

Generalization Error of Machine Learning Algorithms

The Method of Gaps

Samir M. Perlaza

samir.perlaza@inria.fr

INRIA, Centre Inria d'Université Côte d'Azur

Information Theory and Tapas Workshop

Universidad Carlos III de Madrid

April 2, 2025. Madrid Spain



Xinying Zou
INRIA
France



Francisco Daunas
*The University of Sheffield,
UK*



Yaiza Bermudez
INRIA
France



Iñaki Esnaola
*The University of Sheffield,
UK*



H. Vincent Poor
*Princeton University,
United States*



Gaetan Bisson
*University of French Polynesia,
French Polynesia*

Table of Contents

Supervised Statistical Learning and **Generalization Error**

Empirical Risk **Optimization** with Relative Entropy Regularization

The Method of Gaps

Explicit Expressions for the Generalization Error

Concluding Remarks

Table of Contents

Supervised Statistical Learning and Generalization Error

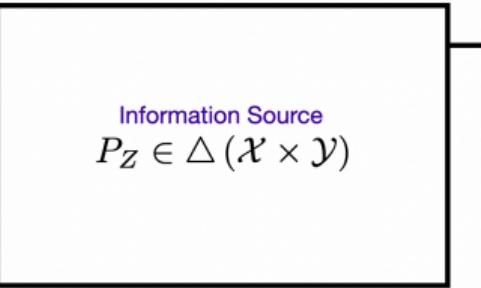
Empirical Risk Optimization with Relative Entropy Regularization

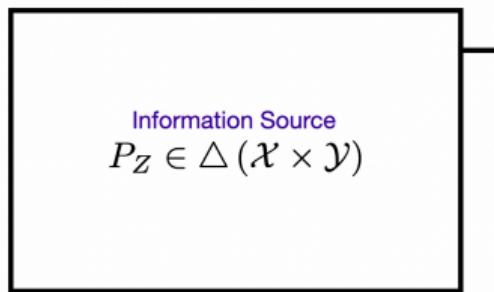
The Method of Gaps

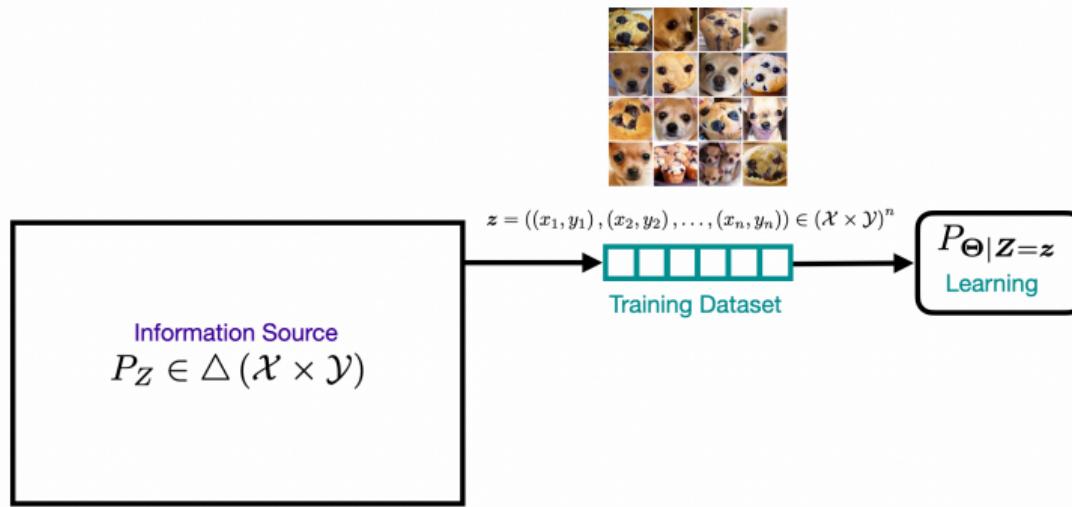
Explicit Expressions for the Generalization Error

