



Empirical Risk Minimization with Relative Entropy Regularization Type-II

Francisco Daunas, Iñaki Esnaola, Samir M. Perlaza, H. Vincent Poor

► To cite this version:

Francisco Daunas, Iñaki Esnaola, Samir M. Perlaza, H. Vincent Poor. Empirical Risk Minimization with Relative Entropy Regularization Type-II. RR-9508, INRIA, Centre Inria d'Université Côte d'Azur, Sophia Antipolis. 2023. hal-04110899v2

HAL Id: hal-04110899

<https://hal.science/hal-04110899v2>

Submitted on 18 Feb 2024

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 Relative Entropy Regularization Type-II

Francisco Daunas, Iñaki Esnaola, Samir M. Perlaza, and
H. Vincent Poor

**RESEARCH
REPORT**

N° 9508

May 2023

Project-Team NEO

ISRN INRIA/RR--9508--FR+ENG

ISSN 0249-6399



Empirical Risk Minimization with Relative Entropy Regularization Type-II

Francisco Daunas, Iñaki Esnaola, Samir M. Perlaza, and
H. Vincent Poor

Project-Team NEO

Research Report n° 9508 — version 2 — initial version May 2023 —
revised version February 2024 — 59 pages

Abstract: The effect of relative entropy asymmetry is analyzed in the context of empirical risk minimization with relative entropy regularization (ERM-RER). A novel regularization is introduced, termed Type-II regularization, that allows for solutions to the ERM-RER problem with a support that extends outside the support of the reference measure. The solution to the ERM-RER Type-II problem is characterized in terms of the Radon-Nikodym derivative of the reference measure with respect to the solution. Analysis of the solution unveils the following properties of relative entropy when it acts as a regularizer in the ERM-RER problem: *i*) relative entropy forces the support of the Type-II solution to collapse into the support of the reference measure, which introduces a strong inductive bias that dominates the evidence provided by the training data; and *ii*) Type-II regularization is equivalent to classical relative entropy regularization with an appropriate transformation of the empirical risk function. Closed-form expressions for the expected empirical risk as a function of the regularization parameters are provided.

Key-words: Supervised Learning, Empirical Risk Minimization; Relative Entropy; Regularization; Gibbs Measure; Inductive Bias; Sensitivity; and Generalization.

Francisco Daunas and Iñaki Esnaola are with the Department of Automatic Control and Systems Engineering at the University of Sheffield, Sheffield, UK. Francisco Daunas is also a member of the NEO Team at INRIA, Centre Inria d'Université Côte d'Azur, Sophia Antipolis.

Samir M. Perlaza is with INRIA, Centre Inria d'Université Côte d'Azur, Sophia Antipolis, France and with the GAATI Laboratory of the Université de la Polynésie Française, Faaa, French Polynesia.

H. Vincent Poor, Iñaki Esnaola, and Samir M. Perlaza are with the ECE Department at Princeton University, Princeton, NJ.

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 Tapez-II

Résumé : L'effet de l'asymétrie de l'entropie relative est analysé dans le problème de la minimisation du risque empirique avec régularisation par entropie relative (MRE-RER). Une nouvelle régularisation est introduite, appelée régularisation de Type-II, qui permet la recherche des solutions au problème de la MRE-RER avec un support qui s'étend en dehors du support de la mesure de référence. La solution au nouveau problème de la Type-II MRE-RER est caractérisée analytiquement en termes de la dérivée de Radon-Nikodym de la mesure de référence par rapport à la solution. L'analyse de la solution dévoile certaines propriétés de l'entropie relative lorsqu'elle agit comme régularisateur du problème de la MRE-RER : (a) l'entropie relative force le support de la solution de Type-II à s'étaler sur tout le support de la mesure de référence, ce qui introduit un fort biais inductif qui domine l'évidence fournie par les données d'entraînement ; (b) La régularisation de Type-II est équivalente à la régularisation d'entropie relative classique avec une transformation appropriée de la fonction du risque empirique. Enfin, une expression sous forme explicite de la valeur espérée du risque empirique en fonction des paramètres de régularisation est présentée.

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

Contents

1	Introduction	5
2	Empirical Risk Minimization	6
3	The Type-I ERM-RER Problem	7
4	The Type-II ERM-RER Problem	8
4.1	The Solution to the Type-II ERM-RER Problem	9
4.2	Properties of the Solution	11
4.2.1	The Normalization Function	11
4.2.2	Bounds on the Radon-Nikodym Derivative	16
4.2.3	Asymptotes of the Radon-Nikodym Derivative	17
5	The Expected Empirical Risk	18
6	(δ, ϵ)-Optimality	20
7	Interplay Between the Relative Entropy Asymmetry and the Empirical Risk	20
7.1	The Connection between Type-I and Type-II	21
7.2	Sensitivity of the Log-Empirical Risk	22
7.3	Type-I and Type-II Optimal Measures	24
8	Final Remarks	25
A	Proof of Lemma 4.1	26
B	Proof of Lemma 4.2	30
C	Proof of Lemma 4.4	31
D	Proof of Lemma 4.5	33
E	Proof of Lemma 4.8	34
F	Proof of Lemma 4.9	35
G	Proof of Lemma 4.10	36
H	Proof of Lemma 4.11	37
I	Proof of Lemma 4.12	37
J	Proof of Lemma 4.13	41
K	Proof of Lemma 5.1	42

L	Proof of Lemma 5.2	43
M	Proof of Lemma 5.3	45
N	Proof of Lemma 5.5	46
O	Proof of Theorem 6.1	47
P	Proof of Lemma 7.1	48
Q	Proof of Lemma 7.2	49
R	Proof of Lemma 7.4	50
S	Example 4.3	52
T	Example 4.2	53

1 Introduction

Empirical risk minimization (ERM) is a central tool in supervised machine learning. Among other uses, it enables the characterization of sample complexity and probably approximately correct (PAC) learning in a wide range of settings [1]. The application of ERM in the study of theoretical guarantees spans related disciplines such as machine learning [2], information theory [3,4] and statistics [5,6]. Classical problems such as classification [7,8], pattern recognition [9,10], regression [11,12], and density estimation [9,13] can be posed as special cases of the ERM problem [13,14]. Unfortunately, ERM is prone to training data memorization, a phenomenon also known as overfitting [15–17]. For that reason, ERM is often regularized in order to provide generalization guarantees [18–21]. Regularization establishes a preference over the models by encoding features of interest that conform to prior knowledge. In different statistical learning frameworks, such as Bayesian learning [22,23] and PAC learning [24–26], the prior knowledge over the set of models can be described by a reference probability measure. More general references can be adapted as proved in [27,28] for the case of σ -finite measures. Prior knowledge on the set of datasets can also be represented by probability measures, e.g., the worst-case data-generating probability measure introduced in [29]. In either case, the solution to the regularized ERM problem can be cast as a probability distribution over the set of models.

A common regularizer of the ERM problem is the relative entropy of the optimization probability measure with respect to a given reference measure over the set of models [13,30–32]. The resulting problem formulation, termed ERM with relative entropy regularization (ERM-RER) has been extensively studied for both the case in which the reference measure is a probability measure [30–33] and the case in which it is a σ -finite measure [27,28,34]. While in both cases the solution is unique and corresponds to a Gibbs probability measure, the existence of the solution is ensured only in the case in which the reference measure is a probability measure [28]. Despite the many merits of the ERM-RER formulation, it has some significant limitations. Firstly, absolute continuity of the optimization measure with respect to the reference measure is required for the existence of the corresponding Radon-Nikodym derivative, which is used by the relative entropy regularization. This absolute continuity sets an insurmountable barrier to the exploration of models outside the support of the reference measure. More specifically, models outside the support of the reference measure exhibit zero probability with respect to the Gibbs probability measure solution to ERM-RER. More importantly, this holds regardless of the evidence provided by the training dataset. Secondly, the choice of relative entropy over alternative divergences often follows arguments based on the simplicity of obtaining generalization guarantees in the form of bounds [18]. Nonetheless, such bounds are often hard to calculate and are not always informative when evaluated in practical settings [35–37].

From all the above, exploring the asymmetry of the relative entropy is of particular interest to advancing the understanding of entropy regularization and

its role in generalization. The problem of ERM with a general f -divergence regularization has been explored in [38] and [39] in the case of a finite countable set of models, and recently extended to uncountable sets of models in [40]. Nonetheless, the authors in [38–40] constrain the optimization domains to sets of measures that are mutually absolutely continuous with respect to the reference probability. The use of the relative entropy of the optimization measure with respect to the reference in ERM-RER is termed Type-I ERM-RER. Alternatively, the use of the relative entropy of the reference measure with respect to the optimization measure, introduced here, is termed Type-II ERM-RER. Interestingly, the existing results in [38–40], which lead to special cases of the Type-I and Type-II ERM-RER problems by assuming that $f(x) = -x \log(x)$ and $f(x) = -\log(x)$, respectively, do not study the impact of the asymmetry of relative entropy in the context of ERM regularization.

This paper presents the solution to Type-II ERM-RER optimization problem using a new method of proof. In particular, mutual absolute continuity of the measures in the optimization domain with respect to the reference measure is not imposed. Nonetheless, such a mutual absolute continuity is exhibited by the solution as a consequence of the structure of the problem. The main properties of the solution are highlighted and an equivalence between Type-I and a Type-II ERM-RER problems is presented by replacing the empirical risk in the Type-I problem by another function, which can be cast as a tunable loss function as in [41–43].

The remainder of the paper is organized as follows. Section 2 presents the standard ERM problem. Section 3 describes the Type-I regularization. The main contribution of this paper, which is the solution to the Type-II ERM-RER problem, is presented in Section 4. Section 7 studies the equivalence between Type-I and Type-II regularization. The conclusions are summarized in Section 8.

2 Empirical Risk Minimization

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

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

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

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

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

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

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

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

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

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

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

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

The PAC and Bayesian frameworks, c.f. in [23] and [25], solve the problem in (4) by constructing probability measures conditioned on the dataset \mathbf{z} , from which models are randomly sampled. In this context, finding probability measures that are minimizers of the ERM problem in (4) over the set of all probability measures that can be defined on the measurable space $(\mathcal{M}, \mathcal{F})$, which is denoted by $\Delta(\mathcal{M})$, requires a metric that enables assessing the goodness of the probability measure. From this perspective, the underlying assumption in the reminder of this work is that the functions f and ℓ in (3) are such that for all $(x, y) \in \mathcal{X} \times \mathcal{Y}$, the function $g_{x,y} : \mathcal{M} \rightarrow [0, \infty]$, such that $g_{x,y}(\theta) = \ell(f_{\theta}(x), y)$, is measurable with respect to the Borel measurable spaces $(\mathcal{M}, \mathcal{F})$ and $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$, where \mathcal{F} and $\mathcal{B}(\mathbb{R})$ are respectively the Borel σ -fields on \mathcal{M} and \mathbb{R} . Under these assumptions, a common metric is the notion of expected empirical risk.

Definition 2.1 (Expected Empirical Risk). *Given the dataset $\mathbf{z} \in (\mathcal{X} \times \mathcal{Y})^n$ in (1), for all probability measures $P \in \Delta(\mathcal{M})$ let the functional $\mathbf{R}_{\mathbf{z}} : \Delta(\mathcal{M}) \rightarrow [0, \infty)$ be such that*

$$\mathbf{R}_{\mathbf{z}}(P) \triangleq \int \mathbf{L}_{\mathbf{z}}(\theta) dP(\theta), \quad (6)$$

where the function $\mathbf{L}_{\mathbf{z}}$ is defined in (3).

In the following section, the Type-I relative entropy regularization is reviewed as it serves as the basis for the analysis of the regularization asymmetry.

3 The Type-I ERM-RER Problem

The Type-I ERM-RER problem is parametrized by a probability measure $Q \in \Delta(\mathcal{M})$ and a real $\lambda \in (0, \infty)$. The measure Q is referred to as the *reference*

measure and λ as the *regularization factor*. The Type-I ERM-RER problem, with parameters Q and λ , is given by the following optimization problem:

$$\min_{P \in \Delta_Q(\mathcal{M})} R_z(P) + \lambda D(P\|Q), \quad (7)$$

where the functional R_z is defined in (6), and the optimization domain is

$$\Delta_Q(\mathcal{M}) \triangleq \{P \in \Delta(\mathcal{M}) : P \ll Q\}. \quad (8)$$

The notation $P \ll Q$ stands for P being absolutely continuous with respect to Q .

The solution to the Type-I ERM-RER problem in (7) is the Gibbs probability measure [27, 30, 31], which is presented in the following lemma.

Lemma 3.1 ([28, Theorem 3.1]). *The solution of the Type-I ERM-RER problem in (7), denoted by $P_{\Theta|Z=z}^{(Q,\lambda)} \in \Delta(\mathcal{M})$, is unique, always exists, and satisfies for all $\theta \in \text{supp } Q$*

$$\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) = \exp\left(-K_{Q,z}\left(-\frac{1}{\lambda}\right) - \frac{1}{\lambda}L_z(\theta)\right), \quad (9)$$

where the function $K_{Q,z} : \mathbb{R} \rightarrow \mathbb{R}$ is such that

$$K_{Q,z}(t) = \log\left(\int \exp(tL_z(\theta)) dQ(\theta)\right). \quad (10)$$

4 The Type-II ERM-RER Problem

The Type-II ERM-RER problem is parametrized by a probability measure $Q \in \Delta(\mathcal{M})$ and a real $\lambda \in (0, \infty)$. As in Type-I ERM-RER problem, the measure Q is referred to as the *reference measure* and λ as the *regularization factor*. Given the dataset $z \in (\mathcal{X} \times \mathcal{Y})^n$ in (1), the Type-II ERM-RER problem, with parameters Q and λ , consists of the following optimization problem:

$$\min_{P \in \nabla_Q(\mathcal{M})} R_z(P) + \lambda D(Q\|P), \quad (11)$$

where the functional R_z is defined in (6), and the optimization domain is

$$\nabla_Q(\mathcal{M}) \triangleq \{P \in \Delta(\mathcal{M}) : Q \ll P\}. \quad (12)$$

The difference between Type-I and Type-II ERM-RER problems lies on the regularization. While the former uses the relative entropy $D(P\|Q)$, the latter uses $D(Q\|P)$. This translates into different optimization domains due to the asymmetry of the relative entropy. More specifically, in the Type-I ERM-RER problem, the optimization domain is the set of probability measures on the Borel measurable space $(\mathcal{M}, \mathcal{F})$ that are absolutely continuous with the reference

measure Q . That is, the set $\Delta_Q(\mathcal{M})$ in (8). Alternatively, in the Type-II ERM-RER problem, the optimization domain consists of probability measures defined on the Borel measurable space $(\mathcal{M}, \mathcal{F})$, with the additional condition that the reference measure Q must be absolutely continuous with respect to them. This corresponds to the set denoted as $\nabla_Q(\mathcal{M})$ in Equation (12). From this perspective, the techniques used in [28] for solving the Type-I ERM-RER no longer hold. As shown in the next section, a new technique is used for solving the Type-II ERM-RER.

The problem in (11) exhibits a trivial solution when the functional R_z is such that for all $P \in \nabla_Q(\mathcal{M})$, it holds that $R_z(P) = c$, for some $c \in [0, \infty)$. In such a case, the solution is unique and equal to the probability measure Q , independently of the parameter λ . In order to avoid this trivial case, which arises under particular conditions over the empirical risk function L_z in (3) and the probability measure Q , the notion of separability of the empirical risk function with respect to the measure Q is borrowed from [28]. A separable empirical risk function with respect to a given probability measure P is defined as follows.

Definition 4.1 (Definition 4.1 in [28]). *The empirical risk function L_z in (3) is said to be separable with respect to the probability measure $P \in \Delta(\mathcal{M})$, if there exist a positive real $c > 0$ and two subsets \mathcal{A} and \mathcal{B} of \mathcal{M} that are nonnegligible with respect to P , and for all $(\theta_1, \theta_2) \in \mathcal{A} \times \mathcal{B}$,*

$$L_z(\theta_1) < c < L_z(\theta_2) < \infty. \quad (13)$$

A nonseparable empirical risk function L_z in (3) with respect to a measure P is a constant almost surely with respect to the measure P . More specifically, there exists a real $a \geq 0$, such that

$$P(\{\theta \in \mathcal{M} : L_z(\theta) = a\}) = 1. \quad (14)$$

When the empirical risk function L_z in (3) is nonseparable with respect to all measures in $P \in \nabla_Q(\mathcal{M})$, the trivial case described above is observed. The notion of separable empirical risk functions would play a central role in the study of the optimization problem in (11).

4.1 The Solution to the Type-II ERM-RER Problem

The solution of the Type-II ERM-RER problem in (11) is presented in the following theorem, where the measure Q is a parameter of the optimization problem in (11) and the function L_z is defined in (3).

Theorem 4.1. *If there exists a real β such that*

$$\beta \in \{t \in \mathbb{R} : \forall \theta \in \text{supp } Q, 0 < t + L_z(\theta)\}, \quad (15a)$$

and

$$\int \frac{\lambda}{\beta + L_z(\theta)} dQ(\theta) = 1, \quad (15b)$$

then the solution to the optimization problem in (11), denoted by $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)} \in \Delta(\mathcal{M})$, is unique and for all $\theta \in \text{supp } Q$, it satisfies

$$\frac{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) = \frac{\lambda}{\beta + L_z(\theta)}. \quad (16)$$

Before introducing the proof of Theorem 4.1, two important results are presented. The first result consists in the solution to the optimization problem in (11) when the optimization domain is restricted to

$$\bigcirc_Q(\mathcal{M}) \triangleq \nabla_Q(\mathcal{M}) \cap \Delta_Q(\mathcal{M}), \quad (17)$$

where the sets $\Delta_Q(\mathcal{M})$ and $\nabla_Q(\mathcal{M})$ are defined in (8) and (12), respectively. Such an ancillary problem can be formulated as follows:

$$\min_{P \in \bigcirc_Q(\mathcal{M})} R_z(P) + \lambda D(Q\|P). \quad (18)$$

The solution to such an ancillary problem is described by the following lemma.

Lemma 4.1. *The solution to the optimization problem in (18) is unique and identical to the probability measure $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ in (16).*

Proof: The proof is presented in Appendix A. ■

The second result consists of comparing the optimal values resulting from the optimizations in (11) and (18). The following lemma formalizes this result.

Lemma 4.2. *The optimization problems in (11) and (18) satisfy*

$$\min_{P \in \nabla_Q} R_z(P) + \lambda D(Q\|P) \geq \min_{P \in \bigcirc_Q} R_z(P) + \lambda D(Q\|P). \quad (19)$$

Proof: The proof is presented in Appendix B. ■

Lemma 4.2 unveils the fact that the objective function in (11) when evaluated at measures whose support extends beyond the support of Q is larger than such an objective function evaluated at measures whose support is identical to the reference measure. This includes the case in which the set $\mathcal{T}(z)$ in (5) lies outside the support of Q .

Using these results, the proof of Theorem 4.1 is as follows.

Proof of Theorem 4.1: The proof follows by observing that from (17), it holds that

$$\bigcirc_Q(\mathcal{M}) \subseteq \nabla_Q(\mathcal{M}). \quad (20)$$

Hence, from (20), it follows that

$$\min_{P \in \nabla_Q} R_z(P) + \lambda D(Q\|P) \leq \min_{P \in \bigcirc_Q} R_z(P) + \lambda D(Q\|P). \quad (21)$$

From the inequalities in (19) and (21), it follows that

$$\min_{P \in \nabla_Q} R_z(P) + \lambda D(Q \| P) = \min_{P \in \bigcirc_Q} R_z(P) + \lambda D(Q \| P). \quad (22)$$

Thus, the measure $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ in (16) is the solution of the optimization problem in (11), which completes the proof of Theorem 4.1. ■

Lemma 4.2 implies that the solution to the optimization problem in (11) is in the set $\bigcirc_Q(\mathcal{M})$ in (16). A consequence of this observation is the following corollary.

Corollary 4.3. *The probability measures Q and $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ in (16) are mutually absolutely continuous.*

Corollary 4.3 also follows from Theorem 4.1 by observing that the solution to the Type-II ERM-RER problem in (11) is expressed in terms of its Radon-Nikodym with respect to Q , which implies the absolute continuity of $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ with respect to Q . The absolute continuity of the measure Q with respect to $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ follows from the optimization domain of the Type-II ERM-RER problem. From this perspective, Corollary 4.3 conveys the fact that there does not exist a dataset that can overcome the inductive bias induced by the reference measure Q . That is, set of models outside the support of Q exhibit zero probability measure with respect to the measure $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$.

This observation is important as at a first glance, the Type-II relative entropy regularization for the ERM problem in (11) does not restrict the solution to be absolutely continuous with respect to the reference measure Q . However, Theorem 4.1 shows that the support of the probability measure $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ in (16) collapses into the support of the reference. A parallel can be established between Type-I and Type-II cases, as in both cases the support of the solution is the support of the reference measure. In a nutshell, the use of relative entropy regularization inadvertently forces the solution to coincide with the support of the reference regardless of the training data.

4.2 Properties of the Solution

This section introduces a function referred to as the *normalization function* and its properties. This function as well as its properties are central for studying the solution to the ERM-RER problem in (11).

