



HAL
open science

Empirical Risk Minimization with Generalized 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 Generalized Relative Entropy Regularization. [Research Report] RR-9454, Inria. 2022. hal-03560072v4

HAL Id: hal-03560072

<https://hal.science/hal-03560072v4>

Submitted on 10 Aug 2022 (v4), 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.



Empirical Risk Minimization with Generalized 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 Generalized Relative Entropy Regularization

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

Project-Team NEO

Research Report n° 9454 — version 3 — initial version February 2022 —
revised version July 2022 — 76 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 σ -finite measure instead of a probability measure. This assumption leads to a generalization of the ERM-RER (g-ERM-RER) problem that allows for a larger degree of flexibility in the incorporation of prior knowledge over the set of models. The solution of the g-ERM-RER problem is shown to be a unique probability measure mutually absolutely continuous with the reference measure and to exhibit a probably-approximately-correct (PAC) guarantee for the ERM problem. For a given dataset, the empirical risk is shown to be a sub-Gaussian random variable when the models are sampled from the solution to the g-ERM-RER problem. Finally, the sensitivity of the expected empirical risk to deviations from the solution of the g-ERM-RER problem is studied. In particular, the expectation of the absolute value of sensitivity is shown to be upper bounded, up to a constant factor, by the square root of the lautum information between the models and the datasets.

Key-words: Supervised Learning, PAC-Learning, Regularization, Relative Entropy, Empirical Risk Minimization, Maximum Entropy Principle, Sub-Gaussian Random Variables, Bayesian Learning.

Samir M. Perlaza and Alain Jean-Marie are with INRIA at Sophia Antipolis. Gaetan Bisson and Samir M. Perlaza are with the GAATI Laboratory, University of French Polynesia. Iñaki Esnaola is with the Department of Automatic Control and Systems Engineering, University of Sheffield. Iñaki Esnaola and Samir M. Perlaza are also with the Department of Electrical and Computer Engineering, Princeton University. Stefano Rini is with the Department of Electrical and Computer Engineering, National Yang-ming Chao-tung University (NYCU).

**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 Généralisée

Résumé : Le problème de minimisation du risque empirique (ERM) avec régularisation par l'entropie relative (ERM-RER) est étudié en considérant que la mesure de référence est une mesure σ -finie au lieu d'une mesure de probabilité. La solution de l'ERM-RER s'avère être une mesure de probabilité unique et son expression explicite est présentée en termes de l'ensemble de données donné et du coefficient de régularisation. Pour un ensemble de données fixe, les ensembles négligeables et la concentration de la mesure optimale (solution à l'ERM-RER) sont caractérisés afin de mettre en évidence l'influence du choix de la mesure de référence et du coefficient de régularisation. Les propriétés de la fonction génératrice de cumulants du risque empirique induite par la mesure optimale sont également étudiées. En utilisant ces propriétés, le risque empirique induit par la mesure optimale s'avère être une variable aléatoire sous-gaussienne. La sensibilité de l'espérance du risque empirique aux déviations de la solution du problème ERM-RER est étudiée. Ensuite, la sensibilité est utilisée pour fournir des bornes supérieures et inférieures sur l'espérance du risque empirique. De plus, il est montré que l'espérance de la sensibilité est majorée, à un facteur constant près, par la racine carrée de l'information lautum entre les modèles et l'ensemble de données.

Mots-clés : Apprentissage Supervisé, Apprentissage PAC, Régularisation, Entropie Relative, Minimisation du Risque Empirique, Principe d'Entropie Maximale, Variables Aléatoires sous-Gaussiennes, Apprentissage Bayésien.

Contents

| | | |
|-----------|--|-----------|
| 1 | Introduction | 5 |
| 2 | Empirical Risk Minimization (ERM) | 6 |
| 2.1 | Main Assumptions | 7 |
| 2.2 | ERM with Relative Entropy Regularization | 8 |
| 2.3 | Special Cases of the G-ERM-RER Problem | 8 |
| 3 | The Solution to the G-ERM-RER Problem | 11 |
| 3.1 | Asymptotic Regimes | 13 |
| 3.2 | Coherent and Consistent Measures | 15 |
| 3.3 | Negligible Sets | 18 |
| 4 | The Log-Partition Function | 18 |
| 5 | Expectation of the Empirical Risk | 21 |
| 6 | Variance of the Empirical Risk | 22 |
| 7 | Concentration of Probability | 24 |
| 7.1 | The Limit Set | 26 |
| 7.2 | The Nonnegligible Limit Set | 28 |
| 8 | Sub-Gaussianity of the Empirical Risk | 29 |
| 9 | (δ, ϵ)-Optimality | 31 |
| 10 | Sensitivity | 33 |
| 10.1 | Dataset-Dependent Bounds | 33 |
| 10.2 | Dataset-Independent Bounds | 35 |
| 11 | Discussion and Final Remarks | 37 |
| | Appendices | 38 |
| A | Proof of Lemma 3.1 | 38 |
| B | Proof of Lemma 3.2 | 38 |
| C | Proof of Theorem 3.1 | 39 |
| D | Proof of Lemma 3.3 | 41 |
| E | Proof of Lemma 3.4 | 41 |
| F | Proof of Lemma 3.6 | 42 |
| G | Proof of Lemma 3.7 | 45 |

| | | |
|----------|------------------------------|-----------|
| H | Proof of Lemma 3.8 | 46 |
| I | Proof of Lemma 3.9 | 47 |
| J | Proof of Lemma 3.12 | 48 |
| K | Proof of Lemma 4.1 | 48 |
| L | Proof of Lemma 4.3 | 49 |
| M | Proof of Lemma 4.4 | 50 |
| N | Proof of Theorem 5.1 | 53 |
| | N.1 Preliminaries | 53 |
| | N.2 The proof | 55 |
| O | Proof of Lemma 5.1 | 58 |
| P | Proof of Theorem 5.2 | 59 |
| Q | Proof of Theorem 7.1 | 59 |
| R | Proof of Theorem 7.2 | 61 |
| S | Proof of Lemma 7.2 | 66 |
| T | Proof of Theorem 7.3 | 66 |
| U | Proof of Lemma 8.1 | 67 |
| V | Proof of Lemma 10.1 | 67 |
| W | Proof of Theorem 10.1 | 68 |
| X | Proof of Theorem 10.2 | 70 |

1 Introduction

In supervised 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 [1–3]. Instead of additional statistical assumptions on the datasets, which are typical in Bayesian methods [4], relative entropy regularization requires a reference probability measure, which is external to the ERM problem. Often, such reference measure represents prior knowledge and is chosen to assign high probability to the models that induce low empirical risks. The ERM-RER problem is known to possess a unique solution, which is a Gibbs probability measure. Such Gibbs probability measure has been studied using information theoretic notions in [5–10]; statistical physics in [1]; PAC (Probably Approximately Correct)-Bayesian learning theory in [11–14]; and proved to be of particular interest in classification problems in [15, 16] and supervised learning with energy-based models [17].

In this report, the ERM-RER is generalized to incorporate a σ -finite measure with arbitrary support as the reference measure. Such problem is referred to as the generalized ERM-RER (g-ERM-RER) problem. The flexibility introduced by the g-ERM-RER 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 [18], or cannot be represented by probability distributions, e.g., equal preferences among elements of infinite countable sets; or equal preferences among the elements of the real numbers.

For some specific choices of the reference measure, the g-ERM-RER boils down to particular cases of special interest: (i) the information-risk minimization problem [8]; (ii) the ERM with differential entropy regularization; and (iii) the ERM with discrete entropy regularization [19]. Hence, the proposed formulation yields a unified mathematical framework that comprises a large class of problems. The g-ERM-RER is shown to possess a unique solution, which is mutually absolutely continuous with the reference measure. Such a solution is recognized as a Gibbs probability measure despite the fact that its partition function is defined with respect to a σ -finite reference measure. Given a dataset, it is shown that the empirical risk observed when models are sampled from the g-ERM-RER optimal probability measure is a sub-Gaussian random variable that exhibits a PAC guarantee for the ERM problem without regularization.

The sensitivity of the expected empirical risk to deviations from the solution of the g-ERM-RER problem is studied. The sensitivity is defined as the difference between two quantities: (a) The expectation of the empirical risk with respect to the solution to the measure that is solution to the g-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 g-ERM-RER-optimal measure. More interestingly,

in a special case, the expectation of the absolute value of the sensitivity with respect to the probability distribution of the data sets is shown to be bounded by a term that is proportional to the squared-root of the lautum information [31] between the models and the datasets. This bound is reminiscent to the result in [10] in which, under certain conditions, the generalization gap of certain machine learning algorithms is upper bounded by a term that is proportional to the squared-root of the mutual information between the models and the datasets.

The reminder of this report is organized as follows. Section 2 introduces three optimization problems: the ERM, the ERM-RER, and the g-ERM-RER. Section 3 presents the solution to the g-ERM-RER problem and introduces its main properties. Section 4 studies the properties of the log-partition function of the g-ERM-RER-optimal probability measure. Section 5 and Section 6 study the properties of the expectation and variance of the empirical risk when the models are sampled from the g-ERM-RER-optimal probability measure. Section 7 describes the monotonic concentration of the g-ERM-RER-optimal probability measure when the regularization factor tends to zero. Section 8 and Section 9 respectively show that the empirical risk when the models are sampled from the g-ERM-RER-optimal probability measure is a sub-Gaussian random variable and exhibits a PAC guarantee with respect to the problem without regularization. Finally, Section 10 studies the sensitivity of the empirical risk function with respect to deviations from the g-ERM-RER optimal measure and shows connections with the lautum information. Section 11 ends this work with conclusions and a discussion on the results.

2 Empirical Risk Minimization (ERM)

Consider three sets \mathcal{M} , \mathcal{X} and \mathcal{Y} , with $\mathcal{M} \subseteq \mathbb{R}^d$ and $d \in \mathbb{N}$. Consider also a function $f : \mathcal{M} \times \mathcal{X} \rightarrow \mathcal{Y}$ such that, for some $\theta^* \in \mathcal{M}$, there exist two random variables X and Y that satisfy,

$$Y = f(\theta^*, X). \quad (1)$$

The random variables X and Y jointly form the probability space

$$(\mathcal{X} \times \mathcal{Y}, \mathcal{F}(\mathcal{X} \times \mathcal{Y}), P_{XY}), \quad (2)$$

where $\mathcal{F}(\mathcal{X} \times \mathcal{Y})$ is a σ -algebra on the set $\mathcal{X} \times \mathcal{Y}$, which is assumed to be fixed in this analysis. The elements of the sets \mathcal{M} , \mathcal{X} and \mathcal{Y} are often referred to as *models*, *patterns*, and *labels*, respectively. A pair $(x, y) \in \mathcal{X} \times \mathcal{Y}$ forms a *labeled pattern* or *data point* under the following condition.

Definition 2.1 (Data Point). *The pair (x, y) is said to be a data point if $(x, y) \in \text{supp } P_{XY}$.*

Several data points form a dataset.

Definition 2.2 (Dataset). Given n data points, with $n \in \mathbb{N}$, denoted by $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$, a dataset is represented by the tuple $((x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)) \in (\mathcal{X} \times \mathcal{Y})^n$.

The model θ^* in (1) is the *ground truth model* and is assumed to be unknown. Given a dataset, the objective is to obtain a model $\theta \in \mathcal{M}$, such that, for all patterns $x \in \mathcal{X}$, the assigned label $f(\theta, x)$ minimizes a notion of *loss* or *risk*. Let the function

$$\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow [0, +\infty), \quad (3)$$

be such that given a data point $(x, y) \in \mathcal{X} \times \mathcal{Y}$, the cost, loss or risk induced by choosing the model $\theta \in \mathcal{M}$ is $\ell(f(\theta, x), y)$. Often, the function ℓ is referred to as the *cost function*, *loss function* or *risk function*. In the following, it is assumed that the function ℓ satisfies that, for all $y \in \mathcal{Y}$, the loss $\ell(y, y) = 0$, which implies that correct labelling induces zero cost. Note that there might exist several models $\theta \in \mathcal{M} \setminus \{\theta^*\}$ such that $\ell(f(\theta, x), y) = 0$, which reveals the need of a large number of labeled patterns for model selection.

The *empirical risk* induced by the model θ , with respect to a dataset

$$z = ((x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)) \in (\mathcal{X} \times \mathcal{Y})^n, \quad (4)$$

with $n \in \mathbb{N}$, is determined by the function $L_z : \mathcal{M} \rightarrow [0, +\infty)$, which satisfies

$$L_z(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(f(\theta, x_i), y_i). \quad (5)$$

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

$$\min_{\theta \in \mathcal{M}} L_z(\theta), \quad (6)$$

whose solutions form the set denoted by

$$\mathcal{T}(z) \triangleq \arg \min_{\theta \in \mathcal{M}} L_z(\theta). \quad (7)$$

The ground truth model θ^* in (1) is one of the solutions to the ERM problem in (6). That is, the model θ^* in (1) satisfies that $\theta^* \in \mathcal{T}(z)$ and $L_z(\theta^*) = 0$. Hence, the ERM problem in (6) is well posed.

2.1 Main Assumptions

The *generalized relative entropy* is defined below as the extension to σ -finite measures of the relative entropy usually defined for probability measures.

Definition 2.3 (Generalized Relative Entropy). Given two σ -finite measures P and Q on the same measurable space, such that Q is absolutely continuous with respect to P , the relative entropy of Q with respect to P is

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

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

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

A fundamental assumption in this work is that the function $\bar{\ell} : \mathcal{M} \times \mathcal{X} \times \mathcal{Y} \rightarrow [0, +\infty)$, such that for all $(\boldsymbol{\theta}, x, y) \in \mathcal{M} \times \mathcal{X} \times \mathcal{Y}$,

$$\bar{\ell}(\boldsymbol{\theta}, x, y) = \ell(f(\boldsymbol{\theta}, x), y), \quad (9)$$

where the functions f and ℓ are those in (1) and (3), is Borel measurable with respect to the measure space $(\mathcal{M} \times \mathcal{X} \times \mathcal{Y}, \mathcal{B}(\mathcal{M}) \times \mathcal{F}(\mathcal{X} \times \mathcal{Y}))$, with $\mathcal{F}(\mathcal{X} \times \mathcal{Y})$ the σ -field in (2).

2.2 ERM with Relative Entropy Regularization

When models are chosen by *sampling* from a probability measure over the measurable space $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$, one of the performance metrics is the *expected empirical risk*, which is introduced hereunder.

Definition 2.4 (Expected Empirical Risk). *Given a dataset $\mathbf{z} \in (\mathcal{X} \times \mathcal{Y})^n$, let the function $R_{\mathbf{z}} : \Delta(\mathcal{M}, \mathcal{B}(\mathcal{M})) \rightarrow [0, +\infty)$ be such that for all σ -finite measures $P \in \Delta(\mathcal{M}, \mathcal{B}(\mathcal{M}))$, it holds that*

$$R_{\mathbf{z}}(P) = \int L_{\mathbf{z}}(\boldsymbol{\theta}) dP(\boldsymbol{\theta}), \quad (10)$$

where the function $L_{\mathbf{z}}$ is in (5). Then, when P is a probability measure, the expected empirical risk induced by P is $R_{\mathbf{z}}(P)$.

The g-ERM-RER problem is parametrized by a σ -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 σ -finite measure and let λ be a positive real. The g-ERM-RER problem, with parameters Q and λ , consists of the following optimization problem:

$$\begin{aligned} \min_{P \in \Delta_Q(\mathcal{M}, \mathcal{B}(\mathcal{M}))} & R_{\mathbf{z}}(P) + \lambda D(P||Q), \\ \text{s. t.} & \int dP(\boldsymbol{\theta}) = 1, \end{aligned} \quad (11)$$

where the dataset \mathbf{z} is in (4); and the function $R_{\mathbf{z}}$ is defined in (10).

2.3 Special Cases of the G-ERM-RER Problem

Particular choices of the set \mathcal{M} and the reference measure Q lead to special cases of the g-ERM-RER problem in (11). Three cases are of particular interest:

(a) The set \mathcal{M} is the set \mathbb{R}^d , with $d \in \mathbb{N}$, and Q is the Lebesgue measure on $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$; (b) The set $\mathcal{M} \subset \mathbb{R}^d$ is countable and the measure Q is a counting measure; and (c) The set \mathcal{M} and the measure Q form a probability measure space $(\mathcal{M}, \mathcal{B}(\mathcal{M}), Q)$.

In the former, the g-ERM-RER in (11) satisfies the following

$$\begin{aligned} & \min_{P \in \Delta_Q(\mathcal{M}, \mathcal{B}(\mathcal{M}))} \int \mathbf{L}_z(\boldsymbol{\theta}) dP(\boldsymbol{\theta}) + \lambda D(P \| Q) \\ & \text{s. t.} \quad \int dP(\boldsymbol{\theta}) = 1 \\ = & \min_{P \in \Delta_Q(\mathcal{M}, \mathcal{B}(\mathcal{M}))} \int \mathbf{L}_z(\boldsymbol{\theta}) \frac{dP}{dQ}(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) + \lambda \int \frac{dP}{dQ}(\boldsymbol{\theta}) \log \left(\frac{dP}{dQ}(\boldsymbol{\theta}) \right) dQ(\boldsymbol{\theta}) \\ & \text{s. t.} \quad \int dP(\boldsymbol{\theta}) = 1 \end{aligned} \tag{12}$$

$$= \min_g \int_{\mathcal{M}} \mathbf{L}_z(\boldsymbol{\theta}) g(\boldsymbol{\theta}) d\boldsymbol{\theta} + \lambda \int_{\mathcal{M}} g(\boldsymbol{\theta}) \log(g(\boldsymbol{\theta})) d\boldsymbol{\theta} \tag{13}$$

$$= \min_g \int_{\mathcal{M}} \mathbf{L}_z(\boldsymbol{\theta}) g(\boldsymbol{\theta}) d\boldsymbol{\theta} - \lambda H(g), \tag{14}$$

where the Radon-Nikodym derivative $\frac{dP}{dQ}$ in (12) is a probability density function (p.d.f.) denoted by g , which implies that the optimization domain in (13) is the set of p.d.f.s on \mathbb{R}^d . In (14), the notation $H(g)$ represents the differential entropy of the p.d.f. g , c.f., Chapter 8 in [20]. In this particular case, denote the p.d.f. solution of the problem in (12) by $g_{\Theta|Z=z}^{(\lambda)}$ and thus, for all $\boldsymbol{\theta} \in \mathcal{M}$, it follows that

$$g_{\Theta|Z=z}^{(\lambda)}(\boldsymbol{\theta}) = \frac{\exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\theta})}{\lambda}\right)}{\int_{\mathbb{R}^d} \exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\nu})}{\lambda}\right) d\boldsymbol{\nu}}, \tag{15}$$

which is the Gibbs measure with respect to the Lebesgue measure, with parameter λ and energy function \mathbf{L}_z in (5).

In case (b), a similar analysis would show that the g-ERM-RER problem in (11) boils down to the ERM with discrete entropy regularization (ERM-DisER).

More specifically,

$$\begin{aligned}
& \min_{P \in \Delta_Q(\mathcal{M}, \mathcal{B}(\mathcal{M}))} \int \mathbf{L}_z(\boldsymbol{\theta}) dP(\boldsymbol{\theta}) + \lambda D(P \| Q) \\
& \text{s. t.} \quad \int dP(\boldsymbol{\theta}) = 1 \\
& = \min_{P \in \Delta_Q(\mathcal{M}, \mathcal{B}(\mathcal{M}))} \int \mathbf{L}_z(\boldsymbol{\theta}) \frac{dP}{dQ}(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) + \lambda \int \frac{dP}{dQ}(\boldsymbol{\theta}) \log\left(\frac{dP}{dQ}(\boldsymbol{\theta})\right) dQ(\boldsymbol{\theta}) \\
& \text{s. t.} \quad \int dP(\boldsymbol{\theta}) = 1 \tag{16}
\end{aligned}$$

$$= \min_p \sum_{\boldsymbol{\theta} \in \mathcal{M}} \mathbf{L}_z(\boldsymbol{\theta}) p(\boldsymbol{\theta}) + \lambda \sum_{\boldsymbol{\theta} \in \mathcal{M}} p(\boldsymbol{\theta}) \log(p(\boldsymbol{\theta})) \tag{17}$$

$$= \min_p \sum_{\boldsymbol{\theta} \in \mathcal{M}} \mathbf{L}_z(\boldsymbol{\theta}) p(\boldsymbol{\theta}) - \lambda H(p), \tag{18}$$

where, the Radon-Nikodym derivative $\frac{dQ}{dP}$ in (16) is a probability mass function (p.m.f.) denoted by p . This implies that the optimization domain in (17) is the set of p.m.f.s on \mathcal{M} ; and the entropy $H(p)$ is that of the p.m.f. p , c.f., Chapter 2 in [20]. In this particular case, denote the p.m.f. solution of the problem in (18) by $p_{\Theta|Z=z}^{(\lambda)}$, and thus, for all $\boldsymbol{\theta} \in \mathcal{M}$, it follows that

$$p_{\Theta|Z=z}^{(\lambda)}(\boldsymbol{\theta}) = \frac{\exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\theta})}{\lambda}\right)}{\sum_{\boldsymbol{\nu} \in \mathcal{M}} \exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\nu})}{\lambda}\right)}, \tag{19}$$

which is the Gibbs measure with respect to a counting measure, with parameter λ and energy function \mathbf{L}_z in (5).

Both, the ERM-DifER and ERM-DisER problems are closely related to those typically arising while using Jayne's maximum entropy principle [21, 22] for classification problems such as those in [15, 16, 23].

Finally, case (c) is a classical optimization problem known as the *information-risk minimization* (IRM) problem in information theory, see for instance [8]. In this case, the g-ERM-RER problem in (11) is known to possess a unique solution consisting in a Gibbs probability measure, see for instance [10, 24] and [25]. Denote such unique solution by $P_{\Theta|Z=z}^{(Q,\lambda)} \in \Delta_Q(\mathcal{M}, \mathcal{B}(\mathcal{M}))$, hence, the Radon-Nikodym derivative of $P_{\Theta|Z=z}^{(Q,\lambda)}$ with respect to Q is

$$\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\boldsymbol{\theta}) = \frac{\exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\theta})}{\lambda}\right)}{\int \exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\theta})}{\lambda}\right) dQ(\boldsymbol{\theta})}. \tag{20}$$

The probability measure $P_{\Theta|Z=z}^{(Q,\lambda)}$ in (20) is a Gibbs measure with respect to the probability measure Q , parameter λ , and energy function \mathbf{L}_z in (5), c.f., [26].

3 The Solution to the G-ERM-RER Problem

The solution to the g-ERM-RER problem in (11) is presented in terms of two objects. First, the function $K_{Q,z} : \mathbb{R} \rightarrow \mathbb{R} \cup \{+\infty\}$ such that for all $t \in \mathbb{R}$,

$$K_{Q,z}(t) = \log \left(\int \exp(t L_z(\boldsymbol{\theta})) dQ(\boldsymbol{\theta}) \right), \quad (21)$$

with L_z in (5). Second, the set $\mathcal{K}_{Q,z} \subset (0, +\infty)$, which is defined by

$$\mathcal{K}_{Q,z} \triangleq \left\{ s \in (0, +\infty) : K_{Q,z} \left(-\frac{1}{s} \right) < +\infty \right\}. \quad (22)$$

The notation for the function $K_{Q,z}$ and the set $\mathcal{K}_{Q,z}$ are chosen such that their parametrization by (or dependence on) the dataset \mathbf{z} in (4) and the σ -finite measure Q in (11) are highlighted. The following lemma describes the set $\mathcal{K}_{Q,z}$.

Lemma 3.1. *The set $\mathcal{K}_{Q,z}$ in (22) is either the empty set or a convex set that satisfies*

$$(0, b) \subset \mathcal{K}_{Q,z}, \quad (23)$$

for some $b \in (0, +\infty]$.

Proof: The proof is presented in Appendix A. ■

In the special case in which the σ -finite measure Q in (11) is a probability measure, the set $\mathcal{K}_{Q,z}$ is the set formed by all positive reals, as shown by the following lemma.

Lemma 3.2. *Assume that the measure Q in (11) is a probability measure. Then, the set $\mathcal{K}_{Q,z}$ in (22) satisfies*

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

Proof: The proof is presented in Appendix B. ■

Using this notation, the solution to the g-ERM-RER problem in (11) is presented by the following theorem.

Theorem 3.1. *The solution to the optimization problem in (11) is a unique measure on $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$, denoted by $P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}$, whose Radon-Nikodym derivative with respect to the σ -finite measure Q satisfies for all $\boldsymbol{\theta} \in \text{supp } Q$,*

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

where, the function L_z is defined in (5) and the function $K_{Q,z}$ is defined in (21).

Proof: The proof is presented in Appendix C. ■

In Theorem 3.1, when the σ -finite measure Q is a probability measure, the probability measure $P_{\Theta|Z=z}^{(Q,\lambda)}$ and the function $K_{Q,z}$ are often referred to as a Gibbs measure and a *log-partition function*, see for instance, [27, Section 7.3.1]. In order not to disrupt with the current nomenclature, in the following, independently of whether Q is a probability measure, the function $K_{Q,z}$ and the probability measure $P_{\Theta|Z=z}^{(Q,\lambda)}$ are respectively recognized as a log-partition function and a Gibbs measure with parameters Q , λ , and energy function L_z .

When the set \mathcal{M} is discrete and the σ -finite measure Q in (11) is the counting measure, the Radon-Nikodym derivative $\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}$ in (25) is equivalent to the p.m.f. $p_{\Theta|Z=z}^{(\lambda)}$ in (19). Alternatively, when $\mathcal{M} \subseteq \mathbb{R}^d$ and the σ -finite measure Q in (11) is the Lebesgue measure, the Radon-Nikodym derivative $\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}$ in (25) is the p.d.f. $g_{\Theta|Z=z}^{(\lambda)}$ in (15). When Q is a probability measure, the expression in (20) and (25) are identical.

The following lemma shows that the maximum of the Radon-Nikodym derivative $\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}$ in (25) is achieved by the models that are at the intersection of the support of the reference measure Q and the set of solutions to the optimization problem in (6), that is, the models in the set $\mathcal{T}(z) \cap \text{supp } Q$, with $\mathcal{T}(z)$ in (7).

Lemma 3.3. *For all $\lambda \in \mathcal{K}_{Q,z}$, with $\mathcal{K}_{Q,z}$ in (22), for all $\theta \in \text{supp } Q$, and for all $(\theta_1, \theta_2) \in (\mathcal{T}(z) \cap \text{supp } Q)^2$, with $\mathcal{T}(z)$ in (7), the Radon-Nikodym derivative $\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}$ in (25) satisfies that*

$$\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) \leq \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta_1) = \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta_2). \quad (26)$$

Proof: The proof is presented in Appendix D. ■

When the Radon-Nikodym derivative $\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}$ in (25) is either the p.d.f. $g_{\Theta|Z=z}^{(\lambda)}$ in (15) or the p.m.f. $p_{\Theta|Z=z}^{(\lambda)}$ in (19), Lemma 3.4 shows that the elements of the set $\mathcal{T}(z) \cap \text{supp } Q$ in (7) are the modes of the corresponding p.d.f. or p.m.f.

The following lemma shows other bounds on the Radon-Nikodym derivative in (25).

Lemma 3.4. *For all $\lambda \in \mathcal{K}_{Q,z}$, with $\mathcal{K}_{Q,z}$ in (22), and for all $\theta \in \text{supp } Q$, the Radon-Nikodym derivative $\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}$ in (25) satisfies that*

$$0 \leq \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) < +\infty. \quad (27)$$

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 E. ■

3.1 Asymptotic Regimes

The following lemma describes the asymptotic behavior of the Radon-Nikodym derivative $\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}$ in (25) when $\lambda \rightarrow +\infty$.