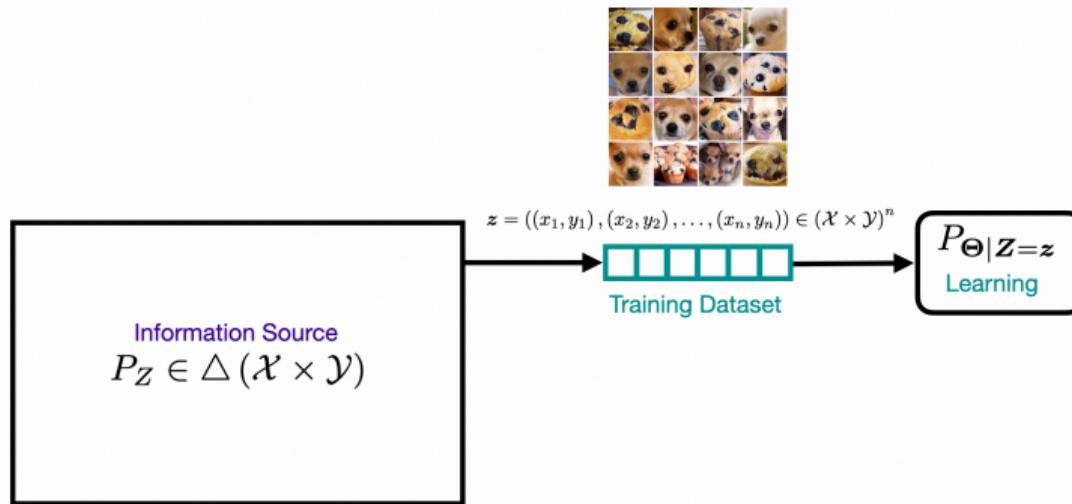
Concluding Remarks

Information Source
 $P_Z \in \Delta(\mathcal{X} \times \mathcal{Y})$



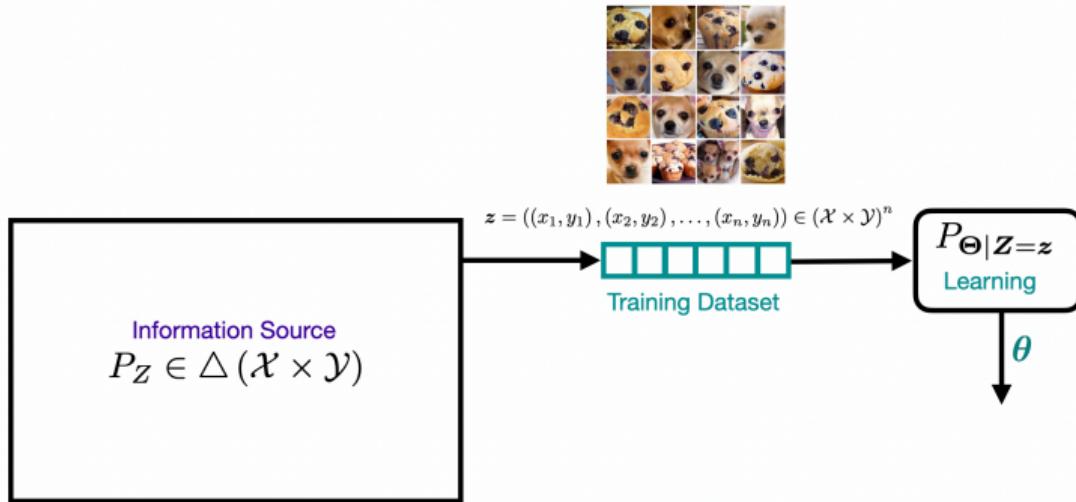


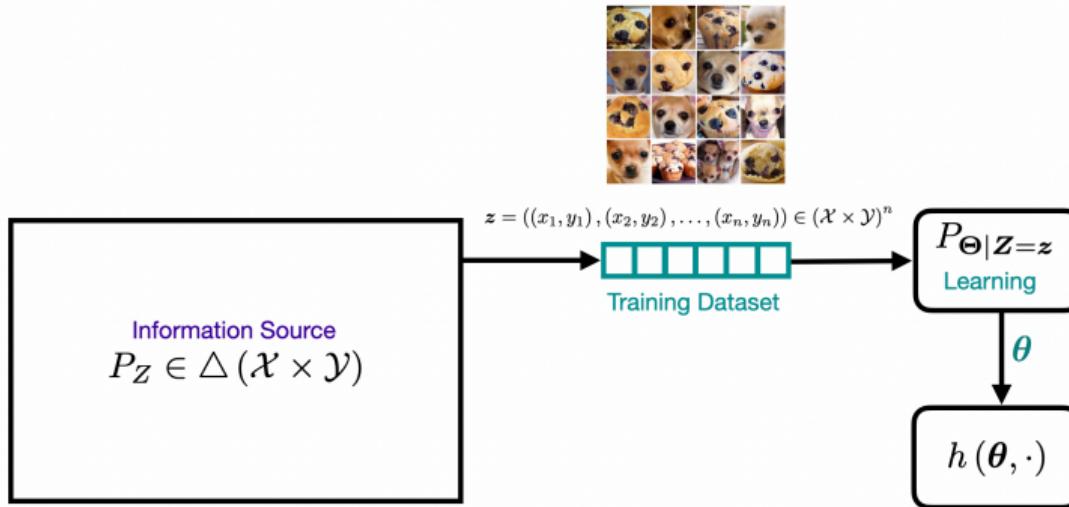


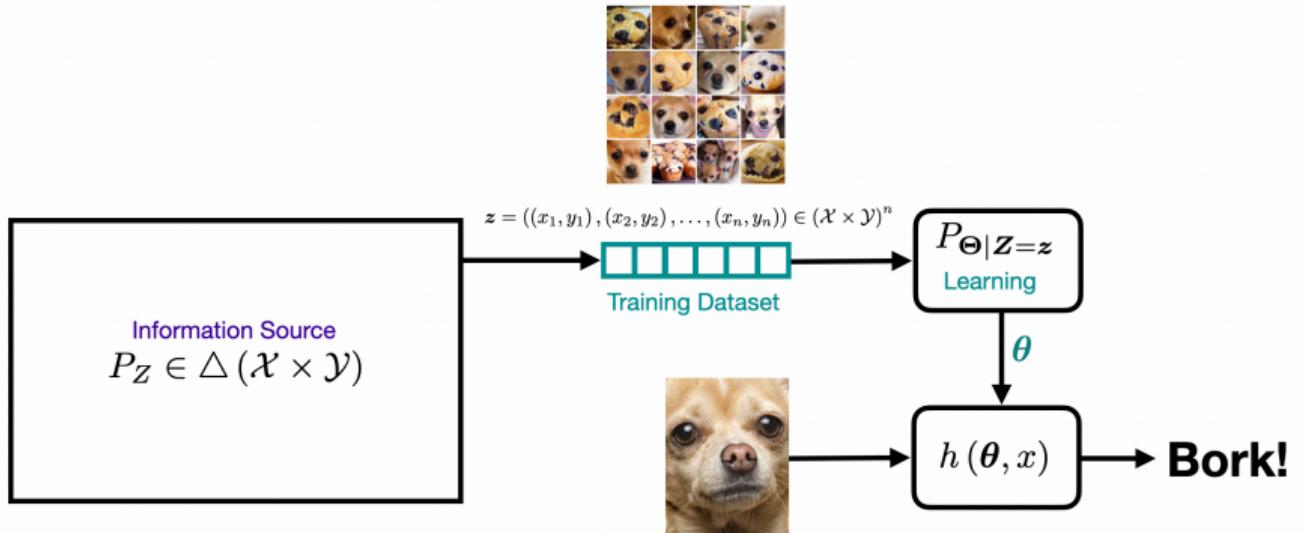


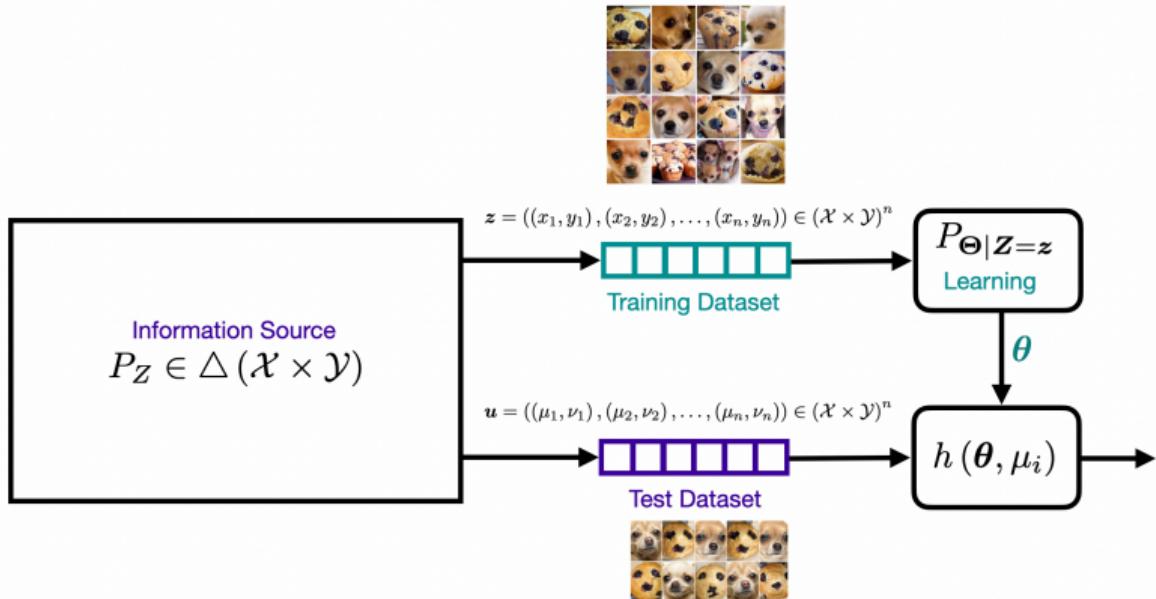
Algorithm

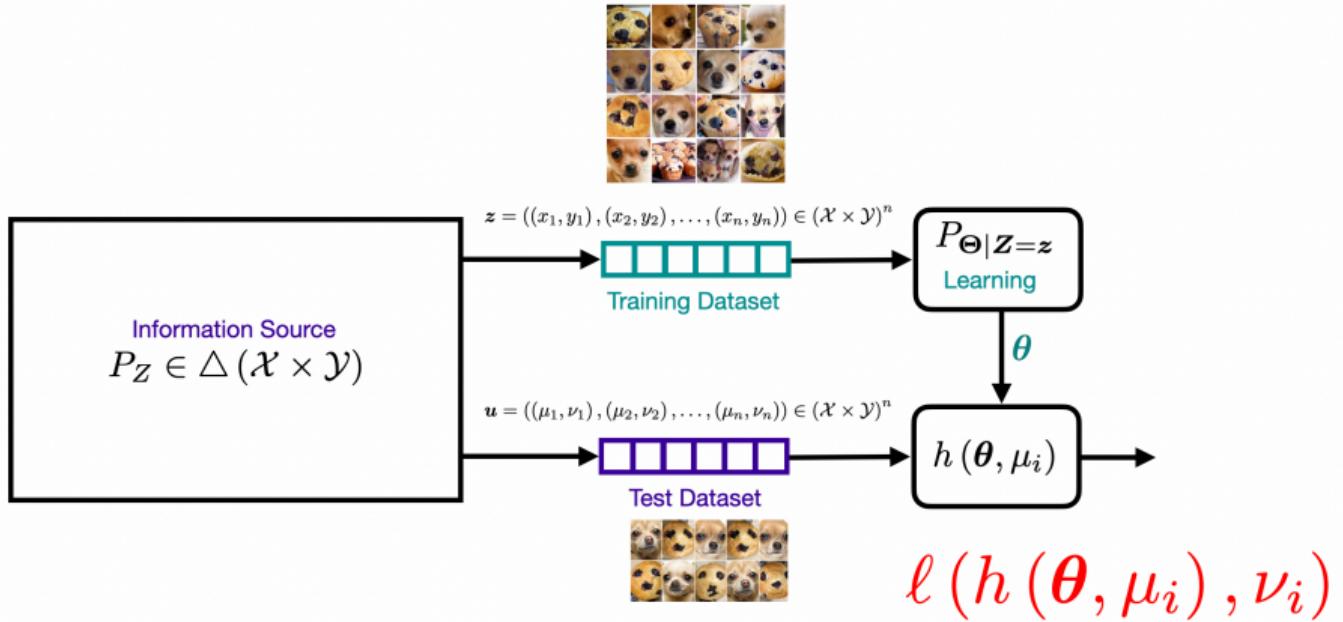
A conditional probability measure $P_{\Theta|Z} \in \Delta(\mathcal{M}|(\mathcal{X} \times \mathcal{Y})^n)$ represents a supervised machine learning algorithm.

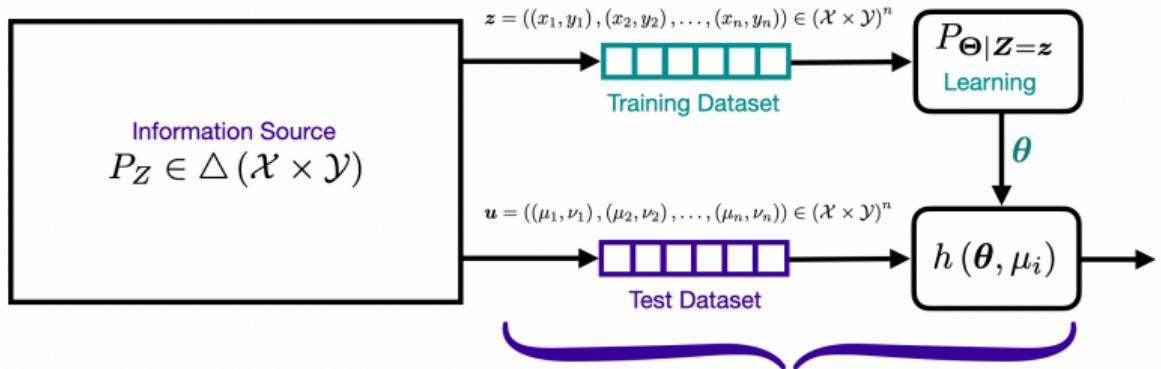






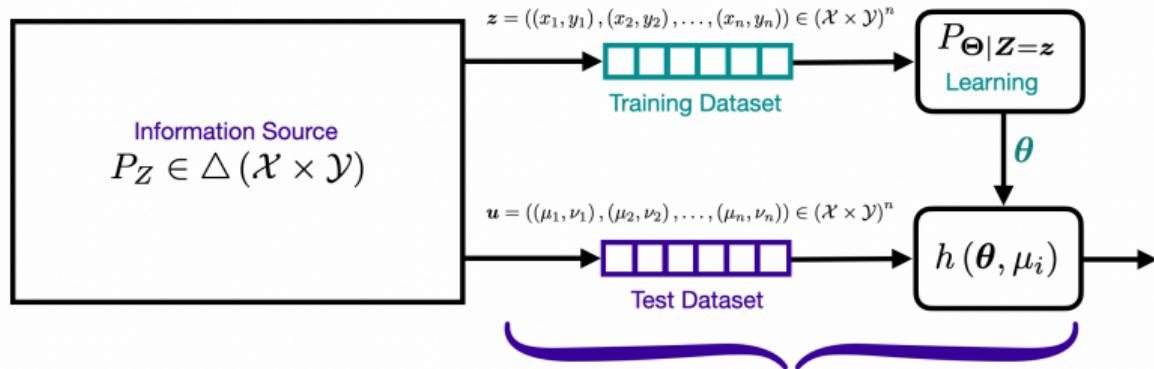






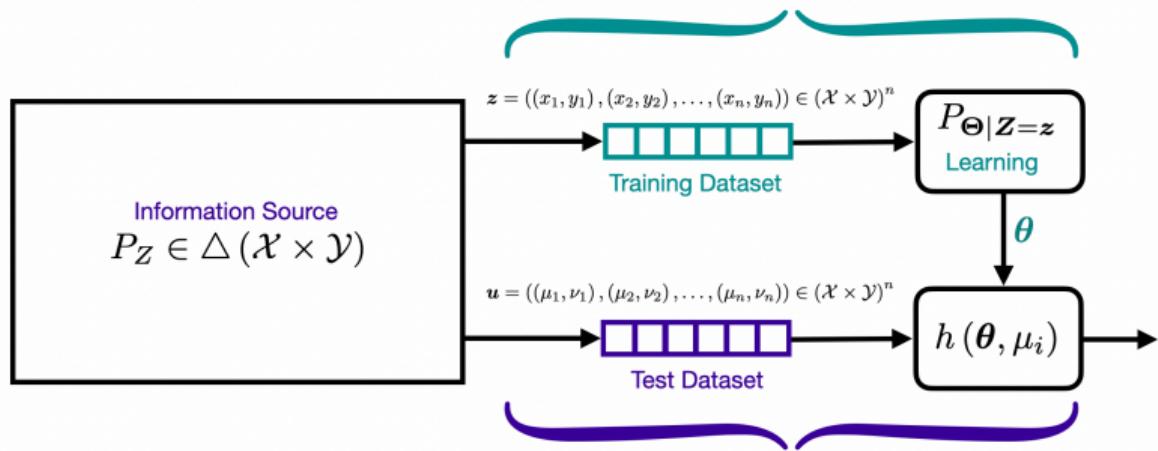
$$L(\mathbf{u}, \boldsymbol{\theta}) = \frac{1}{n} \sum_{t=1}^n \ell(h(\boldsymbol{\theta}, \mu_t), \nu_t)$$

