



**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-03560072v1

**HAL Id: hal-03560072**

**<https://hal.science/hal-03560072v1>**

Submitted on 7 Feb 2022 (v1), 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<sup>(1,2,3)</sup>, Gaetan Bisson<sup>(2)</sup>, Iñaki Esnaola<sup>(3,4)</sup>,  
Alain Jean-Marie<sup>(1)</sup>, and Stefano Rini<sup>(5)</sup>

**RESEARCH  
REPORT**

**N° 9454**

February 7, 2022

Project-Team NEO





# Empirical Risk Minimization with Generalized Relative Entropy Regularization

Samir M. Perlaza<sup>(1,2,3)</sup>, Gaetan Bisson<sup>(2)</sup>, Iñaki Esnaola<sup>(3,4)</sup>,  
Alain Jean-Marie<sup>(1)</sup>, and Stefano Rini<sup>(5)</sup>

Project-Team NEO

Research Report n° 9454 — February 7, 2022 — 66 pages

**Abstract:** The empirical risk minimization problem with relative entropy regularization (ERM-RER) is investigated considering that the reference measure is a  $\sigma$ -finite measure instead of a probability measure. This generalization allows for a larger degree of flexibility in the incorporation of prior knowledge over the set of models. In this setting, the interplay of the regularization parameter, the reference measure, the risk measure, and the expected empirical risk induced by the solution of the ERM-RER problem, which is proved to be unique, is characterized. For a fixed dataset, the empirical risk is shown to be a sub-Gaussian random variable, when the models follow the probability measure that is the solution to the ERM-RER problem. The sensitivity of the expected empirical risk to deviations from the solution of the ERM-RER problem is studied and upper and lower bounds on the expected empirical risk are provided. Finally, it is shown that the expectation of the sensitivity is 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.

Most of this work was done while Samir M. Perlaza was visiting the GAATI Laboratory of the University of French Polynesia between May 3 and October 31, 2021.

<sup>(1)</sup> INRIA, Centre de Recherche de Sophia Antipolis, Sophia Antipolis, France.

<sup>(2)</sup> GAATI Laboratory, University of French Polynesia, BP 6570, 98702 Faaa, French Polynesia.

<sup>(3)</sup> Department of Electrical and Computer Engineering, Princeton University, 08544 Princeton, N.J.

<sup>(4)</sup> Department of Automatic Control and Systems Engineering, University of Sheffield, Sheffield, United Kingdom.

<sup>(5)</sup> Department of Electrical and Computer Engineering, National Chao Tung University. Hsinchu, Taiwan.

**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  $\sigma$ -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

<b>1</b>	<b>Introduction</b>	<b>5</b>
1.1	Notation . . . . .	6
1.2	The Problem of Empirical Risk Minimization (ERM) . . . . .	6
1.3	Stochastic Solvers . . . . .	7
1.3.1	Expected Empirical Risk . . . . .	9
1.3.2	Expected Risk . . . . .	9
1.3.3	Generalization Gap . . . . .	10
1.4	Generalization Guarantees . . . . .	11
1.4.1	Gibbs Stochastic Solvers . . . . .	12
1.4.2	Empirical Risk Minimization with Relative Entropy Regularization . . . . .	13
1.4.3	Special Cases of the ERM-RER Problem . . . . .	13
1.5	Contributions . . . . .	14
<b>2</b>	<b>The Solution to the ERM Problem with Relative Entropy Regularization</b>	<b>15</b>
2.1	Negligible Sets and Coherent Measures . . . . .	18
2.2	The Partition Function, Cumulants and Separability . . . . .	19
2.2.1	The Mean of the Empirical Risk . . . . .	20
2.2.2	Second and Third Cumulants of the Empirical Risk . . . . .	22
2.3	Concentration of Probability . . . . .	23
2.3.1	Preliminaries . . . . .	24
2.3.2	The Limiting Set . . . . .	24
2.3.3	Probability of the Limiting Set . . . . .	26
2.4	Cumulant Generating Function of the empirical risk . . . . .	27
<b>3</b>	<b><math>(\delta, \epsilon)</math>-Optimality and Sensitivity</b>	<b>28</b>
3.1	$(\delta, \epsilon)$ -Optimality of Gibbs Stochastic Solvers . . . . .	29
3.2	Sensitivity of the Expected Empirical Risk . . . . .	30
<b>4</b>	<b>Discussion and Final Remarks</b>	<b>33</b>
	<b>Appendices</b>	<b>34</b>
<b>A</b>	<b>Proof of Lemma 2.1</b>	<b>34</b>
<b>B</b>	<b>Proof of Lemma 2.2</b>	<b>34</b>
<b>C</b>	<b>Proof of Theorem 2.1</b>	<b>35</b>
<b>D</b>	<b>Proof of Lemma 2.3</b>	<b>36</b>
<b>E</b>	<b>Proof of Lemma 2.5</b>	<b>37</b>
<b>F</b>	<b>Proof of Lemma 2.6</b>	<b>40</b>
<b>G</b>	<b>Proof of Lemma 2.7</b>	<b>40</b>
<b>H</b>	<b>Proof of Lemma 2.8</b>	<b>41</b>
<b>I</b>	<b>Proof of Lemma 2.9</b>	<b>42</b>

<b>J</b>	<b>Proof of Lemma 2.10</b>	<b>42</b>
<b>K</b>	<b>Proof of Theorem 2.2</b>	<b>45</b>
<b>L</b>	<b>Proof of Lemma 2.11</b>	<b>48</b>
<b>M</b>	<b>Proof of Lemma 2.12</b>	<b>49</b>
<b>N</b>	<b>Proof of Theorem 2.3</b>	<b>50</b>
<b>O</b>	<b>Proof of Lemma 2.13</b>	<b>50</b>
<b>P</b>	<b>Proof of Theorem 2.4</b>	<b>51</b>
<b>Q</b>	<b>Proof of Theorem 2.5</b>	<b>53</b>
<b>R</b>	<b>Proof of Lemma 2.15</b>	<b>57</b>
<b>S</b>	<b>Proof of Theorem 2.6</b>	<b>57</b>
<b>T</b>	<b>Proof of Lemma 2.16</b>	<b>58</b>
<b>U</b>	<b>Proof of Lemma 3.1</b>	<b>58</b>
<b>V</b>	<b>Proof of Theorem 3.3</b>	<b>59</b>
<b>W</b>	<b>Proof of Theorem 3.4</b>	<b>61</b>

## 1 Introduction

The problem of empirical risk minimization (ERM), also known as  $M$ -estimation [42], minimum contrast estimation [4, 29], sample average approximation [26], among other appellations, appears in numerous central problems in machine learning [45], information theory [32], statistics [47, 48], and operations research [10, 13], etc. Depending on the field, the ERM problem might take different forms and certain mathematical objects might take different names. In the following, ERM and the main results in this work, are presented in the context of supervised machine learning using the notation and nomenclature of this application field.

The empirical risk minimization (ERM) problem 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 and/or the models [2, 10, 51]. Instead of additional statistical assumptions, which is typical in Bayesian methods [36], relative entropy regularization requires a reference probability measure, which is external to the ERM problem. Often, such reference 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 and has been studied using information theoretic notions in [39, 50, 53]; statistical physics [10]; PAC (Probably Approximately Correct)-Bayesian learning theory [19, 20, 31, 41]; and proved to be of particular interest in classification problems in [22, 54].

In this report, the ERM with relative entropy regularization (ERM-RER) is generalized to incorporate a  $\sigma$ -finite measure with arbitrary support as the reference measure. That is, in this setting, the prior knowledge incorporated by the reference measure need not contain in its support the set of solutions of the ERM problem without regularization. The flexibility introduced by this formulation 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, as is often the case in practical scenarios [37]. Moreover, as shown below, for some specific choices of the reference measure, the ERM-RER boils down to particular cases of special interest: the information-risk minimization problem, the ERM with differential entropy regularization; and the ERM with discrete entropy regularization. Hence, the proposed formulation yields a unified mathematical framework that comprises a large class of problems. Nonetheless, despite its general character, the ERM-RER with an arbitrary  $\sigma$ -finite reference measure is shown to possess a unique solution. Indeed, the solution is a Gibbs probability measure whose partition function is defined with respect to the  $\sigma$ -finite reference measure. More importantly, the solution does not depend on the probability distribution of the dataset and can be computed for a particular dataset realization.

Finally, using the mathematical framework described above, the sensitivity of the ERM-RER problem is studied at a given dataset. The sensitivity is defined as the absolute value of the difference between two quantities: (a) The expectation of the empirical risk at a given dataset with respect to the measure that is the solution to the ERM-RER problem; and (b) the expectation of the empirical risk at such dataset with respect to another measure. The sensitivity is upper bounded by a term that is proportional to the squared-root of the relative entropy of the alternative measure with respect to the measure that is a solution to the ERM problem. More interestingly, the expectation of the sensitivity with respect to the probability distribution of the data sets turns out to be bounded by a term that is proportional to the squared-root of the lautum information between the models and the datasets. This bound is reminiscent to the result in [50] in which, under certain conditions, the generalization gap is upper bounded by a term that is proportional to the squared-root of the mutual information between the models and



the datasets. Together, the existing characterisations of the generalization gap and the proposed characterization of the sensitivity provide a solid theoretical framework to study the trade-off between generalization error and empirical risk.

## 1.1 Notation

In this work, sets are denoted by calligraphic letters. Given a set  $\mathcal{A}$ , the notation  $\mathcal{F}(\mathcal{A})$  represents a sigma-field ( $\sigma$ -field) on  $\mathcal{A}$ . When  $\mathcal{A} \subset \mathbb{R}^d$ , for some  $d \in \mathbb{N}$ , the Borel sigma-field on  $\mathcal{A}$  is denoted by  $\mathcal{B}(\mathcal{A})$ . The interior of  $\mathcal{A}$  is denoted  $\text{int}\mathcal{A}$ . The set of all measures that might be defined over a measurable space denoted by the tuple  $(\mathcal{A}, \mathcal{F}(\mathcal{A}))$  is denoted by  $\Delta(\mathcal{A}, \mathcal{F}(\mathcal{A}))$ . The support of the measure  $P$  is denoted by  $\text{supp } P$ . The *generalized relative entropy* is defined below as the extension to  $\sigma$ -finite measures of the relative entropy usually defined for probability measures.

**Definition 1.1** (Relative Entropy). *Given two  $\sigma$ -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*

$$D(Q\|P) = \int \frac{dQ}{dP}(x) \log \left( \frac{dQ}{dP}(x) \right) dP(x), \quad (1)$$

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

**Definition 1.2** (Mutual and Lautum Information). *Consider two random variables  $X$  and  $Y$  jointly inducing the probability measure  $P_{XY}$  in  $(\mathbb{R}^2, \mathcal{B}(\mathbb{R}^2))$ . Let the measure  $P_X$  and the measure  $P_Y$ , both on  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ , be the marginal probability measures of  $P_{XY}$ . The mutual information between  $X$  and  $Y$  is*

$$I(P_{XY}) = D(P_{XY}\|P_Y P_X); \quad (2)$$

and the lautum information [33] between  $X$  and  $Y$  is

$$L(P_{XY}) = D(P_Y P_X\|P_{XY}). \quad (3)$$

## 1.2 The Problem of 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,

$$Y = f(\theta^*, X). \quad (4)$$

The elements of the sets  $\mathcal{M}$ ,  $\mathcal{X}$  and  $\mathcal{Y}$  are often referred to as *models*, *patterns* and *labels*, respectively. A pair  $(x, y) \in \mathcal{X} \times \mathcal{Y}$  is referred to as a *labeled pattern*. In practice, patterns are high dimensional structures such as digital photos, medical images, digital text, speech signals, DNA sequences, etc. Labels are often countable sets, e.g., proper names, or uncountable sets, e.g., the reals numbers. A typical example of a labeled pattern  $(x, y)$  is the case in which  $x$  is a representation of a medical image and  $y$  is the representation of an indication of whether an anomaly is observed in such image.

The model  $\theta^*$  in (4), which is not necessarily unique and often referred to as the *ground truth model*, is unknown. Given a number of labelled patterns, the objective consists in obtaining a model  $\hat{\theta}$ , such that given a new pattern  $u \in \mathcal{X}$ , the estimated label  $f(\hat{\theta}, u)$  is equal, or as close as possible, to the true label  $f(\theta^*, u)$ . Establishing a distance between the labels  $f(\hat{\theta}, u)$  and

$f(\boldsymbol{\theta}^*, u)$  is made possible by the assumption that the set  $\mathcal{Y}$  can be equipped with a metric to form a metric space.

The *loss* or *risk* of adopting a model  $\boldsymbol{\nu} \in \mathcal{M}$  for a particular labelled pattern  $(x, y) \in \mathcal{X} \times \mathcal{Y}$  is determined by the function  $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow [0, +\infty)$ , which is often referred to as the *risk function*. More specifically, the model  $\boldsymbol{\theta}$  induces the risk  $\ell(f(\boldsymbol{\theta}, x), y)$  with respect to the labelled pattern  $(x, y) \in \mathcal{X} \times \mathcal{Y}$ . The risk function is chosen such that  $\ell(y, y) = 0$ . Nonetheless, it might exist other models  $\boldsymbol{\theta} \in \mathcal{M} \setminus \{\boldsymbol{\theta}^*\}$  such that  $\ell(f(\boldsymbol{\theta}, x'), y') = 0$  for some specific data points  $(x', y')$ .

The *empirical risk* induced by the model  $\boldsymbol{\theta}$ , with respect to a data set

$$\mathbf{z} = ((x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)) \in (\mathcal{X} \times \mathcal{Y})^n, \quad (5a)$$

is determined by the function  $L_{\mathbf{z}} : \mathcal{M} \rightarrow [0, +\infty)$ , which satisfies

$$L_{\mathbf{z}}(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \ell(f(\boldsymbol{\theta}, x_i), y_i). \quad (5b)$$

Thus, given the data set  $\mathbf{z}$  in (5a), the model  $\boldsymbol{\theta}$  is preferred against another model  $\boldsymbol{\mu} \in \mathcal{M}$ , if  $L_{\mathbf{z}}(\boldsymbol{\theta}) < L_{\mathbf{z}}(\boldsymbol{\mu})$ . Using this notation, the ERM problem, with respect to the data set  $\mathbf{z}$ , can be formulated as an optimization problem of the form [45]:

$$\min_{\boldsymbol{\theta} \in \mathcal{M}} L_{\mathbf{z}}(\boldsymbol{\theta}), \quad (5c)$$

whose solutions form the set denoted by

$$\mathcal{T}(\mathbf{z}) \triangleq \arg \min_{\boldsymbol{\theta} \in \mathcal{M}} L_{\mathbf{z}}(\boldsymbol{\theta}). \quad (5d)$$

The ground truth model  $\boldsymbol{\theta}^*$  in (4) satisfies that  $\boldsymbol{\theta}^* \in \mathcal{T}(\mathbf{z})$  and  $L_{\mathbf{z}}(\boldsymbol{\theta}^*) = 0$ . That is, the model  $\boldsymbol{\theta}^*$  is one of the solutions to the ERM problem in (5), and thus, the set  $\mathcal{T}(\mathbf{z})$  is not empty.

Despite the apparent simplicity of the ERM problem in (5), it is a difficult problem [46]. Often, the difficulty is due in part to the representation of the patterns, which are collections of heterogeneous classes of data with large dimensions; the function  $f$  in (4), which might be computationally difficult to calculate; the choice of the risk function, which is left at the discretion of the statistician; and the large number of data points whose processing requires unprecedented computing capabilities; among many others reasons. See for instance [5] and [21] for a detailed discussion.

### 1.3 Stochastic Solvers

Among the most popular methods for solving the ERM problem in (5) are those based on the gradient. The stochastic gradient descent algorithm [35] and its variants fall within this class of methods. See, for instance, the literature reviews in [8, 45, 49] and references therein. Other methods are based on constructing probability measures on  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$  that are used to choose a particular model from the set  $\mathcal{M}$ . This operation is known as *sampling* [37]. Methods based on the construction of probability measures include Bayesian methods [29, 36] and PAC-Bayesian methods [18, 40]. In order to describe these methods, it is assumed that independently of the choice of representation of patterns and labels, it is possible to form the following probability space:

$$(\mathcal{X} \times \mathcal{Y}, \mathcal{F}(\mathcal{X} \times \mathcal{Y}), P_{XY}). \quad (6)$$

Given the probability space in (6), a more formal definition of labeled pattern or data point is the following.

**Definition 1.3** (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 data set, as shown hereunder.

**Definition 1.4** (Data Set). *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 data set is represented by the tuple  $((x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)) \in (\mathcal{X} \times \mathcal{Y})^n$ .*

In the following, a central idea is that of *stochastic solvers*, also known as *stochastic estimators* in the realm of estimation theory [34]; and *samplers* in the realm of learning theory [10, 37].

A *stochastic solver* to the ERM problem in (5) is a probability measure on  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$  conditioned on an element of the set  $(\mathcal{X} \times \mathcal{Y})^n$ , i.e., a dataset.

**Definition 1.5** (Stochastic Solver). *A function  $g : \mathcal{B}(\mathcal{M}) \times (\mathcal{X} \times \mathcal{Y})^n \rightarrow [0, 1]$  is said to be a stochastic solver to the ERM problem in (5) if it satisfies the following conditions:*

(a) *For all  $\mathbf{u} \in (\mathcal{X} \times \mathcal{Y})^n$ , the function  $h_{\mathbf{u}} : \mathcal{B}(\mathcal{M}) \rightarrow [0, 1]$  that satisfies for all  $\mathcal{A} \in \mathcal{B}(\mathcal{M})$  that*

$$h_{\mathbf{u}}(\mathcal{A}) = g(\mathcal{A}, \mathbf{u}), \quad (7)$$

*is a probability measure on  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ ; and*

(b) *For all  $\mathcal{A} \in \mathcal{B}(\mathcal{M})$ , the function  $h_{\mathcal{A}} : (\mathcal{X} \times \mathcal{Y})^n \rightarrow [0, 1]$  that satisfies for all  $\mathbf{u} \in (\mathcal{X} \times \mathcal{Y})^n$  that*

$$h_{\mathcal{A}}(\mathbf{u}) = g(\mathcal{A}, \mathbf{u}), \quad (8)$$

*is measurable with respect to the measurable spaces  $((\mathcal{X} \times \mathcal{Y})^n, \mathcal{F}((\mathcal{X} \times \mathcal{Y})^n))$  and  $([0, 1], \mathcal{B}([0, 1]))$ .*

In the following, stochastic solvers are denoted by  $P_{\Theta|\mathbf{Z}}$ . Hence, for all  $\mathbf{u} \in (\mathcal{X} \times \mathcal{Y})^n$ , the notation  $P_{\Theta|\mathbf{Z}=\mathbf{u}}$  represents the probability measure from which models can be sampled from to approximate the solution to the ERM problem in (5) given the dataset  $\mathbf{u}$ . Moreover, given a set  $\mathcal{A} \in \mathcal{B}(\mathcal{M})$ , the function  $P_{\Theta|\mathbf{Z}}(\mathcal{A}|\cdot) : (\mathcal{X} \times \mathcal{Y})^n \rightarrow [0, 1]$  is Borel measurable with respect to the product measurable space  $((\mathcal{X} \times \mathcal{Y})^n, \mathcal{F}((\mathcal{X} \times \mathcal{Y})^n))$ . Hence, the notations  $P_{\Theta|\mathbf{Z}}(\mathcal{A}|\mathbf{u})$  or  $P_{\Theta|\mathbf{Z}=\mathbf{u}}(\mathcal{A})$  are indifferently used and represent the probability assigned by the stochastic solver  $P_{\Theta|\mathbf{Z}}$  to the set  $\mathcal{A}$  conditioned on the observation of the data set  $\mathbf{u}$ .

Using this notation, this work focuses in the case in which given a stochastic solver  $P_{\Theta|\mathbf{Z}}$  to the ERM problem in (5), a model  $\hat{\theta}$  is randomly selected from  $\mathcal{M}$  according to the probability measure  $P_{\Theta|\mathbf{Z}=\mathbf{z}}$ , with  $\mathbf{z}$  being the data set  $\mathbf{z}$  in (5a). Other choices for model selection might be, for instance, the first cumulant of the measure  $P_{\Theta|\mathbf{Z}=\mathbf{z}}$ , i.e.,

$$\hat{\theta} = \int \theta dP_{\Theta|\mathbf{Z}=\mathbf{z}}(\theta). \quad (9)$$

Another choice might be the maximum of the Radon-Nikodym derivative of the measure  $P_{\Theta|\mathbf{Z}=\mathbf{z}}$  with respect to a given measure  $P$  on the measurable space  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ , i.e.,

$$\hat{\theta} = \arg \max_{\theta \in \mathcal{M}} \frac{dP_{\Theta|\mathbf{Z}=\mathbf{z}}}{dP}(\theta), \quad (10)$$

if the Radon-Nikodym derivative and the maximum exist. The choice in (10) is reminiscent to *maximum à posteriori* estimation and *maximum likelihood* estimation, when  $P$  is chosen to be the Lebesgue measure on  $\mathcal{M}$ . See for instance, [34, Chapter IV]. Other choices of  $\hat{\theta}$  obtained from  $P_{\Theta|\mathbf{Z}=\mathbf{z}}$ , and the corresponding implications, are described in [18] and references therein.

In order to study the case in which a model  $\hat{\theta}$  is randomly chosen according to  $P_{\Theta|\mathbf{Z}=\mathbf{z}}$ , it is assumed that the function  $\bar{\ell} : \mathcal{M} \times \mathcal{X} \times \mathcal{Y} \rightarrow [0, +\infty)$ , which satisfies

$$\bar{\ell}(\theta, x, y) = \ell(f(\theta, x), y), \quad (11)$$

where the functions  $f$  and  $\ell$  are those in (4) and (5b), is measurable with respect to the measurable spaces  $(\mathcal{M} \times \mathcal{X} \times \mathcal{Y}, \mathcal{F}(\mathcal{M} \times \mathcal{X} \times \mathcal{Y}))$  and  $([0, +\infty), \mathcal{B}([0, +\infty)))$ . This assumption implies that for all  $\mathbf{u} \in (\mathcal{X} \times \mathcal{Y})^n$ , the empirical risk function  $\mathbf{L}_{\mathbf{u}}$  in (5b) is measurable with respect to the measurable spaces  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$  and  $([0, +\infty), \mathcal{B}([0, +\infty)))$ .

Under these assumptions, let the expected empirical risk induced by the stochastic solver  $P_{\Theta|Z}$  be defined as follows.

### 1.3.1 Expected Empirical Risk

The expected empirical risk is defined with respect to a dataset, as follows.

**Definition 1.6** (Expected Empirical Risk). *Given a data set  $\mathbf{u} \in (\mathcal{X} \times \mathcal{Y})^n$ , let the function  $\mathbf{R}_{\mathbf{u}} : \Delta(\mathcal{M}, \mathcal{B}(\mathcal{M})) \rightarrow [0, +\infty)$  be such that for all probability measures  $Q \in \Delta(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ ,*

$$\mathbf{R}_{\mathbf{u}}(Q) = \int \mathbf{L}_{\mathbf{u}}(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}), \quad (12)$$

where the function  $\mathbf{L}_{\mathbf{u}}$  is in (5b). The empirical risk induced by the measure  $Q$  at dataset  $\mathbf{u}$  is  $\mathbf{R}_{\mathbf{u}}(Q)$ .

Consider a stochastic solver  $P_{\Theta|Z}$  to the ERM problem with respect to the data set  $\mathbf{z}$  in (5), and assume that

$$P_{\Theta|Z=\mathbf{z}}(\mathcal{T}(\mathbf{z})) = 1, \quad (13)$$

with the set  $\mathcal{T}(\mathbf{z})$  in (5d). Hence,  $\mathbf{R}_{\mathbf{z}}(P_{\Theta|Z=\mathbf{z}}) = 0$ , i.e., the stochastic solver  $P_{\Theta|Z}$  achieves zero expected empirical risk with respect to the data set  $\mathbf{z}$  in (5). Nonetheless, given a data set  $\mathbf{u}$  different from  $\mathbf{z}$ , the stochastic solver  $P_{\Theta|Z}$  might induce a nonzero expected empirical risk. That is,  $\mathbf{R}_{\mathbf{u}}(P_{\Theta|Z=\mathbf{z}}) \geq \mathbf{R}_{\mathbf{z}}(P_{\Theta|Z=\mathbf{z}}) = 0$ . This implies that a desired property for any stochastic solver  $P_{\Theta|Z}$  is that for all data sets  $\mathbf{u} \in (\mathcal{X} \times \mathcal{Y})^n$ , the expected empirical risk  $\mathbf{R}_{\mathbf{u}}(P_{\Theta|Z=\mathbf{z}})$  be small and exhibit small changes with respect to changes in the data set  $\mathbf{u}$ . In order to formalize this desired property, the expected empirical risk is compared to the expected risk.

### 1.3.2 Expected Risk

The expected risk is defined as follows.

**Definition 1.7** (Expected Risk). *Let the function  $\mathbf{R} : \Delta(\mathcal{M}, \mathcal{B}(\mathcal{M})) \rightarrow [0, +\infty)$  be such that for all probability measures  $Q \in \Delta(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ ,*

$$\mathbf{R}(Q) = \int \left( \int \bar{\ell}(\boldsymbol{\theta}, u, v) dP_{XY}(u, v) \right) dQ(\boldsymbol{\theta}), \quad (14)$$

where the function  $\bar{\ell}$  is in (11) and the measure  $P_{XY}$  is specified in (6). The expected risk induced by the measure  $Q$  is  $\mathbf{R}(Q)$ .

The interest in the expected risk, i.e., the function  $\mathbf{R}$  in (14), lies on the fact that it establishes a metric of preference among stochastic solvers. More specifically, let  $P_{\Theta|Z}$  and  $Q_{\Theta|Z}$  be two stochastic solvers to the ERM in (5). The stochastic solver  $P_{\Theta|Z}$  is preferred against  $Q_{\Theta|Z}$ , if  $\mathbf{R}(P_{\Theta|Z=\mathbf{z}}) < \mathbf{R}(Q_{\Theta|Z=\mathbf{z}})$ . In particular, consider a stochastic solver  $P_{\Theta|Z}^*$  to the ERM problem with respect to the data set  $\mathbf{z}$  in (5) for which for all  $\mathbf{u} \in (\mathcal{X} \times \mathcal{Y})^n$ ,

$$P_{\Theta|Z=\mathbf{u}}^*(\mathcal{T}) = 1, \quad (15)$$

with  $\mathcal{T}$  being the following set,

$$\mathcal{T} = \arg \min_{\boldsymbol{\theta} \in \mathcal{M}} \int \bar{\ell}(\boldsymbol{\theta}, u, v) dP_{XY}(u, v), \quad (16)$$

where the function  $\bar{\ell}$  is in (11). The optimal model  $\boldsymbol{\theta}^*$  in (4) satisfies  $\boldsymbol{\theta}^* \in \mathcal{T}$ ; and the stochastic solver  $P_{\boldsymbol{\Theta}|Z}^*$  satisfies that for all  $\mathbf{u} \in (\mathcal{X} \times \mathcal{Y})^n$ ,

$$\mathbf{R}_{\mathbf{u}}(P_{\boldsymbol{\Theta}|Z=\mathbf{u}}^*) = \mathbf{R}(P_{\boldsymbol{\Theta}|Z=\mathbf{u}}^*) = 0. \quad (17)$$

In the context of the ERM problem in (5), this justifies the preference for  $P_{\boldsymbol{\Theta}|Z}^*$  against stochastic solvers  $P_{\boldsymbol{\Theta}|Z}$  for which  $\mathbf{R}(P_{\boldsymbol{\Theta}|Z=\mathbf{z}}) > 0$ , with  $\mathbf{z}$  being the data set in (5a). Unfortunately, such preference metric is purely theoretical, as in practice, the function  $\mathbf{R}$  in (14) and the stochastic solver in (15) cannot be calculated due to their dependence on the unknown probability measure  $P_{XY}$  in (6).

Despite the fact that the function  $\mathbf{R}$  in (14) cannot be calculated, it draws particular interest due to the following equality, which holds for all measures  $Q \in \Delta(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ ,

$$\mathbf{R}(Q) = \int \mathbf{R}_{\mathbf{u}}(Q) dP_Z(\mathbf{u}), \quad (18)$$

where the probability measure  $P_Z$  is a product measure in the product measurable space  $((\mathcal{X} \times \mathcal{Y})^n, \mathcal{F}((\mathcal{X} \times \mathcal{Y})^n))$ . That is, for all  $\mathcal{A} \in \mathcal{F}((\mathcal{X} \times \mathcal{Y})^n)$  of the form  $\mathcal{A} = \mathcal{A}_1 \times \mathcal{A}_2 \times \dots \times \mathcal{A}_n$ , such that for all  $i \in \{1, 2, \dots, n\}$ ,  $\mathcal{A}_i \in \mathcal{F}(\mathcal{X} \times \mathcal{Y})$ , it holds that

$$P_Z(\mathcal{A}) = \prod_{t=1}^n P_{XY}(\mathcal{A}_t). \quad (19)$$

More specifically, under the assumption that data points are independent and identically distributed according to  $P_{XY}$  in (6), the expected value with respect to  $P_Z$  of the expected empirical risk in (18) is equal to the expected risk in (14). Given a stochastic solver, this observation allows studying the difference between the expected empirical risk with respect to a given data set  $\mathbf{u}$  and the expected risk, which is defined as the *generalization gap* (GG).

### 1.3.3 Generalization Gap

The generalization gap is defined as follows.

**Definition 1.8** (Generalization Gap). *Let  $P_{\boldsymbol{\Theta}|Z}$  be a stochastic solver to the ERM problem in (5). Let also the function  $\mathbf{G} : (\mathcal{X} \times \mathcal{Y})^n \rightarrow \mathbb{R}$  be such that for all data sets  $\mathbf{u} \in (\mathcal{X} \times \mathcal{Y})^n$ ,*

$$\mathbf{G}(\mathbf{u}) = \mathbf{R}_{\mathbf{u}}(P_{\boldsymbol{\Theta}|Z=\mathbf{u}}) - \mathbf{R}(P_{\boldsymbol{\Theta}|Z=\mathbf{u}}), \quad (20)$$

where the functions  $\mathbf{R}_{\mathbf{u}}$  and  $\mathbf{R}$  are in (12) and (14), respectively. The generalization gap induced by the stochastic solver  $P_{\boldsymbol{\Theta}|Z}$  at the dataset  $\mathbf{u}$  is  $\mathbf{G}(\mathbf{u})$ .

The analysis of the generalization gap relies on the fact that the function  $\mathbf{G}$  in (20) is measurable with respect to the measurable spaces  $((\mathcal{X} \times \mathcal{Y})^n, \mathcal{F}((\mathcal{X} \times \mathcal{Y})^n))$  and  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ . This is a consequence of the fact that the function  $\bar{\ell}$  in (11) is assumed to be measurable with respect to the measurable spaces  $(\mathcal{M} \times \mathcal{X} \times \mathcal{Y}, \mathcal{F}(\mathcal{M} \times \mathcal{X} \times \mathcal{Y}))$  and  $([0, +\infty), \mathcal{B}([0, +\infty)))$ . The measurability of the function  $\mathbf{G}$  allows leveraging existing results in concentration inequalities for obtaining generalization guarantees.