4.2.1 The Normalization Function

Let the function

$$\bar{K}_{Q,z} : (0, \infty) \rightarrow \mathcal{A}, \quad (23a)$$

with $\mathcal{A} \subseteq \mathbb{R}$, and Q and z in (11), be such that for all $t \in (0, \infty)$

$$\bar{K}_{Q,z}(t) = \alpha, \quad (23b)$$

where α satisfies

$$\int \frac{t}{\alpha + \mathbf{L}_{\mathbf{z}}(\boldsymbol{\theta})} dQ(\boldsymbol{\theta}) = 1, \quad (24)$$

and the function $\mathbf{L}_{\mathbf{z}}$ is defined in (3). The function $\bar{K}_{Q,\mathbf{z}}$ is referred to as the *normalization function*. This is essentially due to the observation that β and λ in (16) satisfy

$$\bar{K}_{Q,\mathbf{z}}(\lambda) = \beta, \quad (25)$$

which guarantees that the measure $\bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)} \in \Delta(\mathcal{M})$ in (16) is a probability measure. Hence, the Radon-Nikodym derivative in (16) satisfies for all $\boldsymbol{\theta} \in \text{supp } Q$

$$\frac{d\bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}}{dQ}(\boldsymbol{\theta}) = \frac{\lambda}{\bar{K}_{Q,\mathbf{z}}(\lambda) + \mathbf{L}_{\mathbf{z}}(\boldsymbol{\theta})}. \quad (26)$$

The properties of the function $\bar{K}_{Q,\mathbf{z}}$ in (23b) are studied using the notion of Rashamon sets [44]. Given a real $\delta \in [0, \infty)$, consider the Rashamon set

$$\mathcal{L}_{\mathbf{z}}(\delta) \triangleq \{\boldsymbol{\theta} \in \mathcal{M} : \mathbf{L}_{\mathbf{z}}(\boldsymbol{\theta}) \leq \delta\}. \quad (27)$$

Consider also the nonnegative real number

$$\delta_{Q,\mathbf{z}}^* \triangleq \inf\{\delta \in [0, \infty) : Q(\mathcal{L}_{\mathbf{z}}(\delta)) > 0\}. \quad (28)$$

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

$$\mathcal{L}_{Q,\mathbf{z}}^* \triangleq \{\boldsymbol{\theta} \in \mathcal{M} : \mathbf{L}_{\mathbf{z}}(\boldsymbol{\theta}) = \delta_{Q,\mathbf{z}}^*\}. \quad (29)$$

Using the objects defined above, the following lemma introduces one of the main properties of the function $\bar{K}_{Q,\mathbf{z}}$ in (23).

Lemma 4.4. *The function $\bar{K}_{Q,\mathbf{z}}$ in (23) is strictly increasing and continuous.*

Proof: The proof is presented in Appendix C. ■

As a consequence of Lemma 4.4, the continuity of function $\bar{K}_{Q,\mathbf{z}}$ in (23) and equality (24), it follows that for all $\alpha \in \mathcal{A}$, with \mathcal{A} in (23a), the functional inverse $\bar{K}_{Q,\mathbf{z}}^{-1} : \mathcal{A} \rightarrow (0, \infty)$ is given by

$$\bar{K}_{Q,\mathbf{z}}^{-1}(\alpha) = \frac{1}{\int \frac{1}{\mathbf{L}_{\mathbf{z}}(\boldsymbol{\theta}) + \alpha} dQ(\boldsymbol{\theta})}. \quad (30)$$

Note that the function $\bar{K}_{Q,\mathbf{z}}^{-1}$ in (30) allows defining the codomain of the function $\bar{K}_{Q,\mathbf{z}}$ in (23b), which is formalized by the following lemma.

Lemma 4.5. *The set \mathcal{A} in (23a) satisfies*

$$(-\delta_{Q,z}^*, \infty) \subseteq \mathcal{A} \subseteq [-\delta_{Q,z}^*, \infty). \quad (31)$$

Moreover, the set \mathcal{A} is identical to $[-\delta_{Q,z}^, \infty)$ if and only if*

$$\int \frac{1}{L_z(\boldsymbol{\theta}) - \delta_{Q,z}^*} dQ(\boldsymbol{\theta}) < \infty, \quad (32)$$

with $\delta_{Q,z}^$ defined in (28).*

Proof: The proof is presented in Appendix D. ■

Observe that for all finite sets of models \mathcal{M} , the subset $\mathcal{L}_{Q,z}^*$ as defined in (29) satisfies $Q(\mathcal{L}_{Q,z}^*) > 0$. This implies that the integral in (32) is not finite. The formalization of this observation is presented in the following corollary.

Corollary 4.6. *If the set \mathcal{M} is finite, then the set \mathcal{A} in (23a) is identical to $(-\delta_{Q,z}^*, \infty)$.*

An additional remark, is that from Lemma 4.4 and Lemma 4.5, if the set \mathcal{A} in (23a) is equal to the closed set $[\delta_{Q,z}^*, \infty)$ in (31), there exists a minimum regularization factor $\lambda^* > 0$, where $\lambda^* = \bar{K}_{Q,z}^{-1}(-\delta_{Q,z}^*)$, with $\bar{K}_{Q,z}^{-1}$ defined in (30). This implies that the mapping of the function $\bar{K}_{Q,z}$ in (23) is such that $\bar{K}_{Q,z} : [\lambda^*, \infty) \rightarrow [-\delta_{Q,z}^*, \infty)$. This minimum regularization implies that the existence of a solution to the optimization problem in (11) only exists for regularization factors $\lambda \geq \lambda^*$.

Based on Lemma 4.5, the existence of the minimum regularization factor depends on whether the integral in (32) is finite. The finiteness of the integral is conditioned by the reference measure Q ; the functions ℓ and f in (3); and the dataset \mathbf{z} in (1). The subsequent examples illustrate situations where the set \mathcal{A} is the open set $(\delta_{Q,z}^*, \infty)$ and closed set $[\delta_{Q,z}^*, \infty)$ in (31).

Example 4.1. *Consider the Type-II ERM-RER problem in (11) and assume that: (a) $\mathcal{M} = \mathcal{X} = \mathcal{Y} = [0, \infty)$; (b) $\mathbf{z} = (1, 0)$ and (c) $Q \ll \mu$, with μ the Lebesgue measure, such that for all $\boldsymbol{\theta} \in \text{supp } Q$,*

$$\frac{dQ}{d\mu}(\boldsymbol{\theta}) = 4\theta^2 \exp(-2\boldsymbol{\theta}). \quad (33)$$

Let also the function $f : \mathcal{M} \times \mathcal{X} \rightarrow \mathcal{Y}$ be

$$f(\boldsymbol{\theta}, x) = x\boldsymbol{\theta}, \quad (34)$$

and the loss function ℓ in (2) be

$$\ell(f(\boldsymbol{\theta}, x), y) = (x\boldsymbol{\theta} - y)^2, \quad (35)$$

which implies

$$L_z(\boldsymbol{\theta}) = (x\boldsymbol{\theta} - y)^2, \quad (36)$$

with the function \mathbb{L}_z defined in (3). Under the current assumptions the objective of this example is to show that $\mathcal{A} = [\delta_{Q,z}^*, \infty)$. For this purpose, it is sufficient to show that the inequality in (32) holds. From Theorem 4.1, it follows that $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ in (16) satisfies for all $\theta \in \text{supp } Q$,

$$\frac{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}{d\mu}(\theta) = \frac{\lambda}{\mathbb{L}_z(\theta) + \beta} 4\theta^2 \exp(-2\theta), \quad (37)$$

with β satisfying (15). Thus,

$$\int \frac{1}{\mathbb{L}_z(\theta) - \delta_{Q,z}^*} dQ(\theta) = \int \frac{1}{\mathbb{L}_z(\theta) - \delta_{Q,z}^*} 4\theta^2 \exp(-2\theta) d\mu(\theta) \quad (38a)$$

$$= \int_0^\infty \frac{4\theta^2 \exp(-2\theta)}{(x\theta - y)^2 - \delta_{Q,z}^*} d\theta \quad (38b)$$

$$= \int_0^\infty \frac{4\theta^2 \exp(-2\theta)}{\theta^2 - \delta_{Q,z}^*} d\theta \quad (38c)$$

$$= \int_0^\infty \frac{4\theta^2 \exp(-2\theta)}{\theta^2} d\theta \quad (38d)$$

$$= \int_0^\infty 4 \exp(-2\theta) d\theta \quad (38e)$$

$$= 2, \quad (38f)$$

where equality (38a) follows from equality (33); equality (38c) follows from the assumption that $(x, y) = (1, 0)$; and equality (38d) follows from the fact that $\delta_{Q,z}^* = 0$. Finally, the function $\bar{K}_{Q,z}$ in (23) satisfies $\bar{K}_{Q,z}(\frac{1}{2}) = 0$, which implies $\delta_{Q,z}^* = 0 \in \mathcal{A}$.

Example 4.2. Consider Example 4.1 with $z = (1, 1)$. Under the current assumptions, the objective of this example is to show that $\mathcal{A} = (\delta_{Q,z}^*, \infty)$. For this purpose, it is sufficient to show that the inequality in (32) does not hold:

$$\int \frac{1}{\mathbb{L}_z(\theta) - \delta_{Q,z}^*} dQ(\theta) = \int \frac{1}{\mathbb{L}_z(\theta) - \delta_{Q,z}^*} 4\theta^2 \exp(-2\theta) d\mu(\theta) \quad (39a)$$

$$= \int_0^\infty \frac{4\theta^2 \exp(-2\theta)}{(x\theta - y)^2 - \delta_{Q,z}^*} d\theta \quad (39b)$$

$$= \int_0^\infty \frac{4\theta^2 \exp(-2\theta)}{(\theta - 1)^2} d\theta \quad (39c)$$

$$= \infty. \quad (39d)$$

where equality (39a) follows from equality (33); equality (39b) follows from the assumption that $(x, y) = (1, 0)$; equality (39c) follows from the fact that $\delta_{Q,z}^* = 0$; and the equality (39d) follows from an algebraic development shown in Appendix T. Finally, the function $\bar{K}_{Q,z}$ in (23) is undefined at zero, which implies $\delta_{Q,z}^* = 0 \notin \mathcal{A}$.

The aforementioned examples demonstrate that even if the reference measure Q and functions ℓ and f in (3) are fixed, the set \mathcal{A} might be either $[\delta_{Q,z}^*, \infty)$ or $(\delta_{Q,z}^*, \infty)$ depending on the dataset \mathbf{z} . This observation underscores that the existence of the minimum regularization factor λ^* is coupled on the specific choices of Q , ℓ , f , and \mathbf{z} . Finally, from Lemma 4.5 it is implied that if the set \mathcal{M} is finite, then the set \mathcal{A} is the open set $(-\delta_{Q,z}^*, \infty)$ in (31), as shown by the following example.

Example 4.3. Consider the Type-II ERM-RER problem in (11) and assume that: (a) \mathcal{B} is a proper subset of \mathcal{M} , and (b) the probability measure Q satisfies

$$Q(\mathcal{B}) = \epsilon, \quad \text{and} \quad (40a)$$

$$Q(\mathcal{M} \setminus \mathcal{B}) = 1 - \epsilon, \quad (40b)$$

with $\epsilon > 0$. Let the empirical risk function \mathbf{L}_z in (3) be

$$\mathbf{L}_z(\theta) = \begin{cases} 0 & \text{if } \theta \in \mathcal{B} \\ c & \text{if } \theta \in \mathcal{M} \setminus \mathcal{B}, \end{cases} \quad (41)$$

with $c > 0$.

Under the current assumptions, the objective of this example is to show that for all $\mathbf{z} \in (\mathcal{X} \times \mathcal{Y})^n$, it holds that $\mathcal{A} = (\delta_{Q,z}^*, \infty)$. For this purpose, it is sufficient to show that for all $c \in (0, \infty)$, the function $\bar{K}_{Q,z}$ in (23b) is strictly greater than $-\delta_{Q,z}^*$. From the current assumptions it follows that $\delta_{Q,z}^* = 0$ and that for all $\lambda \in (0, +\infty)$ the function $\bar{K}_{Q,z}$ in (23b) satisfies

$$\bar{K}_{Q,z}(\lambda) = -\frac{(c-\lambda)}{2} + \sqrt{\left(\frac{c-\lambda}{2}\right)^2 + \lambda c Q(\mathcal{B})}. \quad (42)$$

The proof of equality (42) is presented in appendix S. From equality (42), it holds that

$$\bar{K}_{Q,z}(\lambda) = -\frac{(c-\lambda)}{2} + \sqrt{\left(\frac{c-\lambda}{2}\right)^2 + \lambda c Q(\mathcal{B})} \quad (43a)$$

$$> -\frac{(c-\lambda)}{2} + \sqrt{\left(\frac{c-\lambda}{2}\right)^2} \quad (43b)$$

$$= -\frac{(c-\lambda)}{2} + \left| \frac{c-\lambda}{2} \right| \quad (43c)$$

$$\geq 0 \quad (43d)$$

$$= -\delta_{Q,z}^*, \quad (43e)$$

which proves that for all $c \in (0, \infty)$, $\mathcal{A} = (-\delta_{Q,z}^*, \infty)$.

Another immediate consequence of Lemma 4.4 is the following corollary.

Corollary 4.7. If the real value $\delta_{Q,z}^* = 0$, with $\delta_{Q,z}^*$ in (28), then the function $\bar{K}_{Q,z}$ in (23b) is strictly positive.

The following lemma characterizes the limit of $\bar{K}_{Q,z}$ as λ approaches zero from the right, under the assumption that the set \mathcal{A} is $(-\delta_{Q,z}^*, \infty)$.

Lemma 4.8. *If the set \mathcal{A} is $(-\delta_{Q,z}^*, \infty)$, the function $\bar{K}_{Q,z}$ in (23b) satisfies*

$$\lim_{\lambda \rightarrow 0^+} \bar{K}_{Q,z}(\lambda) = -\delta_{Q,z}^*, \quad (44)$$

where $\delta_{Q,z}^*$ is defined in (28).

Proof: The proof is presented in Appendix E. ■

Studying cases in which the choice of the reference measure Q , empirical risk function \mathbb{L}_z and dataset z induce the set \mathcal{A} to be equal $(-\delta_{Q,z}^*, \infty)$ is relevant in the context of learning algorithm, from the fact that if \mathcal{M} is countable the set \mathcal{A} is equal $(-\delta_{Q,z}^*, \infty)$. Furthermore, as the number of data points in datasets increases, choosing a prior Q such that (32) is satisfied becomes less likely. Hence, in the paper, it will be assumed that Q , \mathbb{L}_z and z induce $\mathcal{A} = (-\delta_{Q,z}^*, \infty)$ unless otherwise stated.

4.2.2 Bounds on the Radon-Nikodym Derivative

Note that from Theorem 4.1 models resulting in lower empirical risks correspond to greater values of the Radon-Nikodym derivative. The following corollary formalizes this observation.

Lemma 4.9. *For all $(\theta_1, \theta_2) \in (\text{supp } Q)^2$, such that $\mathbb{L}_z(\theta_1) \leq \mathbb{L}_z(\theta_2)$, with \mathbb{L}_z in (3), the Radon-Nikodym derivative in (16) satisfies*

$$\frac{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta_2) \leq \frac{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta_1), \quad (45)$$

with equality if and only if $\mathbb{L}_z(\theta_1) = \mathbb{L}_z(\theta_2)$.

Proof: The proof is presented in Appendix F. ■

The Radon-Nikodym derivative $\frac{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}$ in (16) is always finite and strictly positive. This observation is formalized in the following lemma.

Lemma 4.10. *The Radon-Nikodym derivative $\frac{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}$ in (16) satisfies for all $\theta \in \text{supp } Q$*

$$0 < \frac{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) \leq \frac{\lambda}{\delta_{Q,z}^* + \bar{K}_{Q,z}(\lambda)} < \infty, \quad (46)$$

where the equality holds if and only if $\theta \in \mathcal{L}_{Q,z}^* \cap \text{supp } Q$.

Proof: The proof is presented in Appendix G. ■

4.2.3 Asymptotes of the Radon-Nikodym Derivative

In the asymptotic regime, when the regularization factor grows to infinity, *i.e.*, $\lambda \rightarrow \infty$, the measure $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ becomes identical to the reference measure Q , as described in the following lemma.

Lemma 4.11. *The Radon-Nikodym derivative $\frac{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}$ in (16) satisfies for all $\theta \in \text{supp } Q$,*

$$\lim_{\lambda \rightarrow \infty} \frac{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) = 1. \quad (47)$$

Proof: The proof is presented in Appendix H. ■

Lemma 4.11 unveils a similarity between Type-I and Type-II regularization as the Type-I measure $P_{\Theta|Z=z}^{(Q,\lambda)}$ in (9) exhibits a similar behavior [28].

In the asymptotic regime, when the regularization factor decreases to zero from the right, *i.e.*, $\lambda \rightarrow 0^+$, the Radon-Nikodym derivative $\frac{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}$ in (16) exhibits the following behavior.

Lemma 4.12. *If $Q(\mathcal{L}_{Q,z}^*) > 0$, with the set $\mathcal{L}_{Q,z}^*$ in (29), then the Radon-Nikodym derivative $\frac{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}$ in (16) satisfies for all $\theta \in \text{supp } Q$,*

$$\lim_{\lambda \rightarrow 0^+} \frac{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) = \frac{1}{Q(\mathcal{L}_{Q,z}^*)} \mathbb{1}_{\{\theta \in \mathcal{L}_{Q,z}^*\}}. \quad (48)$$

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

$$\lim_{\lambda \rightarrow 0^+} \frac{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) = \begin{cases} \infty & \text{if } \theta \in \mathcal{L}_{Q,z}^* \\ 0 & \text{otherwise.} \end{cases} \quad (49)$$

Proof: The proof is presented in Appendix I. ■

Lemma 4.12 highlights that in the asymptotic regime, when the regularization factor decreases to zero from the right, *i.e.*, $\lambda \rightarrow 0^+$, the value $\frac{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta)$ does not depend on the exact model θ but rather on whether $\theta \in \text{supp } Q \cap \mathcal{L}_{Q,z}^*$. In the case in which $\theta \in \text{supp } Q \cap \mathcal{L}_{Q,z}^*$, it holds that $\lim_{\lambda \rightarrow 0^+} \frac{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) > 0$. Otherwise, $\lim_{\lambda \rightarrow 0^+} \frac{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) = 0$. In the special case in which $\delta_{Q,z}^* = 0$, with $\delta_{Q,z}^*$ in (28), the set $\mathcal{L}_{Q,z}^*$ satisfies $\mathcal{L}_{Q,z}^* = \mathcal{T}(z)$, where $\mathcal{T}(z)$ is defined in (5). This implies a concentration of probability over $\mathcal{T}(z) \cap \text{supp } Q$, which establishes a connection with the ERM problem without regularization in (4).

Furthermore, in the asymptotic regime, when the regularization factor decreases to zero from the right, the solutions to the Type-I and Type-II ERM-RER problems exhibit the same asymptotic behavior, as shown in [28]. This aligns with

the observation that as λ decreases, the optimization problems in (7) and (11) exhibit a weaker relative entropy constraint. A stronger result follows from Lemma 4.12 and is presented in the following lemma.

Lemma 4.13. *The measure $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ in (16) and the set $\mathcal{L}_{Q,z}^*$ in (29) satisfy*

$$\lim_{\lambda \rightarrow 0^+} \bar{P}_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_{Q,z}^*) = 1. \quad (50)$$

Proof: The proof is presented in Appendix J. ■

Lemma 4.13 shows that indeed when the regularization factor approaches zero from the right, the probability of the measure $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ in (16), in the asymptotic regime, concentrates in the set of models that induce the minimum empirical risk in $\text{supp } Q$.

5 The Expected Empirical Risk

This section focuses on the expected empirical risk induced by the probability measure $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ in (16). Some important properties of the functional R_z in (6) and the value $R_z(\bar{P}_{\Theta|Z=z}^{(Q,\lambda)})$ are introduced.

The expected empirical risk can be calculated in terms of the regularization parameter λ and the function $\bar{K}_{Q,z}$ defined in (23b), as shown by the following lemma.

Lemma 5.1. *The probability measure $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ in (16) satisfies*

$$R_z(\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}) = \lambda - \bar{K}_{Q,z}(\lambda), \quad (51)$$

where the functional R_z , the function $\bar{K}_{Q,z}$ and the measure $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ are defined in (6), (23b), and (26), respectively.

Proof: The proof is presented in Appendix K. ■

Lemma 5.1 highlights that the function $r : (0, \infty) \rightarrow [0, \infty)$ such that $r(\lambda) = R_z(\bar{P}_{\Theta|Z=z}^{(Q,\lambda)})$, with Q and z fixed, inherits all properties of the function $\bar{K}_{Q,z}$ in (23b). The following lemma formalizes this observation.

Lemma 5.2. *The expected empirical risk $R_z(\bar{P}_{\Theta|Z=z}^{(Q,\lambda)})$, with the functional R_z in (6) and the measure $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ in (16), is nondecreasing with respect to λ . Moreover, it is strictly increasing if and only if the empirical risk function L_z in (3) is separable with respect to the probability measure Q .*

Proof: The proof is presented in Appendix L. ■

The following lemma highlights a connection existing between the expected empirical risks $R_z(Q)$ and $R_z(\bar{P}_{\Theta|Z=z}^{(Q,\lambda)})$; and the relative entropy $D(Q \| \bar{P}_{\Theta|Z=z}^{(Q,\lambda)})$.