$$L(z, \theta) = \frac{1}{n} \sum_{t=1}^n \ell(h(\theta, x_t), y_t)$$



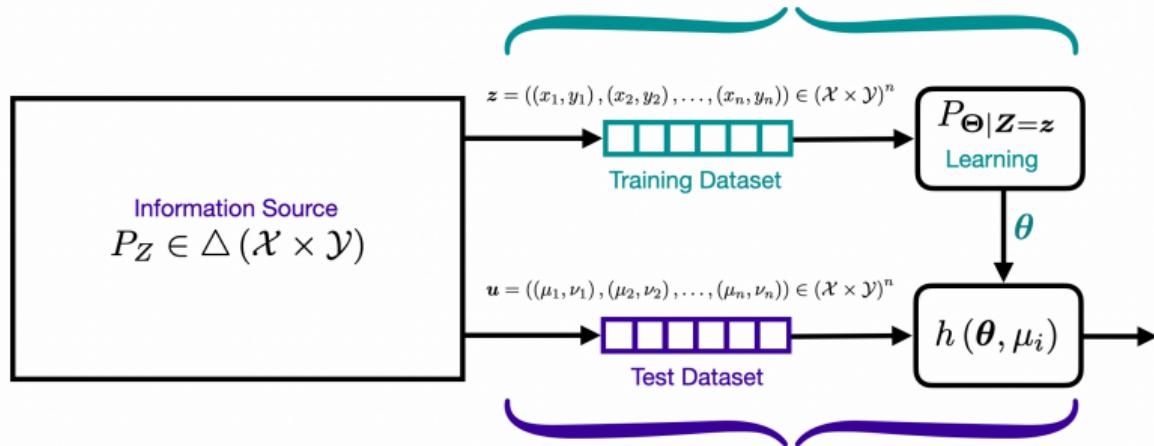
$$L(u, \theta) = \frac{1}{n} \sum_{t=1}^n \ell(h(\theta, \mu_t), \nu_t)$$

$$R_z(P_{\Theta|Z=z}) = \int L(z, \theta) dP_{\Theta|Z=z}(\theta)$$



$$R_u(P_{\Theta|Z=z}) = \int L(u, \theta) dP_{\Theta|Z=z}(\theta)$$

$$R_z(P_{\Theta|Z=z}) = \int L(z, \theta) dP_{\Theta|Z=z}(\theta)$$



$$R_u(P_{\Theta|Z=z}) = \int L(u, \theta) dP_{\Theta|Z=z}(\theta)$$

Training (Expected) Risk and Test (Expected) Risk

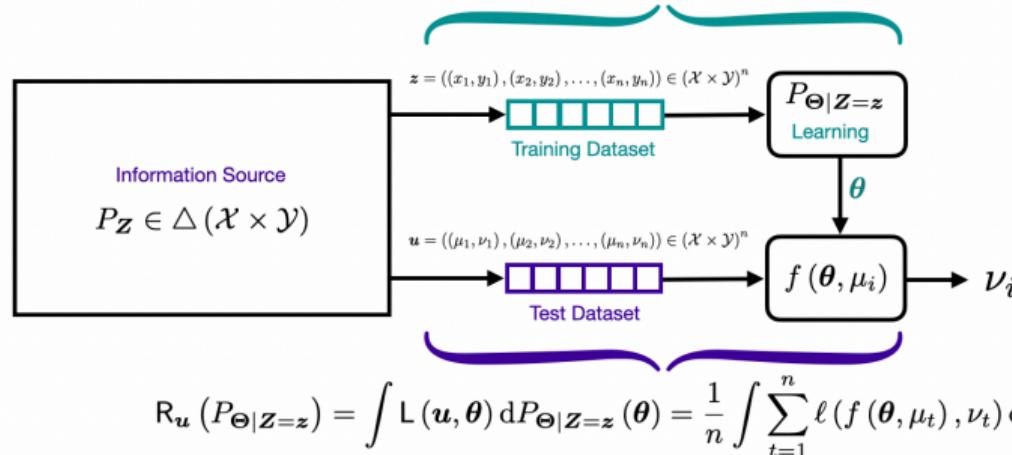
$$\underbrace{R_u(P_{\Theta|Z=z})}_{\text{Test Expected Risk}} - \underbrace{R_z(P_{\Theta|Z=z})}_{\text{Training Expected Risk}}$$

Assumption:

Training datasets and **test datasets** are independent and identically distributed:

- ▶ z is drawn from $P_z \in \Delta((\mathcal{X} \times \mathcal{Y})^n)$; and
- ▶ u is drawn from P_z .

$$R_z(P_{\Theta|Z=z}) = \int L(z, \theta) dP_{\Theta|Z=z}(\theta) = \frac{1}{n} \int \sum_{t=1}^n \ell(f(\theta, x_t), y_t) dP_{\Theta|Z=z}(\theta)$$



Generalization Error

The generalization error of the algorithm $P_{\Theta|Z}$ is

$$\overline{\overline{G}}(P_{\Theta|Z}, P_Z) \triangleq \int \int (R_u(P_{\Theta|Z=z}) - R_z(P_{\Theta|Z=z})) dP_Z(u) dP_Z(z).$$

Table of Contents

Supervised Statistical Learning and **Generalization Error**

Empirical Risk Optimization with Relative Entropy Regularization

The Method of Gaps

Explicit Expressions for the Generalization Error

Concluding Remarks

Empirical Risk **Minimization** with Relative Entropy Regularization

The Gibbs Algorithm

- ▶ Given a **fixed dataset** $z \in (\mathcal{X} \times \mathcal{Y})^n$; and
- ▶ given a **reference measure** $Q \in \Delta(\mathcal{M})$ and a **real** $\lambda > 0$

Problem 1: ERM with Relative Entropy Regularization

$$\min_{P \in \Delta_Q(\mathcal{M})} \int L(z, \theta) dP(\theta) + \lambda D(P\|Q),$$

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

[Perlaza-2024a] Samir M. Perlaza, Gaetan Bisson, Iñaki Esnaola, Alain Jean-Marie, and Stefano Rini. "Empirical Risk Minimization with Relative Entropy Regularization". IEEE Transactions on Information Theory, vol. 70, no. 7, pp. 5122 – 5161, July, 2024.

Empirical Risk **Minimization** with Relative Entropy Regularization

The Gibbs Algorithm

- ▶ Given a **fixed dataset** $z \in (\mathcal{X} \times \mathcal{Y})^n$; and
- ▶ given a **reference measure** $Q \in \Delta(\mathcal{M})$ and a **real** $\lambda > 0$

Problem 1: ERM with Relative Entropy Regularization

$$\min_{P \in \Delta_Q(\mathcal{M})} \int L(z, \theta) dP(\theta) + \lambda D(P\|Q),$$

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

Problem 1a: ERM within a Neighborhood

$$\begin{aligned} \min_{P \in \Delta_Q(\mathcal{M})} & \int L(z, \theta) dP(\theta) \\ \text{s.t. } & D(P\|Q) \leq \gamma. \end{aligned}$$

[Perlaza-2024a] Samir M. Perlaza, Gaetan Bisson, Iñaki Esnaola, Alain Jean-Marie, and Stefano Rini. "Empirical Risk Minimization with Relative Entropy Regularization". IEEE Transactions on Information Theory, vol. 70, no. 7, pp. 5122 – 5161, July, 2024.

Empirical Risk **Minimization** with Relative Entropy Regularization

The Gibbs Algorithm

- ▶ Given a **fixed dataset** $z \in (\mathcal{X} \times \mathcal{Y})^n$; and
- ▶ given a **reference measure** $Q \in \Delta(\mathcal{M})$ and a **real** $\lambda > 0$

Problem 1: ERM with Relative Entropy Regularization

$$\min_{P \in \Delta_Q(\mathcal{M})} \int L(z, \theta) dP(\theta) + \lambda D(P \| Q),$$

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

Problem 1b: Relative Entropy Minimization with First Order Constraint

$$\begin{aligned} & \min_{P \in \Delta_Q(\mathcal{M})} D(P \| Q) \\ \text{s.t. } & \int L(z, \theta) dP(\theta) \leq \xi. \end{aligned}$$

[Perlaza-2024a] Samir M. Perlaza, Gaetan Bisson, Iñaki Esnaola, Alain Jean-Marie, and Stefano Rini. "Empirical Risk Minimization with Relative Entropy Regularization". IEEE Transactions on Information Theory, vol. 70, no. 7, pp. 5122 – 5161, July, 2024.

Empirical Risk **Maximization** with Relative Entropy Regularization

Worst-Case Data-Generating Probability Measure

- ▶ Given a **fixed model** $\theta \in \mathcal{M}$; and
- ▶ Given a **reference measure** $P_S \in \Delta(\mathcal{X} \times \mathcal{Y})$ and a **real** $\beta > 0$

Problem 2: Loss Maximization with Relative Entropy Regularization

$$\max_{P \in \Delta_{P_S}(\mathcal{X} \times \mathcal{Y})} \int \ell(h(\theta, x), y) dP(x, y) - \beta D(P \| P_S),$$

with $\Delta_{P_S}(\mathcal{X} \times \mathcal{Y}) \triangleq \{P \in \Delta(\mathcal{X} \times \mathcal{Y}) : P \ll P_S\}$.

Empirical Risk **Maximization** with Relative Entropy Regularization

Worst-Case Data-Generating Probability Measure

- ▶ Given a **fixed model** $\theta \in \mathcal{M}$; and
- ▶ Given a **reference measure** $P_S \in \Delta(\mathcal{X} \times \mathcal{Y})$ and a **real** $\beta > 0$

Problem 2: Loss Maximization with Relative Entropy Regularization

$$\max_{P \in \Delta_{P_S}(\mathcal{X} \times \mathcal{Y})} \int \ell(h(\theta, x), y) dP(x, y) - \beta D(P \| P_S),$$

with $\Delta_{P_S}(\mathcal{X} \times \mathcal{Y}) \triangleq \{P \in \Delta(\mathcal{X} \times \mathcal{Y}) : P \ll P_S\}$.