## 1.4 Generalization Guarantees

Generalization guarantees are upper-bounds on the generalization gap, which can be of two types [1]. The first type consists in upper bounds on the expected value with respect to  $P_{\mathbf{Z}}$  in (19) of the absolute value of  $\mathbf{G}$ . See for instance, [39, 50, 52, 53]. The second type consists in upper bounds on the function  $\mathbf{G}$  that hold with high probability with respect to the measure  $P_{\mathbf{Z}}$ . Bounds of this type are described in [11, 18–20, 30, 52, 53] and are often referred to as *probably approximately correct* (PAC) guarantees for generalization. This appellation follows from the fact that probability measures satisfying such bounds, concentrate on sets containing models that induce small generalization gaps. Therefore, models sampled according to such measures are *approximately correct*, in the sense of small generalization gaps, with high probability. The notion of PAC can be understood as a generalization of the concept first introduced in [44] and the development of bounds of this kind is a central objective in the so-called PAC-Learning [40, Chapter 3].

The following examples illustrate these two types of guarantees.

**Example 1.1** (Theorem 1 in [50]). *Assume that there exists a real  $\sigma > 0$  such that for all  $\boldsymbol{\theta} \in \mathcal{M}$  and for all  $t \in \mathbb{R}$ ,*

$$\log \left( \int \exp(t \bar{\ell}(\boldsymbol{\theta}, u, v)) dP_{XY}(u, v) \right) \leq \frac{1}{2} \sigma^2 t^2 + t \int \bar{\ell}(\boldsymbol{\theta}, u, v) dP_{XY}(u, v), \quad (21)$$

with the function  $\bar{\ell}$  in (11). Then, for all stochastic solvers  $P_{\boldsymbol{\Theta}|\mathbf{Z}}$  to the ERM problem in (5), it holds that

$$\int |\mathbf{G}(\mathbf{u})| dP_{\mathbf{Z}}(\mathbf{u}) \leq \sqrt{\frac{2\sigma^2}{n} \int D(P_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{u}} \| P_{\boldsymbol{\Theta}}) dP_{\mathbf{Z}}(\mathbf{u})}, \quad (22)$$

where the data set  $\mathbf{z}$  is in (5a); the probability measure  $P_{\boldsymbol{\Theta}}$  is such that for all  $A \in \mathcal{B}(\mathcal{M})$ ,

$$P_{\boldsymbol{\Theta}}(A) = \int P_{\boldsymbol{\Theta}|\mathbf{Z}}(A|\mathbf{u}) dP_{\mathbf{Z}}(\mathbf{u}); \quad (23)$$

and the measure  $P_{\mathbf{Z}}$  is in (19).

**Example 1.2** (Theorem 2 in [30]). *Given a real  $\delta > 0$ , a stochastic solver  $P_{\boldsymbol{\Theta}|\mathbf{Z}}$  to the ERM problem in (5), and a probability measure  $P$  on  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$  such that  $P_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}$  is absolutely continuous with  $P$ , let  $\mathcal{R}(\delta, P_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}, P)$  be the set*

$$\mathcal{R}(\delta, P_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}, P) \triangleq \left\{ \mathbf{u} \in (\mathcal{X} \times \mathcal{Y})^n : \mathbf{G}(\mathbf{u}) < \sqrt{\frac{D(P_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}} \| P) + \ln \frac{2\sqrt{n}}{\delta}}{2n}} \right\}, \quad (24a)$$

where the dataset  $\mathbf{z}$  is in (5a). Then, for all  $\delta > 0$ , and for all stochastic solvers  $P_{\boldsymbol{\Theta}|\mathbf{Z}}$  to the ERM problem in (5) for which the measure  $P_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}$  is absolutely continuous with  $P$ , it holds that

$$P_{\mathbf{Z}}(\mathcal{R}(\delta, P_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}, P)) > 1 - \delta, \quad (24b)$$

where the product probability measure  $P_{\mathbf{Z}}$  is defined in (19).

In Example 1.1, the inequality in (21) implies that, for all  $\boldsymbol{\theta} \in \mathcal{M}$ , the random variable  $\bar{\ell}(\boldsymbol{\theta}, X, Y)$ , with  $X$  and  $Y$  being the random variables that induce the probability measure  $P_{XY}$  in (6), is  $\sigma$  sub-Gaussian. In such a case, for all stochastic solvers  $P_{\boldsymbol{\Theta}|\mathbf{Z}}$  for which

$$\int D(P_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{u}} \| \bar{P}_{\boldsymbol{\Theta}}) dP_{\mathbf{Z}}(\mathbf{u}) \leq \gamma, \quad (25)$$

with  $\gamma \in (0, +\infty)$ , it is always possible to arbitrarily reduce to zero the expectation with respect to  $P_{\mathbf{Z}}$  of the absolute value of  $\mathbf{G}$  in (20) by increasing the number  $n$  of data points in the data set  $\mathbf{z}$  in (5a). This observation is central as it justifies the fact that a stochastic solver that minimizes  $R_{\mathbf{z}}$  in (12) and satisfies the condition

$$\int |\mathbf{G}(\mathbf{u})| dP_{\mathbf{Z}}(\mathbf{u}) \leq \epsilon, \quad (26)$$

for some  $\epsilon > 0$  arbitrarily small, can be accepted as an approximation to a stochastic solver that minimizes the expected risk  $R$  in (14). The relevance of this observation stems from the fact that the function  $R$  cannot be directly minimized as it depends on the unknown probability measure  $P_{XY}$  in (6).

Alternatively, the expectation with respect to  $P_{\mathbf{Z}}$  of the absolute value of  $\mathbf{G}$  in (20) can be made arbitrarily small by reducing the dependence of the measure  $P_{\Theta|\mathbf{Z}=\mathbf{z}}$  on the data set  $\mathbf{z}$  in (5a). See, for instance [39] and [50] in which such dependence is studied using information measures. Consider for instance a stochastic solver  $P_{\Theta|\mathbf{Z}}$  that is independent of the data set  $\mathbf{z}$ . That is, for all data sets  $\mathbf{u} \in (\mathcal{X} \times \mathcal{Y})^n$ ,

$$P_{\Theta|\mathbf{Z}=\mathbf{u}} = Q, \quad (27)$$

for some probability measure  $Q$  on  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ . In this case,  $D(P_{\Theta|\mathbf{Z}=\mathbf{u}} || \bar{P}_{\Theta}) = 0$ , and thus, equality is met in (25) with  $\gamma = 0$ . Together with (22), this implies that for all data sets  $\mathbf{u} \in (\mathcal{X} \times \mathcal{Y})^n$ , the generalization gap at data set  $\mathbf{u}$  is zero, i.e.,  $\mathbf{G}(\mathbf{u}) = 0$ .

In a nutshell, two desired properties for a stochastic solver  $P_{\Theta|\mathbf{Z}}$  to the ERM problem in (5) are:

- (a) The stochastic solver  $P_{\Theta|\mathbf{Z}}$  induces a small expected empirical risk  $R_{\mathbf{z}}(P_{\Theta|\mathbf{Z}=\mathbf{z}})$ ; and
- (b) the relative entropy of  $P_{\Theta|\mathbf{Z}=\mathbf{z}}$  with respect to a probability measure independent of the data set  $\mathbf{z}$  is close to zero.

As shown by the previous examples, properties (a) and (b) can be simultaneously satisfied in the asymptotic regime of large data sets, under some specific conditions. This implies arbitrarily small generalization gap and small expected empirical risk. Nonetheless, in the case in which data sets contain only a finite number of labelled patterns, satisfying condition (a) and satisfying condition (b) are mutually opposing objectives. Hence, observing small generalization gaps does not imply small expected empirical risk. A class of stochastic solvers that is widely known to handle this tradeoff is that of the Gibbs stochastic solvers [16].

#### 1.4.1 Gibbs Stochastic Solvers

A Gibbs stochastic solver to the ERM problem in (5) is often defined with respect to another probability measure on the measurable space  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ . Nonetheless, in this work, a more general definition is adopted and Gibbs stochastic solvers are defined with respect to  $\sigma$ -finite measures. The following definition strengthens this observation.

**Definition 1.9** (Gibbs Stochastic Solver). *Let  $P$  be a  $\sigma$ -finite measure on the measurable space  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$  and let  $\lambda$  be a positive real. A stochastic solver of the ERM problem with respect to the data set  $\mathbf{z}$  in (5), denoted by  $P_{\Theta|\mathbf{Z}}^{(\lambda)}$ , is said to be a  $P$ -Gibbs stochastic solver with parameter  $\lambda$ , if for all  $\theta \in \mathcal{M}$ , it holds that,*

$$\frac{dP_{\Theta|\mathbf{Z}=\mathbf{z}}^{(\lambda)}}{dP}(\theta) = \frac{\exp\left(-\frac{L_{\mathbf{z}}(\theta)}{\lambda}\right)}{\int \exp\left(-\frac{L_{\mathbf{z}}(\nu)}{\lambda}\right) dP(\nu)}, \quad (28)$$

where the function  $L_z$  is defined in (5b).

Let  $P_{\Theta|Z}^{(\lambda)}$  be a Gibbs stochastic solver with respect to a given probability measure  $P$  with parameter  $\lambda > 0$ . Then, the main feature of  $P_{\Theta|Z}^{(\lambda)}$  is that the measure  $P_{\Theta|Z=z}^{(\lambda)}$ , with  $z$  in (5a), is a Gibbs measure with respect to the measure  $P$ , c.f., [6, 7, 17] and [14, Chapter 4].

### 1.4.2 Empirical Risk Minimization with Relative Entropy Regularization

One of the main properties of the  $P$ -Gibbs stochastic solver  $P_{\Theta|Z}^{(\lambda)}$  in (28), when  $P$  is a probability measure, is that the measure  $P_{\Theta|Z=z}^{(\lambda)}$  is the unique solution to the following optimization problem [11, 18, 50],

$$\min_{Q \in \Delta(\mathcal{M}, \mathcal{B}(\mathcal{M}))} R_z(Q) + \lambda D(Q \| P), \quad (29)$$

where the function  $R_z$  is specified in (12). In the following, the optimization problem in (29) is referred to as ERM with relative entropy regularization (ERM-RER) problem.

From the objective function of the optimization problem in (29), it is easy to see that the  $P$ -Gibbs stochastic solver  $P_{\Theta|Z}^{(\lambda)}$  in (28) trades off via the parameter  $\lambda$ , the two competing objectives described above: (a) minimizing the expected empirical risk with respect to the data set  $z$ , i.e., the function  $R_z$ ; and (b) approaching, in the sense of relative entropy, the measure  $P_{\Theta|Z=z}^{(\lambda)}$  to another measure, i.e.,  $P$ , which is independent of the data sets. This is in part, one of the reasons that justify the popularity of Gibbs stochastic solvers. Nonetheless, a prominent difficulty for adopting Gibbs stochastic solvers is the choice of the reference probability measure  $P$ , which is often restricted to probability measures.

### 1.4.3 Special Cases of the ERM-RER Problem

The nature of the set  $\mathcal{M}$  and the choice of the reference measure  $P$  lead to special cases of the ERM-RER problem in (29). Three cases are of particular interest: (a) The set  $\mathcal{M}$  satisfies  $\mathcal{M} \subseteq \mathbb{R}^d$ , with  $d \in \mathbb{N}$ ; and  $P$  is the Lebesgue measure on  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ ; (b) The set  $\mathcal{M}$  is countable; and the measure  $P$  is a counting measure; and (c) The set  $\mathcal{M}$  and the measure  $P$  form a probability measure space  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ .

In the former, the ERM-RER in (29) satisfies the following

$$\begin{aligned} \min_{Q \in \Delta(\mathcal{M}, \mathcal{B}(\mathcal{M}))} R_z(Q) + \lambda D(Q \| P) &= \min_{Q \in \Delta(\mathcal{M}, \mathcal{B}(\mathcal{M}))} \int L_z(\nu) \frac{dQ}{dP}(\nu) dP(\nu) \\ &\quad + \lambda \int \frac{dQ}{dP}(x) \log \left( \frac{dQ}{dP}(x) \right) dP(x) \end{aligned} \quad (30)$$

$$= \min_g \int_{\mathcal{M}} L_z(\nu) g(\nu) d\nu + \lambda \int_{\mathcal{M}} g(\theta) \log(g(\theta)) d\theta \quad (31)$$

$$= \min_g \int_{\mathcal{M}} L_z(\nu) g(\nu) d\nu - \lambda H(g), \quad (32)$$

where the Radon-Nikodym derivative  $\frac{dQ}{dP}$  in (30) is a probability density function (pdf) denoted by  $g$ , which implies that the optimization domain in (32) is the set of pdfs on  $\mathcal{M}$ . In (32), the notation  $H(g)$  represents the differential entropy of the pdf  $g$ , c.f., Chapter 8 in [12].

In case (b), a similar analysis would show that the ERM-RER problem in (29) boils down to the ERM with discrete entropy regularization (ERM-DisER). More specifically, the Radon-Nikodym



derivative  $\frac{dQ}{dP}$  in (30) is a probability mass function (pmf) denoted by  $p$ . This implies that the optimization domain is the set of pmfs on  $\mathcal{M}$ ; and the entropy is that of the pmf  $p$ , c.f., Chapter 2 in [12]. Both, the ERM-DifER and ERM-DisER problems are closely related to those typically arising while using Jayne’s maximum entropy principle [23–25]. In the case of classification problems, see for instance, [16, 22, 54].

Finally, case (c) is the classical formulation of ERM-RER and often referred to as the *information-risk minimization* (IRM) problem in information theory [53].

## 1.5 Contributions

In this report, the problem of ERM-RER with respect to the data set  $\mathbf{z}$  in (29) is studied by assuming that the reference measure  $P$  is a  $\sigma$ -finite measure on  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ , and not necessarily a probability measure. The motivation for this generalization is that some specific choices of the reference measure lead to special cases of the ERM-RER problem that are central in learning theory, e.g., the ERM with (differential or discrete) entropy; or the IRM problem. Moreover, this generalization allows for a larger degree of flexibility in the incorporation of prior knowledge over the set of models. In a nutshell, the proposed formulation yields a unified mathematical framework that comprises a large class of problems.

Under these assumptions, the most relevant results in this work are described hereunder:

- (i) The optimal solution to the ERM-RER with respect to the dataset  $\mathbf{z}$  in (29), when  $P$  is a  $\sigma$ -finite measure, is unique and forms a  $P$ -Gibbs stochastic solver  $P_{\Theta|\mathbf{Z}}^{(\lambda)}$ , where  $\lambda \in (0, b)$ , for some  $b \in \mathbb{R} \cup \{+\infty\}$ .
- (ii) The expected empirical risk  $R_{\mathbf{z}}(P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(\lambda)})$  is nondecreasing with  $\lambda$  (Theorem 2.2). The class of loss functions for which  $R_{\mathbf{z}}(P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(\lambda)})$  is monotone an increasing with  $\lambda$  is introduced.
- (iii) In the limit as  $\lambda$  tends to zero, the probability measure  $P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(\lambda)}$  is proved to concentrate on a set, often referred to as the *limiting set*. Depending on the choice of the reference measure  $P$ , the limiting set might be empty; identical to the set  $\mathcal{T}(\mathbf{z})$  in (5d); or exhibit an empty intersection with  $\mathcal{T}(\mathbf{z})$ . This is due to the flexibility in the choice of the reference measure in this formulation.
- (iv) The class of reference measures for which the limiting set is identical to the set of solutions to the ERM problem in (5), i.e., the set  $\mathcal{T}(\mathbf{z})$  in (5d), is characterized.
- (v) The empirical risk, when the models are distributed according to the optimal probability measure, is shown to be a sub-Gaussian random variable.
- (vi) Necessary and sufficient conditions for the existence of a regularization parameter that achieves an arbitrarily small empirical risk with arbitrarily high probability are presented.
- (vii) The sensitivity of the expected empirical risk to deviations from the solution of the ERM-RER problem is studied. The sensitivity is then used to provide upper and lower bounds on the expected empirical risk. Moreover, it is shown that the expectation of the sensitivity is upper bounded, up to a constant factor, by the square root of the lautum information between the models and the datasets.

These findings, together with other results described in the following sections, provide a deeper understanding of the ERM-RER problem and set a theoretical base for designing algorithms

based on Gibbs stochastic solvers whose reference measures are not necessarily probability measures.

## 2 The Solution to the ERM Problem with Relative Entropy Regularization

This section studies the ERM-RER problem with respect to the data set  $\mathbf{z}$  in (29) assuming that the reference measure  $P$  is a  $\sigma$ -finite measure on  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ ; the parameter  $\lambda > 0$  is fixed; and the dataset  $\mathbf{z}$  is the one in (5a). Under these assumptions, the solution is presented in terms of the function  $K_{\mathbf{z}} : \mathbb{R} \rightarrow \mathbb{R} \cup \{+\infty\}$  and the set  $\mathcal{K}_{\mathbf{z}} \subset (0, +\infty)$ , which are both parametrized by the data set  $\mathbf{z}$  and the measure  $P$ . The former satisfies that for all  $t \in \mathbb{R}$ ,

$$K_{\mathbf{z}}(t) = \log \left( \int \exp(t \mathbf{L}_{\mathbf{z}}(\boldsymbol{\theta})) dP(\boldsymbol{\theta}) \right), \quad (33)$$

with  $\mathbf{L}_{\mathbf{z}}$  in (5b), whereas, the latter satisfies that

$$\mathcal{K}_{\mathbf{z}} \triangleq \left\{ s > 0 : K_{\mathbf{z}} \left( -\frac{1}{s} \right) < +\infty \right\}. \quad (34)$$

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

**Lemma 2.1.** *The set  $\mathcal{K}_{\mathbf{z}}$  in (34) is either the empty set or a convex set that satisfies*

$$(0, b) \subset \mathcal{K}_{\mathbf{z}}, \quad (35)$$

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

*Proof.* The proof of Lemma 2.1 is presented in Appendix A.  $\square$

In the special case in which  $P$  is chosen to be a probability measure over  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ , it holds that for all  $t > 0$ , the value  $K_{\mathbf{z}} \left( -\frac{1}{t} \right)$  is upper bounded by zero. The following corollary formalizes this observation.

**Lemma 2.2.** *Assume that the measure  $P$  in (33) is a probability measure. Then, the set  $\mathcal{K}_{\mathbf{z}}$  in (34) satisfies*

$$\mathcal{K}_{\mathbf{z}} = (0, +\infty). \quad (36)$$

*Proof.* The proof of Lemma 2.2 is presented in Appendix B.  $\square$

Using this notation, the solution to the ERM-RER problem in (29) is presented by the following theorem.

**Theorem 2.1.** *Consider the ERM-RER problem with respect to the data set  $\mathbf{z}$  in (29), when  $P$  is a  $\sigma$ -finite measure on  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$  and  $\lambda \in \mathcal{K}_{\mathbf{z}}$ , with  $\mathcal{K}_{\mathbf{z}}$  in (34). Then, the solution to such problem is a unique measure on  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ , denoted by  $P_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(\lambda)}$ , whose Radon-Nikodym derivative with respect to  $P$  satisfies for all  $\boldsymbol{\theta} \in \text{supp } P$ ,*

$$\frac{dP_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(\lambda)}}{dP}(\boldsymbol{\theta}) = \exp \left( -K_{\mathbf{z}} \left( -\frac{1}{\lambda} \right) - \frac{1}{\lambda} \mathbf{L}_{\mathbf{z}}(\boldsymbol{\theta}) \right), \quad (37)$$

where, the function  $\mathbf{L}_{\mathbf{z}}$  is defined in (5b) and the function  $K_{\mathbf{z}}$  is defined in (33).

*Proof:* The proof of Theorem 2.1 is presented in Appendix C. ■

In Theorem 2.1, when the measure  $P$  is a probability measure, the measure  $P_{\Theta|Z=z}^{(\lambda)}$  is a Gibbs measure with parameter  $\lambda$  with respect to the function  $L_z$  and the measure  $P$ . Moreover, in such a case, the function  $K_z$  is often referred to as the *partition function*. See for instance, [13, Section 7.3.1]. In order not to disrupt with the current nomenclature, in the following, independently on whether  $P$  is a probability measure, both the function  $K_z$  and the probability measure  $P_{\Theta|Z=z}^{(\lambda)}$  are referred to as the partition function and the Gibbs measure with parameter  $\lambda$  with respect to the function  $L_z$  and the measure  $P$ . Similarly, in Definition 1.9, the notion of stochastic solver was introduced by considering that the measure  $P$  in the ERM-RER problem in (29) was a probability measure. In the following, as long as  $P$  is a  $\sigma$ -finite measure, a stochastic solver  $P_{\Theta|Z}^{(\lambda)}$  is said to be a  $P$ -Gibbs stochastic solver if the measure  $P_{\Theta|Z=z}^{(\lambda)}$  satisfies the equality in (37).

The Radon-Nikodym derivative  $\frac{dP_{\Theta|Z=z}^{(\lambda)}}{dP}$  in (37) is strictly positive and bounded, as shown by the following lemma.

**Lemma 2.3.** *Let the set  $\mathcal{T}(z)$  in (5d) and the  $\sigma$ -finite measure  $P$  in (37) be such that  $\mathcal{T}(z) \cap \text{supp } P \neq \emptyset$ . Then, for all  $\theta \in \text{supp } P$ , the Radon-Nikodym derivative  $\frac{dP_{\Theta|Z=z}^{(\lambda)}}{dP}$  in (37) satisfies for all  $(\theta_1, \theta_2) \in \mathcal{T}(z) \times \mathcal{T}(z)$  that*

$$\frac{dP_{\Theta|Z=z}^{(\lambda)}}{dP}(\theta) \leq \frac{dP_{\Theta|Z=z}^{(\lambda)}}{dP}(\theta_1) = \frac{dP_{\Theta|Z=z}^{(\lambda)}}{dP}(\theta_2). \quad (38)$$

*Proof:* The proof of Lemma 2.3 is presented in Appendix D. ■

The following lemmas describe the asymptotic behaviour of the Radon-Nikodym derivative  $\frac{dP_{\Theta|Z=z}^{(\lambda)}}{dP}$  in (37) when  $\lambda \rightarrow 0$  or  $\lambda \rightarrow +\infty$ .

**Lemma 2.4.** *Let the measure  $P$  in (37) be a probability measure. Then, for all  $\theta \in \text{supp } P$ , the Radon-Nikodym derivative  $\frac{dP_{\Theta|Z=z}^{(\lambda)}}{dP}$  in (37) satisfies*

$$\lim_{\lambda \rightarrow +\infty} \frac{dP_{\Theta|Z=z}^{(\lambda)}}{dP}(\theta) = 1. \quad (39)$$

*Proof:* From Theorem 2.1, it follows that for all  $\theta \in \text{supp } P$ ,

$$\lim_{\lambda \rightarrow +\infty} \frac{dP_{\Theta|Z=z}^{(\lambda)}}{dP}(\theta) = \lim_{\lambda \rightarrow +\infty} \frac{\exp\left(-\frac{L_z(\theta)}{\lambda}\right)}{\int \exp\left(-\frac{L_z(\nu)}{\lambda}\right) dP(\nu)} \quad (40)$$

$$= \frac{1}{\int dP(\nu)} \quad (41)$$

$$= 1, \quad (42)$$

which completes the proof. ■

Lemma 2.4 unveils the fact that for all  $\nu \in \text{supp } P$ ,

$$\lim_{\lambda \rightarrow +\infty} P_{\Theta|Z=z}^{(\lambda)}(\nu) = P(\nu). \quad (43)$$

This is consistent with the fact that when  $\lambda$  is arbitrarily large, the optimization problem in (29) boils down to exclusively minimizing the relative entropy. Such minimum is zero and is observed when both measures  $P_{\Theta|Z=z}^{(\lambda)}$  and  $P$  are identical.

From Lemma 2.1, it follows that, when  $P$  is not a probability measure, the set  $\mathcal{K}_z$  might be an interval of the form  $(0, b)$ , with  $b < \infty$ . Hence, in such a case, the analysis in which  $\lambda$  tends to infinity is void.

In the limit when  $\lambda$  tends to zero, the Radon-Nikodym derivative  $\frac{dP_{\Theta|Z=z}^{(\lambda)}}{dP}$  in (37) can be presented in terms of the set

$$\mathcal{L}_z(\delta) \triangleq \{\boldsymbol{\theta} \in \mathcal{M} : L_z(\boldsymbol{\theta}) \leq L_z(\hat{\boldsymbol{\theta}}) + \delta, \text{ with } \hat{\boldsymbol{\theta}} \in \mathcal{T}(z)\} \quad (44)$$

$$= \{\boldsymbol{\theta} \in \mathcal{M} : L_z(\boldsymbol{\theta}) \leq \delta\}, \quad (45)$$

with  $\delta \in [0, +\infty)$ . In particular consider the positive real

$$\delta^* \triangleq \inf \{\delta \in [0, +\infty) : P(\mathcal{L}_z(\delta)) > 0\}, \quad (46)$$

and let  $\mathcal{L}_z^*$  be the following set:

$$\mathcal{L}_z^* \triangleq \{\boldsymbol{\theta} \in \mathcal{M} : L_z(\boldsymbol{\theta}) = L_z(\boldsymbol{\theta}^*) + \delta^*, \text{ with } \boldsymbol{\theta}^* \in \mathcal{T}(z)\} \quad (47)$$

$$= \{\boldsymbol{\theta} \in \mathcal{M} : L_z(\boldsymbol{\theta}) = \delta^*\}. \quad (48)$$

**Lemma 2.5.** *If  $P(\mathcal{L}_z^*) > 0$ , with the set  $\mathcal{L}_z^*$  in (47) and  $P$  the  $\sigma$ -finite measure in (37), then for all  $\boldsymbol{\theta} \in \text{supp } P$ , the Radon-Nikodym derivative  $\frac{dP_{\Theta|Z=z}^{(\lambda)}}{dP}$  in (37) satisfies*

$$\lim_{\lambda \rightarrow 0^+} \frac{dP_{\Theta|Z=z}^{(\lambda)}}{dP}(\boldsymbol{\theta}) = \frac{1}{P(\mathcal{L}_z^*)} \mathbb{1}_{\{\boldsymbol{\theta} \in \mathcal{L}_z^*\}}. \quad (49)$$

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

$$\lim_{\lambda \rightarrow 0^+} \frac{dP_{\Theta|Z=z}^{(\lambda)}}{dP}(\boldsymbol{\theta}) = \begin{cases} +\infty & \text{if } \boldsymbol{\theta} \in \mathcal{L}_z^* \\ 0 & \text{otherwise.} \end{cases} \quad (50)$$

*Proof:* The proof of Lemma 2.5 is presented in Appendix E. ■

Lemma 2.5 describes a concentration phenomenon of the measure  $P_{\Theta|Z=z}^{(\lambda)}$  in (37) on the set  $\mathcal{L}_z^*$  in (47) as  $\lambda$  tends to zero. This is evident from the fact that in the limit, the Radon-Nikodym derivative  $\frac{dP_{\Theta|Z=z}^{(\lambda)}}{dP}$  is nonzero only over the elements of  $\mathcal{L}_z^*$ . This implication is formally proved in Lemma 2.12 below. Nonetheless, before a formal proof of the concentration of probability, it is interesting to study the set  $\mathcal{L}_z^*$  in some special cases.

Consider for instance the case in which for all  $\delta > 0$  it holds that  $P(\mathcal{L}_z(\delta)) > 0$ . Then,  $\delta^* = 0$  and  $\mathcal{L}_z^* = \mathcal{T}(z)$ , with  $\mathcal{T}(z)$  in (5d). Hence, in this particular case, the probability measure  $P_{\Theta|Z=z}^{(\lambda)}$  asymptotically concentrates on the set of solutions to the problem (5), as  $\lambda$  tends to zero. Interestingly, independently of whether  $P(\mathcal{T}(z)) = 0$  or  $P(\mathcal{T}(z)) > 0$ , the concentration phenomenon takes place. Note that measures that satisfy that for all  $\delta > 0$ , it holds that  $P(\mathcal{L}_z(\delta)) > 0$ , always exist. For instance, the Lebesgue measure when  $\mathcal{M}$  is uncountable; or counting measures when  $\mathcal{M}$  is countable. Another case of interest arises when the concentration

of probability occurs over a set that does not contain the set  $\mathcal{T}(\mathbf{z})$ . This occurs when  $\delta^* > 0$  and thus, for all  $\delta < \delta^*$ , it holds that  $P(\mathcal{L}_{\mathbf{z}}(\delta)) = 0$ . In this case, as  $\lambda$  tends to zero, the measure  $P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(\lambda)}$  concentrates on the set  $\mathcal{L}_{\mathbf{z}}^*$ , for which  $\mathcal{L}_{\mathbf{z}}^* \cap \mathcal{T}(\mathbf{z}) = \emptyset$ .

Finally, it is interesting to highlight a case in which the concentration of probability does not take place. Consider for instance, the case in which the  $\mathcal{L}_{\mathbf{z}}^*$  is empty. In such a case, for all  $\theta \in \text{supp } P$ , it holds that  $\lim_{\lambda \rightarrow 0^+} \frac{dP_{\Theta|\mathbf{Z}=\mathbf{z}}^{(\lambda)}}{dP}(\theta) = 0$ . This implies that, in the limit as  $\lambda$  tends to zero, the measure  $P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(\lambda)}$  is not a probability measure, but the trivial measure. This follows from the fact that in the context of the ERM-RER in (29), a  $\lambda$  arbitrarily small leads to exclusively minimizing the empirical risk, which in this case, does not have an optimum in  $\text{supp } P$ .

These observations highlight the influence of the reference measure  $P$  on the probability measure Gibbs measure  $P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(\lambda)}$ . A more in-depth analysis of this influence is presented in the following subsection via the negligible sets with respect to the measure  $P$ .

## 2.1 Negligible Sets and Coherent Measures

A central observation in the choice of the reference measure  $P$  in the ERM-RER in (29) is that it determines the negligible sets with respect to the measure  $P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(\lambda)}$  in (37). More specifically, for all sets  $\mathcal{C} \in \mathcal{B}(\mathcal{M})$ , it follows from Theorem 2.1 that if  $P(\mathcal{C}) = 0$ , then  $P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(\lambda)}(\mathcal{C}) = 0$ . The following lemma shows that the converse is also true.

**Lemma 2.6.** *The probability measures  $P$  and  $P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(\lambda)}$  in (37) are mutually absolutely continuous.*

*Proof:* The proof of Lemma 2.6 is presented in Appendix F. ■

The relevance of Lemma 2.6 is that it proves that for all  $\lambda \in \mathcal{K}_{\mathbf{z}}$ , the collection of negligible sets with respect to the measure  $P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(\lambda)}$  in (37) is identical to the collection of negligible sets with respect to the measure  $P$ . That is, for all subsets  $\mathcal{C} \in \mathcal{B}(\mathcal{M})$ ,

$$P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(\lambda)}(\mathcal{C}) > 0 \quad \text{if and only if} \quad P(\mathcal{C}) > 0. \quad (51)$$

An immediate consequence of these observations is the following.

**Lemma 2.7.** *For all  $(\alpha, \beta) \in \mathcal{K}_{\mathbf{z}} \times \mathcal{K}_{\mathbf{z}}$ , with  $\mathcal{K}_{\mathbf{z}}$  in (34), assume that the measures  $P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(\alpha)}$  and  $P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(\beta)}$  satisfy (37) with  $\lambda = \alpha$  and  $\lambda = \beta$ , respectively. Then,  $P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(\alpha)}$  and  $P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(\beta)}$  are mutually absolutely continuous.*

*Proof:* The proof of Lemma 2.7 is presented in Appendix G. ■

Lemma 2.7 proves that for all  $(\alpha, \beta) \in \mathcal{K}_{\mathbf{z}} \times \mathcal{K}_{\mathbf{z}}$ , the collection of negligible sets with respect to  $P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(\alpha)}$  and the collection of negligible sets with respect to  $P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(\beta)}$  are identical. This implies that the negligible sets with respect to the measure  $P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(\lambda)}$  in (37) are invariant with respect to the choice of  $\lambda \in \mathcal{K}_{\mathbf{z}}$ .