Lemma 5.3. *The functional R_z defined in (6) and the measures Q and $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ in (16) satisfy*

$$R_z(Q) - R_z\left(\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}\right) \geq \lambda \left(\exp\left(D\left(Q\|\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}\right)\right) - 1 \right). \quad (52)$$

Proof: The proof is presented in Appendix M. ■

Note that $D\left(Q\|\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}\right) \geq 0$, which leads to the observation that

$$\left(\exp\left(D\left(Q\|\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}\right)\right) - 1 \right) \geq 0, \quad (53)$$

in (52). Hence, from Lemma 5.3, it follows that the solution to the Type-II ERM-RER problem induces an expected empirical risk that is smaller than the one induced by reference measure Q . This is formalized by the following corollary.

Corollary 5.4. *The probability measures Q and $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ in (16) satisfy*

$$R_z\left(\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}\right) \leq R_z(Q), \quad (54)$$

where the functional R_z is defined in (6) and equality holds if and only if the empirical risk function L_z in (3) is nonseparable.

The following lemma presents a lower bound and an upper bound on the expected empirical risk $R_z\left(\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}\right)$ in which the regularization parameter plays a central role.

Lemma 5.5. *The probability measure $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ in (16) satisfies*

$$\delta_{Q,z}^* \leq R_z\left(\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}\right) < \lambda + \delta_{Q,z}^*, \quad (55)$$

where the functional R_z is defined in (6) and $\delta_{Q,z}^*$ is defined in (28). Moreover, equality holds if and only if the empirical risk function L_z in (3) is nonseparable.

Proof: The proof is presented in Appendix N. ■

The bounds presented in Lemma 5.5 highlight that the regularization parameter λ in (11) governs the increase of the expected empirical risk $R_z\left(\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}\right)$ with respect to its minimum, i.e., $\delta_{Q,z}^*$ in (28). Moreover, the lower bound is tight for the probability measure $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ in (16) in the asymptotic regime when λ decreases to zero from right, as shown hereunder.

Lemma 5.6. *The probability measure $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ in (16) satisfies*

$$\lim_{\lambda \rightarrow 0^+} R_z\left(\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}\right) = \delta_{Q,z}^*, \quad (56)$$

where $\delta_{Q,z}^*$ is defined in (28) and the functional R_z is defined in (6).

Proof: From Lemma 5.1, it holds that

$$\lim_{\lambda \rightarrow 0^+} R_z \left(\bar{P}_{\Theta|Z=z}^{(Q,\lambda)} \right) = \lim_{\lambda \rightarrow 0^+} \lambda - \lim_{\lambda \rightarrow 0^+} \bar{K}_{Q,z}(\lambda) \quad (57a)$$

$$= - \lim_{\lambda \rightarrow 0^+} \bar{K}_{Q,z}(\lambda) \quad (57b)$$

$$= \delta_{Q,z}^*, \quad (57c)$$

where equality (57c) follows from Lemma 4.8. This completes the proof. ■

Finally, note that the functional R_z in (6) is nonnegative. This observation together with Lemma 5.1 lead to a new property for the function $\bar{K}_{Q,z}$ in (23b), which is stated by the following corollary

Corollary 5.7. *The function $\bar{K}_{Q,z}$ in (23b) satisfies, for all $t \in (0, \infty)$,*

$$\bar{K}_{Q,z}(t) \leq t. \quad (58)$$

6 (δ, ϵ) -Optimality

This section presents a PAC guarantee of optimality for models sampled from the probability measure $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ in (16) within the setting of the Type-II ERM-RER problem in (11). Such guarantee is presented using the notion of (δ, ϵ) -optimality introduced in [28].

Definition 6.1 ((δ, ϵ) -Optimality). *Given a pair of positive reals (δ, ϵ) , with $\epsilon < 1$, the probability measure $P \in \Delta(\mathcal{M})$ is said to be (δ, ϵ) -optimal if the set $\mathcal{L}_z(\delta)$ in (27) satisfies*

$$P(\mathcal{L}_z(\delta)) > 1 - \epsilon. \quad (59)$$

The (δ, ϵ) -optimality guarantee ensures that with probability at least $1 - \epsilon$, sampling models from $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ in (16) yields models that induce empirical risks not greater than δ . The following theorem presents an (δ, ϵ) -optimality guarantee for the Type-II ERM-RER solution.

Theorem 6.1. *For all $(\delta, \epsilon) \in (\delta_{Q,z}^*, \infty) \times (0, 1)$, with $\delta_{Q,z}^*$ in (28), there always exists a $\lambda \in (0, \infty)$, such that the probability measure $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ in (16) is (δ, ϵ) -optimal.*

Proof: The proof is presented in Appendix O. ■

7 Interplay Between the Relative Entropy Asymmetry and the Empirical Risk

This section presents a connection between the Type-I ERM-RER in (7) and Type-II ERM-RER problems in (11) established via a transformation of the

empirical risk function. The transformation used is termed *log-empirical risk* and is defined hereunder.

Definition 7.1 (Log-Empirical Risk). *Consider the empirical risk function L_z in (3) and let the function $V_{z,\lambda} : \mathcal{M} \rightarrow \mathbb{R}$, with Q and λ in (11), be such that*

$$V_{z,\lambda}(\theta) \triangleq \log(\bar{K}_{Q,z}(\lambda) + L_z(\theta)), \quad (60)$$

with the function $\bar{K}_{Q,z}$ defined in (23b). The value $V_{z,\lambda}(\theta)$ is said to be the log-empirical risk induced by the model $\theta \in \mathcal{M}$.

The notion of log-empirical risk in (60) leads to the *expected log-empirical risk*, which is denoted as follows.

Definition 7.2 (Expected Log-Empirical Risk). *Consider the log-empirical risk function $V_{z,\lambda}$ in (60) and let the functional $\bar{R}_{z,Q,\lambda} : \Delta(\mathcal{M}) \rightarrow \mathbb{R}$ be such that for all probability measures $P \in \Delta(\mathcal{M})$, it holds that*

$$\bar{R}_{z,Q,\lambda}(P) \triangleq \int V_{z,\lambda}(\theta) dP(\theta). \quad (61)$$

The value $\bar{R}_{z,Q,\lambda}(P)$ is the expected log-empirical risk induced by the measure P .

7.1 The Connection between Type-I and Type-II

Using the objects defined above, a new Type-I ERM-RER problem is presented:

$$\min_{P \in \Delta_Q(\mathcal{M})} \bar{R}_{z,Q,\lambda}(P) + D(P \| Q), \quad (62)$$

with λ and Q being problems of the Type-I and Type-II ERM-RER problems in (7) and (11). The main result of this section is presented in the following theorem.

Theorem 7.1. *The solution to the optimization problem in (62) is unique and identical to the probability measure $\hat{P}_{\Theta|Z=z}^{(Q,\lambda)}$ in (16).*

Proof: Denote by $\hat{P}_{\Theta|Z=z}^{(Q,\lambda)}$ the solution to the optimization problem in (62). From Lemma 3.1, for all $\theta \in \text{supp } Q$, it follows that

$$\frac{d\hat{P}_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) = \frac{\exp(-V_{z,\lambda}(\theta))}{\int \exp(-V_{z,\lambda}(\nu)) dQ(\nu)} \quad (63a)$$

$$= \frac{\exp\left(\log\left(\frac{1}{L_z(\theta) + \bar{K}_{Q,z}(\lambda)}\right)\right)}{\int \exp\left(\log\left(\frac{1}{L_z(\nu) + \bar{K}_{Q,z}(\lambda)}\right)\right) dQ(\nu)} \quad (63b)$$

$$= \frac{\left(\int \frac{1}{L_z(\nu) + \bar{K}_{Q,z}(\lambda)} dQ(\nu)\right)^{-1}}{L_z(\theta) + \bar{K}_{Q,z}(\lambda)} \quad (63c)$$

$$= \frac{\lambda}{\mathbf{L}_z(\boldsymbol{\theta}) + \bar{K}_{Q,z}(\lambda)} \quad (63d)$$

$$= \frac{d\bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}}{dQ}(\boldsymbol{\theta}), \quad (63e)$$

where the equality in (63b) follows from the definition of log-empirical risk in (60); the equality in (63d) follows from (23b) and (15b); and the equality in (63e) follows from Theorem 4.1, which completes the proof. ■

Theorem 7.1 establishes an equivalence between the regularization of Type-I and Type-II. More specifically, Theorem 7.1 highlights that by transforming the empirical risk function \mathbf{L}_z in (60) into the log-empirical risk $\mathbf{V}_{z,\lambda}$ in (60), the Type-II ERM-RER problem in (11) can be solved by solving the Type-I ERM-RER problem in (62). In view of this, it is interesting to note that Type-I regularization forces the support of the solution to be included into the support of the reference measure as a consequence of the optimization domain $\triangle_Q(\mathcal{M})$. Hence, $\text{supp } \bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)} \subseteq \text{supp } Q$. Moreover, as a consequence of the solution, all the models in the support of the reference measure are in the support of the solution. That is, $\text{supp } Q \subseteq \text{supp } \bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}$. Alternatively, Type-II regularization forces the support of the solution to be included into the support of the reference measure as a consequence of the optimization domain $\nabla_Q(\mathcal{M})$. Hence, $\text{supp } Q \subseteq \text{supp } \bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}$. Moreover, as a consequence of the solution, all the models in the support of the solution are in the support of the reference measure. That is, $\text{supp } \bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)} \subseteq \text{supp } Q$. Hence, $\text{supp } \bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)} = \text{supp } Q$, regardless of whether the regularization is of Type-I or Type-II.

Finally, the Type-I - Type-II relation can be used to establish an equality involving the relative entropies $D\left(Q \parallel \bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}\right)$ and $D\left(\bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)} \parallel Q\right)$; and the expected log-empirical risks $\bar{R}_{z,Q,\lambda}\left(\bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}\right)$ and $\bar{R}_{z,Q,\lambda}(Q)$, as shown hereunder.

Lemma 7.1. *The functional $\bar{R}_{z,Q,\lambda}$ in (61) and the probability measures $\bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}$ and Q in (16) satisfy*

$$\bar{R}_{z,Q,\lambda}\left(\bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}\right) + D\left(\bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)} \parallel Q\right) = \log(\lambda), \quad (64)$$

and

$$\bar{R}_{z,Q,\lambda}(Q) - D\left(Q \parallel \bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}\right) = \log(\lambda). \quad (65)$$

Proof: The proof is presented in Appendix P. ■

7.2 Sensitivity of the Log-Empirical Risk

The sensitivity of the expected empirical risk, as presented in [28, Definition 10.1], is defined as follows.

Definition 7.3 (Sensitivity of the Expected Empirical Risk). *Consider the functional R_z defined in (6) and let $S_{Q,\lambda} : (\mathcal{X} \times \mathcal{Y})^n \times \Delta_Q(\mathcal{M}) \rightarrow \mathbb{R}$ be a functional such that*

$$S_{Q,\lambda}(z, P) = R_z(P) - R_z\left(P_{\Theta|Z=z}^{(Q,\lambda)}\right), \quad (66)$$

where the probability measure $P_{\Theta|Z=z}^{(Q,\lambda)}$ is in (9). The sensitivity of the expected empirical risk due to a deviation from $P_{\Theta|Z=z}^{(Q,\lambda)}$ to P is $S_{Q,\lambda}(z, P)$.

Following the same idea, the sensitivity of the expected log-empirical risk is defined as follows.

Definition 7.4 (Sensitivity of the Expected Log-Empirical Risk). *Consider the functional $\bar{R}_{z,Q,\lambda}$ in (61) and let $\bar{S}_{Q,\lambda} : (\mathcal{X} \times \mathcal{Y})^n \times \nabla_Q(\mathcal{M}) \rightarrow \mathbb{R}$ be a functional such that*

$$\bar{S}_{Q,\lambda}(z, P) = \bar{R}_{z,Q,\lambda}(P) - \bar{R}_{z,Q,\lambda}\left(\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}\right), \quad (67)$$

where the probability measure $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ is in (16). The sensitivity of the expected log-empirical risk due to a deviation from $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ to P is $\bar{S}_{Q,\lambda}(z, P)$.

The sensitivity of the expected log-empirical risk due to a deviation from $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ to P exhibits the following closed-form expression.

Lemma 7.2. *The sensitivity $\bar{S}_{Q,\lambda}$ in (67) satisfies for all probability measures $P \in \bigcirc_Q(\mathcal{M})$ that*

$$\bar{S}_{Q,\lambda}(z, P) = D\left(P \parallel \bar{P}_{\Theta|Z=z}^{(Q,\lambda)}\right) - D(P \parallel Q) + D\left(\bar{P}_{\Theta|Z=z}^{(Q,\lambda)} \parallel Q\right), \quad (68)$$

where the probability measures Q and $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ are defined in (16).

Proof: The proof is presented in Appendix Q. ■

An interesting interpretation of Lemma 7.2 follows from rewriting (68) using the objective function of the Type-I ERM-RER problem in (62) as follows:

$$\begin{aligned} D\left(P \parallel \bar{P}_{\Theta|Z=z}^{(Q,\lambda)}\right) &= \bar{R}_{z,Q,\lambda}(P) + D(P \parallel Q) - \bar{R}_{z,Q,\lambda}\left(\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}\right) \\ &\quad - D\left(\bar{P}_{\Theta|Z=z}^{(Q,\lambda)} \parallel Q\right). \end{aligned} \quad (69)$$

That is, the relative entropy $D\left(P \parallel \bar{P}_{\Theta|Z=z}^{(Q,\lambda)}\right)$ represents the variation of the objective function of the Type-I ERM-RER problem in (62) due to a deviation from the solution $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ to an alternative probability measure P .

In Lemma 7.2, when P is chosen to be identical to the reference measure Q , it follows that

$$\bar{S}_{Q,\lambda}(z, Q) = D\left(Q \parallel \bar{P}_{\Theta|Z=z}^{(Q,\lambda)}\right) + D\left(\bar{P}_{\Theta|Z=z}^{(Q,\lambda)} \parallel Q\right), \quad (70)$$

where the right-hand side is a Jeffrey's divergence, also known as the symmetrized Kullback-Liebler divergence, between the measures $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ and Q . Furthermore, by observing that $D(\bar{P}_{\Theta|Z=z}^{(Q,\lambda)} \| Q) \geq 0$, and $D(Q \| \bar{P}_{\Theta|Z=z}^{(Q,\lambda)}) \geq 0$, Lemma 7.2 leads to the following corollary.

Corollary 7.3. *The probability measures Q and $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ in (16) satisfy*

$$\bar{R}_{z,Q,\lambda}(\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}) \leq \bar{R}_{z,Q,\lambda}(Q), \quad (71)$$

where the functional $\bar{R}_{z,Q,\lambda}$ is defined in (61).

7.3 Type-I and Type-II Optimal Measures

The probability measures $P_{\Theta|Z=z}^{(Q,\alpha)}$ and $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ in (9) and in (16), respectively, exhibit the following property.

Lemma 7.4. *The probability measures $P_{\Theta|Z=z}^{(Q,\alpha)}$ and $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ in (9) and in (16), respectively, satisfy*

$$D(P_{\Theta|Z=z}^{(Q,\lambda)} \| Q) - D(\bar{P}_{\Theta|Z=z}^{(Q,\lambda)} \| Q) = \log(\alpha) + K_{Q,z}\left(-\frac{1}{\lambda}\right), \quad (72)$$

where the function $K_{Q,z}$ is defined in (10).

Proof: The proof is presented in Appendix R. ■

Finally, two important properties of the Type-I and Type-II optimal measures are presented. The first one quantifies the variation of the expected log-empirical risk due to a deviation from the probability measure $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ in (9) to the probability measure $P_{\Theta|Z=z}^{(Q,\alpha)}$ in (16) via the sensitivity $\bar{S}_{Q,\alpha}\left(z, P_{\Theta|Z=z}^{(Q,\lambda)}\right)$. The second result quantifies the variation of the expected empirical risk due to a deviation from the probability measure $P_{\Theta|Z=z}^{(Q,\alpha)}$ in (16) to the probability measure $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ in (9) via the sensitivity $S_{Q,\alpha}\left(z, \bar{P}_{\Theta|Z=z}^{(Q,\lambda)}\right)$.

Corollary 7.5. *The probability measures $P_{\Theta|Z=z}^{(Q,\alpha)}$ and $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ in (9) and in (16), respectively, satisfy*

$$\bar{S}_{Q,\alpha}\left(z, P_{\Theta|Z=z}^{(Q,\lambda)}\right) = D\left(P_{\Theta|Z=z}^{(Q,\alpha)} \| \bar{P}_{\Theta|Z=z}^{(Q,\lambda)}\right) - \left(\log(\alpha) + K_{Q,z}\left(-\frac{1}{\lambda}\right)\right) \quad (73)$$

and

$$\frac{1}{\lambda} S_{Q,\lambda}\left(z, \bar{P}_{\Theta|Z=z}^{(Q,\lambda)}\right) = D\left(\bar{P}_{\Theta|Z=z}^{(Q,\lambda)} \| P_{\Theta|Z=z}^{(Q,\alpha)}\right) + \log(\alpha) + K_{Q,z}\left(-\frac{1}{\lambda}\right), \quad (74)$$

where the functionals $S_{Q,\lambda}$ and $\bar{S}_{Q,\alpha}$ are respectively defined in (66) and in (67); and the function $K_{Q,z}$ is defined in (10).

8 Final Remarks

This work has introduced the Type-II ERM-RER problem and has presented its solution through Theorem 4.1. The solution highlights that regardless of the direction in which relative entropy is used as a regularizer, the models that are considered by the solution are necessarily in the support of the reference measure. In that sense, the restriction over the models introduced by the reference measure cannot be bypassed by the training data when relative entropy is used as the regularizer. This limitation has been shown to be a consequence of the equivalence that can be established between Type-I and Type-II regularization. Remarkably, the direction of the relative entropy regularizer can be switched by a logarithmic transformation of the risk. The mutual absolute continuity of both Type-I and Type-II ERM-RER solutions relative to the reference measure can be understood in light of the equivalence between both types of regularization. The analytical results have also enabled us to provide an operationally meaningful characterization of the expected empirical risk induced by the Type-II solution in terms of the regularization parameters. This, in turn, reduces the computational burden of bounding the expected empirical risk. Moreover, the insight provided by the bounds on the expected empirical risk can be distilled into guidelines for selecting the regularization parameter.

A Proof of Lemma 4.1

Proof: The optimization problem in (11) can be re-written in terms of the Radon-Nikodym derivative of the optimization measure P with respect to the measure Q , which yields:

$$\min_{P \in \mathcal{O}_Q(\mathcal{M})} \int \mathbf{L}_z(\boldsymbol{\theta}) \frac{dP}{dQ}(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) + \lambda \int \frac{dQ}{dP}(\boldsymbol{\theta}) \log\left(\frac{dQ}{dP}(\boldsymbol{\theta})\right) dP(\boldsymbol{\theta}), \quad (75a)$$

$$\text{s.t.} \quad \int \frac{dP}{dQ}(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) = 1. \quad (75b)$$

The remainder of the proof focuses on the problem in which the optimization is over the function $\frac{dP}{dQ} : \mathcal{M} \rightarrow \mathbb{R}$, instead of optimizing the measure P . This is due to the fact that for all $P \in \mathcal{O}_Q(\mathcal{M})$, the Radon-Nikodym derivative $\frac{dP}{dQ}$ is unique up to sets of zero measure with respect to Q . Let \mathcal{M} be the set of measurable functions $\mathcal{M} \rightarrow \mathbb{R}$ with respect to the measurable spaces $(\mathcal{M}, \mathcal{F})$ and $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ that are absolutely integrable with respect to Q . That is, for all $\hat{g} \in \mathcal{M}$, it holds that

$$\int |\hat{g}(\boldsymbol{\theta})| dQ(\boldsymbol{\theta}) < \infty. \quad (76)$$

Hence, the optimization problem of interest is:

$$\min_{g \in \mathcal{M}} \int \mathbf{L}_z(\boldsymbol{\theta}) g(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) - \lambda \int \log(g(\boldsymbol{\theta})) dQ(\boldsymbol{\theta}) \quad (77a)$$

$$\text{s.t.} \quad \int g(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) = 1. \quad (77b)$$

Let the Lagrangian of the optimization problem in (77) be $L : \mathcal{M} \times \mathbb{R} \rightarrow \mathbb{R}$ such that

$$\begin{aligned} L(g, \beta) &= \int \mathbf{L}_z(\boldsymbol{\theta}) g(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) - \lambda \int \log(g(\boldsymbol{\theta})) dQ(\boldsymbol{\theta}) \\ &\quad + \beta \left(\int g(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) - 1 \right) \end{aligned} \quad (78)$$

$$= \int \left(g(\boldsymbol{\theta}) (\mathbf{L}_z(\boldsymbol{\theta}) + \beta) - \lambda \log(g(\boldsymbol{\theta})) \right) dQ(\boldsymbol{\theta}) - \beta, \quad (79)$$

where β is a real that acts as a Lagrange multiplier due to the constraint (77b). Let $\hat{g} : \mathcal{M} \rightarrow \mathbb{R}$ be a function in \mathcal{M} . The Gateaux differential of the functional L in (78) at $(g, \beta) \in \mathcal{M} \times \mathbb{R}$ in the direction of \hat{g} , if it exists, is

$$\partial L(g, \beta; \hat{g}) \triangleq \left. \frac{d}{d\gamma} L(g + \gamma \hat{g}, \beta) \right|_{\gamma=0}. \quad (80)$$