Problem 2: Loss Maximization within a Neighbourhood

$$\begin{aligned} & \max_{P \in \Delta_{P_S}(\mathcal{X} \times \mathcal{Y})} \int \ell(h(\theta, x), y) dP(x, y) \\ \text{s.t. } & D(P \| P_S) \leq \gamma. \end{aligned}$$

[Zou-2024] Xinying Zou, Samir M. Perlaza, Iñaki Esnaola, Eitan Altman, and H. Vincent Poor. "The Worst-Case Data-Generating Probability Measure in Statistical Learning". IEEE Journal on Selected Areas in Information Theory, vol. 5, pp. 175–189, Apr., 2024.

Problem 1:
Empirical Risk **Minimization** with Relative Entropy Regularization

Empirical Risk **Minimization** with Relative Entropy Regularization

Problem 1: ERM with Relative Entropy Regularization

$$\min_{P \in \Delta_Q(\mathcal{M})} \int L(z, \theta) dP(\theta) + \lambda D(P\|Q),$$

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

Notation:

$$K_{Q,z}(t) = \log \left(\int \exp(t L(z, \theta)) dQ(\theta) \right) \text{ and } \mathcal{K}_{Q,z} \triangleq \left\{ s \in (0, +\infty) : K_{Q,z}\left(-\frac{1}{s}\right) < +\infty \right\}$$

Theorem (Theorem 3 in [Perlaza-2024a])

If $\lambda \in \mathcal{K}_{Q,z}$, the solution to Problem 1 is unique, denoted by $P_{\Theta|z=z}^{(Q,\lambda)}$, and satisfies for all $\theta \in \text{supp } Q$,

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

[Perlaza-2024a] Samir M. Perlaza, Gaetan Bisson, Iñaki Esnaola, Alain Jean-Marie, and Stefano Rini. "Empirical Risk Minimization with Relative Entropy Regularization". IEEE Transactions on Information Theory, vol. 70, no. 7, pp. 5122 – 5161, July, 2024.

Empirical Risk **Minimization** with Relative Entropy Regularization

Problem 1: ERM with Relative Entropy Regularization

$$\min_{P \in \Delta_Q(\mathcal{M})} \int L(z, \theta) dP(\theta) + \lambda D(P\|Q),$$

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

Notation:

$$K_{Q,z}(t) = \log \left(\int \exp(t L(z, \theta)) dQ(\theta) \right) \text{ and } \mathcal{K}_{Q,z} \triangleq \left\{ s \in (0, +\infty) : K_{Q,z}\left(-\frac{1}{s}\right) < +\infty \right\}$$

Theorem (Equation (28) in [Perlaza-2024a])

If $\lambda \in \mathcal{K}_{Q,z}$, the solution to Problem 1 is a unique, denoted by $P_{\Theta|z=z}^{(Q,\lambda)}$, and satisfies for all $\theta \in \sup Q$,

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

[Perlaza-2024a] Samir M. Perlaza, Gaetan Bisson, Iñaki Esnaola, Alain Jean-Marie, and Stefano Rini. "Empirical Risk Minimization with Relative Entropy Regularization". IEEE Transactions on Information Theory, vol. 70, no. 7, pp. 5122 – 5161, July, 2024.

Empirical Risk **Minimization** with Relative Entropy Regularization

Problem 1: ERM with Relative Entropy Regularization

$$\min_{P \in \Delta_Q(\mathcal{M})} \underbrace{\int L(z, \theta) dP(\theta)}_{R_z(P)} + \lambda D(P\|Q).$$

Solution:

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

Sensitivity to **deviations from the Optimal Measure**:

Lemma (Lemma 33 in [Perlaza-2024b])

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

Empirical Risk **Minimization** with Relative Entropy Regularization

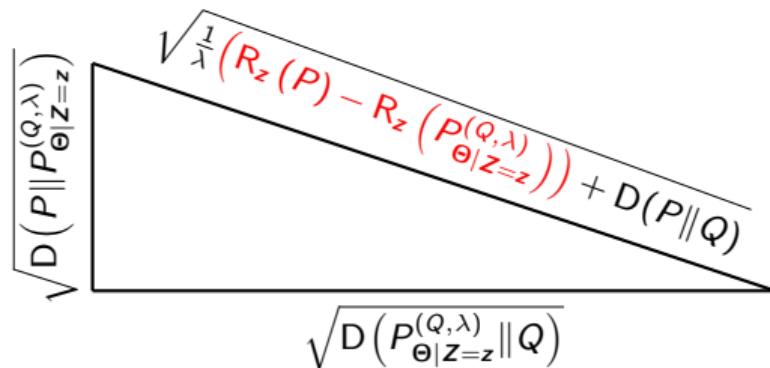


Figure: Geometric interpretation of the gap $R_z(P) - R_z(P_{\Theta|Z=z}^{(Q,\lambda)})$.

Empirical Risk **Minimization** with Relative Entropy Regularization

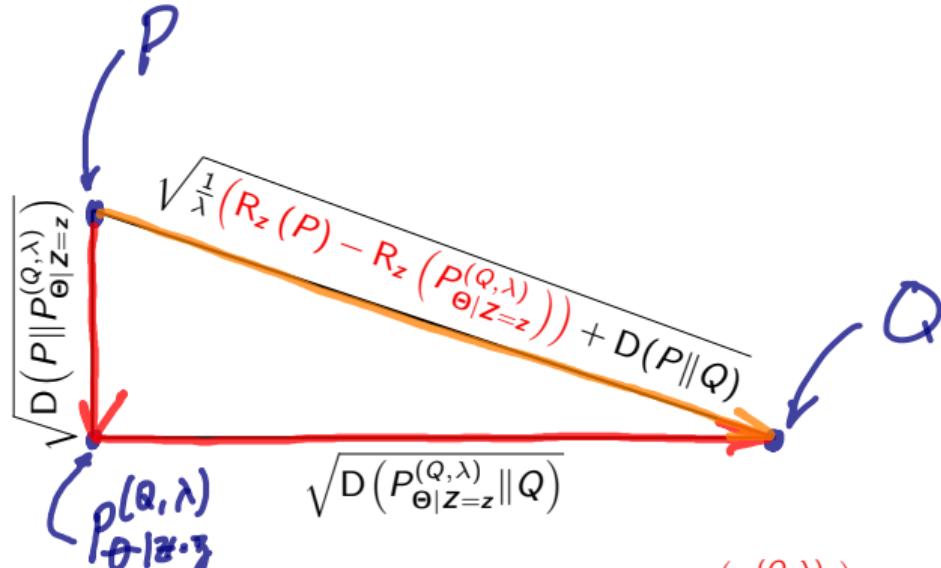


Figure: Geometric interpretation of the gap $R_z(P) - R_z(P_{\Theta|Z=z}^{(Q,\lambda)})$.

Empirical Risk **Minimization** with Relative Entropy Regularization

Problem 1: ERM with Relative Entropy Regularization

$$\min_{P \in \Delta_Q(\mathcal{M})} \underbrace{\int L(z, \theta) dP(\theta)}_{R_z(P)} + \lambda D(P\|Q).$$

Solution:

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

Theorem (Theorem 37 in [Perlaza-2024b])

For all $P_1 \in \Delta_Q(\mathcal{M})$ and $P_2 \in \Delta_Q(\mathcal{M})$,

$$R_z(P_1) - R_z(P_2) = \lambda \left(D(P_1 \| P_{\Theta|z=z}^{(Q,\lambda)}) - D(P_2 \| P_{\Theta|z=z}^{(Q,\lambda)}) + D(P_2 \| Q) - D(P_1 \| Q) \right).$$

Problem 2:
Loss Maximization with Relative Entropy Regularization

Loss Maximization with Relative Entropy Regularization

Problem 2: Loss Maximization with Relative Entropy Regularization

$$\max_{P \in \Delta_{P_S}(\mathcal{X} \times \mathcal{Y})} \int \ell(h(\theta, x), y) dP(x, y) - \beta D(P \| P_S),$$

with $\Delta_{P_S}(\mathcal{X} \times \mathcal{Y}) \triangleq \{P \in \Delta(\mathcal{X} \times \mathcal{Y}) : P \ll P_S\}$.

Notation:

$$J_{P_S, \theta}(t) = \log \left(\int \exp(t\ell(\theta, x, y)) dP_S(x, y) \right) \text{ and } \mathcal{J}_{P_S, \theta} \triangleq \left\{ t \in (0, +\infty) : J_{P_S, \theta} \left(\frac{1}{t} \right) < +\infty \right\}$$

Theorem (Theorem 1 in [Zou-2024])

If $\beta \in \mathcal{J}_{P_S, \theta}$, the solution to Problem 2 is unique, denoted by $P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}$, and satisfies for all $(x, y) \in \text{supp } P_S$,

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

[Zou-2024] Xinying Zou, Samir M. Perlaza, Iñaki Esnaola, Eitan Altman, and H. Vincent Poor. "The Worst-Case Data-Generating Probability Measure in Statistical Learning". IEEE Journal on Selected Areas in Information Theory, vol. 5, pp. 175–189, Apr., 2024.

Loss Maximization with Relative Entropy Regularization

Problem 2: Loss Maximization with Relative Entropy Regularization

$$\max_{P \in \Delta_{P_S}(\mathcal{X} \times \mathcal{Y})} \underbrace{\int \ell(h(\theta, x), y) dP(x, y) - \beta D(P \| P_S)}_{R_\theta(P)},$$

Solution:

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

Assumption: $P_Z(\mathcal{A}_1 \times \mathcal{A}_2 \times \dots \times \mathcal{A}_n) = \prod_{t=1}^n P_Z(\mathcal{A}_t)$

Lemma (Theorem 6 in [Zou-2024])

$$R_\theta(P) - R_\theta(P_{Z|\Theta=\theta}^{(P_S, \beta)}) = \beta \left(D(P \| P_S) - D\left(P \| P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}\right) - D\left(P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)} \| P_S\right) \right)$$

[Zou-2024] Xinying Zou, Samir M. Perlaza, Iñaki Esnaola, Eitan Altman, and H. Vincent Poor. "The Worst-Case Data-Generating Probability Measure in Statistical Learning". IEEE Journal on Selected Areas in Information Theory, vol. 5, pp. 175–189, Apr., 2024.