In the context of the ERM-RER problem with respect to the data set  $\mathbf{z}$  in (29), a desired condition for a  $P$ -Gibbs stochastic solver is that for all  $\delta > 0$ , the set  $\mathcal{L}_{\mathbf{z}}(\delta)$  in (44) possesses nonzero probability measure, i.e.,  $P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(\lambda)}(\mathcal{L}_{\mathbf{z}}(\delta)) > 0$ . More explicitly, it is desired that the models that induce empirical risks smaller than  $\delta$  concentrate most of the probability with respect to  $P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(\lambda)}$ . This is motivated by the fact that if  $P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(\lambda)}(\mathcal{L}_{\mathbf{z}}(\delta)) = 0$ , the expected empirical risk

$R_{\mathbf{z}}\left(P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(\lambda)}\right)$  in (12) is bounded away from zero, c.f., Lemma 2.11. In this context, Lemma 2.6 and Lemma 2.7 rise the need of choosing the reference measure  $P$  in (37) such that  $P(\mathcal{L}_{\mathbf{z}}(\delta)) > 0$ , for all  $\delta > 0$  arbitrarily small. Measures satisfying this condition are referred to as *coherent measures*.

**Definition 2.1.** *A measure  $P$  is said to be coherent with the ERM-RER problem with respect to the data set  $\mathbf{z}$  in (29), if for all  $\delta > 0$ ,*

$$P(\mathcal{L}_{\mathbf{z}}(\delta)) > 0, \quad (52)$$

where the set  $\mathcal{L}_{\mathbf{z}}(\delta)$  is defined in (44).

In the case in which  $P$  is coherent, the probability measure  $P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(\lambda)}$  in (37) satisfies for all  $\delta > 0$

$$P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(\lambda)}(\mathcal{L}_{\mathbf{z}}(\delta)) > 0, \quad (53)$$

as desired.

## 2.2 The Partition Function, Cumulants and Separability

This section introduces some properties of the partition function  $K_{\mathbf{z}}$  in (33), which are central in this work. The first property is concerned with the continuity and differentiability of the function  $K_{\mathbf{z}}$ .

**Lemma 2.8.** *The function  $K_{\mathbf{z}}$  in (33) is continuous and differentiable infinitely many times in  $(-\infty, 0)$ .*

*Proof:* The proof of Lemma 2.8 is presented in Appendix H. ■

More specific properties for the function  $K_{\mathbf{z}}$  in (33) can be stated for the case in which the empirical risk function  $L_{\mathbf{z}}$  in (5b) is separable with respect to the measure  $P$  in (37).

**Definition 2.2** (Separable Empirical Risk Function). *The empirical risk function  $L_{\mathbf{z}}$  in (5b) is said to be separable with respect to the  $\sigma$ -finite measure  $P$  in (37), if there exist a positive real  $c > 0$  and two subsets  $\mathcal{A}$  and  $\mathcal{B}$  of  $\mathcal{M}$  that are nonnegligible with respect to  $P$ , such that for all  $(\theta_1, \theta_2) \in \mathcal{A} \times \mathcal{B}$ ,*

$$L_{\mathbf{z}}(\theta_1) > c > L_{\mathbf{z}}(\theta_2). \quad (54)$$

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

$$P(\{\theta \in \mathcal{M} : L_{\mathbf{z}}(\theta) = a\}) = 1. \quad (55)$$

From this perspective, nonseparable empirical risk functions exhibit little practical interest. This follows from observing that models sampled from the probability measure  $P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(\lambda)}$  in (37) induce the same empirical risk.

The definition of separability in Definition 2.2 is in terms of the reference measure  $P$  in (37). Nonetheless, Lemma 2.6 provides an alternative definition in terms of the measure  $P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(\lambda)}$  in (37), with  $\lambda \in \mathcal{K}_{\mathbf{z}}$ .

**Definition 2.3.** The empirical risk function  $L_z$  in (5b) is said to be separable with respect to a  $\sigma$ -finite measure  $P$ , if and only if there exist a real  $\lambda \in \mathcal{K}_z$ , with  $\mathcal{K}_z$  in (34); a positive real  $c > 0$ ; and two subsets  $\mathcal{A}$  and  $\mathcal{B}$  of  $\mathcal{M}$  nonnegligible with respect to the probability measure  $P_{\Theta|Z=z}^{(\lambda)}$  in (37), such that for all  $(\nu_1, \nu_2) \in \mathcal{A} \times \mathcal{B}$ , the inequality in (54) holds.

The following lemma presents a general property of the function  $K_z$  in (33) for the case of separable empirical risk functions.

**Lemma 2.9.** The function  $K_z$  in (33) is convex. If the empirical risk function  $L_z$  in (5b) is separable with respect to the  $\sigma$ -finite measure  $P$  in (33), then the function  $K_z$  is strictly convex.

*Proof:* The proof of Lemma 2.9 is presented in Appendix I. ■

Let the  $m$ -th derivative of the function  $K_z$  in (33) be denoted by  $K_z^{(m)} : \mathbb{R} \rightarrow \mathbb{R}$ , with  $m \in \mathbb{N}$ . Hence, for all  $s \in \mathcal{K}_z$ ,

$$K_z^{(m)}\left(-\frac{1}{s}\right) \triangleq \frac{d^m}{dt^m} K_z(t) \Big|_{t=-\frac{1}{s}}. \quad (56)$$

The following lemma provides explicit expressions for the first, second and third derivatives of the function  $K_z$  in (33).

**Lemma 2.10.** The first, second and third derivatives of the function  $K_z$  in (33) evaluated at  $-\frac{1}{\lambda}$ , with  $\lambda \in \text{int}\mathcal{K}_z$  and  $\mathcal{K}_z$  in (34), satisfy the following equalities,

$$K_z^{(1)}\left(-\frac{1}{\lambda}\right) = \int L_z(\theta) dP_{\Theta|Z=z}^{(\lambda)}(\theta), \quad (57)$$

$$K_z^{(2)}\left(-\frac{1}{\lambda}\right) = \int \left( L_z(\theta) - K_z^{(1)}\left(-\frac{1}{\lambda}\right) \right)^2 dP_{\Theta|Z=z}^{(\lambda)}(\theta), \text{ and} \quad (58)$$

$$K_z^{(3)}\left(-\frac{1}{\lambda}\right) = \int \left( L_z(\theta) - K_z^{(1)}\left(-\frac{1}{\lambda}\right) \right)^3 dP_{\Theta|Z=z}^{(\lambda)}(\theta), \quad (59)$$

where the function  $L_z$  is defined in (5b) and the measure  $P_{\Theta|Z=z}^{(\lambda)}$  satisfies (37).

*Proof:* The proof of Lemma 2.10 is presented in Appendix J. ■

From Lemma 2.10, it follows that if  $\Theta$  is the random vector that induces the measure  $P_{\Theta|Z=z}^{(\lambda)}$  in (37), with  $\lambda \in \mathcal{K}_z$ , the empirical risk  $L_z$  in (5b) becomes the random variable

$$W \triangleq L_z(\Theta), \quad (60)$$

whose mean, variance, and third cumulant are respectively  $K_z^{(1)}\left(-\frac{1}{\lambda}\right)$ ,  $K_z^{(2)}\left(-\frac{1}{\lambda}\right)$ , and  $K_z^{(3)}\left(-\frac{1}{\lambda}\right)$ .

### 2.2.1 The Mean of the Empirical Risk

The mean of the random variable  $W$  in (60) is equivalent to the expected empirical risk with respect to the data set  $z$  induced by  $P_{\Theta|Z=z}^{(\lambda)}$ , i.e., the value  $R_z\left(P_{\Theta|Z=z}^{(\lambda)}\right)$  in (12), as shown by the following corollary.

**Corollary 2.1.** The probability measure  $P_{\Theta|Z=z}^{(\lambda)}$  in (37) verifies that

$$R_z\left(P_{\Theta|Z=z}^{(\lambda)}\right) = K_z^{(1)}\left(-\frac{1}{\lambda}\right), \quad (61)$$

where the functions  $R_z$  and  $K_z^{(1)}$  are defined in (12) and (57), respectively.

The following theorem shows an interesting property of the value  $R_z \left( P_{\Theta|Z=z}^{(\lambda)} \right)$ .

**Theorem 2.2.** *The value  $R_z \left( P_{\Theta|Z=z}^{(\lambda)} \right)$  in (61) is nondecreasing with  $\lambda \in \mathcal{K}_z$ . Moreover, the function  $L_z$  in (5b) is separable with respect to the measure  $P$  if and only if the value  $R_z \left( P_{\Theta|Z=z}^{(\lambda)} \right)$  is strictly increasing with  $\lambda \in \mathcal{K}_z$ .*

*Proof:* The proof of Theorem 2.2 is presented in Appendix K. ■

A question that arises from Theorem 2.2 is whether the value  $R_z \left( P_{\Theta|Z=z}^{(\lambda)} \right)$  in (61) can be made arbitrarily close to  $L_z(\theta^*) = 0$ , with  $\theta^*$  in (4), by making  $\lambda$  arbitrarily small. The following lemma shows that there exist cases in which the value  $R_z \left( P_{\Theta|Z=z}^{(\lambda)} \right)$  is bounded away from zero, even for arbitrarily small values of  $\lambda$ .

**Lemma 2.11.** *For all  $\lambda \in \mathcal{K}_z$ , with  $\mathcal{K}_z$  in (34), the function  $K_z^{(1)}$  in (57) satisfies,*

$$K_z^{(1)} \left( -\frac{1}{\lambda} \right) \geq \delta^*, \quad (62)$$

where  $\delta^*$  is defined in (46). Moreover, the function  $L_z$  in (5b) is separable with respect to the measure  $P$  if and only if the inequality in (62) is strict.

*Proof:* The proof of Lemma 2.11 is presented in Appendix L. ■

In the limit as  $\lambda$  tends to zero, a probability concentration phenomena takes place. This phenomena is studied in detail in Section 2.3. Nonetheless, the following lemma leads to preliminary intuitions.

**Lemma 2.12.** *The measure  $P_{\Theta|Z=z}^{(\lambda)}$  in (37) and the set  $\mathcal{L}_z^*$  in (47) satisfy,*

$$\lim_{\lambda \rightarrow 0^+} P_{\Theta|Z=z}^{(\lambda)}(\mathcal{L}_z^*) = \begin{cases} 0 & \text{if } \mathcal{L}_z^* = \emptyset \\ 1 & \text{if } \mathcal{L}_z^* \neq \emptyset. \end{cases} \quad (63)$$

*Proof:* The proof of Lemma 2.12 is presented in Appendix M. ■

Lemma 2.12 reveals the fact that as  $\lambda$  tends to zero, the probability measure  $P_{\Theta|Z=z}^{(\lambda)}$  concentrates in the set  $\mathcal{L}_z^*$ . An immediate consequence of this observation is the following theorem.

**Theorem 2.3.** *The function  $K_z^{(1)}$  in (57) satisfies,*

$$\lim_{\lambda \rightarrow 0^+} K_z^{(1)} \left( -\frac{1}{\lambda} \right) = \begin{cases} 0 & \text{if } \mathcal{L}_z^* = \emptyset \\ \delta^* & \text{if } \mathcal{L}_z^* \neq \emptyset. \end{cases} \quad (64)$$

where  $\delta^*$  is defined in (46).

*Proof:* The proof of Theorem 2.3 is presented in Appendix N. ■

Theorem 2.3 unveils the relevance of coherent measures (Definition 2.1). More specifically, the expected empirical risk with respect to the dataset  $z$  induced by  $P_{\Theta|Z=z}^{(\lambda)}$ , i.e., the value  $R_z \left( P_{\Theta|Z=z}^{(\lambda)} \right)$  in (12), can be made arbitrarily close to zero, if the reference measure  $P$  is coherent.

The anomalous case in which  $\mathcal{L}_z^*$  is empty leads to a zero expected empirical risk in the limit as  $\lambda$  tends to zero, i.e.,  $\lim_{\lambda \rightarrow 0^+} R_z \left( P_{\Theta|Z=z}^{(\lambda)} \right) = 0$ . Nonetheless, in such a case, in the limit



as  $\lambda$  tends to zero, the measure  $P_{\Theta|Z=z}^{(\lambda)}$  is the trivial measure. That is, all sets in  $\mathcal{F}(\mathcal{M})$  exhibit zero measure with respect to  $P_{\Theta|Z=z}^{(\lambda)}$ . Hence, this is an anomaly instead of an optimality guarantee.

## 2.2.2 Second and Third Cumulants of the Empirical Risk

The monotonicity of the function  $K_z^{(1)}$ , stated by Theorem 2.2, is not a property exhibit by the functions  $K_z^{(2)}$  and  $K_z^{(3)}$ , respectively specified by (58) and (59). This section highlights this observation via the following example.

**Example 2.1.** Consider the ERM-RER problem with respect to the data set  $\mathbf{z}$  in (29), under the assumption that  $P$  is a probability measure and the empirical risk function is of the form

$$\mathsf{L}_z(\boldsymbol{\theta}) = \begin{cases} 0 & \text{if } \boldsymbol{\theta} \in \mathcal{A} \\ 1 & \text{if } \boldsymbol{\theta} \in \mathcal{M} \setminus \mathcal{A}, \end{cases} \quad (65)$$

where, the sets  $\mathcal{A}$  and  $\mathcal{M} \setminus \mathcal{A}$  are nonnegligible with respect to the reference probability measure  $P$  in (29). Hence, from (33), the following holds for all  $\lambda > 0$ ,

$$K_z\left(-\frac{1}{\lambda}\right) = \log\left(P(\mathcal{A}) + \exp\left(-\frac{1}{\lambda}\right)(1 - P(\mathcal{A}))\right). \quad (66)$$

The derivatives  $K_z^{(1)}$ ,  $K_z^{(2)}$ , and  $K_z^{(3)}$  in (56) of the function  $K_z$  in (66) are:

$$K_z^{(1)}\left(-\frac{1}{\lambda}\right) = \frac{\exp\left(-\frac{1}{\lambda}\right)(1 - P(\mathcal{A}))}{P(\mathcal{A}) + \exp\left(-\frac{1}{\lambda}\right)(1 - P(\mathcal{A}))}; \quad (67)$$

$$K_z^{(2)}\left(-\frac{1}{\lambda}\right) = \frac{P(\mathcal{A})(1 - P(\mathcal{A}))\exp\left(-\frac{1}{\lambda}\right)}{\left(P(\mathcal{A}) + \exp\left(-\frac{1}{\lambda}\right)(1 - P(\mathcal{A}))\right)^2}; \text{ and} \quad (68)$$

$$K_z^{(3)}\left(-\frac{1}{\lambda}\right) = \frac{P(\mathcal{A})(1 - P(\mathcal{A}))\exp\left(-\frac{1}{\lambda}\right)\left(P(\mathcal{A}) - (1 - P(\mathcal{A}))\exp\left(-\frac{1}{\lambda}\right)\right)}{\left(P(\mathcal{A}) + \exp\left(-\frac{1}{\lambda}\right)(1 - P(\mathcal{A}))\right)^3}. \quad (69)$$

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

$$P(\mathcal{A}) - (1 - P(\mathcal{A}))\exp\left(-\frac{1}{\lambda}\right) > 0. \quad (70)$$

Assume that  $P(\mathcal{A}) \geq \frac{1}{2}$ . Thus, it holds that for all  $\lambda > 0$ , the inequality in (70) is always satisfied. This follows from observing that for all  $\lambda > 0$ ,

$$\exp\left(-\frac{1}{\lambda}\right) < 1 \leq \frac{P(\mathcal{A})}{1 - P(\mathcal{A})}. \quad (71)$$

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

$$\frac{1}{4} \geq K_z^{(2)}\left(-\frac{1}{\lambda_1}\right) > K_z^{(2)}\left(-\frac{1}{\lambda_2}\right) > \dots > 0. \quad (72)$$

Alternatively, assume that  $P(\mathcal{A}) < \frac{1}{2}$ . In this case, the inequality in (70) is satisfied if and only if

$$\lambda < \left(\log\left(\frac{1 - P(\mathcal{A})}{P(\mathcal{A})}\right)\right)^{-1}. \quad (73)$$

Hence, if  $P(\mathcal{A}) < \frac{1}{2}$ , then for all decreasing sequences of positive reals  $(\log(\frac{1-P(\mathcal{A})}{P(\mathcal{A})}))^{-1} > \lambda_1 > \lambda_2 > \dots > 0$ , it holds that

$$\frac{1}{4} > K_{\mathbf{z}}^{(2)}\left(-\frac{1}{\lambda_1}\right) > K_{\mathbf{z}}^{(2)}\left(-\frac{1}{\lambda_2}\right) > \dots > 0. \quad (74)$$

Moreover, for all decreasing sequences of positive reals  $\lambda_1 > \lambda_2 > \dots > (\log(\frac{1-P(\mathcal{A})}{P(\mathcal{A})}))^{-1}$ , it holds that

$$K_{\mathbf{z}}^{(2)}\left(-\frac{1}{\lambda_1}\right) < K_{\mathbf{z}}^{(2)}\left(-\frac{1}{\lambda_2}\right) < \dots < \frac{1}{4}. \quad (75)$$

The upperbound by  $\frac{1}{4}$  in (72), (74) and (75) follows by noticing that the value  $K_{\mathbf{z}}^{(2)}(-\frac{1}{\lambda})$  is maximized when  $\lambda = (\log(\frac{1-P(\mathcal{A})}{P(\mathcal{A})}))^{-1}$  and  $K_{\mathbf{z}}^{(2)}(-\frac{1}{\lambda}) = \frac{1}{4}$ .

Example 2.1 provides important insights on the choice of the reference measure  $P$ . Note for instance that the optimal set of models is the set  $\mathcal{A}$ . That is,  $\mathcal{T}(\mathbf{z}) = \mathcal{A}$ , with  $\mathcal{T}(\mathbf{z})$  in (5d). When the reference measure assigns a probability to the set of optimal models  $\mathcal{T}(\mathbf{z})$  that is bigger than or equal to the probability of suboptimal models  $\mathcal{M} \setminus \mathcal{T}(\mathbf{z})$ , i.e.,  $P(\mathcal{T}(\mathbf{z})) \geq \frac{1}{2}$ , the variance is strictly decreasing to zero when  $\lambda$  decreases.

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

In general, these observations suggest that reference measures  $P$  that allocate small measures to the sets containing the set  $\mathcal{T}(\mathbf{z})$  might require reducing the value  $\lambda$  beyond a small threshold in order to observe small values of  $K_{\mathbf{z}}^{(2)}(-\frac{1}{\lambda})$ , e.g., the variance of the random variable  $W$  in (60). These observations are central to understanding the concentration of probability that occurs when  $\lambda$  decreases, as discussed in the following section.

## 2.3 Concentration of Probability

This section describes two phenomena concerning the properties of the measure  $P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(\lambda)}$  in (37). First, this section shows that it is possible to determine a subset of  $\mathcal{M}$ , denoted by  $\mathcal{N}(\lambda)$  and defined later in (81), over which the measure  $P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(\lambda)}$  concentrates most of the probability as a consequence of reducing  $\lambda \in \mathcal{K}_{\mathbf{z}}$ , with  $\mathcal{K}_{\mathbf{z}}$  in (34). As  $\lambda$  tends to zero, the set  $\mathcal{N}(\lambda)$  decreases to a set that is independent of  $\lambda$ , but strongly dependent on the reference measure  $P$  in (37). Second, this section shows that at the same time that the set  $\mathcal{N}(\lambda)$  decreases as a consequence of reducing  $\lambda \in \mathcal{K}_{\mathbf{z}}$ , the probability  $P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(\lambda)}(\mathcal{N}(\lambda))$  increases to one.

### 2.3.1 Preliminaries

Given a real  $\lambda \in \mathcal{K}_z$ , with  $\mathcal{K}_z$  in (34), 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

$$\mathcal{R}_0(\lambda) \triangleq \left\{ \nu \in \mathcal{M} : L_z(\nu) = \int L_z(\alpha) dP_{\Theta|Z=z}^{(\lambda)}(\alpha) \right\}, \quad (76a)$$

$$\mathcal{R}_1(\lambda) \triangleq \left\{ \nu \in \mathcal{M} : L_z(\nu) < \int L_z(\alpha) dP_{\Theta|Z=z}^{(\lambda)}(\alpha) \right\}, \text{ and} \quad (76b)$$

$$\mathcal{R}_2(\lambda) \triangleq \left\{ \nu \in \mathcal{M} : L_z(\nu) > \int L_z(\alpha) dP_{\Theta|Z=z}^{(\lambda)}(\alpha) \right\}, \quad (76c)$$

where the function  $L_z$  is in (5b).

These sets exhibit several properties that are central for proving the main results of this section.

**Lemma 2.13.** *The measure  $P_{\Theta|Z=z}^{(\lambda)}$  in (37), satisfies*

$$P_{\Theta|Z=z}^{(\lambda)}(\mathcal{R}_1(\lambda)) > 0, \quad (77)$$

if and only if

$$P_{\Theta|Z=z}^{(\lambda)}(\mathcal{R}_2(\lambda)) > 0, \quad (78)$$

where the sets  $\mathcal{R}_1(\alpha)$  and  $\mathcal{R}_2(\alpha)$  are in (76b) and (76c), respectively.

*Proof:* The proof of Lemma 2.13 is presented in Appendix O. ■

Lemma 2.13 shows that given a real  $\gamma \in \mathcal{K}_z$ , if there exists a nonnegligible set with respect to  $P_{\Theta|Z=z}^{(\gamma)}$  whose elements induce an empirical risk that is smaller than the expected empirical risk  $K_z^{(1)}\left(-\frac{1}{\gamma}\right)$ , then there exists a nonnegligible set with respect to  $P_{\Theta|Z=z}^{(\gamma)}$  whose elements induce an empirical risk that is bigger than the expected empirical risk  $K_z^{(1)}\left(-\frac{1}{\gamma}\right)$ . Moreover, the converse also holds.

A more general result can be immediately obtained by combining Lemma 2.7 and Lemma 2.13.

**Lemma 2.14.** *For all  $\alpha \in \mathcal{K}_z$ , with  $\mathcal{K}_z$  in (34), the measure  $P_{\Theta|Z=z}^{(\lambda)}$  in (37), satisfies*

$$P_{\Theta|Z=z}^{(\lambda)}(\mathcal{R}_1(\alpha)) > 0, \quad (79)$$

if and only if

$$P_{\Theta|Z=z}^{(\lambda)}(\mathcal{R}_2(\alpha)) > 0, \quad (80)$$

where the sets  $\mathcal{R}_1(\alpha)$  and  $\mathcal{R}_2(\alpha)$  are in (76b) and (76c), respectively.

### 2.3.2 The Limiting Set

Consider the following set, with  $\lambda \in \mathcal{K}_z$ ,

$$\mathcal{N}(\lambda) \triangleq \left\{ \nu \in \mathcal{M} : L_z(\nu) \leq K_z^{(1)}\left(-\frac{1}{\lambda}\right) \right\}, \quad (81)$$

where the function  $K_{\mathbf{z}}^{(1)}$  is defined by (57). This section shows that, when  $\lambda$  decreases, the set  $\mathcal{N}(\lambda)$  forms a monotonic sequence of sets that decreases to the set

$$\mathcal{N}^* \triangleq \mathcal{L}_{\mathbf{z}}(\delta^*), \quad (82)$$

where,  $\delta^*$  is defined in (46); and the set  $\mathcal{L}_{\mathbf{z}}(\cdot)$  is defined in (44).

In the context of the ERM-RER problem in (29), when the reference measure  $P$  is a  $\sigma$ -finite measure on  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$  and the data set  $\mathbf{z}$  is the one in (5a), the set  $\mathcal{N}(\lambda)$ , with  $\lambda \in \mathcal{K}_{\mathbf{z}}$ , contains all the models that induce an empirical risk that is smaller than the expected empirical risk (with respect to the dataset  $\mathbf{z}$ ) induced by the measure  $P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(\lambda)}$  in (37). This observation unveils the existence of a relation between the set  $\mathcal{N}^*$  in (82) and the set  $\mathcal{T}(\mathbf{z})$  in (5d), as shown by the following corollary.

**Corollary 2.2.** *The set  $\mathcal{N}^*$  in (82) satisfies*

$$\mathcal{T}(\mathbf{z}) \subseteq \mathcal{N}^*, \quad (83)$$

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

$$\mathcal{T}(\mathbf{z}) = \mathcal{N}^*, \quad (84)$$

if and only if the reference measure  $P$  in (37) is coherent with the ERM-RER problem with respect to the data set  $\mathbf{z}$  in (29).

Corollary 2.2 shows that  $\mathcal{N}^*$  is not empty. This follows from the fact that the set  $\mathcal{T}(\mathbf{z})$ , which is not empty, is a subset of  $\mathcal{N}^*$ . This observation turns out to be particularly important at the light of the following theorem.

**Theorem 2.4.** *For all  $(\lambda_1, \lambda_2) \in \mathcal{K}_{\mathbf{z}} \times \mathcal{K}_{\mathbf{z}}$ , with  $\mathcal{K}_{\mathbf{z}}$  in (34) and  $\lambda_1 > \lambda_2$ , the sets  $\mathcal{N}(\lambda_1)$  and  $\mathcal{N}(\lambda_2)$  in (81) satisfy*

$$\mathcal{M} \supseteq \mathcal{N}(\lambda_1) \supseteq \mathcal{N}(\lambda_2) \supseteq \mathcal{N}^*, \quad (85)$$

with  $\mathcal{M}$  and  $\mathcal{N}^*$  the sets defined in (5d) and (82). Moreover, if the empirical risk function  $\mathbf{L}_{\mathbf{z}}$  in (5b) is continuous on  $\mathcal{M}$  and separable (Definition 2.2), then,

$$\mathcal{M} \supset \mathcal{N}(\lambda_1) \supset \mathcal{N}(\lambda_2) \supset \mathcal{N}^*. \quad (86)$$

*Proof:* The proof of Theorem 2.4 is presented in Appendix P. ■

Theorem 2.4 shows that, when  $\lambda$  decreases, the set  $\mathcal{N}(\lambda)$  monotonically decreases to the set  $\mathcal{N}^*$ , which always contains the set of solutions  $\mathcal{T}(\mathbf{z})$  to the ERM problem in (5). Nonetheless, for all  $\lambda \in \mathcal{K}_{\mathbf{z}}$ , with  $\mathcal{K}_{\mathbf{z}}$  in (34), depending on the choice of the  $\sigma$ -finite measure  $P$  in (37), only a subset of  $\mathcal{N}(\lambda)$  might exhibit nonzero probability with respect to the measure  $P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(\lambda)}$  in (37). Consider for instance a noncoherent measure  $P$ . That is, there exists a  $\delta > 0$ , such that for all  $\gamma < \delta$ , it holds that  $P(\mathcal{L}_{\mathbf{z}}(\gamma)) = 0$ . Therefore, for all  $\lambda \in \left\{ \alpha \in \mathcal{K}_{\mathbf{z}} : K_{\mathbf{z}}^{(1)}\left(-\frac{1}{\alpha}\right) > \delta \right\}$ , it holds that  $\mathcal{L}_{\mathbf{z}}(\gamma) \subseteq \mathcal{N}(\lambda)$ . In this case, note also that the set  $\mathcal{T}(\mathbf{z})$  in (5d) is a subset of a zero-measure set, i.e.,  $\mathcal{T}(\mathbf{z}) \subseteq \mathcal{L}_{\mathbf{z}}(\gamma)$ .

### 2.3.3 Probability of the Limiting Set

The following theorem shows that the probability measure  $P_{\Theta|Z=z}^{(\lambda)}$  over the set  $\mathcal{N}(\lambda)$  does not increase when  $\lambda$  increases.

**Theorem 2.5.** *For all  $(\lambda_1, \lambda_2) \in \mathcal{K}_z \times \mathcal{K}_z$ , with  $\mathcal{K}_z$  in (34) and  $\lambda_1 > \lambda_2$ , assume that the measures  $P_{\Theta|Z=z}^{(\lambda_1)}$  and  $P_{\Theta|Z=z}^{(\lambda_2)}$  satisfy (37) with  $\lambda = \lambda_1$  and  $\lambda = \lambda_2$ , respectively. Then, the set  $\mathcal{N}(\lambda_2)$  in (81) satisfies*

$$P_{\Theta|Z=z}^{(\lambda_1)}(\mathcal{N}(\lambda_2)) \leq P_{\Theta|Z=z}^{(\lambda_2)}(\mathcal{N}(\lambda_2)). \quad (87)$$

Moreover, the function  $\mathbb{L}_z$  is separable with respect to  $P$ , with  $P$  in (37), (Definition 2.2) if and only if for all pairs  $(\lambda_1, \lambda_2) \in \mathcal{K}_z \times \mathcal{K}_z$ , with  $\lambda_1 > \lambda_2$ , it holds that

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

*Proof:* The proof of Theorem 2.5 is presented in Appendix Q. ■

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

**Lemma 2.15.** *Let the function  $\mathbb{L}_z$  in (5b) be separable and consider two positive reals  $(\lambda_1, \lambda_2) \in \mathcal{K}_z \times \mathcal{K}_z$ , with  $\mathcal{K}_z$  in (34) and  $\lambda_1 > \lambda_2$ . Assume that*

$$P\left(\mathcal{N}(\lambda_1) \cap \mathcal{R}_2(\lambda_2)\right) = 0. \quad (89)$$

Then, the measures  $P_{\Theta|Z=z}^{(\lambda_1)}$  and  $P_{\Theta|Z=z}^{(\lambda_2)}$ , which satisfy (37) with  $\lambda = \lambda_1$  and  $\lambda = \lambda_2$  respectively, verify

$$P_{\Theta|Z=z}^{(\lambda_1)}(\mathcal{N}(\lambda_1)) < P_{\Theta|Z=z}^{(\lambda_2)}(\mathcal{N}(\lambda_2)), \quad (90)$$

where for all  $t \in \{1, 2\}$ , the sets  $\mathcal{R}_2(\lambda_t)$  and  $\mathcal{N}(\lambda_t)$  are in (76c) and (81), respectively.

*Proof:* The proof of Lemma 2.15 is presented in Appendix R. ■

The following example shows the relevance of Lemma 2.15 in the case in which the empirical risk function  $\mathbb{L}_z$  in (5b) is a simple function.

**Example 2.2.** *Consider Example 2.1. Note that, for all  $\lambda > 0$ ,*

$$0 < K_z^{(1)}\left(-\frac{1}{\lambda}\right) < 1, \quad (91)$$

which implies that given two reals  $\lambda_1$  and  $\lambda_2$  such that  $\lambda_1 > \lambda_2 > 0$ , it holds that,

$$\mathcal{N}(\lambda_1) \cap \mathcal{R}_2(\lambda_2) = \left\{ \nu \in \mathcal{M} : K_z^{(1)}\left(-\frac{1}{\lambda_2}\right) < \mathbb{L}_z(\nu) \leq K_z^{(1)}\left(-\frac{1}{\lambda_1}\right) \right\} \quad (92)$$

$$= \emptyset, \quad (93)$$

and moreover,  $\mathcal{N}(\lambda_1) = \mathcal{N}(\lambda_2)$ . Finally, from Lemma 2.15,

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

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

**Theorem 2.6.** *The probability measure  $P_{\Theta|Z=z}^{(\lambda)}$  in (37), with the measure  $P$  being a  $\sigma$ -finite measure  $P$  on  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ , satisfies*

$$\lim_{\lambda \rightarrow 0^+} P_{\Theta|Z=z}^{(\lambda)}(\mathcal{N}(\lambda)) = \begin{cases} 0 & \text{if } \mathcal{L}_z^* = \emptyset \\ 1 & \text{if } \mathcal{L}_z^* \neq \emptyset, \end{cases} \quad (95)$$

where, the sets  $\mathcal{N}(\cdot)$  and  $\mathcal{L}_z^*$  are respectively defined in (81) and (47).

