



HAL
open science

An Exact Characterization of the Generalization Error of Machine Learning Algorithms

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

► **To cite this version:**

Samir M. Perlaza, Xinying Zou, Iñaki Esnaola, Eitan Altman, H. Vincent Poor. An Exact Characterization of the Generalization Error of Machine Learning Algorithms. 2024. hal-04426500v2

HAL Id: hal-04426500

<https://inria.hal.science/hal-04426500v2>

Preprint submitted on 14 May 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.



Distributed under a Creative Commons Attribution 4.0 International License

An Exact Characterization of the Generalization Error of Machine Learning Algorithms

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

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

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

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

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

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

Abstract—In this paper, exact expressions for the generalization error of general machine learning algorithms are presented in terms of information measures. These expressions can be broadly divided into two classes. The first one involves a worst-case data generating (WCDG) probability measure, while the second one involves a Gibbs algorithm. These expressions are formed by the sum of three terms. The first term is either a mutual or lautum information induced by the algorithm between the models and the datasets. The second and third terms are different for the two classes. In the first class, the second term compares via Kullback-Leibler (KL) divergence the posterior probability distribution induced by the algorithm on the datasets with a WCDG probability measure; and the third term compares via KL divergence the marginal of such posterior with the same WCDG probability measure. In the second class, the second term compares via KL divergence the distributions induced on the models by the algorithm and a Gibbs algorithm when both have been training with the same dataset; the third term compares via KL divergence the marginal of the distribution on the models induced by the algorithm with the same Gibbs algorithm. For Gibbs algorithms, the second and third terms jointly become either mutual or lautum information, recovering existing results.

I. INTRODUCTION

An exact expression for the generalization error is known only for the Gibbs algorithm [1]–[4]. For machine learning algorithms, other than the Gibbs algorithm, the generalization error is often characterized via upper bounds. Complete literature reviews of these bounds are presented in [5] and [6]. Nonetheless, it has been pointed out that such upper bounds are often loose, see for instance [7] and [8].

This paper introduces exact expressions for the generalization error that apply to any algorithm. Due to space constraints, the proofs of the main results are presented in [9].

Notation

Sets are denoted by caligraphic uppercase letters, except sets of probability measures. Given a set \mathcal{M} and a sigma-field \mathcal{F} on \mathcal{M} , the set of all probability measures that can be defined on the measurable space $(\mathcal{M}, \mathcal{F})$ is denoted by $\Delta(\mathcal{M}, \mathcal{F})$ or simply $\Delta(\mathcal{M})$, when the sigma-field is fixed in the analysis. If $\mathcal{M} \subseteq \mathbb{R}$, then the Borel sigma-field on \mathcal{M} is denoted by

$\mathcal{B}(\mathcal{M})$. The subset of measures in $\Delta(\mathcal{M})$ that are absolutely continuous with a probability measure Q is denoted by $\Delta_Q(\mathcal{M})$. Given a measure $P \in \Delta_Q(\mathcal{M})$, the Radon-Nikodym derivative of P with respect to Q is denoted by $\frac{dP}{dQ}$. Given a set \mathcal{N} , the set of all probability measures defined on $(\mathcal{M}, \mathcal{F})$ conditioned on an element of \mathcal{N} is denoted by $\Delta(\mathcal{M}|\mathcal{N})$. More specifically, given a measure $P_{M|N} \in \Delta(\mathcal{M}|\mathcal{N})$ and $n \in \mathcal{N}$, the measure $P_{M|N=n}$ is in $\Delta(\mathcal{M})$. Given a probability measure $P_N \in \Delta(\mathcal{N})$, the mutual information [10] and lautum information [11] induced by the measures $P_{M|N}$ and P_N are denoted respectively by

$$I(P_{M|N}; P_N) \triangleq \int D(P_{M|N=n} \| P_M) dP_N(n), \text{ and} \quad (1)$$

$$L(P_{M|N}; P_N) \triangleq \int D(P_M \| P_{M|N=n}) dP_N(n), \quad (2)$$

with $P_M \in \Delta(\mathcal{M})$ the marginal induced by $P_{M|N}$ and P_N ; and $D(\cdot \| \cdot)$ the Kullback-Leibler divergence.

