{\theta})}{dQ} \left( \frac{dP_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\beta)}(\boldsymbol{\theta})}{dQ} \right)^{-1} dP_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\beta)}(\boldsymbol{\theta}) \quad (302)$$

$$= \int_{\mathcal{A}} \frac{dP_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\alpha)}(\boldsymbol{\theta})}{dQ} \left( \frac{dP_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\beta)}(\boldsymbol{\theta})}{dQ} \right)^{-1} \frac{dP_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\beta)}(\boldsymbol{\theta})}{dQ} dQ(\boldsymbol{\theta}) \quad (303)$$

$$= \int_{\mathcal{A}} \frac{dP_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\alpha)}(\boldsymbol{\theta})}{dQ} dQ(\boldsymbol{\theta}) \quad (304)$$

$$= \int_{\mathcal{A}} dP_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\alpha)}(\boldsymbol{\theta}) \quad (305)$$

$$= P_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\alpha)}(\mathcal{A}), \quad (306)$$

which implies that the probability measures  $P$  in (301) and  $P_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\alpha)}$  are identical. Thus,  $P_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\alpha)}$  is absolutely continuous with respect to  $P_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\beta)}$ . The proof that  $P_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\beta)}$  is absolutely continuous with respect to  $P_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\alpha)}$  follows the same argument. This completes the proof.

## G Proof of Lemma 3.6

From Theorem 3.1, the probability measure  $P_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}$  in (25) satisfies for all  $\boldsymbol{\theta} \in \text{supp } Q$ ,

$$\frac{dP_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}(\boldsymbol{\theta})}{dQ} = \frac{\exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\theta})}{\lambda}\right)}{\int \exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\nu})}{\lambda}\right) dQ(\boldsymbol{\nu})} \quad (307)$$

$$= \left( \exp\left(\frac{\mathbf{L}_z(\boldsymbol{\theta})}{\lambda}\right) \int \exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\nu})}{\lambda}\right) dQ(\boldsymbol{\nu}) \right)^{-1} \quad (308)$$

$$= \left( \int \exp\left(\frac{1}{\lambda}(\mathbf{L}_z(\boldsymbol{\theta}) - \mathbf{L}_z(\boldsymbol{\nu}))\right) dQ(\boldsymbol{\nu}) \right)^{-1}. \quad (309)$$

Given  $\boldsymbol{\theta} \in \text{supp } Q$ , consider the partition of  $\text{supp } Q$  formed by the sets  $\mathcal{A}_0(\boldsymbol{\theta})$ ,  $\mathcal{A}_1(\boldsymbol{\theta})$ , and  $\mathcal{A}_2(\boldsymbol{\theta})$ , which satisfy the following:

$$\mathcal{A}_0(\boldsymbol{\theta}) \triangleq \{\boldsymbol{\nu} \in \text{supp } Q : \mathbf{L}_z(\boldsymbol{\theta}) - \mathbf{L}_z(\boldsymbol{\nu}) = 0\}, \quad (310a)$$

$$\mathcal{A}_1(\boldsymbol{\theta}) \triangleq \{\boldsymbol{\nu} \in \text{supp } Q : \mathbf{L}_z(\boldsymbol{\theta}) - \mathbf{L}_z(\boldsymbol{\nu}) < 0\}, \text{ and} \quad (310b)$$

$$\mathcal{A}_2(\boldsymbol{\theta}) \triangleq \{\boldsymbol{\nu} \in \text{supp } Q : \mathbf{L}_z(\boldsymbol{\theta}) - \mathbf{L}_z(\boldsymbol{\nu}) > 0\}. \quad (310c)$$

Using the sets  $\mathcal{A}_0(\boldsymbol{\theta})$ ,  $\mathcal{A}_1(\boldsymbol{\theta})$ , and  $\mathcal{A}_2(\boldsymbol{\theta})$  in (309), the following holds for all  $\boldsymbol{\theta} \in \text{supp } Q$ ,

$$\begin{aligned} \frac{dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}(\boldsymbol{\theta})}{dQ} &= \left( \int_{\mathcal{A}_0(\boldsymbol{\theta})} \exp\left(\frac{1}{\lambda}(\mathbf{L}_z(\boldsymbol{\theta}) - \mathbf{L}_z(\boldsymbol{\nu}))\right) dQ(\boldsymbol{\nu}) \right. \\ &\quad + \int_{\mathcal{A}_1(\boldsymbol{\theta})} \exp\left(\frac{1}{\lambda}(\mathbf{L}_z(\boldsymbol{\theta}) - \mathbf{L}_z(\boldsymbol{\nu}))\right) dQ(\boldsymbol{\nu}) \\ &\quad \left. + \int_{\mathcal{A}_2(\boldsymbol{\theta})} \exp\left(\frac{1}{\lambda}(\mathbf{L}_z(\boldsymbol{\theta}) - \mathbf{L}_z(\boldsymbol{\nu}))\right) dQ(\boldsymbol{\nu}) \right)^{-1} \end{aligned} \quad (311)$$

$$\begin{aligned} &= \left( Q(\mathcal{A}_0(\boldsymbol{\theta})) + \int_{\mathcal{A}_1(\boldsymbol{\theta})} \exp\left(\frac{1}{\lambda}(\mathbf{L}_z(\boldsymbol{\theta}) - \mathbf{L}_z(\boldsymbol{\nu}))\right) dQ(\boldsymbol{\nu}) \right. \\ &\quad \left. + \int_{\mathcal{A}_2(\boldsymbol{\theta})} \exp\left(\frac{1}{\lambda}(\mathbf{L}_z(\boldsymbol{\theta}) - \mathbf{L}_z(\boldsymbol{\nu}))\right) dQ(\boldsymbol{\nu}) \right)^{-1}. \end{aligned} \quad (312)$$

Note that the sets

$$\{\boldsymbol{\nu} \in \text{supp } Q : \mathbf{L}_z(\boldsymbol{\nu}) = \delta_{Q,z}^*\}, \quad (313)$$

$$\{\boldsymbol{\nu} \in \text{supp } Q : \mathbf{L}_z(\boldsymbol{\nu}) > \delta_{Q,z}^*\}, \text{ and} \quad (314)$$

$$\{\boldsymbol{\nu} \in \text{supp } Q : \mathbf{L}_z(\boldsymbol{\nu}) < \delta_{Q,z}^*\}, \quad (315)$$

with  $\delta_{Q,z}^*$  in (38), form a partition of the set  $\text{supp } Q$ . Following this observation, the rest of the proof is divided into three parts. The first part evaluates  $\lim_{\lambda \rightarrow 0^+} \frac{dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}(\boldsymbol{\theta})}{dQ}$ , with  $\boldsymbol{\theta} \in \{\boldsymbol{\nu} \in \mathcal{M} : \mathbf{L}_z(\boldsymbol{\nu}) = \delta_{Q,z}^*\}$ . The second part considers the case in which  $\boldsymbol{\theta} \in \{\boldsymbol{\nu} \in \mathcal{M} : \mathbf{L}_z(\boldsymbol{\nu}) > \delta_{Q,z}^*\}$ . The third part considers the remaining case.

The first part is as follows. Consider that  $\boldsymbol{\theta} \in \{\boldsymbol{\nu} \in \mathcal{M} : \mathbf{L}_z(\boldsymbol{\nu}) = \delta_{Q,z}^*\}$  and note that  $\{\boldsymbol{\nu} \in \mathcal{M} : \mathbf{L}_z(\boldsymbol{\nu}) = \delta_{Q,z}^*\} = \mathcal{L}_{Q,z}^*$ . Hence, the sets  $\mathcal{A}_0(\boldsymbol{\theta})$ ,  $\mathcal{A}_1(\boldsymbol{\theta})$ , and  $\mathcal{A}_2(\boldsymbol{\theta})$  in (310) satisfy the following:

$$\mathcal{A}_0(\boldsymbol{\theta}) = \mathcal{L}_{Q,z}^*, \quad (316a)$$

$$\mathcal{A}_1(\boldsymbol{\theta}) = \{\boldsymbol{\mu} \in \text{supp } Q : \mathbf{L}_z(\boldsymbol{\mu}) > \delta_{Q,z}^*\}, \text{ and} \quad (316b)$$

$$\mathcal{A}_2(\boldsymbol{\theta}) = \{\boldsymbol{\mu} \in \text{supp } Q : \mathbf{L}_z(\boldsymbol{\mu}) < \delta_{Q,z}^*\}. \quad (316c)$$

From the definition of  $\delta_{Q,z}^*$  in (38), it follows that  $Q(\mathcal{A}_2(\boldsymbol{\theta})) = 0$ . Plugging the equalities in (316) in (312) yields for all  $\boldsymbol{\theta} \in \{\boldsymbol{\nu} \in \mathcal{M} : \mathbf{L}_z(\boldsymbol{\nu}) = \delta_{Q,z}^*\}$ ,

$$\frac{dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}(\boldsymbol{\theta})}{dQ} = \left( Q(\mathcal{L}_{Q,z}^*) + \int_{\mathcal{A}_1(\boldsymbol{\theta})} \exp\left(\frac{1}{\lambda}(\mathbf{L}_z(\boldsymbol{\theta}) - \mathbf{L}_z(\boldsymbol{\nu}))\right) dQ(\boldsymbol{\nu}) \right)^{-1}. \quad (317)$$

The equality in (317) implies that for all  $\theta \in \{\nu \in \mathcal{M} : \mathbb{L}_z(\nu) = \delta_{Q,z}^*\}$ ,

$$\lim_{\lambda \rightarrow 0^+} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) = \left( \lim_{\lambda \rightarrow 0^+} \int_{\mathcal{A}_1(\theta)} \exp\left(\frac{1}{\lambda}(\mathbb{L}_z(\theta) - \mathbb{L}_z(\nu))\right) dQ(\nu) + \right. \quad (318)$$

$$\left. + Q(\mathcal{L}_{Q,z}^*) \right)^{-1} \quad (319)$$

$$= \begin{cases} +\infty & \text{if } Q(\mathcal{L}_{Q,z}^*) = 0 \\ \frac{1}{Q(\mathcal{L}_{Q,z}^*)} & \text{otherwise.} \end{cases} \quad (320)$$

where the equality in (320) follows from verifying that the dominated convergence theorem [55, Theorem 2.6.9] holds. That is,

- (a) For all  $\nu \in \mathcal{A}_1(\theta)$ , it holds that  $\exp\left(\frac{1}{\lambda}(\mathbb{L}_z(\theta) - \mathbb{L}_z(\nu))\right) < 1$ ; and
- (b) For all  $\nu \in \mathcal{A}_1(\theta)$ , it holds that

$$\lim_{\lambda \rightarrow 0^+} \exp\left(\frac{1}{\lambda}(\mathbb{L}_z(\theta) - \mathbb{L}_z(\nu))\right) = 0. \quad (321)$$

This completes the first part of the proof.

The second part is as follows. For all  $\delta > \delta_{Q,z}^*$  and for all  $\theta \in \{\nu \in \text{supp } Q : \mathbb{L}_z(\nu) = \delta\}$ , the sets  $\mathcal{A}_0(\theta)$ ,  $\mathcal{A}_1(\theta)$ , and  $\mathcal{A}_2(\theta)$  in (310) satisfy the following:

$$\mathcal{A}_0(\theta) = \{\mu \in \text{supp } Q : \mathbb{L}_z(\mu) = \delta\}, \quad (322a)$$

$$\mathcal{A}_1(\theta) = \{\mu \in \text{supp } Q : \mathbb{L}_z(\mu) > \delta\}, \text{ and} \quad (322b)$$

$$\mathcal{A}_2(\theta) = \{\mu \in \text{supp } Q : \mathbb{L}_z(\mu) < \delta\}. \quad (322c)$$

Consider the sets

$$\mathcal{A}_{2,1}(\theta) \triangleq \{\mu \in \mathcal{A}_2(\theta) : \mathbb{L}_z(\mu) < \delta_{Q,z}^*\}, \text{ and} \quad (323)$$

$$\mathcal{A}_{2,2}(\theta) \triangleq \{\mu \in \mathcal{A}_2(\theta) : \delta_{Q,z}^* \leq \mathbb{L}_z(\mu) < \delta\}, \quad (324)$$

and note that  $\mathcal{A}_{2,1}(\theta)$  and  $\mathcal{A}_{2,2}(\theta)$  form a partition of  $\mathcal{A}_2(\theta)$ . Moreover, from the definition of  $\delta_{Q,z}^*$  in (38), it holds that

$$Q(\mathcal{A}_{2,1}(\theta)) = 0. \quad (325)$$

Hence, plugging the equalities in (322) and (325) in (312) yields, for all  $\delta > \delta_{Q,z}^*$  and for all  $\theta \in \{\nu \in \mathcal{M} : \mathbb{L}_z(\nu) = \delta\}$ ,

$$\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) = \left( Q(\mathcal{A}_0(\theta)) + \int_{\mathcal{A}_1(\theta)} \exp\left(\frac{1}{\lambda}(\mathbb{L}_z(\theta) - \mathbb{L}_z(\nu))\right) dQ(\nu) \right. \quad (326)$$

$$\left. + \int_{\mathcal{A}_{2,2}(\theta)} \exp\left(\frac{1}{\lambda}(\mathbb{L}_z(\theta) - \mathbb{L}_z(\nu))\right) dQ(\nu) \right)^{-1}. \quad (327)$$

The equality in (327) implies that for all  $\delta > \delta_{Q,z}^*$  and for all  $\theta \in \{\nu \in \mathcal{M} : \mathbb{L}_z(\nu) = \delta\}$ ,

$$\lim_{\lambda \rightarrow 0^+} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) = \left( \lim_{\lambda \rightarrow 0^+} \int_{\mathcal{A}_1(\theta)} \exp\left(\frac{1}{\lambda}(\mathbb{L}_z(\theta) - \mathbb{L}_z(\nu))\right) dQ(\nu) \right. \quad (328)$$

$$\left. + \lim_{\lambda \rightarrow 0^+} \int_{\mathcal{A}_{2,2}(\theta)} \exp\left(\frac{1}{\lambda}(\mathbb{L}_z(\theta) - \mathbb{L}_z(\nu))\right) dQ(\nu) \right. \quad (329)$$

$$\left. + Q(\mathcal{A}_0(\theta)) \right)^{-1} \quad (330)$$

$$= \left( \lim_{\lambda \rightarrow 0^+} \int_{\mathcal{A}_{2,2}(\theta)} \exp\left(\frac{1}{\lambda}(\mathbb{L}_z(\theta) - \mathbb{L}_z(\nu))\right) dQ(\nu) \right. \quad (331)$$

$$\left. + Q(\mathcal{A}_0(\theta)) \right)^{-1} \quad (332)$$

$$= 0, \quad (333)$$

where the equality in (331) follows by verifying that the dominated convergence theorem [55, Theorem 2.6.9] holds. That is,

(a) For all  $\nu \in \mathcal{A}_1(\theta)$ , it holds that  $\exp\left(\frac{1}{\lambda}(\mathbb{L}_z(\theta) - \mathbb{L}_z(\nu))\right) < 1$ ; and

(b) For all  $\nu \in \mathcal{A}_1(\theta)$ , it holds that

$$\lim_{\lambda \rightarrow 0^+} \exp\left(\frac{1}{\lambda}(\mathbb{L}_z(\theta) - \mathbb{L}_z(\nu))\right) = 0. \quad (334)$$

This completes the second part.

The third part of the proof follows by noticing that the set  $\{\nu \in \text{supp } Q : \mathbb{L}_z(\nu) < \delta_{Q,z}^*\}$  is a negligible set with respect to  $Q$  and thus, for all  $\theta \in \{\nu \in \text{supp } Q :$

$\mathbb{L}_z(\nu) < \delta_{Q,z}^*\}$ , the value  $\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta)$  is immaterial. Hence, it is arbitrarily assumed that for all  $\theta \in \{\nu \in \text{supp } Q : \mathbb{L}_z(\nu) < \delta_{Q,z}^*\}$ , it holds that

$$\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) = 0. \quad (335)$$

This completes the third part and completes the proof.