Lemma 3.5. *Let the measure Q in (25) 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*

$$\lim_{\lambda \rightarrow +\infty} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) = 1. \tag{28}$$

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

$$\lim_{\lambda \rightarrow +\infty} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) = \lim_{\lambda \rightarrow +\infty} \frac{\exp\left(-\frac{L_z(\theta)}{\lambda}\right)}{\int \exp\left(-\frac{L_z(\nu)}{\lambda}\right) dQ(\nu)} = \frac{1}{\int dQ(\nu)} = 1, \tag{29}$$

where the function L_z is defined in (5). This completes the proof. ■

Lemma 3.5 unveils the fact that, when Q is a probability measure, for all measurable sets $\mathcal{A} \subseteq \text{supp } Q$,

$$\lim_{\lambda \rightarrow +\infty} P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{A}) = Q(\mathcal{A}). \tag{30}$$

This is consistent with the fact that when λ tends to infinity, the optimization problem in (11) 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. Alternatively, from Lemma 3.1, it follows that, when Q is not a probability measure, the set $\mathcal{K}_{Q,z}$ might be an interval of the form $(0, b)$, with $b < \infty$. Hence, in such a case, the analysis in which λ tends to infinity is void.

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

$$\mathcal{L}_z(\delta) = \{\theta \in \mathcal{M} : L_z(\theta) \leq \delta\}, \tag{31}$$

where the function L_z is defined in (5) and $\delta \in [0, +\infty)$. In particular consider the nonnegative real

$$\delta_{Q,z}^* \triangleq \inf \{\delta \in [0, +\infty) : Q(\mathcal{L}_z(\delta)) > 0\}. \tag{32}$$

Let also $\mathcal{L}_{Q,z}^*$ be the following level set of the empirical risk function L_z in (5):

$$\mathcal{L}_{Q,z}^* = \{\boldsymbol{\theta} \in \mathcal{M} : L_z(\boldsymbol{\theta}) = \delta_{Q,z}^*\}. \quad (33)$$

Using this notation, the limit of the Radon-Nikodym derivative $\frac{dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}}{dQ}$ in (25), when λ 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 (33) and Q the σ -finite measure in (25), then for all $\boldsymbol{\theta} \in \text{supp } Q$, the Radon-Nikodym derivative $\frac{dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}}{dQ}$ in (25) satisfies*

$$\lim_{\lambda \rightarrow 0^+} \frac{dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}}{dQ}(\boldsymbol{\theta}) = \frac{1}{Q(\mathcal{L}_{Q,z}^*)} \mathbb{1}_{\{\boldsymbol{\theta} \in \mathcal{L}_{Q,z}^*\}}. \quad (34)$$

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

$$\lim_{\lambda \rightarrow 0^+} \frac{dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}}{dQ}(\boldsymbol{\theta}) = \begin{cases} +\infty & \text{if } \boldsymbol{\theta} \in \mathcal{L}_{Q,z}^* \\ 0 & \text{otherwise.} \end{cases} \quad (35)$$

Proof: The proof is presented in Appendix F. ■

Consider that $Q(\mathcal{L}_{Q,z}^*) > 0$, with $\mathcal{L}_{Q,z}^*$ in (33). Under this assumption, from Lemma 3.6, it holds that the Radon-Nikodym derivative $\frac{dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}}{dQ}$ asymptotically concentrates on the set $\mathcal{L}_{Q,z}^*$ when λ tends to zero from the right. The same observation holds for the probability measure $P_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}$. More specifically, note that for all measurable sets $\mathcal{A} \subseteq \mathcal{L}_{Q,z}^* \cap \text{supp } Q$, it holds that

$$\lim_{\lambda \rightarrow 0^+} P_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}(\mathcal{A}) = \lim_{\lambda \rightarrow 0^+} \int_{\mathcal{A}} \frac{dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}}{dQ}(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) \quad (36)$$

$$= \lim_{\lambda \rightarrow 0^+} \int \frac{dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}}{dQ}(\boldsymbol{\theta}) \mathbb{1}_{\{\boldsymbol{\theta} \in \mathcal{A}\}} dQ(\boldsymbol{\theta}) \quad (37)$$

$$= \int \lim_{\lambda \rightarrow 0^+} \frac{dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}}{dQ}(\boldsymbol{\theta}) \mathbb{1}_{\{\boldsymbol{\theta} \in \mathcal{A}\}} dQ(\boldsymbol{\theta}) \quad (38)$$

$$= \int \frac{1}{Q(\mathcal{L}_{Q,z}^*)} \mathbb{1}_{\{\boldsymbol{\theta} \in \mathcal{L}_{Q,z}^*\}} \mathbb{1}_{\{\boldsymbol{\theta} \in \mathcal{A}\}} dQ(\boldsymbol{\theta}) \quad (39)$$

$$= \frac{1}{Q(\mathcal{L}_{Q,z}^*)} \int \mathbb{1}_{\{\boldsymbol{\theta} \in \mathcal{A}\}} dQ(\boldsymbol{\theta}) \quad (40)$$

$$= \frac{Q(\mathcal{A})}{Q(\mathcal{L}_{Q,z}^*)}, \quad (41)$$

where the equality in (38) follows from Lemma 3.4 and the dominated convergence theorem [28, Theorem 2.6.9]; and the equality in (39) follows from

Lemma 3.6. In the particular case in which $\mathcal{A} = \mathcal{L}_{Q,z}^*$ in (41), 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 equality in (41) implies a uniform distribution of the probability among the elements of the set $\mathcal{L}_{Q,z}^*$ in the limit when λ tends to zero from the right. 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 λ tends to zero from the right is $\frac{1}{|\mathcal{L}_{Q,z}^*|}$.

Alternatively, consider that $Q(\mathcal{L}_{Q,z}^*) = 0$, with $\mathcal{L}_{Q,z}^*$ in (33). Under this assumption, for all $\lambda \in \mathcal{K}_{Q,z}$, with $\mathcal{K}_{Q,z}$ in (22), the probability measure $P_{\Theta|Z=z}^{(Q,\lambda)}$ is absolutely continuous with respect to the measure Q , and thus, $P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_{Q,z}^*) = 0$. This is typically the case in which $\mathcal{M} = \mathbb{R}^d$, the measure Q is absolutely continuous with the Lebesgue measure, and the solution to the ERM problem in (6) has a unique solution, i.e., $|\mathcal{T}(z)| = 1$.

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 λ 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 (33) satisfy,*

$$\lim_{\lambda \rightarrow 0^+} P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_{Q,z}^*) = 1. \tag{42}$$

Proof: The proof is presented in Appendix G. ■

In Lemma 3.7, in the case in which $\delta_{Q,z}^* = 0$, with $\delta_{Q,z}^*$ in (32), it holds that $\mathcal{L}_{Q,z}^* = \mathcal{T}(z)$ and thus, when λ tends to zero from the right, the probability measure $P_{\Theta|Z=z}^{(Q,\lambda)}$ asymptotically concentrates on the set of solutions to the ERM problem in (6). Alternatively, in the case in which $\delta_{Q,z}^* > 0$, it holds that $\mathcal{L}_{Q,z}^* \cap \mathcal{T}(z) = \emptyset$. Hence, when λ 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 (6).

3.2 Coherent and Consistent Measures

Given a $\lambda \in \mathcal{K}_{Q,z}$, with $\mathcal{K}_{Q,z}$ in (22), the support of the probability measure $P_{\Theta|Z=z}^{(Q,\lambda)}$ is of particular interest. In this regard, note that for all sets $\mathcal{C} \in \mathcal{B}(\mathcal{M})$, it follows from Theorem 3.1 that if $Q(\mathcal{C}) = 0$, then $P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{C}) = 0$. That is, the probability measure $P_{\Theta|Z=z}^{(Q,\lambda)}$ is absolutely continuous with the measure Q . The following lemma shows that the converse is also true, under a certain condition.

Lemma 3.8. For all $\lambda \in \mathcal{K}_{Q,z}$, with $\mathcal{K}_{Q,z}$ in (22), under the assumption that

$$Q(\{\boldsymbol{\theta} \in \mathcal{M} : \mathcal{L}_z(\boldsymbol{\theta}) = +\infty\}) = 0, \quad (43)$$

the σ -finite measure Q and the probability measure $P_{\Theta|Z=z}^{(Q,\lambda)}$ in (25) are mutually absolutely continuous.

Proof: The proof is presented in Appendix H. ■

The relevance of Lemma 3.8 is that it proves that for all $\lambda \in \mathcal{K}_{Q,z}$, 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 (11), under the assumption in (43). Under such assumption, for all subsets $\mathcal{C} \in \mathcal{B}(\mathcal{M})$,

$$P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{C}) > 0 \quad \text{if and only if} \quad Q(\mathcal{C}) > 0. \quad (44)$$

The assumption in (43) is trivially true when the function ℓ in (3) is finite.

At the light of Lemma 3.8, a class of reference measures of particular importance in the following sections is that of *coherent* measures.

Definition 3.1 (Coherent Measures). *The σ -finite measure Q in (11) is said to be coherent if, for all $\delta \in (0, +\infty)$, it holds that*

$$Q(\mathcal{L}_z(\delta)) > 0, \quad (45)$$

where the set $\mathcal{L}_z(\delta)$ is defined in (31).

The following lemma highlights the relevance of coherence measures.

Lemma 3.9. *The probability measure $P_{\Theta|Z=z}^{(Q,\lambda)}$ in (25) satisfies for all $\delta \in (0, +\infty)$ that*

$$P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_z(\delta)) > 0, \quad (46)$$

with $\mathcal{L}_z(\delta)$ in (31), if and only if the σ -finite measure Q in (11) is coherent.

Proof: The proof is presented in Appendix I. ■

Reference measures that are coherent also exhibit an interesting property in the asymptotic regime when λ tends to zero from the right, as shown by the following lemma.

Lemma 3.10. *The probability measure $P_{\Theta|Z=z}^{(Q,\lambda)}$ in (25) and the set $\mathcal{T}(z)$ in (7) satisfy*

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

if and only if the σ -finite measure Q in (11) is coherent.

Proof: When the σ -finite measure Q in (11) is coherent, then $\delta_{Q,z}^* = 0$, with $\delta_{Q,z}^*$ in (32) and $\mathcal{L}_{Q,z}^* = \mathcal{T}(z)$, with $\mathcal{L}_{Q,z}^*$ in (33). Thus, from Lemma 3.7, the equality in (47) holds. Alternatively, when the measure Q in (11) is noncoherent, then $\delta_{Q,z}^* > 0$ and thus, $\mathcal{L}_{Q,z}^* \cap \mathcal{T}(z) = \emptyset$, which implies

$$\lim_{\lambda \rightarrow 0^+} P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{T}(z)) = 0, \quad (48)$$

and completes the proof. \blacksquare

Another reference measure Q of particular interest is the one in which for a given $\lambda \in \mathcal{K}_{Q,z}$, with $\mathcal{K}_{Q,z}$ in (22), the set $\mathcal{L}_{Q,z}^*$ in (33) satisfies $Q(\mathcal{L}_{Q,z}^*) > 0$.

Definition 3.2 (Consistent Measure). *The σ -finite measure Q in (11) is said to be consistent if $Q(\mathcal{L}_{Q,z}^*) > 0$, with $\mathcal{L}_{Q,z}^*$ in (33).*

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

Lemma 3.11. *For all $\lambda \in \mathcal{K}_{Q,z}$, with $\mathcal{K}_{Q,z}$ in (22), the probability measure $P_{\Theta|Z=z}^{(Q,\lambda)}$ in (25) and the set $\mathcal{L}_{Q,z}^*$ in (33) satisfy*

$$P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_{Q,z}^*) > 0, \quad (49)$$

if and only if the σ -finite measure Q in (11) 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 Q , it holds that for all $\lambda \in \mathcal{K}_{Q,z}$, $P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_{Q,z}^*) = 0$. 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 3.4, it follows that for all $\lambda \in \mathcal{K}_{Q,z}$, $\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) > 0$. Hence,

$$P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_{Q,z}^*) = \int_{\mathcal{L}_{Q,z}^*} dP_{\Theta|Z=z}^{(Q,\lambda)}(\theta) \quad (50)$$

$$= \int_{\mathcal{L}_{Q,z}^*} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) dQ(\theta) \quad (51)$$

$$> 0, \quad (52)$$

which completes the proof. \blacksquare

The case in which the reference measure Q is simultaneously coherent and consistent exhibits a particular property, as shown by the following corollary of Lemma 3.11.

Corollary 3.1. *For all $\lambda \in \mathcal{K}_{Q,z}$, with $\mathcal{K}_{Q,z}$ in (22), the probability measure $P_{\Theta|Z=z}^{(Q,\lambda)}$ in (25) and the set $\mathcal{T}(z)$ in (7) satisfy*

$$P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{T}(z)) > 0, \quad (53)$$

if and only if the σ -finite measure Q in (11) is coherent and consistent.

3.3 Negligible Sets

Lemma 3.8 shows that, under the assumption in (43), 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 σ -measure Q . 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 the choice of $\lambda \in \mathcal{K}_{Q,z}$.

Lemma 3.12. *For all $(\alpha, \beta) \in \mathcal{K}_{Q,z} \times \mathcal{K}_{Q,z}$, with $\mathcal{K}_{Q,z}$ in (22), assume that the measures $P_{\Theta|Z=z}^{(Q,\lambda)}$ 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 J. ■

4 The Log-Partition Function

The function $K_{Q,z}$ in (25), often referred to as the log-partition function, exhibits a number of properties, which are relevant in the following sections.

Lemma 4.1. *The function $K_{Q,z}$ in (21) is continuous and differentiable infinitely many times in $(-\infty, 0)$.*

Proof: The proof is presented in Appendix K. ■

More specific properties for the function $K_{Q,z}$ in (21) can be stated for the case in which the empirical risk function L_z in (5) is separable with respect to the measure Q in (25).

Definition 4.1 (Separable Empirical Risk Function). *The empirical risk function L_z in (5) is said to be separable with respect to the σ -finite measure Q in (25), 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 Q , such that for all $(\theta_1, \theta_2) \in \mathcal{A} \times \mathcal{B}$,*

$$L_z(\theta_1) < c < L_z(\theta_2) < +\infty. \quad (54)$$

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

$$Q(\{\theta \in \mathcal{M} : L_z(\theta) = a\}) = 1. \quad (55)$$

From this perspective, nonseparable empirical risk functions exhibit little practical interest. This follows from observing that models sampled from a nonseparable probability measure $P_{\Theta|Z=z}^{(Q,\lambda)}$ in (25) induce the same empirical risk.

The definition of separability in Definition 4.1 and Lemma 3.8 lead to the following lemma.

Lemma 4.2. *The empirical risk function L_z in (5) is separable with respect to the σ -finite measure Q in (25) 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 σ -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 (54) holds. Hence, from (54) the following inequalities hold for all $\lambda \in \mathcal{K}_{Q,z}$, with $\mathcal{K}_{Q,z}$ in (22),

$$-\frac{1}{\lambda}L_z(\theta_1) > -\frac{c}{\lambda} > -\frac{1}{\lambda}L_z(\theta_2) > -\infty, \quad (56)$$

$$\exp\left(-\frac{1}{\lambda}L_z(\theta_1)\right) > \exp\left(-\frac{c}{\lambda}\right) > \exp\left(-\frac{1}{\lambda}L_z(\theta_2)\right) > 0, \quad (57)$$

and finally,

$$\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta_1) > \exp\left(-K_{Q,z}\left(-\frac{1}{\lambda}\right) - \frac{c}{\lambda}\right) > \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta_2) > 0. \quad (58)$$

Using the inequality in (58) and the facts that $Q(\mathcal{A}) > 0$ and $Q(\mathcal{B}) > 0$, the following holds

$$P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{A}) = \int_{\mathcal{A}} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) dQ(\theta) > 0, \quad (59)$$

and

$$P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{B}) = \int_{\mathcal{B}} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) dQ(\theta) > 0. \quad (60)$$

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 (54) holds. More specifically,

$$0 < P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{A}) = \int_{\mathcal{A}} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) dQ(\theta) \quad (61)$$

and

$$0 < P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{B}) = \int_{\mathcal{B}} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) dQ(\theta). \quad (62)$$

From Lemma 3.4 and the inequality in (54), 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 (61) and (62),

respectively, it follows that $Q(\mathcal{A}) > 0$ and $Q(\mathcal{B}) > 0$, which implies that the function \mathbb{L}_z is separable with respect to the σ -finite measure Q . This completes the proof. ■

The following lemma presents a general property of the function $K_{Q,z}$ in (21) for the case of separable empirical risk functions.

Lemma 4.3. *The function $K_{Q,z}$ in (21) is convex. The function $K_{Q,z}$ in (21) is strictly convex if and only if the empirical risk function \mathbb{L}_z in (5) is separable with respect to the σ -finite measure Q in (11).*

Proof: The proof is presented in Appendix L. ■

Let the m -th derivative of the function $K_{Q,z}$ in (21) be denoted by $K_{Q,z}^{(m)} : \mathbb{R} \rightarrow \mathbb{R}$, with $m \in \mathbb{N}$. Hence, for all $s \in \mathcal{K}_{Q,z}$,

$$K_{Q,z}^{(m)}\left(-\frac{1}{s}\right) \triangleq \frac{d^m}{dt^m} K_{Q,z}(t) \Big|_{t=-\frac{1}{s}}. \quad (63)$$

The following lemma provides explicit expressions for the first, second and third derivatives of the function $K_{Q,z}$ in (21).

Lemma 4.4. *The first, second and third derivatives of the function $K_{Q,z}$ in (21), 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 (22),*

$$K_{Q,z}^{(1)}\left(-\frac{1}{\lambda}\right) = \int \mathbb{L}_z(\boldsymbol{\theta}) dP_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}(\boldsymbol{\theta}), \quad (64)$$

$$K_{Q,z}^{(2)}\left(-\frac{1}{\lambda}\right) = \int \left(\mathbb{L}_z(\boldsymbol{\theta}) - K_{Q,z}^{(1)}\left(-\frac{1}{\lambda}\right)\right)^2 dP_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}(\boldsymbol{\theta}), \text{ and} \quad (65)$$

$$K_{Q,z}^{(3)}\left(-\frac{1}{\lambda}\right) = \int \left(\mathbb{L}_z(\boldsymbol{\theta}) - K_{Q,z}^{(1)}\left(-\frac{1}{\lambda}\right)\right)^3 dP_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}(\boldsymbol{\theta}), \quad (66)$$

where the function \mathbb{L}_z is defined in (5) and the measure $P_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}$ satisfies (25).

Proof: The proof is presented in Appendix M. ■

From Lemma 4.4, it follows that if $\boldsymbol{\Theta}$ is the random vector that induces the measure $P_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}$ in (25), with $\lambda \in \mathcal{K}_{Q,z}$, the empirical risk function \mathbb{L}_z in (5) becomes the random variable

$$W \triangleq \mathbb{L}_z(\boldsymbol{\Theta}), \quad (67)$$

whose mean, variance, and third cumulant are $K_{Q,z}^{(1)}\left(-\frac{1}{\lambda}\right)$ in (64), $K_{Q,z}^{(2)}\left(-\frac{1}{\lambda}\right)$ in (65), and $K_{Q,z}^{(3)}\left(-\frac{1}{\lambda}\right)$ in (66), respectively.

5 Expectation of the Empirical Risk

The mean of the random variable W in (67) is equivalent to the expectation of the empirical risk function L_z with respect to the probability measure $P_{\Theta|Z=z}^{(Q,\lambda)}$ in (25), which is equal to $R_z(P_{\Theta|Z=z}^{(Q,\lambda)})$, with the function R_z in (10). Often, $R_z(P_{\Theta|Z=z}^{(Q,\lambda)})$ is referred to as the g-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 g-ERM-RER problem in (11). The following corollary of Lemma 4.4 formalizes this observation.

Corollary 5.1. *For all $\lambda \in \mathcal{K}_{Q,z}$, with $\mathcal{K}_{Q,z}$ in (22), the probability measure $P_{\Theta|Z=z}^{(Q,\lambda)}$ in (25) verifies that*

$$R_z(P_{\Theta|Z=z}^{(Q,\lambda)}) = K_{Q,z}^{(1)}\left(-\frac{1}{\lambda}\right), \quad (68)$$

where the functions R_z and $K_{Q,z}^{(1)}$ are defined in (10) and (64), respectively.

The expected empirical risk $R_z(P_{\Theta|Z=z}^{(Q,\lambda)})$ in (68) exhibits the following property.

Theorem 5.1. *The expected empirical risk $R_z(P_{\Theta|Z=z}^{(Q,\lambda)})$ in (68) is nondecreasing with $\lambda \in \mathcal{K}_{Q,z}$, with $\mathcal{K}_{Q,z}$ in (22). Moreover, the function L_z in (5) is separable with respect to the measure Q if and only if $R_z(P_{\Theta|Z=z}^{(Q,\lambda)})$ is strictly increasing with $\lambda \in \mathcal{K}_{Q,z}$.*

Proof: The proof is presented in Appendix N. ■

A question that arises from Theorem 5.1 is whether the value $R_z(P_{\Theta|Z=z}^{(Q,\lambda)})$ in (68) can be made arbitrarily close to $L_z(\theta^*) = 0$, with θ^* in (1), by making λ arbitrarily small. The following lemma shows that there exist cases in which the value $R_z(P_{\Theta|Z=z}^{(Q,\lambda)})$ is bounded away from zero, even for arbitrarily small values of λ .

Lemma 5.1. *For all $\lambda \in \mathcal{K}_{Q,z}$, with $\mathcal{K}_{Q,z}$ in (22), the expected empirical risk $R_z(P_{\Theta|Z=z}^{(Q,\lambda)})$ in (68) satisfies,*

$$R_z(P_{\Theta|Z=z}^{(Q,\lambda)}) \geq \delta_{Q,z}^*, \quad (69)$$

where $\delta_{Q,z}^*$ is defined in (32). Moreover, the function L_z in (5) is separable with respect to the measure Q in (11) if and only if the inequality in (69) is strict.

Proof: The proof is presented in Appendix O. ■

In the asymptotic regime when λ tends to zero, the expected empirical risk $R_z(P_{\Theta|Z=z}^{(Q,\lambda)})$ in (68) is equal to $\delta_{Q,z}^*$, as shown by the following lemma.

Theorem 5.2. *The expected empirical risk $R_z \left(P_{\Theta|Z=z}^{(Q,\lambda)} \right)$ in (68) satisfies,*

$$\lim_{\lambda \rightarrow 0^+} R_z \left(P_{\Theta|Z=z}^{(Q,\lambda)} \right) = \delta_{Q,z}^*, \quad (70)$$

where $\delta_{Q,z}^*$ is defined in (32).

Proof: The proof is presented in Appendix P. ■

6 Variance of the Empirical Risk

The monotonicity of the expectation of the random variable W in (67), stated by Theorem 5.1, is not a property exhibited by the variance nor the third cumulant. This section highlights this observation via the following example.

Example 6.1. *Consider the g-ERM-RER problem in (11), under the assumption that Q is a probability measure and the empirical risk function L_z in (5) is such that for all $\theta \in \mathcal{M}$,*

$$L_z(\theta) = \begin{cases} 0 & \text{if } \theta \in \mathcal{A} \\ 1 & \text{if } \theta \in \mathcal{M} \setminus \mathcal{A}, \end{cases} \quad (71)$$

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 (21) satisfies for all $\lambda > 0$,

$$K_{Q,z} \left(-\frac{1}{\lambda} \right) = \log \left(Q(\mathcal{A}) + \exp \left(-\frac{1}{\lambda} \right) (1 - Q(\mathcal{A})) \right). \quad (72)$$

The derivatives $K_{Q,z}^{(1)}$, $K_{Q,z}^{(2)}$, and $K_{Q,z}^{(3)}$ in (63) of the function $K_{Q,z}$ in (72) satisfy:

$$K_{Q,z}^{(1)} \left(-\frac{1}{\lambda} \right) = \frac{\exp \left(-\frac{1}{\lambda} \right) (1 - Q(\mathcal{A}))}{Q(\mathcal{A}) + \exp \left(-\frac{1}{\lambda} \right) (1 - Q(\mathcal{A}))}; \quad (73)$$

$$K_{Q,z}^{(2)} \left(-\frac{1}{\lambda} \right) = \frac{Q(\mathcal{A}) (1 - Q(\mathcal{A})) \exp \left(-\frac{1}{\lambda} \right)}{\left(Q(\mathcal{A}) + \exp \left(-\frac{1}{\lambda} \right) (1 - Q(\mathcal{A})) \right)^2}; \text{ and} \quad (74)$$

$$K_{Q,z}^{(3)} \left(-\frac{1}{\lambda} \right) = \frac{Q(\mathcal{A}) (1 - Q(\mathcal{A})) \exp \left(-\frac{1}{\lambda} \right) \left(Q(\mathcal{A}) - (1 - Q(\mathcal{A})) \exp \left(-\frac{1}{\lambda} \right) \right)}{\left(Q(\mathcal{A}) + \exp \left(-\frac{1}{\lambda} \right) (1 - Q(\mathcal{A})) \right)^3}. \quad (75)$$

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

$$Q(\mathcal{A}) - (1 - Q(\mathcal{A})) \exp \left(-\frac{1}{\lambda} \right) > 0. \quad (76)$$

Assume that $Q(\mathcal{A}) \geq \frac{1}{2}$. Thus, it holds that for all $\lambda > 0$, the inequality in (76) is always satisfied. This follows from observing that for all $\lambda > 0$,

$$\exp \left(-\frac{1}{\lambda} \right) < 1 \leq \frac{Q(\mathcal{A})}{1 - Q(\mathcal{A})}. \quad (77)$$

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

$$\frac{1}{4} \geq K_{Q,z}^{(2)}\left(-\frac{1}{\lambda_1}\right) > K_{Q,z}^{(2)}\left(-\frac{1}{\lambda_2}\right) > \dots > 0. \quad (78)$$

Alternatively, assume that $Q(\mathcal{A}) < \frac{1}{2}$. In this case, the inequality in (76) is satisfied if and only if

$$\lambda < \left(\log\left(\frac{1-Q(\mathcal{A})}{Q(\mathcal{A})}\right)\right)^{-1}. \quad (79)$$

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

$$\frac{1}{4} > K_{Q,z}^{(2)}\left(-\frac{1}{\lambda_1}\right) > K_{Q,z}^{(2)}\left(-\frac{1}{\lambda_2}\right) > \dots > 0. \quad (80)$$

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

$$K_{Q,z}^{(2)}\left(-\frac{1}{\lambda_1}\right) < K_{Q,z}^{(2)}\left(-\frac{1}{\lambda_2}\right) < \dots < \frac{1}{4}. \quad (81)$$

The upperbound by $\frac{1}{4}$ in (78), (80) and (81) 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 6.1 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 $\mathcal{T}(z)$ 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 λ 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 (6), the variance is monotone with respect to the parameter λ .

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 λ 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 λ

decreases beyond the value $\left(\log\left(\frac{1-Q(\mathcal{A})}{Q(\mathcal{A})}\right)\right)^{-1}$. Otherwise, reducing λ 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 λ 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 (67). These observations are central to understanding the concentration of probability that occurs when λ decreases, as discussed in the following section.

