



HAL
open science

An Exact Characterization of the Generalization Error of Machine Learning Algorithms

Xinying Zou, Samir M. Perlaza, Iñaki Esnaola, Eitan Altman, H. Vincent Poor

► **To cite this version:**

Xinying Zou, Samir M. Perlaza, Iñaki Esnaola, Eitan Altman, H. Vincent Poor. An Exact Characterization of the Generalization Error of Machine Learning Algorithms. 2024. hal-04426500v1

HAL Id: hal-04426500

<https://inria.hal.science/hal-04426500v1>

Preprint submitted on 2 Feb 2024 (v1), last revised 14 May 2024 (v2)

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

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



Distributed under a Creative Commons Attribution 4.0 International License

An Exact Characterization of the Generalization Error of Machine Learning Algorithms

Xinying Zou*, Samir M. Perlaza*^{‡§}, Iñaki Esnaola^{†‡}, Eitan Altman*[¶] and H. Vincent Poor[‡]

*INRIA, Centre Inria d'Université Côte d'Azur, Sophia Antipolis 06902, France.

[†]Department of Automatic Control and Systems Engineering, University of Sheffield, Sheffield S1 3JD, UK.

[‡]Department of Electrical and Computer Engineering, Princeton University, Princeton NJ 08544, USA.

[§]GAATI Laboratory, Université de la Polynésie Française, Faaa 98702, French Polynesia.

[¶]The Laboratoire d'Informatique d'Avignon (LIA), Université d'Avignon, France.

Abstract—In this paper, the notion of worst-case data-generating (WCDG) probability measure, recently introduced in [1], is leveraged to obtain an exact expression of the expected generalization gap (or generalization error) for any machine learning algorithm. This exact expression is provided in terms of information measures involving the WCDG probability measure and leads to an upper bound on the generalization error that is equal to the sum of the mutual information and the lautum information between the models and the datasets, up to a constant factor. This upper bound is achieved by a Gibbs algorithm whose parameters satisfy particular conditions. Finally, given a fixed model, it is shown that the empirical risk is a sub-Gaussian random variable when datasets are sampled from the WCDG probability measure. This observation leads to the construction of new generalization guarantees coined (ϵ, δ) -robustness.

I. INTRODUCTION

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 $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$, a 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 $\theta \in \mathcal{M}$ is $y = f(\theta, x)$. Let also the function $\hat{\ell} : \mathcal{Y} \times \mathcal{Y} \rightarrow [0, +\infty]$ be such that given a data point $(x, y) \in \mathcal{X} \times \mathcal{Y}$, the loss induced by a model $\theta \in \mathcal{M}$ is $\hat{\ell}(f(\theta, x), y)$. In the following, the loss function $\hat{\ell}$ is assumed to be nonnegative and such that for all $y \in \mathcal{Y}$, $\hat{\ell}(y, y) = 0$. For ease of notation, let the function $\ell : \mathcal{M} \times \mathcal{X} \times \mathcal{Y} \rightarrow [0, +\infty]$ be such that

$$\ell(\theta, x, y) = \hat{\ell}(f(\theta, x), y). \quad (2)$$

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

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

where the function ℓ is defined in (2). Given a model $\theta \in \mathcal{M}$, the expected loss induced by a measure $P \in \Delta(\mathcal{X} \times \mathcal{Y})$ is defined as follows.

Definition 1: Let P be a probability measure in $\Delta(\mathcal{X} \times \mathcal{Y})$. The expected loss with respect to a fixed model $\theta \in \mathcal{M}$ induced by the measure P is

$$R_\theta(P) = \int \ell(\theta, x, y) dP(x, y), \quad (4)$$

where the function ℓ is defined in (2).

Given a probability measure $P_S \in \Delta(\mathcal{X} \times \mathcal{Y})$, which can be interpreted as a *prior* on the set of data points, and a model $\theta \in \mathcal{M}$, studying the worst-case data-generating (WCDG) probability measure leads to two different optimization problems. Firstly, a *neighborhood* of candidate data-generating probability measures is constructed as the set of probability measures that satisfy a constraint expressed in terms of a *statistical distance* $d : \Delta(\mathcal{X} \times \mathcal{Y}) \times \Delta(\mathcal{X} \times \mathcal{Y}) \rightarrow [0, +\infty]$. This *neighborhood* describes the set of plausible probability measures from which test datasets might be sampled from. This work focuses on the case in which the statistical distance is the relative entropy with respect to P_S . Thus, the WCDG measure is the probability measure that maximizes the expectation of the loss within this neighborhood. This view leads to the following problem, c.f. [1] and [2]:

$$\max_{P \in \Delta_{P_S}(\mathcal{X} \times \mathcal{Y})} R_\theta(P) \quad (5a)$$

$$\text{s.t.} \quad D(P \| P_S) \leq \gamma, \quad (5b)$$

where $\gamma > 0$ determines the neighborhood around P_S as the set $\{P \in \Delta(\mathcal{X} \times \mathcal{Y}) : D(P \| P_S) \leq \gamma\}$ and the functional R_θ is defined in (4).

Secondly, the WCDG probability measure can also be interpreted as the probability measure that trades off the maximization of the expectation of the loss and the minimization of the statistical distance with respect to P_S . This alternative perspective yields the optimization problem:

$$\max_{P \in \Delta_{P_S}(\mathcal{X} \times \mathcal{Y})} R_\theta(P) - \beta D(P \| P_S), \quad (6a)$$

where the functional R_θ is defined in (4); and $\beta > 0$ determines the trade-off between maximization of the expectation of the loss and the statistical distance to P_S .