The proof continues under the assumption that the function g and \hat{g} are such that the Gateaux differential in (80) exists. Under such an assumption, let

the function $r : \mathbb{R} \rightarrow \mathbb{R}$ satisfy for all $\gamma \in (-\epsilon, \epsilon)$, with ϵ arbitrarily small, that

$$r(\gamma) = L(g + \gamma \hat{g}, \beta) \quad (81)$$

$$= \int \mathbf{L}_z(\boldsymbol{\theta})(g(\boldsymbol{\theta}) + \gamma \hat{g}(\boldsymbol{\theta})) \, dQ(\boldsymbol{\theta}) - \lambda \int \log(g(\boldsymbol{\theta}) + \gamma \hat{g}(\boldsymbol{\theta})) \, dQ(\boldsymbol{\theta}) \\ + \beta \left(\int (g(\boldsymbol{\theta}) + \gamma \hat{g}(\boldsymbol{\theta})) \, dQ(\boldsymbol{\theta}) - 1 \right) \quad (82)$$

$$= -\beta + \int g(\boldsymbol{\theta})(\mathbf{L}_z(\boldsymbol{\theta}) + \beta) \, dQ(\boldsymbol{\theta}) + \gamma \int \hat{g}(\boldsymbol{\theta})(\mathbf{L}_z(\boldsymbol{\theta}) + \beta) \, dQ(\boldsymbol{\theta}) \\ + \lambda \int \log(g(\boldsymbol{\theta}) + \gamma \hat{g}(\boldsymbol{\theta})) \, dQ(\boldsymbol{\theta}), \quad (83)$$

where the last equality is simply an algebraic re-arrangement of terms. From the assumption that the function g and \hat{g} are such that the Gateaux differential in (80) exists, it follows that the function r in (81) is differentiable at zero. Note that the first two terms in (83) are independent of γ ; the third term is linear with γ ; and the fourth term can be written using the function $\hat{r} : \mathbb{R} \rightarrow \mathbb{R}$ such that for all $\gamma \in (-\epsilon, \epsilon)$, with ϵ arbitrarily small, satisfies

$$\hat{r}(\gamma) = -\lambda \int \log(g(\boldsymbol{\theta}) + \gamma \hat{g}(\boldsymbol{\theta})) \, dQ(\boldsymbol{\theta}). \quad (84)$$

Under the same assumption, it follows that the function \hat{r} in (84) is differentiable at zero. That is, the limit

$$\lim_{\delta \rightarrow 0} \frac{1}{\delta} (\hat{r}(\gamma + \delta) - \hat{r}(\gamma)) \quad (85)$$

exists for all $\gamma \in (-\epsilon, \epsilon)$, with ϵ arbitrarily small. The proof of the existence of such a limit relies on the fact that \log is strictly concave and continuous. This implies that $-\log$ is also Lipschitz continuous, which implies that for all $\boldsymbol{\theta} \in \mathcal{M}$ and for all $\gamma \in (-\epsilon, \epsilon)$, with $\epsilon > 0$ arbitrarily small, it holds that

$$|\log(g(\boldsymbol{\theta}) + \gamma \hat{g}(\boldsymbol{\theta})) - \log(g(\boldsymbol{\theta}) + (\gamma + \delta) \hat{g}(\boldsymbol{\theta}))| \leq c |\hat{g}(\boldsymbol{\theta}) \delta|, \quad (86)$$

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

$$\left| \frac{\log(g(\boldsymbol{\theta}) + \gamma \hat{g}(\boldsymbol{\theta})) - \log(g(\boldsymbol{\theta}) + (\gamma + \delta) \hat{g}(\boldsymbol{\theta}))}{\delta} \right| \leq c |\hat{g}(\boldsymbol{\theta})|. \quad (87)$$

Using these arguments, the limit in (85) satisfies for all $\gamma \in (-\epsilon, \epsilon)$, with $\epsilon > 0$ arbitrarily small, that

$$\lim_{\delta \rightarrow 0} \frac{1}{\delta} (\hat{r}(\gamma + \delta) - \hat{r}(\gamma)) = \lambda \lim_{\delta \rightarrow 0} \frac{1}{\delta} \left(\int \log(g(\boldsymbol{\theta}) + \gamma \hat{g}(\boldsymbol{\theta})) \, dQ(\boldsymbol{\theta}) \right. \\ \left. - \int \log(g(\boldsymbol{\theta}) + (\gamma + \delta) \hat{g}(\boldsymbol{\theta})) \, dQ(\boldsymbol{\theta}) \right) \quad (88)$$

$$\begin{aligned}
&= \lambda \lim_{\delta \rightarrow 0} \int \frac{\log(g(\boldsymbol{\theta}) + \gamma \hat{g}(\boldsymbol{\theta}))}{\delta} \\
&\quad - \frac{\log(g(\boldsymbol{\theta}) + (\gamma + \delta) \hat{g}(\boldsymbol{\theta}))}{\delta} dQ(\boldsymbol{\theta}) \quad (89)
\end{aligned}$$

$$= \lambda \int \frac{d}{d\gamma} \log(g(\boldsymbol{\theta}) + \gamma \hat{g}(\boldsymbol{\theta})) dQ(\boldsymbol{\theta}) \quad (90)$$

$$< \infty, \quad (91)$$

where both the equality in (90) and the inequality in (91) follow from noticing that the conditions for the dominated convergence theorem hold [45, Theorem 1.6.9], namely:

- For all $\gamma \in (-\epsilon, \epsilon)$, with $\epsilon > 0$, the inequality in (87) holds;
- The function \hat{g} in (87) satisfies the inequality in (76); and
- For all $\boldsymbol{\theta} \in \mathcal{M}$ and for all $\gamma \in (-\epsilon, \epsilon)$, with $\epsilon > 0$ arbitrarily small, it holds that

$$\begin{aligned}
&\lim_{\delta \rightarrow 0} \frac{\log(g(\boldsymbol{\theta}) + (\gamma + \delta) \hat{g}(\boldsymbol{\theta})) - \log(g(\boldsymbol{\theta}) + \gamma \hat{g}(\boldsymbol{\theta}))}{\delta} \\
&= \frac{d}{d\gamma} \log(g(\boldsymbol{\theta}) + \gamma \hat{g}(\boldsymbol{\theta})) \quad (92)
\end{aligned}$$

$$= \frac{\hat{g}(\boldsymbol{\theta})}{(g(\boldsymbol{\theta}) + \gamma \hat{g}(\boldsymbol{\theta}))}, \quad (93)$$

which follows from the fact that \log is differentiable.

Hence, the derivative of the real function r in (83) is

$$\frac{d}{d\gamma} r(\gamma) = \int \hat{g}(\boldsymbol{\theta}) (\mathbf{L}_{\mathbf{z}}(\boldsymbol{\theta}) + \beta) dQ(\boldsymbol{\theta}) - \lambda \int \frac{\hat{g}(\boldsymbol{\theta})}{(g(\boldsymbol{\theta}) + \gamma \hat{g}(\boldsymbol{\theta}))} dQ(\boldsymbol{\theta}) \quad (94)$$

$$= \int \hat{g}(\boldsymbol{\theta}) \left(\mathbf{L}_{\mathbf{z}}(\boldsymbol{\theta}) + \beta - \frac{\lambda}{(g(\boldsymbol{\theta}) + \gamma \hat{g}(\boldsymbol{\theta}))} \right) dQ(\boldsymbol{\theta}). \quad (95)$$

From (80) and (95), it follows that

$$\partial L(g, \beta; \hat{g}) = \int \hat{g}(\boldsymbol{\theta}) \left(\mathbf{L}_{\mathbf{z}}(\boldsymbol{\theta}) + \beta - \frac{\lambda}{g(\boldsymbol{\theta})} \right) dQ(\boldsymbol{\theta}). \quad (96)$$

The relevance of the Gateaux differential in (96) stems from [46, Theorem 1, page 178], which unveils the fact that a necessary condition for the functional L in (78) to have a stationary point at $\left(\frac{dP_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q, \lambda)}}{dQ}, \beta \right) \in \mathcal{M} \times [0, \infty)$ is that for all functions $\hat{g} \in \mathcal{M}$,

$$\partial L \left(\frac{dP_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q, \lambda)}}{dQ}, \beta; \hat{g} \right) = 0. \quad (97)$$

From (96) and (97), it follows that $\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}$ must satisfy for all functions \hat{g} in \mathcal{M} that

$$\int \hat{g}(\boldsymbol{\theta}) \left(\mathbb{L}_z(\boldsymbol{\theta}) + \beta - \lambda \left(\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\boldsymbol{\theta}) \right)^{-1} \right) dQ(\boldsymbol{\theta}) = 0. \quad (98)$$

This implies that for all $\boldsymbol{\theta} \in \text{supp } Q$,

$$\mathbb{L}_z(\boldsymbol{\theta}) + \beta - \lambda \left(\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\boldsymbol{\theta}) \right)^{-1} = 0, \quad (99)$$

and thus,

$$\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\boldsymbol{\theta}) = \frac{\lambda}{\beta + \mathbb{L}_z(\boldsymbol{\theta})}, \quad (100)$$

where β is chosen to satisfy (77b) and guarantee that for all $\boldsymbol{\theta} \in \text{supp } Q$, it holds that $\frac{dP_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\boldsymbol{\theta}) \in (0, \infty)$. That is,

$$\beta \in \left\{ t \in \mathbb{R} : \forall \boldsymbol{\theta} \in \text{supp } Q, 0 < \frac{\lambda}{t + \mathbb{L}_z(\boldsymbol{\theta})} \right\}, \text{ and} \quad (101)$$

$$1 = \int \frac{\lambda}{\mathbb{L}_z(\boldsymbol{\theta}) + \beta} dQ(\boldsymbol{\theta}). \quad (102)$$

which is an assumption of the theorem.

The proof continues by verifying that the measure $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ that satisfies (100) is the unique solution to the optimization problem in (75). Such verification is done by showing that the objective function in (75) is strictly convex with the optimization variable. Let P_1 and P_2 be two different probability measures in $(\mathcal{M}, \mathcal{F})$ and let α be in $(0, 1)$. Hence,

$$\mathbb{R}_z(\alpha P_1 + (1 - \alpha)P_2) + \lambda \mathbb{D}(\alpha P_1 + (1 - \alpha)P_2 \| Q) \quad (103)$$

$$= \mathbb{R}_z(\alpha P_1) + \mathbb{R}_z((1 - \alpha)P_2) + \lambda \mathbb{D}(\alpha P_1 + (1 - \alpha)P_2 \| Q) \quad (104)$$

$$> \mathbb{R}_z(\alpha P_1) + \mathbb{R}_z((1 - \alpha)P_2) + \lambda \mathbb{D}(\alpha P_1 \| Q) + \lambda \mathbb{D}((1 - \alpha)P_2 \| Q) \quad (105)$$

$$(106)$$

where the functional \mathbb{R}_z is defined in (6). The equality above follows from the properties of the Lebesgue integral, while the inequality follows from [28, Theorem 2.2., page]. This proves that the solution is unique due to the strict concavity of the objective function, which completes the proof. ■

B Proof of Lemma 4.2

Proof: Given a probability measure $V \in \nabla_Q(\mathcal{M})$, with $\nabla_Q(\mathcal{M})$ in (12), let V_0 and V_1 be two probability measures on the measurable space $(\mathcal{M}, \mathcal{F})$ such that for all $\mathcal{A} \in \mathcal{F}$, it holds that

$$V_0(\mathcal{A}) = \frac{V(\mathcal{A} \setminus \text{supp } Q)}{V(\mathcal{M} \setminus \text{supp } Q)}, \quad (107a)$$

and

$$V_1(\mathcal{A}) = \frac{V(\mathcal{A} \cap \text{supp } Q)}{V(\mathcal{M} \cap \text{supp } Q)}. \quad (107b)$$

Let the real value α be

$$\alpha = V(\mathcal{M} \cap \text{supp } Q). \quad (108)$$

Hence, for all $\mathcal{A} \in \mathcal{F}$ the measure V satisfies that

$$V(\mathcal{A}) = (1 - \alpha)V_0(\mathcal{A}) + \alpha V_1(\mathcal{A}). \quad (109)$$

Moreover, from (109) it holds that: *i)* If $V(\mathcal{A}) = 0$, then $V_0(\mathcal{A}) = 0$, which implies that V_0 is absolutely continuous with respect to V ; *ii)* If $V(\mathcal{A}) = 0$, then $V_1(\mathcal{A}) = 0$, which implies that V_1 is absolutely continuous with respect to V . Furthermore, from the definition of $\nabla_Q(\mathcal{M})$ in (12) the probability measure V is absolutely continuous with respect to Q . Hence, for all $\mathcal{A} \in \mathcal{F}$, it follows that

$$Q(\mathcal{A}) = \int_{\mathcal{A}} dQ(\boldsymbol{\theta}) \quad (110)$$

$$= \int_{\mathcal{A}} \frac{dQ}{dV}(\boldsymbol{\theta}) dV(\boldsymbol{\theta}) \quad (111)$$

$$= \int_{\mathcal{A}} \frac{dQ}{dV}(\boldsymbol{\theta}) d((1 - \alpha)V_0 + \alpha V_1)(\boldsymbol{\theta}) \quad (112)$$

$$= (1 - \alpha) \int_{\mathcal{A}} \frac{dQ}{dV}(\boldsymbol{\theta}) dV_0(\boldsymbol{\theta}) + \alpha \int_{\mathcal{A}} \frac{dQ}{dV}(\boldsymbol{\theta}) dV_1(\boldsymbol{\theta}) \quad (113)$$

$$= \int_{\mathcal{A}} \alpha \frac{dQ}{dV}(\boldsymbol{\theta}) dV_1(\boldsymbol{\theta}). \quad (114)$$

Hence, from (114) and the Radon-Nikodym Theorem in [45, Theorem 2.2.1, page 65] the probability measure Q is absolutely continuous with respect to V_1 . This implies that for all $\mathcal{A} \in \mathcal{F}$, it holds that

$$Q(\mathcal{A}) = \int_{\mathcal{A}} \frac{dQ}{dV_1}(\boldsymbol{\theta}) dV_1(\boldsymbol{\theta}), \quad (115)$$

where, for all $\boldsymbol{\theta} \in \text{supp } V$,

$$\frac{dQ}{dV_1}(\boldsymbol{\theta}) = \alpha \frac{dQ}{dV}(\boldsymbol{\theta}). \quad (116)$$

From (116), the following holds:

$$D(Q\|V) = \int \log\left(\frac{dQ}{dV}(\boldsymbol{\theta})\right) dQ(\boldsymbol{\theta}) \quad (117a)$$

$$= \int \log\left(\frac{1}{\alpha} \frac{dQ}{dV_1}(\boldsymbol{\theta})\right) dQ(\boldsymbol{\theta}) \quad (117b)$$

$$= \int \log\left(\frac{dQ}{dV_1}(\boldsymbol{\theta})\right) dQ(\boldsymbol{\theta}) - \int \log(\alpha) dQ(\boldsymbol{\theta}) \quad (117c)$$

$$= D(Q\|V_1) - \log(\alpha). \quad (117d)$$

From (117) it follows that

$$R_{\mathbf{z}}(V) + \lambda D(Q\|V) = R_{\mathbf{z}}((1-\alpha)V_0 + \alpha V_1) + \lambda D(Q\|V_1) - \lambda \log(\alpha) \quad (118a)$$

$$= (1-\alpha)R_{\mathbf{z}}(V_0) + \alpha R_{\mathbf{z}}(V_1) + \lambda D(Q\|V_1) - \lambda \log(\alpha) \quad (118b)$$

$$\geq \alpha R_{\mathbf{z}}(V_1) + \lambda D(Q\|V_1) - \lambda \log(\alpha), \quad (118c)$$

with equality if and only if $\alpha = 1$, which implies that for all $\mathcal{A} \in \mathcal{F}$, it holds that

$$V(\mathcal{A}) = V_1(\mathcal{A}) \quad (119a)$$

$$= V(\mathcal{A} \cap \text{supp } Q), \quad (119b)$$

where the equality in (119b) follows from (107b). This implies that the equality in (118c) holds if and only if

$$\text{supp } Q = \text{supp } V, \quad (120)$$

which implies that the equality in (118c) holds if and only if the measure V is mutually absolutely continuous with respect to the reference measure Q . Finally, the above leads to

$$\min_{P \in \nabla_Q(\mathcal{M}) \setminus \bigcirc_Q(\mathcal{M})} R_{\mathbf{z}}(P) + \lambda D(Q\|P) > \min_{P \in \nabla_Q(\mathcal{M})} R_{\mathbf{z}}(P) + \lambda D(Q\|P), \quad (121)$$

which completes the proof. \blacksquare

C Proof of Lemma 4.4

Proof: The properties of the function $\bar{K}_{Q,\mathbf{z}}$ in (23b), for which an explicit expression is unknown, are proven by studying the functional inverse and the continuous inverse theorem [47, Theorem 5.6]. Hence, the proof is divided into two parts. The first part introduces the functional inverse of the function $\bar{K}_{Q,\mathbf{z}}$ in (23b) and properties. The second part proves the continuity of the function

$\bar{K}_{Q,\mathbf{z}}$ in (23b) is proved using its functional inverse via the continuous inverse theorem [47, Theorem 5.6]

The first part is as follows. For $\bar{K}_{Q,\mathbf{z}}$ defined in (23), assume that $t \in (0, \infty)$ and $\gamma \in \mathcal{A}$, with \mathcal{A} defined in (23a) satisfy that

$$\bar{K}_{Q,\mathbf{z}}(t) = \gamma, \quad (122)$$

which implies that

$$1 = \int \frac{d\bar{P}_{\Theta|\mathbf{Z}=\mathbf{z}}^{(Q,t)}(\boldsymbol{\theta})}{dQ} dQ(\boldsymbol{\theta}) \quad (123a)$$

$$= \int \frac{t}{\mathbb{L}_{\mathbf{z}}(\boldsymbol{\theta}) + \gamma} dQ(\boldsymbol{\theta}). \quad (123b)$$

Let the function $\bar{K}_{Q,\mathbf{z}}^{-1} : \mathcal{A} \rightarrow (0, \infty)$ be the functional inverse of $\bar{K}_{Q,\mathbf{z}}$ in (23b) given by

$$\bar{K}_{Q,\mathbf{z}}^{-1}(\gamma) = \frac{1}{\int \frac{1}{\mathbb{L}_{\mathbf{z}}(\boldsymbol{\theta}) + \gamma} dQ(\boldsymbol{\theta})}, \quad (124)$$

which follows from the constraint in (123) and the equality in (122).

From (123) and the fact that $t \in (0, \infty)$, it holds that

$$0 < \int \frac{1}{\mathbb{L}_{\mathbf{z}}(\boldsymbol{\theta}) + \gamma} dQ(\boldsymbol{\theta}) < \infty, \quad (125)$$

which implies that for all $\gamma \in \mathcal{A}$, with \mathcal{A} in (23a), the function $\bar{K}_{Q,\mathbf{z}}^{-1}$ in (124) satisfies

$$\infty > \bar{K}_{Q,\mathbf{z}}^{-1}(\gamma) > 0. \quad (126)$$

Note that from (125), for all $(\gamma_1, \gamma_2) \in \mathcal{A}^2$, such that $\gamma_1 < \gamma_2$, for all $\boldsymbol{\theta} \in \text{supp } Q$ it holds that

$$\frac{1}{\mathbb{L}_{\mathbf{z}}(\boldsymbol{\theta}) + \gamma_1} > \frac{1}{\mathbb{L}_{\mathbf{z}}(\boldsymbol{\theta}) + \gamma_2}, \quad (127)$$

which implies that

$$\int \frac{1}{\mathbb{L}_{\mathbf{z}}(\boldsymbol{\theta}) + \gamma_1} dQ(\boldsymbol{\theta}) > \int \frac{1}{\mathbb{L}_{\mathbf{z}}(\boldsymbol{\theta}) + \gamma_2} dQ(\boldsymbol{\theta}). \quad (128)$$

Therefore, from equality (124) and inequality (128) for all $(\gamma_1, \gamma_2) \in \mathcal{A}^2$, such that $\gamma_1 < \gamma_2$, the function $\bar{K}_{Q,\mathbf{z}}^{-1}$ in (124) satisfies

$$\bar{K}_{Q,\mathbf{z}}^{-1}(\gamma_1) < \bar{K}_{Q,\mathbf{z}}^{-1}(\gamma_2), \quad (129)$$

which implies that the function $\bar{K}_{Q,z}^{-1}$ in (124) is strictly increasing. Furthermore, from inequality (129) and the fact that $\bar{P}_{\Theta|Z=z}^{(Q,t)}$ is unique (see Theorem 4.1), it follows that $\bar{K}_{Q,z}^{-1}$ is bijective, which completes the proof of the first part.