## H Proof of Lemma 3.7

Consider the following partition of the set  $\mathcal{M}$  formed by the sets

$$\mathcal{A}_0 \triangleq \{\theta \in \mathcal{M} : \mathbb{L}_z(\theta) = \delta_{Q,z}^*\}, \quad (336a)$$

$$\mathcal{A}_1 \triangleq \{\theta \in \mathcal{M} : \mathbb{L}_z(\theta) < \delta_{Q,z}^*\}, \text{ and} \quad (336b)$$

$$\mathcal{A}_2 \triangleq \{\theta \in \mathcal{M} : \mathbb{L}_z(\theta) > \delta_{Q,z}^*\}, \quad (336c)$$

with  $\delta_{Q,z}^*$  in (38) and the function  $L_z$  in (3). Note that  $\mathcal{A}_0 = \mathcal{L}_{Q,z}^*$ , with  $\mathcal{L}_{Q,z}^*$  in (39) and

$$1 = P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{A}_0) + P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{A}_1) + P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{A}_2) \quad (337)$$

$$= P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{A}_0) + P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{A}_2) \quad (338)$$

$$= P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{A}_0) + \int_{\mathcal{A}_2} dP_{\Theta|Z=z}^{(Q,\lambda)}(\theta), \quad (339)$$

where, the equality in (338) follows from noticing that  $P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{A}_1) = 0$ , which follows from the definition of  $\delta_{Q,z}^*$  in (38) and the fact that the probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  is absolutely continuous with respect to the measure  $Q$ .

The above implies that

$$1 = \lim_{\lambda \rightarrow 0^+} P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{A}_0) + \lim_{\lambda \rightarrow 0^+} \int_{\mathcal{A}_2} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) dQ(\theta) \quad (340)$$

$$= \lim_{\lambda \rightarrow 0^+} P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{A}_0) + \int_{\mathcal{A}_2} \lim_{\lambda \rightarrow 0^+} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) dQ(\theta) \quad (341)$$

$$= \lim_{\lambda \rightarrow 0^+} P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{A}_0), \quad (342)$$

where, the equality in (341) follows from the dominated convergence theorem [55, Theorem 1.6.9], given that the Radon-Nikodym derivative  $\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}$  is positive and finite (Lemma 3.2); and the inequality in (342) holds from the fact that for all  $\theta \in \mathcal{A}_2$ , it holds that  $\lim_{\lambda \rightarrow 0^+} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) = 0$  (Lemma 3.6). Hence, it finally holds that

$$\lim_{\lambda \rightarrow 0^+} P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_{Q,z}^*) = 1, \quad (343)$$

which completes the proof.

## I Proof of Lemma 4.1

The proof is presented in two parts. The first part shows that if for all  $\delta \in (\rho^*, +\infty)$ , the inequality in (50) holds, then,  $Q$  is coherent. The second part shows that if  $Q$  is not coherent, then there exists a  $\delta \in (\rho^*, +\infty)$  such that

$$P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_z(\delta)) = 0. \quad (344)$$

The first part is as follows. Note that for all  $\delta \in (\rho^*, +\infty)$  and for all  $\theta \in \mathcal{L}_z(\delta) \cap \text{supp } Q$ , it holds from Lemma 3.2 that

$$\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) > 0. \quad (345)$$

Hence, if for all  $\delta \in (\rho^*, +\infty)$ , the inequality in (50) holds, then

$$0 < P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_z(\delta)) \quad (346)$$

$$= \int_{\mathcal{L}_z(\delta)} dP_{\Theta|Z=z}^{(Q,\lambda)}(\theta) \quad (347)$$

$$= \int_{\mathcal{L}_z(\delta)} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) dQ(\theta), \quad (348)$$

which, together with (345), implies that for all  $\delta \in (\rho^*, +\infty)$ ,  $Q(\mathcal{L}_z(\delta)) > 0$ . Hence,  $Q$  is coherent.

The second part is as follows. Assume that  $Q$  is not coherent. Then, there exists a  $\delta \in (\rho^*, +\infty)$  such that  $Q(\mathcal{L}_z(\delta)) = 0$ . Hence, from the fact that  $P_{\Theta|Z=z}^{(Q,\lambda)}$  is absolutely continuous with respect to  $Q$ , it follows that  $P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_z(\delta)) = 0$ . This completes the proof.

## J Proof of Theorem 4.1

The optimization problem in (83) can be re-written in terms of the Radon-Nikodym derivative of the optimization measure  $P$  with respect to the measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$ , denoted by  $\frac{dP}{dP_{\Theta|Z=z}^{(Q,\lambda)}} : \mathcal{M} \rightarrow [0, +\infty)$ , which yields:

$$\min_{P \in \Delta_Q(\mathcal{M}, \mathcal{B}(\mathcal{M}))} \int \mathbf{L}_z(\nu) \frac{dP}{dP_{\Theta|Z=z}^{(Q,\lambda)}}(\nu) dP_{\Theta|Z=z}^{(Q,\lambda)}(\nu), \quad (349a)$$

subject to:

$$\int \frac{dP}{dP_{\Theta|Z=z}^{(Q,\lambda)}}(\nu) \log \left( \frac{dP}{dP_{\Theta|Z=z}^{(Q,\lambda)}}(\nu) \right) dP_{\Theta|Z=z}^{(Q,\lambda)}(\nu) \leq c, \text{ and} \quad (349b)$$

$$\int \frac{dP}{dP_{\Theta|Z=z}^{(Q,\lambda)}}(\theta) dP_{\Theta|Z=z}^{(Q,\lambda)}(\theta) = 1. \quad (349c)$$

The remainder of the proof focuses on the problem in which the optimization is over the function  $\frac{dP}{dP_{\Theta|Z=z}^{(Q,\lambda)}}$  instead of the measure  $P$ . This is due to the fact that for all  $P \in \Delta_Q(\mathcal{M})$ , the Radon-Nikodym derivate  $\frac{dP}{dP_{\Theta|Z=z}^{(Q,\lambda)}}$  is unique up to sets of zero measure with respect to the measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$ . Let  $\mathcal{M}$  be the set of measurable functions  $\mathcal{M} \rightarrow \mathbb{R}$  with respect to the measurable spaces  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$  and  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$  that are absolutely integrable with respect to  $P_{\Theta|Z=z}^{(Q,\lambda)}$ . That is, for all  $\hat{g} \in \mathcal{M}$ , it holds that

$$\int |\hat{g}(\theta)| dP_{\Theta|Z=z}^{(Q,\lambda)}(\theta) < \infty. \quad (350)$$



Hence, the optimization problem of interest is:

$$\min_{g \in \mathcal{M}} \int \mathbf{L}_z(\boldsymbol{\nu}) g(\boldsymbol{\nu}) dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}(\boldsymbol{\nu}) \quad (351a)$$

$$\text{s.t: } \int g(\boldsymbol{\nu}) \log(g(\boldsymbol{\nu})) dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}(\boldsymbol{\nu}) \leq c, \text{ and} \quad (351b)$$

$$\int g(\boldsymbol{\theta}) dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}(\boldsymbol{\theta}) = 1. \quad (351c)$$

The Lagrangian of the optimization problem in (351) is a functional  $L : \mathcal{M} \times [0, +\infty)^2 \rightarrow \mathbb{R}$  of the form

$$\begin{aligned} L(g, \alpha, \beta) = & \int \mathbf{L}_z(\boldsymbol{\nu}) g(\boldsymbol{\nu}) dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}(\boldsymbol{\nu}) \\ & + \alpha \left( \int g(\boldsymbol{\nu}) \log(g(\boldsymbol{\nu})) dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}(\boldsymbol{\nu}) - c \right) \\ & + \beta \left( \int g(\boldsymbol{\nu}) dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}(\boldsymbol{\nu}) - 1 \right), \end{aligned} \quad (352)$$

where the reals  $\alpha$  and  $\beta$  are both nonnegative and act as Lagrangian multipliers due to the constraints (351b) and (351c), respectively.

Let  $h : \mathcal{M} \rightarrow \mathbb{R}$  be a function in  $\mathcal{M}$ . The Gateaux differential of the functional  $L$  in (352) at  $(g, \alpha, \beta) \in \mathcal{M} \times [0, +\infty)^2$  in the direction of  $h$ , if it exists, is

$$\partial L(g, \alpha, \beta; h) \triangleq \left. \frac{d}{d\gamma} r(\gamma) \right|_{\gamma=0}, \quad (353)$$

where the real function  $r : \mathbb{R} \rightarrow \mathbb{R}$  is such that for all  $\gamma \in (-\epsilon, \epsilon)$ , with some  $\epsilon > 0$ , satisfies

$$\begin{aligned} r(\gamma) = & \int \mathbf{L}_z(\boldsymbol{\nu}) (g(\boldsymbol{\nu}) + \gamma h(\boldsymbol{\nu})) dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}(\boldsymbol{\nu}) \\ & + \alpha \left( \int (g(\boldsymbol{\nu}) + \gamma h(\boldsymbol{\nu})) \log(g(\boldsymbol{\nu}) + \gamma h(\boldsymbol{\nu})) dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}(\boldsymbol{\nu}) - c \right) \\ & + \beta \left( \int (g(\boldsymbol{\nu}) + \gamma h(\boldsymbol{\nu})) dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}(\boldsymbol{\nu}) - 1 \right). \end{aligned} \quad (354)$$

The proof continues under the assumption that the functions  $g$  and  $h$  are such that the Gateaux differential in (353) exists. That is, the function  $r$  in (354) is differentiable in  $(-\epsilon, \epsilon)$ , with some  $\epsilon > 0$ . Using the same arguments as in the proof of Theorem 3.1, it follows that the derivative of the real function  $r$  in (354) is

$$\begin{aligned} \frac{d}{d\gamma} r(\gamma) = & \int \mathbf{L}_z(\boldsymbol{\nu}) dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}(\boldsymbol{\nu}) + \alpha \int h(\boldsymbol{\nu}) dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}(\boldsymbol{\nu}) \\ & + \alpha \int h(\boldsymbol{\nu}) \log(g(\boldsymbol{\nu}) + \gamma h(\boldsymbol{\nu})) dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}(\boldsymbol{\nu}) \\ & + \beta \int h(\boldsymbol{\nu}) dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}(\boldsymbol{\nu}). \end{aligned} \quad (355)$$

From (353) and (355), it follows that

$$\partial L(g, \alpha, \beta; h) = \int h(\boldsymbol{\nu}) (\mathbf{L}_{\mathbf{z}}(\boldsymbol{\nu}) + \alpha(1 + \log g(\boldsymbol{\nu})) + \beta) dP_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q, \lambda)}(\boldsymbol{\nu}). \quad (356)$$

From [69, Theorem 1, page 217], it holds that a necessary condition for the functional  $L$  in (352) to have a minimum at  $(g, \alpha, \beta) \in \mathcal{M} \times [0, +\infty)^2$  is that for all functions  $h \in \mathcal{M}$ ,

$$\partial L(g, \alpha, \beta; h) = 0, \quad (357)$$

which implies that for all  $\boldsymbol{\nu} \in \mathcal{M}$ ,

$$\mathbf{L}_{\mathbf{z}}(\boldsymbol{\nu}) + \alpha(1 + \log g(\boldsymbol{\nu})) + \beta = 0. \quad (358)$$

Thus,

$$g(\boldsymbol{\nu}) = \exp\left(-\frac{\mathbf{L}_{\mathbf{z}}(\boldsymbol{\nu})}{\alpha}\right) \exp\left(-\frac{\beta + \alpha}{\alpha}\right), \quad (359)$$

where  $\alpha$  and  $\beta$  are chosen to satisfy their corresponding constraints with equality. Denote by  $P^*$  the solution of the optimization problem in (83). Hence, from (359), it follows that

$$\frac{dP^*}{dP_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q, \lambda)}}(\boldsymbol{\nu}) = \frac{\exp\left(-\frac{\mathbf{L}_{\mathbf{z}}(\boldsymbol{\nu})}{\alpha}\right)}{\int \exp\left(-\frac{\mathbf{L}_{\mathbf{z}}(\boldsymbol{\theta})}{\alpha}\right) dP_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q, \lambda)}(\boldsymbol{\theta})}, \quad (360)$$

where  $\alpha$  is chosen to satisfy

$$D(P^* \| P_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q, \lambda)}) = c. \quad (361)$$

From Lemma 3.3, it follows that the probability measure  $P^*$  and the  $\sigma$ -finite

measure  $Q$  satisfy,

$$\frac{dP^*}{dQ}(\nu) = \frac{dP^*}{dP_{\Theta|Z=z}^{(Q,\lambda)}}(\nu) \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\nu) \quad (362)$$

$$= \left( \frac{\exp\left(-\frac{\mathbf{L}_z(\nu)}{\alpha}\right)}{\int \exp\left(-\frac{\mathbf{L}_z(\theta)}{\alpha}\right) dP_{\Theta|Z=z}^{(Q,\lambda)}(\theta)} \right) \left( \frac{\exp\left(-\frac{\mathbf{L}_z(\nu)}{\lambda}\right)}{\int \exp\left(-\frac{\mathbf{L}_z(\theta)}{\lambda}\right) dQ(\theta)} \right) \quad (363)$$

$$= \left( \frac{\exp\left(-\frac{\mathbf{L}_z(\nu)}{\alpha}\right)}{\int \frac{\exp\left(-\frac{\mathbf{L}_z(\theta)}{\alpha}\right) \exp\left(-\frac{\mathbf{L}_z(\theta)}{\lambda}\right)}{\int \exp\left(-\frac{\mathbf{L}_z(\alpha)}{\lambda}\right) dQ(\alpha)} dQ(\theta)} \right) \left( \frac{\exp\left(-\frac{\mathbf{L}_z(\nu)}{\lambda}\right)}{\int \exp\left(-\frac{\mathbf{L}_z(\theta)}{\lambda}\right) dQ(\theta)} \right) \quad (364)$$

$$= \frac{\exp\left(-\left(\frac{1}{\alpha} + \frac{1}{\lambda}\right) \mathbf{L}_z(\nu)\right)}{\int \exp\left(-\left(\frac{1}{\alpha} + \frac{1}{\lambda}\right) \mathbf{L}_z(\nu)\right) dQ(\theta)}, \quad (365)$$

which implies that  $P^*$  is a Gibbs probability measure on  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ , with energy function  $\mathbf{L}_z$ , reference measure  $Q$ , and regularization parameter  $\frac{1}{\frac{1}{\alpha} + \frac{1}{\lambda}}$ , where  $\alpha$  is chosen to satisfy (361). Let the positive real  $\omega$  be  $\omega \triangleq \frac{\alpha\lambda}{\alpha+\lambda}$  and note that  $\omega \in (0, \lambda]$  and satisfies  $D(P_{\Theta|Z=z}^{(Q,\omega)}(\nu) \| P_{\Theta|Z=z}^{(Q,\lambda)}) = c$ . The proof ends by verifying that the objective function in (352) is strictly convex, and thus, the measure  $P_{\Theta|Z=z}^{(Q,\omega)}$  is the unique minimizer. This completes the proof.

## K Proof of Lemma 5.2

Note that for all  $(\lambda_1, \lambda_2) \in \{x \in \mathbb{R} : K_{Q,z}(x) < +\infty\}^2$ , such that  $\lambda_1 > \lambda_2$ , it follows that for all  $\theta \in \text{supp } Q$ , the inequality  $\exp(\lambda_2 \mathbf{L}_z(\theta)) \leq \exp(\lambda_1 \mathbf{L}_z(\theta))$  holds. This implies that  $K_{Q,z}(\lambda_2) \leq K_{Q,z}(\lambda_1) < +\infty$ , which proves that the function is nondecreasing.

The proof of continuity of the function  $K_{Q,z}$  follows from observing that for all

$\alpha \in \{x \in \mathbb{R} : K_{Q,z}(x) < +\infty\}$ , it holds that

$$\lim_{t \rightarrow \alpha} K_{Q,z}(t) = \lim_{t \rightarrow \alpha} \log \left( \int \exp(t \mathbb{L}_z(\boldsymbol{\theta})) dQ(\boldsymbol{\theta}) \right) \quad (366)$$

$$= \log \left( \lim_{t \rightarrow \alpha} \int \exp(t \mathbb{L}_z(\boldsymbol{\theta})) dQ(\boldsymbol{\theta}) \right) \quad (367)$$

$$= \log \left( \int \lim_{t \rightarrow \alpha} \exp(t \mathbb{L}_z(\boldsymbol{\theta})) dQ(\boldsymbol{\theta}) \right) \quad (368)$$

$$= \log \left( \int \exp(\alpha \mathbb{L}_z(\boldsymbol{\theta})) dQ(\boldsymbol{\theta}) \right) \quad (369)$$

$$= K_{Q,z}(\alpha), \quad (370)$$

where (367) and (369) follow from the fact that both the logarithmic and exponential functions are continuous; and the equality in (368) follows from the monotone convergence theorem [55, Theorem 1.6.2]. This shows that the function  $K_{Q,z}$  is continuous in  $\{x \in \mathbb{R} : K_{Q,z}(x) < +\infty\}$ .