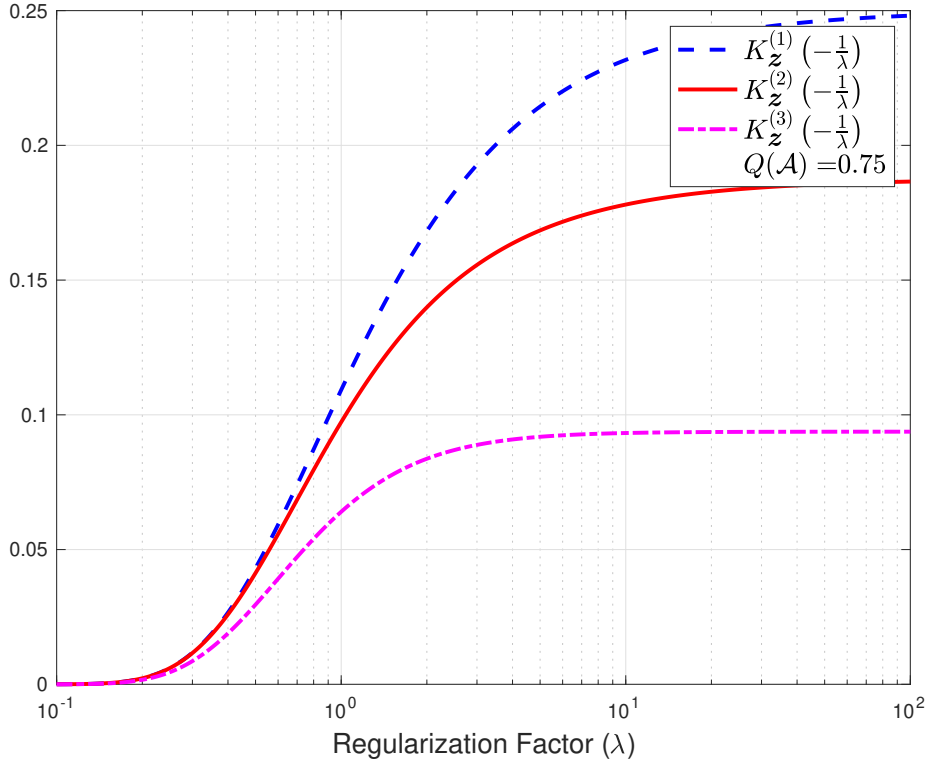


Figure 1: Mean $K_{Q,z}^{(1)}\left(-\frac{1}{\lambda}\right)$, variance $K_{Q,z}^{(2)}\left(-\frac{1}{\lambda}\right)$, and third central moment $K_{Q,z}^{(3)}\left(-\frac{1}{\lambda}\right)$ of the empirical risk in Example 6.1, with $Q(\mathcal{A}) = \frac{3}{4}$

7 Concentration of Probability

Given a positive real $\lambda \in \mathcal{K}_{Q,z}$, with $\mathcal{K}_{Q,z}$ in (22), consider the following set,

$$\mathcal{N}_{Q,z}(\lambda) \triangleq \left\{ \boldsymbol{\theta} \in \mathcal{M} : L_z(\boldsymbol{\theta}) \leq R_z\left(P_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}\right) \right\}, \quad (82)$$

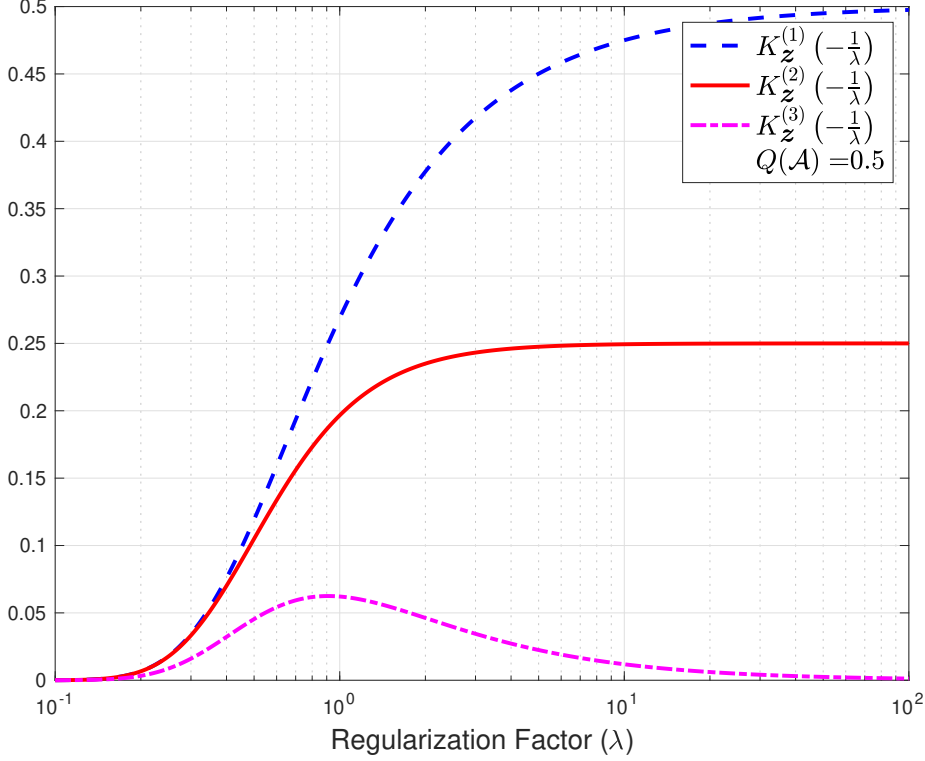


Figure 2: Mean $K_{Q,z}^{(1)}(-\frac{1}{\lambda})$, variance $K_{Q,z}^{(2)}(-\frac{1}{\lambda})$, and third central moment $K_{Q,z}^{(3)}(-\frac{1}{\lambda})$ of the empirical risk in Example 6.1, with $Q(\mathcal{A}) = \frac{1}{2}$

where the function L_z is defined by (5); the function R_z is defined by (10); and the probability measure $P_{\Theta|Z=z}^{(Q,\lambda)}$ is in (25). This section introduces two main results. First, in Theorem 7.1, it is shown that when λ tends to zero, the set $\mathcal{N}_{Q,z}(\lambda)$ forms a monotonic sequence of sets that decreases to the set

$$\mathcal{N}_{Q,z}^* \triangleq \mathcal{L}_z(\delta_{Q,z}^*), \quad (83)$$

where, $\delta_{Q,z}^*$ is defined in (32); and the set $\mathcal{L}_z(\cdot)$ is defined in (31). Second, in Theorem 7.2, it is shown that the sequence formed by $P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{N}_{Q,z}(\lambda))$ when λ tends to zero is increasing and monotone. More importantly, in Theorem 7.3, it is shown that the limit of such sequence is equal to one. These observations justify referring to the set $\mathcal{N}_{Q,z}^*$ as the *limit set*. 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 (33) 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 σ -finite measure Q in (25) is coherent, the sets $\mathcal{N}_{Q,z}^*$ and $\mathcal{L}_{Q,z}^*$ are identical.

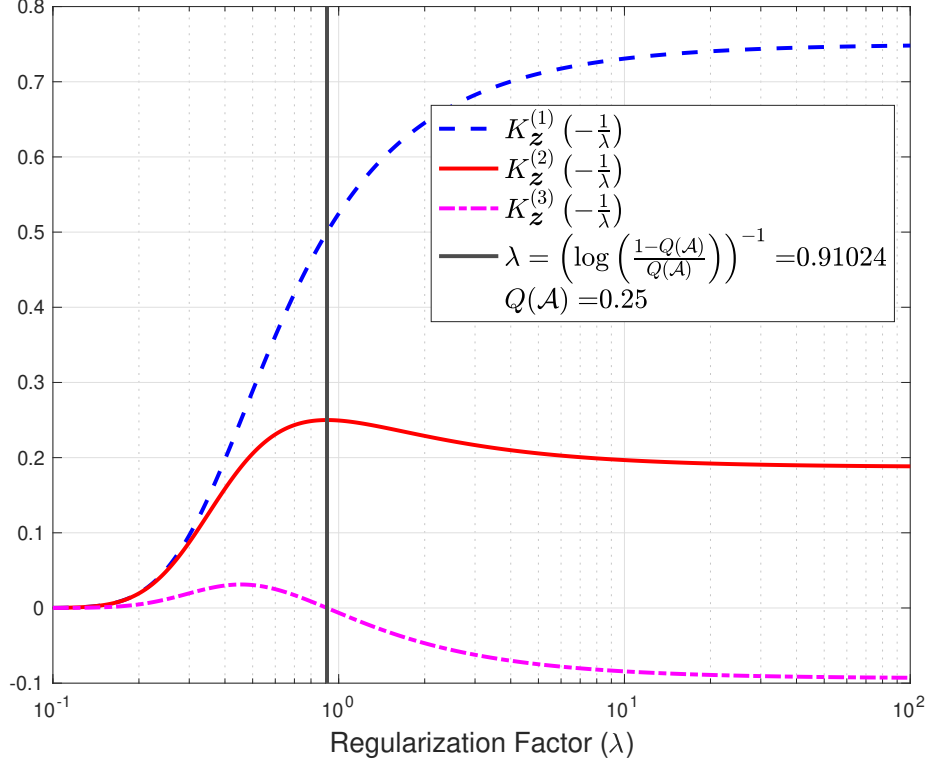


Figure 3: Mean $K_{Q,z}^{(1)}(-\frac{1}{\lambda})$, variance $K_{Q,z}^{(2)}(-\frac{1}{\lambda})$, and third central moment $K_{Q,z}^{(3)}(-\frac{1}{\lambda})$ of the empirical risk in Example 6.1, with $Q(\mathcal{A}) = \frac{1}{4}$

7.1 The Limit Set

The set $\mathcal{N}_{Q,z}(\lambda)$ in (82), with $\lambda \in \mathcal{K}_{Q,z}$ and $\mathcal{K}_{Q,z}$ in (22), contains all the models that induce an empirical risk that is smaller than or equal to $R_z(P_{\Theta|Z=z}^{(Q,\lambda)})$, i.e., the g-ERM-RER-optimal expected empirical risk in (68). This observation unveils the existence of a relation between the set $\mathcal{N}_{Q,z}^*$ in (83) and the set $\mathcal{T}(z)$ in (7), as shown by the following lemma.

Lemma 7.1. *The set $\mathcal{N}_{Q,z}^*$ in (83) satisfies*

$$\mathcal{T}(z) \subseteq \mathcal{N}_{Q,z}^*, \quad (84)$$

where the set $\mathcal{T}(z)$ is in (7). Moreover,

$$\mathcal{T}(z) = \mathcal{N}_{Q,z}^*, \quad (85)$$

if and only if the reference measure Q in (25) is coherent.

Proof: The proof of the inclusion in (84) follows from observing that for all $\theta \in$

$\mathcal{T}(\mathbf{z})$, it holds that $\mathbf{L}_{\mathbf{z}}(\boldsymbol{\theta}) = 0 \leq \delta_{Q,\mathbf{z}}^*$, with $\delta_{Q,\mathbf{z}}^*$ in (32). Hence, $\boldsymbol{\theta} \in \mathcal{N}_{Q,\mathbf{z}}^*$. This completes the proof of the inclusion in (84).

The proof of the second statement is presented in two parts. In the first part, it is proved that if (85) holds, then the measure Q is coherent. The second part proves the converse.

The proof of the first part is as follows. Under the assumption that $\mathcal{T}(\mathbf{z}) = \mathcal{N}_{Q,\mathbf{z}}^*$ holds, it follows that $\mathcal{T}(\mathbf{z}) \supseteq \mathcal{N}_{Q,\mathbf{z}}^*$, which implies that $\delta_{Q,\mathbf{z}}^* = 0$, and thus, for all $\delta \in (0, +\infty)$, it holds that $Q(\mathcal{L}_{\mathbf{z}}(\delta)) > 0$. This 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 measure Q is coherent, it follows that $\delta_{Q,\mathbf{z}}^* = 0$. Hence, $\mathcal{T}(\mathbf{z}) \supseteq \mathcal{N}_{Q,\mathbf{z}}^*$, which together with the inclusion in (84) leads to the equality in (85), which completes the proof of the second part. ■

Lemma 7.1 shows that the limit set $\mathcal{N}_{Q,\mathbf{z}}^*$ in (83) is not empty. This follows from the fact that the set $\mathcal{T}(\mathbf{z})$, which is not empty, is a subset of $\mathcal{N}_{Q,\mathbf{z}}^*$. This observation turns out to be particularly important at the light of the fact that when λ decreases, the set $\mathcal{N}_{Q,\mathbf{z}}(\lambda)$ decreases to the set $\mathcal{N}_{Q,\mathbf{z}}^*$. This observation is formalized by following theorem.

Theorem 7.1. *For all $(\lambda_1, \lambda_2) \in \mathcal{K}_{Q,\mathbf{z}} \times \mathcal{K}_{Q,\mathbf{z}}$, with $\mathcal{K}_{Q,\mathbf{z}}$ in (22) and $\lambda_1 > \lambda_2$, the sets $\mathcal{N}_{Q,\mathbf{z}}(\lambda_1)$ and $\mathcal{N}_{Q,\mathbf{z}}(\lambda_2)$ in (82) satisfy*

$$\mathcal{M} \supseteq \mathcal{N}_{Q,\mathbf{z}}(\lambda_1) \supseteq \mathcal{N}_{Q,\mathbf{z}}(\lambda_2) \supseteq \mathcal{N}_{Q,\mathbf{z}}^*, \quad (86)$$

with $\mathcal{N}_{Q,\mathbf{z}}^*$ the set defined in (83). Moreover, if the empirical risk function $\mathbf{L}_{\mathbf{z}}$ in (5) is continuous on \mathcal{M} and separable with respect to the measure Q in (11), then,

$$\mathcal{M} \supset \mathcal{N}_{Q,\mathbf{z}}(\lambda_1) \supset \mathcal{N}_{Q,\mathbf{z}}(\lambda_2) \supset \mathcal{N}_{Q,\mathbf{z}}^*. \quad (87)$$

Proof: The proof is presented in Appendix Q. ■

An interesting observation is that for all $\lambda \in \mathcal{K}_{Q,\mathbf{z}}$, with $\mathcal{K}_{Q,\mathbf{z}}$ in (22), only a subset of $\mathcal{N}_{Q,\mathbf{z}}(\lambda)$ might exhibit nonzero probability with respect to the measure $P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}$ in (25). Consider for instance that the measure Q in (25) is non-coherent (Definition 3.1). That is, $\delta_{Q,\mathbf{z}}^* > 0$, with $\delta_{Q,\mathbf{z}}^*$ in (32), and thus, for all $\gamma \in [0, \delta_{Q,\mathbf{z}}^*)$, it holds that $Q(\mathcal{L}_{\mathbf{z}}(\gamma)) = 0$, with the set $\mathcal{L}_{\mathbf{z}}(\cdot)$ in (31). Then, for all $\lambda \in \{\alpha \in \mathcal{K}_{Q,\mathbf{z}} : \mathbf{R}_{\mathbf{z}}(P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(Q,\alpha)}) > \delta_{Q,\mathbf{z}}^*\}$ and for all $\gamma < \delta_{Q,\mathbf{z}}^*$, it holds that $\mathcal{L}_{\mathbf{z}}(\gamma) \subseteq \mathcal{N}_{Q,\mathbf{z}}(\lambda)$, while verifying that $Q(\mathcal{L}_{\mathbf{z}}(\gamma)) = 0$, which implies that $P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}(\mathcal{L}_{\mathbf{z}}(\gamma)) = 0$ (Lemma 3.8). Hence, in this case, the set $\mathcal{N}_{Q,\mathbf{z}}(\lambda)$ possesses a subset $\mathcal{L}_{\mathbf{z}}(\gamma)$ that is negligible with respect to the probability measure $P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}$. These observations lead to the analysis of the asymptotic concentration of probability in the following section.

7.2 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 λ tends to zero, as shown by the following theorem.

Theorem 7.2. *For all $(\lambda_1, \lambda_2) \in \mathcal{K}_{Q,z} \times \mathcal{K}_{Q,z}$, with $\mathcal{K}_{Q,z}$ in (22) 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 (82) satisfies*

$$0 < P_{\Theta|Z=z}^{(Q,\lambda_1)}(\mathcal{N}_{Q,z}(\lambda_2)) \leq P_{\Theta|Z=z}^{(Q,\lambda_2)}(\mathcal{N}_{Q,z}(\lambda_2)). \quad (88)$$

Moreover, the function L_z is separable with respect to the σ -finite measure Q , with Q in (25), if and only if for all pairs $(\lambda_1, \lambda_2) \in \mathcal{K}_{Q,z} \times \mathcal{K}_{Q,z}$, with $\lambda_1 > \lambda_2$, it holds that

$$0 < P_{\Theta|Z=z}^{(Q,\lambda_1)}(\mathcal{N}_{Q,z}(\lambda_2)) < P_{\Theta|Z=z}^{(Q,\lambda_2)}(\mathcal{N}_{Q,z}(\lambda_2)). \quad (89)$$

Proof: The proof is presented in Appendix R. ■

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

Lemma 7.2. *Let the function L_z in (5) be separable with respect to the σ -finite measure Q in (25), and consider two positive reals $(\lambda_1, \lambda_2) \in \mathcal{K}_{Q,z} \times \mathcal{K}_{Q,z}$, with $\mathcal{K}_{Q,z}$ in (22) and $\lambda_1 > \lambda_2$. Assume that*

$$Q\left(\mathcal{N}_{Q,z}(\lambda_1) \cap (\mathcal{N}_{Q,z}(\lambda_2))^c\right) = 0. \quad (90)$$

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

$$P_{\Theta|Z=z}^{(Q,\lambda_1)}(\mathcal{N}_{Q,z}(\lambda_1)) < P_{\Theta|Z=z}^{(Q,\lambda_2)}(\mathcal{N}_{Q,z}(\lambda_2)), \quad (91)$$

where, the set $\mathcal{N}_{Q,z}(\cdot)$ is defined in (82).

Proof: The proof is presented in Appendix S. ■

The following example shows the relevance of Lemma 7.2 in the case in which the empirical risk function L_z in (5) is a simple function and separable with respect to the σ -finite measure Q in (25).

Example 7.1. *Consider Example 6.1. Note that, for all $\lambda > 0$,*

$$0 < R_z\left(P_{\Theta|Z=z}^{(Q,\lambda)}\right) < 1, \quad (92)$$

where $R_z\left(P_{\Theta|Z=z}^{(Q,\lambda)}\right)$ is the g -ERM-RER-optimal expected empirical risk in (68). The equality in (92) implies that given two reals λ_1 and λ_2 such that $\lambda_1 > \lambda_2 >$

0, it holds that,

$$\begin{aligned} \mathcal{N}_{Q,z}(\lambda_1) \cap (\mathcal{N}_{Q,z}(\lambda_2))^c &= \{\nu \in \mathcal{M} : R_z(P_{\Theta|Z=z}^{(Q,\lambda_2)}) < L_z(\nu) \leq R_z(P_{\Theta|Z=z}^{(Q,\lambda_1)})\} \\ &= \emptyset, \end{aligned} \quad (93)$$

and moreover, $\mathcal{N}_{Q,z}(\lambda_1) = \mathcal{N}_{Q,z}(\lambda_2)$. Finally, from Lemma 7.2,

$$P_{\Theta|Z=z}^{(Q,\lambda_1)}(\mathcal{N}_{Q,z}(\lambda_1)) < P_{\Theta|Z=z}^{(Q,\lambda_2)}(\mathcal{N}_{Q,z}(\lambda_2)). \quad (94)$$

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

Theorem 7.3. *The probability measure $P_{\Theta|Z=z}^{(Q,\lambda)}$ in (25), with Q being a σ -finite measure on $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$, satisfies*

$$\lim_{\lambda \rightarrow 0^+} P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{N}_{Q,z}(\lambda)) = 1, \quad (95)$$

where, the set $\mathcal{N}_{Q,z}(\lambda)$ is defined in (82).

Proof: The proof is presented in Appendix T. ■

Note that Theorem 7.3 and Lemma 3.7 lead to the following conclusion

$$\lim_{\lambda \rightarrow 0^+} P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{N}_{Q,z}(\lambda) \setminus \mathcal{L}_{Q,z}^*) = 0, \quad (96)$$

which follows from the fact that $\mathcal{L}_{Q,z}^* \subset \mathcal{N}_{Q,z}(\lambda)$, with $\mathcal{L}_{Q,z}^*$ in (33). This justifies referring to the set $\mathcal{L}_{Q,z}^*$ as the nonnegligible limit set.

8 Sub-Gaussianity of the Empirical Risk

Let λ be a real in $\mathcal{K}_{Q,z}$, with $\mathcal{K}_{Q,z}$ in (22), and consider the transport of the measure $P_{\Theta|Z=z}^{(Q,\lambda)}$ in (25) from $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ to $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ through the function L_z in (5). Denote the resulting probability measure in $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ by $P_{W|Z=z}^{(Q,\lambda)}$ in $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$. That is, for all $\mathcal{A} \in \mathcal{B}(\mathbb{R})$,

$$P_{W|Z=z}^{(Q,\lambda)}(\mathcal{A}) = P_{\Theta|Z=z}^{(Q,\lambda)}(L_z^{-1}(\mathcal{A})), \quad (97)$$

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

$$L_z^{-1}(\mathcal{A}) \triangleq \{\nu \in \mathcal{M} : L_z(\nu) \in \mathcal{A}\}. \quad (98)$$

Note that the random variable W in (67) induces the probability measure $P_{W|Z=z}^{(Q,\lambda)}$ in $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$. The objective of this section is to prove that the random variable W is a sub-Gaussian random variable. For this purpose, note that the cumulant generating function induced by the measure $P_{W|Z=z}^{(Q,\lambda)}$, denoted by $J_{z,Q,\lambda}$:

$\mathbb{R} \rightarrow \mathbb{R} \cup \{+\infty\}$, satisfies for all $t \in \mathbb{R}$,

$$J_{\mathbf{z},Q,\lambda}(t) = \log \left(\int \exp(tw) dP_{W|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}(w) \right) \quad (99)$$

$$= \log \left(\int \exp(t\mathbf{L}_{\mathbf{z}}(\mathbf{u})) dP_{\Theta|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}(\mathbf{u}) \right). \quad (100)$$

For all $\lambda \in \mathcal{K}_{Q,\mathbf{z}}$, with $\mathcal{K}_{Q,\mathbf{z}}$ in (22), the following lemma provides an expression for $J_{\mathbf{z},Q,\lambda}(t)$ in terms of the cumulant generating function $K_{Q,\mathbf{z}}$ in (21), for all $t \in (-\infty, \frac{1}{\lambda})$.

Lemma 8.1. *Given a real $\lambda \in \mathcal{K}_{Q,\mathbf{z}}$, with $\mathcal{K}_{Q,\mathbf{z}}$ in (22), the cumulant generating function $J_{\mathbf{z},Q,\lambda}$ in (100), verifies the following equality for all $t \in (-\infty, \frac{1}{\lambda})$,*

$$J_{\mathbf{z},Q,\lambda}(t) = K_{Q,\mathbf{z}} \left(t - \frac{1}{\lambda} \right) - K_{Q,\mathbf{z}} \left(-\frac{1}{\lambda} \right) < +\infty, \quad (101)$$

with the function $K_{Q,\mathbf{z}}$ in (21).

Proof: The proof is presented in Appendix U. ■

Denote by $J_{\mathbf{z},Q,\lambda}^{(m)} : \mathbb{R} \rightarrow \mathbb{R} \cup \{+\infty\}$, with $m \in \mathbb{N}$, the m -th derivative of the function $J_{\mathbf{z},Q,\lambda}$ in (100). That is, for all $s \in (-\infty, \frac{1}{\lambda})$,

$$J_{\mathbf{z},Q,\lambda}^{(m)}(s) = \frac{d^m}{dt^m} J_{\mathbf{z},Q,\lambda}(t) \Big|_{t=s}. \quad (102)$$

From Lemma 8.1, it follows that for all $m \in \mathbb{N}$, and for all $\alpha \in (-\infty, \frac{1}{\lambda})$, the following holds,

$$J_{\mathbf{z},Q,\lambda}^{(m)}(\alpha) = K_{Q,\mathbf{z}}^{(m)} \left(\alpha - \frac{1}{\lambda} \right), \quad (103)$$

where the function $K_{Q,\mathbf{z}}^{(m)}$ denotes the m -th derivative of the function $K_{Q,\mathbf{z}}$ in (21). See for instance, Lemma 4.4. The equality in (103) establishes a relation between the cumulant generating function $J_{\mathbf{z},Q,\lambda}$ and the function $K_{Q,\mathbf{z}}$. The following lemma leverages these observations and presents the main result of this section.

Theorem 8.1. *The cumulant generating function $J_{\mathbf{z},Q,\lambda}$ in (100) verifies the following inequality for all $\alpha \in (-\infty, \frac{1}{\lambda})$,*

$$J_{\mathbf{z},Q,\lambda}(\alpha) \leq \alpha K_{Q,\mathbf{z}}^{(1)} \left(-\frac{1}{\lambda} \right) + \frac{1}{2} \alpha^2 B_{\mathbf{z}}^2, \quad (104)$$

wheret the constant $B_{\mathbf{z}} > 0$ satisfies

$$B_{\mathbf{z}}^2 = \sup_{\gamma \in \mathcal{K}_{Q,\mathbf{z}}} K_{Q,\mathbf{z}}^{(2)} \left(-\frac{1}{\gamma} \right), \quad (105)$$

with $\mathcal{K}_{Q,\mathbf{z}}$ in (22); and the functions $K_{Q,\mathbf{z}}^{(1)}$ and $K_{Q,\mathbf{z}}^{(2)}$ are respectively defined in (64) and (65).

Proof: The function $K_{Q,z}$ in (21) is differentiable infinitely many times over the interior of the set $\mathcal{K}_{Q,z}$ (Lemma 4.1). Thus, from the Taylor-Lagrange theorem, c.f., [29, Theorem 2.5.4], it follows that for all $\lambda \in \mathcal{K}_{Q,z}$ and for all $\alpha \in (-\infty, \frac{1}{\lambda})$, there exists a real $\xi \in (-\infty, \frac{1}{\lambda})$ such that

$$K_{Q,z} \left(\alpha - \frac{1}{\lambda} \right) = K_{Q,z} \left(-\frac{1}{\lambda} \right) + \alpha K_{Q,z}^{(1)} \left(-\frac{1}{\lambda} \right) + \frac{\alpha^2}{2} K_{Q,z}^{(2)}(\xi). \quad (106)$$

From (106) and Lemma 8.1, it holds that

$$J_{z,Q,\lambda}(\alpha) = \alpha K_{Q,z}^{(1)} \left(-\frac{1}{\lambda} \right) + \frac{\alpha^2}{2} K_{Q,z}^{(2)}(\xi). \quad (107)$$

Finally, the inequality in (104) follows from the maximization of the function $K_{Q,z}^{(2)}$ on the set $\mathcal{K}_{Q,z}$, which completes the proof. ■

The main implication of Theorem 8.1 is that the random variable W in (67) is a sub-Gaussian random variable with parameter B_z in (105).

9 (δ, ϵ) -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 (6).

Definition 9.1 ((δ, ϵ) -Optimality). *Given a pair of positive reals (δ, ϵ) , with $\epsilon < 1$, the probability measure $P_{\Theta|Z=z}^{(Q,\lambda)}$ in (25) is said to be (δ, ϵ) -optimal, if the set $\mathcal{L}_z(\delta)$ in (31) satisfies*

$$P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_z(\delta)) > 1 - \epsilon. \quad (108)$$

For all $\delta > 0$, it holds that $\mathcal{T}(z) \subset \mathcal{L}_z(\delta)$, with the sets $\mathcal{T}(z)$ and \mathcal{L}_z in (7) and (31), respectively. Hence, from Definition 9.1, it follows that if the probability measure $P_{\Theta|Z=z}^{(Q,\lambda)}$ in (25) is (δ, ϵ) -optimal, then it assigns a probability that is always greater than $1 - \epsilon$ to the set that contains models that induce an empirical risk that is smaller than δ . From this perspective, particular interest is given to the smallest δ and ϵ for which $P_{\Theta|Z=z}^{(Q,\lambda)}$ is (δ, ϵ) -optimal.