*Proof:* The proof of Theorem 2.6 is presented in Appendix S. ■

## 2.4 Cumulant Generating Function of the empirical risk

Let  $\lambda$  be a real in  $\mathcal{K}_z$ , with  $\mathcal{K}_z$  in (34), and consider the transport of the measure  $P_{\Theta|Z=z}^{(\lambda)}$  in (37) from  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$  to  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$  through the function  $L_z$  in (5b). Denote the resulting probability measure in  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$  by  $P_{W|Z=z}^{(\lambda)}$  in  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ . That is, for all  $\mathcal{A} \in \mathcal{B}(\mathbb{R})$ ,

$$P_{W|Z=z}^{(\lambda)}(\mathcal{A}) = P_{\Theta|Z=z}^{(\lambda)}(L_z^{-1}(\mathcal{A})), \quad (96)$$

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

$$L_z^{-1}(\mathcal{A}) \triangleq \{\nu \in \mathcal{M} : L_z(\nu) \in \mathcal{A}\}. \quad (97)$$

Note that the random variable  $W$  in (60) induces the probability measure  $P_{W|Z=z}^{(\lambda)}$ . 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  $P_{W|Z=z}^{(\lambda)}$ , denoted by  $J_{z,\lambda} : \mathbb{R} \rightarrow \mathbb{R} \cup \{+\infty\}$ , satisfies for all  $t \in \mathbb{R}$ ,

$$J_{z,\lambda}(t) = \log \left( \int \exp(tz) dP_{W|Z=z}^{(\lambda)}(z) \right) = \log \left( \int \exp(tL_z(\mathbf{u})) dP_{\Theta|Z=z}^{(\lambda)}(\mathbf{u}) \right). \quad (98)$$

For all  $\lambda \in \mathcal{K}_z$ , with  $\mathcal{K}_z$  in (34), the following lemma provides an expression for  $J_{z,\lambda}(t)$  in terms of the cumulant generating function  $K_z$  in (33), for all  $t \in (-\infty, \frac{1}{\lambda})$ .

**Lemma 2.16.** *Given a real  $\lambda \in \mathcal{K}_z$ , with  $\mathcal{K}_z$  in (34), the cumulant generating function  $J_{z,\lambda}$  in (98), verifies the following equality for all  $t \in (-\infty, \frac{1}{\lambda})$ ,*

$$J_{z,\lambda}(t) = K_z \left( t - \frac{1}{\lambda} \right) - K_z \left( -\frac{1}{\lambda} \right) < +\infty. \quad (99)$$

*Proof:* The proof of Lemma 2.16 is presented in Appendix T. ■

Denote by  $J_{z,\lambda}^{(m)} : \mathbb{R} \rightarrow \mathbb{R} \cup \{+\infty\}$ , with  $m \in \mathbb{N}$ , the  $m$ -th derivative of the function  $J_{z,\lambda}$  in (98). That is, for all  $s \in (-\infty, \frac{1}{\lambda})$ ,

$$J_{z,\lambda}^{(m)}(s) = \frac{d^m}{dt^m} J_{z,\lambda}(t) \Big|_{t=s}. \quad (100)$$

From Lemma 2.16, it follows that for all  $m \in \mathbb{N}$ , and for all  $\alpha \in (-\infty, \frac{1}{\lambda})$ , the following holds,

$$J_{z,\lambda}^{(m)}(\alpha) = K_z^{(m)} \left( \alpha - \frac{1}{\lambda} \right). \quad (101)$$

The equality in (101) establishes a relation between the cumulant generating functions  $J_{\mathbf{z},\lambda}$  and  $K_{\mathbf{z}}$  respectively induced by the measures  $P_{W|Z=\mathbf{z}}^{(\lambda)}$  in (96) and  $P_{\Theta|Z=\mathbf{z}}^{(\lambda)}$  in (37), with  $\lambda \in \mathcal{K}_{\mathbf{z}}$ . The following lemma leverages these observations and presents the main result of this section.

**Theorem 2.7.** *The cumulant generating function  $J_{\mathbf{z},\lambda}$  in (98) verifies the following inequality for all  $\alpha \in (-\infty, \frac{1}{\lambda})$ ,*

$$J_{\mathbf{z},\lambda}(\alpha) \leq \alpha K_{\mathbf{z}}^{(1)}\left(-\frac{1}{\lambda}\right) + \frac{1}{2}\alpha^2 B_{\mathbf{z}}^2 \quad (102)$$

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

$$B_{\mathbf{z}}^2 = \sup_{\gamma \in \mathcal{K}_{\mathbf{z}}} K_{\mathbf{z}}^{(2)}\left(-\frac{1}{\gamma}\right) < +\infty, \quad (103)$$

with  $\mathcal{K}_{\mathbf{z}}$  in (34); and the functions  $K_{\mathbf{z}}^{(1)}$  and  $K_{\mathbf{z}}^{(2)}$  are respectively defined in (57), and (58).

*Proof:* The function  $K_{\mathbf{z}}$  in (33) is differentiable infinitely many times over the interior of the set  $\mathcal{K}_{\mathbf{z}}$  (Lemma 2.8). Thus, from the Taylor-Lagrange theorem, c.f., [43, Theorem 2.5.4], it follows that for all  $\lambda \in \mathcal{K}_{\mathbf{z}}$  and for all  $\alpha \in (-\infty, \frac{1}{\lambda})$ , there exists a real  $\xi \in (-\infty, \frac{1}{\lambda})$  such that

$$K_{\mathbf{z}}\left(\alpha - \frac{1}{\lambda}\right) = K_{\mathbf{z}}\left(-\frac{1}{\lambda}\right) + \alpha K_{\mathbf{z}}^{(1)}\left(-\frac{1}{\lambda}\right) + \frac{\alpha^2}{2} K_{\mathbf{z}}^{(2)}(\xi). \quad (104)$$

From (104) and Lemma 2.16, it holds that

$$J_{\mathbf{z},\lambda}(\alpha) = \alpha K_{\mathbf{z}}^{(1)}\left(-\frac{1}{\lambda}\right) + \frac{\alpha^2}{2} K_{\mathbf{z}}^{(2)}(\xi). \quad (105)$$

Finally, the inequality in (102) follows from the maximization of the function  $K_{\mathbf{z}}^{(2)}$  on the set  $\mathcal{K}_{\mathbf{z}}$ , which completes the proof.  $\blacksquare$

The main implication of Theorem 2.7 is that the random variable  $W$  in (60) is a sub-Gaussian random variable with parameter  $B$ . This result is central for studying of the generalization error of Gibbs stochastic solvers. Based on Theorem 2.7 and previous results in [50], generalization guarantees can be immediately obtained. Nonetheless, such analysis is left out of the scope of this work.

### 3 $(\delta, \epsilon)$ -Optimality and Sensitivity

This section introduces a notion of optimality in probability, which is reminiscent to the existing PAC-generalization guarantees. In this case, the optimality is in the sense of low expected empirical risks instead of small generalization gaps.

**Definition 3.1** ( $(\delta, \epsilon)$ -Optimality). *Given a pair of positive reals  $(\delta, \epsilon)$ , with  $\epsilon < 1$ , a stochastic solver  $P_{\Theta|Z}$  to the ERM problem with respect to the data set  $\mathbf{z}$  in (5) is said to be  $(\delta, \epsilon)$ -optimal, if the set  $\mathcal{L}_{\mathbf{z}}(\delta)$  in (44) satisfies*

$$P_{\Theta|Z=\mathbf{z}}(\mathcal{L}_{\mathbf{z}}(\delta)) > 1 - \epsilon. \quad (106)$$

For all  $\delta > 0$ , it holds that  $\mathcal{T}(\mathbf{z}) \subset \mathcal{L}_{\mathbf{z}}(\delta)$ , with the sets  $\mathcal{T}(\mathbf{z})$  and  $\mathcal{L}_{\mathbf{z}}$  in (5d) and (44), respectively. Hence, from Definition 3.1, it follows that if the stochastic solver  $P_{\Theta|Z}$  is  $(\delta, \epsilon)$ -optimal, then it assigns a probability that is always bigger than  $1 - \epsilon$  to the set  $\mathcal{L}_{\mathbf{z}}(\delta)$ . That is,

the measure  $P_{\Theta|Z=z}$  concentrates on a set that contains models that induce an empirical risk that is smaller than  $\delta$ . From this perspective, particular interest is given to the smallest  $\delta$  and the smallest  $\epsilon$  when choosing a  $(\delta, \epsilon)$ -optimal stochastic solver.

### 3.1 $(\delta, \epsilon)$ -Optimality of Gibbs Stochastic Solvers

This section shows that Gibbs stochastic solvers to the ERM-RER in (29) might be  $(\delta, \epsilon)$ -optimal for some values of  $\delta$  and  $\epsilon$ , when the reference measure is *consistent*.

**Definition 3.2** (Consistent Measure). *A measure  $P$  is said to be consistent with the ERM-RER problem with respect to the data set  $\mathbf{z}$  in (29), if the set  $\mathcal{L}_{\mathbf{z}}^*$  in (47) is not empty.*

The relation between coherent and consistent measures is described by the following corollary.

**Corollary 3.1.** *Every coherent measure is consistent.*

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

**Theorem 3.1.** *Consider the ERM-RER problem with respect to the data set  $\mathbf{z}$  in (29), when  $P$  is a  $\sigma$ -finite measure on  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ . If the measure  $P$  is consistent, for all  $(\delta, \epsilon) \in (\delta^*, +\infty) \times (0, 1)$ , with  $\delta^*$  in (46), then, there exists a  $\lambda \in \mathcal{K}_{\mathbf{z}}$ , with  $\mathcal{K}_{\mathbf{z}}$  in (34), such that the  $P$ -Gibbs stochastic solver  $P_{\Theta|Z}^{(\lambda)}$  is  $(\delta, \epsilon)$ -optimal.*

*Proof of Theorem 3.1:* Let  $\delta$  be a real in  $(\delta^*, +\infty)$ , with  $\delta^*$  in (46). Let also  $\lambda \in \mathcal{K}_{\mathbf{z}}$  satisfy the following equality:

$$K_{\mathbf{z}}^{(1)}\left(-\frac{1}{\lambda}\right) \leq \delta. \quad (107)$$

Note that under the assumption that  $P$  is consistent, such a  $\lambda$  in (107) always exists. This follows from the fact that the value  $K_{\mathbf{z}}^{(1)}\left(-\frac{1}{\lambda}\right)$  is nondecreasing with  $\lambda$  (Theorem 2.2) and in the limit as  $\lambda$  tends to zero (Theorem 2.3), the value  $K_{\mathbf{z}}^{(1)}\left(-\frac{1}{\lambda}\right)$  tends to  $\delta^*$ . Moreover, from (44) and (81), it holds that

$$\mathcal{N}(\lambda) \subseteq \mathcal{L}_{\mathbf{z}}(\delta), \quad (108)$$

and thus,

$$P_{\Theta|Z=z}^{(\lambda)}(\mathcal{L}_{\mathbf{z}}(\delta)) \geq P_{\Theta|Z=z}^{(\lambda)}(\mathcal{N}(\lambda)). \quad (109)$$

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

$$P_{\Theta|Z=z}^{(\gamma)}(\mathcal{N}(\gamma)) > 1 - \epsilon. \quad (110)$$

The existence of such a positive real  $\gamma$  follows from Theorem 2.6. Hence, from (110), it holds that,

$$1 - \epsilon < P_{\Theta|Z=z}^{(\gamma)}(\mathcal{N}(\gamma)) \quad (111)$$

$$\leq P_{\Theta|Z=z}^{(\gamma)}(\mathcal{L}_{\mathbf{z}}(\delta)), \quad (112)$$

where the inequality in (112) follows from the fact that  $\mathcal{N}(\gamma) \subseteq \mathcal{N}(\lambda) \subseteq \mathcal{L}_{\mathbf{z}}(\delta)$ . Finally, the inequality in (112) implies that the stochastic solver  $P_{\Theta|Z=z}^{(\gamma)}$  is a  $(\delta, \epsilon)$ -optimal probability measure (Definition 3.1). This completes the proof.  $\blacksquare$

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



**Theorem 3.2.** *Consider the ERM-RER problem with respect to the data set  $\mathbf{z}$  in (29), when  $P$  is a  $\sigma$ -finite measure on  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ . Then, for all  $(\delta, \epsilon) \in (0, +\infty) \times (0, 1)$ , there always exists a  $\lambda \in \mathcal{K}_{\mathbf{z}}$ , with  $\mathcal{K}_{\mathbf{z}}$  in (34), such that the  $P$ -Gibbs stochastic solver  $P_{\Theta|\mathbf{Z}}^{(\lambda)}$  is  $(\delta, \epsilon)$ -optimal if and only if the reference measure  $P$  is coherent.*

*Proof of Theorem 3.2:* The proof is divided into two parts. The first part deals with the case in which the measure  $P$  is coherent. The second part deals with the converse case.

The first part is as follows. Under the assumption that the measure  $P$  is coherent it follows that  $\delta^* = 0$ . Moreover, it follows from Corollary 3.1 that the measure  $P$  is also consistent. Then, from Theorem 3.1, it follows that for all  $(\delta, \epsilon) \in (0, +\infty) \times (0, 1)$ , there always exists a  $\lambda \in \mathcal{K}_{\mathbf{z}}$ , with  $\mathcal{K}_{\mathbf{z}}$  in (34), such that the  $P$ -Gibbs stochastic solver  $P_{\Theta|\mathbf{Z}}^{(\lambda)}$  is  $(\delta, \epsilon)$ -optimal.

The second part is as follows. Under the assumption that the measure  $P$  is not coherent, there always exists a  $\delta > 0$ , such that for all  $\gamma < \delta$ , it holds that  $P(\mathcal{L}_{\mathbf{z}}(\gamma)) = 0$ . Hence, the  $P$ -Gibbs stochastic solver  $P_{\Theta|\mathbf{Z}}^{(\lambda)}$  is not  $(\gamma, \epsilon)$ -optimal, for all  $\gamma < \delta$ , with  $\epsilon \in (0, 1)$ . This completes the proof.  $\blacksquare$

### 3.2 Sensitivity of the Expected Empirical Risk

The sensitivity due to a change of measure is a performance metric defined as follows.

**Definition 3.3.** *Given a  $\sigma$ -finite measure  $P$  and a real  $\lambda > 0$ , let*

$$D_{\lambda} : (\mathcal{X} \times \mathcal{Y})^n \times \Delta_P(\mathcal{M}, \mathcal{B}(\mathcal{M})) \rightarrow [0, +\infty]$$

be a function such that

$$D_{\lambda}(\mathbf{u}, Q) = \begin{cases} \left| R_{\mathbf{u}}(Q) - R_{\mathbf{u}}\left(P_{\Theta|\mathbf{Z}=\mathbf{u}}^{(\lambda)}\right) \right| & \text{if } \lambda \in \mathcal{K}_{\mathbf{u}} \\ +\infty & \text{otherwise,} \end{cases} \quad (113)$$

where the function  $R_{\mathbf{u}}$  is defined in (12) and the measure  $P_{\Theta|\mathbf{Z}=\mathbf{u}}^{(\lambda)}$  is the solution to the ERM-RER problem in (29) with respect to the data set  $\mathbf{u}$ . The sensitivity of the expected empirical risk at dataset  $\mathbf{u}$  when the measure changes from  $P_{\Theta|\mathbf{Z}=\mathbf{u}}^{(\lambda)}$  to  $Q$  is  $D_{\lambda}(\mathbf{u}, Q)$ .

The sensitivity  $D_{\lambda}(\mathbf{z}, Q)$ , with  $\mathbf{z}$  the data set in (5a), is a positive real that indicates the variation of the expected empirical risk with respect to  $\mathbf{z}$ , i.e., the function  $R_{\mathbf{z}}$  in (61), when the measure is changed from the optimal solution to the ERM-RER problem in (29), i.e., the measure  $P_{\Theta|\mathbf{Z}=\mathbf{z}}^{(\lambda)}$ , to an alternative measure  $Q$ . That is, the sensitivity  $D_{\lambda}(\mathbf{z}, Q)$  is a means to quantify the change of the expected empirical risk with respect to deviations from the optimal solution of the ERM-RER problem. In the aim of characterizing the sensitivity  $D_{\lambda}(\mathbf{z}, Q)$ , consider the following lemma.

**Lemma 3.1.** *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  $\mathbf{u} \in (\mathcal{X} \times \mathcal{Y})^n$ ,*

$$\int \mathbf{L}_{\mathbf{u}}(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) - \int \mathbf{L}_{\mathbf{u}}(\boldsymbol{\theta}) dP(\boldsymbol{\theta}) \leq \inf_{t \in (-\infty, 0)} \left( \frac{D(Q\|P) + \log \left( \int \exp(t(\mathbf{L}_{\mathbf{u}}(\boldsymbol{\theta}) - \mu)) dP(\boldsymbol{\theta}) \right)}{t} \right) \quad (114)$$

where  $\mu = \int \mathbf{L}_{\mathbf{u}}(\boldsymbol{\theta}) dP(\boldsymbol{\theta})$ , and the function  $\mathbf{L}_{\mathbf{u}}$  is defined in (5b).

*Proof:* The proof of Lemma 3.1 is presented in Appendix U. ■

Lemma 3.1 together with Theorem 2.7 lead to an upper bound on the sensitivity to a change of measure.

**Theorem 3.3.** *Consider the solution  $P_{\Theta|Z=z}^{(\lambda)}$  in (37) to ERM-RER problem in (29) with  $\lambda \in \mathcal{K}_z$  and  $\mathcal{K}_z$  in (34). Let  $Q$  be a probability measure over  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$  absolutely continuous with  $P$ . Then,*

$$D_\lambda(z, Q) \leq \sqrt{2B_z^2 D(Q \| P_{\Theta|Z=z}^{(\lambda)})} \quad (115)$$

where the function  $D_\lambda$  is defined in (113); and the constant  $B_z$  is defined in (103).

*Proof:* The proof of Theorem 3.3 is presented in Appendix V. ■

In the context of the ERM problem in (5), Theorem 3.3 establishes an upper and a lower bound on the increase or decrease on the expected empirical risk minimization that can be obtained by deviating from the optimal solution of the ERM-RER in (29). More specifically, note that for all probability measures  $Q \in \Delta(\mathcal{M}, \mathcal{B}(\mathcal{M}))$  that are absolutely continuous with  $P$ , it holds that,

$$R_z(P_{\Theta|Z=z}^{(\lambda)}) - \sqrt{2B_z^2 D(Q \| P_{\Theta|Z=z}^{(\lambda)})} \leq R_z(Q) \leq R_z(P_{\Theta|Z=z}^{(\lambda)}) + \sqrt{2B_z^2 D(Q \| P_{\Theta|Z=z}^{(\lambda)})}. \quad (116)$$

The probability measure  $Q$  over  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$  that achieves the lower bound (or the upper bound) in (116) can be explicitly calculated, as shown by the following theorem.

**Theorem 3.4.** *Consider the solution  $P_{\Theta|Z=z}^{(\lambda)}$  in (37) to ERM-RER problem in (29) with  $\lambda \in \mathcal{K}_z$  and  $\mathcal{K}_z$  in (34). Consider also the following optimization problem,*

$$\min_{Q \in \Delta_P(\mathcal{M}, \mathcal{B}(\mathcal{M}))} \int \mathbb{L}_z(\theta) dQ(\theta), \quad (117a)$$

$$\text{subject to: } D(Q \| P_{\Theta|Z=z}^{(\lambda)}) \leq c, \quad (117b)$$

where the function  $\mathbb{L}_z$  is in (5b); and the constant  $c$  is nonnegative. Then, the solution to optimization problem in (117) is a Gibbs probability measure  $P_{\Theta|Z=z}^{(\omega)}$  satisfying (37) with  $\omega \leq \lambda$  such that

$$D(P_{\Theta|Z=z}^{(\omega)} \| P_{\Theta|Z=z}^{(\lambda)}) = c. \quad (118)$$

*Proof:* The proof of Theorem 3.4 is presented in Appendix W. ■

The relevance of Theorem 3.4 is that it shows that the measure  $Q$  that reduces the expected empirical risk beyond the expected empirical risk induced by the Gibbs measure  $P_{\Theta|Z=z}^{(\lambda)}$  in (37), i.e., the value  $R_z(P_{\Theta|Z=z}^{(\lambda)})$  in (61), subject to the constraint  $D(Q \| P_{\Theta|Z=z}^{(\lambda)}) \leq c$ , with  $c > 0$ , is also a Gibbs measure  $P_{\Theta|Z=z}^{(\omega)}$ , with  $\omega < \lambda$  chosen to satisfy (118). More specifically, the inequality in the left hand side in (116) is observed with equality when  $Q = P_{\Theta|Z=z}^{(\omega)}$ . Similarly, it can be shown that the measure  $Q$  that increases the expected empirical risk beyond the expected empirical risk induced by the Gibbs measure  $P_{\Theta|Z=z}^{(\lambda)}$  in (37), i.e., the value  $R_z(P_{\Theta|Z=z}^{(\lambda)})$  in (61), subject to the constraint  $D(Q \| P_{\Theta|Z=z}^{(\lambda)}) \leq c$ , with  $c > 0$ , is also a Gibbs measure  $P_{\Theta|Z=z}^{(\omega)}$ , with  $\omega > \lambda$  chosen to satisfy (118), provided that such  $\omega$  exists in  $\mathcal{K}_z$ , with  $\mathcal{K}_z$  in (34), c.f., Lemma 2.1.

The following corollary of Theorem 3.3 studies the expectation of the sensitivity with respect to the probability measure  $P_{\mathbf{Z}}$  in (19).

**Corollary 3.2.** *For all*

$$\lambda \in \bigcap_{\mathbf{u} \in (\mathcal{X} \times \mathcal{Y})^n} \mathcal{K}_{\mathbf{u}}, \quad (119)$$

with  $\mathcal{K}_{\mathbf{u}}$  in (34), given a probability measure  $Q$  over  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$  absolutely continuous with  $P$ , it holds that

$$\int \mathsf{D}_{\lambda}(\mathbf{u}, Q) \, dP_{\mathbf{Z}}(\mathbf{u}) \leq \int \sqrt{2B_{\mathbf{u}}^2 D(Q \| P_{\Theta|\mathbf{Z}=\mathbf{u}}^{(\lambda)})} \, dP_{\mathbf{Z}}(\mathbf{u}), \quad (120)$$

where  $B_{\mathbf{u}}$  is defined in (103); the probability measure  $P_{\Theta|\mathbf{Z}=\mathbf{u}}^{(\lambda)}$  is the solution to the ERM-RER problem in (29) with respect to the data set  $\mathbf{u}$ ; and the probability measure  $P_{\mathbf{Z}}$  is defined in (19).

The set  $\bigcap_{\mathbf{u} \in (\mathcal{X} \times \mathcal{Y})^n} \mathcal{K}_{\mathbf{u}}$  in (122) can be empty for some choices of the  $\sigma$ -finite measure  $P$  and empirical loss function  $\mathsf{L}_{\mathbf{z}}$  in (5b). Nonetheless, from Lemma 2.2, it follows that when  $P$  is a probability measure,

$$\bigcap_{\mathbf{u} \in (\mathcal{X} \times \mathcal{Y})^n} \mathcal{K}_{\mathbf{u}} = (0, +\infty). \quad (121)$$

In the following, the expectation of the sensitivity with respect to the measure  $P_{\mathbf{Z}}$  in (19) is to shown to have an upper-bound that can be expressed in terms of the lautum information between the models and the data sets.

**Theorem 3.5.** *For all*

$$\lambda \in \bigcap_{\mathbf{u} \in (\mathcal{X} \times \mathcal{Y})^n} \mathcal{K}_{\mathbf{u}}, \quad (122)$$

with  $\mathcal{K}_{\mathbf{u}}$  in (34), it follows that

$$\int \mathsf{D}_{\lambda}(\mathbf{u}, P_{\Theta}^{(\lambda)}) \, dP_{\mathbf{Z}}(\mathbf{u}) \leq \sqrt{2B^2 \int D(P_{\Theta}^{(\lambda)} \| P_{\Theta|\mathbf{Z}=\mathbf{u}}^{(\lambda)}) \, dP_{\mathbf{Z}}(\mathbf{u})}, \quad (123)$$

where the probability measure  $P_{\Theta|\mathbf{Z}=\mathbf{u}}^{(\lambda)}$  is the solution to the ERM-RER problem in (29) with respect to the data set  $\mathbf{u}$ ; the probability measure  $P_{\mathbf{Z}}$  is defined in (19); the probability measure  $P_{\Theta}^{(\lambda)}$  is such that for all  $\mathcal{A} \in \mathcal{B}(\mathcal{M})$ ,

$$P_{\Theta}^{(\lambda)}(\mathcal{A}) = \int P_{\Theta|\mathbf{Z}=\mathbf{u}}^{(\lambda)}(\mathcal{A}) \, dP_{\mathbf{Z}}(\mathbf{u}); \quad (124)$$

and the constant  $B$  satisfies  $B^2 = \sup_{\mathbf{u} \in (\mathcal{X} \times \mathcal{Y})^n} B_{\mathbf{u}}^2$ , with  $B_{\mathbf{u}}$  defined in (103).

*Proof:* The proof follows from Corollary 3.5. In particular, from (120), for all probability measures  $Q$  over  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$  absolutely continuous with  $P$ , it holds that

$$\int \mathsf{D}_{\lambda}(\mathbf{u}, Q) \, dP_{\mathbf{Z}}(\mathbf{u}) \leq \int \sqrt{2B_{\mathbf{u}}^2 D(Q \| P_{\Theta|\mathbf{Z}=\mathbf{u}}^{(\lambda)})} \, dP_{\mathbf{Z}}(\mathbf{u}) \quad (125)$$

$$\leq \int \sqrt{2B^2 D(Q \| P_{\Theta|\mathbf{Z}=\mathbf{u}}^{(\lambda)})} \, dP_{\mathbf{Z}}(\mathbf{u}) \quad (126)$$

$$\leq \sqrt{2B^2 \int D(Q \| P_{\Theta|\mathbf{Z}=\mathbf{u}}^{(\lambda)}) \, dP_{\mathbf{Z}}(\mathbf{u})}. \quad (127)$$

where the inequality in (126) follows from (3.5); and the inequality in (127) follows from Jensen's inequality [3, Theorem 6.3.5].

Finally, from (124) and (127), it holds that

$$\int \mathbb{D}_\lambda(\mathbf{u}, P_{\Theta}^{(\lambda)}) dP_{\mathbf{Z}}(\mathbf{u}) \leq \sqrt{2 B^2 \mathbb{L}(P_{\Theta|\mathbf{Z}}^{(\lambda)} P_{\mathbf{Z}})}, \quad (128)$$

where the function  $\mathbb{L}$  in (128) is the lautum information (Definition 1.2). This completes the proof. ■

The right-hand side in (123) can be written in terms of the lautum information [33] between the models and the data sets, denoted by  $\mathbb{L}(P_{\Theta|\mathbf{Z}}^{(\lambda)} P_{\mathbf{Z}})$ , by observing that

$$\mathbb{L}(P_{\Theta|\mathbf{Z}}^{(\lambda)} P_{\mathbf{Z}}) = \int D(P_{\Theta}^{(\lambda)} \| P_{\Theta|\mathbf{Z}=\mathbf{u}}^{(\lambda)}) dP_{\mathbf{Z}}(\mathbf{u}). \quad (129)$$

In a nutshell, it can be concluded that as the expectation of the absolute value of the generalization gap with respect to the measure  $P_{\mathbf{Z}}$  in (19) is upper bounded in terms of the mutual information between the models and the datasets, c.f., Example 1.1; the expectation of the sensitivity with respect to the measure  $P_{\mathbf{Z}}$  is upper bounded by the lautum information between the models and the datasets, c.f., Theorem 3.5.

## 4 Discussion and Final Remarks

The ERM-RER problem in (29) has been studied under the assumption that the reference measure  $P$  is a  $\sigma$ -finite measure. In this context, it has been shown that the solution exists and is the unique probability measure whose Radon-Nikodym derivative with respect to  $P$  is in Theorem 2.1. Interestingly, the optimal measure is shown to belong to an extended class of Gibbs measures for which the partition function is with respect to a  $\sigma$ -finite measure instead of a probability measure. An important remark is that the choice of the reference measure  $P$  plays a central role in the concentration of the probability of the optimal measure; and the mean and variance of the empirical risk when the models are distributed according to the optimal measure. A class of reference measures (coherent measures) for which the highest probability is always allocated to the set of models that are solutions ERM problem (lowest empirical risk) is introduced. More interestingly, it is shown that when the reference measure does not belong to this class, the models that induce the lowest empirical risk are observed with zero probability. Finally, the empirical risk when the models are distributed according to the optimal probability measure is shown to be a sub-Gaussian random variable. This observation is leveraged to study the sensitivity of the ERM-RER problems and unveil the connections between the sensitivity and the lautum information between the models and the data sets.

# Appendices

## A Proof of Lemma 2.1

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

The first part is as follows. Under the assumption that the set  $\mathcal{K}_z$  is not empty, there always exists a real  $b \in \mathcal{K}_z$ , such that  $K_z(-\frac{1}{b}) < +\infty$ . Note that for all  $\theta \in \mathcal{M}$ ,

$$\frac{d}{dt} \exp\left(-\frac{1}{t} L_z(\theta)\right) = \frac{1}{t^2} L_z(\theta) \exp\left(-\frac{1}{t} L_z(\theta)\right) \geq 0, \quad (130)$$

and thus, from (33), it follows that  $K_z(-\frac{1}{b})$  is nondecreasing with  $b$ . This implies that  $(0, b] \subset \mathcal{K}_z$ . This proves the convexity of  $\mathcal{K}_z$ .

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

$$b^* = \sup \mathcal{K}_z. \quad (131)$$

Hence, if  $b^* = +\infty$ , it follows from (34) that

$$\mathcal{K}_z = (0, +\infty). \quad (132)$$

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

$$(0, b^*) \subset \mathcal{K}_z. \quad (133)$$

This completes the first part.

The second part is trivial. Under the assumption that the set  $\mathcal{K}_z$  in (34) is empty, there is nothing to prove.

This completes the proof.

## B Proof of Lemma 2.2

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

$$\exp\left(-\frac{1}{t} L_z(\theta)\right) \leq 1, \quad (134)$$

and thus,

$$K_z\left(-\frac{1}{t}\right) = \log\left(\int \exp\left(-\frac{1}{t} L_z(\mathbf{u})\right) dP(\mathbf{u})\right) \quad (135)$$

$$\leq \log\left(\int dP(\mathbf{u})\right) \quad (136)$$

$$\leq 0, \quad (137)$$

which implies that  $(0, +\infty) \subseteq \mathcal{K}_z$ . Thus, from (34), it holds that  $\mathcal{K}_z = (0, +\infty)$ , which completes the proof.