The proof of differentiability follows by considering the transport of the  $\sigma$ -finite measure  $Q$  in (22) from the measure space  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$  to the measure space  $([0, +\infty), \mathcal{B}([0, +\infty)))$  through the function  $\mathbb{L}_z$  in (3). Denote the resulting measure in  $([0, +\infty), \mathcal{B}([0, +\infty)))$  by  $P$ . More specifically, for all  $\mathcal{A} \in \mathcal{B}([0, +\infty))$ , it holds that  $P(\mathcal{A}) = Q(\{\boldsymbol{\theta} \in \mathcal{M} : \mathbb{L}_z(\boldsymbol{\theta}) \in \mathcal{A}\})$ . Hence, the function  $K_{Q,z}$  satisfies for all  $t \in \{x \in \mathbb{R} : K_{Q,z}(x) < +\infty\}$ ,

$$K_{Q,z}(t) = \log \left( \int \exp(t \mathbb{L}_z(\boldsymbol{\theta})) dQ(\boldsymbol{\theta}) \right) \quad (371)$$

$$= \log \left( \int \exp(t w) dP(w) \right), \quad (372)$$

where the equality (372) follows from [55, Theorem 1.6.12]. Denote by  $\phi$  the Laplace transform of the measure  $P$ . That is, for all  $t \in \{x \in \mathbb{R} : K_{Q,z}(x) < +\infty\}$ ,

$$\phi(t) = \int \exp(t v) dP(v). \quad (373)$$

Hence,  $\phi(t) = \exp(K_{Q,z}(t))$ . From [70, Theorem 1a (page 439)], it follows that the function  $\phi$  has derivatives of all orders in  $\{x \in \mathbb{R} : K_{Q,z}(x) < +\infty\}$ , and thus, so does the function  $K_{Q,z}$  in the interior of  $\{x \in \mathbb{R} : K_{Q,z}(x) < +\infty\}$ . This completes the proof.

## L Proof of Lemma 5.3

Let  $(\gamma_1, \gamma_2) \in \mathbb{R}^2$ , with  $\gamma_1 \neq \gamma_2$  and  $\alpha \in [0, 1]$  be fixed. Assume that  $K_{Q,z}(\gamma_1) < +\infty$  and  $K_{Q,z}(\gamma_2) < +\infty$ . Then, for all  $\alpha \in (0, 1)$ , the following holds

$$\begin{aligned} & \alpha K_{Q,z}(\gamma_1) + (1 - \alpha) K_{Q,z}(\gamma_2) \\ &= \alpha \log \left( \int \exp(\gamma_1 \mathbf{L}_z(\boldsymbol{\theta})) dQ(\boldsymbol{\theta}) \right) \\ & \quad + (1 - \alpha) \log \left( \int \exp(\gamma_2 \mathbf{L}_z(\boldsymbol{\theta})) dQ(\boldsymbol{\theta}) \right) \end{aligned} \quad (374)$$

$$\begin{aligned} &= \log \left( \left( \int \exp(\gamma_1 \mathbf{L}_z(\boldsymbol{\theta})) dQ(\boldsymbol{\theta}) \right)^\alpha \right) \\ & \quad + \log \left( \left( \int \exp(\gamma_2 \mathbf{L}_z(\boldsymbol{\theta})) dQ(\boldsymbol{\theta}) \right)^{(1-\alpha)} \right) \end{aligned} \quad (375)$$

$$\begin{aligned} &= \log \left( \left( \int \exp(\gamma_1 \mathbf{L}_z(\boldsymbol{\theta})) dQ(\boldsymbol{\theta}) \right)^\alpha \right. \\ & \quad \left. \left( \int \exp(\gamma_2 \mathbf{L}_z(\boldsymbol{\theta})) dQ(\boldsymbol{\theta}) \right)^{(1-\alpha)} \right) \end{aligned} \quad (376)$$

$$\begin{aligned} &= \log \left( \left( \int \exp(\gamma_1 \alpha \mathbf{L}_z(\boldsymbol{\theta}))^p dQ(\boldsymbol{\theta}) \right)^{\frac{1}{p}} \right. \\ & \quad \left. \left( \int \exp(\gamma_2 (1 - \alpha) \mathbf{L}_z(\boldsymbol{\theta}))^q dQ(\boldsymbol{\theta}) \right)^{\frac{1}{q}} \right) \end{aligned} \quad (377)$$

$$\geq \log \left( \int \exp(\gamma_1 \alpha \mathbf{L}_z(\boldsymbol{\theta})) \exp(\gamma_2 (1 - \alpha) \mathbf{L}_z(\boldsymbol{\theta})) dQ(\boldsymbol{\theta}) \right) \quad (378)$$

$$= \log \left( \int \exp \left( (\gamma_1 \alpha + \gamma_2 (1 - \alpha)) \mathbf{L}_z(\boldsymbol{\theta}) \right) dQ(\boldsymbol{\theta}) \right) \quad (379)$$

$$= K_{Q,z}(\gamma_1 \alpha + \gamma_2 (1 - \alpha)), \quad (380)$$

where the inequality in (377) follows with  $\alpha \triangleq \frac{1}{p}$  and  $1 - \alpha \triangleq \frac{1}{q}$ ; the inequality in (378) follows from Hölder's inequality. Hence, equality in (378) holds if and only if there exist two constants  $\beta_1$  and  $\beta_2$ , not simultaneously equal to zero, such that the set

$$\begin{aligned} \mathcal{A} &\triangleq \{\boldsymbol{\theta} \in \mathcal{M} : \beta_1 \exp(\gamma_1 \mathbf{L}_z(\boldsymbol{\theta})) = \beta_2 \exp(\gamma_2 \mathbf{L}_z(\boldsymbol{\theta}))\} \\ &= \left\{ \boldsymbol{\theta} \in \mathcal{M} : \exp((\gamma_1 - \gamma_2) \mathbf{L}_z(\boldsymbol{\theta})) = \frac{\beta_2}{\beta_1} \right\} \end{aligned} \quad (381)$$

$$= \left\{ \boldsymbol{\theta} \in \mathcal{M} : \mathbf{L}_z(\boldsymbol{\theta}) = \frac{\log \frac{\beta_2}{\beta_1}}{(\gamma_1 - \gamma_2)} \right\}, \quad (382)$$

satisfies  $Q(\mathcal{A}) = 1$ . That is, strict inequality in (378) holds if and only if the function  $\mathbf{L}_z$  is separable with respect to the  $\sigma$ -finite measure  $Q$ . When  $\alpha = 0$  or  $\alpha = 1$ , the proof is trivial. This completes the proof.

## M Proof of Lemma 5.4

For all  $s \in \mathcal{K}_{Q,z}$ , with  $\mathcal{K}_{Q,z}$  in (23), the equality in (96) implies the following,

$$K_{Q,z}^{(1)}\left(-\frac{1}{s}\right) = \frac{d}{dt} \log \left( \int \exp(t \mathbf{L}_z(\boldsymbol{\theta})) dQ(\boldsymbol{\theta}) \right) \Big|_{t=-\frac{1}{s}} \quad (383)$$

$$= \int \frac{\mathbf{L}_z(\boldsymbol{\theta}) \exp(t \mathbf{L}_z(\boldsymbol{\theta}))}{\int \exp(t \mathbf{L}_z(\mathbf{v})) dQ(\mathbf{v})} dQ(\boldsymbol{\theta}) \Big|_{t=-\frac{1}{s}} \quad (384)$$

$$= \int \frac{\mathbf{L}_z(\boldsymbol{\theta}) \exp\left(-\frac{1}{s} \mathbf{L}_z(\boldsymbol{\theta})\right)}{\int \exp\left(-\frac{1}{s} \mathbf{L}_z(\mathbf{v})\right) dQ(\mathbf{v})} dQ(\boldsymbol{\theta}) \quad (385)$$

$$= \exp\left(-K_{Q,z}\left(-\frac{1}{s}\right)\right) \int \mathbf{L}_z(\boldsymbol{\theta}) \exp\left(-\frac{1}{s} \mathbf{L}_z(\boldsymbol{\theta})\right) dQ(\boldsymbol{\theta}) \quad (386)$$

$$= \int \mathbf{L}_z(\boldsymbol{\theta}) \exp\left(-K_{Q,z}\left(-\frac{1}{s}\right) - \frac{1}{s} \mathbf{L}_z(\boldsymbol{\theta})\right) dQ(\boldsymbol{\theta}) \quad (387)$$

$$= \int \mathbf{L}_z(\boldsymbol{\theta}) dP_{\boldsymbol{\Theta}|Z=z}^{(Q,s)}(\boldsymbol{\theta}), \quad (388)$$

where the equality in (384) holds from the dominated convergence theorem [55, Theorem 1.6.9]; the equality in (386) follows from (22); and the equality in (388) follows from (25).

For all  $s \in \mathcal{K}_{Q,z}$ , with  $\mathcal{K}_{Q,z}$  in (23), the equalities in (96) and (387) imply that

$$K_{Q,z}^{(2)}\left(-\frac{1}{s}\right) = \frac{d}{dt} \int \mathbf{L}_z(\boldsymbol{\theta}) \exp(-K_{Q,z}(t) + t \mathbf{L}_z(\boldsymbol{\theta})) dQ(\boldsymbol{\theta}) \Big|_{t=-\frac{1}{s}} \quad (389)$$

$$= \int \mathbf{L}_z(\boldsymbol{\theta}) \left(-K_{Q,z}^{(1)}(t) + \mathbf{L}_z(\boldsymbol{\theta})\right) \exp(-K_{Q,z}(t) + t \mathbf{L}_z(\boldsymbol{\theta})) dQ(\boldsymbol{\theta}) \Big|_{t=-\frac{1}{s}} \quad (390)$$

$$= \int \mathbf{L}_z(\boldsymbol{\theta}) \left(-K_{Q,z}^{(1)}\left(-\frac{1}{s}\right) + \mathbf{L}_z(\boldsymbol{\theta})\right) \exp\left(-K_{Q,z}\left(-\frac{1}{s}\right) - \frac{1}{s} \mathbf{L}_z(\boldsymbol{\theta})\right) dQ(\boldsymbol{\theta}) \quad (391)$$

$$= \int \mathbf{L}_z(\boldsymbol{\theta}) \left(-K_{Q,z}^{(1)}\left(-\frac{1}{s}\right) + \mathbf{L}_z(\boldsymbol{\theta})\right) dP_{\boldsymbol{\Theta}|Z=z}^{(Q,s)}(\boldsymbol{\theta}) \quad (392)$$

$$= -K_{Q,z}^{(1)}\left(-\frac{1}{s}\right) \int \mathbf{L}_z(\boldsymbol{\theta}) dP_{\boldsymbol{\Theta}|Z=z}^{(Q,s)}(\boldsymbol{\theta}) + \int (\mathbf{L}_z(\boldsymbol{\theta}))^2 dP_{\boldsymbol{\Theta}|Z=z}^{(Q,s)}(\boldsymbol{\theta}) \quad (393)$$

$$= - \left( K_{Q,z}^{(1)} \left( -\frac{1}{s} \right) \right)^2 + \int (\mathbf{L}_z(\boldsymbol{\theta}))^2 dP_{\boldsymbol{\Theta}|\mathbf{Z}=z}^{(Q,s)}(\boldsymbol{\theta}) \quad (394)$$

$$= \int \left( \mathbf{L}_z(\boldsymbol{\theta}) - K_{Q,z}^{(1)} \left( -\frac{1}{s} \right) \right)^2 dP_{\boldsymbol{\Theta}|\mathbf{Z}=z}^{(Q,s)}(\boldsymbol{\theta}), \quad (395)$$

where the equality in (390) follows from the dominated convergence theorem [55, Theorem 1.6.9]; the equality in (392) is due to a change of measure through the Radon-Nikodym derivative in (25); and the equality in (394) follows from (388).

For all  $s \in \mathcal{K}_{Q,z}$ , with  $\mathcal{K}_{Q,z}$  in (23), the equalities in (96) and (394) imply that

$$K_{Q,z}^{(3)} \left( -\frac{1}{s} \right) = \frac{d}{dt} \left( \int (\mathbf{L}_z(\boldsymbol{\theta}))^2 dP_{\boldsymbol{\Theta}|\mathbf{Z}=z}^{(Q,-\frac{1}{t})}(\boldsymbol{\theta}) - \left( K_{Q,z}^{(1)}(t) \right)^2 \right) \Big|_{t=-\frac{1}{s}} \quad (396)$$

$$= \frac{d}{dt} \left( \int \left( (\mathbf{L}_z(\boldsymbol{\theta}))^2 \exp(-K_{Q,z}(t) + t\mathbf{L}_z(\boldsymbol{\theta})) \right) dQ(\boldsymbol{\theta}) - \left( K_{Q,z}^{(1)}(t) \right)^2 \right) \Big|_{t=-\frac{1}{s}} \quad (397)$$

$$= \int (\mathbf{L}_z(\boldsymbol{\theta}))^2 \left( \frac{d}{dt} \exp(-K_{Q,z}(t) + t\mathbf{L}_z(\boldsymbol{\theta})) \right) \Big|_{t=-\frac{1}{s}} dQ(\boldsymbol{\theta}) - 2K_{Q,z}^{(1)}(t) K_{Q,z}^{(2)}(t) \Big|_{t=-\frac{1}{s}} \quad (398)$$

$$= \int (\mathbf{L}_z(\boldsymbol{\theta}))^2 \left( \mathbf{L}_z(\boldsymbol{\theta}) - K_{Q,z}^{(1)}(t) \right) \exp(-K_{Q,z}(t) + t\mathbf{L}_z(\boldsymbol{\theta})) \Big|_{t=-\frac{1}{s}} dQ(\boldsymbol{\theta}) - 2K_{Q,z}^{(1)}(t) K_{Q,z}^{(2)}(t) \Big|_{t=-\frac{1}{s}} \quad (399)$$

$$= \int (\mathbf{L}_z(\boldsymbol{\theta}))^2 \left( \mathbf{L}_z(\boldsymbol{\theta}) - K_{Q,z}^{(1)} \left( -\frac{1}{s} \right) \right) \exp \left( -K_{Q,z} \left( -\frac{1}{s} \right) - \frac{1}{s} \mathbf{L}_z(\boldsymbol{\theta}) \right) dQ(\boldsymbol{\theta}) - 2K_{Q,z}^{(1)} \left( -\frac{1}{s} \right) K_{Q,z}^{(2)} \left( -\frac{1}{s} \right) \quad (400)$$

$$= \int (\mathbf{L}_z(\boldsymbol{\theta}))^2 \left( \mathbf{L}_z(\boldsymbol{\theta}) - K_{Q,z}^{(1)} \left( -\frac{1}{s} \right) \right) dP_{\boldsymbol{\Theta}|\mathbf{Z}=z}^{(Q,s)}(\boldsymbol{\theta}) - 2K_{Q,z}^{(1)} \left( -\frac{1}{s} \right) K_{Q,z}^{(2)} \left( -\frac{1}{s} \right) \quad (401)$$

$$\begin{aligned}
&= \int (\mathbb{L}_{\mathbf{z}}(\boldsymbol{\theta}))^3 dP_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,s)}(\boldsymbol{\theta}) \\
&\quad - K_{Q,\mathbf{z}}^{(1)}\left(-\frac{1}{s}\right) \int (\mathbb{L}_{\mathbf{z}}(\boldsymbol{\theta}))^2 dP_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,s)}(\boldsymbol{\theta}) \\
&\quad - 2K_{Q,\mathbf{z}}^{(1)}\left(-\frac{1}{s}\right) K_{Q,\mathbf{z}}^{(2)}\left(-\frac{1}{s}\right) \tag{402}
\end{aligned}$$

$$\begin{aligned}
&= \int (\mathbb{L}_{\mathbf{z}}(\boldsymbol{\theta}))^3 dP_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,s)}(\boldsymbol{\theta}) \\
&\quad - K_{Q,\mathbf{z}}^{(1)}\left(-\frac{1}{s}\right) \left( K_{Q,\mathbf{z}}^{(2)}\left(-\frac{1}{s}\right) + \left( K_{Q,\mathbf{z}}^{(1)}\left(-\frac{1}{s}\right) \right)^2 \right) \\
&\quad - 2K_{Q,\mathbf{z}}^{(1)}\left(-\frac{1}{s}\right) K_{Q,\mathbf{z}}^{(2)}\left(-\frac{1}{s}\right) \tag{403}
\end{aligned}$$

$$= \int (\mathbb{L}_{\mathbf{z}}(\boldsymbol{\theta}))^3 dP_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,s)}(\boldsymbol{\theta}) - K_{Q,\mathbf{z}}^{(1)}\left(-\frac{1}{s}\right)^3 \tag{404}$$

$$- 3K_{Q,\mathbf{z}}^{(1)}\left(-\frac{1}{s}\right) K_{Q,\mathbf{z}}^{(2)}\left(-\frac{1}{s}\right) \tag{405}$$

$$= \int \left( \mathbb{L}_{\mathbf{z}}(\boldsymbol{\theta}) - K_{Q,\mathbf{z}}^{(1)}\left(-\frac{1}{s}\right) \right)^3 dP_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,s)}(\boldsymbol{\theta}), \tag{406}$$

where the equality in (397) follows from (25); and the equality in (398) follows from the dominated convergence theorem [55, Theorem 1.6.9]; the equality in (401) follows from (25); and the equality in (403) follows from (394).