II. PROBLEM FORMULATION

Let \mathcal{M} , \mathcal{X} and \mathcal{Y} , with $\mathcal{M} \subseteq \mathbb{R}^d$, 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*. Let the probability measure

$$P_Z \in \Delta(\mathcal{X} \times \mathcal{Y}) \quad (3)$$

be the one from which data points are independently sampled. While the probability measure P_Z is unknown in a practical setting, a finite set of datapoints is made available via a *training dataset*. A training dataset $\mathbf{z} \in (\mathcal{X} \times \mathcal{Y})^n$ is a tuple of n datapoints of the form:

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

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 $f(\theta, x)$. The function $f(\theta, \cdot): \mathcal{X} \rightarrow \mathcal{Y}$ might for instance represent a neural network whose weights are vectorized into θ . Let the function $\hat{\ell}: \mathcal{Y} \times \mathcal{Y} \rightarrow [0, +\infty]$, be the risk or loss function. Such a function is often of the form $\hat{\ell}(\mu, \nu) = |\mu - \nu|^p$, with $p \geq 1$; or $\hat{\ell}(\mu, \nu) = \mathbf{1}_{\{\mu = \nu\}}$. In general, given a labelled

pattern $(x, y) \in \mathcal{X} \times \mathcal{Y}$, the risk induced by a model $\theta \in \mathcal{M}$ is $\hat{\ell}(f(\theta, x), y)$. In the following, the risk function $\hat{\ell}$ is assumed to be a general nonnegative function and for all $y \in \mathcal{Y}$, $\hat{\ell}(y, y) = 0$. For the ease of notation, consider the function $\ell : \mathcal{X} \times \mathcal{Y} \times \mathcal{M} \rightarrow [0, +\infty]$ such that

$$\ell(x, y, \theta) = \hat{\ell}(f(\theta, x), y). \quad (5)$$

A. The Problem of Supervised Machine Learning

The problem of supervised machine learning is essentially devising an algorithm that takes as input a training dataset, as z in (4), and outputs a model θ , with certain probability. This notion is formally defined as follows.

Definition 1 (Algorithm): A conditional probability measure $P_{\Theta|Z} \in \Delta(\mathcal{M} | (\mathcal{X} \times \mathcal{Y})^n)$ is said to represent a supervised machine learning algorithm if it satisfies the following conditions:

- (a) For all $z \in (\mathcal{X} \times \mathcal{Y})^n$, the set-function $P_{\Theta|Z=z} : \mathcal{B}(\mathcal{M}) \rightarrow [0, 1]$ is a probability measure in $\Delta(\mathcal{M})$.
- (b) For all $\mathcal{A} \in \mathcal{B}(\mathcal{M})$, the function $h_{\mathcal{A}} : (\mathcal{X} \times \mathcal{Y})^n \rightarrow [0, 1]$ that satisfies $h_{\mathcal{A}}(z) = P_{\Theta|Z=z}(\mathcal{A})$ is measurable with respect to the measurable spaces $((\mathcal{X} \times \mathcal{Y})^n, \mathcal{Z}_n)$ and $([0, 1], \mathcal{B}([0, 1]))$, where \mathcal{Z}_n is a σ -algebra on $(\mathcal{X} \times \mathcal{Y})^n$.

Let $P_{\Theta|Z} \in \Delta(\mathcal{M} | (\mathcal{X} \times \mathcal{Y})^n)$ be an algorithm. Hence, the instance of such an algorithm trained upon the dataset z in (4) is denoted by $P_{\Theta|Z=z}$, which is simply a probability measure in $\Delta(\mathcal{M})$.

B. Empirical Risk and True Risk

The empirical risk is the average risk induced by a model among all datapoints of a given dataset. Let the function $L : (\mathcal{X} \times \mathcal{Y})^n \times \mathcal{M} \rightarrow [0, +\infty]$ be such that

$$L(z, \theta) = \frac{1}{n} \sum_{i=1}^n \ell(x_i, y_i, \theta), \quad (6)$$

where z is a dataset of the form in (4) and the function ℓ is defined in (5). Using this notation, the empirical risk is defined as follows.

Definition 2 (Empirical Risk): The *empirical risk* induced by a model $\theta \in \mathcal{M}$ with respect to the dataset z in (4) is $L(z, \theta)$ in (6).

In the following, the expectation of ℓ in (5) with respect to a probability measure P in $\Delta(\mathcal{X} \times \mathcal{Y})$, for a fixed model θ is denoted using the functional $R_{\theta} : \Delta(\mathcal{X} \times \mathcal{Y}) \rightarrow [0, +\infty]$, which satisfies

$$R_{\theta}(P) = \int \ell(x, y, \theta) dP(x, y). \quad (7)$$

Using this notation, the true risk is defined as follows.

Definition 3 (True Risk): The *true risk* induced by a model $\theta \in \mathcal{M}$, under the assumption that datasets are formed by i.i.d. datapoints sampled from P_Z in (3), is $R_{\theta}(P_Z)$, where the functional R_{θ} is defined (7).