The second part is as follows. The continuity of the function $\bar{K}_{Q,z}^{-1}$ in (124) is shown by considering an arbitrary $\epsilon > 0$ and a pair $(\gamma_1, \gamma_2) \in \mathcal{A}^2$ under the condition

$$\epsilon > \left| \bar{K}_{Q,z}^{-1}(\gamma_2) - \bar{K}_{Q,z}^{-1}(\gamma_1) \right|. \quad (130)$$

From the fact that $\bar{K}_{Q,z}^{-1}$ is bijective and strictly increasing, let the pair γ_1 , and γ_2 satisfy $\gamma_1 < \gamma_2$, without loss of generality. Thus, it holds that

$$\left| \bar{K}_{Q,z}^{-1}(\gamma_2) - \bar{K}_{Q,z}^{-1}(\gamma_1) \right| = \bar{K}_{Q,z}^{-1}(\gamma_2) - \bar{K}_{Q,z}^{-1}(\gamma_1). \quad (131)$$

Substituting equality (131) into (130) yields

$$\bar{K}_{Q,z}^{-1}(\gamma_1) + \epsilon > \bar{K}_{Q,z}^{-1}(\gamma_2). \quad (132)$$

From the fact that $\bar{K}_{Q,z}^{-1}$ is bijective, evaluating inequality (132) with $\bar{K}_{Q,z}$ in (23a) leads to

$$\bar{K}_{Q,z}(\bar{K}_{Q,z}^{-1}(\gamma_1) + \epsilon) > \bar{K}_{Q,z}(\bar{K}_{Q,z}^{-1}(\gamma_2)) \quad (133a)$$

$$= \gamma_2. \quad (133b)$$

Subtracting γ_1 from both sides in (133), results in

$$\bar{K}_{Q,z}(\bar{K}_{Q,z}^{-1}(\gamma_1) + \epsilon) - \gamma_1 > \gamma_2 - \gamma_1 \quad (134a)$$

$$= |\gamma_2 - \gamma_1|, \quad (134b)$$

where equality (134b) follows from the condition that $\gamma_2 > \gamma_1$. Thus, from (134) it follows that for all $\epsilon > 0$, there exist a $\delta > 0$ that satisfies

$$\delta = \bar{K}_{Q,z}(\bar{K}_{Q,z}^{-1}(\gamma_1) + \epsilon) - \gamma_1, \quad (135)$$

which implies the function $\bar{K}_{Q,z}^{-1}$ in (124) is continuous for all $\gamma \in \mathcal{A}$. From (135) the function $\bar{K}_{Q,z}^{-1}$ is continuous and strictly increasing. Hence, from the continuous inverse theorem in [47, Theorem 5.6], the inverse of the function $\bar{K}_{Q,z}^{-1}$ in (124), which is given by $\bar{K}_{Q,z}$ in (23) is continuous and strictly increasing, which completes the proof. ■

D Proof of Lemma 4.5

Proof: The proof shows by exhaustion that the codomain of the function $\bar{K}_{Q,z}$ in (23) is an interval of \mathbb{R} by evaluating β in (25). Three cases are considered: *i*) $\beta < -\delta_{Q,z}^*$; *ii*) $\beta > -\delta_{Q,z}^*$; and *iii*) $\beta = -\delta_{Q,z}^*$, with $\delta_{Q,z}^*$ in (28).

In the first case, from the definition of $\delta_{Q,z}^*$ in (28), it follows from (15b) that if $\beta < -\delta_{Q,z}^*$, then for all $\theta \in \{\theta \in \text{supp } Q : \mathbb{L}_z(\theta) < -\beta\}$, it holds that

$$\frac{d\bar{P}_{\Theta|Z=z}^{(Q,t)}}{dQ}(\theta) = \frac{\lambda}{\mathbb{L}_z(\theta) + \beta} \quad (136a)$$

$$< 0, \quad (136b)$$

which contradicts the fact that the function $\frac{d\bar{P}_{\Theta|Z=z}^{(Q,t)}}{dQ}$ in (16) is nonnegative [45, Corollary 2.2.2]. The above implies that for all $t \in (0, \infty)$, the function $\bar{K}_{Q,z}$ in (23) satisfies that $\bar{K}_{Q,z}(t) \notin (-\infty, -\delta_{Q,z}^*)$.

In the second case, if $\beta > -\delta_{Q,z}^*$, it holds that

$$\int \frac{1}{\mathbb{L}_z(\theta) + \beta} dQ(\theta) \leq \int \frac{1}{\delta_{Q,z}^* + \beta} dQ(\theta) \quad (137a)$$

$$= \frac{1}{\delta_{Q,z}^* + \beta}. \quad (137b)$$

Moreover, from (137) for all $\beta > -\delta_{Q,z}^*$, there exists a $\lambda \in (0, \infty)$ such that the constraint in (24) holds, which implies $\bar{K}_{Q,z}^{-1}(\beta) \in (0, \infty)$.

Finally, under the assumption that $\beta = -\delta_{Q,z}^*$ and $\mathcal{L}_{Q,z}^*$ defined in (29), two cases are considered: (a) $Q(\mathcal{L}_{Q,z}^*) > 0$; and (b) $Q(\mathcal{L}_{Q,z}^*) = 0$. In case (a), if $\beta = -\delta_{Q,z}^*$ and $Q(\mathcal{L}_{Q,z}^*) > 0$, then the integral in the denominator of (30) is undefined, contradicting (24). In the alternative case (b), if $\beta = -\delta_{Q,z}^*$ and $Q(\mathcal{L}_{Q,z}^*) = 0$, then the integral in (30) is either

$$\int \frac{1}{\mathbb{L}_z(\theta) - \delta_{Q,z}^*} dQ(\theta) < \infty, \quad (138a)$$

which implies that $-\delta_{Q,z}^* \in \mathcal{A}$, with \mathcal{A} defined in (23a), or the integral is

$$\int \frac{1}{\mathbb{L}_z(\theta) - \delta_{Q,z}^*} dQ(\theta) = \infty, \quad (138b)$$

which implies that, $-\delta_{Q,z}^* \notin \mathcal{A}$. Note that the integral is always positive from the fact that for all $\theta \in \text{supp } Q$, it holds that $\mathbb{L}_z(\theta) > \delta_{Q,z}^*$. Moreover, the integral in (138) is never zero as it would be a contradiction with (2). Hence, from (136), (137) and (138) the set \mathcal{A} in (23a) is either the open set $(-\delta_{Q,z}^*, \infty)$ or the closed set $[-\delta_{Q,z}^*, \infty)$, which completes the proof. ■

E Proof of Lemma 4.8

Proof: From (15b), for all $\lambda \in (0, \infty)$, it follows that

$$\lambda = \frac{1}{\int \frac{1}{\mathbb{L}_z(\nu) + \beta} dQ(\nu)} \quad (139a)$$

$$= \frac{1}{\int \frac{1}{\mathbb{L}_z(\boldsymbol{\nu}) + \bar{K}_{Q,z}(\lambda)} dQ(\boldsymbol{\nu})}, \quad (139b)$$

where equality (139a) follows from (15b), and equality (139b) follows from the nominal function in (23b). From the definition of $\delta_{Q,z}^*$ in (28), for all $\boldsymbol{\nu} \in \text{supp } Q$, the empirical risk satisfies

$$\mathbb{L}_z(\boldsymbol{\nu}) \geq \delta_{Q,z}^*. \quad (140)$$

Taking the limit of (139) as λ approaches zero from the right satisfies

$$\lim_{\lambda \rightarrow 0^+} \lambda = \lim_{\lambda \rightarrow 0^+} \frac{1}{\int \frac{1}{\mathbb{L}_z(\boldsymbol{\nu}) + \bar{K}(\lambda)} dQ(\boldsymbol{\nu})} \quad (141a)$$

$$\geq \lim_{\lambda \rightarrow 0^+} \frac{1}{\int \frac{1}{\delta_{Q,z}^* + \bar{K}(\lambda)} dQ(\boldsymbol{\nu})} \quad (141b)$$

$$= \lim_{\lambda \rightarrow 0^+} \frac{1}{\frac{1}{\delta_{Q,z}^* + \bar{K}(\lambda)}} \quad (141c)$$

$$= \lim_{\lambda \rightarrow 0^+} (\delta_{Q,z}^* + \bar{K}(\lambda)) \quad (141d)$$

$$= \delta_{Q,z}^* + \lim_{\lambda \rightarrow 0^+} \bar{K}(\lambda), \quad (141e)$$

where inequality (141b) follows from (140). Hence, from (141) it holds that

$$0 \geq \delta_{Q,z}^* + \lim_{\lambda \rightarrow 0^+} \bar{K}(\lambda). \quad (142)$$

From (15a) and Lemma 4.4, for all $\lambda \in (0, \infty)$, it holds that

$$0 \leq \delta_{Q,z}^* + \bar{K}_{Q,z}(\lambda). \quad (143)$$

Taking the limit when λ approaches zero from the right in (143) yields

$$0 \leq \delta_{Q,z}^* + \lim_{\lambda \rightarrow 0^+} \bar{K}_{Q,z}(\lambda). \quad (144)$$

Then, from (142) and (144) it follows that

$$\lim_{\lambda \rightarrow 0^+} \bar{K}_{Q,z}(\lambda) = -\delta_{Q,z}^*, \quad (145)$$

which completes the proof. \blacksquare

F Proof of Lemma 4.9

Proof: For all $\boldsymbol{\theta}_1 \in \text{supp } Q$ and for all $\boldsymbol{\theta}_2 \in \mathcal{L}_{Q,z}^*$, it follows that

$$\mathbb{L}_z(\boldsymbol{\theta}_1) \geq \mathbb{L}_z(\boldsymbol{\theta}_2), \quad (146)$$

and thus, for all $\lambda \in (0, \infty)$, it holds that

$$\frac{1}{\bar{K}_{Q,z}(\lambda) + \mathbb{L}_z(\boldsymbol{\theta}_1)} \leq \frac{1}{\bar{K}_{Q,z}(\lambda) + \mathbb{L}_z(\boldsymbol{\theta}_2)}, \quad (147a)$$

which implies

$$\frac{(\bar{K}_{Q,z}(\lambda) + \mathbf{L}_z(\theta_1))^{-1}}{\int (\bar{K}_{Q,z}(\lambda) + \mathbf{L}_z(\nu))^{-1} dQ(\nu)} \leq \frac{(\bar{K}_{Q,z}(\lambda) + \mathbf{L}_z(\theta_2))^{-1}}{\int (\bar{K}_{Q,z}(\lambda) + \mathbf{L}_z(\nu))^{-1} dQ(\nu)}. \quad (148)$$

Hence, under the assumption that $\mathcal{L}_{Q,z}^* \cap \text{supp } Q \neq \emptyset$, for all $\theta_1 \in \text{supp } Q$ and for all $\theta_2 \in \mathcal{L}_{Q,z}^* \cap \text{supp } Q$, it holds that

$$\frac{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta_1) \leq \frac{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta_2), \quad (149)$$

with equality if and only if $\theta_1 \in \mathcal{L}_{Q,z}^* \cap \text{supp } Q$, which completes the proof. \blacksquare

G Proof of Lemma 4.10

Proof: From Lemma 4.9, it follows that for all $\lambda \in (0, \infty)$, for all $\theta \in \text{supp } Q$, and for all $\phi \in \mathcal{L}_{Q,z}^* \cap \text{supp } Q$, it holds that

$$\frac{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) \leq \frac{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\phi) \quad (150a)$$

$$= \frac{\lambda}{\mathbf{L}_z(\phi) + \bar{K}_{Q,z}(\lambda)} \quad (150b)$$

$$= \frac{\lambda}{\delta_{Q,z}^* + \bar{K}_{Q,z}(\lambda)} \quad (150c)$$

$$< \infty, \quad (150d)$$

where equality (150a) follows from (16); equality (150b) follows from the fact that $\mathbf{L}_z(\phi) \geq \delta_{Q,z}^*$; and the equality in (150d) follows from the fact that for all $\lambda > 0$, the function $\bar{K}_{Q,z}(\lambda) < \infty$. From the definition of $\delta_{Q,z}^*$ in (28) and $\mathcal{L}_{Q,z}^*$ in (29) equality in (150a) holds if and only if $\theta \in \mathcal{L}_{Q,z}^* \cap \text{supp } Q$. This completes the proof of finiteness.

For the proof of positivity, observe that from Lemma 4.4 for all $\lambda \in (0, \infty)$, it holds that

$$-\delta_{Q,z}^* < \bar{K}_{Q,z}(\lambda) < \infty, \quad (151)$$

which implies

$$0 < \delta_{Q,z}^* + \bar{K}_{Q,z}(\lambda). \quad (152)$$

From (151), it follows that

$$\frac{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) = \frac{\lambda}{\bar{K}_{Q,z}(\lambda) + \mathbf{L}_z(\theta)} \quad (153a)$$

$$= \frac{\int \frac{1}{\bar{K}_{Q,z}(\lambda) + \mathbf{L}_z(\nu)} dQ(\nu)}{\bar{K}_{Q,z}(\lambda) + \mathbf{L}_z(\theta)} \quad (153b)$$

$$> 0, \quad (153c)$$

which completes the proof. \blacksquare

H Proof of Lemma 4.11

Proof: From Theorem 4.1, the Radon-Nikodym derivative of the measure $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ with respect to Q , for all $\theta \in \text{supp } Q$, satisfies that

$$\frac{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) = \frac{\lambda}{\bar{K}_{Q,z}(\lambda) + \mathbf{L}_z(\theta)} \quad (154a)$$

$$= \frac{\int \frac{1}{\bar{K}_{Q,z}(\lambda) + \mathbf{L}_z(\nu)} dQ(\nu)}{\bar{K}_{Q,z}(\lambda) + \mathbf{L}_z(\theta)} \quad (154b)$$

$$= \left(\frac{\bar{K}_{Q,z}(\lambda) + \mathbf{L}_z(\theta)}{\int (\bar{K}_{Q,z}(\lambda) + \mathbf{L}_z(\nu))^{-1} dQ(\nu)} \right)^{-1} \quad (154c)$$

$$= \frac{\exp(-\log(\bar{K}_{Q,z}(\lambda) + \mathbf{L}_z(\theta)))}{\int \exp\left(\log\left(\frac{1}{\bar{K}_{Q,z}(\lambda) + \mathbf{L}_z(\nu)}\right)\right) dQ(\nu)} \quad (154d)$$

From Lemma 4.4, the case in which $\lambda \rightarrow \infty$ implies $\bar{K}_{Q,z}(\lambda) \rightarrow \infty$, and from (154), it follows that

$$\lim_{\lambda \rightarrow \infty} \frac{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) = \lim_{\lambda \rightarrow \infty} \frac{\exp(-\log(\bar{K}_{Q,z}(\lambda) + \mathbf{L}_z(\theta)))}{\int \exp\left(\log\left(\frac{1}{\bar{K}_{Q,z}(\lambda) + \mathbf{L}_z(\nu)}\right)\right) dQ(\nu)} \quad (155)$$

$$= \frac{1}{\int dQ(\nu)} \quad (156)$$

$$= 1, \quad (157)$$

where the function \mathbf{L}_z is defined in (3). This completes the proof. \blacksquare

I Proof of Lemma 4.12

Proof: Assume that $t \in (0, \infty)$ and $\gamma \in \mathbb{R}$ satisfy that

$$1 = \int \frac{t}{\mathbf{L}_z(\theta) + \gamma} dQ(\theta), \quad (158)$$

which implies

$$\bar{K}_{Q,z}(t) = \gamma. \quad (159)$$

Let the function $\bar{K}_{Q,z}^{-1} : \mathbb{R} \rightarrow (0, \infty)$ be the functional inverse of $\bar{K}_{Q,z}$ in (159). Thus,

$$\bar{K}_{Q,z}^{-1}(\gamma) = t, \quad (160)$$

with γ and t in (158). From (158), the function $\bar{K}_{Q,z}^{-1}$ in (160) satisfies

$$\bar{K}_{Q,z}^{-1}(\gamma) = \frac{1}{\int \frac{1}{L_z(\theta) + \gamma} dQ(\theta)}. \quad (161)$$

From Theorem 4.1, for all $\lambda \in (0, \infty)$ and for all $\theta \in \text{supp } Q$, it follows that

$$\frac{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}(\theta)}{dQ} = \frac{\lambda}{L_z(\theta) + \bar{K}_{Q,z}(\lambda)} \quad (162a)$$

$$= \frac{\bar{K}_{Q,z}^{-1}(\bar{K}_{Q,z}(\lambda))}{L_z(\theta) + \bar{K}_{Q,z}(\lambda)} \quad (162b)$$

$$= \frac{\int \frac{1}{L_z(\nu) + \bar{K}_{Q,z}(\lambda)} dQ(\nu)}{L_z(\theta) + \bar{K}_{Q,z}(\lambda)} \quad (162c)$$

$$= \left(\int \frac{L_z(\theta) + \bar{K}_{Q,z}(\lambda)}{L_z(\nu) + \bar{K}_{Q,z}(\lambda)} dQ(\nu) \right)^{-1}. \quad (162d)$$

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

$$\mathcal{A}_0 = \{\nu \in \mathcal{M} : L_z(\nu) = \delta_{Q,z}^*\}, \quad (163a)$$

$$\mathcal{A}_1 = \{\nu \in \mathcal{M} : L_z(\nu) > \delta_{Q,z}^*\}, \quad (163b)$$

$$\mathcal{A}_2 = \{\nu \in \mathcal{M} : L_z(\nu) < \delta_{Q,z}^*\}, \quad (163c)$$

with $\delta_{Q,z}^*$ defined in (28). Using the sets \mathcal{A}_0 , \mathcal{A}_1 , and \mathcal{A}_2 in (162), for all $\lambda \in (0, \infty)$ and for all $\theta \in \text{supp } Q$ the following holds

$$\begin{aligned} & \frac{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}(\theta)}{dQ} \\ &= \left(\int_{\mathcal{A}_0} \frac{L_z(\theta) + \bar{K}_{Q,z}(\lambda)}{L_z(\nu) + \bar{K}_{Q,z}(\lambda)} dQ(\nu) + \int_{\mathcal{A}_1} \frac{L_z(\theta) + \bar{K}_{Q,z}(\lambda)}{L_z(\nu) + \bar{K}_{Q,z}(\lambda)} dQ(\nu) \right. \\ & \quad \left. + \int_{\mathcal{A}_2} \frac{L_z(\theta) + \bar{K}_{Q,z}(\lambda)}{L_z(\nu) + \bar{K}_{Q,z}(\lambda)} dQ(\nu) \right)^{-1} \end{aligned} \quad (164a)$$

$$= \left(\int_{\mathcal{L}_{Q,z}^*} \frac{L_z(\theta) + \bar{K}_{Q,z}(\lambda)}{L_z(\nu) + \bar{K}_{Q,z}(\lambda)} dQ(\nu) + \int_{\mathcal{A}_1} \frac{L_z(\theta) + \bar{K}_{Q,z}(\lambda)}{L_z(\nu) + \bar{K}_{Q,z}(\lambda)} dQ(\nu) \right)^{-1} \quad (164b)$$

$$= \left(\int_{\mathcal{L}_{Q,z}^*} \frac{L_z(\theta) + \bar{K}_{Q,z}(\lambda)}{\delta_{Q,z}^* + \bar{K}_{Q,z}(\lambda)} dQ(\nu) + \int_{\mathcal{A}_1} \frac{L_z(\theta) + \bar{K}_{Q,z}(\lambda)}{L_z(\nu) + \bar{K}_{Q,z}(\lambda)} dQ(\nu) \right)^{-1} \quad (164c)$$

where equality (164b) follows from the fact that the \mathcal{A}_2 has measure zero with respect to Q and that the set $\mathcal{A}_0 = \mathcal{L}_{Q,z}^*$, with $\mathcal{L}_{Q,z}^*$ in (29); and equality

(164b) follows from the definition of $\delta_{Q,z}^*$ in (28). Following this observation, the rest of the proof is divided into three parts. The first part evaluates $\lim_{\lambda \rightarrow \infty} \frac{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta)$, with $\theta \in \{\nu \in \mathcal{M} : L_z(\nu) = \delta_{Q,z}^*\}$. The second part considers the case in which $\theta \in \{\nu \in \mathcal{M} : L_z(\nu) > \delta_{Q,z}^*\}$. The third part considers the remaining case.

The first part is as follows. Consider that $\theta \in \{\nu \in \mathcal{M} : L_z(\nu) = \delta_{Q,z}^*\}$ and note that $\mathcal{A}_0 = \mathcal{L}_{Q,z}^*$. From (164), for all $\theta \in \{\nu \in \mathcal{M} : L_z(\nu) = \delta_{Q,z}^*\}$ it holds that

$$\frac{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) = \left(\int_{\mathcal{L}_{Q,z}^*} \frac{\delta_{Q,z}^* + \bar{K}_{Q,z}(\lambda)}{\delta_{Q,z}^* + \bar{K}_{Q,z}(\lambda)} dQ(\nu) + \int_{\mathcal{A}_1} \frac{\delta_{Q,z}^* + \bar{K}_{Q,z}(\lambda)}{L_z(\nu) + \bar{K}_{Q,z}(\lambda)} dQ(\nu) \right)^{-1} \quad (165a)$$

$$= \left(Q(\mathcal{L}_{Q,z}^*) + \int_{\mathcal{A}_1} \frac{\delta_{Q,z}^* + \bar{K}_{Q,z}(\lambda)}{L_z(\nu) + \bar{K}_{Q,z}(\lambda)} dQ(\nu) \right)^{-1}. \quad (165b)$$

The equality in (165) implies that for all $\theta \in \{\nu \in \mathcal{M} : L_z(\nu) = \delta_{Q,z}^*\}$ the limit satisfies

$$\lim_{\lambda \rightarrow 0^+} \frac{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) = \left(Q(\mathcal{L}_{Q,z}^*) + \lim_{\lambda \rightarrow 0^+} \int_{\mathcal{A}_1} \frac{\delta_{Q,z}^* + \bar{K}_{Q,z}(\lambda)}{L_z(\nu) + \bar{K}_{Q,z}(\lambda)} dQ(\nu) \right)^{-1} \quad (166a)$$

$$= \left(Q(\mathcal{L}_{Q,z}^*) + \lim_{\lambda \rightarrow 0^+} (\delta_{Q,z}^* + \bar{K}_{Q,z}(\lambda)) \lim_{\lambda \rightarrow 0^+} \int_{\mathcal{A}_1} \frac{1}{L_z(\nu) + \bar{K}_{Q,z}(\lambda)} dQ(\nu) \right)^{-1} \quad (166b)$$

$$= \left(Q(\mathcal{L}_{Q,z}^*) + \left(\delta_{Q,z}^* + \lim_{\lambda \rightarrow 0^+} \bar{K}_{Q,z}(\lambda) \right) \int_{\mathcal{A}_1} \frac{1}{L_z(\nu) + \lim_{\lambda \rightarrow 0^+} \bar{K}_{Q,z}(\lambda)} dQ(\nu) \right)^{-1} \quad (166c)$$

$$= \left(Q(\mathcal{L}_{Q,z}^*) + 0 \int_{\mathcal{A}_1} \frac{1}{L_z(\nu) - \delta_{Q,z}^*} dQ(\nu) \right)^{-1} \quad (166d)$$

$$= (Q(\mathcal{L}_{Q,z}^*))^{-1}, \quad (166e)$$

where equality (166d) follows from Lemma 4.8.