Loss Maximization with Relative Entropy Regularization

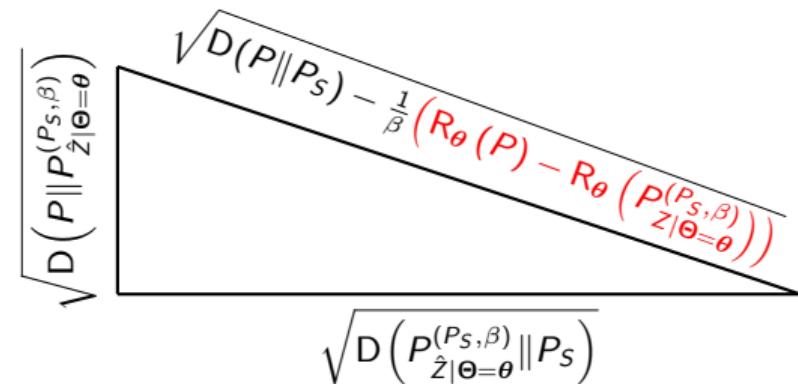
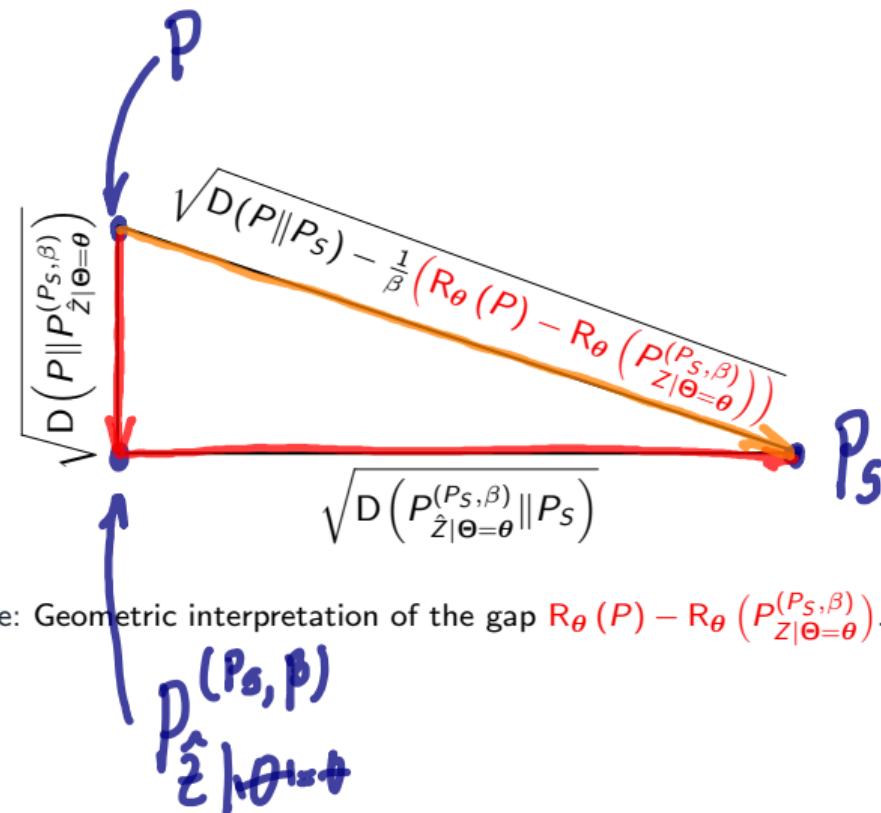


Figure: Geometric interpretation of the gap $R_\theta(P) - R_\theta(P_{Z|\Theta=\theta}^{(P_S, \beta)})$.

Loss Maximization with Relative Entropy Regularization



Loss Maximization with Relative Entropy Regularization

Problem 2: Loss Maximization with Relative Entropy Regularization

$$\max_{P \in \Delta_{P_S}(\mathcal{X} \times \mathcal{Y})} \underbrace{\int \ell(h(\theta, x), y) dP(x, y) - \beta D(P \| P_S)}_{R_\theta(P)},$$

Solution:

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

Assumption: $P_Z(\mathcal{A}_1 \times \mathcal{A}_2 \times \dots \times \mathcal{A}_n) = \prod_{t=1}^n P_Z(\mathcal{A}_t)$

Theorem (Theorem 8 in [Zou-2024])

For all $P_1 \in \Delta_{P_S}(\mathcal{X} \times \mathcal{Y})$ and for all $P_2 \in \Delta_{P_S}(\mathcal{X} \times \mathcal{Y})$,

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

So far...

Theorem (Theorem 37 in [Perlaza-2024b])

For all $P_1 \in \Delta_Q(\mathcal{M})$ and $P_2 \in \Delta_Q(\mathcal{M})$,

$$R_z(P_1) - R_z(P_2) = \lambda \left(D\left(P_1 \| P_{\Theta|Z=z}^{(Q,\lambda)}\right) - D\left(P_2 \| P_{\Theta|Z=z}^{(Q,\lambda)}\right) + D(P_2 \| Q) - D(P_1 \| Q) \right).$$

Theorem (Theorem 8 in [Zou-2024])

For all $P_1 \in \Delta_{P_S}(\mathcal{X} \times \mathcal{Y})$ and for all $P_2 \in \Delta_{P_S}(\mathcal{X} \times \mathcal{Y})$,

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

[Perlaza-2024b] Samir M. Perlaza and Xinying Zou. "The Generalization Error of Machine Learning Algorithms". November, 2024.

[Zou-2024] Xinying Zou, Samir M. Perlaza, Iñaki Esnaola, Eitan Altman, and H. Vincent Poor. "The Worst-Case Data-Generating Probability Measure in Statistical Learning". IEEE Journal on Selected Areas in Information Theory, vol. 5, pp. 175–189, Apr., 2024.

Table of Contents

Supervised Statistical Learning and **Generalization Error**

Empirical Risk Optimization with Relative Entropy Regularization

The Method of Gaps

Explicit Expressions for the Generalization Error

Concluding Remarks

The Method of Gaps

Definition (Expected Empirical Risk)

$$\begin{aligned} R_z(P) &= \int L(z, \theta) dP(\theta) \\ R_\theta(Q) &= \int \ell(h(\theta, x), y) dQ(x, y). \end{aligned}$$

The Method of Gaps

Definition (Expected Empirical Risk)

$$\begin{aligned} R_z(P) &= \int L(z, \theta) dP(\theta) \\ R_\theta(Q) &= \int \ell(h(\theta, x), y) dQ(x, y). \end{aligned}$$

Two **essential** observations:

- ▶ The generalization error is an expectation of the **variations of R_z or R_θ** ; and

The Method of Gaps

Definition (Expected Empirical Risk)

$$R_z(P) = \int L(z, \theta) dP(\theta)$$
$$R_\theta(Q) = \int \ell(h(\theta, x), y) dQ(x, y).$$

Two **essential** observations:

- ▶ The generalization error is an expectation of the **variations of R_z or R_θ** ; and
- ▶ These variations, a.k.a. **gaps**, exhibit closed-form expressions in terms of **information measures**.

The Method of Gaps

Definition (Expected Empirical Risk)

$$R_z(P) = \int L(z, \theta) dP(\theta)$$
$$R_\theta(Q) = \int \ell(h(\theta, x), y) dQ(x, y).$$

Two **essential** observations:

- ▶ The generalization error is an expectation of the **variations of R_z or R_θ** ; and
- ▶ These variations, a.k.a. **gaps**, exhibit closed-form expressions in terms of **information measures**.

Two-step Method:

- ▶ To express the generalization error as an expectation of a gap; and

The Method of Gaps

Definition (Expected Empirical Risk)

$$R_z(P) = \int L(z, \theta) dP(\theta)$$
$$R_\theta(Q) = \int \ell(h(\theta, x), y) dQ(x, y).$$

Two **essential** observations:

- ▶ The generalization error is an expectation of the **variations of R_z or R_θ** ; and
- ▶ These variations, a.k.a. **gaps**, exhibit closed-form expressions in terms of **information measures**.

Two-step Method:

- ▶ To express the generalization error as an expectation of a gap; and
- ▶ To leverage the properties of gaps to obtain closed-form expressions.

The Method of Gaps

Expected-Empirical-Risk Gaps

Definition (Expected Empirical Risk)

$$R_z(P) = \int L(z, \theta) dP(\theta)$$
$$R_\theta(Q) = \int \ell(h(\theta, x), y) dQ(x, y).$$

Definition (Expected-Empirical-Risk Gaps)

Let functionals $G : (\mathcal{X} \times \mathcal{Y})^n \times \Delta(\mathcal{M}) \times \Delta(\mathcal{M}) \rightarrow \mathbb{R}$ and $G : \mathcal{M} \times \Delta(\mathcal{X} \times \mathcal{Y}) \times \Delta(\mathcal{X} \times \mathcal{Y}) \rightarrow \mathbb{R}$ be

$$G(z, P_1, P_2) = R_z(P_1) - R_z(P_2), \text{ **Algorithm-driven Gap**}$$

and

$$G(\theta, P_1, P_2) = R_\theta(P_1) - R_\theta(P_2). \text{ **Data-driven Gap**}$$

The Method of Gaps

Two variants:

- ▶ The Method of **Algorithm-driven** Gaps
 - ▶ Central building-block: **The Gibbs Algorithm**
 - ▶ No assumptions on P_Z (probability distribution of the datasets)
- ▶ The Method of **Data-driven** Gaps
 - ▶ Central building-block: **The Worst-Case Data-Generating** (WCDG) probability measure
 - ▶ I.I.D assumption on P_Z :

$$P_Z (\mathcal{A}_1 \times \mathcal{A}_2 \times \dots \times \mathcal{A}_n) = \prod_{t=1}^n P_Z (\mathcal{A}_t)$$

The Method of **Algorithm-driven** Gaps

Generalization Error

The generalization error of the algorithm $P_{\Theta|Z}$ is

$$\overline{\overline{G}}(P_{\Theta|Z}, P_Z) \triangleq \int \int (R_u(P_{\Theta|Z=z}) - R_z(P_{\Theta|Z=z})) dP_Z(u) dP_Z(z).$$

The Method of **Algorithm-driven** Gaps

Generalization Error

The generalization error of the algorithm $P_{\Theta|Z}$ is

$$\overline{\overline{G}}(P_{\Theta|Z}, P_Z) \triangleq \int \int (R_u(P_{\Theta|Z=z}) - R_z(P_{\Theta|Z=z})) dP_Z(u) dP_Z(z).$$

Step 1:

Lemma (Lemma 3 in [Perlaza-2024b])

The generalization error $\overline{\overline{G}}(P_{\Theta|Z}, P_Z)$ satisfies

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

where for all measurable subsets \mathcal{C} of \mathcal{M} ,

$$P_\Theta(\mathcal{C}) = \int P_{\Theta|Z=z}(\mathcal{C}) dP_Z(z).$$

The Method of **Algorithm-driven** Gaps

Generalization Error

The generalization error of the algorithm $P_{\Theta|Z}$ is

$$\overline{\overline{G}}(P_{\Theta|Z}, P_Z) \triangleq \int \int (R_u(P_{\Theta|Z=z}) - R_z(P_{\Theta|Z=z})) dP_Z(u) dP_Z(z).$$

Step 2:

The Method of **Algorithm-driven** Gaps

Generalization Error

The generalization error of the algorithm $P_{\Theta|Z}$ is

$$\overline{\overline{G}}(P_{\Theta|Z}, P_Z) \triangleq \int \int (R_u(P_{\Theta|Z=z}) - R_z(P_{\Theta|Z=z})) dP_Z(u) dP_Z(z).$$

Step 2:

Lemma (Lemma 4 in [Perlaza-2024b])

The generalization error $\overline{\overline{G}}(P_{\Theta|Z}, P_Z)$ satisfies

$$\overline{\overline{G}}(P_{\Theta|Z}, P_Z) = \lambda \int \left(D(P_\Theta \| P_{\Theta|Z=z}^{(Q,\lambda)}) - D(P_{\Theta|Z=z} \| P_{\Theta|Z=z}^{(Q,\lambda)}) + D(P_{\Theta|Z=z} \| Q) - D(P_\Theta \| Q) \right) dP_Z(z).$$

The Method of Data-driven Gaps

Generalization Error

The generalization error of the algorithm $P_{\Theta|Z}$ is

$$\overline{\overline{G}}(P_{\Theta|Z}, P_Z) \triangleq \int \int (R_u(P_{\Theta|Z=z}) - R_z(P_{\Theta|Z=z})) dP_Z(u) dP_Z(z).$$

Step 1:

- ▶ **Assumption:** $P_Z(\mathcal{A}_1 \times \mathcal{A}_2 \times \dots \times \mathcal{A}_n) = \prod_{t=1}^n P_Z(\mathcal{A}_t)$.