The optimization problems in (5) and (6) exhibit a major difference in the following case. If for all $P \in \Delta_{P_S}$, the functional R_θ in (4) is such that $R_\theta(P) = c$, for some fixed $c \geq 0$, then all the measures in $\{P \in \Delta(\mathcal{X} \times \mathcal{Y}) : D(P||P_S) \leq \gamma\}$ are solutions to the optimization problem in (5). Alternatively, in this case P_S is the unique solution to the problem in (6). Although this difference holds substantial mathematical importance, its practical relevance is limited. This observation is explained by the notion of nonseparable loss functions, which is analogous to [3, Definition 4.1]. When the function ℓ in (2) is nonseparable with respect to the reference measure P_S , it is a constant almost surely with respect to such a measure. More specifically, there exists a real $a \geq 0$, such that $P_S(\{(x, y) \in \mathcal{X} \times \mathcal{Y} : \ell(\theta, x, y) = a\}) = 1$. As a consequence, for all probability measures $P \in \Delta_{P_S}(\mathcal{X} \times \mathcal{Y})$, it holds that $P(\{(x, y) \in \mathcal{X} \times \mathcal{Y} : \ell(\theta, x, y) = a\}) = 1$. If this pathological case of nonseparable loss functions is disregarded, the optimization problems in (5) and (6) exhibit considerable similarity. To recognize this similarity, consider the following. Given a model $\theta \in \mathcal{M}$, let $J_{P_S, \theta} : \mathbb{R} \rightarrow \mathbb{R}$ be the function

$$J_{P_S, \theta}(t) = \log \left(\int \exp(t\ell(\theta, x, y)) dP_S(x, y) \right), \quad (7)$$

where the function ℓ is defined in (2). Let also the set $\mathcal{J}_{P_S, \theta}$ be the set

$$\mathcal{J}_{P_S, \theta} \triangleq \left\{ t \in (0, +\infty) : J_{P_S, \theta} \left(\frac{1}{t} \right) < +\infty \right\}. \quad (8)$$

Using these elements, the following theorem proves that the optimization problems in (5) and (6) exhibit the same unique solution.

Theorem 1: If the function ℓ in (2) is separable with respect to the measure P_S and $\beta \in \mathcal{J}_{P_S, \theta}$, with $\mathcal{J}_{P_S, \theta}$ in (8), then the probability measure $P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)} \in \Delta_{P_S}(\mathcal{X} \times \mathcal{Y})$ that satisfies for all $(x, y) \in \text{supp } P_S$,

$$\frac{dP_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}}{dP_S}(x, y) = \exp \left(\frac{1}{\beta} \ell(\theta, x, y) - J_{P_S, \theta} \left(\frac{1}{\beta} \right) \right), \quad (9)$$

where the function $J_{P_S, \theta}$ is defined in (7) and

$$D \left(P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)} || P_S \right) = \gamma, \quad (10)$$

is the unique solution to the optimization problems in (5) and (6).

Proof: The proof is presented in [4]. \blacksquare

Theorem 1 presents an alternative characterization of the WCDG probability measure introduced in [1], which unveils properties that hold significant importance in the subsequent sections.

II. CONTRIBUTIONS AND RELATED WORKS

The expected generalization gap (EGG) is a central performance metric for machine learning algorithms, see for instance, [5]–[14] and [15]. Analytical formulations for the

exact EGG have been established solely for the Gibbs algorithm when the reference measure is a probability measure in [5] and a σ -finite measure in [3]. For other machine learning algorithms, only upper bounds in terms of information measures have been obtained for the EGG, see for instance, [8], [16]–[26], and references therein. The first contribution of this paper is a closed-form expression for EGG for any statistical machine learning algorithm. This expression is formulated using information measures that incorporate the WCDG probability measure introduced in [1]. The derived expression establishes an upper bound on the EGG for any statistical machine learning algorithm. Such a bound is the sum of the mutual information and the lautum information between models and datasets, up to a constant for which a closed-form expression is provided. Notably, this bound is tight for a Gibbs algorithm, whose parameters are selected to satisfy a specific condition. The idea of studying the *worst-case* is also adopted in minimax distributed robust optimization, where the aim is to minimize the risk over a large ambiguity set, see for instance [27], [28] and [29]. This problem can also be cast as an instance of the maximum entropy principle [30]. Finally, this work has privileged the mathematical formulation whose solution exhibits properties that are useful for the analysis of the generalization. In particular, given a fixed model, the empirical risk is shown to be a sub-Gaussian random variable when datasets are sampled from the WCDG probability measure. This has led to the construction of new generalization guarantees coined (ϵ, δ) -robustness.

III. AN EXPRESSION OF THE EXPECTED GENERALIZATION GAP

Given a model $\theta \in \mathcal{M}$, the variation of the expected loss when the probability measure from which data points are sampled varies to another measure is referred to as the *sensitivity* of the expected loss. Such a sensitivity can be quantified using the functional $G : \mathcal{M} \times \Delta(\mathcal{X} \times \mathcal{Y}) \times \Delta(\mathcal{X} \times \mathcal{Y}) \rightarrow \mathbb{R}$,

$$G(\theta, P_1, P_2) = R_\theta(P_1) - R_\theta(P_2), \quad (11)$$

where the functional R_θ is defined in (4) and $\Delta(\mathcal{X} \times \mathcal{Y})$ is the set of all probability measures over $\mathcal{X} \times \mathcal{Y}$. In order to introduce the definition of the generalization gap, the notion of type induced by a dataset [31] is presented below.