Using the notion of *type induced by a dataset*, the empirical risk $L(z, \theta)$ in (6) can be written in terms of the function R_{θ} in (7). A type is defined as follows.

Definition 4 (Type of a dataset): The type induced by the dataset z in (4) is a probability measure in $\Delta(\mathcal{X} \times \mathcal{Y})$, denoted by P_z , such that for all measurable subsets \mathcal{A} of $\mathcal{X} \times \mathcal{Y}$,

$$P_z(\mathcal{A}) = \frac{1}{n} \sum_{t=1}^n \mathbb{1}_{\{(x_t, y_t) \in \mathcal{A}\}}. \quad (8)$$

The empirical risk $L(z, \theta)$ in (6) satisfies

$$L_z(\theta) = R_{\theta}(P_z), \quad (9)$$

where the functional R_{θ} is defined in (7); and the probability measure P_z is the type in (8) induced by the dataset z . See for instance, [4, Lemma 6.1].

C. Generalization Error

Generalization metrics, including the generalization error, are often defined in terms of differences of certain quantities. In the following, those differences are referred to as *gaps* and are quantified via the following functional: $G : \mathcal{M} \times \Delta(\mathcal{X} \times \mathcal{Y}) \times \Delta(\mathcal{X} \times \mathcal{Y}) \rightarrow \mathbb{R}$, such that

$$G(\theta, P_1, P_2) = R_{\theta}(P_1) - R_{\theta}(P_2). \quad (10)$$

Let the functional $\overline{G} : \Delta(\mathcal{M} | (\mathcal{X} \times \mathcal{Y})^n) \times \Delta(\mathcal{X} \times \mathcal{Y})^n \rightarrow \mathbb{R}$ be such that

$$\overline{G}(P_{\Theta|Z}, P_Z) = \int \int G(\theta, P_Z, P_Z) dP_{\Theta|Z=z}(\theta) dP_Z(z), \quad (11)$$

where P_z is the type induced by the dataset z ; the functional G is defined in (10); the probability measure $P_Z \in \Delta((\mathcal{X} \times \mathcal{Y})^n)$ is a product distribution formed by P_Z in (3); and $P_{\Theta|Z}$ is an algorithm (Definition 1).

The term $G(\theta, P_Z, P_Z)$ in (11) is the difference between the true risk $R_{\theta}(P_Z)$ and the empirical risk $R_{\theta}(P_z)$ induced by a model θ with respect to a dataset z . Hence, $\overline{G}(P_{\Theta|Z}, P_Z)$ is the expectation of such a difference with respect to a joint probability measure in $\Delta(\mathcal{M} \times (\mathcal{X} \times \mathcal{Y})^n)$ formed by the algorithm $P_{\Theta|Z}$ and the product measure P_Z . Using this notation, the *generalization error* of a machine learning algorithm $P_{\Theta|Z} \in \Delta(\mathcal{M} | (\mathcal{X} \times \mathcal{Y})^n)$ is defined as follows.

Definition 5 (Generalization Error): The generalization error induced by an algorithm $P_{\Theta|Z} \in \Delta(\mathcal{M} | (\mathcal{X} \times \mathcal{Y})^n)$, under the assumption that datasets are formed by i.i.d. datapoints sampled from P_Z in (3), is $\overline{G}(P_{\Theta|Z}, P_Z)$ in (11), with P_Z a product distribution formed by P_Z .

A more useful expression for the generalization error $\overline{G}(P_{\Theta|Z}, P_Z)$ in (11) can be obtained in terms of the conditional measure $P_{Z|\Theta} \in \Delta((\mathcal{X} \times \mathcal{Y})^n | \mathcal{M})$ and the measure $P_{\Theta} \in \Delta(\mathcal{M})$, which are particular for the algorithm $P_{\Theta|Z}$ and the product measure P_Z , as described hereunder. The measure P_{Θ} satisfies for all measurable subsets \mathcal{A} of \mathcal{M} ,

$$P_{\Theta}(\mathcal{A}) = \int P_{\Theta|Z=z}(\mathcal{A}) dP_Z(z); \quad (12)$$

and the conditional probability measure $P_{Z|\Theta}$ satisfies for all measurable subsets \mathcal{B} of $(\mathcal{X} \times \mathcal{Y})^n$,

$$P_Z(\mathcal{B}) = \int P_{Z|\Theta=\theta}(\mathcal{B}) dP_{\Theta}(\theta). \quad (13)$$