This completes the proof.

## N Proof of Theorem 6.1

The proof is based on the analysis of the derivative of  $K_{Q,\mathbf{z}}^{(1)}\left(-\frac{1}{\lambda}\right)$  with respect to  $\lambda$  in  $\text{int}\mathcal{K}_{Q,\mathbf{z}}$ . This is due to Corollary 6.1. For instance, note that

$$\frac{d}{d\lambda} \mathbb{R}_{\mathbf{z}}\left(P_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}\right) = \frac{d}{d\lambda} K_{Q,\mathbf{z}}^{(1)}\left(-\frac{1}{\lambda}\right) \tag{407}$$

$$= \frac{1}{\lambda^2} K_{Q,\mathbf{z}}^{(2)}\left(-\frac{1}{\lambda}\right) \tag{408}$$

$$\geq 0, \tag{409}$$

where the equality in (408) follows from Lemma 5.4. The inequality in (409) implies that the expected empirical risk  $\mathbb{R}_{\mathbf{z}}\left(P_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}\right) = K_{Q,\mathbf{z}}^{(1)}\left(-\frac{1}{\lambda}\right)$  in (105) is nondecreasing with respect to  $\lambda$ . The rest of the proof consists in showing that for all  $\alpha \in \mathcal{K}_{Q,\mathbf{z}}$ , the function  $K_{Q,\mathbf{z}}^{(2)}$  in (96) satisfies  $K_{Q,\mathbf{z}}^{(2)}\left(-\frac{1}{\alpha}\right) > 0$  if and only if the function  $\mathbb{L}_{\mathbf{z}}$  in (3) is separable. For doing so, a handful of preliminary results are described in the following subsection. The proof of Theorem 6.1 resumes in Subsection N.2



## N.1 Preliminaries

Given a positive real  $\lambda \in \mathcal{K}_{Q,z}$ , with  $\mathcal{K}_{Q,z}$  in (23), consider a partition of  $\mathcal{M}$  formed by the sets  $\mathcal{R}_0(\lambda)$ ,  $\mathcal{R}_1(\lambda)$  and  $\mathcal{R}_2(\lambda)$ , such that

$$\mathcal{R}_0(\lambda) \triangleq \left\{ \nu \in \mathcal{M} : \mathbb{L}_z(\nu) = \mathbb{R}_z \left( P_{\Theta|Z=z}^{(Q,\lambda)} \right) \right\}, \quad (410a)$$

$$\mathcal{R}_1(\lambda) \triangleq \left\{ \nu \in \mathcal{M} : \mathbb{L}_z(\nu) < \mathbb{R}_z \left( P_{\Theta|Z=z}^{(Q,\lambda)} \right) \right\}, \text{ and} \quad (410b)$$

$$\mathcal{R}_2(\lambda) \triangleq \left\{ \nu \in \mathcal{M} : \mathbb{L}_z(\nu) > \mathbb{R}_z \left( P_{\Theta|Z=z}^{(Q,\lambda)} \right) \right\}, \quad (410c)$$

where the functional  $\mathbb{R}_z$  is in (18) and the probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  is in (25). The sets in (410) exhibit several properties that are central for proving the main results of this section.

**Lemma N.1.** *The probability measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  in (25), satisfies*

$$P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{R}_1(\lambda)) > 0, \quad (411)$$

*if and only if*

$$P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{R}_2(\lambda)) > 0, \quad (412)$$

*where the sets  $\mathcal{R}_1(\cdot)$  and  $\mathcal{R}_2(\cdot)$  are in (410b) and (410c), respectively.*

*Proof:* The proof is divided into two parts. In the first part, given a real  $\alpha \in \mathcal{K}_{Q,z}$ , it is proven that if the set  $\mathcal{R}_1(\alpha)$  is nonnegligible with respect to  $P_{\Theta|Z=z}^{(Q,\alpha)}$ , then the set  $\mathcal{R}_2(\alpha)$  is nonnegligible with respect to  $P_{\Theta|Z=z}^{(Q,\alpha)}$ . The second part proves the converse.

The first part is proved by contradiction. Assume that set  $\mathcal{R}_2(\alpha)$  is negligible

with respect to  $P_{\Theta|Z=z}^{(Q,\alpha)}$ . Hence, from Lemma 5.4, it holds that

$$K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right) = \int_{\mathcal{R}_0(\alpha)} \mathbb{L}_z(\nu) dP_{\Theta|Z=z}^{(Q,\alpha)}(\nu) + \int_{\mathcal{R}_1(\alpha)} \mathbb{L}_z(\nu) dP_{\Theta|Z=z}^{(Q,\alpha)}(\nu) \quad (413)$$

$$+ \int_{\mathcal{R}_2(\alpha)} \mathbb{L}_z(\nu) dP_{\Theta|Z=z}^{(Q,\alpha)}(\nu) \quad (414)$$

$$= \int_{\mathcal{R}_0(\alpha)} \mathbb{L}_z(\nu) dP_{\Theta|Z=z}^{(Q,\alpha)}(\nu) + \int_{\mathcal{R}_1(\alpha)} \mathbb{L}_z(\nu) dP_{\Theta|Z=z}^{(Q,\alpha)}(\nu) \quad (415)$$

$$= K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right) P_{\Theta|Z=z}^{(Q,\alpha)}(\mathcal{R}_0(\alpha)) \quad (416)$$

$$+ \int_{\mathcal{R}_1(\alpha)} \mathbb{L}_z(\nu) dP_{\Theta|Z=z}^{(Q,\alpha)}(\nu) \quad (417)$$

$$< K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right) P_{\Theta|Z=z}^{(Q,\alpha)}(\mathcal{R}_0(\alpha)) \quad (418)$$

$$+ K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right) P_{\Theta|Z=z}^{(Q,\alpha)}(\mathcal{R}_1(\alpha)) \quad (418)$$

$$= K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right) (P_{\Theta|Z=z}^{(Q,\alpha)}(\mathcal{R}_0(\alpha)) + P_{\Theta|Z=z}^{(Q,\alpha)}(\mathcal{R}_1(\alpha))) \quad (419)$$

$$= K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right), \quad (420)$$

which is a contradiction.

The second part of the proof follows the same arguments as in the first part. Assume that the set  $\mathcal{R}_1(\alpha)$  is negligible with respect to  $P_{\Theta|Z=z}^{(Q,\alpha)}$ . Hence, from Lemma 5.4, it holds that

$$K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right) = \int_{\mathcal{R}_0(\alpha)} \mathbb{L}_z(\nu) dP_{\Theta|Z=z}^{(Q,\alpha)}(\nu) + \int_{\mathcal{R}_1(\alpha)} \mathbb{L}_z(\nu) dP_{\Theta|Z=z}^{(Q,\alpha)}(\nu) \quad (421)$$

$$+ \int_{\mathcal{R}_2(\alpha)} \mathbb{L}_z(\nu) dP_{\Theta|Z=z}^{(Q,\alpha)}(\nu) \quad (422)$$

$$= \int_{\mathcal{R}_0(\alpha)} \mathbb{L}_z(\nu) dP_{\Theta|Z=z}^{(Q,\alpha)}(\nu) + \int_{\mathcal{R}_2(\alpha)} \mathbb{L}_z(\nu) dP_{\Theta|Z=z}^{(Q,\alpha)}(\nu) \quad (423)$$

$$= K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right) P_{\Theta|Z=z}^{(Q,\alpha)}(\mathcal{R}_0(\alpha)) \quad (424)$$

$$+ \int_{\mathcal{R}_2(\alpha)} \mathbb{L}_z(\nu) dP_{\Theta|Z=z}^{(Q,\alpha)}(\nu) \quad (424)$$

$$> K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right) P_{\Theta|Z=z}^{(Q,\alpha)}(\mathcal{R}_0(\alpha)) + K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right) P_{\Theta|Z=z}^{(Q,\alpha)}(\mathcal{R}_2(\alpha)) \quad (425)$$

$$= K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right) (P_{\Theta|Z=z}^{(Q,\alpha)}(\mathcal{R}_0(\alpha)) + P_{\Theta|Z=z}^{(Q,\alpha)}(\mathcal{R}_2(\alpha))) \quad (426)$$

$$= K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right), \quad (427)$$

which is also a contradiction. This completes the proof.  $\blacksquare$

A more general result can be immediately obtained by combining Lemma 3.4 and Lemma N.1.

**Lemma N.2.** *For all  $\alpha \in \mathcal{K}_{Q,z}$ , with  $\mathcal{K}_{Q,z}$  in (23), the measure  $P_{\Theta|Z=z}^{(Q,\lambda)}$  in (25), satisfies*

$$P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{R}_1(\alpha)) > 0, \quad (428)$$

*if and only if*

$$P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{R}_2(\alpha)) > 0, \quad (429)$$

*where the sets  $\mathcal{R}_1(\alpha)$  and  $\mathcal{R}_2(\alpha)$  are in (410b) and (410c), respectively.*

## N.2 The proof

The rest of the proof of Theorem 6.1 is divided into two parts. In the first part, it is shown that if for all  $\alpha \in \mathcal{K}_{Q,z}$ ,  $K_{Q,z}^{(2)}(-\frac{1}{\alpha}) > 0$ , then the function  $\mathbf{L}_z$  in (3) is separable. The second part of the proof, consists in showing that if the function  $\mathbf{L}_z$  is separable, then, for all  $\alpha \in \mathcal{K}_{Q,z}$ ,  $K_{Q,z}^{(2)}(-\frac{1}{\alpha}) > 0$ .

The first part is as follows. From Lemma 5.4, it holds that for all  $\alpha \in \mathcal{K}_{Q,z}$ ,

$$K_{Q,z}^{(2)}\left(-\frac{1}{\alpha}\right) = \int \left( \mathbf{L}_z(\boldsymbol{\theta}) - K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right) \right)^2 dP_{\Theta|Z=z}^{(Q,\alpha)}(\boldsymbol{\theta}) \quad (430)$$

$$= \int_{\mathcal{R}_0(\alpha)} \left( \mathbf{L}_z(\boldsymbol{\theta}) - K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right) \right)^2 dP_{\Theta|Z=z}^{(Q,\alpha)}(\boldsymbol{\theta}) \quad (431)$$

$$+ \int_{\mathcal{R}_1(\alpha)} \left( \mathbf{L}_z(\boldsymbol{\theta}) - K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right) \right)^2 dP_{\Theta|Z=z}^{(Q,\alpha)}(\boldsymbol{\theta}) \quad (432)$$

$$+ \int_{\mathcal{R}_2(\alpha)} \left( \mathbf{L}_z(\boldsymbol{\theta}) - K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right) \right)^2 dP_{\Theta|Z=z}^{(Q,\alpha)}(\boldsymbol{\theta}), \quad (433)$$

where the sets  $\mathcal{R}_0(\alpha)$ ,  $\mathcal{R}_1(\alpha)$ , and  $\mathcal{R}_2(\alpha)$  are respectively defined in (410). Hence,

$$K_{Q,z}^{(2)}\left(-\frac{1}{\alpha}\right) = \int_{\mathcal{R}_1(\alpha)} \left( \mathbf{L}_z(\boldsymbol{\theta}) - K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right) \right)^2 dP_{\Theta|Z=z}^{(Q,\alpha)}(\boldsymbol{\theta}) \quad (434)$$

$$+ \int_{\mathcal{R}_2(\gamma)} \left( \mathbf{L}_z(\boldsymbol{\theta}) - K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right) \right)^2 dP_{\Theta|Z=z}^{(Q,\alpha)}(\boldsymbol{\theta}). \quad (435)$$

Under the assumption that for all  $\alpha \in \mathcal{K}_{Q,z}$  the function  $K_{Q,z}^{(2)}$  in (96) satisfies  $K_{Q,z}^{(2)}(-\frac{1}{\alpha}) > 0$ , it follows that at least one of the following claims is true:

- (a)  $P_{\Theta|Z=z}^{(Q,\alpha)}(\mathcal{R}_1(\alpha)) > 0$ ; and
- (b)  $P_{\Theta|Z=z}^{(Q,\alpha)}(\mathcal{R}_2(\alpha)) > 0$ .

Nonetheless, from Lemma N.1, it follows that both claims (a) and (b) hold simultaneously. Hence, the sets  $\mathcal{R}_1(\alpha)$  and  $\mathcal{R}_2(\alpha)$  are both nonnegligible with

respect to  $P_{\Theta|Z=z}^{(Q,\alpha)}$  and moreover, it holds that for all  $(\nu_1, \nu_2) \in \mathcal{R}_1(\alpha) \times \mathcal{R}_2(\alpha)$ ,

$$+\infty > \mathbb{L}_z(\nu_1) > K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right) > \mathbb{L}_z(\nu_2), \quad (436)$$

where  $\mathbb{L}_z(\nu_1) < +\infty$  follows from the fact that  $P_{\Theta|Z=z}^{(Q,\lambda)}(\{\theta \in \mathcal{M} : \mathbb{L}_z(\theta) = +\infty\}) = 0$  (Lemma 3.2). This proves that under the assumption that for all  $\alpha \in \mathcal{K}_{Q,z}$ ,  $K_{Q,z}^{(2)}(-\frac{1}{\alpha}) > 0$ , the function  $\mathbb{L}_z$  in (3) is separable with respect to  $P_{\Theta|Z=z}^{(Q,\alpha)}$ . From Lemma 5.1, it holds that the function  $\mathbb{L}_z$  is separable with respect to  $Q$ . This completes the first part of the proof.

The second part of the proof is simpler. Assume that the empirical risk function  $\mathbb{L}_z$  in (3) is separable with respect to  $P_{\Theta|Z=z}^{(Q,\alpha)}$ . That is, for all  $\gamma \in \mathcal{K}_{Q,z}$ , there exist a positive real  $c_\gamma > 0$ ; and two subsets  $\mathcal{A}(\gamma)$  and  $\mathcal{B}(\gamma)$  of  $\mathcal{M}$  that are nonnegligible with respect to  $P_{\Theta|Z=z}^{(Q,\gamma)}$  in (25) and verify that for all  $(\nu_1, \nu_2) \in \mathcal{A}(\gamma) \times \mathcal{B}(\gamma)$ ,

$$+\infty > \mathbb{L}_z(\nu_1) > c_\gamma > \mathbb{L}_z(\nu_2). \quad (437)$$

From Lemma 5.4, it holds that

$$K_{Q,z}^{(2)}\left(-\frac{1}{\gamma}\right) = \int \left(\mathbb{L}_z(\theta) - K_{Q,z}^{(1)}\left(-\frac{1}{\gamma}\right)\right)^2 dP_{\Theta|Z=z}^{(Q,\gamma)}(\theta) \quad (438)$$

$$= \int_{\mathcal{A}(\gamma)} \left(\mathbb{L}_z(\theta) - K_{Q,z}^{(1)}\left(-\frac{1}{\gamma}\right)\right)^2 dP_{\Theta|Z=z}^{(Q,\gamma)}(\theta) \quad (439)$$

$$+ \int_{\mathcal{B}(\gamma)} \left(\mathbb{L}_z(\theta) - K_{Q,z}^{(1)}\left(-\frac{1}{\gamma}\right)\right)^2 dP_{\Theta|Z=z}^{(Q,\gamma)}(\theta) \quad (440)$$

$$+ \int_{\mathcal{M} \setminus (\mathcal{A}(\gamma) \cup \mathcal{B}(\gamma))} \left(\mathbb{L}_z(\theta) - K_{Q,z}^{(1)}\left(-\frac{1}{\gamma}\right)\right)^2 dP_{\Theta|Z=z}^{(Q,\gamma)}(\theta) \quad (441)$$

$$> 0, \quad (442)$$

where the inequality (442) follows from the following facts. First, if  $c_\gamma < K_{Q,z}^{(1)}(-\frac{1}{\gamma})$ , with  $c_\gamma$  in (437), then for all  $\nu \in \mathcal{B}(\gamma)$ , it holds that  $K_{Q,z}^{(1)}(-\frac{1}{\gamma}) > c_\gamma > \mathbb{L}_z(\nu)$ , and thus,

$$\left(\mathbb{L}_z(\nu) - K_{Q,z}^{(1)}\left(-\frac{1}{\gamma}\right)\right)^2 > \left(c_\gamma - K_{Q,z}^{(1)}\left(-\frac{1}{\gamma}\right)\right)^2, \quad (443)$$

which implies,

$$\int_{\mathcal{B}(\gamma)} \left(\mathbb{L}_z(\theta) - K_{Q,z}^{(1)}\left(-\frac{1}{\gamma}\right)\right)^2 dP_{\Theta|Z=z}^{(Q,\gamma)}(\theta) > \left(c_\gamma - K_{Q,z}^{(1)}\left(-\frac{1}{\gamma}\right)\right)^2 P_{\Theta|Z=z}^{(Q,\gamma)}(\mathcal{B}(\gamma)) \quad (444)$$

$$> 0. \quad (445)$$

Second, if  $c_\gamma \geq K_{Q,z}^{(1)}\left(-\frac{1}{\gamma}\right)$  then for all  $\nu \in \mathcal{A}(\gamma)$ , it holds that  $L_z(\nu) > c_\gamma \geq K_{Q,z}^{(1)}\left(-\frac{1}{\gamma}\right)$ , and thus,

$$\left(L_z(\nu) - K_{Q,z}^{(1)}\left(-\frac{1}{\gamma}\right)\right)^2 > \left(c_\gamma - K_{Q,z}^{(1)}\left(-\frac{1}{\gamma}\right)\right)^2, \quad (446)$$

which implies,

$$\int_{\mathcal{A}(\gamma)} \left(L_z(\theta) - K_{Q,z}^{(1)}\left(-\frac{1}{\gamma}\right)\right)^2 dP_{\Theta|Z=z}^{(Q,\gamma)}(\theta) > \left(c_\gamma - K_{Q,z}^{(1)}\left(-\frac{1}{\gamma}\right)\right)^2 P_{\Theta|Z=z}^{(Q,\gamma)}(\mathcal{A}(\gamma)) \quad (447)$$

$$> 0. \quad (448)$$

Hence, under the assumption that the empirical risk function  $L_z$  in (3) is separable, it holds that for all  $\gamma \in \mathcal{K}_{Q,z}$ ,  $K_{Q,z}^{(2)}\left(-\frac{1}{\gamma}\right) > 0$ . This completes the proof.