## C Proof of Theorem 2.1

The objective function in the optimization problem in (29) can be written as follows:

$$\min_Q \left[ \int \mathbf{L}_z(\boldsymbol{\nu}) \frac{dQ}{dP}(\boldsymbol{\nu}) dP(\boldsymbol{\nu}) + \lambda \int \frac{dQ}{dP}(\boldsymbol{\nu}) \log \left( \frac{dQ}{dP}(\boldsymbol{\nu}) \right) dP(\boldsymbol{\nu}) \right], \quad (138)$$

where the optimization is over all measures  $Q$  on  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$  that are absolutely continuous with  $P$  and satisfy

$$\int \frac{dQ}{dP}(\boldsymbol{\nu}) dP(\boldsymbol{\nu}) = 1, \quad (139)$$

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

Let  $\mathcal{M}$  be the set of nonnegative measurable functions with respect to the measurable spaces  $(\text{supp } P, \mathcal{B}(\text{supp } P))$  and  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ . The Lagrangian of the optimization problem in (138) can be constructed in terms of a function in  $\mathcal{M}$ , instead of a measure in  $\Delta(\text{supp } P, \mathcal{B}(\text{supp } P))$ . Let such Lagrangian be  $L : \mathcal{M} \times [0, +\infty) \rightarrow \mathbb{R}$  of the form

$$\begin{aligned} L \left( \frac{dQ}{dP}, \beta \right) &= \int \mathbf{L}_z(\boldsymbol{\nu}) \frac{dQ}{dP}(\boldsymbol{\nu}) dP(\boldsymbol{\nu}) + \lambda \int \frac{dQ}{dP}(\boldsymbol{\nu}) \log \left( \frac{dQ}{dP}(\boldsymbol{\nu}) \right) dP(\boldsymbol{\nu}) \\ &\quad + \beta \left( \int \frac{dQ}{dP}(\boldsymbol{\nu}) dP(\boldsymbol{\nu}) - 1 \right), \end{aligned} \quad (140)$$

where  $\beta$  is a positive real that acts as a Lagrangian multiplier due to the constraint (139).

Let  $g : \mathbb{R}^k \rightarrow \mathbb{R}$  be a function in  $\mathcal{M}$ . The Gateaux differential of the functional  $L$  in (140) at  $\left( \frac{dQ}{dP}, \beta \right) \in \mathcal{M} \times [0, +\infty)$  in the direction of  $g$  is

$$\partial L \left( \frac{dQ}{dP}, \beta; g \right) \triangleq \left. \frac{d}{d\alpha} r(\alpha) \right|_{\alpha=0}, \quad (141)$$

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

$$\begin{aligned} r(\alpha) &= \int \mathbf{L}_z(\boldsymbol{\nu}) \left( \frac{dQ}{dP}(\boldsymbol{\nu}) + \alpha g(\boldsymbol{\nu}) \right) dP(\boldsymbol{\nu}) + \beta \left( \int \left( \frac{dQ}{dP}(\boldsymbol{\nu}) + \alpha g(\boldsymbol{\nu}) \right) dP(\boldsymbol{\nu}) - 1 \right) \\ &\quad + \lambda \int \left( \frac{dQ}{dP}(\boldsymbol{\nu}) + \alpha g(\boldsymbol{\nu}) \right) \log \left( \frac{dQ}{dP}(\boldsymbol{\nu}) + \alpha g(\boldsymbol{\nu}) \right) dP(\boldsymbol{\nu}). \end{aligned} \quad (142)$$

Note that the derivative of the real function  $r$  in (142) is

$$\begin{aligned} \frac{d}{d\alpha} r(\alpha) &= \int \mathbf{L}_z(\boldsymbol{\nu}) g(\boldsymbol{\nu}) dP(\boldsymbol{\nu}) + \beta \int g(\boldsymbol{\nu}) dP(\boldsymbol{\nu}) \\ &\quad + \lambda \int g(\boldsymbol{\nu}) \left( 1 + \log \left( \frac{dQ}{dP}(\boldsymbol{\nu}) + \alpha g(\boldsymbol{\nu}) \right) \right) dP(\boldsymbol{\nu}). \end{aligned} \quad (143)$$

From (141) and (143), it follows that

$$\partial L \left( \frac{dQ}{dP}, \beta; g \right) = \int g(\boldsymbol{\nu}) \left( \mathbf{L}_z(\boldsymbol{\nu}) + \lambda \left( 1 + \log \left( \frac{dQ}{dP}(\boldsymbol{\nu}) \right) \right) + \beta \right) dP(\boldsymbol{\nu}). \quad (144)$$

The relevance of the Gateaux differential in (144) stems from [27, Theorem 1, page 178], which unveils the fact that a necessary condition for the functional  $L$  in (140) to have a minimum at  $\left(\frac{dP_{\Theta|Z=z}^{(\lambda)}}{dP}, \beta\right) \in \mathcal{M} \times [0, +\infty)$  is that for all functions  $g \in \mathcal{M}$ ,

$$\partial L \left( \frac{dP_{\Theta|Z=z}^{(\lambda)}}{dP}; g \right) = 0. \quad (145)$$

From (145), it follows that  $\frac{dP_{\Theta|Z=z}^{(\lambda)}}{dP}$  must satisfy for all functions  $g$  in  $\mathcal{M}$  that

$$\int g(\nu) \left( \mathsf{L}_z(\nu) + \lambda \left( 1 + \log \left( \frac{dP_{\Theta|Z=z}^{(\lambda)}}{dP}(\nu) \right) \right) + \beta \right) dP(\nu) = 0, \quad (146)$$

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

$$\mathsf{L}_z(\nu) + \lambda \left( 1 + \log \left( \frac{dP_{\Theta|Z=z}^{(\lambda)}}{dP}(\nu) \right) \right) + \beta = 0, \quad (147)$$

and thus,

$$\frac{dP_{\Theta|Z=z}^{(\lambda)}}{dP}(\nu) = \exp\left(-\frac{\beta + \lambda}{\lambda}\right) \exp\left(-\frac{\mathsf{L}_z(\nu)}{\lambda}\right), \quad (148)$$

with  $\beta$  chosen to satisfy (139). That is,

$$\frac{dP_{\Theta|Z=z}^{(\lambda)}}{dP}(\nu) = \frac{\exp\left(-\frac{\mathsf{L}_z(\nu)}{\lambda}\right)}{\int \exp\left(-\frac{\mathsf{L}_z(\theta)}{\lambda}\right) dP(\theta)} \quad (149)$$

$$= \exp\left(-K_z \left(-\frac{1}{\lambda}\right) - \frac{1}{\lambda} \mathsf{L}_z(\nu)\right). \quad (150)$$

The proof continues by verifying that the objective function in (138) is strictly convex, and thus, the measure  $P_{\Theta|Z=z}^{(\lambda)}$  that satisfies (149) is the unique minimizer. More specifically, note that the objective function in (138) is the sum of two terms. The first one, i.e.,  $\int \mathsf{L}_z(\nu) \frac{dQ}{dP}(\nu) dP(\nu)$ , is linear in  $\frac{dQ}{dP}$ . The second, i.e.,  $\int \frac{dQ}{dP}(\nu) \log\left(\frac{dQ}{dP}(\nu)\right) dP(\nu)$ , is strictly convex with  $\frac{dQ}{dP}$ . Hence, given that  $\lambda > 0$ , the sum of both terms is strictly convex with  $\frac{dQ}{dP}$ . This implies the uniqueness of  $P_{\Theta|Z=z}^{(\lambda)}$  and completes the proof.

## D Proof of Lemma 2.3

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

$$\mathsf{L}_z(\theta) \geq \mathsf{L}_z(\nu) = \mathsf{L}_z(\mu), \quad (151)$$

and thus, for all  $\lambda \in \mathcal{K}_z$ , with  $\mathcal{K}_z$  in (34), it holds that

$$\exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\theta})}{\lambda}\right) \leq \exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\nu})}{\lambda}\right) = \exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\mu})}{\lambda}\right), \quad (152)$$

which implies

$$\frac{\exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\theta})}{\lambda}\right)}{\int \exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\alpha})}{\lambda}\right) dP(\boldsymbol{\alpha})} \leq \frac{\exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\nu})}{\lambda}\right)}{\int \exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\alpha})}{\lambda}\right) dP(\boldsymbol{\alpha})} = \frac{\exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\mu})}{\lambda}\right)}{\int \exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\alpha})}{\lambda}\right) dP(\boldsymbol{\alpha})}. \quad (153)$$

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

$$\frac{dP^{(\lambda)}_{\boldsymbol{\Theta}|Z=z}}{dP}(\boldsymbol{\theta}) \leq \frac{dP^{(\lambda)}_{\boldsymbol{\Theta}|Z=z}}{dP}(\boldsymbol{\mu}) = \frac{dP^{(\lambda)}_{\boldsymbol{\Theta}|Z=z}}{dP}(\boldsymbol{\nu}), \quad (154)$$

which completes the proof.

## E Proof of Lemma 2.5

From Theorem 2.1, it follows that for all  $\lambda \in \mathcal{K}_z$  and for all  $\boldsymbol{\theta} \in \text{supp } P$ ,

$$\frac{dP^{(\lambda)}_{\boldsymbol{\Theta}|Z=z}}{dP}(\boldsymbol{\theta}) = \frac{\exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\theta})}{\lambda}\right)}{\int \exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\nu})}{\lambda}\right) dP(\boldsymbol{\nu})} \quad (155)$$

$$= \left( \exp\left(\frac{\mathbf{L}_z(\boldsymbol{\theta})}{\lambda}\right) \int \exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\nu})}{\lambda}\right) dP(\boldsymbol{\nu}) \right)^{-1} \quad (156)$$

$$= \left( \int \exp\left(\frac{1}{\lambda}(\mathbf{L}_z(\boldsymbol{\theta}) - \mathbf{L}_z(\boldsymbol{\nu}))\right) dP(\boldsymbol{\nu}) \right)^{-1}. \quad (157)$$

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

$$\mathcal{A}_0(\boldsymbol{\theta}) \triangleq \{\boldsymbol{\nu} \in \mathcal{M} : \mathbf{L}_z(\boldsymbol{\theta}) - \mathbf{L}_z(\boldsymbol{\nu}) = 0\}, \quad (158a)$$

$$\mathcal{A}_1(\boldsymbol{\theta}) \triangleq \{\boldsymbol{\nu} \in \mathcal{M} : \mathbf{L}_z(\boldsymbol{\theta}) - \mathbf{L}_z(\boldsymbol{\nu}) < 0\}, \text{ and} \quad (158b)$$

$$\mathcal{A}_2(\boldsymbol{\theta}) \triangleq \{\boldsymbol{\nu} \in \mathcal{M} : \mathbf{L}_z(\boldsymbol{\theta}) - \mathbf{L}_z(\boldsymbol{\nu}) > 0\}. \quad (158c)$$



Using the sets  $\mathcal{A}_0(\boldsymbol{\theta})$ ,  $\mathcal{A}_1(\boldsymbol{\theta})$ , and  $\mathcal{A}_2(\boldsymbol{\theta})$  in (157), the following holds for all  $\lambda \in \mathcal{K}_z$  and for all  $\boldsymbol{\theta} \in \text{supp } P$ ,

$$\begin{aligned} \frac{dP_{\Theta|Z=z}^{(\lambda)}}{dP}(\boldsymbol{\theta}) &= \left( \int_{\mathcal{A}_0(\boldsymbol{\theta})} \exp\left(\frac{1}{\lambda}(\mathbb{L}_z(\boldsymbol{\theta}) - \mathbb{L}_z(\boldsymbol{\nu}))\right) dP(\boldsymbol{\nu}) \right. \\ &\quad + \int_{\mathcal{A}_1(\boldsymbol{\theta})} \exp\left(\frac{1}{\lambda}(\mathbb{L}_z(\boldsymbol{\theta}) - \mathbb{L}_z(\boldsymbol{\nu}))\right) dP(\boldsymbol{\nu}) \\ &\quad \left. + \int_{\mathcal{A}_2(\boldsymbol{\theta})} \exp\left(\frac{1}{\lambda}(\mathbb{L}_z(\boldsymbol{\theta}) - \mathbb{L}_z(\boldsymbol{\nu}))\right) dP(\boldsymbol{\nu}) \right)^{-1} \end{aligned} \quad (159)$$

$$\begin{aligned} &= \left( P(\mathcal{A}_0(\boldsymbol{\theta})) + \int_{\mathcal{A}_1(\boldsymbol{\theta})} \exp\left(\frac{1}{\lambda}(\mathbb{L}_z(\boldsymbol{\theta}) - \mathbb{L}_z(\boldsymbol{\nu}))\right) dP(\boldsymbol{\nu}) \right. \\ &\quad \left. + \int_{\mathcal{A}_2(\boldsymbol{\theta})} \exp\left(\frac{1}{\lambda}(\mathbb{L}_z(\boldsymbol{\theta}) - \mathbb{L}_z(\boldsymbol{\nu}))\right) dP(\boldsymbol{\nu}) \right)^{-1}. \end{aligned} \quad (160)$$

Note that the sets  $\{\boldsymbol{\nu} \in \mathcal{M} : \mathbb{L}_z(\boldsymbol{\nu}) > \delta^*\}$  and  $\{\boldsymbol{\nu} \in \mathcal{M} : \mathbb{L}_z(\boldsymbol{\nu}) = \delta^*\}$ , with  $\delta^*$  in (46), form a partition of the set  $\text{supp } P$ . Following this observation, the rest of the proof is divided into two parts. The first part evaluates  $\lim_{\lambda \rightarrow 0^+} \frac{dP_{\Theta|Z=z}^{(\lambda)}}{dP}(\boldsymbol{\theta})$ , with  $\boldsymbol{\theta} \in \{\boldsymbol{\nu} \in \mathcal{M} : \mathbb{L}_z(\boldsymbol{\nu}) > \delta^*\}$ . The second part considers the case in which  $\boldsymbol{\theta} \in \{\boldsymbol{\nu} \in \mathcal{M} : \mathbb{L}_z(\boldsymbol{\nu}) = \delta^*\}$ .

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

$$\mathcal{A}_0(\boldsymbol{\theta}) = \{\boldsymbol{\mu} \in \mathcal{M} : \mathbb{L}_z(\boldsymbol{\mu}) = \delta\}, \quad (161a)$$

$$\mathcal{A}_1(\boldsymbol{\theta}) = \{\boldsymbol{\mu} \in \mathcal{M} : \mathbb{L}_z(\boldsymbol{\mu}) > \delta\}, \text{ and} \quad (161b)$$

$$\mathcal{A}_2(\boldsymbol{\theta}) = \{\boldsymbol{\mu} \in \mathcal{M} : \mathbb{L}_z(\boldsymbol{\mu}) < \delta\}. \quad (161c)$$

Consider the sets  $\mathcal{A}_{2,1}(\boldsymbol{\theta}) \triangleq \{\boldsymbol{\mu} \in \mathcal{M} : \mathbb{L}_z(\boldsymbol{\mu}) < \delta^*\}$  and  $\mathcal{A}_{2,2}(\boldsymbol{\theta}) = \{\boldsymbol{\mu} \in \mathcal{M} : \delta^* \leq \mathbb{L}_z(\boldsymbol{\mu}) < \delta\}$ , and note that  $\mathcal{A}_{2,1}(\boldsymbol{\theta})$  and  $\mathcal{A}_{2,2}(\boldsymbol{\theta})$  form a partition of  $\mathcal{A}_2(\boldsymbol{\theta})$ , and

$$P(\mathcal{A}_{2,1}(\boldsymbol{\theta})) = 0. \quad (162)$$

Hence, plugging the equalities in (161) and (162) in (160) yields,

$$\begin{aligned} \frac{dP_{\Theta|Z=z}^{(\lambda)}}{dP}(\boldsymbol{\theta}) &= \left( P(\mathcal{A}_0(\boldsymbol{\theta})) + \int_{\mathcal{A}_1(\boldsymbol{\theta})} \exp\left(\frac{1}{\lambda}(\mathbb{L}_z(\boldsymbol{\theta}) - \mathbb{L}_z(\boldsymbol{\nu}))\right) dP(\boldsymbol{\nu}) \right. \\ &\quad \left. + \int_{\mathcal{A}_{2,2}(\boldsymbol{\theta})} \exp\left(\frac{1}{\lambda}(\mathbb{L}_z(\boldsymbol{\theta}) - \mathbb{L}_z(\boldsymbol{\nu}))\right) dP(\boldsymbol{\nu}) \right)^{-1}. \end{aligned} \quad (163)$$

The equality in (163) implies that for all  $\delta > \delta^*$  and for all  $\theta \in \{\nu \in \mathcal{M} : \mathbb{L}_z(\nu) = \delta\}$ ,

$$\begin{aligned} \lim_{\lambda \rightarrow 0^+} \frac{dP_{\Theta|Z=z}^{(\lambda)}}{dP}(\theta) &= \left( P(\mathcal{A}_0(\theta)) + \lim_{\lambda \rightarrow 0^+} \int_{\mathcal{A}_1(\theta)} \exp\left(\frac{1}{\lambda}(\mathbb{L}_z(\theta) - \mathbb{L}_z(\nu))\right) dP(\nu) \right. \\ &\quad \left. + \lim_{\lambda \rightarrow 0^+} \int_{\mathcal{A}_2(\theta)} \exp\left(\frac{1}{\lambda}(\mathbb{L}_z(\theta) - \mathbb{L}_z(\nu))\right) dP(\nu) \right)^{-1} \end{aligned} \quad (164)$$

$$= \left( P(\mathcal{A}_0(\theta)) + \lim_{\lambda \rightarrow 0^+} \int_{\mathcal{A}_2(\theta)} \exp\left(\frac{1}{\lambda}(\mathbb{L}_z(\theta) - \mathbb{L}_z(\nu))\right) dP(\nu) \right)^{-1} \quad (165)$$

$$= (P(\mathcal{A}_0(\theta)) + \infty)^{-1} \quad (166)$$

$$= 0, \quad (167)$$

where the equality in (165) follows verifying that the dominated convergence theorem [3, 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

$$\lim_{\lambda \rightarrow 0^+} \exp\left(\frac{1}{\lambda}(\mathbb{L}_z(\theta) - \mathbb{L}_z(\nu))\right) = 0. \quad (168)$$

This completes the first part of the proof.

The second part is as follows. Consider that  $\theta \in \{\nu \in \mathcal{M} : \mathbb{L}_z(\nu) = \delta^*\}$ . Hence, the sets  $\mathcal{A}_0(\theta)$ ,  $\mathcal{A}_1(\theta)$ , and  $\mathcal{A}_2(\theta)$  satisfy the following:

$$\mathcal{A}_0(\theta) = \{\mu \in \mathcal{M} : \mathbb{L}_z(\mu) = \delta^*\}, \quad (169a)$$

$$\mathcal{A}_1(\theta) = \{\mu \in \mathcal{M} : \mathbb{L}_z(\mu) > \delta^*\}, \text{ and} \quad (169b)$$

$$\mathcal{A}_2(\theta) = \{\mu \in \mathcal{M} : \mathbb{L}_z(\mu) < \delta^*\}. \quad (169c)$$

Observing that  $P(\mathcal{A}_2(\theta)) = 0$ , plugging the equalities in (169) in (160) yields,

$$\frac{dP_{\Theta|Z=z}^{(\lambda)}}{dP}(\theta) = \left( P(\mathcal{A}_0(\theta)) + \int_{\mathcal{A}_1(\theta)} \exp\left(\frac{1}{\lambda}(\mathbb{L}_z(\theta) - \mathbb{L}_z(\nu))\right) dP(\nu) \right)^{-1}. \quad (170)$$

The equality in (170) implies that for all  $\theta \in \mathcal{L}_z(\delta^*)$ ,

$$\begin{aligned} \lim_{\lambda \rightarrow 0^+} \frac{dP_{\Theta|Z=z}^{(\lambda)}}{dP}(\theta) &= \left( P(\mathcal{A}_0(\theta)) + \lim_{\lambda \rightarrow 0^+} \int_{\mathcal{A}_1(\theta)} \exp\left(\frac{1}{\lambda}(\mathbb{L}_z(\theta) - \mathbb{L}_z(\nu))\right) dP(\nu) \right)^{-1} \\ &= \frac{1}{P(\mathcal{A}_0(\theta))}, \end{aligned} \quad (171)$$

where the equality in (171) follows from the same arguments as in (165). This completes the second part.

Finally, from (167) and (171), it follows that for all  $\theta \in \text{supp } P$ ,

$$\lim_{\lambda \rightarrow 0^+} \frac{dP_{\Theta|Z=z}^{(\lambda)}}{dP}(\theta) = \frac{1}{P(\mathcal{L}_z^*)} \mathbb{1}_{\{\theta \in \mathcal{L}_z^*\}}, \quad (172)$$

which completes the proof.

## F Proof of Lemma 2.6

For all  $\lambda \in \mathcal{K}_z$ , with  $\mathcal{K}_z$  in (34), and for all  $\mathcal{C} \in \mathcal{B}(\mathcal{M})$ ,

$$P_{\Theta|Z=z}^{(\lambda)}(\mathcal{C}) = \int_{\mathcal{C}} \frac{dP_{\Theta|Z=z}^{(\lambda)}}{dP}(\theta) dP(\theta), \quad (173)$$

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

$$P_{\Theta|Z=z}^{(\lambda)}(\mathcal{C}) = 0, \quad (174)$$

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

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

$$0 = P_{\Theta|Z=z}^{(\lambda)}(\mathcal{C}) \quad (175)$$

$$= \int_{\mathcal{C}} \frac{dP_{\Theta|Z=z}^{(\lambda)}}{dP}(\theta) dP(\theta). \quad (176)$$

From Theorem 2.1, it holds that for all  $\alpha \in \mathcal{K}_z$  and for all  $\nu \in \text{supp } P$ ,

$$\frac{dP_{\Theta|Z=z}^{(\alpha)}}{dP}(\nu) > 0, \quad (177)$$

which implies that

$$\int_{\mathcal{C}} \frac{dP_{\Theta|Z=z}^{(\lambda)}}{dP}(\theta) dP(\theta) = 0, \quad (178)$$

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

## G Proof of Lemma 2.7

The proof that the measure  $P_{\Theta|Z=z}^{(\alpha)}$  is absolutely continuous with respect to  $P_{\Theta|Z=z}^{(\beta)}$  is an immediate consequence of Lemma 2.6 and the Radon-Nikodym theorem [3, Theorem 2.2.1]. More specifically, from Lemma 2.6 it holds that for all  $(\alpha, \beta) \in \mathcal{K}_z \times \mathcal{K}_z$ ,

- (i) the measure  $P_{\Theta|Z=z}^{(\alpha)}$  is absolutely continuous with respect to  $P$ ; and
- (ii) the measure  $P$  is absolutely continuous with respect to  $P_{\Theta|Z=z}^{(\beta)}$ .

From (i) and the Radon-Nikodym theorem [3, Theorem 2.2.1], it follows that Radon-Nikodym derivative  $\frac{dP_{\Theta|Z=z}^{(\alpha)}}{dP}$  exists.

From (ii) and the Radon-Nikodym theorem [3, Theorem 2.2.1], it follows that the Radon-Nikodym derivative  $\frac{dP}{dP_{\Theta|Z=z}^{(\beta)}}$  exists.

Hence, from [28, Theorem 3.41], the Radon-Nikodym derivative  $\frac{dP_{\Theta|Z=z}^{(\alpha)}}{dP_{\Theta|Z=z}^{(\beta)}}$  exists and for all  $\theta \in \text{supp } P$ ,

$$\frac{dP_{\Theta|Z=z}^{(\alpha)}}{dP_{\Theta|Z=z}^{(\beta)}}(\theta) = \frac{dP_{\Theta|Z=z}^{(\alpha)}}{dP}(\theta) \left( \frac{dP_{\Theta|Z=z}^{(\beta)}}{dP}(\theta) \right)^{-1}, \quad (179)$$

which implies the absolute continuity of the measure  $P_{\Theta|Z=z}^{(\alpha)}$  with respect to  $P_{\Theta|Z=z}^{(\beta)}$  ([3, Theorem 2.2.1]).

The proof that the measure  $P_{\Theta|Z=z}^{(\beta)}$  is absolutely continuous with respect to  $P_{\Theta|Z=z}^{(\alpha)}$  is obtained by exchanging the choices of  $\alpha$  and  $\beta$  and completes the proof.

## H Proof of Lemma 2.8

Consider the transport of the measure  $P$  from  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$  to  $([0, +\infty), \mathcal{B}([0, +\infty)))$  through the function  $L_z$  in (5b). Denote the resulting measure in  $([0, +\infty), \mathcal{B}([0, +\infty)))$  by  $P_V$ . More specifically, for all  $\mathcal{A} \in \mathcal{B}([0, +\infty))$ , it holds that  $P_V(\mathcal{A}) = P(\{\theta \in \mathcal{M} : L_z(\theta) \in \mathcal{A}\})$ . Hence, the function  $K_z$  in (33) can be written for all  $t \in \mathbb{R}$  in terms of the measure  $P_V$  as follows

$$K_z(t) = \log \left( \int \exp(tv) dP_V(v) \right). \quad (180)$$

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

$$\phi_V(t) = \int \exp(tv) dP_V(v) = \exp(K_z(-t)). \quad (181)$$

From [15, Theorem 1a (page 439)], it follows that the function  $\phi_V$  has derivatives of all orders in  $(0, +\infty)$ , and thus, so does the function  $K_z$  in  $(-\infty, 0)$ . This implies the continuity of  $K_z$  in  $(-\infty, 0)$ , and completes the proof.

## I Proof of Lemma 2.9

Let  $(\gamma_1, \gamma_2) \in \mathcal{K}_z \times \mathcal{K}_z$  and  $\alpha \in [0, 1]$  be fixed. When  $\alpha = 0$  or  $\alpha = 1$ , the proof is trivial. Then, for all  $\alpha \in (0, 1)$ , the following holds

$$\begin{aligned} & \alpha K_z(\gamma_1) + (1 - \alpha) K_z(\gamma_2) \\ &= \alpha \log \left( \int \exp(\gamma_1 \mathbf{L}_z(\mathbf{u})) dP(\mathbf{u}) \right) + (1 - \alpha) \log \left( \int \exp(\gamma_2 \mathbf{L}_z(\mathbf{u})) dP(\mathbf{u}) \right) \end{aligned} \quad (182)$$

$$= \log \left( \left( \int \exp(\gamma_1 \mathbf{L}_z(\mathbf{u})) dP(\mathbf{u}) \right)^\alpha \right) + \log \left( \left( \int \exp(\gamma_2 \mathbf{L}_z(\mathbf{u})) dP(\mathbf{u}) \right)^{(1-\alpha)} \right) \quad (183)$$

$$= \log \left( \left( \int \exp(\gamma_1 \mathbf{L}_z(\mathbf{u})) dP(\mathbf{u}) \right)^\alpha \left( \int \exp(\gamma_2 \mathbf{L}_z(\mathbf{u})) dP(\mathbf{u}) \right)^{(1-\alpha)} \right) \quad (184)$$

$$= \log \left( \left( \int \exp(\alpha \gamma_1 \mathbf{L}_z(\mathbf{u}))^{\frac{1}{\alpha}} dP(\mathbf{u}) \right)^\alpha \left( \int \exp((1 - \alpha) \gamma_2 \mathbf{L}_z(\mathbf{u}))^{\frac{1}{1-\alpha}} dP(\mathbf{u}) \right)^{(1-\alpha)} \right) \quad (185)$$

$$\geq \log \left( \int \exp(\alpha \gamma_1 \mathbf{L}_z(\mathbf{u})) \exp((1 - \alpha) \gamma_2 \mathbf{L}_z(\mathbf{u})) dP(\mathbf{u}) \right) \quad (186)$$

$$= \log \left( \int \exp((\alpha \gamma_1 + (1 - \alpha) \gamma_2) \mathbf{L}_z(\mathbf{u})) dP(\mathbf{u}) \right) \quad (187)$$

$$= K_z(\alpha \gamma_1 + (1 - \alpha) \gamma_2), \quad (188)$$

where the inequality in (186) follows from Hölder's inequality. Note that if the function  $\mathbf{L}_z$  in (5b) is a constant, then (186) holds with equality. Alternatively, under the assumption that there exist at least two disjoint subsets  $\mathcal{A}$  and  $\mathcal{B}$  of  $\mathcal{M}$  that are nonnegligible with respect to  $P$  and satisfy (54) for some positive real  $c > 0$ , the inequality is strict. This completes the proof.

## J Proof of Lemma 2.10

For all  $s \in \mathcal{K}_z$ , with  $\mathcal{K}_z$  in (34), the equality in (56) implies the following,

$$K_z^{(1)} \left( -\frac{1}{s} \right) = \frac{d}{dt} \log \left( \int \exp(t \mathbf{L}_z(\mathbf{u})) dP(\mathbf{u}) \right) \Big|_{t=-\frac{1}{s}} \quad (189)$$

$$= \frac{1}{\int \exp(t \mathbf{L}_z(\mathbf{v})) dP(\mathbf{v})} \int \mathbf{L}_z(\mathbf{u}) \exp(t \mathbf{L}_z(\mathbf{u})) dP(\mathbf{u}) \Big|_{t=-\frac{1}{s}} \quad (190)$$

$$= \frac{1}{\int \exp\left(-\frac{1}{s} \mathbf{L}_z(\mathbf{v})\right) dP(\mathbf{v})} \int \mathbf{L}_z(\mathbf{u}) \exp\left(-\frac{1}{s} \mathbf{L}_z(\mathbf{u})\right) dP(\mathbf{u}) \quad (191)$$

$$= \exp\left(-K_z\left(-\frac{1}{s}\right)\right) \int \mathbf{L}_z(\mathbf{u}) \exp\left(-\frac{1}{s} \mathbf{L}_z(\mathbf{u})\right) dP(\mathbf{u}) \quad (192)$$

$$= \int \mathbf{L}_z(\mathbf{u}) \exp\left(-K_z\left(-\frac{1}{s}\right) - \frac{1}{s} \mathbf{L}_z(\mathbf{u})\right) dP(\mathbf{u}) \quad (193)$$

$$= \int \mathbf{L}_z(\boldsymbol{\theta}) dP_{\Theta|Z=z}^{(s)}(\boldsymbol{\theta}), \quad (194)$$

where the equality in (190) holds from the dominated convergence theorem [3]; the equality in (192) follows from (33); and the equality in (194) follows from (37).