The second part is as follows. Consider that $\theta \in \{\nu \in \mathcal{M} : L_z(\nu) > \delta_{Q,z}^*\}$. Hence, for all $\theta \in \{\nu \in \mathcal{M} : L_z(\nu) > \delta_{Q,z}^*\}$, there exists a real value $\epsilon > 0$ such

that the function L_z satisfies

$$L_z(\theta) - \delta_{Q,z}^* = \epsilon. \quad (167)$$

From (164), for all $\theta \in \{\nu \in \mathcal{M} : L_z(\nu) > \delta_{Q,z}^*\}$ it holds that

$$\begin{aligned} \frac{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) &= \left(\int_{\mathcal{L}_{Q,z}^*} \frac{\delta_{Q,z}^* + \epsilon + \bar{K}_{Q,z}(\lambda)}{\delta_{Q,z}^* + \bar{K}_{Q,z}(\lambda)} dQ(\nu) \right. \\ &\quad \left. + \int_{\mathcal{A}_1} \frac{\delta_{Q,z}^* + \epsilon + \bar{K}_{Q,z}(\lambda)}{L_z(\nu) + \bar{K}_{Q,z}(\lambda)} dQ(\nu) \right)^{-1} \end{aligned} \quad (168a)$$

$$\begin{aligned} &= \left(\frac{\delta_{Q,z}^* + \epsilon + \bar{K}_{Q,z}(\lambda)}{\delta_{Q,z}^* + \bar{K}_{Q,z}(\lambda)} Q(\mathcal{L}_{Q,z}^*) \right. \\ &\quad \left. + \int_{\mathcal{A}_1} \frac{\delta_{Q,z}^* + \bar{K}_{Q,z}(\lambda)}{L_z(\nu) + \bar{K}_{Q,z}(\lambda)} dQ(\nu) \right)^{-1}. \end{aligned} \quad (168b)$$

The equality in (168) implies that for all $\theta \in \{\nu \in \mathcal{M} : L_z(\nu) = \delta_{Q,z}^*\}$, the limit satisfies

$$\begin{aligned} \lim_{\lambda \rightarrow 0^+} \frac{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) &= \left(\lim_{\lambda \rightarrow 0^+} \frac{\delta_{Q,z}^* + \epsilon + \bar{K}_{Q,z}(\lambda)}{\delta_{Q,z}^* + \bar{K}_{Q,z}(\lambda)} Q(\mathcal{L}_{Q,z}^*) \right. \\ &\quad \left. + \lim_{\lambda \rightarrow 0^+} \int_{\mathcal{A}_1} \frac{\delta_{Q,z}^* + \epsilon + \bar{K}_{Q,z}(\lambda)}{L_z(\nu) + \bar{K}_{Q,z}(\lambda)} dQ(\nu) \right)^{-1} \end{aligned} \quad (169a)$$

$$\begin{aligned} &= \left(\frac{\delta_{Q,z}^* + \epsilon + \lim_{\lambda \rightarrow 0^+} \bar{K}_{Q,z}(\lambda)}{\delta_{Q,z}^* + \lim_{\lambda \rightarrow 0^+} \bar{K}_{Q,z}(\lambda)} Q(\mathcal{L}_{Q,z}^*) \right. \\ &\quad \left. + \lim_{\lambda \rightarrow 0^+} (\delta_{Q,z}^* + \bar{K}_{Q,z}(\lambda)) \right. \\ &\quad \left. \lim_{\lambda \rightarrow 0^+} \int_{\mathcal{A}_1} \frac{1}{L_z(\nu) + \bar{K}_{Q,z}(\lambda)} dQ(\nu) \right)^{-1} \end{aligned} \quad (169b)$$

$$= \left(\infty Q(\mathcal{L}_{Q,z}^*) + \epsilon \int_{\mathcal{A}_1} \frac{1}{L_z(\nu) - \delta_{Q,z}^*} dQ(\nu) \right)^{-1} \quad (169c)$$

The proof continues by considering the cases in which $Q(\mathcal{A}_0) > 0$ and $Q(\mathcal{A}_0) = 0$ for the limit in (169). Under the assumption that $Q(\mathcal{A}_0) > 0$, from (169b) the limit satisfies

$$\lim_{\lambda \rightarrow 0^+} \frac{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) = 0. \quad (170a)$$

Under the assumption that $Q(\mathcal{A}_0) = 0$, from the definition of $\delta_{Q,z}^*$ in (28), if $Q(\mathcal{A}_0) = 0$, then it is implied that there is at least an infinite countable sequence of models ν_1, ν_2, \dots in $\text{supp } Q$ that satisfies

$$L_z(\nu_1) > L_z(\nu_2) > \dots, \quad (171)$$

such that

$$\lim_{i \rightarrow \infty} \mathsf{L}_{\mathbf{z}}(\boldsymbol{\nu}_i) = \delta_{Q,\mathbf{z}}^*. \quad (172)$$

From (172) the function inside the integral in (169c) is unbounded over the set \mathcal{A}_1 . Therefore,

$$\lim_{\lambda \rightarrow 0^+} \frac{d\bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}}{dQ}(\boldsymbol{\theta}) = \frac{1}{0 + (\infty)\epsilon} \quad (173a)$$

$$= 0. \quad (173b)$$

The third part of the proof follows by noticing that the set $\left\{ \boldsymbol{\nu} \in \text{supp } Q : \mathsf{L}_{\mathbf{z}}(\boldsymbol{\nu}) < \delta_{Q,\mathbf{z}}^* \right\}$ has measure zero with respect to Q and thus, for all $\boldsymbol{\theta} \in \left\{ \boldsymbol{\nu} \in \text{supp } Q : \mathsf{L}_{\mathbf{z}}(\boldsymbol{\nu}) < \delta_{Q,\mathbf{z}}^* \right\}$, the value $\frac{d\bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}}{dQ}(\boldsymbol{\theta})$ is unchanged. Hence, it is assumed that for all $\boldsymbol{\theta} \in \left\{ \boldsymbol{\nu} \in \text{supp } Q : \mathsf{L}_{\mathbf{z}}(\boldsymbol{\nu}) < \delta_{Q,\mathbf{z}}^* \right\}$, it holds that

$$\frac{d\bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}}{dQ}(\boldsymbol{\theta}) = 0. \quad (174)$$

This completes the third part and the entire proof. ■

J Proof of Lemma 4.13

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

$$\mathcal{A}_0 = \left\{ \boldsymbol{\theta} \in \mathcal{M} : \mathsf{L}_{\mathbf{z}}(\boldsymbol{\theta}) = \delta_{Q,\mathbf{z}}^* \right\}, \quad (175a)$$

$$\mathcal{A}_1 = \left\{ \boldsymbol{\theta} \in \mathcal{M} : \mathsf{L}_{\mathbf{z}}(\boldsymbol{\theta}) > \delta_{Q,\mathbf{z}}^* \right\}, \quad (175b)$$

$$\mathcal{A}_2 = \left\{ \boldsymbol{\theta} \in \mathcal{M} : \mathsf{L}_{\mathbf{z}}(\boldsymbol{\theta}) < \delta_{Q,\mathbf{z}}^* \right\}, \quad (175c)$$

with $\delta_{Q,\mathbf{z}}^*$ in (28) and the function $\mathsf{L}_{\mathbf{z}}$ in (3). Note that $\mathcal{A}_0 = \mathcal{L}_{Q,\mathbf{z}}^*$, with $\mathcal{L}_{Q,\mathbf{z}}^*$ in (29).

For all $\lambda \in (0, \infty)$, it holds that

$$1 = \bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}(\mathcal{A}_0) + \bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}(\mathcal{A}_1) + \bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}(\mathcal{A}_2) \quad (176a)$$

$$= \bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}(\mathcal{A}_0) + \bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}(\mathcal{A}_1) \quad (176b)$$

$$= \bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}(\mathcal{A}_0) + \int_{\mathcal{A}_1} d\bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}(\boldsymbol{\theta}), \quad (176c)$$

where equality (176b) follows from the fact that $\bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}(\mathcal{A}_2) = 0$, which follows from the definition of $\delta_{Q,\mathbf{z}}^*$ in (28) and the fact that the probability measure

$\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ is mutually absolutely continuous with respect to the reference measure Q . The above implies that

$$\begin{aligned} \lim_{\lambda \rightarrow 0^+} \left(\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{A}_0) + \int_{\mathcal{A}_1} d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}(\theta) \right) \\ = \lim_{\lambda \rightarrow 0^+} \bar{P}_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{A}_0) + \lim_{\lambda \rightarrow 0^+} \int_{\mathcal{A}_1} \frac{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) dQ(\theta) \end{aligned} \quad (177a)$$

$$= \lim_{\lambda \rightarrow 0^+} \bar{P}_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{A}_0) + \int_{\mathcal{A}_1} \lim_{\lambda \rightarrow 0^+} \frac{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) dQ(\theta) \quad (177b)$$

$$= \lim_{\lambda \rightarrow 0^+} \bar{P}_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{A}_0), \quad (177c)$$

$$= 1, \quad (177d)$$

where, the equality in (177b) follows from the dominated convergence theorem [45, Theorem 1.6.9 page 50], given that from Lemma 4.10 for all $\lambda \in \bar{K}_{Q,z}$, the Radon-Nikodym derivative $\frac{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}$ is positive and finite; and the equality in (177c) follows from the fact that for all $\theta \in \mathcal{A}_1$, it holds that $\lim_{\lambda \rightarrow 0^+} \frac{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) = 0$ from Lemma 4.12.

Hence, it holds that

$$\lim_{\lambda \rightarrow 0^+} \bar{P}_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_{Q,z}^*) = 1, \quad (178)$$

which completes the proof. \blacksquare

K Proof of Lemma 5.1

Proof: From Lemma 4.1 and the fact that the measures $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ and Q are mutually absolutely continuous, it holds that for all $\theta \in \text{supp } Q$,

$$\frac{dQ}{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}(\theta) = \frac{\bar{K}_{Q,z}(\lambda) + \mathbf{L}_z(\theta)}{\lambda}, \quad (179)$$

where the functions \mathbf{L}_z and $\bar{K}_{Q,z}$ are in (3) and (23b), respectively. From (179), it follows that for all $\theta \in \text{supp } Q$,

$$0 = \lambda \frac{dQ}{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}(\theta) - \mathbf{L}_z(\theta) - \bar{K}_{Q,z}(\lambda). \quad (180)$$

Integrating both sides of (180) with respect to the probability measure $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ yields

$$0 = \int \left(\mathbf{L}_z(\theta) - \lambda \left(\frac{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) \right)^{-1} + \bar{K}_{Q,z}(\lambda) \right) d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}(\theta) \quad (181a)$$

$$\begin{aligned}
 &= \int \mathbb{L}_z(\boldsymbol{\theta}) d\bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=z}^{(Q,\lambda)}(\boldsymbol{\theta}) - \lambda \int \left(\frac{d\bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=z}^{(Q,\lambda)}}{dQ}(\boldsymbol{\theta}) \right)^{-1} d\bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=z}^{(Q,\lambda)}(\boldsymbol{\theta}) \\
 &\quad + \int \bar{K}_{Q,z}(\lambda) d\bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=z}^{(Q,\lambda)}(\boldsymbol{\theta}) \tag{181b}
 \end{aligned}$$

$$= R_z\left(\bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=z}^{(Q,\lambda)}\right) - \lambda \int dQ(\boldsymbol{\theta}) + \bar{K}_{Q,z}(\lambda) \tag{181c}$$

$$= R_z\left(\bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=z}^{(Q,\lambda)}\right) - \lambda + \bar{K}_{Q,z}(\lambda). \tag{181d}$$

From (181d), it holds that

$$R_z\left(\bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=z}^{(Q,\lambda)}\right) = \lambda - \bar{K}_{Q,z}(\lambda), \tag{182}$$

which completes the proof. \blacksquare

L Proof of Lemma 5.2

Proof: The proof is divided into four parts. The first part characterizes the functional inverse of the function $\bar{K}_{Q,z}$ in (23b). The second part provides necessary conditions for the differentiation lemma [48, Theorem 6.28, page 160] to hold, which is used in the third and fourth parts of the proof. The third part presents the first derivative of the functional inverse and shows that its derivative is strictly positive. The forth part shows the expected empirical risk $R_z\left(\bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=z}^{(Q,\lambda)}\right)$ decreases with lambdas decreasing.

The first part is as follows. Under the assumption that $t \in (0, \infty)$ and $\gamma \in \mathbb{R}$ satisfy that

$$1 = \int \frac{d\bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=z}^{(Q,t)}}{dQ}(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) \tag{183a}$$

$$= \int \frac{t}{\mathbb{L}_z(\boldsymbol{\theta}) + \gamma} dQ(\boldsymbol{\theta}), \tag{183b}$$

which implies that

$$\bar{K}_{Q,z}(t) = \gamma, \tag{184}$$

where $\bar{K}_{Q,z}$ is defined in (23b). From Lemma 4.4, let the functional $\bar{K}_{Q,z}^{-1} : (-\delta_{Q,z}^*, \infty) \rightarrow (0, \infty)$ be the functional inverse of $\bar{K}_{Q,z}$ in (23b) given by

$$\bar{K}_{Q,z}^{-1}(\beta) = \frac{1}{\int \frac{1}{\mathbb{L}_z(\boldsymbol{\theta}) + \beta} dQ(\boldsymbol{\theta})}, \tag{185}$$

which follows from the constraint in (15b) and completes the first part of the proof.

The second part is as follows. From Lemma 4.4, for all $\lambda \in (0, \infty)$, it holds that

$$\bar{K}_{Q,z}(\lambda) > -\delta_{Q,z}^*. \quad (186)$$

Then, for the third case, if $\gamma > -\delta_{Q,z}^*$, it follows that

$$\int \frac{1}{L_z(\theta) + \gamma} dQ(\theta) \leq \int \frac{1}{\delta_{Q,z}^* + \gamma} dQ(\theta) \quad (187a)$$

$$= \frac{1}{\delta_{Q,z}^* + \gamma}. \quad (187b)$$

From the equality in (187b) and the fact that Q is a probability measure, it follows that the integral on the right-hand side in (187a) is bounded. Hence, from the dominated convergence theorem [45, Theorem 1.6.9, page 50] the left-hand side of (187a) is finite for all $\beta \in (-\delta_{Q,z}^*, \infty)$. Furthermore, for all $\theta \in \mathcal{M}$ the partial derivative of $\frac{1}{L_z(\theta) + \beta}$ in (187a) with respect to β yields

$$\frac{\partial}{\partial \beta} \left(\frac{1}{L_z(\theta) + \beta} \right) = -\frac{1}{(\beta + L_z(\theta))^2}, \quad (188)$$

which exists for all $\beta \in (-\delta_{Q,z}^*, \infty)$. Then, from the differentiation lemma [48, Theorem 6.28, page 160], the interchange of the integral with the derivative on the right-hand side of (187a) is possible. Hence,

$$\frac{d}{d\beta} \int \frac{1}{L_z(\theta) + \beta} dQ(\theta) = \int \frac{\partial}{\partial \beta} \frac{1}{L_z(\theta) + \beta} dQ(\theta), \quad (189)$$

which completes the second part of the proof.

The third part is as follows.

For all $\beta \in (-\delta_{Q,z}^*, \infty)$, the derivative of the function $\bar{K}_{Q,z}^{-1}$ in (185) satisfies:

$$\bar{K}_{Q,z}^{-1(1)}(\beta) = \frac{d}{d\beta} \left(\int \frac{1}{\beta + L_z(\theta)} dQ(\theta) \right)^{-1} \quad (190a)$$

$$= - \left(\int \frac{1}{\beta + L_z(\theta)} dQ(\theta) \right)^{-2} \frac{d}{d\beta} \left(\bar{K}_{Q,z}^{-1}(\beta) \right)^{-1} \quad (190b)$$

$$= - \left(\int \frac{1}{\beta + L_z(\theta)} dQ(\theta) \right)^2 \frac{d}{d\beta} \left(\bar{K}_{Q,z}^{-1}(\beta) \right)^{-1} \quad (190c)$$

$$= - \frac{\frac{d}{d\beta} \left(\bar{K}_{Q,z}^{-1}(\beta) \right)^{-1}}{\left(\int \frac{1}{\beta + L_z(\theta)} dQ(\theta) \right)^2} \quad (190d)$$

$$= - \frac{\int -\frac{1}{(\beta + L_z(\theta))^2} dQ(\theta)}{\left(\int \frac{1}{\beta + L_z(\theta)} dQ(\theta) \right)^2} \quad (190e)$$

$$= \frac{\int \frac{1}{(\beta + \mathbf{L}_{\mathbf{z}}(\boldsymbol{\theta}))^2} dQ(\boldsymbol{\theta})}{\left(\int \frac{1}{\beta + \mathbf{L}_{\mathbf{z}}(\boldsymbol{\theta})} dQ(\boldsymbol{\theta}) \right)^2}. \quad (190f)$$

where the equality in (190b) follows from (189); and the equality in (190d) follows from (188).

Note that from Jensen's inequality, it follows that

$$\left(\int \frac{1}{\beta + \mathbf{L}_{\mathbf{z}}(\boldsymbol{\theta})} dQ(\boldsymbol{\theta}) \right)^2 \leq \int \frac{1}{(\beta + \mathbf{L}_{\mathbf{z}}(\boldsymbol{\theta}))^2} dQ(\boldsymbol{\theta}). \quad (191)$$

Then, from (190f) and (191), for all $\beta \in (\delta_{Q,\mathbf{z}}^*, \infty)$ it holds that

$$\bar{K}_{Q,\mathbf{z}}^{-1(1)}(\beta) \geq 1. \quad (192)$$

which completes the third part of the proof.