## O Proof of Lemma 6.1

Consider the partition of the set  $\mathcal{M}$  formed by the sets  $\mathcal{A}_0$ ,  $\mathcal{A}_1$ , and  $\mathcal{A}_2$  in (336). From (97), for all  $\lambda \in \mathcal{K}_{Q,z}$ , with  $\mathcal{K}_{Q,z}$  in (23), it holds that,

$$K_{Q,z}^{(1)}\left(-\frac{1}{\lambda}\right) = \int_{\mathcal{A}_0} L_z(\theta) dP_{\Theta|Z=z}^{(Q,\lambda)}(\theta) + \int_{\mathcal{A}_1} L_z(\theta) dP_{\Theta|Z=z}^{(Q,\lambda)}(\theta) \quad (449)$$

$$+ \int_{\mathcal{A}_2} L_z(\theta) dP_{\Theta|Z=z}^{(Q,\lambda)}(\theta) \quad (450)$$

$$= \int_{\mathcal{A}_0} L_z(\theta) dP_{\Theta|Z=z}^{(Q,\lambda)}(\theta) + \int_{\mathcal{A}_2} L_z(\theta) dP_{\Theta|Z=z}^{(Q,\lambda)}(\theta) \quad (451)$$

$$= \delta_{Q,z}^* P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_{Q,z}^*) + \int_{\mathcal{A}_2} L_z(\theta) dP_{\Theta|Z=z}^{(Q,\lambda)}(\theta) \quad (452)$$

$$\geq \delta_{Q,z}^* P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_{Q,z}^*) + \delta_{Q,z}^* P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{A}_2) \quad (453)$$

$$= \delta_{Q,z}^*, \quad (454)$$

where the equality in (451) follows by noticing that  $Q(\mathcal{A}_1) = 0$ , which implies that  $P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{A}_1) = 0$  (Lemma 3.3); the equality in (452) follows from noticing that  $\mathcal{A}_0 = \mathcal{L}_{Q,z}^*$ , with  $\mathcal{L}_{Q,z}^*$  in (39); and the equality in (453) follows from (336c). This completes the proof.

## P Proof of Theorem 6.2

From (452) in the proof of Lemma 6.1, it holds that

$$\lim_{\lambda \rightarrow 0^+} K_{Q,z}^{(1)} \left( -\frac{1}{\lambda} \right) = \lim_{\lambda \rightarrow 0^+} \delta_{Q,z}^* P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_{Q,z}^*) + \lim_{\lambda \rightarrow 0^+} \int_{\mathcal{A}_2} \mathbf{L}_z(\boldsymbol{\theta}) dP_{\Theta|Z=z}^{(Q,\lambda)}(\boldsymbol{\theta}) \quad (455)$$

$$\begin{aligned} &= \lim_{\lambda \rightarrow 0^+} \delta_{Q,z}^* P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_{Q,z}^*) \\ &\quad + \lim_{\lambda \rightarrow 0^+} \int_{\mathcal{A}_2} \mathbf{L}_z(\boldsymbol{\theta}) \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) \end{aligned} \quad (456)$$

$$\begin{aligned} &= \lim_{\lambda \rightarrow 0^+} \delta_{Q,z}^* P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_{Q,z}^*) \\ &\quad + \int_{\mathcal{A}_2} \mathbf{L}_z(\boldsymbol{\theta}) \lim_{\lambda \rightarrow 0^+} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) \end{aligned} \quad (457)$$

$$= \delta_{Q,z}^* \lim_{\lambda \rightarrow 0^+} P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_{Q,z}^*) \quad (458)$$

$$= \delta_{Q,z}^*, \quad (459)$$

where, the equality in (457) follows from noticing two facts: (a) For all  $\lambda \in \mathcal{K}_{Q,z}$ , the Randon-Nikodym derivative  $\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}$  is positive and finite (Lemma 3.2); and (b) For all  $\boldsymbol{\theta} \in \mathcal{A}_2$ , it holds that  $\lim_{\lambda \rightarrow 0^+} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\boldsymbol{\theta}) = 0$  (Lemma 3.6). Hence, the dominated convergence theorem [55, Theorem 1.6.9] holds. The inequality in (458) follows from Lemma 3.7. This completes the proof.

## Q Proof of Theorem 9.1

From Theorem 6.1, it follows that for all  $(\lambda_1, \lambda_2) \in \mathcal{K}_{Q,z} \times \mathcal{K}_{Q,z}$  with  $\lambda_1 > \lambda_2$ ,

$$\int \mathbf{L}_z(\boldsymbol{\alpha}) \frac{dP_{\Theta|Z=z}^{(Q,\lambda_1)}}{dQ}(\boldsymbol{\alpha}) dQ(\boldsymbol{\alpha}) \geq \int \mathbf{L}_z(\boldsymbol{\alpha}) \frac{dP_{\Theta|Z=z}^{(Q,\lambda_2)}}{dQ}(\boldsymbol{\alpha}) dQ(\boldsymbol{\alpha}),$$

which implies the following inclusions:

$$\mathcal{R}_1(\lambda_2) \subseteq \mathcal{R}_1(\lambda_1), \text{ and} \quad (460a)$$

$$\mathcal{R}_2(\lambda_1) \subseteq \mathcal{R}_2(\lambda_2), \quad (460b)$$

with the sets  $\mathcal{R}_1(\cdot)$  and  $\mathcal{R}_2(\cdot)$  in (410). From (193), it holds that for all  $i \in \{1, 2\}$ ,

$$\mathcal{N}_{Q,z}(\lambda_i) = \mathcal{R}_2(\lambda_i)^c, \quad (461)$$

where the complement is with respect to  $\mathcal{M}$ . Thus, the inclusion in (460b) and the equality in (461) yields,

$$\mathcal{N}_{Q,z}(\lambda_1) \supseteq \mathcal{N}_{Q,z}(\lambda_2). \quad (462)$$

The inclusion  $\mathcal{M} \supseteq \mathcal{N}_{Q,z}(\lambda_1)$  follows from (193). Alternatively, the inclusion  $\mathcal{N}_{Q,z}(\lambda_2) \supseteq \mathcal{N}_{Q,z}^*$ , follows from Lemma 6.1 and from observing that for all  $\nu \in \mathcal{N}_{Q,z}^*$ ,

$$R_z \left( P_{\Theta|Z=z}^{(Q,\lambda_2)} \right) \geq \delta_{Q,z}^* = L_z(\nu), \quad (463)$$

which implies that  $\nu \in \mathcal{N}_{Q,z}(\lambda_2)$ . This completes the proof of (197).

The proof of (198) is as follows. From the intermediate value theorem [66, Theorem 4.23] and the assumption that the empirical risk function  $L_z$  in (3) is continuous on  $\mathcal{M}$ , it follows that for all  $\lambda \in \mathcal{K}_{Q,z}$ , there always exists a model  $\theta \in \mathcal{M}$ , such that

$$L_z(\theta) = \int L_z(\alpha) dP_{\Theta|Z=z}^{(Q,\lambda)}(\alpha), \quad (464)$$

which implies that  $\mathcal{R}_0(\lambda)$  is not empty, and as a consequence,  $\mathcal{N}_{Q,z}(\lambda) = \mathcal{R}_0(\lambda) \cup \mathcal{R}_1(\lambda)$  is not empty. Hence, for all  $\theta \in \mathcal{R}_0(\lambda_1)$  it holds that  $\theta \notin \mathcal{N}_{Q,z}(\lambda_2)$ . This proves that the elements of  $\mathcal{R}_0(\lambda_1)$  are in  $\mathcal{N}_{Q,z}(\lambda_1)$  but not in  $\mathcal{N}_{Q,z}(\lambda_2)$ . This, together with (462), verifies that

$$\mathcal{N}_{Q,z}(\lambda_1) \supset \mathcal{N}_{Q,z}(\lambda_2). \quad (465)$$

The strict inclusion  $\mathcal{M} \supset \mathcal{N}_{Q,z}(\lambda_1)$  is proved by contradiction. Assume that there exists a  $\lambda \in \mathcal{K}_{Q,z}$  such that  $\mathcal{M} = \mathcal{N}_{Q,z}(\lambda)$ . Then,  $\mathcal{R}_2(\lambda) = \emptyset$  and thus,  $P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{R}_2(\lambda)) = 0$ , which together with Lemma N.1, implies that  $P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{R}_1(\lambda)) = 0$  and consequently,

$$P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{R}_0(\lambda)) = 1. \quad (466)$$

This contradicts the assumption that the function  $L_z$  is separable (Definition 5.1). Hence,  $\mathcal{M} \supset \mathcal{N}_{Q,z}(\lambda_1)$ .

Finally, the strict inclusion  $\mathcal{N}_{Q,z}(\lambda_2) \supset \mathcal{N}_{Q,z}^*$  is proved by contradiction. Assume that there exists a  $\lambda \in \mathcal{K}_{Q,z}$  such that  $\mathcal{N}_{Q,z}^* = \mathcal{N}_{Q,z}(\lambda)$ . That is,

$$\{\theta \in \mathcal{M} : L_z(\theta) \leq \delta_{Q,z}^*\} = \mathcal{N}_{Q,z}^* \quad (467)$$

$$= \mathcal{N}_{Q,z}(\lambda) \quad (468)$$

$$= \left\{ \theta \in \mathcal{M} : L_z(\theta) \leq K_{Q,z}^{(1)} \left( -\frac{1}{\lambda} \right) \right\}. \quad (469)$$

Hence, three cases might arise:

(a) there exists a  $\lambda \in \mathcal{K}_{Q,z}$ , such that  $\delta_{Q,z}^* < K_{Q,z}^{(1)} \left( -\frac{1}{\lambda} \right)$  and it holds that

$$\left\{ \nu \in \mathcal{M} : \delta_{Q,z}^* < L_z(\nu) \leq K_{Q,z}^{(1)} \left( -\frac{1}{\lambda} \right) \right\} = \emptyset;$$

(b) there exists a  $\lambda \in \mathcal{K}_{Q,z}$ , such that  $\delta_{Q,z}^* > K_{Q,z}^{(1)}(-\frac{1}{\lambda})$  and it holds that

$$\left\{ \nu \in \mathcal{M} : K_{Q,z}^{(1)}\left(-\frac{1}{\lambda}\right) < \mathbb{L}_z(\nu) \leq \delta_{Q,z}^* \right\} = \emptyset;$$

or (c) there exists a  $\lambda \in \mathcal{K}_{Q,z}$ , such that  $\delta_{Q,z}^* = K_{Q,z}^{(1)}(-\frac{1}{\lambda})$ .

The cases (a) and (b) are absurd. Hence, the proof is complete only by considering the case (c). In the case (c), it holds that,

$$\mathcal{R}_1(\lambda) = \left\{ \nu \in \mathcal{M} : \mathbb{L}_z(\nu) < \delta_{Q,z}^* \right\}, \quad (470)$$

and from the definition of  $\delta_{Q,z}^*$  in (38), it holds that

$$P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{R}_1(\lambda)) = 0. \quad (471)$$

From Lemma N.1 and (471), it follows that,

$$P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{R}_2(\lambda)) = 0. \quad (472)$$

Finally, by noticing that

$$1 = P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{R}_0(\lambda)) + P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{R}_1(\lambda)) + P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{R}_2(\lambda)) \quad (473)$$

$$= P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{R}_0(\lambda)), \quad (474)$$

reveals a contradiction to the assumption that the function  $\mathbb{L}_z$  is separable with respect to  $P_{\Theta|Z=z}^{(Q,\lambda)}$  (and thus, separable with respect to  $Q$  by Lemma 5.1). This completes the proof of (198).

## R Proof of Theorem 9.2

The proof of (199) is based on the analysis of the derivative of  $P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{A})$  with respect to  $\lambda$ , for some fixed set  $\mathcal{A} \subseteq \mathcal{B}(\mathcal{M})$ . More specifically, given a  $\gamma \in \mathcal{K}_{Q,z}$ , it holds that

$$P_{\Theta|Z=z}^{(Q,\gamma)}(\mathcal{A}) = \int_{\mathcal{A}} \frac{dP_{\Theta|Z=z}^{(Q,\gamma)}}{dQ}(\alpha) dQ(\alpha), \quad (475)$$

and from the fundamental theorem of calculus [66, Theorem 6.21], it follows that for all  $(\lambda_1, \lambda_2) \in \mathcal{K}_{Q,z} \times \mathcal{K}_{Q,z}$  with  $\lambda_1 > \lambda_2$ ,

$$P_{\Theta|Z=z}^{(Q,\lambda_1)}(\mathcal{A}) - P_{\Theta|Z=z}^{(Q,\lambda_2)}(\mathcal{A}) = \int_{\lambda_2}^{\lambda_1} \frac{d}{d\gamma} P_{\Theta|Z=z}^{(Q,\gamma)}(\mathcal{A}) d\gamma \quad (476)$$

$$= \int_{\lambda_2}^{\lambda_1} \frac{d}{d\gamma} \int_{\mathcal{A}} \frac{dP_{\Theta|Z=z}^{(Q,\gamma)}}{dQ}(\alpha) dQ(\alpha) d\gamma \quad (477)$$

$$= \int_{\lambda_2}^{\lambda_1} \int_{\mathcal{A}} \frac{d}{d\gamma} \frac{dP_{\Theta|Z=z}^{(Q,\gamma)}}{dQ}(\alpha) dQ(\alpha) d\gamma, \quad (478)$$



where the equality in (477) follows from (475); and the equality in (478) holds from Lemma 3.2 and the dominated convergence theorem [55, Theorem 1.6.9].