For all  $s \in \mathcal{K}_z$ , with  $\mathcal{K}_z$  in (34), the equalities in (56) and (193) imply that

$$K_z^{(2)}\left(-\frac{1}{s}\right) = \frac{d}{dt} \int \mathbf{L}_z(\mathbf{u}) \exp(-K_z(t) + t\mathbf{L}_z(\mathbf{u})) dP(\mathbf{u}) \Big|_{t=-\frac{1}{s}} \quad (195)$$

$$= \int \mathbf{L}_z(\mathbf{u}) \left(-K_z^{(1)}(t) + \mathbf{L}_z(\mathbf{u})\right) \exp(-K_z(t) + t\mathbf{L}_z(\mathbf{u})) dP(\mathbf{u}) \Big|_{t=-\frac{1}{s}} \quad (196)$$

$$= \int \mathbf{L}_z(\mathbf{u}) \left(-K_z^{(1)}\left(-\frac{1}{s}\right) + \mathbf{L}_z(\mathbf{u})\right) \exp\left(-K_z\left(-\frac{1}{s}\right) - \frac{1}{s}\mathbf{L}_z(\mathbf{u})\right) dP(\mathbf{u}) \quad (197)$$

$$= \int \mathbf{L}_z(\mathbf{u}) \left(-K_z^{(1)}\left(-\frac{1}{s}\right) + \mathbf{L}_z(\mathbf{u})\right) dP_{\Theta|Z=z}^{(s)}(\mathbf{u}) \quad (198)$$

$$= -K_z^{(1)}\left(-\frac{1}{s}\right) \int \mathbf{L}_z(\mathbf{u}) dP_{\Theta|Z=z}^{(s)}(\mathbf{u}) + \int (\mathbf{L}_z(\mathbf{u}))^2 dP_{\Theta|Z=z}^{(s)}(\mathbf{u}) \quad (199)$$

$$= -\left(K_z^{(1)}\left(-\frac{1}{s}\right)\right)^2 + \int (\mathbf{L}_z(\mathbf{u}))^2 dP_{\Theta|Z=z}^{(s)}(\mathbf{u}) \quad (200)$$

$$= \int \left(\mathbf{L}_z(\mathbf{u}) - K_z^{(1)}\left(-\frac{1}{s}\right)\right)^2 dP_{\Theta|Z=z}^{(s)}(\mathbf{u}), \quad (201)$$

where the equality in (196) follows from the dominated convergence theorem [3]; the equality in (198) is due to a change of measure through the Radon-Nikodym derivative in (37); and the equality in (200) follows from (194).

For all  $s \in \mathcal{K}_z$ , with  $\mathcal{K}_z$  in (34), the equalities in (56) and (200) imply that

$$K_z^{(3)}\left(-\frac{1}{s}\right) = \frac{d}{dt} \left( \int (\mathbf{L}_z(\mathbf{u}))^2 dP_{\Theta|Z=z}^{(-\frac{1}{t})}(\mathbf{u}) - (K_z^{(1)}(t))^2 \right) \Big|_{t=-\frac{1}{s}} \quad (202)$$

$$= \frac{d}{dt} \left( \int (\mathbf{L}_z(\mathbf{u}))^2 \exp(-K_z(t) + t\mathbf{L}_z(\mathbf{u})) dP(\mathbf{u}) - (K_z^{(1)}(t))^2 \right) \Big|_{t=-\frac{1}{s}} \quad (203)$$

$$= \int (\mathbf{L}_z(\mathbf{u}))^2 \left( \frac{d}{dt} \exp(-K_z(t) + t\mathbf{L}_z(\mathbf{u})) \Big|_{t=-\frac{1}{s}} \right) dP(\mathbf{u}) - 2K_z^{(1)}(t) K_z^{(2)}(t) \Big|_{t=-\frac{1}{s}} \quad (204)$$

$$= \int (\mathbf{L}_z(\mathbf{u}))^2 \left( (\mathbf{L}_z(\mathbf{u}) - K_z^{(1)}(t)) \exp(-K_z(t) + t\mathbf{L}_z(\mathbf{u})) \Big|_{t=-\frac{1}{s}} \right) dP(\mathbf{u}) - 2K_z^{(1)}(t) K_z^{(2)}(t) \Big|_{t=-\frac{1}{s}} \quad (205)$$

$$= \int (\mathbf{L}_z(\mathbf{u}))^2 \left( \mathbf{L}_z(\mathbf{u}) - K_z^{(1)}\left(-\frac{1}{s}\right) \right) \exp\left(-K_z\left(-\frac{1}{s}\right) - \frac{1}{s}\mathbf{L}_z(\mathbf{u})\right) dP(\mathbf{u}) - 2K_z^{(1)}\left(-\frac{1}{s}\right) K_z^{(2)}\left(-\frac{1}{s}\right) \quad (206)$$

$$= \int (\mathbf{L}_z(\mathbf{u}))^2 \left( \mathbf{L}_z(\mathbf{u}) - K_z^{(1)}\left(-\frac{1}{s}\right) \right) dP_{\Theta|Z=z}^{(s)}(\mathbf{u}) - 2K_z^{(1)}\left(-\frac{1}{s}\right) K_z^{(2)}\left(-\frac{1}{s}\right) \quad (207)$$

$$= \int (\mathbf{L}_z(\mathbf{u}))^3 dP_{\Theta|Z=z}^{(s)}(\mathbf{u}) - K_z^{(1)}\left(-\frac{1}{s}\right) \int (\mathbf{L}_z(\mathbf{u}))^2 dP_{\Theta|Z=z}^{(s)}(\mathbf{u}) - 2K_z^{(1)}\left(-\frac{1}{s}\right) K_z^{(2)}\left(-\frac{1}{s}\right) \quad (208)$$

$$= \int (\mathbf{L}_z(\mathbf{u}))^3 dP_{\Theta|Z=z}^{(s)}(\mathbf{u}) - K_z^{(1)}\left(-\frac{1}{s}\right) \left( K_z^{(2)}\left(-\frac{1}{s}\right) + \left( K_z^{(1)}\left(-\frac{1}{s}\right) \right)^2 \right) - 2K_z^{(1)}\left(-\frac{1}{s}\right) K_z^{(2)}\left(-\frac{1}{s}\right) \quad (209)$$

$$= \int (\mathbf{L}_z(\mathbf{u}))^3 dP_{\Theta|Z=z}^{(s)}(\mathbf{u}) - K_z^{(1)}\left(-\frac{1}{s}\right)^3 - 3K_z^{(1)}\left(-\frac{1}{s}\right) K_z^{(2)}\left(-\frac{1}{s}\right) \quad (210)$$

$$= \int \left( \mathbf{L}_z(\mathbf{u}) - K_z^{(1)}\left(-\frac{1}{s}\right) \right)^3 dP_{\Theta|Z=z}^{(s)}(\mathbf{u}), \quad (211)$$

where the equality in (203) follows from (37); the equality in (204) follows from the dominated convergence theorem [3]; the equality in (207) follows from (37); and the equality in (209) follows from (200). This completes the proof.

## K Proof of Theorem 2.2

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

$$\frac{d}{d\lambda} \int \mathbb{L}_{\mathbf{z}}(\boldsymbol{\theta}) \frac{dP_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(\lambda)}(\boldsymbol{\theta})}{dP}(\boldsymbol{\theta}) dP(\boldsymbol{\theta}) \quad (212)$$

$$= \frac{d}{d\lambda} \int \frac{\mathbb{L}_{\mathbf{z}}(\boldsymbol{\theta}) \exp\left(-\frac{\mathbb{L}_{\mathbf{z}}(\boldsymbol{\theta})}{\lambda}\right)}{\int \exp\left(-\frac{\mathbb{L}_{\mathbf{z}}(\boldsymbol{\nu})}{\lambda}\right) dP(\boldsymbol{\nu})} dP(\boldsymbol{\theta}) \quad (213)$$

$$= \int \mathbb{L}_{\mathbf{z}}(\boldsymbol{\theta}) \frac{d}{d\lambda} \frac{\exp\left(-\frac{\mathbb{L}_{\mathbf{z}}(\boldsymbol{\theta})}{\lambda}\right)}{\int \exp\left(-\frac{\mathbb{L}_{\mathbf{z}}(\boldsymbol{\nu})}{\lambda}\right) dP(\boldsymbol{\nu})} dP(\boldsymbol{\theta}) \quad (214)$$

$$= \frac{1}{\lambda^2} \int (\mathbb{L}_{\mathbf{z}}(\boldsymbol{\theta}))^2 \frac{\exp\left(-\frac{\mathbb{L}_{\mathbf{z}}(\boldsymbol{\theta})}{\lambda}\right)}{\int \exp\left(-\frac{\mathbb{L}_{\mathbf{z}}(\boldsymbol{\nu})}{\lambda}\right) dP(\boldsymbol{\nu})} dP(\boldsymbol{\theta})$$

$$- \frac{1}{\lambda^2} \int \frac{\mathbb{L}_{\mathbf{z}}(\boldsymbol{\theta}) \exp\left(-\frac{\mathbb{L}_{\mathbf{z}}(\boldsymbol{\theta})}{\lambda}\right) \int \mathbb{L}_{\mathbf{z}}(\boldsymbol{\alpha}) \exp\left(-\frac{\mathbb{L}_{\mathbf{z}}(\boldsymbol{\alpha})}{\lambda}\right) dP(\boldsymbol{\alpha})}{\left(\int \exp\left(-\frac{\mathbb{L}_{\mathbf{z}}(\boldsymbol{\nu})}{\lambda}\right) dP(\boldsymbol{\nu})\right)^2} dP(\boldsymbol{\theta}) \quad (215)$$

$$= \frac{1}{\lambda^2} \int (\mathbb{L}_{\mathbf{z}}(\boldsymbol{\theta}))^2 \frac{dP_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(\lambda)}(\boldsymbol{\theta})}{dP}(\boldsymbol{\theta}) dP(\boldsymbol{\theta})$$

$$- \frac{1}{\lambda^2} \left( \int \mathbb{L}_{\mathbf{z}}(\boldsymbol{\theta}) \frac{dP_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(\lambda)}(\boldsymbol{\theta})}{dP}(\boldsymbol{\theta}) dP(\boldsymbol{\theta}) \right)^2 \quad (216)$$

$$= \frac{1}{\lambda^2} \int \left( \mathbb{L}_{\mathbf{z}}(\boldsymbol{\theta}) - \int \mathbb{L}_{\mathbf{z}}(\boldsymbol{\nu}) dP_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(\lambda)}(\boldsymbol{\nu}) \right)^2 dP_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(\lambda)}(\boldsymbol{\theta}) \quad (217)$$

$$\geq 0, \quad (218)$$

where the equality in (214) follows from the dominated convergence theorem [3].

The inequality in (218) implies that the the expected empirical risk  $K_{\mathbf{z}}^{(1)}\left(-\frac{1}{\lambda}\right)$  in (33) is non-decreasing with respect to  $\lambda$ . Hence, given two reals  $\lambda_1$  and  $\lambda_2$ , with  $\lambda_1 > \lambda_2 > 0$ , it holds that

$$K_{\mathbf{z}}^{(1)}\left(-\frac{1}{\lambda_1}\right) \geq K_{\mathbf{z}}^{(1)}\left(-\frac{1}{\lambda_2}\right). \quad (219)$$

The rest of the proof consists in showing that for all  $\alpha \in \mathcal{K}_{\mathbf{z}}$ , the function  $K_{\mathbf{z}}^{(2)}$  in (56) satisfies  $K_{\mathbf{z}}^{(2)}\left(-\frac{1}{\alpha}\right) > 0$  if and only if the function  $\mathbb{L}_{\mathbf{z}}$  in (5b) is separable. Hence, the proof is divided into two parts. In the first part, it is shown that if for all  $\alpha \in \mathcal{K}_{\mathbf{z}}$ ,  $K_{\mathbf{z}}^{(2)}\left(-\frac{1}{\alpha}\right) > 0$ , then the function  $\mathbb{L}_{\mathbf{z}}$  in (5b) is separable. The second part of the proof, consists in showing that if the function  $\mathbb{L}_{\mathbf{z}}$  is separable, then, for all  $\alpha \in \mathcal{K}_{\mathbf{z}}$ ,  $K_{\mathbf{z}}^{(2)}\left(-\frac{1}{\alpha}\right) > 0$ .



The first part is as follows. From Lemma 2.10, it holds that for all  $\alpha \in \mathcal{K}_z$ ,

$$K_z^{(2)}\left(-\frac{1}{\alpha}\right) = \int \left( L_z(\boldsymbol{\theta}) - K_z^{(1)}\left(-\frac{1}{\alpha}\right) \right)^2 dP_{\boldsymbol{\Theta}|Z=z}^{(\alpha)}(\boldsymbol{\theta}) \quad (220)$$

$$= \int_{\mathcal{R}_0(\alpha)} \left( L_z(\boldsymbol{\theta}) - K_z^{(1)}\left(-\frac{1}{\alpha}\right) \right)^2 dP_{\boldsymbol{\Theta}|Z=z}^{(\alpha)}(\boldsymbol{\theta}) \quad (221)$$

$$+ \int_{\mathcal{R}_1(\alpha)} \left( L_z(\boldsymbol{\theta}) - K_z^{(1)}\left(-\frac{1}{\alpha}\right) \right)^2 dP_{\boldsymbol{\Theta}|Z=z}^{(\alpha)}(\boldsymbol{\theta}) \quad (222)$$

$$+ \int_{\mathcal{R}_2(\alpha)} \left( L_z(\boldsymbol{\theta}) - K_z^{(1)}\left(-\frac{1}{\alpha}\right) \right)^2 dP_{\boldsymbol{\Theta}|Z=z}^{(\alpha)}(\boldsymbol{\theta}), \quad (223)$$

where the sets  $\mathcal{R}_0(\alpha)$ ,  $\mathcal{R}_1(\alpha)$ , and  $\mathcal{R}_2(\alpha)$  are respectively defined in (76). Hence,

$$K_z^{(2)}\left(-\frac{1}{\alpha}\right) = \int_{\mathcal{R}_1(\alpha)} \left( L_z(\boldsymbol{\theta}) - K_z^{(1)}\left(-\frac{1}{\alpha}\right) \right)^2 dP_{\boldsymbol{\Theta}|Z=z}^{(\alpha)}(\boldsymbol{\theta}) \\ + \int_{\mathcal{R}_2(\alpha)} \left( L_z(\boldsymbol{\theta}) - K_z^{(1)}\left(-\frac{1}{\alpha}\right) \right)^2 dP_{\boldsymbol{\Theta}|Z=z}^{(\alpha)}(\boldsymbol{\theta}) \quad (224)$$

$$\leq \left( \inf_{\boldsymbol{\theta}} L_z(\boldsymbol{\theta}) - K_z^{(1)}\left(-\frac{1}{\alpha}\right) \right)^2 P_{\boldsymbol{\Theta}|Z=z}^{(\alpha)}(\mathcal{R}_1(\alpha)) \\ + \left( \sup_{\boldsymbol{\theta}} L_z(\boldsymbol{\theta}) - K_z^{(1)}\left(-\frac{1}{\alpha}\right) \right)^2 P_{\boldsymbol{\Theta}|Z=z}^{(\alpha)}(\mathcal{R}_2(\alpha)). \quad (225)$$

Under the assumption that for all  $\alpha \in \mathcal{K}_z$  the function  $K_z^{(2)}$  in (56) satisfies  $K_z^{(2)}\left(-\frac{1}{\alpha}\right) > 0$ , it follows from (225) that

$$0 < \left( \inf_{\boldsymbol{\theta}} L_z(\boldsymbol{\theta}) - K_z^{(1)}\left(-\frac{1}{\alpha}\right) \right)^2 P_{\boldsymbol{\Theta}|Z=z}^{(\alpha)}(\mathcal{R}_1(\alpha)) \\ + \left( \sup_{\boldsymbol{\theta}} L_z(\boldsymbol{\theta}) - K_z^{(1)}\left(-\frac{1}{\alpha}\right) \right)^2 P_{\boldsymbol{\Theta}|Z=z}^{(\alpha)}(\mathcal{R}_2(\alpha)). \quad (226)$$

Note that if

$$P_{\boldsymbol{\Theta}|Z=z}^{(\alpha)}(\mathcal{R}_1(\alpha)) > 0, \quad (227)$$

then,  $\inf_{\boldsymbol{\theta}} L_z(\boldsymbol{\theta}) \neq K_z^{(1)}\left(-\frac{1}{\alpha}\right)$ . Moreover, if

$$P_{\boldsymbol{\Theta}|Z=z}^{(\alpha)}(\mathcal{R}_2(\alpha)) > 0, \quad (228)$$

then,  $\sup_{\boldsymbol{\theta}} L_z(\boldsymbol{\theta}) \neq K_z^{(1)}\left(-\frac{1}{\alpha}\right)$ . Therefore, the inequality in (226) implies that at least one of the following claims is true:

(a)  $P_{\boldsymbol{\Theta}|Z=z}^{(\alpha)}(\mathcal{R}_1(\alpha)) > 0$ ; and

(b)  $P_{\boldsymbol{\Theta}|Z=z}^{(\alpha)}(\mathcal{R}_2(\alpha)) > 0$ .

Nonetheless, from Lemma 2.13, 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  $P_{\boldsymbol{\Theta}|Z=z}^{(\alpha)}$  and moreover, it holds that for all  $(\boldsymbol{\mu}, \boldsymbol{\nu}) \in \mathcal{R}_1(\alpha) \times \mathcal{R}_2(\alpha)$ ,

$$L_z(\boldsymbol{\nu}_1) > K_z^{(1)}\left(-\frac{1}{\alpha}\right) > L_z(\boldsymbol{\nu}_2). \quad (229)$$

This proves that under the assumption that for all  $\alpha \in \mathcal{K}_z$ ,  $K_z^{(2)}\left(-\frac{1}{\alpha}\right) > 0$ , the function  $L_z$  in (5b) is separable. This completes the first part of the proof.

The second part of the proof is simpler. Assume that the empirical risk function  $L_z$  in (5b) is separable. That is, for all  $\gamma \in \mathcal{K}_z$ , 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  $P_{\Theta|Z=z}^{(\gamma)}$  in (37) and verify that for all  $(\nu_1, \nu_2) \in \mathcal{A}(\gamma) \times \mathcal{B}(\gamma)$ ,

$$L_z(\nu_1) > c_\gamma > L_z(\nu_2). \quad (230)$$

From Lemma 2.10, it holds that for all  $\gamma \in \mathcal{K}_z$ ,

$$K_z^{(2)}\left(-\frac{1}{\gamma}\right) = \int \left(L_z(\theta) - K_z^{(1)}\left(-\frac{1}{\gamma}\right)\right)^2 dP_{\Theta|Z=z}^{(\gamma)}(\theta) \quad (231)$$

$$= \int_{\mathcal{A}(\gamma)} \left(L_z(\theta) - K_z^{(1)}\left(-\frac{1}{\gamma}\right)\right)^2 dP_{\Theta|Z=z}^{(\gamma)}(\theta) \quad (232)$$

$$+ \int_{\mathcal{B}(\gamma)} \left(L_z(\theta) - K_z^{(1)}\left(-\frac{1}{\gamma}\right)\right)^2 dP_{\Theta|Z=z}^{(\gamma)}(\theta) \quad (233)$$

$$+ \int_{(\mathcal{A}(\gamma) \cup \mathcal{B}(\gamma))^c} \left(L_z(\theta) - K_z^{(1)}\left(-\frac{1}{\gamma}\right)\right)^2 dP_{\Theta|Z=z}^{(\gamma)}(\theta) \quad (234)$$

$$> 0, \quad (235)$$

where the inequality (235) follows from the following facts. First, if  $c_\gamma < K_z^{(1)}\left(-\frac{1}{\gamma}\right)$ , with  $c_\gamma$  in (230), then for all  $\nu \in \mathcal{B}(\gamma)$ , it holds that

$$\left(L_z(\theta) - K_z^{(1)}\left(-\frac{1}{\gamma}\right)\right)^2 > \left(c_\gamma - K_z^{(1)}\left(-\frac{1}{\gamma}\right)\right)^2, \quad (236)$$

and thus,

$$\int_{\mathcal{B}(\gamma)} \left(L_z(\theta) - K_z^{(1)}\left(-\frac{1}{\gamma}\right)\right)^2 dP_{\Theta|Z=z}^{(\gamma)}(\theta) > \left(c_\gamma - K_z^{(1)}\left(-\frac{1}{\gamma}\right)\right)^2 P_{\Theta|Z=z}^{(\gamma)}(\mathcal{B}(\gamma)) \quad (237)$$

$$> 0. \quad (238)$$

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

$$\left(L_z(\theta) - K_z^{(1)}\left(-\frac{1}{\gamma}\right)\right)^2 > \left(c_\gamma - K_z^{(1)}\left(-\frac{1}{\gamma}\right)\right)^2, \quad (239)$$

and thus,

$$\int_{\mathcal{A}(\gamma)} \left(L_z(\theta) - K_z^{(1)}\left(-\frac{1}{\gamma}\right)\right)^2 dP_{\Theta|Z=z}^{(\gamma)}(\theta) > \left(c_\gamma - K_z^{(1)}\left(-\frac{1}{\gamma}\right)\right)^2 P_{\Theta|Z=z}^{(\gamma)}(\mathcal{A}(\gamma)) \quad (240)$$

$$\geq 0. \quad (241)$$

Hence, under the assumption that the empirical risk function  $L_z$  in (5b) is separable, it holds that for all  $\gamma \in \mathcal{K}_z$ ,  $K_z^{(2)}\left(-\frac{1}{\gamma}\right) > 0$ . This completes the proof.

## L Proof of Lemma 2.11

Consider the following sets

$$\mathcal{A}_0 \triangleq \{\boldsymbol{\nu} \in \mathcal{M} : \mathsf{L}_z(\boldsymbol{\nu}) = \delta^*\}, \quad (242)$$

$$\mathcal{A}_1 \triangleq \{\boldsymbol{\nu} \in \mathcal{M} : \mathsf{L}_z(\boldsymbol{\nu}) < \delta^*\}, \text{ and} \quad (243)$$

$$\mathcal{A}_2 \triangleq \{\boldsymbol{\nu} \in \mathcal{M} : \mathsf{L}_z(\boldsymbol{\nu}) > \delta^*\}, \quad (244)$$

with  $\delta^*$  in (46) and the function  $\mathsf{L}_z$  in (5b).

From (57), for all  $\lambda \in \mathcal{K}_z$ , with  $\mathcal{K}_z$  in (34), it holds that,

$$K_z^{(1)}\left(-\frac{1}{\lambda}\right) = \int_{\mathcal{A}_0} \mathsf{L}_z(\boldsymbol{\theta}) \, dP_{\Theta|Z=z}^{(\lambda)}(\boldsymbol{\theta}) + \int_{\mathcal{A}_1} \mathsf{L}_z(\boldsymbol{\theta}) \, dP_{\Theta|Z=z}^{(\lambda)}(\boldsymbol{\theta}) + \int_{\mathcal{A}_2} \mathsf{L}_z(\boldsymbol{\theta}) \, dP_{\Theta|Z=z}^{(\lambda)}(\boldsymbol{\theta}) \quad (245)$$

$$= \int_{\mathcal{A}_0} \mathsf{L}_z(\boldsymbol{\theta}) \, dP_{\Theta|Z=z}^{(\lambda)}(\boldsymbol{\theta}) + \int_{\mathcal{A}_2} \mathsf{L}_z(\boldsymbol{\theta}) \, dP_{\Theta|Z=z}^{(\lambda)}(\boldsymbol{\theta}) \quad (246)$$

$$= \delta^* P_{\Theta|Z=z}^{(\lambda)}(\mathcal{L}_z^*) + \int_{\mathcal{A}_2} \mathsf{L}_z(\boldsymbol{\theta}) \, dP_{\Theta|Z=z}^{(\lambda)}(\boldsymbol{\theta}), \quad (247)$$

where, the equality in (246) follows from noticing that  $P(\mathcal{A}_1) = 0$  and thus,  $P_{\Theta|Z=z}^{(\lambda)}(\mathcal{A}_1) = 0$  (Lemma 2.6); and the equality in (247) follows from noticing that  $\mathcal{A}_0 = \mathcal{L}_z^*$ , with  $\mathcal{L}_z^*$  in (47).

The rest of the proof is divided into two parts. The first part considers the case in which the function  $\mathsf{L}_z$  in (5b) is separable with respect to the measure  $P$ . The second part considers the converse case.

The first part is as follows. Under the assumption that the function  $\mathsf{L}_z$  in (5b) is separable with respect to the measure  $P$ , if  $P(\mathcal{A}_0) > 0$ , it follows that there exists a real  $c > \delta^*$ , such that  $P(\{\boldsymbol{\nu} \in \mathcal{M} : \mathsf{L}_z(\boldsymbol{\nu}) > c\}) > 0$ . By observing that  $\{\boldsymbol{\nu} \in \mathcal{M} : \mathsf{L}_z(\boldsymbol{\nu}) > c\} \subset \mathcal{A}_2$ , it follows that  $P(\mathcal{A}_2) > 0$ . Moreover, if  $P(\mathcal{A}_0) = 0$ , it holds that  $P(\mathcal{A}_2) = 1$ . Hence, from (247), for all  $\lambda \in \mathcal{K}_z$ , it holds that,

$$K_z^{(1)}\left(-\frac{1}{\lambda}\right) > \delta^* P_{\Theta|Z=z}^{(\lambda)}(\mathcal{L}_z^*) + \delta^* P_{\Theta|Z=z}^{(\lambda)}(\mathcal{A}_2) \quad (248)$$

$$= \delta^*, \quad (249)$$

where, the inequality in (248) follows from noticing that for all  $\boldsymbol{\theta} \in \mathcal{A}_2$ , it holds that  $\mathsf{L}_z(\boldsymbol{\theta}) > \delta^*$ .

The second part of the proof is as follows. Under the assumption that the function  $\mathsf{L}_z$  in (5b) is not separable with respect to the measure  $P$ , it holds that  $P(\mathcal{A}_0) = 1$  and  $P(\mathcal{A}_1) = P(\mathcal{A}_2) = 0$ . Hence, from (247), it follows that

$$K_z^{(1)}\left(-\frac{1}{\lambda}\right) = \delta^*. \quad (250)$$

This completes the proof.

## M Proof of Lemma 2.12

Consider the following sets

$$\mathcal{A}_0 \triangleq \{\boldsymbol{\nu} \in \mathcal{M} : \mathbf{L}_z(\boldsymbol{\nu}) = \delta^*\}, \quad (251)$$

$$\mathcal{A}_1 \triangleq \{\boldsymbol{\nu} \in \mathcal{M} : \mathbf{L}_z(\boldsymbol{\nu}) < \delta^*\}, \text{ and} \quad (252)$$

$$\mathcal{A}_2 \triangleq \{\boldsymbol{\nu} \in \mathcal{M} : \mathbf{L}_z(\boldsymbol{\nu}) > \delta^*\}, \quad (253)$$

with  $\delta^*$  in (46) and the function  $\mathbf{L}_z$  in (5b).

For all  $\lambda \in \mathcal{K}_z$ , the following holds,

$$1 = P_{\Theta|Z=z}^{(\lambda)}(\mathcal{A}_0) + P_{\Theta|Z=z}^{(\lambda)}(\mathcal{A}_1) + P_{\Theta|Z=z}^{(\lambda)}(\mathcal{A}_2) \quad (254)$$

$$= P_{\Theta|Z=z}^{(\lambda)}(\mathcal{A}_0) + P_{\Theta|Z=z}^{(\lambda)}(\mathcal{A}_2) \quad (255)$$

$$= P_{\Theta|Z=z}^{(\lambda)}(\mathcal{A}_0) + \int_{\mathcal{A}_2} dP_{\Theta|Z=z}^{(\lambda)}(\boldsymbol{\theta}), \quad (256)$$

where, the equality in (255) follows from noticing that  $P(\mathcal{A}_1) = 0$  and thus,  $P_{\Theta|Z=z}^{(\lambda)}(\mathcal{A}_1) = 0$  (Lemma 2.6).

The above implies that

$$1 = \lim_{\lambda \rightarrow 0^+} P_{\Theta|Z=z}^{(\lambda)}(\mathcal{A}_0) + \lim_{\lambda \rightarrow 0^+} \int_{\mathcal{A}_2} \frac{dP_{\Theta|Z=z}^{(\lambda)}(\boldsymbol{\theta})}{dP}(\boldsymbol{\theta}) dP(\boldsymbol{\theta}) \quad (257)$$

$$= \lim_{\lambda \rightarrow 0^+} P_{\Theta|Z=z}^{(\lambda)}(\mathcal{A}_0) + \int_{\mathcal{A}_2} \lim_{\lambda \rightarrow 0^+} \frac{dP_{\Theta|Z=z}^{(\lambda)}(\boldsymbol{\theta})}{dP}(\boldsymbol{\theta}) dP(\boldsymbol{\theta}) \quad (258)$$

$$= \lim_{\lambda \rightarrow 0^+} P_{\Theta|Z=z}^{(\lambda)}(\mathcal{A}_0), \quad (259)$$

where, the equality in (258) follows from noticing two facts: (a) For all  $\lambda \in \mathcal{K}_z$ , the Radon-Nikodym derivative  $\frac{dP_{\Theta|Z=z}^{(\lambda)}}{dP}$  is positive and bounded (Lemma 2.3); and (b) For all  $\boldsymbol{\theta} \in \mathcal{A}_2$ , it holds that  $\lim_{\lambda \rightarrow 0^+} \frac{dP_{\Theta|Z=z}^{(\lambda)}(\boldsymbol{\theta})}{dP}(\boldsymbol{\theta}) = 0$ . Hence, the dominated convergence theorem [3, Theorem 1.6.9] holds.

Finally, by noticing that  $\mathcal{A}_0 = \mathcal{L}_z^*$ , it holds that

$$\lim_{\lambda \rightarrow 0^+} P_{\Theta|Z=z}^{(\lambda)}(\mathcal{L}_z^*) = 1, \quad (260)$$

which completes the proof.

## N Proof of Theorem 2.3

From (247) in the proof of Lemma 2.11, it holds that

$$\lim_{\lambda \rightarrow 0^+} K_z^{(1)} \left( -\frac{1}{\lambda} \right) = \lim_{\lambda \rightarrow 0^+} \delta^* P_{\Theta|Z=z}^{(\lambda)}(\mathcal{L}_z^*) + \lim_{\lambda \rightarrow 0^+} \int_{\mathcal{A}_2} \mathbb{L}_z(\boldsymbol{\theta}) dP_{\Theta|Z=z}^{(\lambda)}(\boldsymbol{\theta}) \quad (261)$$

$$= \lim_{\lambda \rightarrow 0^+} \delta^* P_{\Theta|Z=z}^{(\lambda)}(\mathcal{L}_z^*) + \lim_{\lambda \rightarrow 0^+} \int_{\mathcal{A}_2} \mathbb{L}_z(\boldsymbol{\theta}) \frac{dP_{\Theta|Z=z}^{(\lambda)}(\boldsymbol{\theta})}{dP}(\boldsymbol{\theta}) dP(\boldsymbol{\theta}) \quad (262)$$

$$= \lim_{\lambda \rightarrow 0^+} \delta^* P_{\Theta|Z=z}^{(\lambda)}(\mathcal{L}_z^*) + \int_{\mathcal{A}_2} \mathbb{L}_z(\boldsymbol{\theta}) \lim_{\lambda \rightarrow 0^+} \frac{dP_{\Theta|Z=z}^{(\lambda)}(\boldsymbol{\theta})}{dP}(\boldsymbol{\theta}) dP(\boldsymbol{\theta}) \quad (263)$$

$$= \delta^* \lim_{\lambda \rightarrow 0^+} P_{\Theta|Z=z}^{(\lambda)}(\mathcal{L}_z^*) \quad (264)$$

$$= \delta^*, \quad (265)$$

where, the equality in (263) follows from noticing two facts: (a) For all  $\lambda \in \mathcal{K}_z$ , the Randon-Nikodym derivative  $\frac{dP_{\Theta|Z=z}^{(\lambda)}}{dP}$  is positive and bounded (Lemma 2.3); and (b) For all  $\boldsymbol{\theta} \in \mathcal{A}_2$ , it holds that  $\lim_{\lambda \rightarrow 0^+} \frac{dP_{\Theta|Z=z}^{(\lambda)}(\boldsymbol{\theta})}{dP}(\boldsymbol{\theta}) = 0$ . Hence, the dominated convergence theorem [3, Theorem 1.6.9] holds. The inequality in (264) follows from Lemma 2.12 This completes the proof.