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

Theorem 9.1. *For all $(\delta, \epsilon) \in (\delta_{Q,z}^*, +\infty) \times (0, 1)$, with $\delta_{Q,z}^*$ in (32), there exists a real $\lambda \in \mathcal{K}_{Q,z}$, with $\mathcal{K}_{Q,z}$ in (22), such that the probability measure $P_{\Theta|Z=z}^{(Q,\lambda)}$ is (δ, ϵ) -optimal.*

Proof of Theorem 9.1: Let δ be a real in $(\delta_{Q,z}^*, +\infty)$, with $\delta_{Q,z}^*$ in (32). Let also $\lambda \in \mathcal{K}_{Q,z}$ satisfy the following equality:

$$K_{Q,z}^{(1)} \left(-\frac{1}{\lambda} \right) \leq \delta. \quad (109)$$

Note that from Lemma 4.1, it follows that the function $K_{Q,z}^{(1)}$ is continuous. Moreover, from Theorem 5.2, it follows that such a λ in (109) always exists. From (31) and (82), it holds that

$$\mathcal{N}_{Q,z}(\lambda) \subseteq \mathcal{L}_z(\delta), \quad (110)$$

and thus,

$$P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_z(\delta)) \geq P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{N}_{Q,z}(\lambda)). \quad (111)$$

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

$$P_{\Theta|Z=z}^{(Q,\gamma)}(\mathcal{N}_{Q,z}(\gamma)) > 1 - \epsilon. \quad (112)$$

The existence of such a positive real γ follows from Theorem 7.3. Hence, from (112), it holds that,

$$1 - \epsilon < P_{\Theta|Z=z}^{(Q,\gamma)}(\mathcal{N}_{Q,z}(\gamma)) \quad (113)$$

$$\leq P_{\Theta|Z=z}^{(Q,\gamma)}(\mathcal{L}_z(\delta)), \quad (114)$$

where the inequality in (114) 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 (114) implies that the probability measure $P_{\Theta|Z=z}^{(Q,\lambda)}$ is (δ, ϵ) -optimal (Definition 9.1). This completes the proof. ■

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

Theorem 9.2. *For all $(\delta, \epsilon) \in (0, +\infty) \times (0, 1)$, there always exists a $\lambda \in \mathcal{K}_{Q,z}$, with $\mathcal{K}_{Q,z}$ in (22), such that the probability measure $P_{\Theta|Z=z}^{(Q,\lambda)}$ is (δ, ϵ) -optimal if and only if the reference measure Q is coherent.*

Proof of Theorem 9.2: The proof is divided into two parts. The first part shows that if for all $(\delta, \epsilon) \in (0, +\infty) \times (0, 1)$, there always exists a $\lambda \in \mathcal{K}_{Q,z}$, with $\mathcal{K}_{Q,z}$ in (22), such that the probability measure $P_{\Theta|Z=z}^{(Q,\lambda)}$ in (25) is (δ, ϵ) -optimal, then, the measure Q is coherent. The second part deals with the converse.

The first part is as follows. If for all $(\delta, \epsilon) \in (0, +\infty) \times (0, 1)$, there always exists a $\lambda \in \mathcal{K}_{Q,z}$, such that

$$P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_z(\delta)) > 1 - \epsilon, \quad (115)$$

then, it follows from Lemma 3.8 that for all $\delta > 0$,

$$Q(\mathcal{L}_z(\delta)) > 0, \quad (116)$$

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 $\delta_{Q,z}^* = 0$. Then, from Theorem 9.1, it follows that for all $(\delta, \epsilon) \in (0, +\infty) \times (0, 1)$, there always exists a $\lambda \in \mathcal{K}_{Q,z}$, with $\mathcal{K}_{Q,z}$ in (22), such that the probability measure $P_{\Theta|Z=z}^{(Q,\lambda)}$ is (δ, ϵ) -optimal. This completes the second part of the proof. ■

10 Sensitivity

This section studies the sensitivity of the expected empirical risk R_z (Definition 2.4) to deviations from the probability measure $P_{\Theta|Z=z}^{(Q,\lambda)}$ in (25) towards an alternative probability measure P . Deviations from the probability measure $P_{\Theta|Z=z}^{(Q,\lambda)}$ towards an alternative probability measure P over the measurable space $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ might arise due to several reasons. For instance, if new datasets become available, a new g-ERM-RER problem can be formulated using a larger dataset obtained by aggregating the old and the new datasets [30]. Similarly, the parameters Q (the reference measure) and λ (the regularization factor) in (11) might be changed based on side-information leading to new g-ERM-RER problems and thus, to new probability measures. Other techniques different from g-ERM-RER might also be used to obtain a probability measure over the measurable space $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$, e.g., Bayesian methods. Within this context, the sensitivity is a performance metric defined as follows.

Definition 10.1 (Sensitivity). *Given a σ -finite measure $Q \in \Delta(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ and a positive real $\lambda > 0$, let $S_{Q,\lambda} : (\mathcal{X} \times \mathcal{Y})^n \times \Delta_Q(\mathcal{M}, \mathcal{B}(\mathcal{M})) \rightarrow (-\infty, +\infty]$ be a function such that for all datasets $z \in (\mathcal{X} \times \mathcal{Y})^n$ and for probability measures $P \in \Delta_Q(\mathcal{M}, \mathcal{B}(\mathcal{M}))$, it holds that*

$$S_{Q,\lambda}(z, P) = \begin{cases} R_z(P) - R_z(P_{\Theta|Z=z}^{(Q,\lambda)}) & \text{if } \lambda \in \mathcal{K}_{Q,z} \\ +\infty & \text{otherwise,} \end{cases} \quad (117)$$

where the function R_z is defined in (10) 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 $S_{Q,\lambda}(z, P)$.

10.1 Dataset-Dependent Bounds

In the aim of characterizing the sensitivity $S_{Q,\lambda}(z, P)$ in (117), consider the following lemma.

Lemma 10.1. *Given two probability measures P and Q over $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$, with Q absolutely continuous with P , the following holds for all $z \in (\mathcal{X} \times \mathcal{Y})^n$,*

$$\begin{aligned} & R_z(Q) - R_z(P) \\ & \leq \inf_{t \in (-\infty, 0)} \left(\frac{D(Q||P) + \log \left(\int \exp(t(L_z(\theta) - R_z(P))) dP(\theta) \right)}{t} \right), \end{aligned} \quad (118)$$

where the functions L_z and R_z is defined in (5) and in (10), respectively.

Proof: The proof is presented in Appendix V. ■

Lemma 10.1 together with Theorem 8.1 lead to an upper bound on the absolute value of the sensitivity $S_{Q,\lambda}(z, P)$ in (117).

Theorem 10.1. *The function $S_{Q,\lambda}$ in (117) satisfies for all probability measures $P \in \Delta_Q(\mathcal{M}, \mathcal{B}(\mathcal{M}))$,*

$$|S_{Q,\lambda}(z, P)| \leq \sqrt{2B_{Q,z}^2 D(P \| P_{\Theta|Z=z}^{(Q,\lambda)})}, \quad (119)$$

where the constant $B_{Q,z}$ is defined in (105).

Proof: The proof is presented in Appendix W. ■

Theorem 10.1 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 g-ERM-RER in (11). More specifically, note that for all probability measures $P \in \Delta_Q(\mathcal{M}, \mathcal{B}(\mathcal{M}))$, it holds that,

$$R_z(P) \geq R_z(P_{\Theta|Z=z}^{(Q,\lambda)}) - \sqrt{2B_z^2 D(P \| P_{\Theta|Z=z}^{(Q,\lambda)})} \quad \text{and} \quad (120a)$$

$$R_z(P) \leq R_z(P_{\Theta|Z=z}^{(Q,\lambda)}) + \sqrt{2B_z^2 D(P \| P_{\Theta|Z=z}^{(Q,\lambda)})}. \quad (120b)$$

The following theorem highlights the fact that the measure that minimizes the expected empirical risk subject to a constraint in the relative entropy with respect to the g-ERM-RER-optimal measure $P_{\Theta|Z=z}^{(Q,\lambda)}$ in (25) is also the solution to an g-ERM-RER problem with parameters Q and ω , for some specific $\omega > 0$.

Theorem 10.2. *Given a σ -finite measure $Q \in \Delta(\mathcal{M}, \mathcal{B}(\mathcal{M}))$, a dataset $z \in (\mathcal{X} \times \mathcal{Y})^n$, and a nonnegative real $\lambda \in \mathcal{K}_{Q,z}$, with $\mathcal{K}_{Q,z}$ in (22), consider the following optimization problem*

$$\min_{P \in \Delta_{P_{\Theta|Z=z}^{(Q,\lambda)}}(\mathcal{M}, \mathcal{B}(\mathcal{M}))} \int L_z(\theta) dP(\theta), \quad (121a)$$

$$\text{subject to: } D(P \| P_{\Theta|Z=z}^{(Q,\lambda)}) \leq c, \quad \text{and} \quad (121b)$$

$$\int dP(\theta) = 1, \quad (121c)$$

with, c a given nonnegative constant; $P_{\Theta|Z=z}^{(Q,\lambda)}$ the probability measure in (25); and L_z the function in (5). Then, the solution to the optimization problem in (121) is a probability measure $P_{\Theta|Z=z}^{(Q,\omega)}$ satisfying for all $\theta \in \text{supp } Q$,

$$\frac{dP_{\Theta|Z=z}^{(Q,\omega)}}{dQ}(\theta) = \exp\left(-K_{Q,z}\left(-\frac{1}{\omega}\right) - \frac{1}{\omega}L_z(\theta)\right), \quad (122)$$

with $\omega \in (0, \lambda]$ such that

$$D(P_{\Theta|Z=z}^{(Q,\omega)} \| P_{\Theta|Z=z}^{(Q,\lambda)}) = c. \quad (123)$$

Proof: The proof is presented in Appendix X. ■

10.2 Dataset-Independent Bounds

Consider a probability measure, denoted by $P_{\mathbf{Z}} \in \Delta\left((\mathcal{X} \times \mathcal{Y})^n, (\mathcal{F}(\mathcal{X} \times \mathcal{Y}))^n\right)$, such that for all $\mathcal{A} \in (\mathcal{F}(\mathcal{X} \times \mathcal{Y}))^n$ of the form $\mathcal{A} = \mathcal{A}_1 \times \mathcal{A}_2 \times \dots \times \mathcal{A}_n$ with $\mathcal{A}_i \in \mathcal{F}(\mathcal{X} \times \mathcal{Y})$ and $i \in \{1, 2, \dots, n\}$, it holds that

$$P_{\mathbf{Z}}(\mathcal{A}) = \prod_{t=1}^n P_{XY}(\mathcal{A}_t), \quad (124)$$

where the probability measure P_{XY} is defined in (2). More specifically, $P_{\mathbf{Z}}(\mathcal{A})$ is the probability measure induced by a random variable

$$\mathbf{Z} = ((X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)),$$

in which the n random variables $(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)$ are independent and identically distributed according to P_{XY} .

Let the set \mathcal{K}_Q , with $Q \in \Delta(\mathcal{M}, \mathcal{B}(\mathcal{M}))$, be

$$\mathcal{K}_Q = \bigcap_{z \in \text{supp } P_{\mathbf{Z}}} \mathcal{K}_{Q,z}, \quad (125)$$

where the set $\mathcal{K}_{Q,z}$ is defined in (21) and the probability measure $P_{\mathbf{Z}}$ is defined in (124). The set \mathcal{K}_Q in (125) can be empty for some choices of the σ -finite measure Q and empirical loss function \mathbf{L}_z in (5). Nonetheless, from Lemma 3.2, it follows that when Q is a probability measure, then,

$$\mathcal{K}_Q = (0, +\infty). \quad (126)$$

Using this notation, the following corollary of Theorem 10.1 provides an upper bound on the expectation of the sensitivity with respect to the probability measure $P_{\mathbf{Z}}$ in (124).

Corollary 10.1. *Given a σ -finite measure $Q \in \Delta(\mathcal{M}, \mathcal{B}(\mathcal{M}))$, for all $\lambda \in \mathcal{K}_Q$, with \mathcal{K}_Q in (125), and for all probability measures $P \in \Delta_Q(\mathcal{M}, \mathcal{B}(\mathcal{M}))$, it holds that*

$$\int |\mathbf{S}_{Q,\lambda}(z, P)| \, dP_{\mathbf{Z}}(z) \leq \int \sqrt{2B_{Q,z}^2 D\left(P \| P_{\Theta|Z=z}^{(Q,\lambda)}\right)} \, dP_{\mathbf{Z}}(z), \quad (127)$$

where $B_{Q,z}$ is defined in (105); the probability measure $P_{\Theta|Z=z}^{(Q,\lambda)}$ is in (25); and the probability measure $P_{\mathbf{Z}}$ is defined in (124).

In the following theorem, the expectation of the sensitivity with respect to the measure $P_{\mathbf{Z}}$ in (124) is shown to have an upper bound that can be expressed in terms of the lautum information between the models and the data sets.

Theorem 10.3. *Given a σ -finite measure $Q \in \Delta(\mathcal{M}, \mathcal{B}(\mathcal{M}))$, for all $\lambda \in \mathcal{K}_Q$, with \mathcal{K}_Q in (125), it holds that*

$$\begin{aligned} & \int |S_{Q,\lambda}(z, P_{\Theta}^{(Q,\lambda)})| dP_{\mathbf{Z}}(z) \\ & \leq \sqrt{2B_Q^2} \int D(P_{\Theta}^{(Q,\lambda)} \| P_{\Theta|\mathbf{Z}=\mathbf{u}}^{(Q,\lambda)}) dP_{\mathbf{Z}}(\mathbf{u}), \end{aligned} \quad (128)$$

where the probability measure $P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}$ is the solution to the g -ERM-RER problem in (11); the probability measure $P_{\mathbf{Z}}$ is defined in (124); the probability measure $P_{\Theta}^{(Q,\lambda)}$ is such that for all $A \in \mathcal{B}(\mathcal{M})$,

$$P_{\Theta}^{(Q,\lambda)}(A) = \int P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}(A) dP_{\mathbf{Z}}(\mathbf{z}); \quad (129)$$

and the constant B_Q satisfies

$$B_Q^2 = \sup_{z \in \text{supp } P_{\mathbf{Z}}} B_{Q,z}^2, \quad (130)$$

with $B_{Q,z}$ defined in (105).

Proof: The proof follows from Corollary 10.1. In particular, from (127), for all probability measures P over $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ absolutely continuous with Q , it holds that

$$\int |S_{Q,\lambda}(z, P)| dP_{\mathbf{Z}}(z) \leq \int \sqrt{2B_{Q,z}^2 D(P \| P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)})} dP_{\mathbf{Z}}(z) \quad (131)$$

$$\leq \int \sqrt{2B_Q^2 D(P \| P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)})} dP_{\mathbf{Z}}(z) \quad (132)$$

$$\leq \sqrt{2B_Q^2} \int D(P \| P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}) dP_{\mathbf{Z}}(z), \quad (133)$$

where the inequality in (132) follows from (130); and the inequality in (133) follows from Jensen's inequality [28, Theorem 6.3.5]. This completes the proof. \blacksquare

Given a σ -finite measure $Q \in \Delta(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ and a positive real $\lambda \in \mathcal{K}_Q$, with \mathcal{K}_Q in (125), let \mathbf{Z} and Θ be the random variables that jointly induce a probability measure $P_{\mathbf{Z}\Theta}^{(Q,\lambda)}$ with marginals $P_{\mathbf{Z}}$ in (124) and $P_{\Theta}^{(Q,\lambda)}$ in (129). Under these assumptions, the right-hand side in (128) can be written in terms of the lautum information [31] between the random variables \mathbf{Z} and Θ , which is denoted by $L(\mathbf{Z}; \Theta)$. More specifically, note that

$$L(\mathbf{Z}; \Theta) = \int D(P_{\Theta}^{(Q,\lambda)} \| P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}) dP_{\mathbf{Z}}(\mathbf{z}). \quad (134)$$

In a nutshell, it can be concluded that the expectation of $|S_{Q,\lambda}(z, P_{\Theta}^{(Q,\lambda)})|$ with respect to the measure $P_{\mathbf{Z}}$ in (124) is upper bounded up to a constant factor by the square root of the lautum information between the random variables \mathbf{Z} and Θ .

11 Discussion and Final Remarks

The classical ERM-RER problem in (11) has been studied under the assumption that the reference measure Q is a σ -finite measure, instead of a probability measure, which leads to a more general problem coined as the g-ERM-RER problem. In particular, it has been highlighted that the g-ERM-RER problem provides a larger flexibility for including prior knowledge on the models. Special cases of the g-ERM-RER problem include the ERM problem with (discrete or differential) entropy regularization and the information-risk minimization problem. The solution to the g-ERM-RER problem has been shown to exist and to be unique. Interestingly, the empirical risk observed when models are sampled from the g-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 δ and ϵ , it is shown that there always exist some parameters for the g-ERM-RER-optimal measure such that the set of models that induce an empirical loss smaller than δ exhibit a probability that is not smaller than $1 - \epsilon$. Interestingly, none of these results relies on statistical assumptions on the datasets. Finally, the sensitivity of the expected empirical risk to deviations from the g-ERM-RER-optimal measure has also been studied. In particular, an upper bound on the expectation of the sensitivity with respect to the dataset is presented. In some particular cases, connections between the sensitivity and the lautum information between the models and the data sets have been established.

Appendices

A Proof of Lemma 3.1

The proof is divided into two parts. The first part develops under the assumption that the set $\mathcal{K}_{Q,z} \subset (0, +\infty)$ is not empty. The second part considers the opposite assumption.

The first part is as follows. 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 $K_{Q,z}(-\frac{1}{b}) < +\infty$. Note that for all $\theta \in \mathcal{M}$,

$$\frac{d}{dt} \exp\left(-\frac{1}{t} L_z(\theta)\right) = \frac{1}{t^2} L_z(\theta) \exp\left(-\frac{1}{t} L_z(\theta)\right) \geq 0, \quad (135)$$

with L_z in (5). Thus, from (21), it follows that $K_{Q,z}(-\frac{1}{b})$ is nondecreasing with b . This implies that $(0, b] \subset \mathcal{K}_{Q,z}$. This proves the convexity of $\mathcal{K}_{Q,z}$.

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

$$b^* = \sup \mathcal{K}_{Q,z}. \quad (136)$$

Hence, if $b^* = +\infty$, it follows from (22) that

$$\mathcal{K}_{Q,z} = (0, +\infty). \quad (137)$$

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

$$(0, b^*) \subset \mathcal{K}_{Q,z}. \quad (138)$$

This completes the first part.

The second part is trivial. Under the assumption that the set $\mathcal{K}_{Q,z}$ in (22) is empty, there is nothing to prove.

This completes the proof.

B Proof of Lemma 3.2

Note that for all $\theta \in \mathcal{M}$ and for all for all $t > 0$, it follows that

$$\exp\left(-\frac{1}{t} L_z(\theta)\right) \leq 1, \quad (139)$$

with L_z in (5). Thus,

$$K_{Q,z}\left(-\frac{1}{t}\right) = \log\left(\int \exp\left(-\frac{1}{t} L_z(\theta)\right) dQ(\theta)\right) \quad (140)$$

$$\leq \log\left(\int dQ(\theta)\right) \quad (141)$$

$$\leq 0, \quad (142)$$

which implies that $(0, +\infty) \subseteq \mathcal{K}_{Q,z}$. Thus, from (22), it holds that $\mathcal{K}_{Q,z} = (0, +\infty)$, which completes the proof.

C Proof of Theorem 3.1

The objective function in the optimization problem in (11) can be written as follows:

$$\min_{P \in \Delta_Q(\mathcal{M}, \mathcal{B}(\mathcal{M}))} \int \mathbf{L}_z(\boldsymbol{\theta}) \frac{dP}{dQ}(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) + \lambda \int \frac{dP}{dQ}(\boldsymbol{\theta}) \log \left(\frac{dP}{dQ}(\boldsymbol{\theta}) \right) dQ(\boldsymbol{\theta}) \quad (143a)$$

$$\text{s. t. } \int dP(\boldsymbol{\theta}) = 1. \quad (143b)$$

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

Let \mathcal{M} be the set of nonnegative measurable functions with respect to the measurable spaces $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ and $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$. The Lagrangian of the optimization problem in (143) can be constructed in terms of a function in \mathcal{M} , instead of a measure in $\Delta(\mathcal{M}, \mathcal{B}(\mathcal{M}))$. Let such Lagrangian be $L : \mathcal{M} \times [0, +\infty) \rightarrow \mathbb{R}$ of the form

$$\begin{aligned} L \left(\frac{dP}{dQ}, \beta \right) &= \int \mathbf{L}_z(\boldsymbol{\nu}) \frac{dP}{dQ}(\boldsymbol{\nu}) dQ(\boldsymbol{\nu}) + \lambda \int \frac{dP}{dQ}(\boldsymbol{\nu}) \log \left(\frac{dP}{dQ}(\boldsymbol{\nu}) \right) dQ(\boldsymbol{\nu}) \\ &\quad + \beta \left(\int \frac{dP}{dQ}(\boldsymbol{\nu}) dP(\boldsymbol{\nu}) - 1 \right), \end{aligned} \quad (144)$$

where β is a positive real that acts as a Lagrangian multiplier due to the constraint (143b).

Let $g : \mathbb{R}^k \rightarrow \mathbb{R}$ be a function in \mathcal{M} . The Gateaux differential of the functional L in (144) at $\left(\frac{dP}{dQ}, \beta \right) \in \mathcal{M} \times [0, +\infty)$ in the direction of g is

$$\partial L \left(\frac{dP}{dQ}, \beta; g \right) \triangleq \left. \frac{d}{d\alpha} r(\alpha) \right|_{\alpha=0}, \quad (145)$$

where the real function $r : \mathbb{R} \rightarrow \mathbb{R}$ is such that for all $\alpha \in \mathbb{R}$,

$$\begin{aligned} r(\alpha) &= \int \mathbf{L}_z(\boldsymbol{\nu}) \left(\frac{dP}{dQ}(\boldsymbol{\nu}) + \alpha g(\boldsymbol{\nu}) \right) dQ(\boldsymbol{\nu}) \\ &\quad + \beta \left(\int \left(\frac{dP}{dQ}(\boldsymbol{\nu}) + \alpha g(\boldsymbol{\nu}) \right) dQ(\boldsymbol{\nu}) - 1 \right) \\ &\quad + \lambda \int \left(\frac{dP}{dQ}(\boldsymbol{\nu}) + \alpha g(\boldsymbol{\nu}) \right) \log \left(\frac{dP}{dQ}(\boldsymbol{\nu}) + \alpha g(\boldsymbol{\nu}) \right) dQ(\boldsymbol{\nu}). \end{aligned} \quad (146)$$

Note that the derivative of the real function r in (146) is

$$\begin{aligned} \frac{d}{d\alpha} r(\alpha) &= \int \mathbf{L}_z(\boldsymbol{\nu}) g(\boldsymbol{\nu}) dQ(\boldsymbol{\nu}) + \beta \int g(\boldsymbol{\nu}) dQ(\boldsymbol{\nu}) \\ &\quad + \lambda \int g(\boldsymbol{\nu}) \left(1 + \log \left(\frac{dP}{dQ}(\boldsymbol{\nu}) + \alpha g(\boldsymbol{\nu}) \right) \right) dQ(\boldsymbol{\nu}). \end{aligned} \quad (147)$$

From (145) and (147), it follows that

$$\partial L\left(\frac{dP}{dQ}, \beta; g\right) = \int g(\boldsymbol{\nu}) \left(\mathbf{L}_z(\boldsymbol{\nu}) + \lambda \left(1 + \log\left(\frac{dP}{dQ}(\boldsymbol{\nu})\right) \right) + \beta \right) dQ(\boldsymbol{\nu}). \quad (148)$$

The relevance of the Gateaux differential in (148) stems from [32, Theorem 1, page 178], which unveils the fact that a necessary condition for the functional L in (144) to have a minimum at $\left(\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}, \beta\right) \in \mathcal{M} \times [0, +\infty)$ is that for all functions $g \in \mathcal{M}$,

$$\partial L\left(\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}; g\right) = 0. \quad (149)$$

From (149), it follows that $\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}$ must satisfy for all functions g in \mathcal{M} that

$$\int g(\boldsymbol{\nu}) \left(\mathbf{L}_z(\boldsymbol{\nu}) + \lambda \left(1 + \log\left(\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\boldsymbol{\nu})\right) \right) + \beta \right) dQ(\boldsymbol{\nu}) = 0, \quad (150)$$

which implies that for all $\boldsymbol{\nu} \in \mathcal{M}$,

$$\mathbf{L}_z(\boldsymbol{\nu}) + \lambda \left(1 + \log\left(\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\boldsymbol{\nu})\right) \right) + \beta = 0, \quad (151)$$

and thus,

$$\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\boldsymbol{\nu}) = \exp\left(-\frac{\beta + \lambda}{\lambda}\right) \exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\nu})}{\lambda}\right), \quad (152)$$

with β chosen to satisfy (143b). That is,

$$\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\boldsymbol{\nu}) = \frac{\exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\nu})}{\lambda}\right)}{\int \exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\theta})}{\lambda}\right) dQ(\boldsymbol{\theta})} \quad (153)$$

$$= \exp\left(-K_{Q,z}\left(-\frac{1}{\lambda}\right) - \frac{1}{\lambda}\mathbf{L}_z(\boldsymbol{\nu})\right). \quad (154)$$

The proof continues by verifying that the objective function in (143) is strictly convex, and thus, the measure $P_{\Theta|Z=z}^{(Q,\lambda)}$ that satisfies (153) is the unique minimizer. More specifically, note that the objective function in (143) is the sum of two terms. The first one, i.e., $\int \mathbf{L}_z(\boldsymbol{\nu}) \frac{dP}{dQ}(\boldsymbol{\nu}) dQ(\boldsymbol{\nu})$, is linear in $\frac{dP}{dQ}$. The second, i.e., $\int \frac{dP}{dQ}(\boldsymbol{\nu}) \log\left(\frac{dP}{dQ}(\boldsymbol{\nu})\right) dQ(\boldsymbol{\nu})$, is strictly convex with $\frac{dP}{dQ}$. Hence, given that $\lambda > 0$, the sum of both terms is strictly convex with $\frac{dP}{dQ}$. This implies the uniqueness of $\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}$ and completes the proof.

D Proof of Lemma 3.3

For all $\theta \in \mathcal{M}$ and for all $(\mu, \nu) \in \mathcal{T}(z) \times \mathcal{T}(z)$, it follows that

$$\mathbb{L}_z(\theta) \geq \mathbb{L}_z(\nu) = \mathbb{L}_z(\mu), \quad (155)$$

and thus, for all $\lambda \in \mathcal{K}_{Q,z}$, with $\mathcal{K}_{Q,z}$ in (22), it holds that

$$\exp\left(-\frac{\mathbb{L}_z(\theta)}{\lambda}\right) \leq \exp\left(-\frac{\mathbb{L}_z(\nu)}{\lambda}\right) = \exp\left(-\frac{\mathbb{L}_z(\mu)}{\lambda}\right), \quad (156)$$

which implies

$$\frac{\exp\left(-\frac{\mathbb{L}_z(\theta)}{\lambda}\right)}{\int \exp\left(-\frac{\mathbb{L}_z(\alpha)}{\lambda}\right) dQ(\alpha)} \leq \frac{\exp\left(-\frac{\mathbb{L}_z(\nu)}{\lambda}\right)}{\int \exp\left(-\frac{\mathbb{L}_z(\alpha)}{\lambda}\right) dQ(\alpha)} \quad (157)$$

$$= \frac{\exp\left(-\frac{\mathbb{L}_z(\mu)}{\lambda}\right)}{\int \exp\left(-\frac{\mathbb{L}_z(\alpha)}{\lambda}\right) dQ(\alpha)}. \quad (158)$$

Hence, under the assumption that $\mathcal{T}(z) \cap \text{supp } Q \neq \emptyset$, for all $\theta \in \text{supp } Q$ and for all $(\mu, \nu) \in (\mathcal{T}(z) \cap \text{supp } Q)^2$, it holds that

$$\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}(\theta)}{dQ} \leq \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}(\mu)}{dQ} = \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}(\nu)}{dQ}. \quad (159)$$

which completes the proof.