Note that under the assumption that datasets are formed by i.i.d. datapoints sampled from P_Z in (3), i.e., P_Z in (11) is a product measure formed by P_Z , it follows that for all $\theta \in \mathcal{M}$, the measure $P_{Z|\Theta=\theta}$ is a product measure formed by a probability measure $P_{Z|\Theta=\theta} \in \Delta(\mathcal{X} \times \mathcal{Y})$, which for all measurable sets of the form $\mathcal{A}_1 \times \mathcal{A}_2 \times \dots \times \mathcal{A}_n$ in $(\mathcal{X} \times \mathcal{Y})^n$, satisfies

$$P_{Z|\Theta=\theta}(\mathcal{A}_1 \times \mathcal{A}_2 \times \dots \times \mathcal{A}_n) = \prod_{t=1}^n P_{Z|\Theta=\theta}(\mathcal{A}_t), \quad (14)$$

and for all $t \in \{1, 2, \dots, n\}$,

$$P_Z(\mathcal{A}_t) = \int P_{Z|\Theta=\theta}(\mathcal{A}_t) dP_{\Theta}(\theta), \quad (15)$$

where the probability measure $P_{\Theta}(\theta)$ is defined in (12). Using this notation, the generalization error $\overline{G}(P_{\Theta|Z}, P_Z)$ in (11) can be re-written as follows:

$$\begin{aligned} & \overline{G}(P_{\Theta|Z}, P_Z) \\ &= \int \int G(\theta, P_Z, P_Z) dP_{Z|\Theta=\theta}(z) dP_{\Theta}(\theta), \end{aligned} \quad (16)$$

$$\begin{aligned} &= \int \int \left(\int \ell(\theta, x, y) dP_Z(x, y) \right) dP_{Z|\Theta=\theta}(z) dP_{\Theta}(\theta) \\ &\quad - \int \int \left(\int \ell(\theta, x, y) dP_Z(x, y) \right) dP_{Z|\Theta=\theta}(z) dP_{\Theta}(\theta) \end{aligned} \quad (17)$$

$$\begin{aligned} &= \int \int \ell(\theta, x, y) dP_Z(x, y) dP_{\Theta}(\theta), \\ &\quad - \int \int \left(\frac{1}{n} \sum_{t=1}^n \ell(\theta, x_t, y_t) \right) dP_{Z|\Theta=\theta}(z) dP_{\Theta}(\theta) \end{aligned} \quad (18)$$

$$\begin{aligned} &= \int \int \ell(\theta, x, y) dP_Z(x, y) dP_{\Theta}(\theta) \\ &\quad - \int \int \ell(\theta, x, y) dP_{Z|\Theta=\theta}(z) dP_{\Theta}(\theta) \end{aligned} \quad (19)$$

$$= \int G(\theta, P_Z, P_{Z|\Theta=\theta}) dP_{\Theta}(\theta), \quad (20)$$

where the equality in (16) follows from (12) and (13); the equality in (17) follows from (10); the equality in (18) follows from (8) and (9); the equality in (19) follows from the fact that the probability measure $P_{Z|\Theta=\theta}$ is a product measure as in (14); and the equality in (20) follows from (10).

The equality in (20) turns out to be instrumental in the remaining of this work.

III. THE GIBBS ALGORITHM

A typical example of a statistical learning algorithm is the Gibbs algorithm, which is parametrized by a positive real λ and by a probability measure $Q \in \Delta(\mathcal{M}, \mathcal{B}(\mathcal{M}))$ [3]. The probability measure representing such an algorithm is denoted by $P_{\Theta|Z}^{(Q, \lambda)} \in \Delta(\mathcal{M} | (\mathcal{X} \times \mathcal{Y})^n)$. When the Gibbs algorithm $P_{\Theta|Z}^{(Q, \lambda)}$ is trained with the training dataset z in (4), it satisfies for all $\theta \in \text{supp } Q$,

$$\frac{dP_{\Theta|Z=z}^{(Q, \lambda)}(\theta)}{dQ} = \exp\left(-K_{Q, z}\left(-\frac{1}{\lambda}\right) - \frac{1}{\lambda}L(z, \theta)\right), \quad (21)$$

where the function $K_{Q, z} : \mathbb{R} \rightarrow \mathbb{R}$ satisfies

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

and the function L is defined in (6).

The generalization error induced by the Gibbs algorithm with parameters Q and λ , 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 is

$$\overline{G}(P_{\Theta|Z}^{(Q, \lambda)}, P_Z), \quad (23)$$

where the functional \overline{G} is defined in (11). Such a generalization error satisfies the following property.