For all  $\theta \in \text{supp } Q$ , the following holds,

$$\frac{d}{d\lambda} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) = \frac{d}{d\lambda} \frac{\exp\left(-\frac{L_z(\theta)}{\lambda}\right)}{\int \exp\left(-\frac{L_z(\nu)}{\lambda}\right) dQ(\nu)} \quad (479)$$

$$\begin{aligned} &= \frac{\frac{1}{\lambda^2} L_z(\theta) \exp\left(-\frac{L_z(\theta)}{\lambda}\right)}{\int \exp\left(-\frac{L_z(\nu)}{\lambda}\right) dQ(\nu)} \\ &= \frac{\frac{1}{\lambda^2} \exp\left(-\frac{L_z(\theta)}{\lambda}\right) \int L_z(\alpha) \exp\left(-\frac{L_z(\alpha)}{\lambda}\right) dQ(\alpha)}{\left(\int \exp\left(-\frac{L_z(\nu)}{\lambda}\right) dQ(\nu)\right)^2} \quad (480) \end{aligned}$$

$$\begin{aligned} &= \frac{1}{\lambda^2} L_z(\theta) \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) \\ &\quad - \frac{1}{\lambda^2} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) \int L_z(\nu) \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\nu) dQ(\nu) \quad (481) \end{aligned}$$

$$= \frac{1}{\lambda^2} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) \left( L_z(\theta) - \int L_z(\nu) dP_{\Theta|Z=z}^{(Q,\lambda)}(\nu) \right). \quad (482)$$

Plugging (482) into (478) yields,

$$\begin{aligned} &P_{\Theta|Z=z}^{(Q,\lambda_1)}(\mathcal{A}) - P_{\Theta|Z=z}^{(Q,\lambda_2)}(\mathcal{A}) \\ &= \int_{\lambda_2}^{\lambda_1} \int_{\mathcal{A}} \frac{1}{\gamma^2} \frac{dP_{\Theta|Z=z}^{(Q,\gamma)}}{dQ}(\alpha) \left( L_z(\alpha) - \int L_z(\nu) dP_{\Theta|Z=z}^{(Q,\gamma)}(\nu) \right) dQ(\alpha) d\gamma \quad (483) \\ &= \int_{\lambda_2}^{\lambda_1} \int_{\mathcal{A}} \frac{1}{\gamma^2} \left( L_z(\alpha) - \int L_z(\nu) dP_{\Theta|Z=z}^{(Q,\gamma)}(\nu) \right) dP_{\Theta|Z=z}^{(Q,\gamma)}(\alpha) d\gamma. \quad (484) \end{aligned}$$

Note that for all  $\alpha \in \mathcal{N}_{Q,z}(\lambda_2)$ , it holds that for all  $\gamma \in (\lambda_2, \lambda_1)$ ,

$$L_z(\alpha) - \int L_z(\nu) dP_{\Theta|Z=z}^{(Q,\gamma)}(\nu) \leq 0, \quad (485)$$

and thus,

$$\int_{\mathcal{N}_{Q,z}(\lambda_2)} \frac{1}{\gamma^2} \left( L_z(\alpha) - \int L_z(\nu) dP_{\Theta|Z=z}^{(Q,\gamma)}(\nu) \right) dP_{\Theta|Z=z}^{(Q,\gamma)}(\alpha) \leq 0. \quad (486)$$

The equalities in (484) and (486), with  $\mathcal{A} = \mathcal{N}_{Q,z}(\lambda)$ , imply that

$$P_{\Theta|Z=z}^{(Q,\lambda_1)}(\mathcal{N}_{Q,z}(\lambda_2)) - P_{\Theta|Z=z}^{(Q,\lambda_2)}(\mathcal{N}_{Q,z}(\lambda_2)) \leq 0. \quad (487)$$

The inequality  $0 < P_{\Theta|Z=z}^{(Q,\lambda_1)}(\mathcal{N}_{Q,z}(\lambda_2))$  in (199) is proved by contradiction. Assume that for some  $\lambda \in \mathcal{K}_{Q,z}$  it holds that  $0 = P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{N}_{Q,z}(\lambda_2))$ . Then,  $P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{R}_0(\lambda_2)) + P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{R}_1(\lambda_2)) = 0$ , which implies that  $P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{R}_2(\lambda_2)) = 1$ , which is a contradiction. See for instance, Lemma N.2. This completes the proof of (199).

The proof of strict inequality in (199) is divided into two parts. The first part shows that if for all pairs  $(\lambda_1, \lambda_2) \in \mathcal{K}_{Q,z} \times \mathcal{K}_{Q,z}$  with  $\lambda_1 > \lambda_2$ ,

$$P_{\Theta|Z=z}^{(Q,\lambda_1)}(\mathcal{N}_{Q,z}(\lambda_2)) < P_{\Theta|Z=z}^{(Q,\lambda_2)}(\mathcal{N}_{Q,z}(\lambda_2)), \quad (488)$$

then the function  $\mathbf{L}_z$  is separable with respect to  $Q$ . The second part of the proof shows that if the function  $\mathbf{L}_z$  is separable with respect to  $Q$ , then, for all pairs  $(\lambda_1, \lambda_2) \in \mathcal{K}_{Q,z} \times \mathcal{K}_{Q,z}$  with  $\lambda_1 > \lambda_2$ , the inequality in (488) holds.

The first part is as follows. In the proof of Theorem 9.1 it is shown (see (484)) that for all pairs  $(\lambda_1, \lambda_2) \in \mathcal{K}_{Q,z} \times \mathcal{K}_{Q,z}$  with  $\lambda_1 > \lambda_2$ ,

$$\begin{aligned} & P_{\Theta|Z=z}^{(Q,\lambda_1)}(\mathcal{N}_{Q,z}(\lambda_2)) - P_{\Theta|Z=z}^{(Q,\lambda_2)}(\mathcal{N}_{Q,z}(\lambda_2)) \\ &= \int_{\lambda_2}^{\lambda_1} \int_{\mathcal{N}_{Q,z}(\lambda_2)} \frac{1}{\gamma^2} \left( \mathbf{L}_z(\alpha) - K_{Q,z}^{(1)}\left(-\frac{1}{\gamma}\right) \right) dP_{\Theta|Z=z}^{(Q,\gamma)}(\alpha) d\gamma. \end{aligned} \quad (489)$$

Assume that for a given pair  $(\lambda_1, \lambda_2) \in \mathcal{K}_{Q,z} \times \mathcal{K}_{Q,z}$ , with  $\lambda_1 > \lambda_2$ , the inequality in (488) holds. Then, from (489),

$$0 > \int_{\lambda_2}^{\lambda_1} \int_{\mathcal{N}_{Q,z}(\lambda_2)} \frac{1}{\gamma^2} \left( \mathbf{L}_z(\alpha) - K_{Q,z}^{(1)}\left(-\frac{1}{\gamma}\right) \right) dP_{\Theta|Z=z}^{(Q,\gamma)}(\alpha) d\gamma \quad (490)$$

$$= \int_{\lambda_2}^{\lambda_1} \int_{\mathcal{R}_1(\lambda_2)} \frac{1}{\gamma^2} \left( \mathbf{L}_z(\alpha) - K_{Q,z}^{(1)}\left(-\frac{1}{\gamma}\right) \right) dP_{\Theta|Z=z}^{(Q,\gamma)}(\alpha) d\gamma, \quad (491)$$

where the equality in (491) follows from noticing that  $\mathcal{R}_0(\lambda_2)$  and  $\mathcal{R}_1(\lambda_2)$  form a partition of  $\mathcal{N}_{Q,z}(\lambda_2)$ , with the sets  $\mathcal{R}_0(\lambda_2)$ ,  $\mathcal{R}_1(\lambda_2)$  and  $\mathcal{N}_{Q,z}(\lambda_2)$  defined in (410a), (410b), and (193), respectively.

The inequality in (491) implies that the set  $\mathcal{R}_1(\lambda_2)$  is nonnegligible with respect to  $P_{\Theta|Z=z}^{(Q,\gamma)}$ , for some  $\gamma \in (\lambda_2, \lambda_1)$ . Hence, from Lemma N.2, it follows that both sets  $\mathcal{R}_1(\lambda_2)$  and  $\mathcal{R}_2(\lambda_2)$  are nonnegligible with respect to  $P_{\Theta|Z=z}^{(Q,\gamma)}$ .

From the arguments above, it has been proved that given a pair  $(\lambda_1, \lambda_2) \in \mathcal{K}_{Q,z} \times \mathcal{K}_{Q,z}$  with  $\lambda_1 > \lambda_2$ , if

$$P_{\Theta|Z=z}^{(Q,\lambda_1)}(\mathcal{N}_{Q,z}(\lambda_2)) < P_{\Theta|Z=z}^{(Q,\lambda_2)}(\mathcal{N}_{Q,z}(\lambda_2)), \quad (492)$$

then there always exists a positive  $\gamma \in (\lambda_1, \lambda_2)$  such that the sets  $\mathcal{R}_1(\lambda_2)$  and  $\mathcal{R}_2(\lambda_2)$  are not negligible with respect to  $P_{\Theta|Z=z}^{(Q,\gamma)}$ . Moreover, such sets  $\mathcal{R}_1(\lambda_2)$

and  $\mathcal{R}_2(\lambda_2)$  satisfy for all  $(\nu_1, \nu_2) \in \mathcal{R}_2(\lambda) \times \mathcal{R}_1(\lambda)$ ,

$$+\infty > \mathbf{L}_z(\nu_1) > K_{Q,z}^{(1)}\left(-\frac{1}{\lambda}\right) > \mathbf{L}_z(\nu_2), \quad (493)$$

which together with Definition 5.1 verify that the function  $\mathbf{L}_z$  is separable with respect to  $P_{\Theta|Z=z}^{(Q,\gamma)}$  (and thus, with respect to  $Q$  by Lemma 5.1). This ends the first part of the proof.

The second part of the proof is under the assumption that the empirical risk function  $\mathbf{L}_z$  in (3) is separable with respect to  $Q$  (and thus, with respect to  $P_{\Theta|Z=z}^{(Q,\gamma)}$  by Lemma 5.1). That is, from Definition 5.1, for all  $\gamma \in \mathcal{K}_{Q,z}$ , there exist a positive real  $c_\gamma > 0$  and two subsets  $\mathcal{A}(\gamma)$  and  $\mathcal{B}(\gamma)$  of  $\mathcal{M}$  nonnegligible with respect to  $P_{\Theta|Z=z}^{(Q,\gamma)}$  in (25) that verify that for all  $(\nu_1, \nu_2) \in \mathcal{A}(\gamma) \times \mathcal{B}(\gamma)$ ,

$$\mathbf{L}_z(\nu_1) > c_\gamma > \mathbf{L}_z(\nu_2). \quad (494)$$

In the proof of Theorem 9.1, cf. (484), it has been proved that given a pair  $(\alpha_1, \alpha_2) \in \mathcal{K}_{Q,z} \times \mathcal{K}_{Q,z}$ , with  $\alpha_1 > \gamma > \alpha_2$ , it holds that for all subsets  $\mathcal{A}$  of  $\mathcal{M}$ ,

$$\begin{aligned} & P_{\Theta|Z=z}^{(Q,\alpha_1)}(\mathcal{A}) - P_{\Theta|Z=z}^{(Q,\alpha_2)}(\mathcal{A}) \\ &= \int_{\alpha_2}^{\alpha_1} \int_{\mathcal{A}} \frac{1}{\lambda^2} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\alpha) \left( \mathbf{L}_z(\alpha) - K_{Q,z}^{(1)}\left(-\frac{1}{\lambda}\right) \right) dP(\alpha) d\lambda \\ &= \int_{\alpha_2}^{\alpha_1} \int_{\mathcal{A}} \frac{1}{\lambda^2} \left( \mathbf{L}_z(\alpha) - K_{Q,z}^{(1)}\left(-\frac{1}{\lambda}\right) \right) dP_{\Theta|Z=z}^{(Q,\lambda)}(\alpha) d\lambda. \end{aligned} \quad (495)$$

Hence, two cases are studied. The first case considers that

$$c_\gamma < K_{Q,z}^{(1)}\left(-\frac{1}{\gamma}\right), \quad (496)$$

with  $c_\gamma$  in (494). The second case considers that

$$c_\gamma \geq K_{Q,z}^{(1)}\left(-\frac{1}{\gamma}\right). \quad (497)$$

In the first case, it follows from (193) that

$$\mathcal{B}(\gamma) \subset \mathcal{N}_{Q,z}(\gamma), \quad (498)$$

which implies that

$$P_{\Theta|Z=z}^{(Q,\gamma)}(\mathcal{N}_{Q,z}(\gamma)) \geq P_{\Theta|Z=z}^{(Q,\gamma)}(\mathcal{B}(\gamma)) \quad (499)$$

$$> 0, \quad (500)$$

where, the inequality in (500) follows from the fact that  $\mathcal{B}(\gamma)$  is nonnegligible with respect to  $P_{\Theta|Z=z}^{(Q,\gamma)}$ . This implies that the set  $\mathcal{N}_{Q,z}(\gamma)$  is not negligible with

respect  $P_{\Theta|Z=z}^{(Q,\gamma)}$ . Moreover, from (193) and (498), it follows that for all  $\alpha \in \mathcal{N}_{Q,z}(\gamma)$  and for all  $\lambda \in (\gamma, \alpha_1)$ ,

$$\mathbb{L}_z(\alpha) - \int \mathbb{L}_z(\nu) dP_{\Theta|Z=z}^{(Q,\lambda)}(\nu) < \mathbb{L}_z(\alpha) - c_\gamma \quad (501)$$

$$< 0, \quad (502)$$

where the inequality in (501) follows from (496); and the inequality in (502) follows from (494). Thus,

$$\int_\gamma^{\alpha_1} \int_{\mathcal{N}_{Q,z}(\gamma)} \frac{1}{\lambda^2} \left( \mathbb{L}_z(\alpha) - K_{Q,z}^{(1)} \left( -\frac{1}{\lambda} \right) \right) dP_{\Theta|Z=z}^{(Q,\lambda)}(\alpha) d\lambda < 0,$$

which implies, from (495), that

$$P_{\Theta|Z=z}^{(Q,\alpha_1)}(\mathcal{N}_{Q,z}(\gamma)) - P_{\Theta|Z=z}^{(Q,\gamma)}(\mathcal{N}_{Q,z}(\gamma)) < 0. \quad (503)$$

Assume now that  $c_\gamma \geq K_{Q,z}^{(1)} \left( -\frac{1}{\gamma} \right)$ . Hence, the following holds

$$\mathcal{A}(\gamma) \subseteq \mathcal{R}_2(\gamma), \quad (504)$$

which implies that

$$P_{\Theta|Z=z}^{(Q,\gamma)}(\mathcal{R}_2(\gamma)) \geq P_{\Theta|Z=z}^{(Q,\gamma)}(\mathcal{A}(\gamma)) \quad (505)$$

$$> 0, \quad (506)$$

where the inequality in (506) follows from the fact that  $\mathcal{A}(\gamma)$  is nonnegligible with respect to  $P_{\Theta|Z=z}^{(Q,\gamma)}$ . This implies that the set  $\mathcal{R}_2(\gamma)$  is not negligible with respect to  $P_{\Theta|Z=z}^{(Q,\gamma)}$ . From Lemma N.1, it follows that both  $\mathcal{R}_1(\gamma)$  and  $\mathcal{R}_2(\gamma)$  are nonnegligible with respect to  $P_{\Theta|Z=z}^{(Q,\gamma)}$ . Using this result, the following holds,

$$P_{\Theta|Z=z}^{(Q,\gamma)}(\mathcal{N}_{Q,z}(\gamma)) \geq P_{\Theta|Z=z}^{(Q,\gamma)}(\mathcal{R}_1(\gamma)) \quad (507)$$

$$> 0, \quad (508)$$

which proves the set  $\mathcal{N}_{Q,z}(\gamma)$  is nonnegligible with respect to  $P_{\Theta|Z=z}^{(Q,\gamma)}$ .

From (193) and Theorem 6.1, it follows that for all  $\alpha \in \mathcal{N}_{Q,z}(\gamma)$  and for all  $\lambda \in (\gamma, \alpha_1)$ ,

$$0 \geq \mathbb{L}_z(\alpha) - \int \mathbb{L}_z(\nu) dP_{\Theta|Z=z}^{(Q,\gamma)}(\nu) \quad (509)$$

$$> \mathbb{L}_z(\alpha) - \int \mathbb{L}_z(\nu) dP_{\Theta|Z=z}^{(Q,\lambda)}(\nu). \quad (510)$$

Thus,

$$\int_\gamma^{\alpha_1} \int_{\mathcal{N}_{Q,z}(\gamma)} \frac{1}{\lambda^2} \left( \mathbb{L}_z(\alpha) - K_{Q,z}^{(1)} \left( -\frac{1}{\lambda} \right) \right) dP_{\Theta|Z=z}^{(Q,\lambda)}(\alpha) d\lambda < 0,$$

which implies, from (495), that

$$P_{\Theta|Z=z}^{(Q,\alpha_1)}(\mathcal{N}_{Q,z}(\gamma)) - P_{\Theta|Z=z}^{(Q,\gamma)}(\mathcal{N}_{Q,z}(\gamma)) < 0. \quad (511)$$

This completes the proof.