Definition 2 (Type of a dataset): The type induced by the dataset z in (1) is a probability measure in $\Delta(\mathcal{X} \times \mathcal{Y})$, denoted by P_z , such that for all measurable subsets \mathcal{A} of $\mathcal{X} \times \mathcal{Y}$, $P_z(\mathcal{A}) = \frac{1}{n} \sum_{t=1}^n \mathbb{1}_{\{(x_t, y_t) \in \mathcal{A}\}}$.

In [1, Lemma 5.1], it is proved that the empirical risk $L(z, \theta)$ in (3) can be written as $L(z, \theta) = \int \ell(\theta, x, y) dP_z(x, y) = R_\theta(P_z)$, with the functional R_θ defined in (4). Given a model $\theta \in \mathcal{M}$ obtained from the training dataset $z \in (\mathcal{X} \times \mathcal{Y})^n$ and assuming that the training and test data points are generated according to an independent and identically distributed probability measure $P_Z \in \Delta(\mathcal{X} \times \mathcal{Y})$, the *generalization gap* induced by the model θ is

$$G(\theta, P_Z, P_z) = R_\theta(P_Z) - R_\theta(P_z). \quad (12)$$

The term $R_\theta(P_z) = L(z, \theta)$ is referred to as the training empirical risk, while the term $R_\theta(P_Z)$ is referred to as *true risk* or *population risk*.

Let $\Delta(\mathcal{M} | (\mathcal{X} \times \mathcal{Y})^n)$ be the set containing all probability measures over \mathcal{M} conditioned over an element of $(\mathcal{X} \times \mathcal{Y})^n$. Let also the function $\overline{G} : \Delta(\mathcal{M} | (\mathcal{X} \times \mathcal{Y})^n) \times \Delta((\mathcal{X} \times \mathcal{Y})^n) \rightarrow \mathbb{R}$ be such that

$$\overline{G}(P_{\Theta|Z}, P_Z) = \int \int G(\theta, P_Z, P_z) dP_{\Theta|Z=z}(\theta) dP_Z(z), \quad (13)$$

where P_z is the type induced by the dataset z ; and the functional G is defined in (12). Consider the random transformation $P_{\Theta|Z}$, such that given a training dataset $z \in (\mathcal{X} \times \mathcal{Y})^n$, the measure $P_{\Theta|Z=z} \in \Delta(\mathcal{M})$ is used to perform model selection. The conditional measure $P_{\Theta|Z}$ is referred to as a statistical learning algorithm. Under the assumption that the training and test datasets are both formed by independent and identically distributed datapoints sampled from the measure P_Z in (12), the expected generalization gap induced by $P_{\Theta|Z}$ is $\overline{G}(P_{\Theta|Z}, P_Z)$ in (13), where $P_Z \in \Delta((\mathcal{X} \times \mathcal{Y})^n)$ is a product distribution formed by P_Z . The main result of this section is presented by the following theorem.

Theorem 2: The expected generalization gap $\overline{G}(P_{\Theta|Z}, P_Z)$ in (13), with $P_Z \in \Delta((\mathcal{X} \times \mathcal{Y})^n)$, satisfies:

$$\begin{aligned} \frac{1}{\beta} \overline{G}(P_{\Theta|Z}, P_Z) &= \frac{1}{n} L(P_{\Theta|Z}; P_Z) \\ &+ \int D\left(P_{Z|\Theta=\theta} \| P_{Z|\Theta=\theta}^{(P_{Z|\Theta=\theta}, \beta)}\right) dP_\Theta(\theta) \\ &- \int D\left(P_Z \| P_{Z|\Theta=\theta}^{(P_{Z|\Theta=\theta}, \beta)}\right) dP_\Theta(\theta) \end{aligned} \quad (14)$$

where the term $L(P_{\Theta|Z}; P_Z)$ is the *lautum information* [32] given by

$$L(P_{\Theta|Z}; P_Z) \triangleq \int D(P_\Theta \| P_{\Theta|Z=\nu}) dP_Z(\nu); \quad (15)$$

and the measure $P_{Z|\Theta=\theta}^{(P_{Z|\Theta=\theta}, \beta)}$ is the WCDG probability measure in (9), with reference measure $P_{Z|\Theta=\theta}$ obtained from the joint measure $P_{\Theta|Z} P_Z$ and $\beta \in \left\{t > 0 : \forall \theta \in \mathcal{M}, \int \exp\left(\frac{1}{t} \ell(\theta, x, y)\right) dP_{Z|\Theta=\theta}(x, y) < +\infty\right\}$.

Proof: The proof is presented in [4]. \blacksquare

Theorem 2 appears to be the first exact characterization of the expected generalization gap induced by an arbitrary algorithm $P_{\Theta|Z}$ in terms of information measures. The following theorem unveils an upper bound closely linked to the celebrated Gibbs algorithm, which adds an interesting insight to the characterization of the generalization gap via the WCDG probability measure.

Theorem 3: The expected generalization gap $\overline{G}(P_{\Theta|Z}, P_Z)$ in (13), with $P_Z \in \Delta((\mathcal{X} \times \mathcal{Y})^n)$, satisfies:

$$\overline{G}(P_{\Theta|Z}, P_Z) \leq \frac{\beta}{n} (I(P_{\Theta|Z}; P_Z) + L(P_{\Theta|Z}; P_Z)) \quad (16)$$

where the term $I(P_{\Theta|Z}; P_Z)$ is the mutual information given by

$$I(P_{\Theta|Z}; P_Z) \triangleq \int D(P_{\Theta|Z=\nu} \| P_\Theta) dP_Z(\nu); \quad (17)$$

$L(P_{\Theta|Z}; P_Z)$ is in (15); and β satisfies for all $\theta \in \mathcal{M}$, $D\left(P_Z \| P_{Z|\Theta=\theta}^{(P_{Z|\Theta=\theta}, \beta)}\right) = 0$, with the measures P_Z , $P_{Z|\Theta=\theta}$, and $P_{Z|\Theta=\theta}^{(P_{Z|\Theta=\theta}, \beta)}$ defined in (14).