E Proof of Lemma 3.4

From Lemma 3.3, it follows that for all $\lambda \in \mathcal{K}_{Q,z}$, for all $\theta \in \text{supp } Q$, and for all $\mu \in \mathcal{T}(z) \cap \text{supp } Q$, it holds that

$$\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}(\theta)}{dQ} \leq \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}(\mu)}{dQ} \quad (160)$$

$$= \exp\left(-K_{Q,z}\left(-\frac{1}{\lambda}\right) - \frac{1}{\lambda}\mathbb{L}_z(\mu)\right) \quad (161)$$

$$= \exp\left(-K_{Q,z}\left(-\frac{1}{\lambda}\right)\right) \quad (162)$$

$$< +\infty, \quad (163)$$

where the equality in (161) follows from (25); the equality in (162) follows from the fact that $\mathbb{L}_z(\mu) = 0$; and the equality in (163) 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 for all $\lambda \in \mathcal{K}_{Q,z}$, it holds that $K_{Q,z}\left(-\frac{1}{\lambda}\right) < +\infty$, and thus, $\exp\left(-K_{Q,z}\left(-\frac{1}{\lambda}\right)\right) > 0$. Moreover, for

all $\lambda \in \mathcal{K}_{Q,z}$ and for all $\boldsymbol{\theta} \in \text{supp } Q$, it holds that $L_z(\boldsymbol{\theta}) \leq +\infty$, which implies that $-\frac{1}{\lambda}L_z(\boldsymbol{\theta}) \geq -\infty$, and thus, $\exp(-\frac{1}{\lambda}L_z(\boldsymbol{\theta})) \geq 0$, with equality if and only if $L_z(\boldsymbol{\theta}) = +\infty$. These two observations put together yield

$$\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}(\boldsymbol{\theta})}{dQ} = \exp\left(-K_{Q,z}\left(-\frac{1}{\lambda}\right) - \frac{1}{\lambda}L_z(\boldsymbol{\theta})\right) \quad (164)$$

$$= \exp\left(-K_{Q,z}\left(-\frac{1}{\lambda}\right)\right) \exp\left(-\frac{1}{\lambda}L_z(\boldsymbol{\theta})\right) \quad (165)$$

$$\geq 0, \quad (166)$$

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

F Proof of Lemma 3.6

From Theorem 3.1, it follows that for all $\lambda \in \mathcal{K}_{Q,z}$ and for all $\boldsymbol{\theta} \in \text{supp } Q$,

$$\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}(\boldsymbol{\theta})}{dQ} = \frac{\exp\left(-\frac{L_z(\boldsymbol{\theta})}{\lambda}\right)}{\int \exp\left(-\frac{L_z(\boldsymbol{\nu})}{\lambda}\right) dQ(\boldsymbol{\nu})} \quad (167)$$

$$= \left(\exp\left(\frac{L_z(\boldsymbol{\theta})}{\lambda}\right) \int \exp\left(-\frac{L_z(\boldsymbol{\nu})}{\lambda}\right) dQ(\boldsymbol{\nu})\right)^{-1} \quad (168)$$

$$= \left(\int \exp\left(\frac{1}{\lambda}(L_z(\boldsymbol{\theta}) - L_z(\boldsymbol{\nu}))\right) dQ(\boldsymbol{\nu})\right)^{-1}. \quad (169)$$

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:

$$\mathcal{A}_0(\boldsymbol{\theta}) \triangleq \{\boldsymbol{\nu} \in \text{supp } Q : L_z(\boldsymbol{\theta}) - L_z(\boldsymbol{\nu}) = 0\}, \quad (170a)$$

$$\mathcal{A}_1(\boldsymbol{\theta}) \triangleq \{\boldsymbol{\nu} \in \text{supp } Q : L_z(\boldsymbol{\theta}) - L_z(\boldsymbol{\nu}) < 0\}, \text{ and} \quad (170b)$$

$$\mathcal{A}_2(\boldsymbol{\theta}) \triangleq \{\boldsymbol{\nu} \in \text{supp } Q : L_z(\boldsymbol{\theta}) - L_z(\boldsymbol{\nu}) > 0\}. \quad (170c)$$

Using the sets $\mathcal{A}_0(\boldsymbol{\theta})$, $\mathcal{A}_1(\boldsymbol{\theta})$, and $\mathcal{A}_2(\boldsymbol{\theta})$ in (169), the following holds for all $\lambda \in \mathcal{K}_{Q,z}$ and for all $\boldsymbol{\theta} \in \text{supp } Q$,

$$\begin{aligned} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}(\boldsymbol{\theta})}{dQ} &= \left(\int_{\mathcal{A}_0(\boldsymbol{\theta})} \exp\left(\frac{1}{\lambda}(L_z(\boldsymbol{\theta}) - L_z(\boldsymbol{\nu}))\right) dQ(\boldsymbol{\nu}) \right. \\ &\quad + \int_{\mathcal{A}_1(\boldsymbol{\theta})} \exp\left(\frac{1}{\lambda}(L_z(\boldsymbol{\theta}) - L_z(\boldsymbol{\nu}))\right) dQ(\boldsymbol{\nu}) \\ &\quad \left. + \int_{\mathcal{A}_2(\boldsymbol{\theta})} \exp\left(\frac{1}{\lambda}(L_z(\boldsymbol{\theta}) - L_z(\boldsymbol{\nu}))\right) dQ(\boldsymbol{\nu}) \right)^{-1} \end{aligned} \quad (171)$$

$$\begin{aligned} &= \left(Q(\mathcal{A}_0(\boldsymbol{\theta})) + \int_{\mathcal{A}_1(\boldsymbol{\theta})} \exp\left(\frac{1}{\lambda}(L_z(\boldsymbol{\theta}) - L_z(\boldsymbol{\nu}))\right) dQ(\boldsymbol{\nu}) \right. \\ &\quad \left. + \int_{\mathcal{A}_2(\boldsymbol{\theta})} \exp\left(\frac{1}{\lambda}(L_z(\boldsymbol{\theta}) - L_z(\boldsymbol{\nu}))\right) dQ(\boldsymbol{\nu}) \right)^{-1}. \end{aligned} \quad (172)$$

Note that the sets

$$\{\boldsymbol{\nu} \in \text{supp } Q : \mathbb{L}_z(\boldsymbol{\nu}) = \delta_{Q,z}^*\}, \quad (173)$$

$$\{\boldsymbol{\nu} \in \text{supp } Q : \mathbb{L}_z(\boldsymbol{\nu}) > \delta_{Q,z}^*\}, \text{ and} \quad (174)$$

$$\{\boldsymbol{\nu} \in \text{supp } Q : \mathbb{L}_z(\boldsymbol{\nu}) < \delta_{Q,z}^*\}, \quad (175)$$

with $\delta_{Q,z}^*$ in (32), 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 $\lim_{\lambda \rightarrow 0^+} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\boldsymbol{\theta})$, with $\boldsymbol{\theta} \in \{\boldsymbol{\nu} \in \mathcal{M} : \mathbb{L}_z(\boldsymbol{\nu}) = \delta_{Q,z}^*\}$. The second part considers the case in which $\boldsymbol{\theta} \in \{\boldsymbol{\nu} \in \mathcal{M} : \mathbb{L}_z(\boldsymbol{\nu}) > \delta_{Q,z}^*\}$. The third part considers the remaining case.

The first part is as follows. Consider that $\boldsymbol{\theta} \in \{\boldsymbol{\nu} \in \mathcal{M} : \mathbb{L}_z(\boldsymbol{\nu}) = \delta_{Q,z}^*\}$ and note that $\{\boldsymbol{\nu} \in \mathcal{M} : \mathbb{L}_z(\boldsymbol{\nu}) = \delta_{Q,z}^*\} = \mathcal{L}_{Q,z}^*$. Hence, the sets $\mathcal{A}_0(\boldsymbol{\theta})$, $\mathcal{A}_1(\boldsymbol{\theta})$, and $\mathcal{A}_2(\boldsymbol{\theta})$ in (170) satisfy the following:

$$\mathcal{A}_0(\boldsymbol{\theta}) = \mathcal{L}_{Q,z}^*, \quad (176a)$$

$$\mathcal{A}_1(\boldsymbol{\theta}) = \{\boldsymbol{\mu} \in \text{supp } Q : \mathbb{L}_z(\boldsymbol{\mu}) > \delta_{Q,z}^*\}, \text{ and} \quad (176b)$$

$$\mathcal{A}_2(\boldsymbol{\theta}) = \{\boldsymbol{\mu} \in \text{supp } Q : \mathbb{L}_z(\boldsymbol{\mu}) < \delta_{Q,z}^*\}. \quad (176c)$$

From the definition of $\delta_{Q,z}^*$ in (32), it follows that $Q(\mathcal{A}_2(\boldsymbol{\theta})) = 0$. Plugging the equalities in (176) in (172) yields for all $\boldsymbol{\theta} \in \{\boldsymbol{\nu} \in \mathcal{M} : \mathbb{L}_z(\boldsymbol{\nu}) = \delta_{Q,z}^*\}$,

$$\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\boldsymbol{\theta}) = \left(Q(\mathcal{L}_{Q,z}^*) + \int_{\mathcal{A}_1(\boldsymbol{\theta})} \exp\left(\frac{1}{\lambda}(\mathbb{L}_z(\boldsymbol{\theta}) - \mathbb{L}_z(\boldsymbol{\nu}))\right) dQ(\boldsymbol{\nu}) \right)^{-1} \quad (177)$$

The equality in (177) implies that for all $\boldsymbol{\theta} \in \{\boldsymbol{\nu} \in \mathcal{M} : \mathbb{L}_z(\boldsymbol{\nu}) = \delta_{Q,z}^*\}$,

$$\begin{aligned} \lim_{\lambda \rightarrow 0^+} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\boldsymbol{\theta}) &= \left(Q(\mathcal{L}_{Q,z}^*) \right. \\ &\quad \left. + \lim_{\lambda \rightarrow 0^+} \int_{\mathcal{A}_1(\boldsymbol{\theta})} \exp\left(\frac{1}{\lambda}(\mathbb{L}_z(\boldsymbol{\theta}) - \mathbb{L}_z(\boldsymbol{\nu}))\right) dQ(\boldsymbol{\nu}) \right)^{-1} \quad (178) \\ &= \frac{1}{Q(\mathcal{L}_{Q,z}^*)} \quad (179) \end{aligned}$$

where the equality in (179) follows from verifying that the dominated convergence theorem [28, Theorem 2.6.9] holds. That is,

- (a) For all $\boldsymbol{\nu} \in \mathcal{A}_1(\boldsymbol{\theta})$, it holds that $\exp\left(\frac{1}{\lambda}(\mathbb{L}_z(\boldsymbol{\theta}) - \mathbb{L}_z(\boldsymbol{\nu}))\right) < 1$; and
- (b) For all $\boldsymbol{\nu} \in \mathcal{A}_1(\boldsymbol{\theta})$, it holds that

$$\lim_{\lambda \rightarrow 0^+} \exp\left(\frac{1}{\lambda}(\mathbb{L}_z(\boldsymbol{\theta}) - \mathbb{L}_z(\boldsymbol{\nu}))\right) = 0. \quad (180)$$

This completes the first part of the proof.

The second part is as follows. For all $\delta > \delta_{Q,z}^*$ and for all $\boldsymbol{\theta} \in \{\boldsymbol{\nu} \in \text{supp } Q : \mathbb{L}_z(\boldsymbol{\nu}) = \delta\}$, the sets $\mathcal{A}_0(\boldsymbol{\theta})$, $\mathcal{A}_1(\boldsymbol{\theta})$, and $\mathcal{A}_2(\boldsymbol{\theta})$ in (170) satisfy the following:

$$\mathcal{A}_0(\boldsymbol{\theta}) = \{\boldsymbol{\mu} \in \text{supp } Q : \mathbb{L}_z(\boldsymbol{\mu}) = \delta\}, \quad (181a)$$

$$\mathcal{A}_1(\boldsymbol{\theta}) = \{\boldsymbol{\mu} \in \text{supp } Q : \mathbb{L}_z(\boldsymbol{\mu}) > \delta\}, \text{ and} \quad (181b)$$

$$\mathcal{A}_2(\boldsymbol{\theta}) = \{\boldsymbol{\mu} \in \text{supp } Q : \mathbb{L}_z(\boldsymbol{\mu}) < \delta\}. \quad (181c)$$

Consider the sets

$$\mathcal{A}_{2,1}(\boldsymbol{\theta}) \triangleq \{\boldsymbol{\mu} \in \mathcal{A}_2(\boldsymbol{\theta}) : \mathbb{L}_z(\boldsymbol{\mu}) < \delta_{Q,z}^*\}, \text{ and} \quad (182)$$

$$\mathcal{A}_{2,2}(\boldsymbol{\theta}) \triangleq \{\boldsymbol{\mu} \in \mathcal{A}_2(\boldsymbol{\theta}) : \delta_{Q,z}^* \leq \mathbb{L}_z(\boldsymbol{\mu}) < \delta\}, \quad (183)$$

and note that $\mathcal{A}_{2,1}(\boldsymbol{\theta})$ and $\mathcal{A}_{2,2}(\boldsymbol{\theta})$ form a partition of $\mathcal{A}_2(\boldsymbol{\theta})$. Moreover, from the definition of $\delta_{Q,z}^*$ in (32), it holds that

$$Q(\mathcal{A}_{2,1}(\boldsymbol{\theta})) = 0. \quad (184)$$

Hence, plugging the equalities in (181) and (184) in (172) yields, for all $\delta > \delta_{Q,z}^*$ and for all $\boldsymbol{\theta} \in \{\boldsymbol{\nu} \in \mathcal{M} : \mathbb{L}_z(\boldsymbol{\nu}) = \delta\}$,

$$\begin{aligned} \frac{dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}}{dQ}(\boldsymbol{\theta}) &= \left(Q(\mathcal{A}_0(\boldsymbol{\theta})) + \int_{\mathcal{A}_1(\boldsymbol{\theta})} \exp\left(\frac{1}{\lambda}(\mathbb{L}_z(\boldsymbol{\theta}) - \mathbb{L}_z(\boldsymbol{\nu}))\right) dQ(\boldsymbol{\nu}) \right. \\ &\quad \left. + \int_{\mathcal{A}_{2,2}(\boldsymbol{\theta})} \exp\left(\frac{1}{\lambda}(\mathbb{L}_z(\boldsymbol{\theta}) - \mathbb{L}_z(\boldsymbol{\nu}))\right) dQ(\boldsymbol{\nu}) \right)^{-1}. \end{aligned} \quad (185)$$

The equality in (185) implies that for all $\delta > \delta_{Q,z}^*$ and for all $\boldsymbol{\theta} \in \{\boldsymbol{\nu} \in \mathcal{M} : \mathbb{L}_z(\boldsymbol{\nu}) = \delta\}$,

$$\begin{aligned} \lim_{\lambda \rightarrow 0^+} \frac{dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}}{dQ}(\boldsymbol{\theta}) &= \left(Q(\mathcal{A}_0(\boldsymbol{\theta})) \right. \\ &\quad \left. + \lim_{\lambda \rightarrow 0^+} \int_{\mathcal{A}_1(\boldsymbol{\theta})} \exp\left(\frac{1}{\lambda}(\mathbb{L}_z(\boldsymbol{\theta}) - \mathbb{L}_z(\boldsymbol{\nu}))\right) dQ(\boldsymbol{\nu}) \right. \\ &\quad \left. + \lim_{\lambda \rightarrow 0^+} \int_{\mathcal{A}_{2,2}(\boldsymbol{\theta})} \exp\left(\frac{1}{\lambda}(\mathbb{L}_z(\boldsymbol{\theta}) - \mathbb{L}_z(\boldsymbol{\nu}))\right) dQ(\boldsymbol{\nu}) \right)^{-1} \end{aligned} \quad (186)$$

$$\begin{aligned} &= \left(Q(\mathcal{A}_0(\boldsymbol{\theta})) \right. \\ &\quad \left. + \lim_{\lambda \rightarrow 0^+} \int_{\mathcal{A}_{2,2}(\boldsymbol{\theta})} \exp\left(\frac{1}{\lambda}(\mathbb{L}_z(\boldsymbol{\theta}) - \mathbb{L}_z(\boldsymbol{\nu}))\right) dQ(\boldsymbol{\nu}) \right)^{-1} \end{aligned} \quad (187)$$

$$= (Q(\mathcal{A}_0(\boldsymbol{\theta})) + \infty)^{-1} \quad (188)$$

$$= 0, \quad (189)$$

where the equality in (187) follows by verifying that the dominated convergence theorem [28, 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

$$\lim_{\lambda \rightarrow 0^+} \exp\left(\frac{1}{\lambda}(\mathbb{L}_z(\theta) - \mathbb{L}_z(\nu))\right) = 0. \quad (190)$$

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

$$\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) = 0. \quad (191)$$

This completes the third part of the proof.

Finally, from (179), (189), and (191), it follows that for all $\theta \in \text{supp } Q$,

$$\lim_{\lambda \rightarrow 0^+} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) = \frac{1}{Q(\mathcal{L}_{Q,z}^*)} \mathbb{1}_{\{\theta \in \mathcal{L}_{Q,z}^*\}}, \quad (192)$$

which completes the proof.

G Proof of Lemma 3.7

Consider the following partition of the set \mathcal{M} formed by the sets

$$\mathcal{A}_0 \triangleq \{\theta \in \mathcal{M} : \mathbb{L}_z(\theta) = \delta_{Q,z}^*\}, \quad (193a)$$

$$\mathcal{A}_1 \triangleq \{\theta \in \mathcal{M} : \mathbb{L}_z(\theta) < \delta_{Q,z}^*\}, \text{ and} \quad (193b)$$

$$\mathcal{A}_2 \triangleq \{\theta \in \mathcal{M} : \mathbb{L}_z(\theta) > \delta_{Q,z}^*\}, \quad (193c)$$

with $\delta_{Q,z}^*$ in (32) and the function \mathbb{L}_z in (5). Note that $\mathcal{A}_0 = \mathcal{L}_{Q,z}^*$, with $\mathcal{L}_{Q,z}^*$ in (33).

For all $\lambda \in \mathcal{K}_{Q,z}$, the following holds,

$$1 = P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{A}_0) + P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{A}_1) + P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{A}_2) \quad (194)$$

$$= P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{A}_0) + P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{A}_2) \quad (195)$$

$$= P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{A}_0) + \int_{\mathcal{A}_2} dP_{\Theta|Z=z}^{(Q,\lambda)}(\theta), \quad (196)$$

where, the equality in (195) follows from noticing that $P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{A}_1) = 0$, which follows from the definition of $\delta_{Q,z}^*$ in (32) and the fact that the probability measure $P_{\Theta|Z=z}^{(Q,\lambda)}$ is absolutely continuous with the measure Q .

The above implies that

$$1 = \lim_{\lambda \rightarrow 0^+} P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{A}_0) + \lim_{\lambda \rightarrow 0^+} \int_{\mathcal{A}_2} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) dQ(\theta) \quad (197)$$

$$= \lim_{\lambda \rightarrow 0^+} P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{A}_0) + \int_{\mathcal{A}_2} \lim_{\lambda \rightarrow 0^+} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) dQ(\theta) \quad (198)$$

$$= \lim_{\lambda \rightarrow 0^+} P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{A}_0), \quad (199)$$

where, the equality in (198) follows from the dominated convergence theorem [28, Theorem 1.6.9], given that 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 3.4); and the inequality in (199) 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

$$\lim_{\lambda \rightarrow 0^+} P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_{Q,z}^*) = 1, \quad (200)$$

which completes the proof.

H Proof of Lemma 3.8

For all $\lambda \in \mathcal{K}_{Q,z}$, with $\mathcal{K}_{Q,z}$ in (22), and for all $\mathcal{C} \in \mathcal{B}(\mathcal{M})$,

$$P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{C}) = \int_{\mathcal{C}} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) dQ(\theta), \quad (201)$$

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

$$P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{C}) = 0, \quad (202)$$

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

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

$$0 = P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{C}) \quad (203)$$

$$= \int_{\mathcal{C}} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) dQ(\theta). \quad (204)$$

From Lemma 3.4, and the assumption $Q(\{\boldsymbol{\theta} \in \mathcal{M} : \mathcal{L}_z(\boldsymbol{\theta}) = +\infty\}) = 0$, it holds that for all $\boldsymbol{\theta} \in \text{supp } Q$,

$$\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}(\boldsymbol{\theta})}{dQ}(\boldsymbol{\theta}) > 0, \quad (205)$$

which implies that

$$\int_{\mathcal{C}} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}(\boldsymbol{\theta})}{dQ}(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) = 0, \quad (206)$$

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.

I Proof of Lemma 3.9

The proof is presented in two parts. The first part shows that if for all $\delta \in (0, +\infty)$, the inequality in (46) holds, then, Q is coherent. The second part shows that if there exists a $\delta \in (0, +\infty)$ such that

$$P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_z(\delta)) = 0, \quad (207)$$

then Q is noncoherent.

The first part is as follows. Note that for all $\delta \in (0, +\infty)$ and for all $\boldsymbol{\theta} \in \mathcal{L}_z(\delta) \cap \text{supp } Q$, it holds from Lemma 3.4 that

$$\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}(\boldsymbol{\theta})}{dQ}(\boldsymbol{\theta}) > 0. \quad (208)$$

Hence, if for all $\delta \in (0, +\infty)$, the inequality in (46) holds, then

$$0 < P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_z(\delta)) \quad (209)$$

$$= \int_{\mathcal{L}_z(\delta)} dP_{\Theta|Z=z}^{(Q,\lambda)}(\boldsymbol{\theta}) \quad (210)$$

$$= \int_{\mathcal{L}_z(\delta)} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}(\boldsymbol{\theta})}{dQ}(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}), \quad (211)$$

which, together with (208), implies that for all $\delta \in (0, +\infty)$, $Q(\mathcal{L}_z(\delta)) > 0$. Hence, Q is coherent.

The second part is as follows. Assume that there exists a $\delta \in (0, +\infty)$, for which it holds that $P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_z(\delta)) = 0$. Hence, the following holds for such δ ,

$$0 = P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_z(\delta)) \quad (212)$$

$$= \int_{\mathcal{L}_z(\delta)} dP_{\Theta|Z=z}^{(Q,\lambda)}(\boldsymbol{\theta}) \quad (213)$$

$$= \int_{\mathcal{L}_z(\delta)} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}(\boldsymbol{\theta})}{dQ}(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}). \quad (214)$$

From Lemma 3.4, it holds that for all $\theta \in \mathcal{L}_z(\delta) \cap \text{supp } Q$, the inequality in (208) holds, which implies that $Q(\mathcal{L}_z(\delta)) = 0$. That is, Q is noncoherent. This completes the proof.

J Proof of Lemma 3.12

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

$$g(\theta) = \frac{dP_{\Theta|Z=z}^{(Q,\alpha)}(\theta)}{dQ}(\theta) \left(\frac{dP_{\Theta|Z=z}^{(Q,\beta)}(\theta)}{dQ}(\theta) \right)^{-1}, \quad (215)$$

and note that from Lemma 3.4, it holds that for all $\theta \in \text{supp } Q \setminus \{\nu \in \mathcal{M} : \mathbb{L}_z(\nu) = +\infty\}$, $g(\theta) > 0$ and for all $\theta \in \{\nu \in \mathcal{M} : \mathbb{L}_z(\nu) = +\infty\}$, $g(\theta) = 0$, 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 $\mathcal{A} \in \mathcal{B}(\mathcal{M})$,

$$P(\mathcal{A}) = \int_{\mathcal{A}} g(\theta) dP_{\Theta|Z=z}^{(Q,\beta)}(\theta), \quad (216)$$

and note that if $P_{\Theta|Z=z}^{(Q,\beta)}(\mathcal{A}) = 0$, then $P(\mathcal{A}) = 0$. This implies that P is absolutely continuous with $P_{\Theta|Z=z}^{(Q,\beta)}(\mathcal{A})$. Moreover, from (216), it follows that

$$P(\mathcal{A}) = \int_{\mathcal{A}} \frac{dP_{\Theta|Z=z}^{(Q,\alpha)}(\theta)}{dQ}(\theta) \left(\frac{dP_{\Theta|Z=z}^{(Q,\beta)}(\theta)}{dQ}(\theta) \right)^{-1} dP_{\Theta|Z=z}^{(Q,\beta)}(\theta) \quad (217)$$

$$= \int_{\mathcal{A}} \frac{dP_{\Theta|Z=z}^{(Q,\alpha)}(\theta)}{dQ}(\theta) \left(\frac{dP_{\Theta|Z=z}^{(Q,\beta)}(\theta)}{dQ}(\theta) \right)^{-1} \frac{dP_{\Theta|Z=z}^{(Q,\beta)}(\theta)}{dQ}(\theta) dQ(\theta) \quad (218)$$

$$= \int_{\mathcal{A}} \frac{dP_{\Theta|Z=z}^{(Q,\alpha)}(\theta)}{dQ}(\theta) dQ(\theta) \quad (219)$$

$$= \int_{\mathcal{A}} dP_{\Theta|Z=z}^{(Q,\alpha)}(\theta) \quad (220)$$

$$= P_{\Theta|Z=z}^{(Q,\alpha)}(\mathcal{A}). \quad (221)$$

This proves that $P_{\Theta|Z=z}^{(Q,\alpha)}$ is absolutely continuous with $P_{\Theta|Z=z}^{(Q,\beta)}$. The proof that $P_{\Theta|Z=z}^{(Q,\beta)}$ is absolutely continuous with $P_{\Theta|Z=z}^{(Q,\alpha)}$ follows the same argument. This completes the proof.

K Proof of Lemma 4.1

Consider the transport of the σ -finite measure Q in (21) from the measure space $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ to the measure space $([0, +\infty), \mathcal{B}([0, +\infty)))$ through the

function L_z in (5). 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 $P(\mathcal{A}) = Q(\{\boldsymbol{\theta} \in \mathcal{M} : L_z(\boldsymbol{\theta}) \in \mathcal{A}\})$. Hence, the function $K_{Q,z}$ in (21) can be written for all $t \in \mathbb{R}$ in terms of the measure P as follows

$$K_{Q,z}(t) = \log \left(\int \exp(tv) dP(v) \right). \quad (222)$$

Denote by ϕ the Laplace transform of the measure P . That is, for all $t \in (0, +\infty)$,

$$\phi(t) = \int \exp(tv) dP(v) = \exp(K_{Q,z}(-t)). \quad (223)$$

From [33, Theorem 1a (page 439)], it follows that the function ϕ has derivatives of all orders in $(0, +\infty)$, and thus, so does the function $K_{Q,z}$ in $(-\infty, 0)$. This implies the continuity of $K_{Q,z}$ in $(-\infty, 0)$, and completes the proof.