The Method of Data-driven Gaps

Generalization Error

The generalization error of the algorithm $P_{\Theta|Z}$ is

$$\overline{\overline{G}}(P_{\Theta|Z}, P_Z) \triangleq \int \int (R_u(P_{\Theta|Z=z}) - R_z(P_{\Theta|Z=z})) dP_Z(u) dP_Z(z).$$

Step 1:

- **Assumption:** $P_Z(\mathcal{A}_1 \times \mathcal{A}_2 \times \dots \times \mathcal{A}_n) = \prod_{t=1}^n P_Z(\mathcal{A}_t)$.

Lemma (Lemma 6 in [Perlaza-2024b])

The generalization error $\overline{\overline{G}}(P_{\Theta|Z}, P_Z)$ satisfies

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

The Method of Data-driven Gaps

Generalization Error

The generalization error of the algorithm $P_{\Theta|Z}$ is

$$\overline{\overline{G}}(P_{\Theta|Z}, P_Z) \triangleq \int \int (R_u(P_{\Theta|Z=z}) - R_z(P_{\Theta|Z=z})) dP_Z(u) dP_Z(z).$$

Step 2:

- ▶ **Assumption:** $P_Z(\mathcal{A}_1 \times \mathcal{A}_2 \times \dots \times \mathcal{A}_n) = \prod_{t=1}^n P_Z(\mathcal{A}_t)$.

The Method of Data-driven Gaps

Generalization Error

The generalization error of the algorithm $P_{\Theta|Z}$ is

$$\overline{\overline{G}}(P_{\Theta|Z}, P_Z) \triangleq \int \int (R_u(P_{\Theta|Z=z}) - R_z(P_{\Theta|Z=z})) dP_Z(u) dP_Z(z).$$

Step 2:

- **Assumption:** $P_Z(\mathcal{A}_1 \times \mathcal{A}_2 \times \dots \times \mathcal{A}_n) = \prod_{t=1}^n P_Z(\mathcal{A}_t)$.

Lemma (Lemma 7 in [Perlaza-2024b])

The generalization error $\overline{\overline{G}}(P_{\Theta|Z}, P_Z)$ satisfies

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

So far...

Generalization Error

The generalization error of the algorithm $P_{\Theta|Z}$ is

$$\overline{\overline{G}}(P_{\Theta|Z}, P_Z) \triangleq \int \int (R_u(P_{\Theta|Z=z}) - R_z(P_{\Theta|Z=z})) dP_Z(u) dP_Z(z).$$

So far...

Generalization Error

The generalization error of the algorithm $P_{\Theta|Z}$ is

$$\overline{\overline{G}}(P_{\Theta|Z}, P_Z) \triangleq \int \int (R_u(P_{\Theta|Z=z}) - R_z(P_{\Theta|Z=z})) dP_Z(u) dP_Z(z).$$

Lemma (Lemma 4 in [Perlaza-2024b])

The generalization error $\overline{\overline{G}}(P_{\Theta|Z}, P_Z)$ satisfies

$$\overline{\overline{G}}(P_{\Theta|Z}, P_Z) = \lambda \int \left(D(P_\Theta \| P_{\Theta|Z=z}^{(Q,\lambda)}) - D(P_{\Theta|Z=z} \| P_{\Theta|Z=z}^{(Q,\lambda)}) + D(P_{\Theta|Z=z} \| Q) - D(P_\Theta \| Q) \right) dP_Z(z).$$

So far...

Generalization Error

The generalization error of the algorithm $P_{\Theta|Z}$ is

$$\overline{\overline{G}}(P_{\Theta|Z}, P_Z) \triangleq \int \int (R_u(P_{\Theta|Z=z}) - R_z(P_{\Theta|Z=z})) dP_Z(u) dP_Z(z).$$

Lemma (Lemma 4 in [Perlaza-2024b])

The generalization error $\overline{\overline{G}}(P_{\Theta|Z}, P_Z)$ satisfies

$$\overline{\overline{G}}(P_{\Theta|Z}, P_Z) = \lambda \int \left(D(P_\Theta \| P_{\Theta|Z=z}^{(Q,\lambda)}) - D(P_{\Theta|Z=z} \| P_{\Theta|Z=z}^{(Q,\lambda)}) + D(P_{\Theta|Z=z} \| Q) - D(P_\Theta \| Q) \right) dP_Z(z).$$

Lemma (Lemma 7 in [Perlaza-2024b])

The generalization error $\overline{\overline{G}}(P_{\Theta|Z}, P_Z)$ satisfies

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

Table of Contents

Supervised Statistical Learning and **Generalization Error**

Empirical Risk Optimization with Relative Entropy Regularization

The Method of Gaps

Explicit Expressions for the Generalization Error

Concluding Remarks

Expressions Obtained Via the Method of Gaps

- ▶ Particular **choices of the parameters**

Expressions Obtained Via the Method of Gaps

- ▶ Particular **choices of the parameters**
- ▶ **Algebraic manipulations** of the closed-form expressions shown before

Expressions Obtained Via the Method of Gaps

- ▶ Particular **choices of the parameters**
- ▶ **Algebraic manipulations** of the closed-form expressions shown before
- ▶ **More manipulations lead to less generality**
- ▶ Additional conditions to allow manipulations are imposed on:
 - ▶ The algorithm; and
 - ▶ The data-generating distribution.

Expressions Obtained Via the Method of Gaps

- ▶ Particular **choices of the parameters**
- ▶ **Algebraic manipulations** of the closed-form expressions shown before
- ▶ **More manipulations lead to less generality**
- ▶ Additional conditions to allow manipulations are imposed on:
 - ▶ The algorithm; and
 - ▶ The data-generating distribution.
- ▶ Some Expressions establish **bridges with other areas**: Hypothesis Testing, Geometry, etc.

Expressions Obtained Via the Method of Gaps

Connections to Information Measures

Theorem (Theorem 14 in [Perlaza-2024b])

The generalization error $\overline{\overline{G}}(P_{\Theta|Z}, P_Z)$ satisfies

$$\begin{aligned}\overline{\overline{G}}(P_{\Theta|Z}, P_Z) &= \lambda (I(P_{\Theta|Z}; P_Z) + L(P_{\Theta|Z}; P_Z)) \\ &\quad + \lambda \int \int \log \frac{dP_{\Theta|Z=z}}{dP_{\Theta|Z=z}^{(Q,\lambda)}}(\theta) dP_{\Theta}(\theta) dP_Z(z) - \lambda \int \int \log \frac{dP_{\Theta|Z=z}}{dP_{\Theta|Z=z}^{(Q,\lambda)}}(\theta) dP_{\Theta|Z=z}(\theta) dP_Z(z).\end{aligned}$$

Expressions Obtained Via the Method of Gaps

Connections to Information Measures

Theorem (Theorem 14 in [Perlaza-2024b])

The generalization error $\overline{\overline{G}}(P_{\Theta|Z}, P_Z)$ satisfies

$$\begin{aligned}\overline{\overline{G}}(P_{\Theta|Z}, P_Z) &= \lambda (I(P_{\Theta|Z}; P_Z) + L(P_{\Theta|Z}; P_Z)) \\ &\quad + \lambda \int \int \log \frac{dP_{\Theta|Z=z}}{dP_{\Theta|Z=z}^{(Q,\lambda)}}(\theta) dP_\Theta(\theta) dP_Z(z) - \lambda \int \int \log \frac{dP_{\Theta|Z=z}}{dP_{\Theta|Z=z}^{(Q,\lambda)}}(\theta) dP_{\Theta|Z=z}(\theta) dP_Z(z).\end{aligned}$$

What if...

$$\lambda \int \int \log \frac{dP_{\Theta|Z=z}}{dP_{\Theta|Z=z}^{(Q,\lambda)}}(\theta) dP_\Theta(\theta) dP_Z(z) - \lambda \int \int \log \frac{dP_{\Theta|Z=z}}{dP_{\Theta|Z=z}^{(Q,\lambda)}}(\theta) dP_{\Theta|Z=z}(\theta) dP_Z(z) = 0.$$

Expressions Obtained Via the Method of Gaps

Connections to Information Measures

Theorem (Theorem 14 in [Perlaza-2024b])

The generalization error $\overline{\overline{G}}(P_{\Theta|Z}, P_Z)$ satisfies

$$\begin{aligned}\overline{\overline{G}}(P_{\Theta|Z}, P_Z) &= \lambda(I(P_{\Theta|Z}; P_Z) + L(P_{\Theta|Z}; P_Z)) \\ &\quad + \lambda \int \int \log \frac{dP_{\Theta|Z=z}(\theta)}{dP_{\Theta|Z=z}^{(Q,\lambda)}(\theta)} dP_{\Theta}(\theta) dP_Z(z) - \lambda \int \int \log \frac{dP_{\Theta|Z=z}(\theta)}{dP_{\Theta|Z=z}^{(Q,\lambda)}(\theta)} dP_{\Theta|Z=z}(\theta) dP_Z(z).\end{aligned}$$

Generalization Error of the **Gibbs Algorithm**:

Corollary (Theorem 1 in [Aminian-2021])

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

[Aminian-2021] G Aminian, Y Bu, L Toni, M Rodrigues, G Wornell. "An exact characterization of the generalization error for the Gibbs algorithm" Advances in Neural Information Processing Systems, vol. 34, pp. 8106-8118, 2021