Proof: The proof is presented in [4]. \blacksquare

A typical example of a statistical learning algorithm is the Gibbs algorithm, which is parametrized by a positive real λ and by a σ -measure $Q \in \Delta(\mathcal{M})$, c.f. [3] and [33]. In the following, the focus is on the case in which Q is a probability measure. Under this assumption, the probability measure representing such an algorithm, which is denoted by $P_{\Theta|Z}^{(Q, \lambda)}$ with $\lambda > 0$, satisfies for all $\theta \in \text{supp } Q$ and for all $z \in (\mathcal{X} \times \mathcal{Y})^n$, $\frac{dP_{\Theta|Z=z}^{(Q, \lambda)}}{dQ}(\theta) = \exp\left(-K_{Q,z}\left(-\frac{1}{\lambda}\right) - \frac{1}{\lambda} L(z, \theta)\right)$, where the dataset z represents the training dataset; and the function $K_{Q,z} : \mathbb{R} \rightarrow \mathbb{R}$, satisfies $K_{Q,z}(t) = \log\left(\int \exp(t L(z, \nu)) dQ(\nu)\right)$. The EGG induced by $P_{\Theta|Z}^{(Q, \lambda)}$, under the assumption that datasets are sampled from a product distribution $P_Z \in \Delta((\mathcal{X} \times \mathcal{Y})^n)$ formed by the measure P_Z , denoted by $\overline{G}(P_{\Theta|Z}^{(Q, \lambda)}, P_Z)$, satisfies the following property [1], [5] and [3].

Lemma 1: The expected generalization gap $\overline{G}(P_{\Theta|Z}^{(Q, \lambda)}, P_Z)$ satisfies

$$\overline{G}(P_{\Theta|Z}^{(Q, \lambda)}, P_Z) = \lambda (I(P_{\Theta|Z}^{(Q, \lambda)}; P_Z) + L(P_{\Theta|Z}^{(Q, \lambda)}; P_Z)), \quad (18)$$

where the terms $I(P_{\Theta|Z}^{(Q, \lambda)}; P_Z)$ and $L(P_{\Theta|Z}^{(Q, \lambda)}; P_Z)$ are, respectively, a mutual information and a lautum information:

$$I(P_{\Theta|Z}^{(Q, \lambda)}; P_Z) \triangleq \int D(P_{\Theta|Z=\nu}^{(Q, \lambda)} \| P_\Theta^{(Q, \lambda)}) dP_Z(\nu); \quad \text{and} \quad (19)$$

$$L(P_{\Theta|Z}^{(Q, \lambda)}; P_Z) \triangleq \int D(P_\Theta^{(Q, \lambda)} \| P_{\Theta|Z=\nu}^{(Q, \lambda)}) dP_Z(\nu), \quad (20)$$

with $P_\Theta^{(Q, \lambda)}$ being a measure such that for all measurable sets $\mathcal{A} \in \mathcal{B}(\mathcal{M})$,

$$P_\Theta^{(Q, \lambda)}(\mathcal{A}) = \int P_{\Theta|Z=\nu}^{(Q, \lambda)}(\mathcal{A}) dP_Z(\nu). \quad (21)$$

Note that $I(P_{\Theta|Z}^{(Q, \lambda)}; P_Z)$ in (19) and $L(P_{\Theta|Z}^{(Q, \lambda)}; P_Z)$ in (20) also satisfy that

$$I(P_{Z|\Theta=\theta}^{(Q, \lambda)}; P_\Theta^{(Q, \lambda)}) = \int D(P_{Z|\Theta=\theta}^{(Q, \lambda)} \| P_Z) dP_\Theta^{(Q, \lambda)}(\nu); \quad (22)$$

$$L(P_{Z|\Theta=\theta}^{(Q, \lambda)}; P_\Theta^{(Q, \lambda)}) = \int D(P_Z \| P_{Z|\Theta=\theta}^{(Q, \lambda)}) dP_\Theta^{(Q, \lambda)}(\nu), \quad (23)$$

where $P_{Z|\Theta}^{(Q, \lambda)} \in \Delta((\mathcal{X} \times \mathcal{Y})^n | \mathcal{M})$ is the conditional probability measure that satisfies for all measurable subsets \mathcal{B} of $(\mathcal{X} \times \mathcal{Y})^n$,

$$P_Z(\mathcal{B}) = \int P_{Z|\Theta=\theta}^{(Q, \lambda)}(\mathcal{B}) dP_\Theta^{(Q, \lambda)}(\theta). \quad (24)$$

Moreover, for all $\theta \in \mathcal{M}$, the measure $P_{Z|\Theta=\theta}^{(Q,\lambda)} \in \Delta((\mathcal{X} \times \mathcal{Y})^n | \mathcal{M})$ is a product measure formed by a measure $P_{Z|\Theta=\theta}^{(Q,\lambda)} \in \Delta(\mathcal{X} \times \mathcal{Y} | \mathcal{M})$ that satisfies for all measurable sets \mathcal{A}_i in $\mathcal{X} \times \mathcal{Y}$, with $i \in \{1, 2, \dots, n\}$,

$$P_{Z|\Theta=\theta}^{(Q,\lambda)}(\mathcal{A}_1 \times \mathcal{A}_2 \times \dots \times \mathcal{A}_n) = \prod_{t=1}^n P_{Z|\Theta=\theta}^{(Q,\lambda)}(\mathcal{A}_t). \quad (25)$$

Using this notation, Theorem 3 and Lemma 1 lead to the following theorem.