L Proof of Lemma 4.3

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

$$\begin{aligned} & \alpha K_{Q,z}(\gamma_1) + (1 - \alpha) K_{Q,z}(\gamma_2) \\ &= \alpha \log \left(\int \exp(\gamma_1 L_z(\boldsymbol{\theta})) dQ(\boldsymbol{\theta}) \right) \\ & \quad + (1 - \alpha) \log \left(\int \exp(\gamma_2 L_z(\boldsymbol{\theta})) dQ(\boldsymbol{\theta}) \right) \end{aligned} \quad (224)$$

$$\begin{aligned} &= \log \left(\left(\int \exp(\gamma_1 L_z(\boldsymbol{\theta})) dQ(\boldsymbol{\theta}) \right)^\alpha \right) \\ & \quad + \log \left(\left(\int \exp(\gamma_2 L_z(\boldsymbol{\theta})) dQ(\boldsymbol{\theta}) \right)^{(1-\alpha)} \right) \end{aligned} \quad (225)$$

$$= \log \left(\left(\int \exp(\gamma_1 L_z(\boldsymbol{\theta})) dQ(\boldsymbol{\theta}) \right)^\alpha \left(\int \exp(\gamma_2 L_z(\boldsymbol{\theta})) dQ(\boldsymbol{\theta}) \right)^{(1-\alpha)} \right) \quad (226)$$

$$\begin{aligned} &= \log \left(\left(\int \exp(\gamma_1 \alpha L_z(\boldsymbol{\theta}))^p dQ(\boldsymbol{\theta}) \right)^{\frac{1}{p}} \right. \\ & \quad \left. \left(\int \exp(\gamma_2 (1 - \alpha) L_z(\boldsymbol{\theta}))^q dQ(\boldsymbol{\theta}) \right)^{\frac{1}{q}} \right) \end{aligned} \quad (227)$$

$$\geq \log \left(\int \exp(\gamma_1 \alpha L_z(\boldsymbol{\theta})) \exp(\gamma_2 (1 - \alpha) L_z(\boldsymbol{\theta})) dQ(\boldsymbol{\theta}) \right) \quad (228)$$

$$= \log \left(\int \exp \left((\gamma_1 \alpha + \gamma_2 (1 - \alpha)) L_z(\boldsymbol{\theta}) \right) dQ(\boldsymbol{\theta}) \right) \quad (229)$$

$$= K_{Q,z}(\gamma_1 \alpha + \gamma_2 (1 - \alpha)), \quad (230)$$

where the inequality in (227) follows with $\alpha = \frac{1}{p}$ and $1 - \alpha = \frac{1}{q}$; the inequality in (228) follows from Hölder's inequality. Hence, equality in (228) holds if and

only if there exist two constants β_1 and β_2 , not simultaneously equal to zero, such that the set

$$\mathcal{A} \triangleq \{\boldsymbol{\theta} \in \mathcal{M} : \beta_1 \exp(\gamma_1 \alpha \mathbf{L}_z(\boldsymbol{\theta}))^p = \beta_2 \exp(\gamma_2(1 - \alpha) \mathbf{L}_z(\boldsymbol{\theta}))^q\} \quad (231)$$

$$= \{\boldsymbol{\theta} \in \mathcal{M} : \beta_1 \exp(\gamma_1 \mathbf{L}_z(\boldsymbol{\theta})) = \beta_2 \exp(\gamma_2 \mathbf{L}_z(\boldsymbol{\theta}))\} \quad (232)$$

$$= \left\{ \boldsymbol{\theta} \in \mathcal{M} : \exp((\gamma_1 - \gamma_2) \mathbf{L}_z(\boldsymbol{\theta})) = \frac{\beta_2}{\beta_1} \right\} \quad (233)$$

$$= \left\{ \boldsymbol{\theta} \in \mathcal{M} : \mathbf{L}_z(\boldsymbol{\theta}) = \frac{\log \frac{\beta_2}{\beta_1}}{(\gamma_1 - \gamma_2)} \right\}, \quad (234)$$

satisfies

$$Q(\mathcal{A}) = 1. \quad (235)$$

That is, strict inequality in (228) holds if and only if the function \mathbf{L}_z is separable with respect to the σ -finite measure Q . When $\alpha = 0$ or $\alpha = 1$, the proof is trivial. This completes the proof.

M Proof of Lemma 4.4

For all $s \in \mathcal{K}_{Q,z}$, with $\mathcal{K}_{Q,z}$ in (22), the equality in (63) implies the following,

$$K_{Q,z}^{(1)}\left(-\frac{1}{s}\right) = \frac{d}{dt} \log \left(\int \exp(t \mathbf{L}_z(\boldsymbol{\theta})) dQ(\boldsymbol{\theta}) \right) \Big|_{t=-\frac{1}{s}} \quad (236)$$

$$= \frac{1}{\int \exp(t \mathbf{L}_z(\boldsymbol{v})) dQ(\boldsymbol{v})} \int \mathbf{L}_z(\boldsymbol{\theta}) \exp(t \mathbf{L}_z(\boldsymbol{\theta})) dQ(\boldsymbol{\theta}) \Big|_{t=-\frac{1}{s}} \quad (237)$$

$$= \frac{1}{\int \exp\left(-\frac{1}{s} \mathbf{L}_z(\boldsymbol{v})\right) dQ(\boldsymbol{v})} \int \mathbf{L}_z(\boldsymbol{\theta}) \exp\left(-\frac{1}{s} \mathbf{L}_z(\boldsymbol{\theta})\right) dQ(\boldsymbol{\theta}) \quad (238)$$

$$= \exp\left(-K_{Q,z}\left(-\frac{1}{s}\right)\right) \int \mathbf{L}_z(\boldsymbol{\theta}) \exp\left(-\frac{1}{s} \mathbf{L}_z(\boldsymbol{\theta})\right) dQ(\boldsymbol{\theta}) \quad (239)$$

$$= \int \mathbf{L}_z(\boldsymbol{\theta}) \exp\left(-K_{Q,z}\left(-\frac{1}{s}\right) - \frac{1}{s} \mathbf{L}_z(\boldsymbol{\theta})\right) dQ(\boldsymbol{\theta}) \quad (240)$$

$$= \int \mathbf{L}_z(\boldsymbol{\theta}) dP_{\boldsymbol{\Theta}|Z=z}^{(Q,s)}(\boldsymbol{\theta}), \quad (241)$$

where the equality in (237) holds from the dominated convergence theorem [28]; the equality in (239) follows from (21); and the equality in (241) follows from (25).

For all $s \in \mathcal{K}_{Q,z}$, with $\mathcal{K}_{Q,z}$ in (22), the equalities in (63) and (240) imply

that

$$K_{Q,z}^{(2)}\left(-\frac{1}{s}\right) = \frac{d}{dt} \int \mathbb{L}_z(\boldsymbol{\theta}) \exp(-K_{Q,z}(t) + t\mathbb{L}_z(\boldsymbol{\theta})) dQ(\boldsymbol{\theta}) \Big|_{t=-\frac{1}{s}} \quad (242)$$

$$= \int \mathbb{L}_z(\boldsymbol{\theta}) \left(-K_{Q,z}^{(1)}(t) + \mathbb{L}_z(\boldsymbol{\theta})\right) \exp(-K_{Q,z}(t) + t\mathbb{L}_z(\boldsymbol{\theta})) dQ(\boldsymbol{\theta}) \Big|_{t=-\frac{1}{s}} \quad (243)$$

$$= \int \mathbb{L}_z(\boldsymbol{\theta}) \left(-K_{Q,z}^{(1)}\left(-\frac{1}{s}\right) + \mathbb{L}_z(\boldsymbol{\theta})\right) \exp\left(-K_{Q,z}\left(-\frac{1}{s}\right) - \frac{1}{s}\mathbb{L}_z(\boldsymbol{\theta})\right) dQ(\boldsymbol{\theta}) \\ = \int \mathbb{L}_z(\boldsymbol{\theta}) \left(-K_{Q,z}^{(1)}\left(-\frac{1}{s}\right) + \mathbb{L}_z(\boldsymbol{\theta})\right) dP_{\boldsymbol{\Theta}|Z=z}^{(Q,s)}(\boldsymbol{\theta}) \quad (244)$$

$$= -K_{Q,z}^{(1)}\left(-\frac{1}{s}\right) \int \mathbb{L}_z(\boldsymbol{\theta}) dP_{\boldsymbol{\Theta}|Z=z}^{(Q,s)}(\boldsymbol{\theta}) + \int (\mathbb{L}_z(\boldsymbol{\theta}))^2 dP_{\boldsymbol{\Theta}|Z=z}^{(Q,s)}(\boldsymbol{\theta}) \quad (245)$$

$$= -\left(K_{Q,z}^{(1)}\left(-\frac{1}{s}\right)\right)^2 + \int (\mathbb{L}_z(\boldsymbol{\theta}))^2 dP_{\boldsymbol{\Theta}|Z=z}^{(Q,s)}(\boldsymbol{\theta}) \quad (246)$$

$$= \int \left(\mathbb{L}_z(\boldsymbol{\theta}) - K_{Q,z}^{(1)}\left(-\frac{1}{s}\right)\right)^2 dP_{\boldsymbol{\Theta}|Z=z}^{(Q,s)}(\boldsymbol{\theta}), \quad (247)$$

where the equality in (243) follows from the dominated convergence theorem [28]; the equality in (244) is due to a change of measure through the Radon-Nikodym derivative in (25); and the equality in (246) follows from (241).

For all $s \in \mathcal{K}_{Q,z}$, with $\mathcal{K}_{Q,z}$ in (22), the equalities in (63) and (246) imply that

$$\begin{aligned}
& K_{Q,z}^{(3)}\left(-\frac{1}{s}\right) \\
&= \frac{d}{dt} \left(\int (\mathbf{L}_z(\boldsymbol{\theta}))^2 dP_{\boldsymbol{\Theta}|\mathbf{Z}=z}^{(Q,-\frac{1}{t})}(\boldsymbol{\theta}) - \left(K_{Q,z}^{(1)}(t)\right)^2 \right) \Big|_{t=-\frac{1}{s}} \quad (248)
\end{aligned}$$

$$\begin{aligned}
&= \frac{d}{dt} \left(\int \left((\mathbf{L}_z(\boldsymbol{\theta}))^2 \exp(-K_{Q,z}(t) + t\mathbf{L}_z(\boldsymbol{\theta})) \right) dQ(\boldsymbol{\theta}) \right. \\
&\quad \left. - \left(K_{Q,z}^{(1)}(t)\right)^2 \right) \Big|_{t=-\frac{1}{s}} \quad (249)
\end{aligned}$$

$$\begin{aligned}
&= \int (\mathbf{L}_z(\boldsymbol{\theta}))^2 \left(\frac{d}{dt} \exp(-K_{Q,z}(t) + t\mathbf{L}_z(\boldsymbol{\theta})) \Big|_{t=-\frac{1}{s}} \right) dQ(\boldsymbol{\theta}) \\
&\quad - 2K_{Q,z}^{(1)}(t) K_{Q,z}^{(2)}(t) \Big|_{t=-\frac{1}{s}} \quad (250)
\end{aligned}$$

$$\begin{aligned}
&= \int (\mathbf{L}_z(\boldsymbol{\theta}))^2 \left((\mathbf{L}_z(\boldsymbol{\theta}) - K_{Q,z}^{(1)}(t)) \exp(-K_{Q,z}(t) + t\mathbf{L}_z(\boldsymbol{\theta})) \Big|_{t=-\frac{1}{s}} \right) dQ(\boldsymbol{\theta}) \\
&\quad - 2K_{Q,z}^{(1)}(t) K_{Q,z}^{(2)}(t) \Big|_{t=-\frac{1}{s}} \quad (251)
\end{aligned}$$

$$\begin{aligned}
&= \int (\mathbf{L}_z(\boldsymbol{\theta}))^2 \left(\mathbf{L}_z(\boldsymbol{\theta}) - K_{Q,z}^{(1)}\left(-\frac{1}{s}\right) \right) \exp\left(-K_{Q,z}\left(-\frac{1}{s}\right) - \frac{1}{s}\mathbf{L}_z(\boldsymbol{\theta})\right) dQ(\boldsymbol{\theta}) \\
&\quad - 2K_{Q,z}^{(1)}\left(-\frac{1}{s}\right) K_{Q,z}^{(2)}\left(-\frac{1}{s}\right) \quad (252)
\end{aligned}$$

$$\begin{aligned}
&= \int (\mathbf{L}_z(\boldsymbol{\theta}))^2 \left(\mathbf{L}_z(\boldsymbol{\theta}) - K_{Q,z}^{(1)}\left(-\frac{1}{s}\right) \right) dP_{\boldsymbol{\Theta}|\mathbf{Z}=z}^{(Q,s)}(\boldsymbol{\theta}) \\
&\quad - 2K_{Q,z}^{(1)}\left(-\frac{1}{s}\right) K_{Q,z}^{(2)}\left(-\frac{1}{s}\right) \quad (253)
\end{aligned}$$

$$\begin{aligned}
&= \int (\mathbf{L}_z(\boldsymbol{\theta}))^3 dP_{\boldsymbol{\Theta}|\mathbf{Z}=z}^{(Q,s)}(\boldsymbol{\theta}) \\
&\quad - K_{Q,z}^{(1)}\left(-\frac{1}{s}\right) \int (\mathbf{L}_z(\boldsymbol{\theta}))^2 dP_{\boldsymbol{\Theta}|\mathbf{Z}=z}^{(Q,s)}(\boldsymbol{\theta}) - 2K_{Q,z}^{(1)}\left(-\frac{1}{s}\right) K_{Q,z}^{(2)}\left(-\frac{1}{s}\right) \quad (254)
\end{aligned}$$

$$\begin{aligned}
&= \int (\mathbf{L}_z(\boldsymbol{\theta}))^3 dP_{\boldsymbol{\Theta}|\mathbf{Z}=z}^{(Q,s)}(\boldsymbol{\theta}) \\
&\quad - K_{Q,z}^{(1)}\left(-\frac{1}{s}\right) \left(K_{Q,z}^{(2)}\left(-\frac{1}{s}\right) + \left(K_{Q,z}^{(1)}\left(-\frac{1}{s}\right)\right)^2 \right) \\
&\quad - 2K_{Q,z}^{(1)}\left(-\frac{1}{s}\right) K_{Q,z}^{(2)}\left(-\frac{1}{s}\right) \quad (255)
\end{aligned}$$

$$\begin{aligned}
&= \int (\mathbf{L}_z(\boldsymbol{\theta}))^3 dP_{\boldsymbol{\Theta}|\mathbf{Z}=z}^{(Q,s)}(\boldsymbol{\theta}) - K_{Q,z}^{(1)}\left(-\frac{1}{s}\right)^3 - 3K_{Q,z}^{(1)}\left(-\frac{1}{s}\right) K_{Q,z}^{(2)}\left(-\frac{1}{s}\right) \quad (256)
\end{aligned}$$

$$\begin{aligned}
&= \int \left(\mathbf{L}_z(\boldsymbol{\theta}) - K_{Q,z}^{(1)}\left(-\frac{1}{s}\right) \right)^3 dP_{\boldsymbol{\Theta}|\mathbf{Z}=z}^{(Q,s)}(\boldsymbol{\theta}), \quad (257) \\
&\quad \text{Inria}
\end{aligned}$$

where the equality in (249) follows from (25); the equality in (250) follows from the dominated convergence theorem [28]; the equality in (253) follows from (25); and the equality in (255) follows from (246). This completes the proof.

N Proof of Theorem 5.1

The proof is based on the analysis of the derivative of $K_{Q,z}^{(1)}\left(-\frac{1}{\lambda}\right)$ with respect to λ in $\text{int}\mathcal{K}_{Q,z}$. That is,

$$\frac{d}{d\lambda}K_{Q,z}^{(1)}\left(-\frac{1}{\lambda}\right) = \frac{1}{\lambda^2}K_{Q,z}^{(2)}\left(-\frac{1}{\lambda}\right) \quad (258)$$

$$\geq 0, \quad (259)$$

where the equality in (259) follows from Lemma 4.4. From Lemma 5.1, the inequality in (259) implies that the expected empirical risk $R_z\left(P_{\Theta|Z=z}^{(Q,\lambda)}\right) = K_{Q,z}^{(1)}\left(-\frac{1}{\lambda}\right)$ in (21) is nondecreasing with respect to λ . That is, given two reals λ_1 and λ_2 , with $\lambda_1 > \lambda_2 > 0$, it holds that

$$R_z\left(P_{\Theta|Z=z}^{(Q,\lambda_1)}\right) = K_{Q,z}^{(1)}\left(-\frac{1}{\lambda_1}\right) \geq K_{Q,z}^{(1)}\left(-\frac{1}{\lambda_2}\right) = R_z\left(P_{\Theta|Z=z}^{(Q,\lambda_2)}\right). \quad (260)$$

The rest of the proof consists in showing that for all $\alpha \in \mathcal{K}_{Q,z}$, the function $K_{Q,z}^{(2)}$ in (63) satisfies $K_{Q,z}^{(2)}\left(-\frac{1}{\alpha}\right) > 0$ if and only if the function L_z in (5) is separable. For doing so, a handful of preliminary results are described in the following subsection. The proof of Theorem 5.1 resumes in Subsection N.2

N.1 Preliminaries

Given a positive real $\lambda \in \mathcal{K}_{Q,z}$, with $\mathcal{K}_{Q,z}$ in (22), 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

$$\mathcal{R}_0(\lambda) \triangleq \left\{ \nu \in \mathcal{M} : L_z(\nu) = R_z\left(P_{\Theta|Z=z}^{(Q,\lambda)}\right) \right\}, \quad (261a)$$

$$\mathcal{R}_1(\lambda) \triangleq \left\{ \nu \in \mathcal{M} : L_z(\nu) < R_z\left(P_{\Theta|Z=z}^{(Q,\lambda)}\right) \right\}, \text{ and} \quad (261b)$$

$$\mathcal{R}_2(\lambda) \triangleq \left\{ \nu \in \mathcal{M} : L_z(\nu) > R_z\left(P_{\Theta|Z=z}^{(Q,\lambda)}\right) \right\}, \quad (261c)$$

where the function R_z is in (10) and the probability measure $P_{\Theta|Z=z}^{(Q,\lambda)}$ is in (25). Note that the value $R_z\left(P_{\Theta|Z=z}^{(Q,\lambda)}\right)$ is the g-ERM-RER-optimal expected empirical risk in (68).

The sets in (261) 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*

$$P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{R}_1(\lambda)) > 0, \quad (262)$$

if and only if

$$P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{R}_2(\lambda)) > 0, \quad (263)$$

where the sets $\mathcal{R}_1(\alpha)$ and $\mathcal{R}_2(\alpha)$ are in (261b) and (261c), 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 of the proof consists in proving that, given a real $\alpha \in \mathcal{K}_{Q,z}$, if the set $\mathcal{R}_2(\alpha)$ is nonnegligible with respect to $P_{\Theta|Z=z}^{(Q,\alpha)}$, then the set $\mathcal{R}_1(\alpha)$ is nonnegligible with respect to $P_{\Theta|Z=z}^{(Q,\alpha)}$.

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 4.4, it holds that

$$\begin{aligned} K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right) &= \int_{\mathcal{R}_0(\alpha)} \mathbf{L}_z(\boldsymbol{\nu}) \, dP_{\Theta|Z=z}^{(Q,\alpha)}(\boldsymbol{\nu}) + \int_{\mathcal{R}_1(\alpha)} \mathbf{L}_z(\boldsymbol{\nu}) \, dP_{\Theta|Z=z}^{(Q,\alpha)}(\boldsymbol{\nu}) \\ &\quad + \int_{\mathcal{R}_2(\alpha)} \mathbf{L}_z(\boldsymbol{\nu}) \, dP_{\Theta|Z=z}^{(Q,\alpha)}(\boldsymbol{\nu}) \end{aligned} \quad (264)$$

$$= \int_{\mathcal{R}_0(\alpha)} \mathbf{L}_z(\boldsymbol{\nu}) \, dP_{\Theta|Z=z}^{(Q,\alpha)}(\boldsymbol{\nu}) + \int_{\mathcal{R}_1(\alpha)} \mathbf{L}_z(\boldsymbol{\nu}) \, dP_{\Theta|Z=z}^{(Q,\alpha)}(\boldsymbol{\nu}) \quad (265)$$

$$= K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right) P_{\Theta|Z=z}^{(Q,\alpha)}(\mathcal{R}_0(\alpha)) + \int_{\mathcal{R}_1(\alpha)} \mathbf{L}_z(\boldsymbol{\nu}) \, dP_{\Theta|Z=z}^{(Q,\alpha)}(\boldsymbol{\nu}) \quad (266)$$

$$< K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right) P_{\Theta|Z=z}^{(Q,\alpha)}(\mathcal{R}_0(\alpha)) + K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right) P_{\Theta|Z=z}^{(Q,\alpha)}(\mathcal{R}_1(\alpha)) \quad (267)$$

$$= K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right) \left(P_{\Theta|Z=z}^{(Q,\alpha)}(\mathcal{R}_0(\alpha)) + P_{\Theta|Z=z}^{(Q,\alpha)}(\mathcal{R}_1(\alpha)) \right) \quad (268)$$

$$= K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right), \quad (269)$$

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 4.4, it holds that

$$\begin{aligned}
 K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right) &= \int_{\mathcal{R}_0(\alpha)} \mathbf{L}_z(\boldsymbol{\nu}) dP_{\Theta|Z=z}^{(Q,\alpha)}(\boldsymbol{\nu}) + \int_{\mathcal{R}_1(\alpha)} \mathbf{L}_z(\boldsymbol{\nu}) dP_{\Theta|Z=z}^{(Q,\alpha)}(\boldsymbol{\nu}) \\
 &\quad + \int_{\mathcal{R}_2(\alpha)} \mathbf{L}_z(\boldsymbol{\nu}) dP_{\Theta|Z=z}^{(Q,\alpha)}(\boldsymbol{\nu}) \tag{270}
 \end{aligned}$$

$$= \int_{\mathcal{R}_0(\alpha)} \mathbf{L}_z(\boldsymbol{\nu}) dP_{\Theta|Z=z}^{(Q,\alpha)}(\boldsymbol{\nu}) + \int_{\mathcal{R}_2(\alpha)} \mathbf{L}_z(\boldsymbol{\nu}) dP_{\Theta|Z=z}^{(Q,\alpha)}(\boldsymbol{\nu}) \tag{271}$$

$$= K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right) P_{\Theta|Z=z}^{(Q,\alpha)}(\mathcal{R}_0(\alpha)) + \int_{\mathcal{R}_2(\alpha)} \mathbf{L}_z(\boldsymbol{\nu}) dP_{\Theta|Z=z}^{(Q,\alpha)}(\boldsymbol{\nu}) \tag{272}$$

$$\begin{aligned}
 &> K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right) P_{\Theta|Z=z}^{(Q,\alpha)}(\mathcal{R}_0(\alpha)) \\
 &\quad + K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right) P_{\Theta|Z=z}^{(Q,\alpha)}(\mathcal{R}_2(\alpha)) \tag{273}
 \end{aligned}$$

$$= K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right) \left(P_{\Theta|Z=z}^{(Q,\alpha)}(\mathcal{R}_0(\alpha)) + P_{\Theta|Z=z}^{(Q,\alpha)}(\mathcal{R}_2(\alpha)) \right) \tag{274}$$

$$= K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right), \tag{275}$$

which is also a contradiction. This completes the proof. \blacksquare

A more general result can be immediately obtained by combining Lemma 3.12 and Lemma N.1.

Lemma N.2. *For all $\alpha \in \mathcal{K}_{Q,z}$, with $\mathcal{K}_{Q,z}$ in (22), the measure $P_{\Theta|Z=z}^{(Q,\lambda)}$ in (25), satisfies*

$$P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{R}_1(\alpha)) > 0, \tag{276}$$

if and only if

$$P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{R}_2(\alpha)) > 0, \tag{277}$$

where the sets $\mathcal{R}_1(\alpha)$ and $\mathcal{R}_2(\alpha)$ are in (261b) and (261c), respectively.

N.2 The proof

The rest of the proof of Theorem 5.1 is divided into two parts. In the first part, it is shown that if for all $\alpha \in \mathcal{K}_{Q,z}$, $K_{Q,z}^{(2)}\left(-\frac{1}{\alpha}\right) > 0$, then the function \mathbf{L}_z in (5) is separable. The second part of the proof, consists in showing that if the function \mathbf{L}_z is separable, then, for all $\alpha \in \mathcal{K}_{Q,z}$, $K_{Q,z}^{(2)}\left(-\frac{1}{\alpha}\right) > 0$.

The first part is as follows. From Lemma 4.4, it holds that for all $\alpha \in \mathcal{K}_{Q,z}$,

$$K_{Q,z}^{(2)}\left(-\frac{1}{\alpha}\right) = \int \left(L_z(\boldsymbol{\theta}) - K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right) \right)^2 dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\alpha)}(\boldsymbol{\theta}) \quad (278)$$

$$= \int_{\mathcal{R}_0(\alpha)} \left(L_z(\boldsymbol{\theta}) - K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right) \right)^2 dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\alpha)}(\boldsymbol{\theta}) \quad (279)$$

$$+ \int_{\mathcal{R}_1(\alpha)} \left(L_z(\boldsymbol{\theta}) - K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right) \right)^2 dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\alpha)}(\boldsymbol{\theta}) \quad (280)$$

$$+ \int_{\mathcal{R}_2(\alpha)} \left(L_z(\boldsymbol{\theta}) - K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right) \right)^2 dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\alpha)}(\boldsymbol{\theta}), \quad (281)$$

where the sets $\mathcal{R}_0(\alpha)$, $\mathcal{R}_1(\alpha)$, and $\mathcal{R}_2(\alpha)$ are respectively defined in (261). Hence,

$$K_{Q,z}^{(2)}\left(-\frac{1}{\alpha}\right) = \int_{\mathcal{R}_1(\alpha)} \left(L_z(\boldsymbol{\theta}) - K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right) \right)^2 dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\alpha)}(\boldsymbol{\theta}) \\ + \int_{\mathcal{R}_2(\alpha)} \left(L_z(\boldsymbol{\theta}) - K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right) \right)^2 dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\alpha)}(\boldsymbol{\theta}) \quad (282)$$

$$\leq \left(0 - K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right) \right)^2 P_{\boldsymbol{\Theta}|Z=z}^{(Q,\alpha)}(\mathcal{R}_1(\alpha)) \\ + \left(\sup_{\boldsymbol{\theta}} L_z(\boldsymbol{\theta}) - K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right) \right)^2 P_{\boldsymbol{\Theta}|Z=z}^{(Q,\alpha)}(\mathcal{R}_2(\alpha)). \quad (283)$$

Under the assumption that for all $\alpha \in \mathcal{K}_{Q,z}$ the function $K_{Q,z}^{(2)}$ in (63) satisfies $K_{Q,z}^{(2)}\left(-\frac{1}{\alpha}\right) > 0$, it follows from (283) that

$$0 < \left(0 - K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right) \right)^2 P_{\boldsymbol{\Theta}|Z=z}^{(Q,\alpha)}(\mathcal{R}_1(\alpha)) \\ + \left(\sup_{\boldsymbol{\theta}} L_z(\boldsymbol{\theta}) - K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right) \right)^2 P_{\boldsymbol{\Theta}|Z=z}^{(Q,\alpha)}(\mathcal{R}_2(\alpha)). \quad (284)$$

Note that if

$$P_{\boldsymbol{\Theta}|Z=z}^{(Q,\alpha)}(\mathcal{R}_1(\alpha)) > 0, \quad (285)$$

then, $0 \neq K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right)$. Moreover, if

$$P_{\boldsymbol{\Theta}|Z=z}^{(Q,\alpha)}(\mathcal{R}_2(\alpha)) > 0, \quad (286)$$

then, $\sup_{\boldsymbol{\theta}} L_z(\boldsymbol{\theta}) \neq K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right)$. Therefore, the inequality in (284) implies that at least one of the following claims is true:

- (a) $P_{\boldsymbol{\Theta}|Z=z}^{(Q,\alpha)}(\mathcal{R}_1(\alpha)) > 0$; and
- (b) $P_{\boldsymbol{\Theta}|Z=z}^{(Q,\alpha)}(\mathcal{R}_2(\alpha)) > 0$.