Lemma 1 (Generalization Gap of the Gibbs Algorithm): The generalization error $\overline{G}(P_{\Theta|Z}^{(Q, \lambda)}, P_Z)$ in (23) satisfies

$$\overline{G}(P_{\Theta|Z}^{(Q, \lambda)}, P_Z) = \lambda \left(I(P_{\Theta|Z}^{(Q, \lambda)}; P_Z) + L(P_{\Theta|Z}^{(Q, \lambda)}; P_Z) \right). \quad (24)$$

Proof: This result has been proved before in the case in which Q is a probability measure in [1]; and in the more general case in which Q is a σ -finite measure in [3]. A proof for the particular case in which Q is a probability measure and datasets are formed by independent and identically distributed datapoints is presented in [4] \blacksquare

IV. THE WORST-CASE DATA-GENERATING PROBABILITY MEASURE AND GENERALIZATION ERROR

The WCDG probability measure was introduced in [4] and further studied in [12]. The WCDG probability measure, which is parametrized by a real $\beta > 0$ and a probability measure P_S in $\Delta(\mathcal{X} \times \mathcal{Y})$, is specific for a given model θ and thus, denoted by $P_{Z|\Theta=\theta}^{(P_S, \beta)}$. Note that $P_{Z|\Theta=\theta}^{(P_S, \beta)}$ is a measure in $\Delta_{P_S}(\mathcal{X} \times \mathcal{Y})$. An explicit expression for $P_{Z|\Theta=\theta}^{(P_S, \beta)}$ requires defining the following function. Let $J_{P_S, \theta} : \mathbb{R} \rightarrow \mathbb{R}$ be the function

$$J_{P_S, \theta}(t) = \log\left(\int \exp(t\ell(\theta, x, y)) dP_S(x, y)\right), \quad (25)$$

where the function ℓ is defined in (5). The probability measure $P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}$ satisfies for all $(x, y) \in \text{supp } P_S$,

$$\frac{dP_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}}{dP_S}(x, y) = \exp\left(\frac{1}{\beta}\ell(\theta, x, y) - J_{P_S, \theta}\left(\frac{1}{\beta}\right)\right), \quad (26)$$

where the function $J_{P_S, \theta}$ is defined in (25). The existence of $P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}$ is ensured as long as $\beta \in \mathcal{J}_{P_S, \theta}$, where

$$\mathcal{J}_{P_S, \theta} \triangleq \left\{t \in (0, +\infty) : J_{P_S, \theta}\left(\frac{1}{t}\right) < +\infty\right\}. \quad (27)$$

The relevance of the WCDG probability measure $P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}$ stems from the fact that a closed-form expression for the term $G(\theta, P_1, P_2)$ in (10) can be obtained in terms of $P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}$. The following theorem shows such an expression.

Lemma 2 (Theorem 8 in [12]): The gap $G(\theta, P_1, P_2)$ in (10) satisfies

$$G(\theta, P_1, P_2) = \beta \left(D\left(P_2 \| P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}\right) - D\left(P_1 \| P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}\right) - D(P_2 \| P_S) + D(P_1 \| P_S) \right), \quad (28)$$

where the parameter β , the model θ , and the measures P_S and $P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}$ satisfy (26); and additionally, $\beta \in \mathcal{J}_{P_S, \theta}$, with the set $\mathcal{J}_{P_S, \theta}$ defined in (27).

In Lemma 2, the choice of β and P_S is flexible up to some constraints. Note that P_S can be any measure in $\Delta(\mathcal{X} \times \mathcal{Y})$, as long as both P_1 and P_2 are absolutely continuous with P_S ; and β is in the set $\mathcal{J}_{P_S, \theta}$ in (27). This flexibility is instrumental to obtain the main results. More specifically, the remaining of this work consists in studying the following expression, which follows immediately from the equality in (20) and Lemma 2,

$$\begin{aligned} \overline{\overline{G}}(P_{\Theta|Z}, P_Z) &= \beta \int \left(D\left(P_{Z|\Theta=\theta} \| P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}\right) - D\left(P_Z \| P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}\right) \right. \\ &\quad \left. - D\left(P_{Z|\Theta=\theta} \| P_S\right) + D\left(P_Z \| P_S\right) \right) dP_{\Theta}(\theta), \end{aligned} \quad (29)$$

where for all $\theta \in \text{supp } P_{\Theta}$, the probability measure $P_{Z|\Theta=\theta}$ is absolutely continuous with P_S ; and for all $\theta \in \text{supp } P_{\Theta}$, $\beta \in \mathcal{J}_{P_S, \theta}$, with the set $\mathcal{J}_{P_S, \theta}$ defined in (27).

V. MAIN RESULTS

The main results are closed-form expressions for the generalization error $\overline{\overline{G}}(P_{\Theta|Z}, P_Z)$ in (11), where $P_{\Theta|Z}$ is an arbitrary algorithm. Some expressions involve a WCDG probability measure, while others involve a Gibbs algorithm.