## O Proof of Lemma 2.13

The proof is divided into two parts. In the first part, given a real  $\alpha \in \mathcal{K}_z$ , it is proven that if the set  $\mathcal{R}_1(\alpha)$  is nonnegligible with respect to  $P_{\Theta|Z=z}^{(\alpha)}$ , then the set  $\mathcal{R}_2(\alpha)$  is nonnegligible with respect to  $P_{\Theta|Z=z}^{(\alpha)}$ . The second part of the proof consists in proving that, given a real  $\alpha \in \mathcal{K}_z$ , if the set  $\mathcal{R}_2(\alpha)$  is nonnegligible with respect to  $P_{\Theta|Z=z}^{(\alpha)}$ , then the set  $\mathcal{R}_1(\alpha)$  is nonnegligible with respect to  $P_{\Theta|Z=z}^{(\alpha)}$ .

The first part is proved by contradiction. Assume that set  $\mathcal{R}_2(\alpha)$  is negligible with respect to  $P_{\Theta|Z=z}^{(\alpha)}$ . Hence, from Lemma 2.10, it holds that

$$\begin{aligned} K_z^{(1)} \left( -\frac{1}{\alpha} \right) &= \int_{\mathcal{R}_0(\alpha)} \mathbb{L}_z(\boldsymbol{\nu}) dP_{\Theta|Z=z}^{(\alpha)}(\boldsymbol{\nu}) + \int_{\mathcal{R}_1(\alpha)} \mathbb{L}_z(\boldsymbol{\nu}) dP_{\Theta|Z=z}^{(\alpha)}(\boldsymbol{\nu}) \\ &\quad + \int_{\mathcal{R}_2(\alpha)} \mathbb{L}_z(\boldsymbol{\nu}) dP_{\Theta|Z=z}^{(\alpha)}(\boldsymbol{\nu}) \end{aligned} \quad (266)$$

$$= \int_{\mathcal{R}_0(\alpha)} \mathbb{L}_z(\boldsymbol{\nu}) dP_{\Theta|Z=z}^{(\alpha)}(\boldsymbol{\nu}) + \int_{\mathcal{R}_1(\alpha)} \mathbb{L}_z(\boldsymbol{\nu}) dP_{\Theta|Z=z}^{(\alpha)}(\boldsymbol{\nu}) \quad (267)$$

$$= K_z^{(1)} \left( -\frac{1}{\alpha} \right) P_{\Theta|Z=z}^{(\alpha)}(\mathcal{R}_0(\alpha)) + \int_{\mathcal{R}_1(\alpha)} \mathbb{L}_z(\boldsymbol{\nu}) dP_{\Theta|Z=z}^{(\alpha)}(\boldsymbol{\nu}) \quad (268)$$

$$< K_z^{(1)} \left( -\frac{1}{\alpha} \right) P_{\Theta|Z=z}^{(\alpha)}(\mathcal{R}_0(\alpha)) + K_z^{(1)} \left( -\frac{1}{\alpha} \right) P_{\Theta|Z=z}^{(\alpha)}(\mathcal{R}_1(\alpha)) \quad (269)$$

$$= K_z^{(1)} \left( -\frac{1}{\alpha} \right) \left( P_{\Theta|Z=z}^{(\alpha)}(\mathcal{R}_0(\alpha)) + P_{\Theta|Z=z}^{(\alpha)}(\mathcal{R}_1(\alpha)) \right) \quad (270)$$

$$= K_z^{(1)} \left( -\frac{1}{\alpha} \right), \quad (271)$$

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  $P_{\Theta|Z=z}^{(\alpha)}$ . Hence, from Lemma 2.10, it holds that

$$K_z^{(1)}\left(-\frac{1}{\alpha}\right) = \int_{\mathcal{R}_0(\alpha)} \mathbf{L}_z(\boldsymbol{\nu}) dP_{\Theta|Z=z}^{(\alpha)}(\boldsymbol{\nu}) + \int_{\mathcal{R}_1(\alpha)} \mathbf{L}_z(\boldsymbol{\nu}) dP_{\Theta|Z=z}^{(\alpha)}(\boldsymbol{\nu}) + \int_{\mathcal{R}_2(\alpha)} \mathbf{L}_z(\boldsymbol{\nu}) dP_{\Theta|Z=z}^{(\alpha)}(\boldsymbol{\nu}) \quad (272)$$

$$= \int_{\mathcal{R}_0(\alpha)} \mathbf{L}_z(\boldsymbol{\nu}) dP_{\Theta|Z=z}^{(\alpha)}(\boldsymbol{\nu}) + \int_{\mathcal{R}_2(\alpha)} \mathbf{L}_z(\boldsymbol{\nu}) dP_{\Theta|Z=z}^{(\alpha)}(\boldsymbol{\nu}) \quad (273)$$

$$= K_z^{(1)}\left(-\frac{1}{\alpha}\right) P_{\Theta|Z=z}^{(\alpha)}(\mathcal{R}_0(\alpha)) + \int_{\mathcal{R}_2(\alpha)} \mathbf{L}_z(\boldsymbol{\nu}) dP_{\Theta|Z=z}^{(\alpha)}(\boldsymbol{\nu}) \quad (274)$$

$$> K_z^{(1)}\left(-\frac{1}{\alpha}\right) P_{\Theta|Z=z}^{(\alpha)}(\mathcal{R}_0(\alpha)) + K_z^{(1)}\left(-\frac{1}{\alpha}\right) P_{\Theta|Z=z}^{(\alpha)}(\mathcal{R}_2(\alpha)) \quad (275)$$

$$= K_z^{(1)}\left(-\frac{1}{\alpha}\right) \left( P_{\Theta|Z=z}^{(\alpha)}(\mathcal{R}_0(\alpha)) + P_{\Theta|Z=z}^{(\alpha)}(\mathcal{R}_2(\alpha)) \right) \quad (276)$$

$$= K_z^{(1)}\left(-\frac{1}{\alpha}\right), \quad (277)$$

which is also a contradiction. This completes the proof.

## P Proof of Theorem 2.4

From Theorem 2.2, it follows that for all  $(\lambda_1, \lambda_2) \in \mathcal{K}_z \times \mathcal{K}_z$  with  $\lambda_1 > \lambda_2$ ,

$$\int \mathbf{L}_z(\boldsymbol{\alpha}) \frac{dP_{\Theta|Z=z}^{(\lambda_1)}}{dP}(\boldsymbol{\alpha}) dP(\boldsymbol{\alpha}) \geq \int \mathbf{L}_z(\boldsymbol{\alpha}) \frac{dP_{\Theta|Z=z}^{(\lambda_2)}}{dP}(\boldsymbol{\alpha}) dP(\boldsymbol{\alpha}), \quad (278)$$

which implies the following inclusions:

$$\mathcal{R}_1(\lambda_2) \subseteq \mathcal{R}_1(\lambda_1), \text{ and} \quad (279a)$$

$$\mathcal{R}_2(\lambda_1) \subseteq \mathcal{R}_2(\lambda_2), \quad (279b)$$

with the sets  $\mathcal{R}_1(\cdot)$  and  $\mathcal{R}_2(\cdot)$  in (76). From (81), it holds that for all  $i \in \{1, 2\}$ ,

$$\mathcal{N}(\lambda_i) = \mathcal{R}_2(\lambda_i)^c, \quad (280)$$

where the complement is with respect to  $\mathcal{M}$ . Thus, the inclusion in (279b) and the equality in (280) yields,

$$\mathcal{N}(\lambda_1) \supseteq \mathcal{N}(\lambda_2). \quad (281)$$

The inclusion  $\mathcal{M} \supseteq \mathcal{N}(\lambda_1)$  follows from (81). Alternatively, the inclusion  $\mathcal{N}(\lambda_2) \supseteq \mathcal{N}^*$ , follows from observing that for all  $\boldsymbol{\nu} \in \mathcal{N}^*$ ,

$$\int \mathbf{L}_z(\boldsymbol{\alpha}) \frac{dP_{\Theta|Z=z}^{(\lambda_2)}}{dP}(\boldsymbol{\alpha}) dP(\boldsymbol{\alpha}) \geq \delta^* \int \frac{dP_{\Theta|Z=z}^{(\lambda_2)}}{dP}(\boldsymbol{\alpha}) dP(\boldsymbol{\alpha}) \quad (282)$$

$$\geq \delta^* \quad (283)$$

$$\geq \mathbf{L}_z(\boldsymbol{\nu}), \quad (284)$$

which implies that  $\nu \in \mathcal{N}(\lambda_2)$ . This completes the proof of (85).

The proof of (86) is as follows. From the mean value theorem [3] and the assumption that the empirical risk function  $L_z$  in (5b) is continuous on  $\mathcal{M}$ , it follows that for all  $\lambda \in \mathcal{K}_z$ , there always exists a model  $\theta \in \mathcal{M}$ , such that

$$L_z(\theta) = \int L_z(\alpha) dP_{\Theta|Z=z}^{(\lambda)}(\alpha), \quad (285)$$

which implies that  $\mathcal{R}_0(\lambda)$  is not empty, and as a consequence,  $\mathcal{N}(\lambda) = \mathcal{R}_0(\lambda) \cup \mathcal{R}_1(\lambda)$  is not empty. Hence, for all  $\theta \in \mathcal{R}_0(\lambda)$  it holds that  $\theta \notin \mathcal{N}(\lambda_2)$ . This proves that the elements of  $\mathcal{R}_0(\lambda)$  are in  $\mathcal{N}(\lambda)$  but not in  $\mathcal{N}(\lambda_2)$ . This, together with (281), verifies that

$$\mathcal{N}(\lambda) \supset \mathcal{N}(\lambda_2). \quad (286)$$

The strict inclusion  $\mathcal{M} \supset \mathcal{N}(\lambda)$  is proved by contradiction. Assume that there exists a  $\lambda \in \mathcal{K}_z$  such that  $\mathcal{M} = \mathcal{N}(\lambda)$ . Then,  $\mathcal{R}_2(\lambda) = \emptyset$  and thus,  $P_{\Theta|Z=z}^{(\lambda)}(\mathcal{R}_2(\lambda)) = 0$ , which together with Lemma 2.13, implies that  $P_{\Theta|Z=z}^{(\lambda)}(\mathcal{R}_1(\lambda)) = 0$  and consequently,

$$P_{\Theta|Z=z}^{(\lambda)}(\mathcal{R}_0(\lambda)) = 1. \quad (287)$$

This contradicts the assumption that the function  $L_z$  is separable (Definition 2.2). Hence,  $\mathcal{M} \supset \mathcal{N}(\lambda)$ .

Finally, the strict inclusion  $\mathcal{N}(\lambda_2) \supset \mathcal{N}^*$  is proved by contradiction. Assume that there exists a  $\lambda \in \mathcal{K}_z$  such that  $\mathcal{N}^* = \mathcal{N}(\lambda)$ . Hence, three cases might arise:

- (a) there exists a  $\lambda \in \mathcal{K}_z$ , such that  $\delta^* < K_z^{(1)}(-\frac{1}{\lambda})$  and the set  $\{\nu \in \mathcal{M} : \delta^* < L_z(\nu) \leq K_z^{(1)}(-\frac{1}{\lambda})\} = \emptyset$ ;
- (b) there exists a  $\lambda \in \mathcal{K}_z$ , such that  $\delta^* > K_z^{(1)}(-\frac{1}{\lambda})$  and the set  $\{\nu \in \mathcal{M} : K_z^{(1)}(-\frac{1}{\lambda}) < L_z(\nu) \leq \delta^*\} = \emptyset$ ; or
- (c) there exists a  $\lambda \in \mathcal{K}_z$ , such that  $\delta^* = K_z^{(1)}(-\frac{1}{\lambda})$ .

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,

$$\mathcal{R}_1(\lambda) = \{\nu \in \mathcal{M} : L_z(\nu) < \delta^*\}, \quad (288)$$

which implies that

$$P_{\Theta|Z=z}^{(\lambda)}(\mathcal{R}_1(\lambda)) = 0. \quad (289)$$

From Lemma 2.13 and (289), it follows that,

$$P_{\Theta|Z=z}^{(\lambda)}(\mathcal{R}_2(\lambda)) = 0. \quad (290)$$

Finally, by noticing that

$$1 = P_{\Theta|Z=z}^{(\lambda)}(\mathcal{R}_0(\lambda)) + P_{\Theta|Z=z}^{(\lambda)}(\mathcal{R}_1(\lambda)) + P_{\Theta|Z=z}^{(\lambda)}(\mathcal{R}_2(\lambda)) \quad (291)$$

$$= P_{\Theta|Z=z}^{(\lambda)}(\mathcal{R}_0(\lambda)), \quad (292)$$

leads to the conclusion that the assumption that there exists a  $\lambda \in \mathcal{K}_z$  such that  $\mathcal{N}^* = \mathcal{N}(\lambda)$  is a contradiction to the assumption that the function  $L_z$  is separable (Definition 2.2). Hence,  $\mathcal{N}(\lambda_2) \supset \mathcal{T}(z)$ , which completes the proof of (86).

## Q Proof of Theorem 2.5

The proof of (87) is based on the analysis of the derivative of  $P_{\Theta|Z=z}^{(\lambda)}(\mathcal{A})$  with respect to  $\lambda$ , for some fixed set  $\mathcal{A} \subseteq \mathcal{M}$ . More specifically, given a  $\gamma \in \mathcal{K}_z$ , it holds that

$$P_{\Theta|Z=z}^{(\gamma)}(\mathcal{A}) = \int_{\mathcal{A}} \frac{dP_{\Theta|Z=z}^{(\gamma)}(\alpha)}{dP}(\alpha) dP(\alpha), \quad (293)$$

and from the fundamental theorem of calculus [38], it follows that for all  $(\lambda_1, \lambda_2) \in \mathcal{K}_z \times \mathcal{K}_z$  with  $\lambda_1 > \lambda_2$ ,

$$P_{\Theta|Z=z}^{(\lambda_1)}(\mathcal{A}) - P_{\Theta|Z=z}^{(\lambda_2)}(\mathcal{A}) = \int_{\lambda_2}^{\lambda_1} \frac{d}{d\gamma} P_{\Theta|Z=z}^{(\gamma)}(\mathcal{A}) d\gamma \quad (294)$$

$$= \int_{\lambda_2}^{\lambda_1} \frac{d}{d\gamma} \int_{\mathcal{A}} \frac{dP_{\Theta|Z=z}^{(\gamma)}(\alpha)}{dP}(\alpha) dP(\alpha) d\gamma \quad (295)$$

$$= \int_{\lambda_2}^{\lambda_1} \int_{\mathcal{A}} \frac{d}{d\gamma} \frac{dP_{\Theta|Z=z}^{(\gamma)}(\alpha)}{dP}(\alpha) dP(\alpha) d\gamma, \quad (296)$$

where the equality in (295) follows from (293); and the equality in (296) holds from Theorem 2.1 and the dominated convergence theorem [3].

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

$$\frac{d}{d\lambda} \frac{dP_{\Theta|Z=z}^{(\lambda)}(\theta)}{dP}(\theta) = \frac{d}{d\lambda} \frac{\exp\left(-\frac{L_z(\theta)}{\lambda}\right)}{\int \exp\left(-\frac{L_z(\nu)}{\lambda}\right) dP(\nu)} \quad (297)$$

$$\begin{aligned} &= \frac{\frac{1}{\lambda^2} L_z(\theta) \exp\left(-\frac{L_z(\theta)}{\lambda}\right)}{\int \exp\left(-\frac{L_z(\nu)}{\lambda}\right) dP(\nu)} \\ &= \frac{\frac{1}{\lambda^2} \exp\left(-\frac{L_z(\theta)}{\lambda}\right) \int L_z(\alpha) \exp\left(-\frac{L_z(\alpha)}{\lambda}\right) dP(\alpha)}{\left(\int \exp\left(-\frac{L_z(\nu)}{\lambda}\right) dP(\nu)\right)^2} \quad (298) \end{aligned}$$

$$\begin{aligned} &= \frac{1}{\lambda^2} L_z(\theta) \frac{dP_{\Theta|Z=z}^{(\lambda)}(\theta)}{dP}(\theta) \\ &\quad - \frac{1}{\lambda^2} \frac{dP_{\Theta|Z=z}^{(\lambda)}(\theta)}{dP}(\theta) \int L_z(\nu) \frac{dP_{\Theta|Z=z}^{(\lambda)}(\nu)}{dP}(\nu) dP(\nu) \quad (299) \end{aligned}$$

$$= \frac{1}{\lambda^2} \frac{dP_{\Theta|Z=z}^{(\lambda)}(\theta)}{dP}(\theta) \left( L_z(\theta) - \int L_z(\nu) dP_{\Theta|Z=z}^{(\lambda)}(\nu) \right). \quad (300)$$

Plugging (300) into (296) yields,

$$\begin{aligned} &P_{\Theta|Z=z}^{(\lambda_1)}(\mathcal{A}) - P_{\Theta|Z=z}^{(\lambda_2)}(\mathcal{A}) \\ &= \int_{\lambda_2}^{\lambda_1} \int_{\mathcal{A}} \frac{1}{\gamma^2} \frac{dP_{\Theta|Z=z}^{(\gamma)}(\alpha)}{dP}(\alpha) \left( L_z(\alpha) - \int L_z(\nu) dP_{\Theta|Z=z}^{(\gamma)}(\nu) \right) dP(\alpha) d\gamma \quad (301) \end{aligned}$$

$$= \int_{\lambda_2}^{\lambda_1} \int_{\mathcal{A}} \frac{1}{\gamma^2} \left( L_z(\alpha) - \int L_z(\nu) dP_{\Theta|Z=z}^{(\gamma)}(\nu) \right) dP_{\Theta|Z=z}^{(\gamma)}(\alpha) d\gamma. \quad (302)$$



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

$$\mathbf{L}_z(\alpha) - \int \mathbf{L}_z(\nu) dP_{\Theta|Z=z}^{(\gamma)}(\nu) \leq 0, \quad (303)$$

and thus,

$$\int_{\mathcal{N}(\lambda_2)} \frac{1}{\gamma^2} \left( \mathbf{L}_z(\alpha) - \int \mathbf{L}_z(\nu) dP_{\Theta|Z=z}^{(\gamma)}(\nu) \right) dP_{\Theta|Z=z}^{(\gamma)}(\alpha) \leq 0. \quad (304)$$

The equalities in (302) and (304), with  $\mathcal{A} = \mathcal{N}(\lambda)$ , imply that

$$P_{\Theta|Z=z}^{(\lambda_1)}(\mathcal{N}(\lambda_2)) - P_{\Theta|Z=z}^{(\lambda_2)}(\mathcal{N}(\lambda_2)) \leq 0. \quad (305)$$

The inequality  $0 < P_{\Theta|Z=z}^{(\lambda_1)}(\mathcal{N}(\lambda_2))$  in (87) is proved by contradiction. Assume that for some  $\lambda \in \mathcal{K}_z$  it holds that  $0 = P_{\Theta|Z=z}^{(\lambda)}(\mathcal{N}(\lambda_2))$ . Then,  $P_{\Theta|Z=z}^{(\lambda)}(\mathcal{R}_0(\lambda_2)) + P_{\Theta|Z=z}^{(\lambda)}(\mathcal{R}_1(\lambda_2)) = 0$ , which implies that  $P_{\Theta|Z=z}^{(\lambda)}(\mathcal{R}_2(\lambda_2)) = 1$ , which is a contradiction. See for instance, Lemma 2.14. This completes the proof of (87).

The proof of (88) is divided into two parts. The first part shows that if for all pairs  $(\lambda_1, \lambda_2) \in \mathcal{K}_z \times \mathcal{K}_z$  with  $\lambda_1 > \lambda_2$ ,

$$P_{\Theta|Z=z}^{(\lambda_1)}(\mathcal{N}(\lambda_2)) < P_{\Theta|Z=z}^{(\lambda_2)}(\mathcal{N}(\lambda_2)), \quad (306)$$

then the function  $\mathbf{L}_z$  is separable. The second part of the proof shows that if the function  $\mathbf{L}_z$  is separable, then, for all pairs  $(\lambda_1, \lambda_2) \in \mathcal{K}_z \times \mathcal{K}_z$  with  $\lambda_1 > \lambda_2$ , the inequality in (306) holds.

The first part is as follows. In the proof of Theorem 2.4 it is shown (see (302)) that for all pairs  $(\lambda_1, \lambda_2) \in \mathcal{K}_z \times \mathcal{K}_z$  with  $\lambda_1 > \lambda_2$ ,

$$\begin{aligned} & P_{\Theta|Z=z}^{(\lambda_1)}(\mathcal{N}(\lambda_2)) - P_{\Theta|Z=z}^{(\lambda_2)}(\mathcal{N}(\lambda_2)) \\ &= \int_{\lambda_2}^{\lambda_1} \int_{\mathcal{N}(\lambda_2)} \frac{1}{\gamma^2} \left( \mathbf{L}_z(\alpha) - \int \mathbf{L}_z(\nu) dP_{\Theta|Z=z}^{(\gamma)}(\nu) \right) dP_{\Theta|Z=z}^{(\gamma)}(\alpha) d\gamma. \end{aligned} \quad (307)$$

Assume that for a given pair  $(\lambda_1, \lambda_2) \in \mathcal{K}_z \times \mathcal{K}_z$ , with  $\lambda_1 > \lambda_2$ , the inequality in (306) holds. Then, from (307),

$$\begin{aligned} 0 &> \int_{\lambda_2}^{\lambda_1} \int_{\mathcal{N}(\lambda_2)} \frac{1}{\gamma^2} \left( \mathbf{L}_z(\alpha) - \int \mathbf{L}_z(\nu) dP_{\Theta|Z=z}^{(\gamma)}(\nu) \right) dP_{\Theta|Z=z}^{(\gamma)}(\alpha) d\gamma \\ &= \int_{\lambda_2}^{\lambda_1} \int_{\mathcal{R}_1(\lambda_2)} \frac{1}{\gamma^2} \left( \mathbf{L}_z(\alpha) - \int \mathbf{L}_z(\nu) dP_{\Theta|Z=z}^{(\gamma)}(\nu) \right) dP_{\Theta|Z=z}^{(\gamma)}(\alpha) d\gamma \end{aligned} \quad (308)$$

where the equality in (308) follows from noticing that  $\mathcal{R}_0(\lambda_2)$  and  $\mathcal{R}_1(\lambda_2)$  form a partition of  $\mathcal{N}(\lambda_2)$ , with the sets  $\mathcal{R}_0(\lambda_2)$ ,  $\mathcal{R}_1(\lambda_2)$  and  $\mathcal{N}(\lambda_2)$  defined in (76a), (76b), and (81), respectively.

The inequality in (308) implies that the set  $\mathcal{R}_1(\lambda_2)$  is nonnegligible with respect to  $P_{\Theta|Z=z}^{(\gamma)}$ , for some  $\gamma \in (\lambda_2, \lambda_1)$ . Hence, from Lemma 2.14, it follows that both sets  $\mathcal{R}_1(\lambda_2)$  and  $\mathcal{R}_2(\lambda_2)$  are nonnegligible with respect to  $P_{\Theta|Z=z}^{(\gamma)}$ .

From the above arguments, it has been proved that given a pair  $(\lambda_1, \lambda_2) \in \mathcal{K}_z \times \mathcal{K}_z$  with  $\lambda_1 > \lambda_2$ , if

$$P_{\Theta|Z=z}^{(\lambda_1)}(\mathcal{N}(\lambda_2)) < P_{\Theta|Z=z}^{(\lambda_2)}(\mathcal{N}(\lambda_2)), \quad (309)$$

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  $P_{\Theta|Z=z}^{(\gamma)}$ . Moreover, such sets  $\mathcal{R}_1(\lambda_2)$  and  $\mathcal{R}_2(\lambda_2)$  satisfy for all  $(\nu_1, \nu_2) \in \mathcal{R}_2(\lambda) \times \mathcal{R}_1(\lambda)$ ,

$$\mathbb{L}_z(\nu_1) > K_z^{(1)} \left( -\frac{1}{\lambda} \right) > \mathbb{L}_z(\nu_2), \quad (310)$$

which together with Definition 2.3 verify that the function  $\mathbb{L}_z$  is separable. This ends the first part of the proof.

The second part of the proof is under the assumption that the empirical risk function  $\mathbb{L}_z$  in (5b) is separable. That is, from Definition 2.3, for all  $\gamma \in \mathcal{K}_z$ , 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  $P_{\Theta|Z=z}^{(\gamma)}$  in (37) that verify that for all  $(\nu_1, \nu_2) \in \mathcal{A}(\gamma) \times \mathcal{B}(\gamma)$ ,

$$\mathbb{L}_z(\nu_1) > c_\gamma > \mathbb{L}_z(\nu_2). \quad (311)$$

In the proof of Theorem 2.4, c.f. (302), it has been proved that given a pair  $(\alpha_1, \alpha_2) \in \mathcal{K}_z \times \mathcal{K}_z$ , with  $\alpha_1 > \gamma > \alpha_2$ , it holds that for all subsets  $\mathcal{A}$  of  $\mathcal{M}$ ,

$$\begin{aligned} & P_{\Theta|Z=z}^{(\alpha_1)}(\mathcal{A}) - P_{\Theta|Z=z}^{(\alpha_2)}(\mathcal{A}) \\ &= \int_{\alpha_2}^{\alpha_1} \int_{\mathcal{A}} \frac{1}{\lambda^2} \frac{dP_{\Theta|Z=z}^{(\lambda)}}{dP}(\alpha) \left( \mathbb{L}_z(\alpha) - \int \mathbb{L}_z(\nu) dP_{\Theta|Z=z}^{(\lambda)}(\nu) \right) dP(\alpha) d\lambda \end{aligned} \quad (312)$$

$$= \int_{\alpha_2}^{\alpha_1} \int_{\mathcal{A}} \frac{1}{\lambda^2} \left( \mathbb{L}_z(\alpha) - \int \mathbb{L}_z(\nu) dP_{\Theta|Z=z}^{(\lambda)}(\nu) \right) dP_{\Theta|Z=z}^{(\lambda)}(\alpha) d\lambda. \quad (313)$$

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

$$c_\gamma < K_z^{(1)} \left( -\frac{1}{\gamma} \right), \quad (314)$$

with  $c_\gamma$  in (311). The second case considers that

$$c_\gamma \geq K_z^{(1)} \left( -\frac{1}{\gamma} \right). \quad (315)$$

In the first case, it follows from (81) that

$$\mathcal{B}(\gamma) \subset \mathcal{N}(\gamma), \quad (316)$$

which implies that

$$P_{\Theta|Z=z}^{(\gamma)}(\mathcal{N}(\gamma)) \geq P_{\Theta|Z=z}^{(\gamma)}(\mathcal{B}(\gamma)) \quad (317)$$

$$> 0, \quad (318)$$

where, the inequality in (318) follows from the fact that  $\mathcal{B}(\gamma)$  is nonnegligible with respect to  $P_{\Theta|Z=z}^{(\gamma)}$ . This implies that the set  $\mathcal{N}_1(\gamma)$  is not negligible with respect to  $P_{\Theta|Z=z}^{(\gamma)}$ . Moreover, from (81) and (316), it follows that for all  $\alpha \in \mathcal{N}(\gamma)$  and for all  $\lambda \in (\gamma, \alpha_1)$ ,

$$\mathbb{L}_z(\alpha) - \int \mathbb{L}_z(\nu) dP_{\Theta|Z=z}^{(\lambda)}(\nu) < \mathbb{L}_z(\alpha) - c_\gamma \quad (319)$$

$$< 0, \quad (320)$$

where the inequality in (319) follows from (314); and the inequality in (320) follows from (311). Thus,

$$\int_\gamma^{\alpha_1} \int_{\mathcal{N}(\gamma)} \frac{1}{\lambda^2} \left( \mathbb{L}_z(\alpha) - \int \mathbb{L}_z(\nu) dP_{\Theta|Z=z}^{(\lambda)}(\nu) \right) dP_{\Theta|Z=z}^{(\lambda)}(\alpha) d\lambda < 0, \quad (321)$$

which implies, from (313), that

$$P_{\Theta|Z=z}^{(\alpha_1)}(\mathcal{N}(\gamma)) - P_{\Theta|Z=z}^{(\gamma)}(\mathcal{N}(\gamma)) < 0. \quad (322)$$

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

$$A(\gamma) \subseteq \mathcal{R}_2(\gamma), \quad (323)$$

which implies that

$$P_{\Theta|Z=z}^{(\gamma)}(\mathcal{R}_2(\gamma)) \geq P_{\Theta|Z=z}^{(\gamma)}(A(\gamma)) \quad (324)$$

$$> 0, \quad (325)$$

where the inequality in (325) follows from the fact that  $A(\gamma)$  is nonnegligible with respect to  $P_{\Theta|Z=z}^{(\gamma)}$ . This implies that the set  $\mathcal{R}_2(\gamma)$  is not negligible with respect to  $P_{\Theta|Z=z}^{(\gamma)}$ . From Lemma 2.13, it follows that both  $\mathcal{R}_1(\gamma)$  and  $\mathcal{R}_2(\gamma)$  are nonnegligible with respect to  $P_{\Theta|Z=z}^{(\gamma)}$ . Using this result, the following holds,

$$P_{\Theta|Z=z}^{(\gamma)}(\mathcal{N}(\gamma)) \geq P_{\Theta|Z=z}^{(\gamma)}(\mathcal{R}_1(\gamma)) \quad (326)$$

$$> 0, \quad (327)$$

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

From (81) and Theorem 2.2, it follows that for all  $\alpha \in \mathcal{N}(\gamma)$  and for all  $\lambda \in (\gamma, \alpha_1)$ ,

$$0 \geq \mathbb{L}_z(\alpha) - \int \mathbb{L}_z(\nu) dP_{\Theta|\underline{X}=z=y}^{(\gamma)}(\nu) \quad (328)$$

$$> \mathbb{L}_z(\alpha) - \int \mathbb{L}_z(\nu) dP_{\Theta|\underline{X}=z=y}^{(\lambda)}(\nu). \quad (329)$$

Thus,

$$\int_\gamma^{\alpha_1} \int_{\mathcal{N}(\gamma)} \frac{1}{\lambda^2} \left( \mathbb{L}_z(\alpha) - \int \mathbb{L}_z(\nu) dP_{\Theta|Z=z}^{(\lambda)}(\nu) \right) dP_{\Theta|Z=z}^{(\lambda)}(\alpha) d\lambda < 0, \quad (330)$$

which implies, from (313), that

$$P_{\Theta|Z=z}^{(\alpha_1)}(\mathcal{N}(\gamma)) - P_{\Theta|Z=z}^{(\gamma)}(\mathcal{N}(\gamma)) < 0. \quad (331)$$

The inequality  $0 < P_{\Theta|Z=z}^{(\lambda_1)}(\mathcal{N}(\lambda_2))$  in (88) has already been proved while proving (87), and thus, this completes the proof of (88).

## R Proof of Lemma 2.15

The proof is based on the observation that

$$\mathcal{N}(\lambda_1) = \mathcal{N}(\lambda_2) \cup (\mathcal{N}(\lambda_1) \cap \mathcal{R}_2(\lambda_2)), \quad (332)$$

and the fact that  $\mathcal{N}(\lambda_2)$  and  $(\mathcal{N}(\lambda_1) \cap \mathcal{R}_2(\lambda_2))$  are disjoint. Hence, for all  $i \in \{1, 2\}$ ,

$$P_{\Theta|Z=z}^{(\lambda_i)}(\mathcal{N}(\lambda_1)) = P_{\Theta|Z=z}^{(\lambda_i)}\left(\mathcal{N}(\lambda_2) \cup (\mathcal{N}(\lambda_1) \cap \mathcal{R}_2(\lambda_2))\right) \quad (333)$$

$$= P_{\Theta|Z=z}^{(\lambda_i)}(\mathcal{N}(\lambda_2)) + P_{\Theta|Z=z}^{(\lambda_i)}(\mathcal{N}(\lambda_1) \cap \mathcal{R}_2(\lambda_2)) \quad (334)$$

$$= P_{\Theta|Z=z}^{(\lambda_i)}(\mathcal{N}(\lambda_2)), \quad (335)$$

where the equality in (334) follows from Lemma 2.6 and the equality in (89).