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 $(\nu_1, \nu_2) \in \mathcal{R}_1(\alpha) \times \mathcal{R}_2(\alpha)$,

$$\mathbb{L}_z(\nu_1) > K_{Q,z}^{(1)}\left(-\frac{1}{\alpha}\right) > \mathbb{L}_z(\nu_2). \quad (287)$$

This proves that under the assumption that for all $\alpha \in \mathcal{K}_{Q,z}$, $K_{Q,z}^{(2)}\left(-\frac{1}{\alpha}\right) > 0$, the function \mathbb{L}_z in (5) is separable with respect to $P_{\Theta|Z=z}^{(Q,\alpha)}$. From Lemma 4.2, 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 (5) is separable with respect to $P_{\Theta|Z=z}^{(Q,\gamma)}$. 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 $(\nu_1, \nu_2) \in \mathcal{A}(\gamma) \times \mathcal{B}(\gamma)$,

$$\mathbb{L}_z(\nu_1) > c_\gamma > \mathbb{L}_z(\nu_2). \quad (288)$$

From Lemma 4.4, it holds that

$$K_{Q,z}^{(2)}\left(-\frac{1}{\gamma}\right) = \int \left(\mathbb{L}_z(\theta) - K_{Q,z}^{(1)}\left(-\frac{1}{\gamma}\right)\right)^2 dP_{\Theta|Z=z}^{(Q,\gamma)}(\theta) \quad (289)$$

$$= \int_{\mathcal{A}(\gamma)} \left(\mathbb{L}_z(\theta) - K_{Q,z}^{(1)}\left(-\frac{1}{\gamma}\right)\right)^2 dP_{\Theta|Z=z}^{(Q,\gamma)}(\theta) \quad (290)$$

$$+ \int_{\mathcal{B}(\gamma)} \left(\mathbb{L}_z(\theta) - K_{Q,z}^{(1)}\left(-\frac{1}{\gamma}\right)\right)^2 dP_{\Theta|Z=z}^{(Q,\gamma)}(\theta) \quad (291)$$

$$+ \int_{\mathcal{M} \setminus (\mathcal{A}(\gamma) \cup \mathcal{B}(\gamma))} \left(\mathbb{L}_z(\theta) - K_{Q,z}^{(1)}\left(-\frac{1}{\gamma}\right)\right)^2 dP_{\Theta|Z=z}^{(Q,\gamma)}(\theta) \quad (292)$$

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

where the inequality (293) follows from the following facts. First, if $c_\gamma < K_{Q,z}^{(1)}\left(-\frac{1}{\gamma}\right)$, with c_γ in (288), then for all $\nu \in \mathcal{B}(\gamma)$, it holds that

$$\left(\mathbb{L}_z(\theta) - K_{Q,z}^{(1)}\left(-\frac{1}{\gamma}\right)\right)^2 > \left(c_\gamma - K_{Q,z}^{(1)}\left(-\frac{1}{\gamma}\right)\right)^2, \quad (294)$$

and thus,

$$\int_{\mathcal{B}(\gamma)} \left(\mathbb{L}_z(\theta) - K_{Q,z}^{(1)}\left(-\frac{1}{\gamma}\right)\right)^2 dP_{\Theta|Z=z}^{(Q,\gamma)}(\theta) > \left(c_\gamma - K_{Q,z}^{(1)}\left(-\frac{1}{\gamma}\right)\right)^2 P_{\Theta|Z=z}^{(Q,\gamma)}(\mathcal{B}(\gamma)) > 0. \quad (295)$$

Second, if $c_\gamma \geq K_{Q,z}^{(1)}\left(-\frac{1}{\gamma}\right)$ then for all $\boldsymbol{\nu} \in \mathcal{A}(\gamma)$, it holds that

$$\left(\mathbf{L}_z(\boldsymbol{\theta}) - K_{Q,z}^{(1)}\left(-\frac{1}{\gamma}\right)\right)^2 > \left(c_\gamma - K_{Q,z}^{(1)}\left(-\frac{1}{\gamma}\right)\right)^2, \quad (297)$$

and thus,

$$\begin{aligned} & \int_{\mathcal{A}(\gamma)} \left(\mathbf{L}_z(\boldsymbol{\theta}) - K_{Q,z}^{(1)}\left(-\frac{1}{\gamma}\right)\right)^2 dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\gamma)}(\boldsymbol{\theta}) \\ & > \left(c_\gamma - K_{Q,z}^{(1)}\left(-\frac{1}{\gamma}\right)\right)^2 P_{\boldsymbol{\Theta}|Z=z}^{(Q,\gamma)}(\mathcal{A}(\gamma)) \\ & > 0. \end{aligned} \quad (298)$$

$$(299)$$

Hence, under the assumption that the empirical risk function \mathbf{L}_z in (5) 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 5.1

Consider the partition of the set \mathcal{M} formed by the sets following sets \mathcal{A}_0 , \mathcal{A}_1 , and \mathcal{A}_2 in (193). From (64), for all $\lambda \in \mathcal{K}_{Q,z}$, with $\mathcal{K}_{Q,z}$ in (22), it holds that,

$$K_{Q,z}^{(1)}\left(-\frac{1}{\lambda}\right) = \int_{\mathcal{A}_0} \mathbf{L}_z(\boldsymbol{\theta}) dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}(\boldsymbol{\theta}) + \int_{\mathcal{A}_1} \mathbf{L}_z(\boldsymbol{\theta}) dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}(\boldsymbol{\theta}) \quad (300)$$

$$+ \int_{\mathcal{A}_2} \mathbf{L}_z(\boldsymbol{\theta}) dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}(\boldsymbol{\theta}) \quad (301)$$

$$= \int_{\mathcal{A}_0} \mathbf{L}_z(\boldsymbol{\theta}) dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}(\boldsymbol{\theta}) + \int_{\mathcal{A}_2} \mathbf{L}_z(\boldsymbol{\theta}) dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}(\boldsymbol{\theta}) \quad (302)$$

$$= \delta_{Q,z}^* P_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}(\mathcal{L}_{Q,z}^*) + \int_{\mathcal{A}_2} \mathbf{L}_z(\boldsymbol{\theta}) dP_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}(\boldsymbol{\theta}) \quad (303)$$

$$\geq \delta_{Q,z}^* P_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}(\mathcal{L}_{Q,z}^*) + \delta_{Q,z}^* P_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}(\mathcal{A}_2) \quad (304)$$

$$= \delta_{Q,z}^*, \quad (305)$$

where, the equality in (302) follows by noticing that $Q(\mathcal{A}_1) = 0$, which implies that $P_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}(\mathcal{A}_1) = 0$ (Lemma 3.8); the equality in (303) follows from noticing that $\mathcal{A}_0 = \mathcal{L}_{Q,z}^*$, with $\mathcal{L}_{Q,z}^*$ in (33); and the equality in (304) follows from (193c).

This completes the proof.

P Proof of Theorem 5.2

From (303) in the proof of Lemma 5.1, it holds that

$$\lim_{\lambda \rightarrow 0^+} K_{Q,z}^{(1)} \left(-\frac{1}{\lambda} \right) = \lim_{\lambda \rightarrow 0^+} \delta_{Q,z}^* P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_{Q,z}^*) + \lim_{\lambda \rightarrow 0^+} \int_{\mathcal{A}_2} \mathbb{L}_z(\boldsymbol{\theta}) dP_{\Theta|Z=z}^{(Q,\lambda)}(\boldsymbol{\theta}) \quad (306)$$

$$\begin{aligned} &= \lim_{\lambda \rightarrow 0^+} \delta_{Q,z}^* P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_{Q,z}^*) \\ &\quad + \lim_{\lambda \rightarrow 0^+} \int_{\mathcal{A}_2} \mathbb{L}_z(\boldsymbol{\theta}) \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}(\boldsymbol{\theta})}{dQ}(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) \end{aligned} \quad (307)$$

$$\begin{aligned} &= \lim_{\lambda \rightarrow 0^+} \delta_{Q,z}^* P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_{Q,z}^*) \\ &\quad + \int_{\mathcal{A}_2} \mathbb{L}_z(\boldsymbol{\theta}) \lim_{\lambda \rightarrow 0^+} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}(\boldsymbol{\theta})}{dQ}(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) \end{aligned} \quad (308)$$

$$= \delta_{Q,z}^* \lim_{\lambda \rightarrow 0^+} P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_{Q,z}^*) \quad (309)$$

$$= \delta_{Q,z}^*, \quad (310)$$

where, the equality in (308) follows from noticing two facts: (a) For all $\lambda \in \mathcal{K}_{Q,z}$, the Radon-Nikodym derivative $\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}$ is positive and finite (Lemma 3.4); and (b) For all $\boldsymbol{\theta} \in \mathcal{A}_2$, it holds that $\lim_{\lambda \rightarrow 0^+} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}(\boldsymbol{\theta})}{dQ}(\boldsymbol{\theta}) = 0$. Hence, the dominated convergence theorem [28, Theorem 1.6.9] holds. The inequality in (309) follows from Lemma 3.7 This completes the proof.

Q Proof of Theorem 7.1

From Theorem 5.1, it follows that for all $(\lambda_1, \lambda_2) \in \mathcal{K}_{Q,z} \times \mathcal{K}_{Q,z}$ with $\lambda_1 > \lambda_2$,

$$\int \mathbb{L}_z(\boldsymbol{\alpha}) \frac{dP_{\Theta|Z=z}^{(Q,\lambda_1)}(\boldsymbol{\alpha})}{dQ}(\boldsymbol{\alpha}) dQ(\boldsymbol{\alpha}) \geq \int \mathbb{L}_z(\boldsymbol{\alpha}) \frac{dP_{\Theta|Z=z}^{(Q,\lambda_2)}(\boldsymbol{\alpha})}{dQ}(\boldsymbol{\alpha}) dQ(\boldsymbol{\alpha}), \quad (311)$$

which implies the following inclusions:

$$\mathcal{R}_1(\lambda_2) \subseteq \mathcal{R}_1(\lambda_1), \text{ and} \quad (312a)$$

$$\mathcal{R}_2(\lambda_1) \subseteq \mathcal{R}_2(\lambda_2), \quad (312b)$$

with the sets $\mathcal{R}_1(\cdot)$ and $\mathcal{R}_2(\cdot)$ in (261). From (82), it holds that for all $i \in \{1, 2\}$,

$$\mathcal{N}_{Q,z}(\lambda_i) = \mathcal{R}_2(\lambda_i)^c, \quad (313)$$

where the complement is with respect to \mathcal{M} . Thus, the inclusion in (312b) and the equality in (313) yields,

$$\mathcal{N}_{Q,z}(\lambda_1) \supseteq \mathcal{N}_{Q,z}(\lambda_2). \quad (314)$$

The inclusion $\mathcal{M} \supseteq \mathcal{N}_{Q,z}(\lambda_1)$ follows from (82). Alternatively, the inclusion $\mathcal{N}_{Q,z}(\lambda_2) \supseteq \mathcal{N}_{Q,z}^*$, follows from Lemma 5.1 and from observing that for all $\nu \in \mathcal{N}_{Q,z}^*$,

$$\mathbb{R}_z \left(P_{\Theta|Z=z}^{(Q,\lambda_2)} \right) \geq \delta_{Q,z}^* \quad (315)$$

$$\geq \mathbb{L}_z(\nu), \quad (316)$$

which implies that $\nu \in \mathcal{N}_{Q,z}(\lambda_2)$. This completes the proof of (86).

The proof of (87) is as follows. From the mean value theorem [28] and the assumption that the empirical risk function \mathbb{L}_z in (5) 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

$$\mathbb{L}_z(\theta) = \int \mathbb{L}_z(\alpha) dP_{\Theta|Z=z}^{(Q,\lambda)}(\alpha), \quad (317)$$

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 (314), verifies that

$$\mathcal{N}_{Q,z}(\lambda_1) \supset \mathcal{N}_{Q,z}(\lambda_2). \quad (318)$$

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,

$$P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{R}_0(\lambda)) = 1. \quad (319)$$

This contradicts the assumption that the function \mathbb{L}_z is separable (Definition 4.1). 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,

$$\{\theta \in \mathcal{M} : \mathbb{L}_z(\theta) \leq \delta_{Q,z}^*\} = \mathcal{N}_{Q,z}^* \quad (320)$$

$$= \mathcal{N}_{Q,z}(\lambda) \quad (321)$$

$$= \left\{ \theta \in \mathcal{M} : \mathbb{L}_z(\theta) \leq K_{Q,z}^{(1)} \left(-\frac{1}{\lambda} \right) \right\}. \quad (322)$$

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

$$\left\{ \nu \in \mathcal{M} : \delta_{Q,z}^* < \mathbb{L}_z(\nu) \leq K_{Q,z}^{(1)} \left(-\frac{1}{\lambda} \right) \right\} = \emptyset;$$

(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,

$$\mathcal{R}_1(\lambda) = \left\{ \nu \in \mathcal{M} : \mathbb{L}_z(\nu) < \delta_{Q,z}^* \right\}, \quad (323)$$

and from the definition of $\delta_{Q,z}^*$ in (32), it holds that

$$P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{R}_1(\lambda)) = 0. \quad (324)$$

From Lemma N.1 and (324), it follows that,

$$P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{R}_2(\lambda)) = 0. \quad (325)$$

Finally, by noticing that

$$1 = P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{R}_0(\lambda)) + P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{R}_1(\lambda)) + P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{R}_2(\lambda)) \quad (326)$$

$$= P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{R}_0(\lambda)), \quad (327)$$

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 4.2). This completes the proof of (87).

R Proof of Theorem 7.2

The proof of (88) is based on the analysis of the derivative of $P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{A})$ with respect to λ , for some fixed set $\mathcal{A} \subseteq \mathcal{M}$. More specifically, given a $\gamma \in \mathcal{K}_{Q,z}$, it holds that

$$P_{\Theta|Z=z}^{(Q,\gamma)}(\mathcal{A}) = \int_{\mathcal{A}} \frac{dP_{\Theta|Z=z}^{(Q,\gamma)}(\alpha)}{dQ}(\alpha) dQ(\alpha), \quad (328)$$

and from the fundamental theorem of calculus [34], it follows that for all $(\lambda_1, \lambda_2) \in \mathcal{K}_{Q,z} \times \mathcal{K}_{Q,z}$ with $\lambda_1 > \lambda_2$,

$$P_{\Theta|Z=z}^{(Q,\lambda_1)}(\mathcal{A}) - P_{\Theta|Z=z}^{(Q,\lambda_2)}(\mathcal{A}) = \int_{\lambda_2}^{\lambda_1} \frac{d}{d\gamma} P_{\Theta|Z=z}^{(Q,\gamma)}(\mathcal{A}) d\gamma \quad (329)$$

$$= \int_{\lambda_2}^{\lambda_1} \frac{d}{d\gamma} \int_{\mathcal{A}} \frac{dP_{\Theta|Z=z}^{(Q,\gamma)}(\alpha)}{dQ}(\alpha) dQ(\alpha) d\gamma \quad (330)$$

$$= \int_{\lambda_2}^{\lambda_1} \int_{\mathcal{A}} \frac{d}{d\gamma} \frac{dP_{\Theta|Z=z}^{(Q,\gamma)}(\alpha)}{dQ}(\alpha) dQ(\alpha) d\gamma, \quad (331)$$

where the equality in (330) follows from (328); and the equality in (331) holds from Lemma 3.4 and the dominated convergence theorem [28].

For all $\theta \in \text{supp } Q$, the following holds,

$$\frac{d}{d\lambda} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) = \frac{d}{d\lambda} \frac{\exp\left(-\frac{L_z(\theta)}{\lambda}\right)}{\int \exp\left(-\frac{L_z(\nu)}{\lambda}\right) dQ(\nu)} \quad (332)$$

$$\begin{aligned} &= \frac{\frac{1}{\lambda^2} L_z(\theta) \exp\left(-\frac{L_z(\theta)}{\lambda}\right)}{\int \exp\left(-\frac{L_z(\nu)}{\lambda}\right) dQ(\nu)} \\ &= \frac{\frac{1}{\lambda^2} \exp\left(-\frac{L_z(\theta)}{\lambda}\right) \int L_z(\alpha) \exp\left(-\frac{L_z(\alpha)}{\lambda}\right) dQ(\alpha)}{\left(\int \exp\left(-\frac{L_z(\nu)}{\lambda}\right) dQ(\nu)\right)^2} \quad (333) \end{aligned}$$

$$\begin{aligned} &= \frac{1}{\lambda^2} L_z(\theta) \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) \\ &\quad - \frac{1}{\lambda^2} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) \int L_z(\nu) \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\nu) dQ(\nu) \quad (334) \end{aligned}$$

$$= \frac{1}{\lambda^2} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) \left(L_z(\theta) - \int L_z(\nu) dP_{\Theta|Z=z}^{(Q,\lambda)}(\nu) \right). \quad (335)$$

Plugging (335) into (331) yields,

$$\begin{aligned} &P_{\Theta|Z=z}^{(Q,\lambda_1)}(\mathcal{A}) - P_{\Theta|Z=z}^{(Q,\lambda_2)}(\mathcal{A}) \\ &= \int_{\lambda_2}^{\lambda_1} \int_{\mathcal{A}} \frac{1}{\gamma^2} \frac{dP_{\Theta|Z=z}^{(Q,\gamma)}}{dQ}(\alpha) \left(L_z(\alpha) - \int L_z(\nu) dP_{\Theta|Z=z}^{(Q,\gamma)}(\nu) \right) dQ(\alpha) d\gamma \quad (336) \end{aligned}$$

$$= \int_{\lambda_2}^{\lambda_1} \int_{\mathcal{A}} \frac{1}{\gamma^2} \left(L_z(\alpha) - \int L_z(\nu) dP_{\Theta|Z=z}^{(Q,\gamma)}(\nu) \right) dP_{\Theta|Z=z}^{(Q,\gamma)}(\alpha) d\gamma. \quad (337)$$

Note that for all $\alpha \in \mathcal{N}_{Q,z}(\lambda_2)$, it holds that for all $\gamma \in (\lambda_2, \lambda_1)$,

$$L_z(\alpha) - \int L_z(\nu) dP_{\Theta|Z=z}^{(Q,\gamma)}(\nu) \leq 0, \quad (338)$$

and thus,

$$\int_{\mathcal{N}_{Q,z}(\lambda_2)} \frac{1}{\gamma^2} \left(L_z(\alpha) - \int L_z(\nu) dP_{\Theta|Z=z}^{(Q,\gamma)}(\nu) \right) dP_{\Theta|Z=z}^{(Q,\gamma)}(\alpha) \leq 0. \quad (339)$$

The equalities in (337) and (339), with $\mathcal{A} = \mathcal{N}_{Q,z}(\lambda)$, imply that

$$P_{\Theta|Z=z}^{(Q,\lambda_1)}(\mathcal{N}_{Q,z}(\lambda_2)) - P_{\Theta|Z=z}^{(Q,\lambda_2)}(\mathcal{N}_{Q,z}(\lambda_2)) \leq 0. \quad (340)$$

The inequality $0 < P_{\Theta|Z=z}^{(Q,\lambda_1)}(\mathcal{N}_{Q,z}(\lambda_2))$ in (88) 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 (88).

The proof of (89) 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 $\lambda_1 > \lambda_2$,

$$P_{\Theta|Z=z}^{(Q,\lambda_1)}(\mathcal{N}_{Q,z}(\lambda_2)) < P_{\Theta|Z=z}^{(Q,\lambda_2)}(\mathcal{N}_{Q,z}(\lambda_2)), \quad (341)$$

then the function L_z is separable with respect to Q . The second part of the proof shows that if the function 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 (341) holds.

The first part is as follows. In the proof of Theorem 7.1 it is shown (see (337)) that for all pairs $(\lambda_1, \lambda_2) \in \mathcal{K}_{Q,z} \times \mathcal{K}_{Q,z}$ with $\lambda_1 > \lambda_2$,

$$\begin{aligned} & P_{\Theta|Z=z}^{(Q,\lambda_1)}(\mathcal{N}_{Q,z}(\lambda_2)) - P_{\Theta|Z=z}^{(Q,\lambda_2)}(\mathcal{N}_{Q,z}(\lambda_2)) \\ &= \int_{\lambda_2}^{\lambda_1} \int_{\mathcal{N}_{Q,z}(\lambda_2)} \frac{1}{\gamma^2} \left(L_z(\alpha) - \int L_z(\nu) dP_{\Theta|Z=z}^{(Q,\gamma)}(\nu) \right) dP_{\Theta|Z=z}^{(Q,\gamma)}(\alpha) d\gamma \end{aligned} \quad (342)$$

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 (341) holds. Then, from (342),

$$\begin{aligned} & 0 > \int_{\lambda_2}^{\lambda_1} \int_{\mathcal{N}_{Q,z}(\lambda_2)} \frac{1}{\gamma^2} \left(L_z(\alpha) - \int L_z(\nu) dP_{\Theta|Z=z}^{(Q,\gamma)}(\nu) \right) dP_{\Theta|Z=z}^{(Q,\gamma)}(\alpha) d\gamma \\ & \geq \int_{\lambda_2}^{\lambda_1} \int_{\mathcal{R}_1(\lambda_2)} \frac{1}{\gamma^2} \left(L_z(\alpha) - \int L_z(\nu) dP_{\Theta|Z=z}^{(Q,\gamma)}(\nu) \right) dP_{\Theta|Z=z}^{(Q,\gamma)}(\alpha) d\gamma \end{aligned} \quad (343)$$

where the equality in (343) 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 (261a), (261b), and (82), respectively.

The inequality in (343) 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 above arguments, 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

$$P_{\Theta|Z=z}^{(Q,\lambda_1)}(\mathcal{N}_{Q,z}(\lambda_2)) < P_{\Theta|Z=z}^{(Q,\lambda_2)}(\mathcal{N}_{Q,z}(\lambda_2)), \quad (344)$$

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 $(\nu_1, \nu_2) \in \mathcal{R}_2(\lambda) \times \mathcal{R}_1(\lambda)$,

$$L_z(\nu_1) > K_{Q,z}^{(1)} \left(-\frac{1}{\lambda} \right) > L_z(\nu_2), \quad (345)$$

which together with Definition 4.2 verify that the function L_z is separable with respect to $P_{\Theta|Z=z}^{(Q,\gamma)}$ (and thus, with respect to Q by Lemma 4.2). This ends the first part of the proof.

The second part of the proof is under the assumption that the empirical risk function L_z in (5) is separable with respect to Q (and thus, with respect to $P_{\Theta|Z=z}^{(Q,\gamma)}$ by Lemma 4.2). That is, from Definition 4.2, 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 $(\nu_1, \nu_2) \in \mathcal{A}(\gamma) \times \mathcal{B}(\gamma)$,

$$L_z(\nu_1) > c_\gamma > L_z(\nu_2). \quad (346)$$

In the proof of Theorem 7.1, c.f. (337), 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} ,

$$\begin{aligned} & P_{\Theta|Z=z}^{(Q,\alpha_1)}(\mathcal{A}) - P_{\Theta|Z=z}^{(Q,\alpha_2)}(\mathcal{A}) \\ &= \int_{\alpha_2}^{\alpha_1} \int_{\mathcal{A}} \frac{1}{\lambda^2} \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}(\alpha)}{dQ} \left(L_z(\alpha) - \int L_z(\nu) dP_{\Theta|Z=z}^{(Q,\lambda)}(\nu) \right) dP(\alpha) d\lambda \end{aligned} \quad (347)$$

$$= \int_{\alpha_2}^{\alpha_1} \int_{\mathcal{A}} \frac{1}{\lambda^2} \left(L_z(\alpha) - \int L_z(\nu) dP_{\Theta|Z=z}^{(Q,\lambda)}(\nu) \right) dP_{\Theta|Z=z}^{(Q,\lambda)}(\alpha) d\lambda. \quad (348)$$

Hence, two cases are studied. The first case considers that

$$c_\gamma < K_{Q,z}^{(1)} \left(-\frac{1}{\gamma} \right), \quad (349)$$

with c_γ in (346). The second case considers that

$$c_\gamma \geq K_{Q,z}^{(1)} \left(-\frac{1}{\gamma} \right). \quad (350)$$

In the first case, it follows from (82) that

$$\mathcal{B}(\gamma) \subset \mathcal{N}_{Q,z}(\gamma), \quad (351)$$

which implies that

$$P_{\Theta|Z=z}^{(Q,\gamma)}(\mathcal{N}_{Q,z}(\gamma)) \geq P_{\Theta|Z=z}^{(Q,\gamma)}(\mathcal{B}(\gamma)) \quad (352)$$

$$> 0, \quad (353)$$

where, the inequality in (353) 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 (82) and (351), it follows that for all $\alpha \in \mathcal{N}_{Q,z}(\gamma)$ and for all $\lambda \in (\gamma, \alpha_1)$,

$$L_z(\alpha) - \int L_z(\nu) dP_{\Theta|Z=z}^{(Q,\lambda)}(\nu) < L_z(\alpha) - c_\gamma \quad (354)$$

$$< 0, \quad (355)$$

where the inequality in (354) follows from (349); and the inequality in (355) follows from (346). Thus,

$$\int_{\gamma}^{\alpha_1} \int_{\mathcal{N}_{Q,z}(\gamma)} \frac{1}{\lambda^2} \left(\mathbb{L}_z(\alpha) - \int \mathbb{L}_z(\nu) dP_{\Theta|Z=z}^{(Q,\lambda)}(\nu) \right) dP_{\Theta|Z=z}^{(Q,\lambda)}(\alpha) d\lambda < 0, \quad (356)$$

which implies, from (348), that

$$P_{\Theta|Z=z}^{(Q,\alpha_1)}(\mathcal{N}_{Q,z}(\gamma)) - P_{\Theta|Z=z}^{(Q,\gamma)}(\mathcal{N}_{Q,z}(\gamma)) < 0. \quad (357)$$

Assume now that $c_{\gamma} \geq K_{Q,z}^{(1)} \left(-\frac{1}{\gamma}\right)$. Hence, the following holds

$$A(\gamma) \subseteq \mathcal{R}_2(\gamma), \quad (358)$$

which implies that

$$P_{\Theta|Z=z}^{(Q,\gamma)}(\mathcal{R}_2(\gamma)) \geq P_{\Theta|Z=z}^{(Q,\gamma)}(A(\gamma)) \quad (359)$$

$$> 0, \quad (360)$$

where the inequality in (360) follows from the fact that $A(\gamma)$ 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,

$$P_{\Theta|Z=z}^{(Q,\gamma)}(\mathcal{N}_{Q,z}(\gamma)) \geq P_{\Theta|Z=z}^{(Q,\gamma)}(\mathcal{R}_1(\gamma)) \quad (361)$$

$$> 0, \quad (362)$$

which proves the set $\mathcal{N}_{Q,z}(\gamma)$ is nonnegligible with respect to $P_{\Theta|Z=z}^{(Q,\gamma)}$.

From (82) and Theorem 5.1, it follows that for all $\alpha \in \mathcal{N}_{Q,z}(\gamma)$ and for all $\lambda \in (\gamma, \alpha_1)$,

$$0 \geq \mathbb{L}_z(\alpha) - \int \mathbb{L}_z(\nu) dP_{\Theta|\underline{X}=z=y}^{(\gamma)}(\nu) \quad (363)$$

$$> \mathbb{L}_z(\alpha) - \int \mathbb{L}_z(\nu) dP_{\Theta|\underline{X}=z=y}^{(\lambda)}(\nu). \quad (364)$$

Thus,

$$\int_{\gamma}^{\alpha_1} \int_{\mathcal{N}_{Q,z}(\gamma)} \frac{1}{\lambda^2} \left(\mathbb{L}_z(\alpha) - \int \mathbb{L}_z(\nu) dP_{\Theta|Z=z}^{(Q,\lambda)}(\nu) \right) dP_{\Theta|Z=z}^{(Q,\lambda)}(\alpha) d\lambda < 0, \quad (365)$$

which implies, from (348), that

$$P_{\Theta|Z=z}^{(Q,\alpha_1)}(\mathcal{N}_{Q,z}(\gamma)) - P_{\Theta|Z=z}^{(Q,\gamma)}(\mathcal{N}_{Q,z}(\gamma)) < 0. \quad (366)$$

The inequality $0 < P_{\Theta|Z=z}^{(Q,\lambda_1)}(\mathcal{N}_{Q,z}(\lambda_2))$ in (89) has already been proved while proving (88), and thus, this completes the proof of (89).

S Proof of Lemma 7.2

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 (261c). Second, note that

$$\mathcal{N}_{Q,z}(\lambda_1) = \mathcal{N}_{Q,z}(\lambda_2) \cup (\mathcal{N}_{Q,z}(\lambda_1) \cap \mathcal{R}_2(\lambda_2)), \quad (367)$$

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 $i \in \{1, 2\}$,

$$P_{\Theta|Z=z}^{(\lambda_i)}(\mathcal{N}_{Q,z}(\lambda_1)) = P_{\Theta|Z=z}^{(\lambda_i)}\left(\mathcal{N}_{Q,z}(\lambda_2) \cup (\mathcal{N}_{Q,z}(\lambda_1) \cap \mathcal{R}_2(\lambda_2))\right) \quad (368)$$

$$\begin{aligned} &= P_{\Theta|Z=z}^{(\lambda_i)}\left(\mathcal{N}_{Q,z}(\lambda_2)\right) \\ &\quad + P_{\Theta|Z=z}^{(\lambda_i)}\left(\mathcal{N}_{Q,z}(\lambda_1) \cap \mathcal{R}_2(\lambda_2)\right) \end{aligned} \quad (369)$$

$$= P_{\Theta|Z=z}^{(\lambda_i)}\left(\mathcal{N}_{Q,z}(\lambda_2)\right), \quad (370)$$

where the equality in (369) follows from Lemma 3.8 and the equality in (90).