Theorem 4: Assume that there exists a $\lambda > 0$ that satisfies

$$\lambda \in \left\{ t > 0 : \forall \theta \in \mathcal{M}, \int \exp\left(\frac{1}{nt} \ell(\theta, x, y)\right) dP_{Z|\Theta=\theta}^{(Q,\lambda)}(x, y) < +\infty \right\}, \quad (26)$$

where the function ℓ is defined in (2); and the measure $P_{Z|\Theta=\theta}^{(Q,\lambda)}$ is defined in (25). Let the measure $P_{\hat{Z}|\Theta=\theta}^{(P_{Z|\Theta=\theta}^{(Q,\lambda)}, n\lambda)}$ be the WCDG probability measure of the form in (9). If λ satisfies for all $\theta \in \mathcal{M}$,

$$D\left(P_Z \| P_{\hat{Z}|\Theta=\theta}^{(P_{Z|\Theta=\theta}^{(Q,\lambda)}, n\lambda)}\right) = 0, \quad (27)$$

then, for all $P_{\Theta|Z} \in \Delta(\mathcal{M} | (\mathcal{X} \times \mathcal{Y})^n)$, the expected generalization gaps $\overline{G}(P_{\Theta|Z}, P_Z)$ in (13) and $\overline{G}(P_{\Theta|Z}^{(Q,\lambda)}, P_Z)$ in (18) satisfy

$$\overline{G}(P_{\Theta|Z}, P_Z) \leq \overline{G}(P_{\Theta|Z}^{(Q,\lambda)}, P_Z). \quad (28)$$

Proof: The proof is presented in [4]. \blacksquare

Theorem 4 shows that, under the assumption that datasets are sampled from P_Z , the EGG of any algorithm $P_{\Theta|Z}$ is upper-bounded by the expected generalization gap induced by a particular Gibbs algorithm $P_{\Theta|Z}^{(Q,\lambda)}$. This Gibbs algorithm induces a posterior for a given model θ , denoted by $P_{Z|\Theta=\theta}^{(Q,\lambda)}$ in (25). When such a posterior is used as a reference measure to build the WCDG probability measure for such a model θ , with parameter $n\lambda$, it leads to the probability measure $P_{\hat{Z}|\Theta=\theta}^{(P_{Z|\Theta=\theta}^{(Q,\lambda)}, n\lambda)}$. If the parameters Q and λ are chosen to satisfy (27), then for all $\theta \in \mathcal{M}$, the WCDG probability measure is identical to the actual ground-truth data-generating probability measure P_Z . This is reminiscent of the principle of *indifference* over which the notion of equilibrium in zero-sum games with noisy observations is built [34].

The inequality in (28) reveals the central role of the Gibbs algorithm in statistical machine learning. Essentially, by studying the Gibbs algorithm $P_{\Theta|Z}^{(Q,\lambda)}$, for which its parameters Q and λ are chosen to satisfy (27), the expected generalization gap of any algorithm facing data sampled from $P_Z = P_{\hat{Z}|\Theta=\theta}^{(P_{Z|\Theta=\theta}^{(Q,\lambda)}, n\lambda)}$ can be upper bounded.

IV. SUB-GAUSSIANITY OF THE LOSS

The function $J_{P_S, \theta}$ in (7) exhibits several properties that enable the analysis of general machine learning algorithms. Let the m -th derivative of the function $J_{P_S, \theta}$ be denoted by $J_{P_S, \theta}^{(m)} : \mathbb{R} \rightarrow \mathbb{R}$, with $m \in \mathbb{N}$. Hence, for all $t \in \mathcal{J}_{P_S, \theta}$, with $\mathcal{J}_{P_S, \theta}$ in (8), $J_{P_S, \theta}^{(m)}(t) \triangleq \frac{d^m}{ds^m} J_{P_S, \theta}(s) \Big|_{s=t}$. It is proved in [4] that for a fixed $\theta \in \mathcal{M}$, it follows that if $(X, Y) \sim P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}$, with $P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}$ in (9), then the random variable

$$W_\theta \triangleq \ell(\theta, X, Y), \quad (29)$$

with the function ℓ in (2), possesses a mean, variance, and third moment given by $J_{P_S, \theta}^{(1)}\left(\frac{1}{\beta}\right)$, $J_{P_S, \theta}^{(2)}\left(\frac{1}{\beta}\right)$, and $J_{P_S, \theta}^{(3)}\left(\frac{1}{\beta}\right)$, respectively. In particular, the cumulant generating function of the random variable W_θ , which is denoted by $J_{P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}, \theta} : \mathbb{R} \rightarrow \mathbb{R}$, satisfies

$$J_{P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}, \theta}(t) = \log \left(\int \exp(t\ell(\theta, x, y)) dP_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}(x, y) \right), \quad (30)$$

where the function ℓ is defined in (2). It is also proved in [4] that $J_{P_S, \theta}^{(1)}\left(\frac{1}{\beta}\right) = R_\theta\left(P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}\right)$, with the functional R_θ defined in (4). The next theorem provides an upper bound on the cumulant generating function $J_{P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}, \theta}$ in (30).