Note that the sets  $\mathcal{N}(\lambda_2)$ ,  $\mathcal{N}(\lambda_1) \cap \mathcal{R}_2(\lambda_2)$  and  $\mathcal{R}_1(\lambda_1)$  form a partition of  $\mathcal{M}$  and thus, the equality in (89) implies that for all  $i \in \{1, 2\}$ ,

$$0 < P_{\Theta|Z=z}^{(\lambda_i)}(\mathcal{N}(\lambda_2)) = 1 - P_{\Theta|Z=z}^{(\lambda_i)}(\mathcal{R}_1(\lambda_1)) \quad (336)$$

$$= P_{\Theta|Z=z}^{(\lambda_i)}(\mathcal{R}_0(\lambda_1)) + P_{\Theta|Z=z}^{(\lambda_i)}(\mathcal{R}_2(\lambda_1)), \quad (337)$$

where the inequality in (336) holds from Theorem 2.5. Finally, under the assumption that the empirical function  $L_z$  in (5b) is separable, it holds from Theorem 2.5 that

$$P_{\Theta|Z=z}^{(\lambda_1)}(\mathcal{N}(\lambda_2)) < P_{\Theta|Z=z}^{(\lambda_2)}(\mathcal{N}(\lambda_2)). \quad (338)$$

Plugging (335) into (338), with  $i = 1$ , yields,

$$P_{\Theta|Z=z}^{(\lambda_1)}(\mathcal{N}(\lambda_1)) < P_{\Theta|Z=z}^{(\lambda_2)}(\mathcal{N}(\lambda_2)), \quad (339)$$

and this completes the proof.

## S Proof of Theorem 2.6

$$\lim_{\lambda \rightarrow 0^+} P_{\Theta|Z=z}^{(\lambda)}(\mathcal{N}(\lambda)) = \lim_{\lambda \rightarrow 0^+} P_{\Theta|Z=z}^{(\lambda)}(\mathcal{L}_z^*) + \lim_{\lambda \rightarrow 0^+} P_{\Theta|Z=z}^{(\lambda)}(\mathcal{N}(\lambda)) \quad (340)$$

$$= \lim_{\lambda \rightarrow 0^+} P_{\Theta|Z=z}^{(\lambda)}(\mathcal{L}_z^*) \quad (341)$$

$$= \begin{cases} 0 & \text{if } \mathcal{L}_z^* = \emptyset \\ 1 & \text{if } \mathcal{L}_z^* \neq \emptyset, \end{cases} \quad (342)$$

where, the equalities in (341) follows Lemma 2.12. This completes the proof.

## T Proof of Lemma 2.16

The cumulant generating function  $J_{\mathbf{z},\lambda}$  in (98) induced by the measure  $P_{W|\mathbf{Z}=\mathbf{z}}^{(\lambda)}$  in (96) evaluated at  $t$ , with  $t \leq \frac{1}{\lambda}$ , is

$$J_{\mathbf{z},\lambda}(t) = \log \left( \int \exp(t \mathbf{L}_{\mathbf{z}}(\mathbf{u})) \frac{dP_{\Theta|\mathbf{Z}=\mathbf{z}}^{(\lambda)}(\mathbf{u})}{dP}(\mathbf{u}) dP(\mathbf{u}) \right) \quad (343)$$

$$= \log \left( \int \exp(t \mathbf{L}_{\mathbf{z}}(\mathbf{u})) \exp \left( -K_{\mathbf{z}} \left( -\frac{1}{\lambda} \right) - \frac{1}{\lambda} \mathbf{L}_{\mathbf{z}}(\mathbf{u}) \right) dP(\mathbf{u}) \right) \quad (344)$$

$$= \log \left( \int \exp \left( \left( t - \frac{1}{\lambda} \right) \mathbf{L}_{\mathbf{z}}(\mathbf{u}) - K_{\mathbf{z}} \left( -\frac{1}{\lambda} \right) \right) dP(\mathbf{u}) \right) \quad (345)$$

$$= \log \left( \int \exp \left( \left( t - \frac{1}{\lambda} \right) \mathbf{L}_{\mathbf{z}}(\mathbf{u}) - K_{\mathbf{z}} \left( t - \frac{1}{\lambda} \right) + K_{\mathbf{z}} \left( t - \frac{1}{\lambda} \right) - K_{\mathbf{z}} \left( -\frac{1}{\lambda} \right) \right) dP(\mathbf{u}) \right) \quad (346)$$

$$= K_{\mathbf{z}} \left( t - \frac{1}{\lambda} \right) - K_{\mathbf{z}} \left( -\frac{1}{\lambda} \right) + \log \left( \int \exp \left( \left( t - \frac{1}{\lambda} \right) \mathbf{L}_{\mathbf{z}}(\mathbf{u}) - K_{\mathbf{z}} \left( t - \frac{1}{\lambda} \right) \right) dP(\mathbf{u}) \right) \quad (347)$$

$$= K_{\mathbf{z}} \left( t - \frac{1}{\lambda} \right) - K_{\mathbf{z}} \left( -\frac{1}{\lambda} \right) + \log \left( \int \frac{dP_{\Theta|\mathbf{Z}=\mathbf{z}} \left( -\frac{1}{t-\frac{1}{\lambda}} \right)}{dP}(\mathbf{u}) dP(\mathbf{u}) \right) \quad (348)$$

$$= K_{\mathbf{z}} \left( t - \frac{1}{\lambda} \right) - K_{\mathbf{z}} \left( -\frac{1}{\lambda} \right), \quad (349)$$

where the equality in (344) follows from Theorem 2.1; and the equality in (348) follows from the fact that  $-\frac{1}{t-\frac{1}{\lambda}} > 0$  for all  $t < \frac{1}{\lambda}$ . This completes the proof.

## U Proof of Lemma 3.1

From [9, 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:

$$D(Q\|P) = \sup_f \int f(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) - \log \int \exp(f(\boldsymbol{\theta})) dP(\boldsymbol{\theta}), \quad (350)$$

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  $\mathbf{L}_{\mathbf{u}}$  in (5b) satisfies that

$$D(Q\|P) \geq \int t \mathbf{L}_{\mathbf{u}}(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) - \log \int \exp(t \mathbf{L}_{\mathbf{u}}(\boldsymbol{\theta})) dP(\boldsymbol{\theta}) \quad (351)$$

$$\geq \int t \mathbf{L}_{\mathbf{u}}(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) - \log \int \exp(t \mathbf{L}_{\mathbf{u}}(\boldsymbol{\theta}) + t\mu - t\mu) dP(\boldsymbol{\theta}). \quad (352)$$

$$= \int t \mathbf{L}_{\mathbf{u}}(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) - \int t \mathbf{L}_{\mathbf{u}}(\boldsymbol{\theta}) dP(\boldsymbol{\theta}) - \log \int \exp(t \mathbf{L}_{\mathbf{u}}(\boldsymbol{\theta}) - t\mu) dP(\boldsymbol{\theta}), \quad (353)$$

which leads to

$$\int \mathbb{L}_{\mathbf{u}}(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) - \int \mathbb{L}_{\mathbf{u}}(\boldsymbol{\theta}) dP(\boldsymbol{\theta}) \leq \frac{D(Q\|P) + \log \int \exp(t(\mathbb{L}_{\mathbf{u}}(\boldsymbol{\theta}) - \mu)) dP(\boldsymbol{\theta})}{t}. \quad (354)$$

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

$$\int \mathbb{L}_{\mathbf{u}}(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) - \int \mathbb{L}_{\mathbf{u}}(\boldsymbol{\theta}) dP(\boldsymbol{\theta}) \leq \inf_{t \in (-\infty, 0)} \frac{D(Q\|P) + \log \int \exp(t(\mathbb{L}_{\mathbf{u}}(\boldsymbol{\theta}) - \mu)) dP(\boldsymbol{\theta})}{t}, \quad (355)$$

which completes the proof.

## V Proof of Theorem 3.3

From Lemma 3.1, it holds that the probability measure  $P_{\Theta|Z=z}^{(\lambda)}$  in (37), satisfies

$$\begin{aligned} & \int \mathbb{L}_{\mathbf{z}}(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) - \int \mathbb{L}_{\mathbf{z}}(\boldsymbol{\theta}) dP_{\Theta|Z=z}^{(\lambda)}(\boldsymbol{\theta}) \\ & \leq \inf_{t \in (-\infty, 0)} \left( \frac{D(Q\|P_{\Theta|Z=z}^{(\lambda)}) + \log \left( \int \exp \left( t \left( \mathbb{L}_{\mathbf{z}}(\boldsymbol{\theta}) - K_{\mathbf{z}}^{(1)} \left( -\frac{1}{\lambda} \right) \right) \right) dP_{\Theta|Z=z}^{(\lambda)}(\boldsymbol{\theta}) \right)}{t} \right), \end{aligned} \quad (356)$$

where the function  $K_{\mathbf{z}}^{(1)}$  is defined in (57). Moreover, for all  $t \in (-\infty, 0)$ ,

$$\begin{aligned} & \log \left( \int \exp \left( t \left( \mathbb{L}_{\mathbf{z}}(\boldsymbol{\theta}) - K_{\mathbf{z}}^{(1)} \left( -\frac{1}{\lambda} \right) \right) \right) dP_{\Theta|Z=z}^{(\lambda)}(\boldsymbol{\theta}) \right) \\ & = \log \left( \int \exp(t \mathbb{L}_{\mathbf{z}}(\boldsymbol{\theta})) dP_{\Theta|Z=z}^{(\lambda)}(\boldsymbol{\theta}) \right) - t K_{\mathbf{z}}^{(1)} \left( -\frac{1}{\lambda} \right) \end{aligned} \quad (357)$$

$$= J_{\mathbf{z}, \lambda}(t) - t K_{\mathbf{z}}^{(1)} \left( -\frac{1}{\lambda} \right) \quad (358)$$

$$\leq \frac{1}{2} t^2 B_{\mathbf{z}}^2, \quad (359)$$

where the inequality in (358) follows from (98); the inequality in (359) follows from Theorem 2.7; and the constant  $B_{\mathbf{z}}$  is defined in (103).

Plugging (359) into (356) yields for all  $t \in (-\infty, 0)$ ,

$$\int \mathbb{L}_{\mathbf{z}}(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) - \int \mathbb{L}_{\mathbf{z}}(\boldsymbol{\theta}) dP_{\Theta|Z=z}^{(\lambda)}(\boldsymbol{\theta}) \leq \inf_{t \in (-\infty, 0)} \frac{D(Q\|P_{\Theta|Z=z}^{(\lambda)}) + \frac{1}{2} t^2 B_{\mathbf{z}}^2}{t}. \quad (360)$$

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

$$c \triangleq \int \mathbb{L}_{\mathbf{z}}(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) - \int \mathbb{L}_{\mathbf{z}}(\boldsymbol{\theta}) dP_{\Theta|Z=z}^{(\lambda)}(\boldsymbol{\theta}). \quad (361)$$

Hence, for all  $t \in (-\infty, 0)$ ,

$$c t - \frac{1}{2} t^2 B_{\mathbf{z}}^2 \leq D(Q\|P_{\Theta|Z=z}^{(\lambda)}). \quad (362)$$

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

$$\phi(\alpha) = \frac{1}{2} \alpha^2 B_{\mathbf{z}}^2, \quad (363)$$

and note that  $\phi$  is a positive and strictly convex function with  $\phi(0) = 0$ . Let the Legendre-Fenchel transform of  $\phi$  be the function  $\phi^* : \mathbb{R} \rightarrow \mathbb{R}$ , and thus for all  $x \in \mathbb{R}$ ,

$$\phi^*(x) = \max_{t \in (-\infty, 0)} xt - \phi(t). \quad (364)$$

In particular, note that

$$\phi^*(c) \leq D \left( Q \| P_{\Theta|Z=z}^{(\lambda)} \right). \quad (365)$$

Note that for all  $x \in \mathbb{R}$  and for all  $t \in (-\infty, 0)$ , the function  $\phi^*$  in (364) satisfies

$$xt - \frac{1}{2}t^2 B_z^2 \leq \phi^*(x) = x\alpha^*(x) - \phi(\alpha^*(x)), \quad (366)$$

where the term  $\alpha^*(x)$  represents the unique solution in  $\alpha \in (-\infty, 0)$  to

$$\frac{d}{d\alpha} (x\alpha - \phi(\alpha)) = x - \alpha B_z^2 = 0. \quad (367)$$

That is,

$$\alpha^*(x) = \frac{x}{B_z^2}. \quad (368)$$

Plugging (368) into (366) yields,

$$\phi^*(x) = \frac{x^2}{2B_z^2}. \quad (369)$$

Hence, from (365) and (366), given  $c$  in (361) for all  $t \in (-\infty, 0)$ ,

$$ct - \frac{1}{2}t^2 B_z^2 \leq \phi^*(c) \leq D \left( Q \| P_{\Theta|Z=z}^{(\lambda)} \right), \quad (370)$$

and thus,

$$\frac{c^2}{2B_z^2} \leq D \left( Q \| P_{\Theta|Z=z}^{(\lambda)} \right). \quad (371)$$

This implies that either,

$$c \leq \sqrt{2B_z^2 D \left( Q \| P_{\Theta|Z=z}^{(\lambda)} \right)} \quad (372)$$

or

$$c \geq -\sqrt{2B_z^2 D \left( Q \| P_{\Theta|Z=z}^{(\lambda)} \right)}, \quad (373)$$

which leads to

$$\left| \int \mathbb{L}_z(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) - \int \mathbb{L}_z(\boldsymbol{\theta}) dP_{\Theta|Z=z}^{(\lambda)}(\boldsymbol{\theta}) \right| \leq \sqrt{2B_z^2 D \left( Q \| P_{\Theta|Z=z}^{(\lambda)} \right)}, \quad (374)$$

and completes the proof.

## W Proof of Theorem 3.4

The optimization problem in (117) can be re-written as follows:

$$\min_Q \int \mathbb{L}_z(\boldsymbol{\nu}) \frac{dQ}{dP_{\Theta|Z=z}^{(\lambda)}}(\boldsymbol{\nu}) dP_{\Theta|Z=z}^{(\lambda)}(\boldsymbol{\nu}), \quad (375a)$$

$$\text{subject to: } \int \frac{dQ}{dP_{\Theta|Z=z}^{(\lambda)}}(\boldsymbol{\nu}) \log \left( \frac{dQ}{dP_{\Theta|Z=z}^{(\lambda)}}(\boldsymbol{\nu}) \right) dP_{\Theta|Z=z}^{(\lambda)}(\boldsymbol{\nu}) \leq c, \quad (375b)$$

where the optimization is over all measures  $Q$  on  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$  that are absolutely continuous with  $P_{\Theta|Z=z}^{(\lambda)}$  and satisfy

$$\int dQ(\boldsymbol{\nu}) = 1. \quad (376)$$

Let  $\mathcal{M}$  be the set of nonnegative measurable functions with respect to the measurable spaces  $(\text{supp } P_{\Theta|Z=z}^{(\lambda)}, \mathcal{B}(\text{supp } P_{\Theta|Z=z}^{(\lambda)}))$  and  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ . The Lagrangian of the optimization problem in (375) can be constructed in terms of a function in  $\mathcal{M}$ , instead of a measure over the measurable space  $(\text{supp } P_{\Theta|Z=z}^{(\lambda)}, \mathcal{B}(\text{supp } P_{\Theta|Z=z}^{(\lambda)}))$ . Let such Lagrangian  $L : \mathcal{M} \times [0, +\infty)^2 \rightarrow \mathbb{R}$  be of the form

$$\begin{aligned} L(g, \alpha, \beta) = & \int \mathbb{L}_z(\boldsymbol{\nu}) g(\boldsymbol{\nu}) dP_{\Theta|Z=z}^{(\lambda)}(\boldsymbol{\nu}) + \alpha \left( \int g(\boldsymbol{\nu}) \log(g(\boldsymbol{\nu})) dP_{\Theta|Z=z}^{(\lambda)}(\boldsymbol{\nu}) - c \right) \\ & + \beta \left( \int g(\boldsymbol{\nu}) dP_{\Theta|Z=z}^{(\lambda)}(\boldsymbol{\nu}) - 1 \right), \end{aligned} \quad (377)$$

where  $g$  is a notation to represent the Radon-Nikodym derivative  $\frac{dQ}{dP_{\Theta|Z=z}^{(\lambda)}}$ ; the reals  $\alpha$  and  $\beta$  are both nonnegative and act as Lagrangian multipliers due to the constraint (375b) and (139), respectively.

Let  $h : \mathbb{R}^k \rightarrow \mathbb{R}$  be a function in  $\mathcal{M}$ . The Gateaux differential of the functional  $L$  in (377) at  $(g, \alpha, \beta) \in \mathcal{M} \times [0, +\infty)^2$  in the direction of  $h$  is

$$\partial L(g, \alpha, \beta; h) \triangleq \left. \frac{d}{d\gamma} r(\gamma) \right|_{\gamma=0}, \quad (378)$$

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

$$\begin{aligned} r(\gamma) = & \int \mathbb{L}_z(\boldsymbol{\nu}) (g(\boldsymbol{\nu}) + \gamma h(\boldsymbol{\nu})) dP_{\Theta|Z=z}^{(\lambda)}(\boldsymbol{\nu}) + \\ & \alpha \left( \int (g(\boldsymbol{\nu}) + \gamma h(\boldsymbol{\nu})) \log(g(\boldsymbol{\nu}) + \gamma h(\boldsymbol{\nu})) dP_{\Theta|Z=z}^{(\lambda)}(\boldsymbol{\nu}) - c \right) \\ & + \beta \left( \int (g(\boldsymbol{\nu}) + \gamma h(\boldsymbol{\nu})) dP_{\Theta|Z=z}^{(\lambda)}(\boldsymbol{\nu}) - 1 \right), \end{aligned} \quad (379)$$

Note that the derivative of the real function  $r$  in (380) is

$$\begin{aligned} \frac{d}{d\gamma} r(\gamma) = & \int \mathbb{L}_z(\boldsymbol{\nu}) h(\boldsymbol{\nu}) dP_{\Theta|Z=z}^{(\lambda)}(\boldsymbol{\nu}) \\ & + \alpha \int h(\boldsymbol{\nu}) dP_{\Theta|Z=z}^{(\lambda)}(\boldsymbol{\nu}) + \alpha \int h(\boldsymbol{\nu}) \log(g(\boldsymbol{\nu}) + \gamma h(\boldsymbol{\nu})) dP_{\Theta|Z=z}^{(\lambda)}(\boldsymbol{\nu}) \\ & + \beta \int h(\boldsymbol{\nu}) dP_{\Theta|Z=z}^{(\lambda)}(\boldsymbol{\nu}). \end{aligned} \quad (380)$$



From (378) and (380), it follows that

$$\partial L(g, \alpha, \beta; h) = \int h(\boldsymbol{\nu}) (\mathbf{L}_z(\boldsymbol{\nu}) + \alpha(1 + \log g(\boldsymbol{\nu})) + \beta) dP_{\Theta|Z=z}^{(\lambda)}(\boldsymbol{\nu}). \quad (381)$$

From [27, Theorem 1, page 178], it holds that a necessary condition for the functional  $L$  in (377) to have a minimum at  $(g, \alpha, \beta) \in \mathcal{M} \times [0, +\infty)^2$  is that for all functions  $h \in C(\mathcal{M})$

$$\partial L(g, \alpha, \beta; h) = 0, \quad (382)$$

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

$$\mathbf{L}_z(\boldsymbol{\nu}) + \alpha(1 + \log g(\boldsymbol{\nu})) + \beta = 0. \quad (383)$$

Thus,

$$g(\boldsymbol{\nu}) = \exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\nu})}{\alpha}\right) \exp\left(-\frac{\beta + \alpha}{\alpha}\right), \quad (384)$$

where  $\alpha$  and  $\beta$  are chosen to satisfy their corresponding constraints. That is,

$$\frac{dQ}{dP_{\Theta|Z=z}^{(\lambda)}}(\boldsymbol{\nu}) = \frac{\exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\nu})}{\alpha}\right)}{\int \exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\theta})}{\alpha}\right) dP_{\Theta|Z=z}^{(\lambda)}(\boldsymbol{\theta})}, \quad (385)$$

where  $\alpha$  is chosen to satisfy

$$D(Q \| P_{\Theta|Z=z}^{(\lambda)}) = c. \quad (386)$$

From Lemma 2.6, it follows that the probability measure  $Q$  and the  $\sigma$ -finite measure  $P$  satisfy,

$$\frac{dQ}{dP}(\boldsymbol{\nu}) = \frac{dQ}{dP_{\Theta|Z=z}^{(\lambda)}}(\boldsymbol{\nu}) \frac{dP_{\Theta|Z=z}^{(\lambda)}}{dP}(\boldsymbol{\nu}) \quad (387)$$

$$= \left( \frac{\exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\nu})}{\alpha}\right)}{\int \exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\theta})}{\alpha}\right) dP_{\Theta|Z=z}^{(\lambda)}(\boldsymbol{\theta})} \right) \left( \frac{\exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\nu})}{\lambda}\right)}{\int \exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\theta})}{\lambda}\right) dP(\boldsymbol{\theta})} \right) \quad (388)$$

$$= \left( \frac{\exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\nu})}{\alpha}\right)}{\int \frac{\exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\theta})}{\alpha}\right) \exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\theta})}{\lambda}\right)}{\int \exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\alpha})}{\lambda}\right) dP(\boldsymbol{\alpha})} dP(\boldsymbol{\theta})} \right) \left( \frac{\exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\nu})}{\lambda}\right)}{\int \exp\left(-\frac{\mathbf{L}_z(\boldsymbol{\theta})}{\lambda}\right) dP(\boldsymbol{\theta})} \right) \quad (389)$$

$$= \frac{\exp\left(-\left(\frac{1}{\alpha} + \frac{1}{\lambda}\right) \mathbf{L}_z(\boldsymbol{\nu})\right)}{\int \exp\left(-\left(\frac{1}{\alpha} + \frac{1}{\lambda}\right) \mathbf{L}_z(\boldsymbol{\nu})\right) dP(\boldsymbol{\theta})}, \quad (390)$$

which implies that  $Q$  is a  $P$ -Gibbs probability measure on  $(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ , with parameter  $\frac{\alpha\lambda}{\alpha+\lambda}$ . That is, for all  $\boldsymbol{\nu} \in \text{supp } P$ ,

$$Q(\boldsymbol{\nu}) = P_{\Theta|Z=z}^{\left(\frac{\alpha\lambda}{\alpha+\lambda}\right)}(\boldsymbol{\nu}), \quad (391)$$

where  $\alpha$  is chosen to satisfy (386). Let the positive real  $\omega$  be  $\omega \triangleq \frac{\alpha\lambda}{\alpha+\lambda}$ . Thus, from Theorem 2.2, it follows that  $\omega \in (0, \lambda]$  and satisfies  $D(P_{\Theta|Z=z}^{(\omega)}(\nu) \| P_{\Theta|Z=z}^{(\lambda)}) = c$ . The proof ends by verifying that the objective function in (377) is strictly convex, and thus, the measure  $P_{\Theta|Z=z}^{(\omega)}$  is the unique minimizer. This completes the proof.

## References

- [1] I. Alabdulmohsin, “An information-theoretic route from generalization in expectation to generalization in probability,” in *Proc. of the 20th International Conference on Artificial Intelligence and Statistics*, Apr. 2017, pp. 92–100.
- [2] 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.
- [3] R. B. Ash and C. A. Doleans-Dade, *Probability and Measure Theory*, 2nd ed. Burlington, MA: Harcourt/Academic Press, 1999.
- [4] L. Birge and P. Massart, “Rates of convergence for minimum contrast estimators,” *Probability Theory and Related Fields*, vol. 97, pp. 113–150, 1993.
- [5] C. M. Bishop, *Pattern Recognition and Machine Learning*, 1st ed. Berlin, Heidelberg: Springer-Verlag, 2006.
- [6] L. Boltzmann, *Leçons sur la théorie des gaz (première partie)*. Paris, France: Gauthier-Villars, Réédition Jacques Gabay, 1902.
- [7] —, *Leçons sur la théorie des gaz (deuxième partie)*. Paris, France: Gauthier-Villars, Réédition Jacques Gabay, 1905.
- [8] L. Bottou, F. E. Curtis, and J. Nocedal, “Optimization methods for large-scale machine learning,” *Siam Review*, vol. 60, no. 2, pp. 223–311, 2018.
- [9] S. Boucheron, G. Lugosi, and P. Massart, *Concentration inequalities: A nonasymptotic theory of independence*, 1st ed. Oxford, UK: Oxford University Press, 2013.
- [10] 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.
- [11] —, *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.
- [12] T. M. Cover and J. A. Thomas, *Elements of Information Theory*, 2nd ed. Hoboken, NJ: Wiley-Interscience, 2006.
- [13] A. Dembo and O. Zeitouni, *Large Deviations Techniques and Applications*, 2nd ed. New York, NY: Springer-Verlag, 2009.
- [14] A. Engel and C. Van den Broeck, *Statistical mechanics of learning*. Cambridge University Press, 2001.
- [15] W. Feller, *An Introduction to Probability Theory and Its Applications*, 2nd ed. New York, NY: Jhon Wiley & Sons, 1971, vol. II.
- [16] 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.
- [17] J. W. Gibbs, *Elementary principles in statistical mechanics*, 1st ed. New Haven, NJ: Yale University Press, 1902.

- [18] B. Guedj, “A primer on PAC-Bayesian learning.” in *Tutorials of the International Conference on Machine Learning (ICML)*, Jun. 2019.
- [19] B. Guedj and L. Pujol, “Still no free lunches: The price to pay for tighter PAC-Bayes bounds,” *Entropy*, vol. 23, no. 11, 2021.
- [20] M. Haddouche, B. Guedj, O. Rivasplata, and J. Shawe-Taylor, “PAC-Bayes unleashed: Generalisation bounds with unbounded losses,” *Entropy*, vol. 23, no. 10, 2021.
- [21] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning: Data Mining, and Prediction*, 2nd ed. New York, NY: Springer, 2009.
- [22] T. Jaakkola, M. Meila, and T. Jebara, “Maximum entropy discrimination,” *Neural Information Processing Systems*, 1999.
- [23] E. T. Jaynes, “Information theory and statistical mechanics I,” *Physical Review Journals*, vol. 106, pp. 620–630, May 1957.
- [24] —, “Information theory and statistical mechanics II,” *Physical Review Journals*, vol. 108, pp. 171–190, Oct. 1957.
- [25] J. N. Kapur, *Maximum Entropy Models in Science and Engineering*, 1st ed. New York, NY: Wiley, 1989.
- [26] A. J. Kleywegt, A. Shapiro, and T. Homem-de Mello, “The sample average approximation method for stochastic discrete optimization,” *SIAM Journal on Optimization*, vol. 12, no. 2, pp. 479–502, 2002.
- [27] D. G. Luenberger, *Optimization by Vector Space Methods*, 1st ed. New York, NY: Wiley, 1997.
- [28] S. M. Perlaza, “INFO5147 Lecture Notes: Selected topics in information theory,” École Normale Supérieure (ENS) de Lyon, August 2020.
- [29] P. Massart, *Concentration inequalities and model selection*, 1st ed. New York, NY: Springer, 2007.
- [30] D. A. McAllester, “Some PAC-Bayesian theorems,” in *Proceedings of the Eleventh Annual Conference on Computational Learning Theory*, Jul. 1998, pp. 230–234.
- [31] —, “PAC-Bayesian stochastic model selection,” *Machine Learning*, vol. 51, no. 1, pp. 5–21, 2003.
- [32] M. Mezard and A. Montanari, *Information, physics, and computation*, 1st ed. New York, NY, USA: Oxford University Press, 2009.
- [33] D. P. Palomar and S. Verdú, “Lautum information,” *IEEE Transactions on Information Theory*, vol. 54, no. 3, pp. 964–975, Mar. 2008.
- [34] H. V. Poor, *An introduction to signal detection and estimation*, 2nd ed. New York, NY, USA: Springer Science & Business Media, 2013.
- [35] H. Robbins and S. Monro, “A stochastic approximation method,” *The annals of mathematical statistics*, vol. 22, no. 3, pp. 400–407, Sep. 1951.
- [36] C. P. Robert, *The Bayesian choice: From decision-theoretic foundations to computational implementation*, 1st ed. New York, NY: Springer, 2007.

- [37] C. P. Robert, G. Casella, and G. Casella, *Monte Carlo statistical methods*, 2nd ed. New York, NY, USA: Springer, 2004.
- [38] W. Rudin, *Principles of mathematical analysis*, 1st ed. New York, NY: McGraw-Hill Book Company, Inc., 1953.
- [39] 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.
- [40] S. Shalev-Shwartz and S. Ben-David, *Understanding machine learning: From theory to algorithms*, 1st ed. New York, NY, USA: Cambridge University Press, 2014.
- [41] 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.
- [42] L. A. Stefanski and D. D. Boos, “The calculus of M-estimation,” *The American Statistician*, vol. 56, no. 1, pp. 29–38, 2002.
- [43] W. F. Trench, *Introduction to Real Analysis*, 1st ed. Hoboken, NJ: Prentice Hall/Pearson Education, 2003.
- [44] L. G. Valiant, “A theory of the learnable,” *Communications of the ACM*, vol. 27, no. 11, pp. 1134–1142, 1984.
- [45] V. Vapnik, “Principles of risk minimization for learning theory,” in *Advances in neural information processing systems*, vol. 4, 1992, pp. 831–838.
- [46] V. N. Vapnik, “An overview of statistical learning theory,” *IEEE transactions on neural networks*, vol. 10, no. 5, pp. 988–999, 1999.
- [47] R. Vershynin, *High-dimensional probability: An introduction with applications in data science*, 1st ed. New York, NY, USA: Cambridge university press, 2018.
- [48] M. J. Wainwright, *High-dimensional statistics: A non-asymptotic viewpoint*, 1st ed. New York, NY, USA: Cambridge University Press, 2019.
- [49] R. Xin, S. Kar, and U. A. Khan, “Decentralized stochastic optimization and machine learning: A unified variance-reduction framework for robust performance and fast convergence,” *IEEE Signal Processing Magazine*, vol. 37, no. 3, pp. 102–113, May 2020.
- [50] 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.
- [51] L. Zdeborová and F. Krzakala, “Statistical physics of inference: Thresholds and algorithms,” *Advances in Physics*, vol. 65, no. 5, pp. 453–552, Aug. 2016.
- [52] T. Zhang, “From  $\epsilon$ -entropy to KL-entropy: Analysis of minimum information complexity density estimation,” *The Annals of Statistics*, vol. 34, no. 5, pp. 2180–2210, 2006.
- [53] —, “Information-theoretic upper and lower bounds for statistical estimation,” *IEEE Transactions on Information Theory*, vol. 52, no. 4, pp. 1307–1321, Apr. 2006.
- [54] J. Zhu and E. P. Xing, “Maximum entropy discrimination Markov networks,” *Journal of Machine Learning Research*, vol. 10, no. 11, 2009.



**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

ISSN 0249-6399