A. Expressions Involving a WCDG Measure

The following theorem shows one closed-form expression for the generalization error $\overline{\overline{G}}(P_{\Theta|Z}, P_Z)$ in (11) involving a mutual information.

Theorem 1: The generalization error $\overline{\overline{G}}(P_{\Theta|Z}, P_Z)$ in (11) satisfies:

$$\begin{aligned} \overline{\overline{G}}(P_{\Theta|Z}, P_Z) &= -\frac{\beta}{n} I(P_{\Theta|Z}; P_Z) \\ &\quad + \beta \int D\left(P_{Z|\Theta=\theta} \| P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}\right) dP_{\Theta}(\theta) \\ &\quad - \beta \int D\left(P_Z \| P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}\right) dP_{\Theta}(\theta) \end{aligned} \quad (30)$$

where the measure P_{Θ} is defined in (12); the conditional measure $P_{Z|\Theta}$ is defined in (15); the measure $P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}$ is the WCDG probability measure in (26), with P_S being such that for all $\theta \in \text{supp } P_{\Theta}$, the measure $P_{Z|\Theta=\theta}$ is absolutely continuous with P_S ; and β being such that for all $\theta \in \text{supp } P_{\Theta}$, $\beta \in \mathcal{J}_{P_S, \theta}$, with the set $\mathcal{J}_{P_S, \theta}$ defined in (27).

From Theorem 1, it follows that the generalization error can be strictly negative for a special class of algorithms. The following lemma formalizes this observation.

Lemma 3: Let $P_{\Theta|Z} \in \Delta(\mathcal{M} | (\mathcal{X} \times \mathcal{Y})^n)$ be an algorithm such that for all $\theta \in \mathcal{M}$, the corresponding probability measure $P_{Z|\Theta=\theta}$ in (15) is identical to the WCDG probability measure $P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}$ in (26), for some reference P_S and some real $\beta > 0$. Then, the expected generalization gap $\overline{\overline{G}}(P_{\Theta|Z}, P_Z)$ in (11) satisfies:

$$\overline{\overline{G}}(P_{\Theta|Z}, P_Z) = -\frac{\beta}{n} (I(P_{\Theta|Z}; P_Z) + L(P_{\Theta|Z}; P_Z)) \leq 0. \quad (31)$$

Note that if the algorithm $P_{\Theta|Z}$ in Theorem 1 is such that models are chosen independently of the training dataset, i.e., for all $(\mathbf{u}, \mathbf{v}) \in (\text{supp } P_Z)^n \times (\text{supp } P_Z)^n$, it holds that $D(P_{\Theta|Z=\mathbf{u}} \| P_{\Theta|Z=\mathbf{v}}) = 0$, then, $I(P_{\Theta|Z}; P_Z) = 0$ and moreover, for all $\theta \in \text{supp } P_{\Theta}$, it holds that

$$D\left(P_{Z|\Theta=\theta} \| P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}\right) = D\left(P_Z \| P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}\right), \quad (32)$$

which implies that $\overline{\overline{G}}(P_{\Theta|Z}, P_Z) = 0$. That is, algorithms that neglect the training datasets induce zero generalization error, which is a well-known result.

The following lemma introduces another expression for the generalization error $\overline{\overline{G}}(P_{\Theta|Z}, P_Z)$ in (11) involving a lautum information.

Lemma 4 (Theorem 10 in [12]): The generalization error $\overline{G}(P_{\Theta|Z}, P_Z)$ in (11) satisfies:

$$\begin{aligned} \overline{G}(P_{\Theta|Z}, P_Z) &= \frac{\beta}{n} L(P_{\Theta|Z}; P_Z) \\ &+ \beta \int D\left(P_{Z|\Theta=\theta} \| P_{\hat{Z}|\Theta=\theta}^{(P_{Z|\Theta=\theta}, \beta)}\right) dP_{\Theta}(\theta) \\ &- \beta \int D\left(P_Z \| P_{\hat{Z}|\Theta=\theta}^{(P_{Z|\Theta=\theta}, \beta)}\right) dP_{\Theta}(\theta) \end{aligned} \quad (33)$$

where the measure P_{Θ} is defined in (12); the conditional measure $P_{Z|\Theta}$ is defined in (15); the measure $P_{Z|\Theta=\theta}^{(P_{Z|\Theta=\theta}, \beta)}$ is the WCDG probability measure in (26), with reference measure $P_{Z|\Theta=\theta}$ and β such that for all $\theta \in \text{supp } P_{\Theta}$, $\beta \in \mathcal{J}_{P_{Z|\Theta=\theta}, \theta}$, with the set $\mathcal{J}_{P_{Z|\Theta=\theta}, \theta}$ defined as in (27).