Theorem 5: For all $t \in \{\alpha \in \mathbb{R} : J_{P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}, \theta}(\alpha) < +\infty\}$, the function $J_{P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}, \theta}$ in (30) satisfies the following inequality:

$$J_{P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}, \theta}(t) \leq t J_{P_S, \theta}^{(1)}\left(\frac{1}{\beta}\right) + \frac{1}{2} t^2 \zeta_{P_S, \theta}^2, \quad (31)$$

where $\zeta_{P_S, \theta}$ is finite, and satisfies

$$\zeta_{P_S, \theta} \triangleq \sup \left\{ \sqrt{J_{P_S, \theta}^{(2)}(\xi)} : \xi \in \left(-\infty, b - \frac{1}{\beta}\right) \right\}, \quad (32)$$

with $b \triangleq \sup \{\alpha \in \mathbb{R} : J_{P_S, \theta}(\alpha) < +\infty\}$, and the functions $J_{P_S, \theta}^{(1)}$ and $J_{P_S, \theta}^{(2)}$ are the first and the second derivatives of the function $J_{P_S, \theta}$ in (7).

Proof: The proof is presented in [4]. \blacksquare

The relevance of Theorem 5 lies on the fact that it implies that the random variable W_θ in (29) is a sub-Gaussian random variable with sub-Gaussian parameter $\zeta_{P_S, \theta}$ in (32). An interesting discussion on the impact of the random variable W_θ being sub-Gaussian is presented in [7, Theorem 1]. This observation leads to insightful guidelines for algorithm design.

V. (ϵ, δ) -ROBUSTNESS

An important observation about the dependence of the WCDG probability distribution $P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}$ in (9) on the parameter γ in (10) is that $P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}$ is the measure that induces the largest expected loss within the neighborhood of P_S , defined as

$$\mathcal{N}(P_S, \gamma) \triangleq \{P \in \Delta(\mathcal{X} \times \mathcal{Y}) : D(P \| P_S) \leq \gamma\}. \quad (33)$$

The following lemma provides insights on how $P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}$ changes according to the expansion or the contraction of the neighborhood.

Lemma 2: The relative entropy $D\left(P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}\|P_S\right)$ in (10) satisfies

$$\frac{d}{d\beta}D\left(P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}\|P_S\right) = -\frac{1}{\beta^3}J_{P_S, \theta}^{(2)}\left(\frac{1}{\beta}\right) \leq 0, \quad (34)$$

where $J_{P_S, \theta}^{(2)}\left(\frac{1}{\beta}\right)$ is the variance of W_θ in (29). Moreover, the inequality is strict if and only if the function ℓ in (2) is separable with respect to P_S .

Proof: The proof is presented in [4]. \blacksquare

From Lemma 2, it is shown that there exists a bijection between γ and β through (10). More interestingly, the variations of $D\left(P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}\|P_S\right)$ and the expected loss $J_{P_S, \theta}^{(1)}\left(\frac{1}{\beta}\right)$ with respect to β have the following connection: $\frac{d}{d\beta}D\left(P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}\|P_S\right) = \frac{1}{\beta}\frac{d}{d\beta}J_{P_S, \theta}^{(1)}\left(\frac{1}{\beta}\right)$, which implies that both $D\left(P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}\|P_S\right)$ and $J_{P_S, \theta}^{(1)}\left(\frac{1}{\beta}\right)$ decrease with β , while their rates of change might differ. For instance, for larger values of β with $\beta > 1$, the variation of $D\left(P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}\|P_S\right)$ is larger than $J_{P_S, \theta}^{(1)}\left(\frac{1}{\beta}\right)$. The following definition describes a performance metric for a given data-generating probability measure $P \in \Delta((\mathcal{X} \times \mathcal{Y})^n)$ and a model $\theta \in \mathcal{M}$ that leverages the observations above and provides guidelines for the choice of the values of β (or γ).

Definition 3 ((ϵ, δ)-Robustness): Given a pair of positive reals (ϵ, δ) with $\epsilon < 1$, a model $\theta \in \mathcal{M}$ is said to be (ϵ, δ) -robust to a probability measure $P \in \Delta((\mathcal{X} \times \mathcal{Y})^n)$, if

$$P(\{z \in (\text{supp } P_S)^n : L(z, \theta) \geq \delta\}) \leq \epsilon. \quad (35)$$

This notion of robustness enables the study of the performance guarantees that a model θ yields when faced with data generated by the WCDG probability measure for specific parameters β and P_S . An important issue that arises from this definition is the characterization of the largest value of γ (or smallest value of β) that achieves (ϵ, δ) -robustness, i.e. how much can the WCDG probability measure deviate from the reference P_S while the guarantee still holds. The following theorem establishes a condition for the (ϵ, δ) -robustness to hold for a given model θ .

Theorem 6: The probability measure $P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}$ in (9) satisfies that for all $\delta > J_{P_S, \theta}^{(1)}\left(\frac{1}{\beta}\right)$, with $J_{P_S, \theta}^{(1)}\left(\frac{1}{\beta}\right)$ being the mean of W_θ in (29),

$$\begin{aligned} & P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}\left(\{z \in (\mathcal{X} \times \mathcal{Y})^n : L(\theta, z) \geq \delta\}\right) \\ & \leq \exp\left(-nD\left(P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta^*)}\|P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}\right)\right), \end{aligned} \quad (36)$$

where $\beta^* \in (0, \beta) \cap \mathcal{J}_{P_S, \theta}$, with $\mathcal{J}_{P_S, \theta}$ in (8), satisfies

$$J_{P_S, \theta}^{(1)}\left(\frac{1}{\beta^*}\right) = \delta; \quad (37)$$

the function L is defined in (3); and the measure $P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)} \in \Delta((\mathcal{X} \times \mathcal{Y})^n)$ is a product measure formed by $P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}$.