Expressions Obtained Via the Method of Gaps

Connections to Information Measures

Theorem (Theorem 29 in [Perlaza-2024b])

The generalization error $\overline{\overline{G}}(P_{\Theta|Z}, P_Z)$ satisfies

$$\begin{aligned}\overline{\overline{G}}(P_{\Theta|Z}, P_Z) &= -\beta(I(P_{Z|\Theta}; P_\Theta) + L(P_{Z|\Theta}; P_\Theta)) \\ &+ \beta \int \int \log \left(\frac{dP_{Z|\Theta=\theta}}{dP_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}}(z) \right) dP_{Z|\Theta=\theta}(z) dP_\Theta(\theta) - \beta \int \int \log \left(\frac{dP_{Z|\Theta=\theta}}{dP_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}}(z) \right) dP_Z(z) dP_\Theta(\theta).\end{aligned}$$

Expressions Obtained Via the Method of Gaps

Connections to Information Measures

Theorem (Theorem 29 in [Perlaza-2024b])

The generalization error $\overline{\overline{G}}(P_{\Theta|Z}, P_Z)$ satisfies

$$\begin{aligned}\overline{\overline{G}}(P_{\Theta|Z}, P_Z) &= -\beta(I(P_{Z|\Theta}; P_\Theta) + L(P_{Z|\Theta}; P_\Theta)) \\ &+ \beta \int \int \log \left(\frac{dP_{Z|\Theta=\theta}}{dP_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}}(z) \right) dP_{Z|\Theta=\theta}(z) dP_\Theta(\theta) - \beta \int \int \log \left(\frac{dP_{Z|\Theta=\theta}}{dP_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}}(z) \right) dP_Z(z) dP_\Theta(\theta).\end{aligned}$$

What if...

$$\beta \int \int \log \left(\frac{dP_{Z|\Theta=\theta}}{dP_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}}(z) \right) dP_{Z|\Theta=\theta}(z) dP_\Theta(\theta) - \beta \int \int \log \left(\frac{dP_{Z|\Theta=\theta}}{dP_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}}(z) \right) dP_Z(z) dP_\Theta(\theta) = 0.$$

Expressions Obtained Via the Method of Gaps

Connections to Euclidian Geometry

Theorem (Theorem 18 in [Perlaza-2024b])

The generalization error $\overline{\overline{G}}(P_{\Theta|Z}, P_Z)$ satisfies

$$\overline{\overline{G}}(P_{\Theta|Z}, P_Z) = \lambda \int \int \left(D(P_{\Theta|Z=z} \| P_{\Theta|Z=u}^{(Q,\lambda)}) - D(P_{\Theta|Z=z} \| P_{\Theta|Z=z}^{(Q,\lambda)}) \right) dP_Z(u) dP_Z(z).$$

Expressions Obtained Via the Method of Gaps

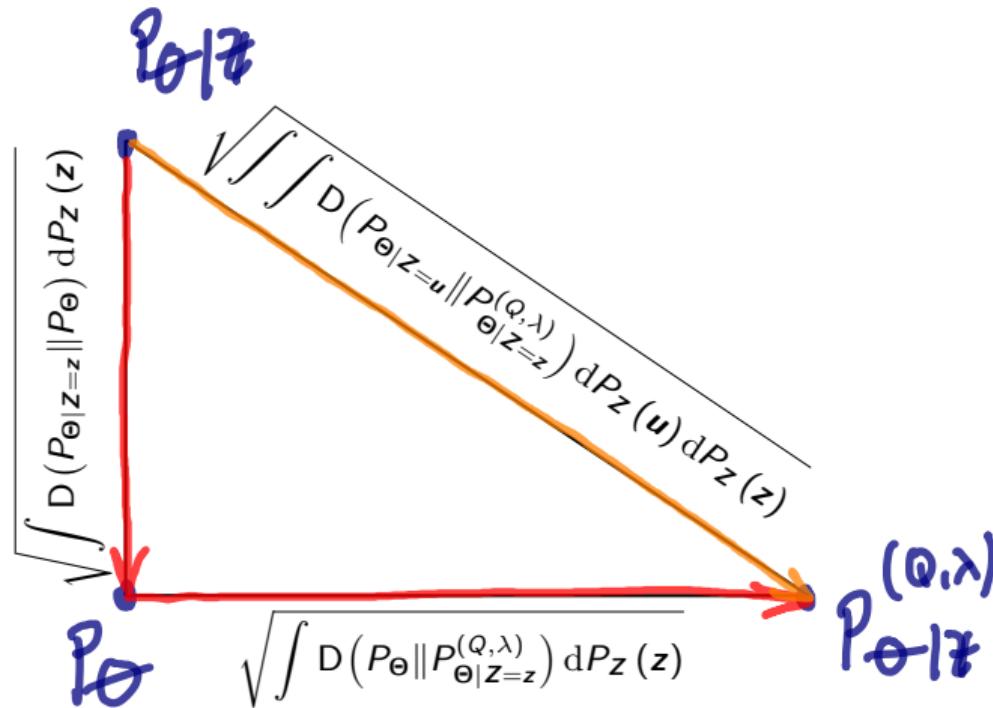
Connections to Euclidian Geometry

$$\frac{\sqrt{\int D(P_{\Theta|Z=u} \| P_{\Theta|Z=z}^{(Q,\lambda)}) dP_Z(u) dP_Z(z)}}{\sqrt{\int D(P_\Theta \| P_{\Theta|Z=z}^{(Q,\lambda)}) dP_Z(z)}}$$

$$\int \int D(P_{\Theta|Z=u} \| P_{\Theta|Z=z}^{(Q,\lambda)}) dP_Z(u) dP_Z(z) = \int D(P_\Theta \| P_{\Theta|Z=z}^{(Q,\lambda)}) dP_Z(z) + \int D(P_{\Theta|Z=z} \| P_\Theta) dP_Z(z).$$

Expressions Obtained Via the Method of Gaps

Connections to Euclidian Geometry



$$\int \int D(P_{\Theta|z=u} \| P_{\Theta|z=z}^{(Q,\lambda)}) dP_Z(u) dP_Z(z) = \int D(P_\Theta \| P_{\Theta|z=z}^{(Q,\lambda)}) dP_Z(z) + \int D(P_{\Theta|z=z} \| P_\Theta) dP_Z(z).$$

Expressions Obtained Via the Method of Gaps

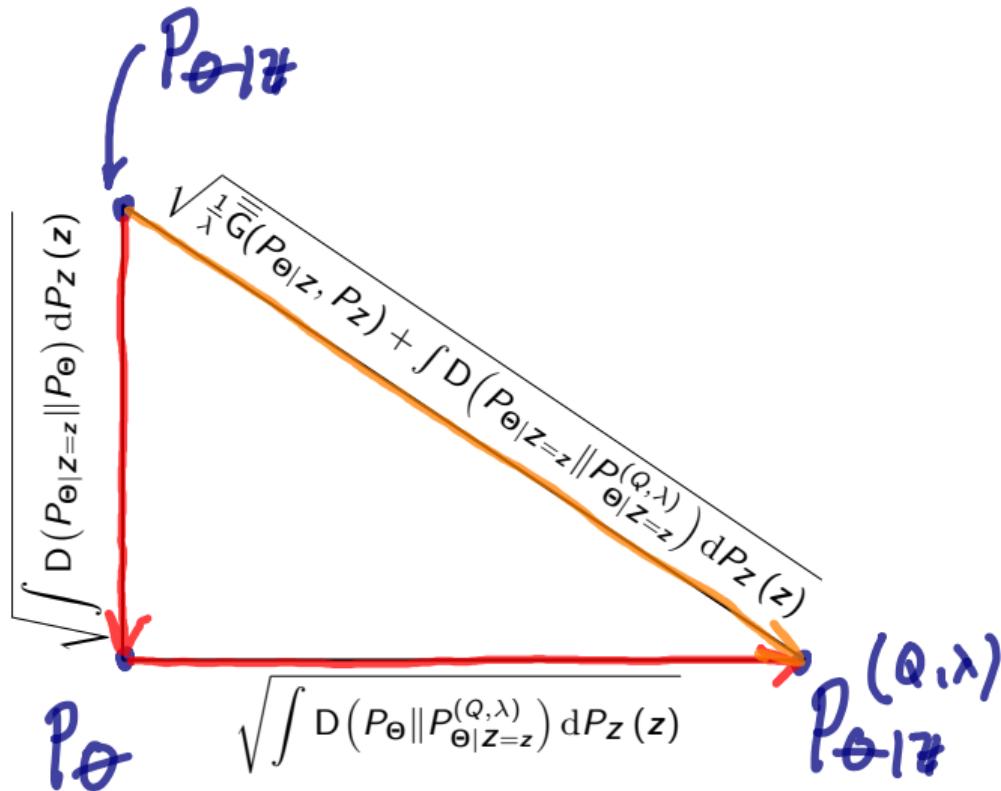
Connections to Euclidian Geometry

$$\sqrt{\int D(P_{\Theta|Z=z} \| P_{\Theta}) dP_Z(z)}$$

$\sqrt{\frac{1}{\lambda} G(P_{\Theta|Z}, P_Z) + \int D(P_{\Theta|Z=z} \| P_{\Theta|Z=z}^{(Q,\lambda)}) dP_Z(z)}$

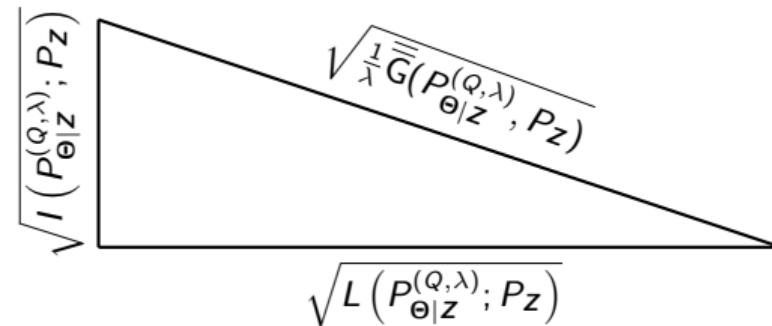
Expressions Obtained Via the Method of Gaps

Connections to Euclidian Geometry



Expressions Obtained Via the Method of Gaps

Connections to Euclidian Geometry



Expressions Obtained Via the Method of Gaps

Connections to Euclidian Geometry

Theorem (Theorem 31 in [Perlaza-2024b])