## S Proof of Lemma 9.2

The proof is based on the following two observations. First, note that  $(\mathcal{N}_{Q,z}(\lambda_2))^c = \mathcal{R}_2(\lambda_2)$ , with the set  $\mathcal{R}_2(\cdot)$  defined in (410c). Second, note that

$$\mathcal{N}_{Q,z}(\lambda_1) = \mathcal{N}_{Q,z}(\lambda_2) \cup (\mathcal{N}_{Q,z}(\lambda_1) \cap \mathcal{R}_2(\lambda_2)), \quad (512)$$

and the fact that the sets  $\mathcal{N}_{Q,z}(\lambda_2)$  and  $(\mathcal{N}_{Q,z}(\lambda_1) \cap \mathcal{R}_2(\lambda_2))$  are disjoint. Hence, for all  $i \in \{1, 2\}$ ,

$$P_{\Theta|Z=z}^{(\lambda_i)}(\mathcal{N}_{Q,z}(\lambda_1)) = P_{\Theta|Z=z}^{(\lambda_i)}(\mathcal{N}_{Q,z}(\lambda_2) \cup (\mathcal{N}_{Q,z}(\lambda_1) \cap \mathcal{R}_2(\lambda_2))) \quad (513)$$

$$= P_{\Theta|Z=z}^{(\lambda_i)}(\mathcal{N}_{Q,z}(\lambda_2)) + P_{\Theta|Z=z}^{(\lambda_i)}(\mathcal{N}_{Q,z}(\lambda_1) \cap \mathcal{R}_2(\lambda_2)) \quad (514)$$

$$= P_{\Theta|Z=z}^{(\lambda_i)}(\mathcal{N}_{Q,z}(\lambda_2)), \quad (515)$$

where the equality in (514) follows from Lemma 3.3 and the equality in (200).

Finally, under the assumption that the empirical function  $\mathbb{L}_z$  in (3) is separable, it holds from Theorem 9.2 that

$$P_{\Theta|Z=z}^{(Q,\lambda_1)}(\mathcal{N}_{Q,z}(\lambda_2)) < P_{\Theta|Z=z}^{(Q,\lambda_2)}(\mathcal{N}_{Q,z}(\lambda_2)). \quad (516)$$

Plugging (515) into (516), with  $i = 1$ , yields,

$$P_{\Theta|Z=z}^{(Q,\lambda_1)}(\mathcal{N}_{Q,z}(\lambda_1)) < P_{\Theta|Z=z}^{(Q,\lambda_2)}(\mathcal{N}_{Q,z}(\lambda_2)), \quad (517)$$

and this completes the proof.

## T Proof of Theorem 11.2

Consider the following lemma.

**Lemma T.1.** *Given two probability measures  $P_1$  and  $P_2$  over  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ , with  $P_2$  absolutely continuous with respect to  $P_1$ , the following holds for all  $\mathbf{z} \in (\mathcal{X} \times \mathcal{Y})^n$ ,*

$$\mathbb{R}_z(P_2) - \mathbb{R}_z(P_1) \leq \inf_{t < 0} \left( \frac{D(P_2 \| P_1) + \log \left( \int \exp(t(\mathbb{L}_z(\boldsymbol{\theta}) - \mathbb{R}_z(P_1))) dP_1(\boldsymbol{\theta}) \right)}{t} \right), \quad (518)$$

where the function  $\mathbb{L}_z$  and the functional  $\mathbb{R}_z$  are defined in (3) and in (18), respectively.

*Proof:* From [53, Corollary 4.15, Page 100], it follows that the probability measures  $P_1$  and  $P_2$  in  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$  satisfy the following equality:

$$D(P_2 \| P_1) = \sup_f \int f(\boldsymbol{\theta}) dP_2(\boldsymbol{\theta}) - \log \int \exp(f(\boldsymbol{\theta})) dP_1(\boldsymbol{\theta}), \quad (519)$$

where the supremum is over the space of all measurable functions  $f$  with respect to  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$  and  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ , such that  $\int \exp(f(\boldsymbol{\theta})) dP_1(\boldsymbol{\theta}) < \infty$ . Hence, for all  $\mathbf{z} \in (\mathcal{X} \times \mathcal{Y})^n$  and for all  $t \in (-\infty, 0)$ , it follows that the empirical risk function  $\mathbf{L}_{\mathbf{z}}$  in (3) satisfies that

$$D(P_2 \| P_1) \geq \int t \mathbf{L}_{\mathbf{z}}(\boldsymbol{\theta}) dP_2(\boldsymbol{\theta}) - \log \int \exp(t \mathbf{L}_{\mathbf{z}}(\boldsymbol{\theta})) dP_1(\boldsymbol{\theta}) \quad (520)$$

$$\begin{aligned} &\geq \int t \mathbf{L}_{\mathbf{z}}(\boldsymbol{\theta}) dP_2(\boldsymbol{\theta}) \\ &\quad - \log \int \exp(t \mathbf{L}_{\mathbf{z}}(\boldsymbol{\theta}) + t \mathbf{R}_{\mathbf{z}}(P_1) - t \mathbf{R}_{\mathbf{z}}(P_1)) dP_1(\boldsymbol{\theta}) \end{aligned} \quad (521)$$

$$\begin{aligned} &= \int t \mathbf{L}_{\mathbf{z}}(\boldsymbol{\theta}) dP_2(\boldsymbol{\theta}) - t \mathbf{R}_{\mathbf{z}}(P_1) \\ &\quad - \log \int \exp(t \mathbf{L}_{\mathbf{z}}(\boldsymbol{\theta}) - t \mathbf{R}_{\mathbf{z}}(P_1)) dP_1(\boldsymbol{\theta}) \end{aligned} \quad (522)$$

$$= t \mathbf{R}_{\mathbf{z}}(P_2) - t \mathbf{R}_{\mathbf{z}}(P_1) - \log \int \exp(t \mathbf{L}_{\mathbf{z}}(\boldsymbol{\theta}) - t \mathbf{R}_{\mathbf{z}}(P_1)) dP_1(\boldsymbol{\theta}), \quad (523)$$

which leads to

$$\mathbf{R}_{\mathbf{z}}(P_2) - \mathbf{R}_{\mathbf{z}}(P_1) \leq \frac{D(P_2 \| P_1) + \log \int \exp(t(\mathbf{L}_{\mathbf{z}}(\boldsymbol{\theta}) - \mathbf{R}_{\mathbf{z}}(P_1))) dP_1(\boldsymbol{\theta})}{t}. \quad (524)$$

Given that  $t$  can be chosen arbitrarily in  $(-\infty, 0)$ , it holds that

$$\mathbf{R}_{\mathbf{z}}(P_2) - \mathbf{R}_{\mathbf{z}}(P_1) \leq \inf_{t \in (-\infty, 0)} \frac{D(P_2 \| P_1) + \log \int \exp(t(\mathbf{L}_{\mathbf{z}}(\boldsymbol{\theta}) - \mathbf{R}_{\mathbf{z}}(P_1))) dP_1(\boldsymbol{\theta})}{t}, \quad (525)$$

which completes the proof.  $\blacksquare$

From Lemma T.1, it holds that the probability measure  $P_{\boldsymbol{\Theta} | \mathbf{Z} = \mathbf{z}}^{(Q, \lambda)}$  in (25), satisfies for all  $P \in \Delta_Q(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ ,

$$\begin{aligned} &\mathbf{R}_{\mathbf{z}}(P) - \mathbf{R}_{\mathbf{z}}(P_{\boldsymbol{\Theta} | \mathbf{Z} = \mathbf{z}}^{(Q, \lambda)}) \\ &\leq \inf_{t \in (-\infty, 0)} \left( \frac{D(P \| P_{\boldsymbol{\Theta} | \mathbf{Z} = \mathbf{z}}^{(Q, \lambda)})}{t} + \frac{\log \left( \int \exp \left( t \left( \mathbf{L}_{\mathbf{z}}(\boldsymbol{\theta}) - K_{Q, \mathbf{z}}^{(1)} \left( -\frac{1}{\lambda} \right) \right) \right) dP_{\boldsymbol{\Theta} | \mathbf{Z} = \mathbf{z}}^{(Q, \lambda)}(\boldsymbol{\theta}) \right)}{t} \right), \end{aligned} \quad (526)$$

where the function  $K_{Q, \mathbf{z}}^{(1)}$  is defined in (97) and satisfies (105). Moreover, for

all  $t \in (-\infty, 0)$ ,

$$\begin{aligned} & \log \left( \int \exp \left( t \left( \mathbf{L}_z(\boldsymbol{\theta}) - K_{Q,z}^{(1)} \left( -\frac{1}{\lambda} \right) \right) \right) dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}(\boldsymbol{\theta}) \right) \\ &= \log \left( \int \exp(t \mathbf{L}_z(\boldsymbol{\theta})) dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}(\boldsymbol{\theta}) \right) - t K_{Q,z}^{(1)} \left( -\frac{1}{\lambda} \right) \end{aligned} \quad (527)$$

$$= J_{z,Q,\lambda}(t) - t K_{Q,z}^{(1)} \left( -\frac{1}{\lambda} \right) \quad (528)$$

$$\leq \frac{1}{2} t^2 \beta_{Q,z}^2, \quad (529)$$

where the equality in (528) follows from (142); the inequality in (529) follows from Theorem 8.1; and the constant  $\beta_{Q,z}$  is defined in (165).

Plugging (529) into (526) yields for all  $t \in (-\infty, 0)$ ,

$$\mathbf{R}_z(P) - \mathbf{R}_z(P_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}) \leq \inf_{t \in (-\infty, 0)} \frac{D(P \| P_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}) + \frac{1}{2} t^2 \beta_{Q,z}^2}{t}. \quad (530)$$

Let the  $c \in \mathbb{R}$  be defined as follows:

$$c \triangleq \mathbf{R}_z(P) - \mathbf{R}_z(P_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}). \quad (531)$$

Hence, from (530), it follows that for all  $t \in (-\infty, 0)$ ,

$$c t - \frac{1}{2} t^2 \beta_{Q,z}^2 \leq D(P \| P_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}). \quad (532)$$

The rest of the proof consists in finding an explicit expression for the absolute value of  $c$  in (532). To this aim, consider the function  $\phi : \mathbb{R} \rightarrow \mathbb{R}$  such that

$$\phi(\alpha) = \frac{1}{2} \alpha^2 \beta_{Q,z}^2, \quad (533)$$

and note that  $\phi$  is a positive and strictly convex function with  $\phi(0) = 0$ . Let the Legendre-Fenchel transform of  $\phi$  be the function  $\phi^* : \mathbb{R} \rightarrow \mathbb{R}$ , and thus for all  $x \in \mathbb{R}$ ,

$$\phi^*(x) = \max_{t \in (-\infty, 0)} x t - \phi(t). \quad (534)$$

In particular, note that

$$\phi^*(c) \leq D(P \| P_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}). \quad (535)$$

Note that for all  $x \in \mathbb{R}$  and for all  $t \in (-\infty, 0)$ , the function  $\phi^*$  in (534) satisfies

$$x t - \frac{1}{2} t^2 \beta_{Q,z}^2 \leq \phi^*(x) = x \alpha^*(x) - \phi(\alpha^*(x)), \quad (536)$$

where the term  $\alpha^*(x)$  represents the unique solution in  $\alpha$  within the interval  $(-\infty, 0)$  to

$$\frac{d}{d\alpha} (x \alpha - \phi(\alpha)) = x - \alpha \beta_{Q,z}^2 = 0. \quad (537)$$

That is,

$$\alpha^*(x) = \frac{x}{\beta_{Q,z}^2}. \quad (538)$$

Plugging (538) into (536) yields,

$$\phi^*(x) = \frac{x^2}{2\beta_{Q,z}^2}. \quad (539)$$

Hence, from (535) and (536), given  $c$  in (531) for all  $t \in (-\infty, 0)$ ,

$$ct - \frac{1}{2}t^2\beta_{Q,z}^2 \leq \phi^*(c) \leq D(P \| P_{\Theta|Z=z}^{(Q,\lambda)}), \quad (540)$$

and thus,

$$\frac{c^2}{2\beta_{Q,z}^2} \leq D(P \| P_{\Theta|Z=z}^{(Q,\lambda)}). \quad (541)$$

This implies that

$$c \leq \sqrt{2\beta_{Q,z}^2 D(P \| P_{\Theta|Z=z}^{(Q,\lambda)})} \quad (542)$$

and

$$c \geq -\sqrt{2\beta_{Q,z}^2 D(P \| P_{\Theta|Z=z}^{(Q,\lambda)})}, \quad (543)$$

which leads to

$$\left| \int \mathbb{L}_z(\theta) dP(\theta) - \int \mathbb{L}_z(\theta) dP_{\Theta|Z=z}^{(Q,\lambda)}(\theta) \right| \leq \sqrt{2\beta_{Q,z}^2 D(P \| P_{\Theta|Z=z}^{(Q,\lambda)})}, \quad (544)$$

and completes the proof.

## U Proof of Theorem 11.4

Under the condition that  $\lambda \in \mathcal{K}_{Q,P_Z}$ , from Theorem 11.1 and Definition 11.2, it follows that the generalization error  $\mathbb{G}_{Q,\lambda}(P_Z)$  in (232) satisfies

$$\begin{aligned} & \mathbb{G}_{Q,\lambda}(P_Z) \\ &= \lambda \int \left( D(P_{\Theta|Z=\nu}^{(Q,\lambda)} \| Q) + D(P_{\Theta}^{(Q,\lambda)} \| P_{\Theta|Z=\nu}^{(Q,\lambda)}) - D(P_{\Theta}^{(Q,\lambda)} \| Q) \right) dP_Z(\nu), \end{aligned} \quad (545)$$

$$= \lambda \left( \int D(P_{\Theta|Z=\nu}^{(Q,\lambda)} \| P_{\Theta}^{(Q,\lambda)}) dP_Z(\nu) + \int D(P_{\Theta}^{(Q,\lambda)} \| P_{\Theta|Z=\nu}^{(Q,\lambda)}) dP_Z(\nu) \right), \quad (546)$$



where the equality in (546) follows from the fact that

$$\begin{aligned} & \int \left( D(P_{\Theta|Z=\nu}^{(Q,\lambda)} \| Q) - D(P_{\Theta}^{(Q,\lambda)} \| Q) \right) dP_Z(\nu) \\ &= \int D(P_{\Theta|Z=\nu}^{(Q,\lambda)} \| Q) dP_Z(\nu) - D(P_{\Theta}^{(Q,\lambda)} \| Q) \end{aligned} \quad (547)$$

$$= \int \left( \int \log \left( \frac{dP_{\Theta|Z=\nu}^{(Q,\lambda)}}{dQ}(\theta) \right) dP_{\Theta|Z=\nu}^{(Q,\lambda)}(\theta) \right) P_Z(\nu) - D(P_{\Theta}^{(Q,\lambda)} \| Q) \quad (548)$$

$$\begin{aligned} &= \int \left( \int \log \left( \frac{dP_{\Theta|Z=\nu}^{(Q,\lambda)}}{dQ}(\theta) \right) dP_{\Theta|Z=\nu}^{(Q,\lambda)}(\theta) \right) dP_Z(\nu) \\ &\quad - \int \log \left( \frac{dP_{\Theta}^{(Q,\lambda)}}{dQ}(\theta) \right) dP_{\Theta}^{(Q,\lambda)}(\theta) \end{aligned} \quad (549)$$

$$\begin{aligned} &= \int \left( \int \log \left( \frac{dP_{\Theta|Z=\nu}^{(Q,\lambda)}}{dQ}(\theta) \right) dP_{\Theta|Z=\nu}^{(Q,\lambda)}(\theta) \right) dP_Z(\nu) \\ &\quad - \int \left( \int \log \left( \frac{dP_{\Theta}^{(Q,\lambda)}}{dQ}(\theta) \right) dP_{\Theta}^{(Q,\lambda)}(\theta) \right) dP_Z(\nu) \end{aligned} \quad (550)$$

$$\begin{aligned} &= \int \left( \int \left( \log \left( \frac{dP_{\Theta|Z=\nu}^{(Q,\lambda)}}{dQ}(\theta) \right) \right. \right. \\ &\quad \left. \left. + \log \left( \frac{dQ}{dP_{\Theta}^{(Q,\lambda)}}(\theta) \right) \right) dP_{\Theta|Z=\nu}^{(Q,\lambda)}(\theta) \right) dP_Z(\nu) \end{aligned} \quad (551)$$

$$= \int \left( \int \log \left( \frac{dP_{\Theta|Z=\nu}^{(Q,\lambda)}}{dP_{\Theta}^{(Q,\lambda)}}(\theta) \right) dP_{\Theta|Z=\nu}^{(Q,\lambda)}(\theta) \right) dP_Z(\nu) \quad (552)$$

$$= \int D(P_{\Theta|Z=\nu}^{(Q,\lambda)} \| P_{\Theta}^{(Q,\lambda)}) dP_Z(\nu). \quad (553)$$

The equality in (550) follows from (230); and the equality in (551) follows from the fact that the measures  $Q$  and  $P_{\Theta|Z=\nu}^{(Q,\lambda)}$ , with  $\nu \in \text{supp } P_Z$ , are mutually absolutely continuous. This completes the proof.