Proof: The proof is presented in [4]. \blacksquare

Theorem 6 describes the (ϵ, δ) -robustness of a model θ in terms of another WCDG probability distribution $P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta^*)}$, where $\beta^* < \beta$. Interestingly, the WCDG probability measure $P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta^*)}$ induces an expected loss that is equal to δ and greater than the expected loss induced by $P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}$.

The following theorem provides a (ϵ, δ) -robust guarantee to any product probability measure formed by a measure in the neighborhood $\mathcal{N}(P_S, \gamma)$ in (33).

Theorem 7: For all $\theta \in \mathcal{M}$ and for all $P_Z \in \mathcal{N}(P_S, \gamma)$, with $\mathcal{N}(P_S, \gamma)$ in (33) and β in (10), it follows that for all $\delta > J_{P_S, \theta}^{(1)}\left(\frac{1}{\beta}\right)$, with $J_{P_S, \theta}^{(1)}\left(\frac{1}{\beta}\right)$ being the mean of W_θ in (29),

$$P_Z(\{z \in (\mathcal{X} \times \mathcal{Y})^n : L(\theta, z) \geq \delta\}) \leq \frac{1}{\delta}J_{P_S, \theta}^{(1)}\left(\frac{1}{\beta}\right), \quad (38)$$

where the function ℓ is defined in (2) and the product probability measure P_Z is formed by P_Z .

Proof: The proof is presented in [4]. \blacksquare

The relevance of Theorem 7 is that given a model θ , it establishes that for all $\delta > J_{P_S, \theta}^{(1)}\left(\frac{1}{\beta}\right)$ and $\epsilon_\delta \triangleq \frac{1}{\delta}J_{P_S, \theta}^{(1)}\left(\frac{1}{\beta}\right)$, such a model θ is $(\epsilon_\delta, \delta)$ -robust to all probability measures in $\mathcal{N}(P_S, \gamma)$, with $\mathcal{N}(P_S, \gamma)$ in (33). This observation raises the question on whether performing model selection based on the minimization of $J_{P_S, \theta}^{(1)}\left(\frac{1}{\beta}\right)$ is an alternative to classical approaches based on empirical risk minimization (ERM) as in [35]; or statistical ERM as in [3] and [36].

VI. CONCLUSIONS AND REMARKS

The first explicit expression in terms of information measures of the EGG for statistical learning algorithms has been presented. This expression has led to an upper bound on the EGG consisting of the sum of the mutual and lautum information between the models and the datasets, up to a constant factor. Such a bound is tight at least for one Gibbs algorithm, whose parameters are chosen to satisfy a particular condition. This observation reveals the central role of the Gibbs algorithm in statistical machine learning. Fundamentally, an exploration into the EGG of the Gibbs algorithm facilitates the derivation of overarching insights applicable to any statistical learning algorithm. Finally, the study of the properties of the WCDG probability measure leads to new machine learning algorithm performance metrics and insightful guidelines for algorithm design.

ACKNOWLEDGMENT

This work is funded in part by the ANR Project PARFAIT under grant ANR-21-CE25-0013 and the INRIA Exploratory Action IDEM.