The generalization error $\overline{\overline{G}}(P_{\Theta|Z}, P_Z)$ satisfies

$$\overline{\overline{G}}(P_{\Theta|Z}, P_Z) = \beta \int \int \left(D(P_{Z|\Theta=\mu} \| P_{\hat{Z}|\Theta=\mu}^{(P_S, \beta)}) - D(P_{Z|\Theta=\mu} \| P_{\hat{Z}|\Theta=\nu}^{(P_S, \beta)}) \right) dP_\Theta(\nu) dP_\Theta(\mu).$$

Expressions Obtained Via the Method of Gaps

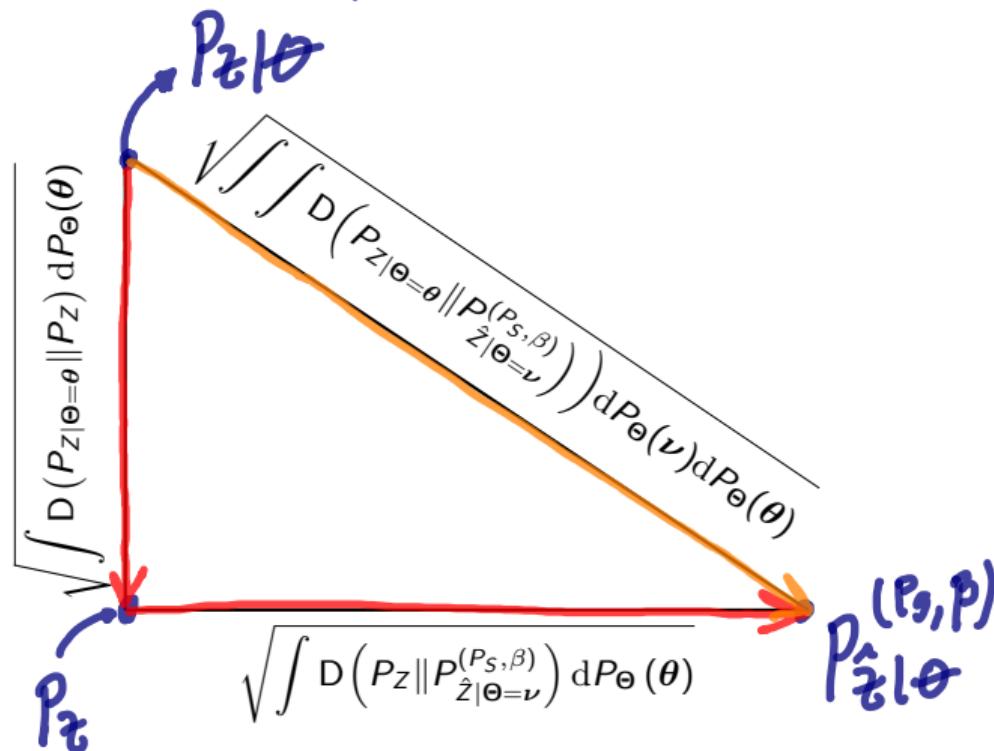
Connections to Euclidian Geometry

$$\sqrt{\int \int D(P_{Z|\Theta=\theta} \| P_Z) dP_\Theta(\theta)} = \sqrt{\int \int D(P_{Z|\Theta=\theta} \| P_{\hat{Z}|\Theta=\nu}^{(P_S, \beta)}) dP_\Theta(\nu) dP_\Theta(\theta)}$$
$$\sqrt{\int D(P_Z \| P_{\hat{Z}|\Theta=\nu}^{(P_S, \beta)}) dP_\Theta(\theta)}$$

$$\int \int D(P_{Z|\Theta=\theta} \| P_{\hat{Z}|\Theta=\nu}^{(P_S, \beta)}) dP_\Theta(\nu) dP_\Theta(\theta) = \int D(P_{Z|\Theta=\theta} \| P_Z) dP_\Theta(\theta) + \int D(P_Z \| P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}) dP_\Theta(\theta),$$

Expressions Obtained Via the Method of Gaps

Connections to Euclidian Geometry



$$\int \int D(P_{Z|\Theta=\theta} \| P_{\hat{Z}|\Theta=\nu}^{(P_S, \beta)}) dP_\Theta(\nu) dP_\Theta(\theta) = \int D(P_{Z|\Theta=\theta} \| P_Z) dP_\Theta(\theta) + \int D(P_Z \| P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}) dP_\Theta(\theta),$$

Expressions Obtained Via the Method of Gaps

Connections to Euclidian Geometry

$$\sqrt{\int D(P_{Z|\Theta=\theta} \| P_Z) dP_\Theta(\theta)} - \frac{1}{\beta} \bar{G}(P_{\Theta|Z}, P_Z)$$
$$\sqrt{\int D(P_Z \| P_{\hat{Z}|\Theta=\theta}^{(P_S, \beta)}) dP_\Theta(\theta)}$$

Expressions Obtained Via the Method of Gaps

Connections to Euclidian Geometry

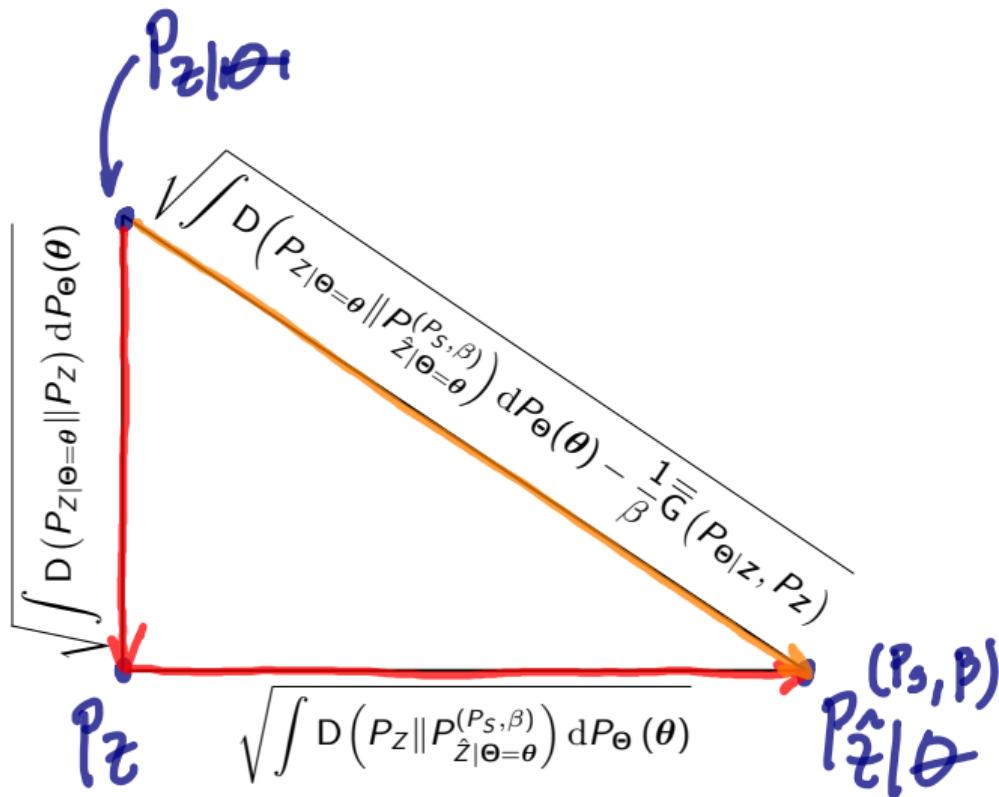


Table of Contents

Supervised Statistical Learning and **Generalization Error**

Empirical Risk Optimization with Relative Entropy Regularization

The Method of Gaps

Explicit Expressions for the Generalization Error

Concluding Remarks

Some Concluding Remarks

- ▶ Solution to the **Empirical Risk Optimization with Relative Entropy regularization**
 - ▶ Maximization → **Worst-Case Data-Generating Probability Measure**
 - ▶ Minimization → **Gibbs Algorithm**

Some Concluding Remarks

- ▶ Solution to the **Empirical Risk Optimization with Relative Entropy regularization**
 - ▶ Maximization → **Worst-Case Data-Generating Probability Measure**
 - ▶ Minimization → **Gibbs Algorithm**
- ▶ Method of Gaps
 - ▶ **Algorithm-driven** gaps → uses **Gibbs Algorithm**
 - ▶ **Data-driven** gaps → uses **Worst-Case Data-Generating Probability Measure**

Some Concluding Remarks

- ▶ Solution to the **Empirical Risk Optimization with Relative Entropy regularization**
 - ▶ Maximization → **Worst-Case Data-Generating Probability Measure**
 - ▶ Minimization → **Gibbs Algorithm**
- ▶ Method of Gaps
 - ▶ **Algorithm-driven** gaps → uses **Gibbs Algorithm**
 - ▶ **Data-driven** gaps → uses **Worst-Case Data-Generating Probability Measure**
- ▶ **Generalization error** obtained via
 - ▶ Expectation of **Algorithm-driven** gaps
 - ▶ Expectation of **Data-driven** gaps

WHAT IS THE LONG-RUN DISTRIBUTION OF STOCHASTIC GRADIENT DESCENT? A LARGE DEVIATIONS ANALYSIS

WAÏSS AZIZIAN^{c,*}, FRANCK IUTZELER[#],
JÉRÔME MALICK^{*}, AND PANAYOTIS MERTIKOPOULOS^o

ABSTRACT. In this paper, we examine the long-run distribution of stochastic gradient descent (SGD) in general, non-convex problems. Specifically, we seek to understand which regions of the problem’s state space are more likely to be visited by SGD, and by how much. Using an approach based on the theory of large deviations and randomly perturbed dynamical systems, we show that the long-run distribution of SGD resembles the Boltzmann–Gibbs distribution of equilibrium thermodynamics with temperature equal to the method’s step-size and energy levels determined by the problem’s objective and the statistics of the noise. In particular, we show that, in the long run, (a) the problem’s critical region is visited exponentially more often than any non-critical region;

Examples

Corollary (What is the long-run Generalization Error of Stochastic Gradient Descent ?)

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

Thank you for your attention!

Questions/Comments/Typos: samir.perlaza@inria.fr

► This work appears in:

- Samir M. Perlaza and Xinying Zou. "**The Generalization Error of Machine Learning Algorithms**". November, 2024.
- Samir M. Perlaza, Gaetan Bisson, Iñaki Esnaola, Alain Jean-Marie, and Stefano Rini. "**Empirical Risk Minimization with Relative Entropy Regularization**". IEEE Transactions on Information Theory, vol. 70, no. 7, pp. 5122 – 5161, July, 2024.
- Xinying Zou, Samir M. Perlaza, Iñaki Esnaola, Eitan Altman, and H. Vincent Poor. "**The Worst-Case Data-Generating Probability Measure in Statistical Learning**". IEEE Journal on Selected Areas in Information Theory, vol. 5, pp. 175–189, Apr., 2024.