## References

- [1] S. M. Perlaza, G. Bisson, I. Esnaola, A. Jean-Marie, and S. Rini, “Empirical risk minimization with relative entropy regularization: Optimality and sensitivity,” in *Proceedings of the IEEE International Symposium on Information Theory (ISIT)*, Espoo, Finland, Jul. 2022, pp. 684–689.
- [2] O. Catoni, *Statistical Learning Theory and Stochastic Optimization: Ecole d’Eté de Probabilités de Saint-Flour, XXXI-2001*, 1st ed. New York, NY, USA: Springer Science & Business Media, 2004, vol. 1851.
- [3] L. Zdeborová and F. Krzakala, “Statistical physics of inference: Thresholds and algorithms,” *Advances in Physics*, vol. 65, no. 5, pp. 453–552, Aug. 2016.
- [4] P. Alquier, J. Ridgway, and N. Chopin, “On the properties of variational approximations of Gibbs posteriors,” *Journal of Machine Learning Research*, vol. 17, no. 1, pp. 8374–8414, Dec. 2016.
- [5] H. P. Young, *Strategic Learning and its Limits*, 1st ed. Oxford, UK: Oxford University Press, 2004.
- [6] C. P. Robert, *The Bayesian Choice: From Decision-theoretic Foundations to Computational Implementation*, 1st ed. New York, NY, USA: Springer, 2007.
- [7] G. Aminian, L. Toni, and M. R. Rodrigues, “Information-theoretic bounds on the moments of the generalization error of learning algorithms,” in *Proceedings of the IEEE International Symposium on Information Theory (ISIT)*, Melbourne, Australia, Jul. 2021, pp. 682–687.
- [8] G. Aminian, Y. Bu, L. Toni, M. Rodrigues, and G. Wornell, “An exact characterization of the generalization error for the Gibbs algorithm,” *Advances in Neural Information Processing Systems*, vol. 34, pp. 8106–8118, Dec. 2021.
- [9] W. Jiang and M. A. Tanner, “Gibbs posterior for variable selection in high-dimensional classification and data mining,” *The Annals of Statistics*, vol. 36, no. 5, pp. 2207–2231, Oct. 2008.
- [10] S. M. Perlaza, I. Esnaola, G. Bisson, and H. V. Poor, “On the validation of Gibbs algorithms: Training datasets, test datasets and their aggregation,” in *Proceedings of the IEEE International Symposium on Information Theory (ISIT)*, Taipei, Taiwan, Jun. 2023.
- [11] X. Zou, S. M. Perlaza, I. Esnaola, and E. Altman, “The worst-case data-generating probability measure,” INRIA, Centre Inria d’Université Côte d’Azur, Sophia Antipolis, France, Tech. Rep. RR-9515, Aug. 2023.
- [12] Y. Bu, G. Aminian, L. Toni, G. W. Wornell, and M. Rodrigues, “Characterizing and understanding the generalization error of transfer learning with

- Gibbs algorithm,” in *Proceedings of the 25th International Conference on Artificial Intelligence and Statistics (AISTATS)*, Virtual Conference, Mar. 2022, pp. 8673–8699.
- [13] H. He, G. Aminian, Y. Bu, M. Rodrigues, and V. Y. Tan, “How does pseudo-labeling affect the generalization error of the semi-supervised Gibbs algorithm?” in *Proceedings of the 26th International Conference on Artificial Intelligence and Statistics (AISTATS)*, Valencia, Spain, Apr. 2023, pp. 8494–8520.
  - [14] F. Hellström, G. Durisi, B. Guedj, and M. Raginsky, “Generalization bounds: Perspectives from information theory and PAC-Bayes,” arXiv preprint arXiv:2309.04381, Sep. 2023.
  - [15] G. Aminian, L. Toni, and M. R. Rodrigues, “Jensen-Shannon information based characterization of the generalization error of learning algorithms,” in *Proceedings of the IEEE Information Theory Workshop (ITW)*, Kanazawa, Japan, Oct. 2021.
  - [16] S. Masiha, A. Gohari, and M. H. Yassaee, “ $f$ -divergences and their applications in lossy compression and bounding generalization error,” *IEEE Transactions on Information Theory*, to appear, arXiv preprint arXiv:2206.11042.
  - [17] P. Alquier, “Non-exponentially weighted aggregation: regret bounds for unbounded loss functions,” in *Proceedings of the 38th International Conference on Machine Learning (ICML)*, vol. 139, Jul. 2021, pp. 207–218.
  - [18] F. Daunas, I. Esnaola, S. M. Perlaza, and H. V. Poor, “Empirical risk minimization with  $f$ -divergence regularization in statistical learning,” INRIA, Centre Inria d’Université Côte d’Azur, Sophia Antipolis, France, Tech. Rep. RR-9521, Oct. 2023.
  - [19] C. P. Robert, G. Casella, and G. Casella, *Monte Carlo statistical methods*, 2nd ed. New York, NY, USA: Springer, 2004.
  - [20] T. Zhang, “Information-theoretic upper and lower bounds for statistical estimation,” *IEEE Transactions on Information Theory*, vol. 52, no. 4, pp. 1307–1321, Apr. 2006.
  - [21] J. N. Kapur, *Maximum Entropy Models in Science and Engineering*, 1st ed. New York, NY, USA: Wiley, 1989.
  - [22] V. De Bortoli and A. Desolneux, “On quantitative Laplace-type convergence results for some exponential probability measures, with two applications,” arXiv preprint arXiv:2110.12922, 2021.
  - [23] K. B. Athreya and C.-R. Hwang, “Gibbs measures asymptotics,” *Sankhyā: The Indian Journal of Statistics, Series A*, vol. 72, no. 1, pp. 191–207, 2010.

- [24] C.-R. Hwang, “Laplace’s method revisited: Weak convergence of probability measures,” *The Annals of Probability*, vol. 8, no. 6, pp. 1177–1182, 1980.
- [25] M. Hasenpflug, D. Rudolf, and B. Sprungk, “Wasserstein convergence rates of increasingly concentrating probability measures,” *arXiv preprint arXiv:2207.08551*, 2022.
- [26] M. Raginsky, A. Rakhlin, M. Tsao, Y. Wu, and A. Xu, “Information-theoretic analysis of stability and bias of learning algorithms,” in *Proceedings of the IEEE Information Theory Workshop (ITW)*, Cambridge, UK, Sep. 2016, pp. 26–30.
- [27] D. Russo and J. Zou, “How much does your data exploration overfit? Controlling bias via information usage,” *IEEE Transactions on Information Theory*, vol. 66, no. 1, pp. 302–323, Jan. 2019.
- [28] A. R. Asadi and E. Abbe, “Chaining meets chain rule: Multilevel entropic regularization and training of neural networks,” *J. Mach. Learn. Res.*, vol. 21, pp. 139–1, Jun. 2020.
- [29] A. Xu and M. Raginsky, “Information-theoretic analysis of generalization capability of learning algorithms,” *Advances in Neural Information Processing Systems*, vol. 30, pp. 1–10, Dec. 2017.
- [30] J. Shawe-Taylor and R. C. Williamson, “A PAC analysis of a Bayesian estimator,” in *Proceedings of the 10th Annual Conference on Computational Learning Theory (COLT)*, Nashville, TN, USA, Jul. 1997, pp. 2–9.
- [31] D. A. McAllester, “PAC-Bayesian stochastic model selection,” *Machine Learning*, vol. 51, no. 1, pp. 5–21, Apr. 2003.
- [32] M. Haddouche, B. Guedj, O. Rivasplata, and J. Shawe-Taylor, “PAC-Bayes unleashed: Generalisation bounds with unbounded losses,” *Entropy*, vol. 23, no. 10, pp. 1–20, Oct. 2021.
- [33] B. Guedj and L. Pujol, “Still no free lunches: The price to pay for tighter PAC-Bayes bounds,” *Entropy*, vol. 23, no. 11, Nov. 2021.
- [34] S. Mazuelas, Y. Shen, and A. Pérez, “Generalized maximum entropy for supervised classification,” *IEEE Transactions on Information Theory*, vol. 68, no. 4, pp. 2530–2550, 2022.
- [35] T. Jaakkola, M. Meila, and T. Jebara, “Maximum entropy discrimination,” *Neural Information Processing Systems*, vol. 12, pp. 470–476, Dec. 1999.
- [36] J. Zhu and E. P. Xing, “Maximum entropy discrimination Markov networks,” *Journal of Machine Learning Research*, vol. 10, no. 11, pp. 2531–2569, Nov. 2009.
- [37] Y. Le Cun, S. Chopra, R. Hadsell, M. Ranzato, and F. Huang, “A tutorial on energy-based learning,” in *Predicting structured data*, 1st ed., G. BakIr,

- T. Hofmann, B. Schölkopf, A. J. Smola, B. Taskar, and S. Vishwanathan, Eds. New York, NY, USA: The MIT Press, 2007, ch. 10, pp. 191–241.
- [38] L. P. Barnes, A. Dytso, and H. V. Poor, “Improved information-theoretic generalization bounds for distributed, federated, and iterative learning,” *Entropy*, vol. 24, no. 9, 2022.
- [39] T. Zhang, “From  $\epsilon$ -entropy to KL-entropy: Analysis of minimum information complexity density estimation,” *The Annals of Statistics*, vol. 34, no. 5, pp. 2180–2210, Oct. 2006.
- [40] J. Jiao, Y. Han, and T. Weissman, “Dependence measures bounding the exploration bias for general measurements,” in *Proceedings of the IEEE International Symposium on Information Theory (ISIT)*, Aachen, Germany, Jun. 2017, pp. 1475–1479.
- [41] H. Wang, M. Diaz, J. C. S. Santos Filho, and F. P. Calmon, “An information-theoretic view of generalization via Wasserstein distance,” in *Proceedings of the IEEE International Symposium on Information Theory (ISIT)*, Paris, France, Jul. 2019, pp. 577–581.
- [42] I. Issa, A. R. Esposito, and M. Gastpar, “Strengthened information-theoretic bounds on the generalization error,” in *Proceedings of the IEEE International Symposium on Information Theory (ISIT)*, Paris, France, Jul. 2019, pp. 582–586.
- [43] Y. Bu, S. Zou, and V. V. Veeravalli, “Tightening mutual information-based bounds on generalization error,” *IEEE Journal on Selected Areas in Information Theory*, vol. 1, no. 1, pp. 121–130, Jan. 2020.
- [44] A. Asadi, E. Abbe, and S. Verdú, “Chaining mutual information and tightening generalization bounds,” *Advances in Neural Information Processing Systems*, pp. 7245–7254, Dec. 2018.
- [45] A. T. Lopez and V. Jog, “Generalization error bounds using Wasserstein distances,” in *Proceedings of the IEEE Information Theory Workshop (ITW)*, Guangzhou, China, Nov. 2018, pp. 1–5.
- [46] H. Hafez-Kolahi, Z. Golgooni, S. Kasaei, and M. Soleymani, “Conditioning and processing: Techniques to improve information-theoretic generalization bounds,” *Advances in Neural Information Processing Systems*, pp. 16 457–16 467, Dec. 2020.
- [47] M. Haghifam, J. Negrea, A. Khisti, D. M. Roy, and G. K. Dziugaite, “Sharpened generalization bounds based on conditional mutual information and an application to noisy, iterative algorithms,” *Advances in Neural Information Processing Systems*, pp. 9925–9935, Dec. 2018.
- [48] B. Rodríguez Gálvez, G. Bassi, R. Thobaben, and M. Skoglund, “Tighter expected generalization error bounds via Wasserstein distance,” *Advances*

- in *Neural Information Processing Systems*, vol. 34, pp. 19 109–19 121, Dec. 2021.
- [49] A. R. Esposito, M. Gastpar, and I. Issa, “Generalization error bounds via Rényi-, f-divergences and maximal leakage,” *IEEE Transactions on Information Theory*, vol. 67, no. 8, pp. 4986–5004, May 2021.
  - [50] G. Aminian, Y. Bu, L. Toni, M. R. D. Rodrigues, and G. W. Wornell, “Information-theoretic characterizations of generalization error for the Gibbs algorithm,” *IEEE Transactions on Information Theory*, to appear, arXiv preprint arXiv:2210.09864.
  - [51] G. Aminian, Y. Bu, G. W. Wornell, and M. R. Rodrigues, “Tighter expected generalization error bounds via convexity of information measures,” in *Proceedings of the IEEE International Symposium on Information Theory (ISIT)*, Aalto, Finland, Jun. 2022, pp. 2481–2486.
  - [52] D. P. Palomar and S. Verdú, “Lautum information,” *IEEE Transactions on Information Theory*, vol. 54, no. 3, pp. 964–975, Mar. 2008.
  - [53] S. Boucheron, G. Lugosi, and P. Massart, *Concentration Inequalities: A Non-asymptotic Theory of Independence*, 1st ed. Oxford, UK: Oxford University Press, 2013.
  - [54] T. M. Cover and J. A. Thomas, *Elements of Information Theory*, 2nd ed. Hoboken, NJ: Wiley-Interscience, 2006.
  - [55] R. B. Ash and C. A. Doleans-Dade, *Probability and Measure Theory*, 2nd ed. Burlington, MA, USA: Harcourt Academic Press, 2000.
  - [56] F. Daunas, I. Esnaola, S. M. Perlaza, and H. V. Poor, “Analysis of the relative entropy asymmetry in the regularization of empirical risk minimization,” in *Proceedings of the IEEE International Symposium on Information Theory (ISIT)*, Taipei, Taiwan, Jun. 2023.
  - [57] H.-O. Georgii, *Gibbs measures and phase transitions*, 2nd ed. New York, NY, USA: De Gruyter, 2011.
  - [58] A. Dembo and O. Zeitouni, *Large Deviations Techniques and Applications*, 2nd ed. New York, NY, USA: Springer-Verlag, 2009.
  - [59] E. T. Jaynes, “Information theory and statistical mechanics I,” *Physical Review Journals*, vol. 106, no. 4, pp. 620–630, May 1957.
  - [60] —, “Information theory and statistical mechanics II,” *Physical Review Journals*, vol. 108, no. 2, pp. 171–190, Oct. 1957.
  - [61] S. Geman and D. Geman, “Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images,” *IEEE Transactions on pattern analysis and machine intelligence*, vol. 11, no. 6, pp. 721–741, Nov. 1984.
  - [62] J. W. Gibbs, *Elementary principles in statistical mechanics*, 1st ed. New Haven, NJ, USA: Yale University Press, 1902.

- [63] O. Catoni, *PAC-Bayesian Supervised Classification: The Thermodynamics of Statistical Learning*, 1st ed. Beachwood, OH, USA: Institute of Mathematical Statistics Lecture Notes - Monograph Series, 2007, vol. 56.
- [64] B. Guedj, “A primer on PAC-Bayesian learning,” in *Proceedings of the International Conference on Machine Learning (ICML)*, Long Beach, CA, USA, May 2019, pp. 2931–2940.
- [65] H. Jeffreys, “An invariant form for the prior probability in estimation problems,” *Proceedings of the Royal Society of London. Series A. Mathematical and Physical Sciences*, vol. 186, no. 1007, pp. 453–461, Sep. 1946.
- [66] W. Rudin, *Principles of mathematical analysis*, 1st ed. New York, NY, USA: McGraw-Hill Book Company, Inc., 1953.
- [67] W. F. Trench, *Introduction to Real Analysis*, 1st ed. Hoboken, NJ, USA: Prentice Hall Pearson Education, 2003.
- [68] S. M. Perlaza, I. Esnaola, G. Bisson, and H. V. Poor, “Sensitivity of the Gibbs algorithm to data aggregation in supervised machine learning,” INRIA, Centre Inria d’Université Côte d’Azur, Sophia Antipolis, France, Tech. Rep. RR-9474, Jun. 2022.
- [69] D. G. Luenberger, *Optimization by Vector Space Methods*, 1st ed. New York, NY, USA: Wiley, 1997.
- [70] W. Feller, *An Introduction to Probability Theory and Its Applications II*, 2nd ed. New York, NY, USA: Jhon Wiley & Sons, 1971.



**RESEARCH CENTRE  
SOPHIA ANTIPOLIS – MÉDITERRANÉE**

2004 route des Lucioles - BP 93  
06902 Sophia Antipolis Cedex

Publisher  
Inria  
Domaine de Voluceau -  
Rocquencourt  
BP 105 - 78153 Le Chesnay  
Cedex  
[inria.fr](http://inria.fr)