In [12, Theorem 11], it has been suggested that there might exist a β that satisfies for all $\theta \in \text{supp } P_{\Theta}$,

$$D\left(P_Z \| P_{\hat{Z}|\Theta=\theta}^{(P_{Z|\Theta=\theta}, \beta)}\right) = 0, \quad (34)$$

which leads to the equality $\overline{G}(P_{\Theta|Z}, P_Z) = \frac{\beta}{n} (I(P_{\Theta|Z}; P_Z) + L(P_{\Theta|Z}; P_Z))$, from the fact that in such a case for all $\theta \in \text{supp } P_{\Theta}$, $D\left(P_{Z|\Theta=\theta} \| P_{\hat{Z}|\Theta=\theta}^{(P_{Z|\Theta=\theta}, \beta)}\right) = D(P_{Z|\Theta=\theta} \| P_Z)$, and thus,

$$\int D\left(P_{Z|\Theta=\theta} \| P_{\hat{Z}|\Theta=\theta}^{(P_{Z|\Theta=\theta}, \beta)}\right) dP_{\Theta}(\theta) = \frac{1}{n} I(P_{\Theta|Z}, P_Z). \quad (35)$$

This observation together with Lemma 2 might lead to conjecturing that the Gibbs algorithm $P_{\Theta|Z}^{(Q, \lambda)} \in \Delta(\mathcal{M} | (\mathcal{X} \times \mathcal{Y})^n)$ in (21), with $\beta = n\lambda$, satisfies the equality in (34). The following lemma shows that it is true under a very restrictive condition. Before presenting such a lemma, note that the measures P_{Θ} in (12); $P_{Z|\Theta=\theta}$ in (13); and $P_{Z|\Theta=\theta}$ in (14) that are specific to the Gibbs algorithm $P_{\Theta|Z}^{(Q, \lambda)}$ in (21) are denoted $P_{\Theta}^{(Q, \lambda)}$; $P_{Z|\Theta=\theta}^{(Q, \lambda)}$; and $P_{Z|\Theta=\theta}^{(Q, \lambda)}$, respectively.

Lemma 5: Let $P_{\Theta|Z}^{(Q, \lambda)} \in \Delta(\mathcal{M} | (\mathcal{X} \times \mathcal{Y})^n)$ be the Gibbs algorithm in (21). Then, for all $\theta \in \text{supp } P_{\Theta}^{(Q, \lambda)}$,

$$D\left(P_Z \| P_{\hat{Z}|\Theta=\theta}^{(P_{Z|\Theta=\theta}^{(Q, \lambda)}, n\lambda)}\right) = 0, \quad (36)$$

if and only if

$$P_Z\left(\left\{z \in (\mathcal{X} \times \mathcal{Y})^n : K_{Q, z}\left(-\frac{1}{\lambda}\right) = c\right\}\right) = 1, \quad (37)$$

for some real $c \leq 0$, where the function $K_{Q, z}$ is defined in (22).

The equality in (37) implies that the function $K_{Q, z}$ is invariant with respect to z , almost surely with P_Z . Note that this is always true if for all $\theta \in \mathcal{M}$, the function $\ell(\cdot, \cdot, \theta)$ in (5) is nonseparable with respect to P_Z [12, Definition 2]. That is, for all $\theta \in \mathcal{M}$,

$$P_Z(\{(x, y) \in \mathcal{X} \times \mathcal{Y} : \ell(x, y, \theta) = a(\theta)\}) = 1, \quad (38)$$

for some function $a : \mathcal{M} \rightarrow [0, +\infty)$. Nonetheless, this case is uninteresting in the context of machine learning, as the loss function ℓ in (5) is invariant with respect to the datapoints almost surely with P_Z in (3). Alternatively, if the function ℓ is separable, it is not clear whether the equality in (37) can be met.

The following lemma presents a property of the Gibbs algorithm that leads to the equality $\overline{G}(P_{\Theta|Z}^{(Q, \lambda)}, P_Z) = \lambda (I(P_{\Theta|Z}^{(Q, \lambda)}; P_Z) + L(P_{\Theta|Z}^{(Q, \lambda)}; P_Z))$, which is the claim of Lemma 1, without the need of verifying (36).