The fourth part is as follows. Let the real values $(\lambda_1, \lambda_2) \in (0, \infty)^2$ be such that $\lambda_2 > \lambda_1$, which implies from Lemma 4.4 that $\bar{K}_{Q,\mathbf{z}}(\lambda_2) > \bar{K}_{Q,\mathbf{z}}(\lambda_1)$. Then, from Lemma 5.1 it follows that

$$\mathbf{R}_{\mathbf{z}}\left(\bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda_2)}\right) - \mathbf{R}_{\mathbf{z}}\left(\bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda_1)}\right) = \lambda_2 - \lambda_1 + \bar{K}_{Q,\mathbf{z}}(\lambda_1) - \bar{K}_{Q,\mathbf{z}}(\lambda_2) \quad (193a)$$

$$= \bar{K}_{Q,\mathbf{z}}^{-1}(\beta_2) - \bar{K}_{Q,\mathbf{z}}^{-1}(\beta_1) + \beta_1 - \beta_2, \quad (193b)$$

where equality (193b) follows from substituting (185) into (193a). Note that (192) implies that

$$\bar{K}_{Q,\mathbf{z}}^{-1}(\beta_2) - \bar{K}_{Q,\mathbf{z}}^{-1}(\beta_1) \geq \beta_2 - \beta_1. \quad (194)$$

Thus, from (193b) and (194) it follows that

$$\mathbf{R}_{\mathbf{z}}\left(\bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda_2)}\right) - \mathbf{R}_{\mathbf{z}}\left(\bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda_1)}\right) \geq 0, \quad (195)$$

Furthermore, from Lemma 4.4 and $\mathbf{L}_{\mathbf{z}}$ in (3), for all $\boldsymbol{\theta} \in \text{supp } Q$ the fraction $(\beta + \mathbf{L}_{\mathbf{z}}(\boldsymbol{\theta}))^{-1}$ is strictly convex. Hence, equality in (192) and in (195) holds if and only if $\mathbf{L}_{\mathbf{z}}$ is nonseparable with respect to Q , which completes the proof. ■

M Proof of Lemma 5.3

Proof: From Theorem 4.1 and the definition of the Type-II relative entropy, it holds that

$$\mathbf{D}\left(Q \parallel \bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}\right) = \int \log \left(\frac{dQ}{d\bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}} \right) dQ(\boldsymbol{\theta}) \quad (196a)$$

$$\leq \log \left(\int \frac{dQ}{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) \right) \quad (196b)$$

$$= \log \left(\int \frac{\bar{K}_{Q,z}(\lambda) + \mathbf{L}_z(\boldsymbol{\theta})}{\lambda} dQ(\boldsymbol{\theta}) \right) \quad (196c)$$

$$= \log \left(\frac{1}{\lambda} \int \bar{K}_{Q,z}(\lambda) + \mathbf{L}_z(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) \right) \quad (196d)$$

$$= \log \left(\frac{1}{\lambda} (\bar{K}_{Q,z}(\lambda) + \mathbf{R}_z(Q)) \right), \quad (196e)$$

where inequality (196b) follows from Jensen's Inequality. From (196e), it follows that

$$\mathbf{R}_z(Q) \geq \exp \left(D \left(Q \| \bar{P}_{\Theta|Z=z}^{(Q,\lambda)} \right) \right) \lambda - \bar{K}_{Q,z}(\lambda). \quad (196f)$$

Hence, the difference between the expected empirical risk of the probability measures $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ and Q , from Lemma 5.1 and (196f), satisfies that

$$\mathbf{R}_z(Q) - \mathbf{R}_z \left(\bar{P}_{\Theta|Z=z}^{(Q,\lambda)} \right) \geq \lambda \left(\exp \left(D \left(Q \| \bar{P}_{\Theta|Z=z}^{(Q,\lambda)} \right) \right) - 1 \right), \quad (197)$$

which completes the proof. \blacksquare

N Proof of Lemma 5.5

Proof: From Lemma 4.4 and Lemma 5.1, it holds that

$$\mathbf{R}_z \left(\bar{P}_{\Theta|Z=z}^{(Q,\lambda)} \right) < \lambda + \delta_{Q,z}^*. \quad (198)$$

Similarly, from the definition of the Rashmon set in (27) and $\delta_{Q,z}^*$ in (28), it follows that

$$\mathbf{R}_z \left(\bar{P}_{\Theta|Z=z}^{(Q,\lambda)} \right) \geq \delta_{Q,z}^*. \quad (199)$$

The proof continues by determining the conditions for which the equality in (199) holds. Assume the empirical risk \mathbf{L}_z in (3) is separable with respect to the probability measure $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ in (16) (see Definition 4.1). Then, there exists a real value $\epsilon > 0$ and two nonnegligible sets \mathcal{A} and \mathcal{B} with respect to the probability measure $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ in (16), such that

$$\mathcal{A} = \{ \boldsymbol{\theta} \in \mathcal{M} : \mathbf{L}_z(\boldsymbol{\theta}) < \delta_{Q,z}^* + \epsilon \}, \text{ and} \quad (200a)$$

$$\mathcal{B} = \{ \boldsymbol{\theta} \in \mathcal{M} : \mathbf{L}_z(\boldsymbol{\theta}) \geq \delta_{Q,z}^* + \epsilon \}. \quad (200b)$$

Under the assumption that \mathbf{L}_z is separable, the expected empirical risk satisfies

$$\mathbf{R}_z \left(\bar{P}_{\Theta|Z=z}^{(Q,\lambda)} \right) = \int_{\mathcal{A}} \mathbf{L}_z(\boldsymbol{\theta}) d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}(\boldsymbol{\theta})$$

$$+ \int_{\mathcal{B}} \mathbf{L}_z(\boldsymbol{\theta}) d\bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}(\boldsymbol{\theta}) \quad (201a)$$

$$\geq \int_{\mathcal{A}} \delta_{Q,z}^* d\bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}(\boldsymbol{\theta}) + \int_{\mathcal{B}} (\delta_{Q,z}^* + \epsilon) d\bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}(\boldsymbol{\theta}) \quad (201b)$$

$$= \delta_{Q,z}^* \bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}(\mathcal{A}) + (\delta_{Q,z}^* + \epsilon) \bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}(\mathcal{B}), \quad (201c)$$

$$= \delta_{Q,z}^* \bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}(\mathcal{A}) + (\delta_{Q,z}^* + \epsilon) \left(1 - \bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}(\mathcal{A})\right) \quad (201d)$$

$$= \delta_{Q,z}^* + \epsilon \left(1 - \bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}(\mathcal{A})\right) \quad (201e)$$

$$> \delta_{Q,z}^*, \quad (201f)$$

where inequality (201b) follows from the fact that $\delta_{Q,z}^*$ and $(\delta_{Q,z}^* + \epsilon)$ are the minimum empirical risk for the set \mathcal{A} and \mathcal{B} respectively; and inequality (201f) follows from the fact that the sets \mathcal{A} and \mathcal{B} are nonnegligible with respect to the probability measure $\bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}$, which implies $\bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}(\mathcal{A}) < 1$. This proves that for separable empirical risk functions \mathbf{L}_z the inequality in (199) is strict.

Considering the case in which the empirical risk \mathbf{L}_z in (3) is not separable with respect to $\bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}$ in (16). Then, for all $\boldsymbol{\theta} \in \text{supp } \bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}$, the empirical risk satisfies $\mathbf{L}_z(\boldsymbol{\theta}) = \delta_{Q,z}^*$, which implies $\mathbf{R}_z\left(\bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}\right) = \delta_{Q,z}^*$. Hence, equality in (199) holds for all nonseparable empirical risk functions \mathbf{L}_z , which completes the proof. ■

O Proof of Theorem 6.1

Proof: Let δ be a real in $(\delta_{Q,z}^*, \infty)$, with $\delta_{Q,z}^*$ in (28). Let also $\gamma \in (0, \infty)$ satisfy the following equality:

$$\mathbf{R}_z\left(\bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\gamma)}\right) \leq \delta. \quad (202)$$

Note that from Lemma 4.4 the function $\bar{K}_{Q,z}$ is continuous and from Lemma 4.4 the set $(0, \infty)$ is convex. Moreover, from Lemma 5.6, it follows that such γ in (202) always exists. From (27), for all $\delta \in (\delta_{Q,z}^*, \infty)$, it holds that

$$\mathcal{L}_z(\delta) \supseteq \mathcal{L}_{Q,z}^*, \quad (203)$$

and thus,

$$\bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\gamma)}(\mathcal{L}_z(\delta)) \geq \bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\gamma)}(\mathcal{L}_{Q,z}^*), \quad (204)$$

with $\mathcal{L}_{Q,z}^*$ defined in (29). Let λ be a positive real such that $\lambda \leq \gamma$, and

$$\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_z(\delta)) > 1 - \epsilon. \quad (205)$$

The existence of such a positive real λ that satisfies (205) follows from Lemma 4.13. From Theorem 4.1 and (203), it follows that

$$\bar{P}_{\Theta|Z=z}^{(Q,\gamma)}(\mathcal{L}_z(\delta)) \geq \bar{P}_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_{Q,z}^*) \quad (206a)$$

Hence, from equality (206a) it holds that

$$1 - \epsilon < \bar{P}_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_{Q,z}^*) \quad (207)$$

$$\leq \bar{P}_{\Theta|Z=z}^{(Q,\lambda)}(\mathcal{L}_z(\delta)). \quad (208)$$

The equality in (208) implies that the probability measure $\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}$ is (δ, ϵ) -optimal (Definition 6.1), which completes the proof. ■

P Proof of Lemma 7.1

Proof: From Theorem 4.1, it follows that for all $\theta \in \mathcal{M}$,

$$\log \left(\frac{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) \right) = \log \left(\frac{\lambda}{\bar{K}_{Q,z}(\lambda) + \mathbf{L}_z(\theta)} \right) \quad (209a)$$

$$= \log(\lambda) - \log(\bar{K}_{Q,z}(\lambda) + \mathbf{L}_z(\theta)) \quad (209b)$$

$$= \log(\lambda) - \mathbf{V}_{z,\lambda}(\theta), \quad (209c)$$

where the function $\mathbf{V}_{z,\lambda}$ is defined in (60). Thus,

$$\mathbf{D}(\bar{P}_{\Theta|Z=z}^{(Q,\lambda)} \| Q) = \int \log \left(\frac{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) \right) d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}(\theta) \quad (210a)$$

$$= \log(\lambda) - \int \mathbf{V}_{z,\lambda}(\theta) d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}(\theta) \quad (210b)$$

$$= \log(\lambda) - \bar{\mathbf{R}}_{z,\lambda}(\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}), \quad (210c)$$

where the functional $\bar{\mathbf{R}}_{z,\lambda}$ is defined in (61). Hence, it follows from (210c) that

$$\log(\lambda) = \bar{\mathbf{R}}_{z,\lambda}(\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}) + \mathbf{D}(\bar{P}_{\Theta|Z=z}^{(Q,\lambda)} \| Q), \quad (211)$$

which completes the proof for (64).

From (209), it follows that

$$\mathbf{D}(Q \| \bar{P}_{\Theta|Z=z}^{(Q,\lambda)}) = - \int \log \left(\frac{d\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}}{dQ}(\theta) \right) dQ(\theta) \quad (212a)$$

$$= -\log(\lambda) + \int \mathbf{V}_{\mathbf{z},\lambda}(\boldsymbol{\theta}) \, dQ(\boldsymbol{\theta}) \quad (212b)$$

$$= -\log(\lambda) + \bar{\mathbf{R}}_{\mathbf{z},\lambda}(Q), \quad (212c)$$

Hence, it follows from (210c) that

$$\log(\lambda) = \bar{\mathbf{R}}_{\mathbf{z},\lambda}(Q) - \mathbf{D}\left(Q \parallel \bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}\right), \quad (213)$$

which completes the proof for (65). \blacksquare

Q Proof of Lemma 7.2

Proof: The proof uses the fact that the probability measure $\bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}$ in (16) is mutually absolutely continuous with the probability measure Q in Theorem 4.1. Hence, the probability measure P is mutually absolutely continuous with $\bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}$, as a consequence of the assumption that $P \in \mathcal{O}_Q(\mathcal{M})$.

The proof follows by noticing that for all $\boldsymbol{\theta} \in \mathcal{M}$,

$$\log\left(\frac{dP}{d\bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}}(\boldsymbol{\theta})\right) = \log\left(\frac{dP}{dQ}(\boldsymbol{\theta}) \frac{dQ}{d\bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}}(\boldsymbol{\theta})\right) \quad (214a)$$

$$= \log\left(\frac{dP}{dQ}(\boldsymbol{\theta})\right) - \log\left(\frac{d\bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}}{dQ}(\boldsymbol{\theta})\right) \quad (214b)$$

$$= \log\left(\frac{dP}{dQ}(\boldsymbol{\theta})\right) - \log\left(\frac{\lambda}{\bar{K}_{Q,\mathbf{z}}(\lambda) + \mathbf{L}_{\mathbf{z}}(\boldsymbol{\theta})}\right) \quad (214c)$$

$$= \log\left(\frac{dP}{dQ}(\boldsymbol{\theta})\right) - \log(\lambda) + \log(\bar{K}_{Q,\mathbf{z}}(\lambda) + \mathbf{L}_{\mathbf{z}}(\boldsymbol{\theta})) \quad (214d)$$

$$= \log\left(\frac{dP}{dQ}(\boldsymbol{\theta})\right) - \log(\lambda) + \mathbf{V}_{\mathbf{z},\lambda}(\boldsymbol{\theta}), \quad (214e)$$

where the functions $\mathbf{L}_{\mathbf{z}}$, $\bar{K}_{Q,\mathbf{z}}$ and $\mathbf{V}_{\mathbf{z},\lambda}$ are defined in (3), (23b) and in (60), respectively; and the equality in (214c) follows from (16). Hence, the relative entropy $\mathbf{D}\left(P \parallel \bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}\right)$ satisfies,

$$\mathbf{D}\left(P \parallel \bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}\right) = \int \log\left(\frac{dP}{d\bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}}(\boldsymbol{\theta})\right) dP(\boldsymbol{\theta}) \quad (215a)$$

$$= \int \left(\log\left(\frac{dP}{dQ}(\boldsymbol{\theta})\right) - \log(\lambda) + \mathbf{V}_{\mathbf{z},\lambda}(\boldsymbol{\theta})\right) dP(\boldsymbol{\theta}) \quad (215b)$$

$$= \int \log\left(\frac{dP}{dQ}(\boldsymbol{\theta})\right) dP(\boldsymbol{\theta}) - \log(\lambda)$$

$$+ \int V_{z,\lambda}(\boldsymbol{\theta}) dP(\boldsymbol{\theta}) \quad (215c)$$

$$= D(P\|Q) - \log(\lambda) + \bar{R}_{z,\lambda}(P) \quad (215d)$$

$$= D(P\|Q) - \bar{R}_{z,\lambda}\left(\bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}\right) - D\left(\bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}\|Q\right) + \bar{R}_{z,\lambda}(P), \quad (215e)$$

where equality (215b) follows from (214); and equality (215e) follows from Lemma 7.1. Thus, from (215), it follows that

$$\begin{aligned} \bar{R}_{z,\lambda}(P) - \bar{R}_{z,\lambda}\left(\bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}\right) &= D\left(P\|\bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}\right) - D(P\|Q) \\ &\quad + D\left(\bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}\|Q\right), \end{aligned} \quad (216)$$

which completes the proof. \blacksquare

R Proof of Lemma 7.4

Proof: The proof is presented in two parts. First, the sensitivity $S_{Q,\lambda}$ ([49, Definition 3.1 page 9]) is evaluated with respect to the probability measure $\bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}$ in (16). Second, the log sensitive defined $\bar{S}_{Q,\lambda}$ in (67) is evaluated with respect to the probability measure $P_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}$ in (9).

For the first part, from Lemma 7.1, for all $\alpha \in (0, \infty)$, it holds that

$$D\left(\bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\alpha)}\|Q\right) = -\bar{R}_{z,\alpha}\left(\bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\alpha)}\right) + \log(\alpha), \quad (217)$$

where the functional $\bar{R}_{z,\alpha}$ is defined in (61).

Similarly, from [49, Lemma 2.2, page 8], for all $\lambda \in (0, \infty)$, it holds that

$$R_z\left(P_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}\right) = -\lambda\left(D\left(P_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}\|Q\right) + K_{Q,z}\left(-\frac{1}{\lambda}\right)\right), \quad (218)$$

with the functional R_z defined in (6).

From [49, Theorem 3.1 page 9] the sensitivity of the probability measure $\bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\alpha)}$ satisfies

$$\begin{aligned} S_{Q,\lambda}\left(z, \bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\alpha)}\right) &= R_z\left(\bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\alpha)}\right) - R_z\left(P_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}\right) \\ &= \lambda\left(D\left(\bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\alpha)}\|P_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}\right) - D\left(\bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\alpha)}\|Q\right) + D\left(P_{\boldsymbol{\Theta}|Z=z}^{(Q,\lambda)}\|Q\right)\right). \end{aligned} \quad (219a)$$

Plugging (217) and (218) into (219) yields

$$\frac{1}{\lambda}R_z\left(\bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\alpha)}\right) - \bar{R}_{z,\alpha}\left(\bar{P}_{\boldsymbol{\Theta}|Z=z}^{(Q,\alpha)}\right)$$

$$= \mathcal{D}\left(\bar{P}_{\Theta|Z=z}^{(Q,\alpha)} \| P_{\Theta|Z=z}^{(Q,\lambda)}\right) - \log(\alpha) - K_{Q,z}\left(-\frac{1}{\lambda}\right), \quad (220)$$

which completes the first part of the proof.

For the second part, from Lemma 7.1, for all $\alpha \in (0, \infty)$, it holds that

$$\bar{R}_{z,\alpha}\left(\bar{P}_{\Theta|Z=z}^{(Q,\alpha)}\right) = -\mathcal{D}\left(\bar{P}_{\Theta|Z=z}^{(Q,\alpha)} \| Q\right) + \log(\alpha), \quad (221)$$

where the functional $\bar{R}_{z,\alpha}$ is defined in (61).

Similarly, from [49, Lemma 2.2, page 8], for all $\lambda \in (0, \infty)$, it holds that

$$\mathcal{D}\left(P_{\Theta|Z=z}^{(Q,\lambda)} \| Q\right) = \frac{1}{\lambda} R_z\left(P_{\Theta|Z=z}^{(Q,\lambda)}\right) + K_{Q,z}\left(-\frac{1}{\lambda}\right), \quad (222)$$

with the functional R_z defined in (6).

From Lemma 7.2 the logarithmic sensitivity of the probability measure $P_{\Theta|Z=z}^{(Q,\lambda)}$ is

$$\begin{aligned} & \bar{R}_{z,\alpha}\left(P_{\Theta|Z=z}^{(Q,\lambda)}\right) - \bar{R}_{z,\alpha}\left(\bar{P}_{\Theta|Z=z}^{(Q,\lambda)}\right) \\ &= \mathcal{D}\left(P_{\Theta|Z=z}^{(Q,\lambda)} \| \bar{P}_{\Theta|Z=z}^{(Q,\alpha)}\right) - \mathcal{D}\left(P_{\Theta|Z=z}^{(Q,\lambda)} \| Q\right) + \mathcal{D}\left(\bar{P}_{\Theta|Z=z}^{(Q,\alpha)} \| Q\right). \end{aligned} \quad (223)$$

Plugging (221) and (222) into (223) yields

$$\begin{aligned} & \frac{1}{\lambda} R_z\left(P_{\Theta|Z=z}^{(Q,\lambda)}\right) - \bar{R}_{z,\alpha}\left(P_{\Theta|Z=z}^{(Q,\lambda)}\right) \\ &= -\mathcal{D}\left(P_{\Theta|Z=z}^{(Q,\lambda)} \| \bar{P}_{\Theta|Z=z}^{(Q,\alpha)}\right) - \left(\log(\alpha) + K_{Q,z}\left(-\frac{1}{\lambda}\right)\right). \end{aligned} \quad (224)$$

The proof proceeds by subtracting (224) from (220), resulting in

$$\begin{aligned} & \frac{1}{\lambda} S_{Q,\lambda}\left(z, \bar{P}_{\Theta|Z=z}^{(Q,\alpha)}\right) - \bar{S}_{Q,\alpha}\left(z, P_{\Theta|Z=z}^{(Q,\lambda)}\right) \\ &= \mathcal{D}\left(\bar{P}_{\Theta|Z=z}^{(Q,\alpha)} \| P_{\Theta|Z=z}^{(Q,\lambda)}\right) - \mathcal{D}\left(P_{\Theta|Z=z}^{(Q,\lambda)} \| \bar{P}_{\Theta|Z=z}^{(Q,\alpha)}\right) \\ &+ 2\left(\log(\alpha) + K_{Q,z}\left(-\frac{1}{\lambda}\right)\right), \end{aligned} \quad (225)$$

where the functions $S_{Q,\lambda}$ and $\bar{S}_{Q,\alpha}$ are respectively defined in [49, Definition 8] and (67). From [34, Theorem 1] and Lemma 7.4, it follows that

$$\begin{aligned} & \frac{1}{\lambda} S_{Q,\lambda}\left(z, \bar{P}_{\Theta|Z=z}^{(Q,\alpha)}\right) - \bar{S}_{Q,\alpha}\left(z, P_{\Theta|Z=z}^{(Q,\lambda)}\right) \\ &= 2\left(\mathcal{D}\left(P_{\Theta|Z=z}^{(Q,\lambda)} \| Q\right) - \mathcal{D}\left(\bar{P}_{\Theta|Z=z}^{(Q,\lambda)} \| Q\right)\right) + \mathcal{D}\left(\bar{P}_{\Theta|Z=z}^{(Q,\alpha)} \| P_{\Theta|Z=z}^{(Q,\lambda)}\right) \\ &- \mathcal{D}\left(P_{\Theta|Z=z}^{(Q,\lambda)} \| \bar{P}_{\Theta|Z=z}^{(Q,\alpha)}\right). \end{aligned} \quad (226)$$

Substituting (226) into (225) yields

$$\mathcal{D}\left(P_{\Theta|Z=z}^{(Q,\lambda)} \| Q\right) - \mathcal{D}\left(\bar{P}_{\Theta|Z=z}^{(Q,\lambda)} \| Q\right) = \log(\alpha) + K_{Q,z}\left(-\frac{1}{\lambda}\right), \quad (227)$$

which completes the proof. \blacksquare

S Example 4.3

Proof. The Lagrangian multiplier β for the optimization problem in (11) satisfies

$$\int \frac{\lambda}{\beta + \mathbf{L}_z(\boldsymbol{\nu})} dQ(\boldsymbol{\nu}) = 1. \quad (228)$$

From the empirical risk function $\mathbf{L}_z : \mathcal{M} \rightarrow \mathbb{R}_0^+$ in (41), which is a simple function, and the probability measure Q in (40a), it holds that

$$\int \frac{\lambda}{\beta + \mathbf{L}_z(\boldsymbol{\nu})} dQ(\boldsymbol{\nu}) = \lambda \left(\frac{1}{\beta + c_0} Q(\mathcal{T}(z)) + \frac{1}{\beta + c_1} Q(\mathcal{M} \setminus \mathcal{T}(z)) \right) \quad (229a)$$

$$= \lambda \left(\frac{1}{\beta + c_0} Q(\mathcal{T}(z)) + \frac{1}{\beta + c_1} (1 - Q(\mathcal{T}(z))) \right) \quad (229b)$$

$$= \lambda \left(\frac{(\beta + c_1)Q(\mathcal{T}(z)) + (\beta + c_0)(1 - Q(\mathcal{T}(z)))}{\beta^2 + \beta(c_0 + c_1) + c_0c_1} \right) \quad (229c)$$

$$= \lambda \left(\frac{(c_1 - c_0)Q(\mathcal{T}(z)) + \beta + c_0}{\beta^2 + \beta(c_0 + c_1) + c_0c_1} \right). \quad (229d)$$

From (228) and (229d), it follows that

$$0 = \beta^2 + \beta(c_0 + c_1) + c_0c_1 - \lambda((c_1 - c_0)Q(\mathcal{T}(z)) + \beta + c_0) \quad (230a)$$

$$= \beta^2 + \beta(c_0 + c_1 - \lambda) + c_0c_1 - \lambda c_0 - \lambda(c_1 - c_0)Q(\mathcal{T}(z)). \quad (230b)$$

From (230b) and the fact that $c_0 = 0$ in equation (41), it holds that

$$0 = \beta^2 + \beta(c_1 - \lambda) - \lambda c_1 Q(\mathcal{T}(z)). \quad (230c)$$

Observe that the equality in (230c) is a quadratic polynomial that has two roots r_1 and r_2 . Hence, (230c) in terms of r_1 and r_2 satisfies

$$0 = \beta^2 - (r_1 + r_2)\beta + r_1r_2 \quad (231a)$$

$$= (\beta - r_1)(\beta - r_2), \quad (231b)$$

where the roots r_1 and r_2 are given by the quadratic formula such that

$$r_1 = -\frac{(c_1 - \lambda)}{2} - \sqrt{\left(\frac{c_1 - \lambda}{2}\right)^2 + \lambda c_1 Q(\mathcal{T}(z))}, \quad (232a)$$

and

$$r_2 = -\frac{(c_1 - \lambda)}{2} + \sqrt{\left(\frac{c_1 - \lambda}{2}\right)^2 + \lambda c_1 Q(\mathcal{T}(z))}. \quad (232b)$$

The proof continues by verifying that the roots in (232a) and (232b) are real and there is only one positive root for all $\lambda \in (0, +\infty)$ and for all $Q(\mathcal{T}(z)) \in [0, 1]$.