REFERENCES

- [1] X. Zou, S. M. Perlaza, I. Esnaola, and E. Altman, “Generalization analysis of machine learning algorithms via the worst-case data-generating probability measure,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, Vancouver, Canada, Feb. 2024.
- [2] —, “The Worst-Case Data-Generating Probability Measure,” INRIA, Centre Inria d’Université Côte d’Azur, Sophia Antipolis, France, Tech. Rep. RR-9515, Aug. 2023.
- [3] S. M. Perlaza, G. Bisson, I. Esnaola, A. Jean-Marie, and S. Rini, “Empirical risk minimization with relative entropy regularization,” INRIA, Centre Inria d’Université Côte d’Azur, Sophia Antipolis, France, Tech. Rep. RR-9454, Feb. 2022.
- [4] X. Zou, S. M. Perlaza, I. Esnaola, E. Altman, and H. V. Poor, “An exact characterization of the generalization error for machine learning algorithms,” INRIA, Centre Inria d’Université Côte d’Azur, Sophia Antipolis, France, Tech. Rep. RR-9539, Jan. 2024.
- [5] 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.
- [6] G. Aminian, Y. Bu, G. W. Wornell, and M. R. Rodrigues, “Tighter expected generalization error bounds via convexity of information measures,” in *Proceedings of the IEEE International Symposium on Information Theory (ISIT)*, Aalto, Finland, Jun. 2022, pp. 2481–2486.
- [7] A. Xu and M. Raginsky, “Information-theoretic analysis of generalization capability of learning algorithms,” *Advances in Neural Information Processing Systems*, vol. 30, pp. 1–10, Dec. 2017.
- [8] Y. Chu and M. Raginsky, “A unified framework for information-theoretic generalization bounds,” arXiv preprint arXiv:2305.11042, May 2023.
- [9] A. Asadi, E. Abbe, and S. Verdú, “Chaining mutual information and tightening generalization bounds,” *Advances in Neural Information Processing Systems*, vol. 31, pp. 7245–7254, Dec. 2018.
- [10] A. R. Asadi and E. Abbe, “Chaining meets chain rule: Multilevel entropic regularization and training of neural networks,” *J. Mach. Learn. Res.*, vol. 21, pp. 139–1, Jun. 2020.
- [11] G. Aminian, L. Toni, and M. R. Rodrigues, “Jensen-Shannon information based characterization of the generalization error of learning algorithms,” in *Proceedings of the IEEE Information Theory Workshop (ITW)*, Kanazawa, Japan, Oct. 2021, pp. 1–5.
- [12] 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.
- [13] F. Daunas, I. Esnaola, S. M. Perlaza, and H. V. Poor, “Analysis of the relative entropy asymmetry in the regularization of empirical risk minimization,” in *Proceedings of the IEEE International Symposium on Information Theory (ISIT)*, Taipei, Taiwan, Jun. 2023.
- [14] —, “Empirical risk minimization with relative entropy regularization Type-II,” INRIA, Centre Inria d’Université Côte d’Azur, Sophia Antipolis, France, Tech. Rep. RR-9575, May. 2023.
- [15] 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.
- [16] Y. Bu, S. Zou, and V. V. Veeravalli, “Tightening mutual information-based bounds on generalization error,” *IEEE Journal on Selected Areas in Information Theory*, vol. 1, no. 1, pp. 121–130, Jan. 2020.
- [17] F. Hellström and G. Durisi, “Generalization bounds via information density and conditional information density,” *IEEE Journal on Selected Areas in Information Theory*, vol. 1, no. 3, pp. 824–839, Nov. 2020.
- [18] H. Hafez-Kolahi, Z. Golgooni, S. Kasaei, and M. Soleymani, “Conditioning and processing: Techniques to improve information-theoretic generalization bounds,” *Advances in Neural Information Processing Systems*, pp. 16 457–16 467, Dec. 2020.
- [19] I. Issa, A. R. Esposito, and M. Gastpar, “Strengthened information-theoretic bounds on the generalization error,” in *Proceedings of the IEEE International Symposium on Information Theory (ISIT)*, Paris, France, Jul. 2019, pp. 582–586.
- [20] A. R. Esposito, M. Gastpar, and I. Issa, “Robust generalization via α -mutual information,” arXiv preprint arXiv:2001.06399, Jan. 2020.
- [21] A. T. Lopez and V. Jog, “Generalization error bounds using Wasserstein distances,” in *Proceedings of the IEEE Information Theory Workshop (ITW)*, Guangzhou, China, Nov. 2018, pp. 1–5.
- [22] H. Wang, M. Diaz, J. C. S. Santos Filho, and F. P. Calmon, “An information-theoretic view of generalization via Wasserstein distance,” in *Proceedings of the IEEE International Symposium on Information Theory (ISIT)*, Paris, France, Jul. 2019, pp. 577–581.
- [23] I. Casado, L. A. Ortega, A. R. Masegosa, and A. Pérez, “Pac-bayes-chernoff bounds for unbounded losses,” arXiv preprint arXiv:2401.01148, Jan. 2024.
- [24] 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.
- [25] L. P. Barnes, A. Dytso, and H. V. Poor, “Improved information-theoretic generalization bounds for distributed, federated, and iterative learning,” *Entropy*, vol. 24, no. 9, 2022.
- [26] G. Aminian, Y. Bu, L. Toni, M. R. D. Rodrigues, and G. W. Wornell, “Information-theoretic characterizations of generalization error for the Gibbs algorithm,” *IEEE Transactions on Information Theory*, vol. 70, no. 1, pp. 632–655, 2024.
- [27] E. Delage and Y. Ye, “Distributionally robust optimization under moment uncertainty with application to data-driven problems,” *Operations research*, vol. 58, no. 3, pp. 595–612, 2010.
- [28] Z. Hu and L. J. Hong, “Kullback-Leibler divergence constrained distributionally robust optimization,” *Optimization Online*, vol. 1, no. 2, p. 9, 2013.
- [29] J. Lee and M. Raginsky, “Minimax Statistical Learning with Wasserstein distances,” *Advances in Neural Information Processing Systems*, vol. 31, pp. 2687–2696, 2018.
- [30] S. Mazuelas, Y. Shen, and A. Pérez, “Generalized maximum entropy for supervised classification,” *IEEE Transactions on Information Theory*, vol. 68, no. 4, pp. 2530–2550, Jan. 2022.
- [31] I. Csiszár, “The method of types,” *IEEE Transactions on Information Theory*, vol. 44, no. 6, pp. 2505–2523, Oct. 1998.
- [32] D. P. Palomar and S. Verdú, “Lautum information,” *IEEE Transactions on Information Theory*, vol. 54, no. 3, pp. 964–975, Mar. 2008.
- [33] 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.
- [34] K. Sun, S. M. Perlaza, and A. Jean-Marie, “ 2×2 Zero-sum games with commitments and Noisy Observations,” in *Proceedings of the IEEE International Symposium on Information Theory (ISIT)*, Taipei, Taiwan, Jun. 2023.
- [35] A. C. Wilson, R. Roelofs, M. Stern, N. Srebro, and B. Recht, “The marginal value of adaptive gradient methods in machine learning,” *Advances in neural information processing systems*, vol. 30, pp. 4151–4161, Dec. 2017.
- [36] F. Daunas, I. Esnaola, S. M. Perlaza, and H. V. Poor, “Empirical risk minimization with f -divergence regularization in statistical learning,” INRIA, Centre Inria d’Université Côte d’Azur, Sophia Antipolis, France, Tech. Rep. RR-9521, Oct. 2023.