Finally, under the assumption that the empirical function L_z in (5) is separable, it holds from Theorem 7.2 that

$$P_{\Theta|Z=z}^{(Q,\lambda_1)}(\mathcal{N}_{Q,z}(\lambda_2)) < P_{\Theta|Z=z}^{(Q,\lambda_2)}(\mathcal{N}_{Q,z}(\lambda_2)). \quad (371)$$

Plugging (370) into (371), with $i = 1$, yields,

$$P_{\Theta|Z=z}^{(Q,\lambda_1)}(\mathcal{N}_{Q,z}(\lambda_1)) < P_{\Theta|Z=z}^{(Q,\lambda_2)}(\mathcal{N}_{Q,z}(\lambda_2)), \quad (372)$$

and this completes the proof.

T Proof of Theorem 7.3

$$\begin{aligned} \lim_{\lambda \rightarrow 0^+} P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{N}_{Q,z}(\lambda)) &= \lim_{\lambda \rightarrow 0^+} P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_{Q,z}^*) \\ &\quad + \lim_{\lambda \rightarrow 0^+} P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{N}_{Q,z}(\lambda) \setminus \mathcal{L}_{Q,z}^*) \end{aligned} \quad (373)$$

$$= \lim_{\lambda \rightarrow 0^+} P_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_{Q,z}^*) \quad (374)$$

$$= 1 \quad (375)$$

where, the equalities in (374) follows Lemma 3.7. This completes the proof.

U Proof of Lemma 8.1

The cumulant generating function $J_{z,Q,\lambda}$ in (100) induced by the measure $P_{W|Z=z}^{(\lambda)}$ in (97) evaluated at t , with $t \leq \frac{1}{\lambda}$, is

$$J_{z,Q,\lambda}(t) = \log \left(\int \exp(t \mathbf{L}_z(\mathbf{u})) \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}(\mathbf{u})}{dQ}(\mathbf{u}) dQ(\mathbf{u}) \right) \quad (376)$$

$$= \log \left(\int \exp(t \mathbf{L}_z(\mathbf{u})) \exp \left(-K_{Q,z} \left(-\frac{1}{\lambda} \right) - \frac{1}{\lambda} \mathbf{L}_z(\mathbf{u}) \right) dQ(\mathbf{u}) \right) \quad (377)$$

$$= \log \left(\int \exp \left(\left(t - \frac{1}{\lambda} \right) \mathbf{L}_z(\mathbf{u}) - K_{Q,z} \left(-\frac{1}{\lambda} \right) \right) dQ(\mathbf{u}) \right) \quad (378)$$

$$= \log \left(\int \exp \left(\left(t - \frac{1}{\lambda} \right) \mathbf{L}_z(\mathbf{u}) - K_{Q,z} \left(t - \frac{1}{\lambda} \right) + K_{Q,z} \left(t - \frac{1}{\lambda} \right) - K_{Q,z} \left(-\frac{1}{\lambda} \right) \right) dQ(\mathbf{u}) \right) \quad (379)$$

$$= K_{Q,z} \left(t - \frac{1}{\lambda} \right) - K_{Q,z} \left(-\frac{1}{\lambda} \right) + \log \left(\int \exp \left(\left(t - \frac{1}{\lambda} \right) \mathbf{L}_z(\mathbf{u}) - K_{Q,z} \left(t - \frac{1}{\lambda} \right) \right) dQ(\mathbf{u}) \right) \quad (380)$$

$$= K_{Q,z} \left(t - \frac{1}{\lambda} \right) - K_{Q,z} \left(-\frac{1}{\lambda} \right) + \log \left(\int \frac{dP_{\Theta|Z=z}^{(Q, -\frac{1}{t-\frac{1}{\lambda}})}(\mathbf{u})}{dQ}(\mathbf{u}) dQ(\mathbf{u}) \right) \quad (381)$$

$$= K_{Q,z} \left(t - \frac{1}{\lambda} \right) - K_{Q,z} \left(-\frac{1}{\lambda} \right), \quad (382)$$

where the equality in (377) follows from Theorem 3.1; and the equality in (381) follows from the fact that $-\frac{1}{t-\frac{1}{\lambda}} > 0$ for all $t < \frac{1}{\lambda}$. This completes the proof.

V Proof of Lemma 10.1

From [35, Corollary 4.15, Page 100], it follows that the probability measures P and Q in $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ satisfy the following equality:

$$D(Q\|P) = \sup_f \int f(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) - \log \int \exp(f(\boldsymbol{\theta})) dP(\boldsymbol{\theta}), \quad (383)$$

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 $\int \exp(f(\boldsymbol{\theta})) dP(\boldsymbol{\theta}) < \infty$. Hence, for

all $\mathbf{u} \in (\mathcal{X} \times \mathcal{Y})^n$ and for all $t \in (-\infty, 0)$, it follows that the empirical risk function L_z in (5) satisfies that

$$D(Q\|P) \geq \int tL_z(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) - \log \int \exp(tL_z(\boldsymbol{\theta})) dP(\boldsymbol{\theta}) \quad (384)$$

$$\begin{aligned} &\geq \int tL_z(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) \\ &\quad - \log \int \exp(tL_z(\boldsymbol{\theta}) + tR_z(P) - tR_z(P)) dP(\boldsymbol{\theta}) \end{aligned} \quad (385)$$

$$\begin{aligned} &= \int tL_z(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) - tR_z(P) \\ &\quad - \log \int \exp(tL_z(\boldsymbol{\theta}) - tR_z(P)) dP(\boldsymbol{\theta}) \end{aligned} \quad (386)$$

$$= tR_z(Q) - tR_z(P) - \log \int \exp(tL_z(\boldsymbol{\theta}) - tR_z(P)) dP(\boldsymbol{\theta}), \quad (387)$$

which leads to

$$R_z(Q) - R_z(P) \leq \frac{D(Q\|P) + \log \int \exp(t(L_z(\boldsymbol{\theta}) - R_z(P))) dP(\boldsymbol{\theta})}{t}. \quad (388)$$

Given that t can be chosen arbitrarily in $(-\infty, 0)$, it holds that

$$R_z(Q) - R_z(P) \leq \inf_{t \in (-\infty, 0)} \frac{D(Q\|P) + \log \int \exp(t(L_z(\boldsymbol{\theta}) - R_z(P))) dP(\boldsymbol{\theta})}{t} \quad (389)$$

which completes the proof.

W Proof of Theorem 10.1

From Lemma 10.1, it holds that the probability measure $P_{\Theta|Z=z}^{(Q,\lambda)}$ in (25), satisfies

$$\begin{aligned} R_z(Q) - R_z(P_{\Theta|Z=z}^{(Q,\lambda)}) &\leq \inf_{t \in (-\infty, 0)} \left(\frac{D(Q\|P_{\Theta|Z=z}^{(Q,\lambda)})}{t} \right. \\ &\quad \left. + \frac{\log \left(\int \exp \left(t \left(L_z(\boldsymbol{\theta}) - K_{Q,z}^{(1)} \left(-\frac{1}{\lambda} \right) \right) \right) dP_{\Theta|Z=z}^{(Q,\lambda)}(\boldsymbol{\theta}) \right)}{t} \right), \end{aligned} \quad (390)$$

where the function $K_{Q,z}^{(1)}$ is defined in (64) and satisfies (68). Moreover, for all $t \in (-\infty, 0)$,

$$\begin{aligned} &\log \left(\int \exp \left(t \left(L_z(\boldsymbol{\theta}) - K_{Q,z}^{(1)} \left(-\frac{1}{\lambda} \right) \right) \right) dP_{\Theta|Z=z}^{(Q,\lambda)}(\boldsymbol{\theta}) \right) \\ &= \log \left(\int \exp(tL_z(\boldsymbol{\theta})) dP_{\Theta|Z=z}^{(Q,\lambda)}(\boldsymbol{\theta}) \right) - tK_{Q,z}^{(1)} \left(-\frac{1}{\lambda} \right) \end{aligned} \quad (391)$$

$$= J_{z,Q,\lambda}(t) - tK_{Q,z}^{(1)} \left(-\frac{1}{\lambda} \right) \quad (392)$$

$$\leq \frac{1}{2} t^2 B_z^2, \quad (393)$$

where the inequality in (392) follows from (100); the inequality in (393) follows from Theorem 8.1; and the constant B_z is defined in (105).

Plugging (393) into (390) yields for all $t \in (-\infty, 0)$,

$$\mathbb{R}_z(Q) - \mathbb{R}_z\left(P_{\Theta|Z=z}^{(Q,\lambda)}\right) \leq \inf_{t \in (-\infty, 0)} \frac{D\left(Q\|P_{\Theta|Z=z}^{(Q,\lambda)}\right) + \frac{1}{2}t^2 B_z^2}{t} : \quad (394)$$

Let the $c \in \mathbb{R}$ be defined as follows:

$$c \triangleq \mathbb{R}_z(Q) - \mathbb{R}_z\left(P_{\Theta|Z=z}^{(Q,\lambda)}\right). \quad (395)$$

Hence, from (394), it follows that for all $t \in (-\infty, 0)$,

$$ct - \frac{1}{2}t^2 B_z^2 \leq D\left(Q\|P_{\Theta|Z=z}^{(Q,\lambda)}\right). \quad (396)$$

The rest of the proof consists in finding an explicit expression for the absolute value of c . To this aim, consider the function $\phi : \mathbb{R} \rightarrow \mathbb{R}$ such that

$$\phi(\alpha) = \frac{1}{2}\alpha^2 B_z^2, \quad (397)$$

and note that ϕ is a positive and strictly convex function with $\phi(0) = 0$. Let the Legendre-Fenchel transform of ϕ be the function $\phi^* : \mathbb{R} \rightarrow \mathbb{R}$, and thus for all $x \in \mathbb{R}$,

$$\phi^*(x) = \max_{t \in (-\infty, 0)} xt - \phi(t). \quad (398)$$

In particular, note that

$$\phi^*(c) \leq D\left(Q\|P_{\Theta|Z=z}^{(Q,\lambda)}\right). \quad (399)$$

Note that for all $x \in \mathbb{R}$ and for all $t \in (-\infty, 0)$, the function ϕ^* in (398) satisfies

$$xt - \frac{1}{2}t^2 B_z^2 \leq \phi^*(x) = x\alpha^*(x) - \phi(\alpha^*(x)), \quad (400)$$

where the term $\alpha^*(x)$ represents the unique solution in α within the interval $(-\infty, 0)$ to

$$\frac{d}{d\alpha}(x\alpha - \phi(\alpha)) = x - \alpha B_z^2 = 0. \quad (401)$$

That is,

$$\alpha^*(x) = \frac{x}{B_z^2}. \quad (402)$$

Plugging (402) into (400) yields,

$$\phi^*(x) = \frac{x^2}{2B_z^2}. \quad (403)$$

Hence, from (399) and (400), given c in (395) for all $t \in (-\infty, 0)$,

$$ct - \frac{1}{2}t^2 B_z^2 \leq \phi^*(c) \leq D \left(Q \| P_{\Theta|Z=z}^{(Q,\lambda)} \right), \quad (404)$$

and thus,

$$\frac{c^2}{2B_z^2} \leq D \left(Q \| P_{\Theta|Z=z}^{(Q,\lambda)} \right). \quad (405)$$

This implies that

$$c \leq \sqrt{2B_z^2 D \left(Q \| P_{\Theta|Z=z}^{(Q,\lambda)} \right)} \quad (406)$$

and

$$c \geq -\sqrt{2B_z^2 D \left(Q \| P_{\Theta|Z=z}^{(Q,\lambda)} \right)}, \quad (407)$$

which leads to

$$\left| \int \mathbb{L}_z(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) - \int \mathbb{L}_z(\boldsymbol{\theta}) dP_{\Theta|Z=z}^{(Q,\lambda)}(\boldsymbol{\theta}) \right| \leq \sqrt{2B_z^2 D \left(Q \| P_{\Theta|Z=z}^{(Q,\lambda)} \right)}, \quad (408)$$

and completes the proof.

X Proof of Theorem 10.2

The optimization problem in (121) can be re-written as follows:

$$\min_{P \in \Delta_{P_{\Theta|Z=z}^{(Q,\lambda)}}(\mathcal{M}, \mathcal{B}(\mathcal{M}))} \int \mathbb{L}_z(\boldsymbol{\nu}) \frac{dP}{dP_{\Theta|Z=z}^{(Q,\lambda)}}(\boldsymbol{\nu}) dP_{\Theta|Z=z}^{(Q,\lambda)}(\boldsymbol{\nu}), \quad (409a)$$

$$\text{subject to: } \int \frac{dP}{dP_{\Theta|Z=z}^{(Q,\lambda)}}(\boldsymbol{\nu}) \log \left(\frac{dP}{dP_{\Theta|Z=z}^{(Q,\lambda)}}(\boldsymbol{\nu}) \right) dP_{\Theta|Z=z}^{(Q,\lambda)}(\boldsymbol{\nu}) \leq c, \quad (409b)$$

$$\int \frac{dP}{dP_{\Theta|Z=z}^{(Q,\lambda)}}(\boldsymbol{\theta}) dP_{\Theta|Z=z}^{(Q,\lambda)}(\boldsymbol{\theta}) = 1. \quad (409c)$$

Let \mathcal{M} be the set of nonnegative measurable functions with respect to the measurable spaces $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ and $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$. The Lagrangian of the optimization problem in (409) can be constructed in terms of a function in \mathcal{M} , instead of a measure over the measurable space $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$. Let such Lagrangian $L : \mathcal{M} \times [0, +\infty)^2 \rightarrow \mathbb{R}$ be of the form

$$\begin{aligned} L(g, \alpha, \beta) &= \int \mathbb{L}_z(\boldsymbol{\nu}) g(\boldsymbol{\nu}) dP_{\Theta|Z=z}^{(Q,\lambda)}(\boldsymbol{\nu}) \\ &\quad + \alpha \left(\int g(\boldsymbol{\nu}) \log(g(\boldsymbol{\nu})) dP_{\Theta|Z=z}^{(Q,\lambda)}(\boldsymbol{\nu}) - c \right) \\ &\quad + \beta \left(\int g(\boldsymbol{\nu}) dP_{\Theta|Z=z}^{(Q,\lambda)}(\boldsymbol{\nu}) - 1 \right), \end{aligned} \quad (410)$$

where g is a notation to represent the Radon-Nikodym derivative $\frac{dP}{dP_{\Theta|Z=z}^{(Q,\lambda)}}$; the reals α and β are both nonnegative and act as Lagrangian multipliers due to the constraint (409b) and (409c), respectively.

Let $h : \mathbb{R}^k \rightarrow \mathbb{R}$ be a function in \mathcal{M} . The Gateaux differential of the functional L in (410) at $(g, \alpha, \beta) \in \mathcal{M} \times [0, +\infty)^2$ in the direction of h is

$$\partial L(g, \alpha, \beta; h) \triangleq \left. \frac{d}{d\gamma} r(\gamma) \right|_{\gamma=0}, \quad (411)$$

where the real function $r : \mathbb{R} \rightarrow \mathbb{R}$ is such that for all $\gamma \in \mathbb{R}$,

$$\begin{aligned} r(\gamma) = & \int \mathbf{L}_z(\boldsymbol{\nu}) (g(\boldsymbol{\nu}) + \gamma h(\boldsymbol{\nu})) dP_{\Theta|Z=z}^{(Q,\lambda)}(\boldsymbol{\nu}) + \\ & \alpha \left(\int (g(\boldsymbol{\nu}) + \gamma h(\boldsymbol{\nu})) \log(g(\boldsymbol{\nu}) + \gamma h(\boldsymbol{\nu})) dP_{\Theta|Z=z}^{(Q,\lambda)}(\boldsymbol{\nu}) - c \right) \\ & + \beta \left(\int (g(\boldsymbol{\nu}) + \gamma h(\boldsymbol{\nu})) dP_{\Theta|Z=z}^{(Q,\lambda)}(\boldsymbol{\nu}) - 1 \right), \end{aligned} \quad (412)$$

Note that the derivative of the real function r in (412) is

$$\begin{aligned} \frac{d}{d\gamma} r(\gamma) = & \int \mathbf{L}_z h(\boldsymbol{\nu}) dP_{\Theta|Z=z}^{(Q,\lambda)}(\boldsymbol{\nu}) \\ & + \alpha \int h(\boldsymbol{\nu}) dP_{\Theta|Z=z}^{(Q,\lambda)}(\boldsymbol{\nu}) + \alpha \int h(\boldsymbol{\nu}) \log(g(\boldsymbol{\nu}) + \gamma h(\boldsymbol{\nu})) dP_{\Theta|Z=z}^{(Q,\lambda)}(\boldsymbol{\nu}) \\ & + \beta \int h(\boldsymbol{\nu}) dP_{\Theta|Z=z}^{(Q,\lambda)}(\boldsymbol{\nu}). \end{aligned} \quad (413)$$

From (411) and (413), it follows that

$$\partial L(g, \alpha, \beta; h) = \int h(\boldsymbol{\nu}) (\mathbf{L}_z(\boldsymbol{\nu}) + \alpha(1 + \log g(\boldsymbol{\nu})) + \beta) dP_{\Theta|Z=z}^{(Q,\lambda)}(\boldsymbol{\nu}). \quad (414)$$

From [32, Theorem 1, page 178], it holds that a necessary condition for the functional L in (410) to have a minimum at $(g, \alpha, \beta) \in \mathcal{M} \times [0, +\infty)^2$ is that for all functions $h \in C(\mathcal{M})$

$$\partial L(g, \alpha, \beta; h) = 0, \quad (415)$$

which implies that for all $\boldsymbol{\nu} \in \mathcal{M}$,

$$\mathbf{L}_z(\boldsymbol{\nu}) + \alpha(1 + \log g(\boldsymbol{\nu})) + \beta = 0. \quad (416)$$

Thus,

$$g(\boldsymbol{\nu}) = \exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\nu})}{\alpha}\right) \exp\left(-\frac{\beta + \alpha}{\alpha}\right), \quad (417)$$

where α and β are chosen to satisfy their corresponding constraints. Denote by P^* the solution of the optimization problem in (121). Hence, from (417), it follows that

$$\frac{dP^*}{dP_{\Theta|Z=z}^{(Q,\lambda)}}(\boldsymbol{\nu}) = \frac{\exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\nu})}{\alpha}\right)}{\int \exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\theta})}{\alpha}\right) dP_{\Theta|Z=z}^{(Q,\lambda)}(\boldsymbol{\theta})}, \quad (418)$$

where α is chosen to satisfy

$$D\left(P^* \parallel P_{\Theta|Z=z}^{(Q,\lambda)}\right) = c. \quad (419)$$

From Lemma 3.8, it follows that the probability measure P^* and the σ -finite measure Q satisfy,

$$\frac{dP^*}{dQ}(\boldsymbol{\nu}) = \frac{dP^*}{dP_{\Theta|Z=z}^{(Q,\lambda)}}(\boldsymbol{\nu}) \frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\boldsymbol{\nu}) \quad (420)$$

$$= \left(\frac{\exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\nu})}{\alpha}\right)}{\int \exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\theta})}{\alpha}\right) dP_{\Theta|Z=z}^{(Q,\lambda)}(\boldsymbol{\theta})} \right) \left(\frac{\exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\nu})}{\lambda}\right)}{\int \exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\theta})}{\lambda}\right) dQ(\boldsymbol{\theta})} \right) \quad (421)$$

$$= \left(\frac{\exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\nu})}{\alpha}\right)}{\int \frac{\exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\theta})}{\alpha}\right) \exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\theta})}{\lambda}\right)}{\int \exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\alpha})}{\lambda}\right) dQ(\boldsymbol{\alpha})} dQ(\boldsymbol{\theta})} \right) \left(\frac{\exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\nu})}{\lambda}\right)}{\int \exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\theta})}{\lambda}\right) dQ(\boldsymbol{\theta})} \right) \quad (422)$$

$$= \frac{\exp\left(-\left(\frac{1}{\alpha} + \frac{1}{\lambda}\right) \mathbf{L}_z(\boldsymbol{\nu})\right)}{\int \exp\left(-\left(\frac{1}{\alpha} + \frac{1}{\lambda}\right) \mathbf{L}_z(\boldsymbol{\nu})\right) dQ(\boldsymbol{\theta})}, \quad (423)$$

which implies that P^* is a Gibbs probability measure on $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$, with reference measure Q , regularization parameter $\frac{1}{\frac{1}{\alpha} + \frac{1}{\lambda}}$, and energy function \mathbf{L}_z . That is, for all $\boldsymbol{\nu} \in \text{supp } Q$,

$$P^*(\boldsymbol{\nu}) = P_{\Theta|Z=z}^{(Q, \frac{\alpha\lambda}{\alpha+\lambda})}(\boldsymbol{\nu}), \quad (424)$$

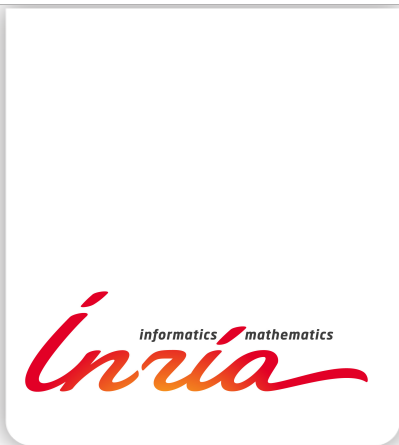
where α is chosen to satisfy (419). Let the positive real ω be $\omega \triangleq \frac{\alpha\lambda}{\alpha+\lambda}$ and note that $\omega \in (0, \lambda]$ and satisfies $D\left(P_{\Theta|Z=z}^{(Q,\omega)}(\nu) \| P_{\Theta|Z=z}^{(Q,\lambda)}\right) = c$. The proof ends by verifying that the objective function in (410) is strictly convex, and thus, the measure $P_{\Theta|Z=z}^{(Q,\omega)}$ is the unique minimizer. This completes the proof.

References

- [1] 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.
- [2] L. Zdeborová and F. Krzakala, “Statistical physics of inference: Thresholds and algorithms,” *Advances in Physics*, vol. 65, no. 5, pp. 453–552, Aug. 2016.
- [3] P. Alquier, J. Ridgway, and N. Chopin, “On the properties of variational approximations of Gibbs posteriors,” *The Journal of Machine Learning Research*, vol. 17, no. 1, pp. 8374–8414, 2016.
- [4] C. P. Robert, *The Bayesian choice: From decision-theoretic foundations to computational implementation*, 1st ed. New York, NY: Springer, 2007.
- [5] G. Aminian, Y. Bu, L. Toni, M. Rodrigues, and G. Wornell, “An exact characterization of the generalization error for the Gibbs algorithm,” in *Proc. Conference on Neural Information Processing Systems (NeurIPS)*, vol. 4, New Orleans, LA, USA, Dec. 2021, pp. 831–838.
- [6] M. Raginsky, A. Rakhlin, M. Tsao, Y. Wu, and A. Xu, “Information-theoretic analysis of stability and bias of learning algorithms,” in *Information Theory Workshop*, Cambridge, United Kingdom, Sep. 2016, pp. 26–30.
- [7] D. Russo and J. Zou, “How much does your data exploration overfit? Controlling bias via information usage,” *Transactions on Information Theory*, vol. 66, no. 1, pp. 302–323, Jan. 2019.
- [8] 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.
- [9] 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, 2020.
- [10] A. Xu and M. Raginsky, “Information-theoretic analysis of generalization capability of learning algorithms,” in *Proc. of the Thirty-first Conference on Neural Information Processing Systems (NeurIPS)*, Dec. 2017.
- [11] J. Shawe-Taylor and R. C. Williamson, “A PAC analysis of a Bayesian estimator,” in *Proceedings of the tenth annual conference on Computational learning theory*, 1997, pp. 2–9.
- [12] D. A. McAllester, “PAC-Bayesian stochastic model selection,” *Machine Learning*, vol. 51, no. 1, pp. 5–21, 2003.

- [13] M. Haddouche, B. Guedj, O. Rivasplata, and J. Shawe-Taylor, “PAC-Bayes unleashed: Generalisation bounds with unbounded losses,” *Entropy*, vol. 23, no. 10, 2021.
- [14] B. Guedj and L. Pujol, “Still no free lunches: The price to pay for tighter PAC-Bayes bounds,” *Entropy*, vol. 23, no. 11, 2021.
- [15] T. Jaakkola, M. Meila, and T. Jebara, “Maximum entropy discrimination,” *Neural Information Processing Systems*, 1999.
- [16] J. Zhu and E. P. Xing, “Maximum entropy discrimination Markov networks,” *Journal of Machine Learning Research*, vol. 10, no. 11, 2009.
- [17] Y. LeCun, 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.: The MIT Press, 2007, ch. 10, pp. 191–241.
- [18] C. P. Robert and G. Casella, *Monte Carlo statistical methods*, 2nd ed. New York, NY, USA: Springer, 2004.
- [19] J. N. Kapur, *Maximum Entropy Models in Science and Engineering*, 1st ed. New York, NY: Wiley, 1989.
- [20] T. M. Cover and J. A. Thomas, *Elements of Information Theory*, 2nd ed. Hoboken, NJ: Wiley-Interscience, 2006.
- [21] E. T. Jaynes, “Information theory and statistical mechanics I,” *Physical Review Journals*, vol. 106, pp. 620–630, May 1957.
- [22] —, “Information theory and statistical mechanics II,” *Physical Review Journals*, vol. 108, pp. 171–190, Oct. 1957.
- [23] 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, 1984.
- [24] 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.
- [25] B. Guedj, “A primer on PAC-Bayesian learning.” in *Tutorials of the International Conference on Machine Learning (ICML)*, Jun. 2019.
- [26] J. W. Gibbs, *Elementary principles in statistical mechanics*, 1st ed. New Haven, NJ: Yale University Press, 1902.
- [27] A. Dembo and O. Zeitouni, *Large Deviations Techniques and Applications*, 2nd ed. New York, NY: Springer-Verlag, 2009.
- [28] R. B. Ash and C. A. Doleans-Dade, *Probability and Measure Theory*, 2nd ed. Burlington, MA: Harcourt/Academic Press, 1999.

- [29] W. F. Trench, *Introduction to Real Analysis*, 1st ed. Hoboken, NJ: Prentice Hall/Pearson Education, 2003.
- [30] S. M. Perlaza, I. Esnaola, and H. V. Poor, “Sensitivity of the Gibbs algorithm to data aggregation in supervised machine learning,” Inria, Centre de Recherche de Sophia Antipolis Méditerranée, Sophia Antipolis, Tech. Rep. RR-9474, Jun. 2022.
- [31] D. P. Palomar and S. Verdú, “Lautum information,” *IEEE Transactions on Information Theory*, vol. 54, no. 3, pp. 964–975, Mar. 2008.
- [32] D. G. Luenberger, *Optimization by Vector Space Methods*, 1st ed. New York, NY: Wiley, 1997.
- [33] W. Feller, *An Introduction to Probability Theory and Its Applications*, 2nd ed. New York, NY: Jhon Wiley & Sons, 1971, vol. II.
- [34] W. Rudin, *Principles of mathematical analysis*, 1st ed. New York, NY: McGraw-Hill Book Company, Inc., 1953.
- [35] S. Boucheron, G. Lugosi, and P. Massart, *Concentration inequalities: A nonasymptotic theory of independence*, 1st ed. Oxford, UK: Oxford University Press, 2013.



**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