Note that for all $c_1 \in (0, \infty)$ and for all $\lambda \in [0, +\infty)$, it holds that

$$-\frac{c_1 - \lambda}{2} \leq \left| \frac{c_1 - \lambda}{2} \right| \quad (233)$$

$$= \sqrt{\left(\frac{c_1 - \lambda}{2} \right)^2} \quad (234)$$

$$\leq \sqrt{\left(\frac{c_1 - \lambda}{2} \right)^2 + \lambda c_1 Q(\mathcal{T}(\mathbf{z}))}. \quad (235)$$

Observe that for all $Q(\mathcal{T}(\mathbf{z})) \in [0, 1)$, $c_1 \in (0, \infty)$ and $\lambda \in [0, \infty)$ the expressions $\left(\frac{c_1 - \lambda}{2} \right)^2$ and $\lambda c_1 Q(\mathcal{T}(\mathbf{z}))$ are always positive. Thus, the square roots in (232a) and (232b) are real, which implies that r_1 and r_2 are real. From (235), for all $\lambda \in [0, +\infty)$ and for all $Q(\mathcal{T}(\mathbf{z})) \in [0, 1)$, it holds

$$r_1 < 0; \quad (236a)$$

and following the same arguments

$$r_2 > 0. \quad (236b)$$

Hence, the solution for the Lagrange Multiplier β that satisfies (228) given the empirical risk function $\mathbf{L}_{\mathbf{z}}$ in (41) and the probability measure Q in (40a) is

$$\beta = -\frac{(c_1 - \lambda)}{2} + \sqrt{\left(\frac{c_1 - \lambda}{2} \right)^2 + \lambda c_1 Q(\mathcal{T}(\mathbf{z}))}, \quad (237)$$

which completes the proof. \square

T Example 4.2

Proof. Consider the Type-II ERM-RER problem in (11) and assume that: (a) $\lambda = 0.5$; (b) $\mathcal{M} = \mathcal{X} = \mathcal{Y} = [0, \infty)$; (c) $\mathbf{z} = ((1, 1))$ and (d) $Q \ll \mu$, with μ the Lebesgue measure, such that for all $\boldsymbol{\theta} \in \text{supp } Q$,

$$\frac{dQ}{d\mu}(\boldsymbol{\theta}) = 4\boldsymbol{\theta}^2 \exp(-2\boldsymbol{\theta}). \quad (238)$$

Let also the function $f : \mathcal{M} \times \mathcal{X} \rightarrow \mathcal{Y}$ be

$$f(\boldsymbol{\theta}, x) = x\boldsymbol{\theta}, \quad (239)$$

and the loss function ℓ in (2) be

$$\ell(f(\boldsymbol{\theta}, x), y) = (x\boldsymbol{\theta} - y)^2. \quad (240)$$

Hence, from the fact that \mathbf{z} is a single data point and the definition of $\mathbf{L}_{\mathbf{z}}$ in (3), it follows that

$$\mathbf{L}_{\mathbf{z}}(\boldsymbol{\theta}) = \ell(f(\boldsymbol{\theta}, x), y) \quad (241a)$$

$$= (x\boldsymbol{\theta} - y)^2. \quad (241b)$$

Based on the previous assumptions regarding the functions ℓ and f , as well as the sets \mathcal{M} , \mathcal{X} , \mathcal{Y} , and the measure Q , the definition of $\delta_{Q,\mathbf{z}}^*$ in (28) implies that $\delta_{Q,\mathbf{z}}^* = 0$.

From Theorem 4.1, the solution to the Type-II ERM-RER problem in (11) is the probability measure $\bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}$ that is mutually absolutely continuous with respect to Q , which implies that under the above assumptions $\bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}$ is absolutely continuous with respect to the Lebesgue measure μ such that

$$\frac{d\bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}}{d\mu}(\boldsymbol{\theta}) = \frac{\lambda}{\mathbf{L}_{\mathbf{z}}(\boldsymbol{\theta}) + \beta} 4\boldsymbol{\theta}^2 \exp(-2\boldsymbol{\theta}). \quad (242)$$

Hence, it follows that

$$\bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}(\mathcal{M}) = \int_{\mathcal{M}} \frac{d\bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}}{dQ}(\boldsymbol{\theta}) dQ(\boldsymbol{\theta}) \quad (243a)$$

$$= \int_{\mathcal{M}} \frac{d\bar{P}_{\boldsymbol{\Theta}|\mathbf{Z}=\mathbf{z}}^{(Q,\lambda)}}{d\mu}(\boldsymbol{\theta}) d\mu(\boldsymbol{\theta}) \quad (243b)$$

$$= \int_0^\infty \frac{\lambda 4\boldsymbol{\theta}^2 \exp(-2\boldsymbol{\theta})}{(x\boldsymbol{\theta} - y)^2 + \beta} d\boldsymbol{\theta} \quad (243c)$$

$$= 4\lambda \int_0^\infty \frac{\boldsymbol{\theta}^2 \exp(-2\boldsymbol{\theta})}{(\boldsymbol{\theta} - 1)^2} d\boldsymbol{\theta} \quad (243d)$$

$$(243e)$$

The integral in (243d) can be rewritten as

$$\begin{aligned} & \int_0^\infty \frac{\boldsymbol{\theta}^2 \exp(-2\boldsymbol{\theta})}{(\boldsymbol{\theta} - 1)^2} d\boldsymbol{\theta} \\ &= \int_0^\infty \left(\frac{\boldsymbol{\theta} \exp(-2\boldsymbol{\theta}) + \frac{1}{2} \exp(-2\boldsymbol{\theta})}{\boldsymbol{\theta} - 1} + \frac{\frac{1}{2} \boldsymbol{\theta} \exp(-2\boldsymbol{\theta}) + \frac{1}{2} \exp(-2\boldsymbol{\theta})}{(\boldsymbol{\theta} - 1)^2} \right) d\boldsymbol{\theta} \quad (244) \\ &= \frac{1}{2} \int_0^\infty \frac{2\boldsymbol{\theta} \exp(-2\boldsymbol{\theta}) + \exp(-2\boldsymbol{\theta})}{\boldsymbol{\theta} - 1} d\boldsymbol{\theta} \\ & \quad + \frac{1}{2} \int_0^\infty \frac{\boldsymbol{\theta} \exp(-2\boldsymbol{\theta}) + \exp(-2\boldsymbol{\theta})}{(\boldsymbol{\theta} - 1)^2} d\boldsymbol{\theta}. \end{aligned} \quad (245)$$

Using integration by parts on the second integral in (245) yields

$$\int_0^\infty \frac{\boldsymbol{\theta} \exp(-2\boldsymbol{\theta}) + \exp(-2\boldsymbol{\theta})}{(\boldsymbol{\theta} - 1)^2} d\boldsymbol{\theta}$$

$$= -\frac{(\boldsymbol{\theta} + 1) \exp(-2\boldsymbol{\theta})}{(\boldsymbol{\theta} - 1)} - \int_0^\infty \frac{2\boldsymbol{\theta} \exp(-2\boldsymbol{\theta}) + \exp(-2\boldsymbol{\theta})}{\boldsymbol{\theta} - 1} d\boldsymbol{\theta}. \quad (246)$$

Plugging (246) into (245) yields

$$\int_0^\infty \frac{\boldsymbol{\theta}^2 \exp(-2\boldsymbol{\theta})}{(\boldsymbol{\theta} - 1)^2} d\boldsymbol{\theta} = \left[-\frac{(\boldsymbol{\theta} + 1) \exp(-2\boldsymbol{\theta})}{2(\boldsymbol{\theta} - 1)} \right]_0^\infty. \quad (247)$$

From (247), it follows that

$$4\lambda \int_0^\infty \frac{\boldsymbol{\theta}^2 \exp(-2\boldsymbol{\theta})}{(\boldsymbol{\theta} - 1)^2} d\boldsymbol{\theta} = 4\lambda \left[-\frac{(\boldsymbol{\theta} + 1) \exp(-2\boldsymbol{\theta})}{2(\boldsymbol{\theta} - 1)} \right]_0^\infty \quad (248a)$$

$$= 4\lambda \left[-\frac{(\boldsymbol{\theta} + 1) \exp(-2\boldsymbol{\theta})}{2(\boldsymbol{\theta} - 1)} \right]_0^1 + 4\lambda \left[-\frac{(\boldsymbol{\theta} + 1) \exp(-2\boldsymbol{\theta})}{2(\boldsymbol{\theta} - 1)} \right]_1^\infty \quad (248b)$$

$$\geq 4\lambda \left[-\frac{(\boldsymbol{\theta} + 1) \exp(-2\boldsymbol{\theta})}{2(\boldsymbol{\theta} - 1)} \right]_0^1 \quad (248c)$$

$$= \lim_{a \rightarrow 1^-} 4\lambda \left[-\frac{(a + 1) \exp(-2a)}{2(a - 1)} \right]_0^a \quad (248d)$$

$$= \lim_{a \rightarrow 1^-} -4\lambda \frac{(a + 1) \exp(-2a)}{2(a - 1)} - 4\lambda \frac{(0 + 1) \exp(-0)}{2(0 - 1)} \quad (248e)$$

$$= \lim_{a \rightarrow 1^-} -4\lambda \frac{(a + 1) \exp(-2a)}{2(a - 1)} - 2\lambda \quad (248f)$$

$$= \infty, \quad (248g)$$

which completes the proof. \square

References

- [1] V. Vapnik, “Principles of risk minimization for learning theory,” *Advances in Neural Information Processing Systems*, vol. 4, pp. 831–838, Jan. 1992.
- [2] V. Vapnik and A. Y. Chervonenkis, “On a perceptron class,” *Avtomatika i Telemekhanika*, vol. 25, no. 1, pp. 112–120, Feb. 1964.
- [3] M. R. Rodrigues and Y. C. Eldar, *Information-theoretic Methods in Data Science*, 1st ed. Cambridge, UK: Cambridge University Press, 2021.
- [4] M. Mezard and A. Montanari, *Information, Physics, and Computation*, 1st ed. New York, NY, USA: Oxford University Press, 2009.
- [5] M. J. Wainwright, *High-dimensional Statistics: A Non-asymptotic Viewpoint*, 1st ed. New York, NY, USA: Cambridge University Press, 2019.

- [6] R. Vershynin, *High-dimensional probability: An Introduction with Applications in Data Science*, 1st ed. New York, NY, USA: Cambridge University Press, 2018.
- [7] A. Blumer, A. Ehrenfeucht, D. Haussler, and M. K. Warmuth, “Learnability and the Vapnik-Chervonenkis dimension,” *Journal of the ACM*, vol. 36, no. 4, pp. 929–965, Oct. 1989.
- [8] I. Guyon, V. Vapnik, B. Boser, L. Bottou, and S. A. Solla, “Structural risk minimization for character recognition,” *Advances in Neural Information Processing Systems*, vol. 4, Dec. 1991.
- [9] G. Lugosi and K. Zeger, “Nonparametric estimation via empirical risk minimization,” *IEEE Transactions on Information Theory*, vol. 41, no. 3, pp. 677–687, May 1995.
- [10] P. L. Bartlett, “The sample complexity of pattern classification with neural networks: the size of the weights is more important than the size of the network,” *IEEE Transactions on Information Theory*, vol. 44, no. 2, pp. 525–536, Mar. 1998.
- [11] V. Vapnik and L. Bottou, “Local algorithms for pattern recognition and dependencies estimation,” *Neural Computation*, vol. 5, no. 6, pp. 893–909, Nov. 1993.
- [12] V. Cherkassky, X. Shao, F. M. Mulier, and V. N. Vapnik, “Model complexity control for regression using VC generalization bounds,” *IEEE Transactions on Neural Networks*, vol. 10, no. 5, pp. 1075–1089, Sep. 1999.
- [13] V. Vapnik, “An overview of statistical learning theory,” *IEEE Transactions on Neural Networks*, vol. 10, no. 5, pp. 988–999, Sep. 1999.
- [14] L. Bottou, F. E. Curtis, and J. Nocedal, “Optimization methods for large-scale machine learning,” *SIAM Review*, vol. 60, no. 2, pp. 223–311, May 2018.
- [15] A. Krzyżak, T. Linder, and C. Lugosi, “Nonparametric estimation and classification using radial basis function nets and empirical risk minimization,” *IEEE Transactions on Neural Networks*, vol. 7, no. 2, pp. 475–487, Mar. 1996.
- [16] W. Deng, Q. Zheng, and L. Chen, “Regularized extreme learning machine,” in *Proceedings of the IEEE Symposium on Computational Intelligence in Data Mining (CIDM)*, Nashville, TN, USA, Apr. 2009, pp. 389–395.
- [17] D. Arpit, S. Jastrzębski, N. Ballas, D. Krueger, E. Bengio, M. S. Kanwal, T. Maharaj, A. Fischer, A. Courville, Y. Bengio, and S. Lacoste-Julien, “A closer look at memorization in deep networks,” in *Proceedings of the 34th International Conference on Machine Learning (ICML)*, vol. 70, Aug. 2017, pp. 233–242.

- [18] F. Hellström, G. Durisi, B. Guedj, and M. Raginsky, “Generalization bounds: Perspectives from information theory and PAC-Bayes,” arXiv preprint arXiv:2309.04381, Sep. 2023.
- [19] O. Bousquet and A. Elisseeff, “Stability and generalization,” *The Journal of Machine Learning Research*, vol. 2, no. 1, pp. 499–526, Mar. 2002.
- [20] V. N. Vapnik and A. Y. Chervonenkis, “On the uniform convergence of relative frequencies of events to their probabilities,” *Measures of Complexity: Festschrift for Alexey Chervonenkis*, vol. 16, no. 2, pp. 11–30, Oct. 2015.
- [21] G. Aminian, Y. Bu, L. Toni, M. Rodrigues, and G. Wornell, “An exact characterization of the generalization error for the Gibbs algorithm,” *Advances in Neural Information Processing Systems*, vol. 34, pp. 8106–8118, Dec. 2021.
- [22] C. P. Robert, *The Bayesian Choice: From Decision-theoretic Foundations to Computational Implementation*, 1st ed. New York, NY, USA: Springer, 2007.
- [23] D. A. McAllester, “Some PAC-Bayesian theorems,” in *Proceedings of the 11th Annual Conference on Computational Learning Theory (COLT)*, Madison, WI, USA, Jul. 1998, pp. 230–234.
- [24] L. G. Valiant, “A theory of the learnable,” *Communications of the ACM*, vol. 27, no. 11, pp. 1134–1142, Nov. 1984.
- [25] J. Shawe-Taylor and R. C. Williamson, “A PAC analysis of a Bayesian estimator,” in *Proceedings of the 10th Annual Conference on Computational Learning Theory (COLT)*, Nashville, TN, USA, Jul. 1997, pp. 2–9.
- [26] D. Cullina, A. N. Bhagoji, and P. Mittal, “PAC-learning in the presence of adversaries,” *Advances in Neural Information Processing Systems*, vol. 31, no. 1, pp. 1–12, Dec. 2018.
- [27] S. M. Perlaza, G. Bisson, I. Esnaola, A. Jean-Marie, and S. Rini, “Empirical risk minimization with relative entropy regularization: Optimality and sensitivity,” in *Proceedings of the IEEE International Symposium on Information Theory (ISIT)*, Espoo, Finland, Jul. 2022, pp. 684–689.
- [28] —, “Empirical risk minimization with relative entropy regularization,” INRIA, Centre Inria d’Université Côte d’Azur, Sophia Antipolis, France, Tech. Rep. RR-9454, Feb. 2022.
- [29] X. Zou, S. M. Perlaza, I. Esnaola, and E. Altman, “The worst-case data-generating probability measure,” INRIA, Centre Inria d’Université Côte d’Azur, Sophia Antipolis, France, Tech. Rep. RR-9515, Aug. 2023.
- [30] M. Raginsky, A. Rakhlin, M. Tsao, Y. Wu, and A. Xu, “Information-theoretic analysis of stability and bias of learning algorithms,” in *Proceedings of the IEEE Information Theory Workshop (ITW)*, Cambridge, UK, Sep. 2016, pp. 26–30.

- [31] D. Russo and J. Zou, “How much does your data exploration overfit? Controlling bias via information usage,” *IEEE Transactions on Information Theory*, vol. 66, no. 1, pp. 302–323, Jan. 2019.
- [32] B. Zou, L. Li, and Z. Xu, “The generalization performance of ERM algorithm with strongly mixing observations,” *Machine Learning*, vol. 75, no. 3, pp. 275–295, Feb. 2009.
- [33] G. Aminian, Y. Bu, L. Toni, M. R. D. Rodrigues, and G. W. Wornell, “Information-theoretic characterizations of generalization error for the Gibbs algorithm,” arXiv preprint arXiv:2210.09864, Oct. 2022.
- [34] S. M. Perlaza, I. Esnaola, G. Bisson, and H. V. Poor, “On the validation of Gibbs algorithms: Training datasets, test datasets and their aggregation,” in *Proceedings of the IEEE International Symposium on Information Theory (ISIT)*, Taipei, Taiwan, Jun. 2023.
- [35] X. Wang and Q. He, “Enhancing generalization capability of SVM classifiers with feature weight adjustment,” in *Proceedings of the Knowledge-Based Intelligent Information and Engineering Systems: 8th International Conference (KES)*, Wellington, New Zealand, Sep. 2004, pp. 1037–1043.
- [36] Q. Lin, Z. Lu, and L. Xiao, “An accelerated proximal coordinate gradient method and its application to regularized empirical risk minimization,” arXiv preprint arXiv:1407.1296, Jul. 2014.
- [37] X. Yang and D. Li, “Estimation of the empirical risk-return relation: A generalized-risk-in-mean model,” *Journal of Time Series Analysis*, vol. 43, no. 6, pp. 938–963, May 2022.
- [38] M. Teboulle, “Entropic proximal mappings with applications to nonlinear programming,” *Mathematics of Operations Research*, vol. 17, no. 3, pp. 670–690, Aug. 1992.
- [39] A. Beck and M. Teboulle, “Mirror descent and nonlinear projected sub-gradient methods for convex optimization,” *Operations Research Letters*, vol. 31, no. 3, pp. 167–175, Jan. 2003.
- [40] P. Alquier, “Non-exponentially weighted aggregation: regret bounds for unbounded loss functions,” in *Proceedings of the 38th International Conference on Machine Learning (ICML)*, vol. 139, Jul. 2021, pp. 207–218.
- [41] J. Liao, O. Kosut, L. Sankar, and F. P. Calmon, “A tunable measure for information leakage,” in *Proceedings of the IEEE International Symposium on Information Theory (ISIT)*, Vail, CO, USA, Jun. 2018, pp. 701–705.
- [42] T. Sypherd, M. Diaz, L. Sankar, and P. Kairouz, “A tunable loss function for binary classification,” in *Proceedings of the IEEE international symposium on information theory (ISIT)*, Paris, France, Jul. 2019, pp. 2479–2483.

- [43] G. R. Kurri, T. Sypherd, and L. Sankar, “Realizing GANs via a tunable loss function,” in *Proceedings of the IEEE Information Theory Workshop (ITW)*, virtual conference, 2021, pp. 1–6.
- [44] H. Hsu and F. Calmon, “Rashomon capacity: A metric for predictive multiplicity in classification,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 28 988–29 000, Dec. 2022.
- [45] R. B. Ash and C. A. Doleans-Dade, *Probability and Measure Theory*, 2nd ed. Burlington, MA, USA: Academic Press, 2000.
- [46] D. G. Luenberger, *Optimization by Vector Space Methods*, 1st ed. New York, NY, USA: Wiley, 1997.
- [47] R. G. Bartle and D. R. Sherbert, *Introduction to Real Analysis*, 2nd ed. New York, NY, USA: Wiley New York, 2000.
- [48] A. Klenke, *Probability Theory: A Comprehensive Course*, 3rd ed. New York, NY, USA: Springer, 2020.
- [49] S. M. Perlaza, I. Esnaola, G. Bisson, and H. V. Poor, “Sensitivity of the Gibbs algorithm to data aggregation in supervised machine learning,” INRIA, Centre Inria d’Université Côte d’Azur, Sophia Antipolis, France, Tech. Rep. RR-9474, Jun. 2022.



**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