Lemma 6: Let $P_{\Theta|Z}^{(Q, \lambda)} \in \Delta(\mathcal{M} | (\mathcal{X} \times \mathcal{Y})^n)$ be the Gibbs algorithm in (21). Let $\beta \triangleq n\lambda$ be such that for all $\theta \in \text{supp } P_{\Theta}^{(Q, \lambda)}$, $\beta \in \mathcal{J}_{P_{Z|\Theta=\theta}^{(Q, \lambda)}, \theta}$, with the set $\mathcal{J}_{P_{Z|\Theta=\theta}^{(Q, \lambda)}, \theta}$ defined as in (27). Then,

$$\begin{aligned} \frac{1}{n} I(P_{\Theta|Z}^{(Q, \lambda)}; P_Z) &= \int D\left(P_{Z|\Theta=\theta}^{(Q, \lambda)} \| P_{\hat{Z}|\Theta=\theta}^{(P_{Z|\Theta=\theta}^{(Q, \lambda)}, \beta)}\right) dP_{\Theta}^{(Q, \lambda)}(\theta) \\ &- \int D\left(P_Z \| P_{\hat{Z}|\Theta=\theta}^{(P_{Z|\Theta=\theta}^{(Q, \lambda)}, \beta)}\right) dP_{\Theta}^{(Q, \lambda)}(\theta). \end{aligned} \quad (39)$$

Note that Lemma 6, together with Theorem 4, form an alternative proof for Lemma 1.

B. Expressions Involving a Gibbs Algorithm

The following lemma introduces another expression for the generalization error $\overline{G}(P_{\Theta|Z}, P_Z)$ in (11) involving a mutual information and a Gibbs algorithm.

Theorem 2: The generalization error $\overline{G}(P_{\Theta|Z}, P_Z)$ in (11) satisfies:

$$\begin{aligned} \overline{G}(P_{\Theta|Z}, P_Z) &= \lambda I(P_{\Theta|Z}; P_Z) \\ &+ \lambda \int \left(D(P_{\Theta} \| P_{\Theta|Z=z}^{(Q, \lambda)}) - D(P_{\Theta|Z=z} \| P_{\Theta|Z=z}^{(Q, \lambda)})\right) dP_Z(z), \end{aligned} \quad (40)$$

where the measure $P_{\Theta|Z}^{(Q, \lambda)} \in \Delta(\mathcal{M} | (\mathcal{X} \times \mathcal{Y})^n)$ is defined in (21), with $\lambda > 0$ and the measure Q being such that for all $z \in \text{supp } P_Z$, the measure $P_{\Theta|Z=z}$ is absolutely continuous with Q .

Theorem 2 allows the calculation of the generalization error of an arbitrary algorithm $P_{\Theta|Z}$ by comparing it with a Gibbs algorithm $P_{\Theta|Z}^{(Q, \lambda)} \in \Delta(\mathcal{M} | (\mathcal{X} \times \mathcal{Y})^n)$, whose reference measure Q is chosen such that all the instances of the algorithm $P_{\Theta|Z=z}$, obtained with all possible training datasets z , are absolutely continuous with Q .

ACKNOWLEDGMENT

This work is funded in part by the ANR Project PARFAIT under grant ANR-21-CE25-0013; in part by the H2020 RISE Project TESTBED2 under EU Grant 872172; in part by the U.S. National Science Foundation under Grant ECCS-2335876; in part by the C3.ai Digital Transformation Institute; and in part by Princeton Language + Intelligence.

REFERENCES

- [1] 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.
- [2] 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.
- [3] S. M. Perlaza, G. Bisson, I. Esnaola, A. Jean-Marie, and S. Rini, "Empirical risk minimization with relative entropy regularization," *IEEE Transactions on Information Theory*, 2024.
- [4] 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.
- [5] 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.
- [6] B. Rodríguez Gálvez, "An information-theoretic approach to generalization theory," Ph.D. dissertation, KTH Royal Institute of Technology, Apr. 2024.
- [7] Y. Jiang, B. Neyshabur, H. Mobahi, D. Krishnan, and S. Bengio, "Fantastic generalization measures and where to find them," in *Proceedings of the International Conference on Learning Representations*, Addis Ababa, Ethiopia, Apr. 2020.
- [8] M. Gastpar, I. Nachum, J. Shafer, and T. Weinberger, "Fantastic generalization measures are nowhere to be found," in *International Conference on Learning Representations (ICLR)*, Vienna, Austria, May 2024.
- [9] 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.
- [10] C. E. Shannon, "A mathematical theory of communication," *The Bell system technical journal*, vol. 27, no. 3, pp. 379–423, 1948.
- [11] D. P. Palomar and S. Verdú, "Lautum information," *IEEE Transactions on Information Theory*, vol. 54, no. 3, pp. 964–975, Mar. 2008.
- [12] X. Zou, S. M. Perlaza, I. Esnaola, E. Altman, and H. V. Poor, "The worst-case data-generating probability measure in statistical learning," *IEEE Journal on Selected Areas in Information Theory*, vol. 5, pp. 175–189, 2024.