

Tutorial

Characterizing the Generalization Error of Machine Learning Algorithms via Information Measures

Gholamali Aminian, Yuheng Bu, Iñaki Esnaola, and Samir M. Perlaza

2024 IEEE Information Theory Workshop

The 24th of November, 2024
Shenzhen, China

Slides for Part II



Table of Contents (Part II)

Generalization Error of Gibbs Algorithm

- Gibbs algorithm

- Exact Characterizations

- Other Properties

Optimal Inverse Temperature

Asymptotic Behavior and Information Criteria for Model Selection

Generalization Error of Transfer learning

Conclusion

Preliminaries

Information Measures

▶ KL divergence: $\text{KL}(P\|Q) \triangleq \int_{\mathcal{X}} \log\left(\frac{dP}{dQ}\right) dP$

▶ Symmetrized KL divergence (Jeffrey's divergence)

$$D_{\text{SKL}}(P\|Q) \triangleq \text{KL}(P\|Q) + \text{KL}(Q\|P).$$

▶ Mutual information: $I(X; Y) \triangleq \text{KL}(P_{X,Y}\|P_X \otimes P_Y)$

▶ Lautum information [Palomar and Verdú, 2008]: $L(X; Y) \triangleq \text{KL}(P_X \otimes P_Y\|P_{X,Y})$

▶ Symmetrized KL information [Aminian et al., 2015]:

$$I_{\text{SKL}}(X; Y) \triangleq D_{\text{SKL}}(P_{X,Y}\|P_X \otimes P_Y) = I(X; Y) + L(X; Y).$$

Information-theoretic Generalization Bounds

Lemma ([Xu and Raginsky, 2017])

Suppose $\ell(w, Z)$ is σ -sub-Gaussian under $Z \sim \mu$ for all $w \in \mathcal{W}$, then

$$|\text{gen}(\mu, P_{W|S})| \leq \sqrt{\frac{2\sigma^2}{n} I(S; W)}.$$

- ▶ Depends on every ingredient in the supervised learning problem
- ▶ Reducing dependence between W and S leads to better generalization bound
- ▶ This bound is only tight if $I(S; W) = 0$ and $\text{gen}(\mu, P_{W|S}) = 0$
- ▶ Multiple techniques to improve this result, including ISMI [Bu et al., 2020], CMI [Steinke and Zakynthinou, 2020], f -CMI [Harutyunyan et al., 2021], ΔL -CMI [Wang and Mao, 2023]

Table of Contents (Part II)

Generalization Error of Gibbs Algorithm

Gibbs algorithm

Exact Characterizations

Other Properties

Optimal Inverse Temperature

Asymptotic Behavior and Information Criteria for Model Selection

Generalization Error of Transfer learning

Conclusion

Regularized ERM problem

- ▶ How can we use this result to develop a better learning algorithm?
- ▶ Regularizing mutual information $I(S; W)$ during ERM

$$P_{W|S}^* = \arg \min_{P_{W|S}} \left(\mathbb{E}_{P_{W,S}} [L_E(W, S)] + \frac{1}{\gamma} I(S; W) \right)$$

- ▶ inverse temperature $\gamma \geq 0$ balances between fitting and generalization
- ▶ Replacing $I(S; W)$ with $\text{KL}(P_{W|S} || \pi(W) | P_S)$ for any prior $\pi(W)$
- ▶ It gives **information risk minimization** (IRM) problem

$$P_{W|S}^* = \arg \min_{P_{W|S}} \left(\mathbb{E}_{P_{W,S}} [L_E(W, S)] + \frac{1}{\gamma} \text{KL}(P_{W|S} || \pi(W) | P_S) \right)$$

Lemma ([Zhang, 2006, Xu and Raginsky, 2017])

Solution to IRM problem is $(\gamma, \pi(w), L_E(w, s))$ -Gibbs distribution

$$P_{W|S}^\gamma(w|s) \triangleq \frac{\pi(w)e^{-\gamma L_E(w,s)}}{V(s,\gamma)}, \quad \gamma \geq 0,$$

where $V(s,\gamma) \triangleq \int \pi(w)e^{-\gamma L_E(w,s)} dw$ is partition function.

Proof.

For any learning algorithm $P_{W|S}$ with fixed $S = s$,

$$\begin{aligned} 0 &\leq \text{KL}(P_{W|S=s} \| P_{W|S=s}^\gamma) \\ &= \mathbb{E}_{P_{W|S=s}} \left[\log \frac{P_{W|S=s} \cdot V(s,\gamma)}{\pi(W) \cdot e^{-\gamma L_E(w,s)}} \right] \\ &= \text{KL}(P_{W|S=s} \| \pi(W)) + \log V(s,\gamma) + \gamma \mathbb{E}_{P_{W|S=s}} [L_E(w,s)]. \end{aligned}$$

$$\min_{P_{W|S}} \mathbb{E}_{P_{W|S=s}} [L_E(W,s)] + \frac{1}{\gamma} \text{KL}(P_{W|S=s} \| \pi) = -\frac{1}{\gamma} \log V(s,\gamma). \quad \square$$

Gibbs Algorithm

We focus on the generalization error of Gibbs algorithm (distribution)

$(\gamma, \pi(w), L_E(w, s))$ -Gibbs distribution:

$$P_{W|S}^\gamma(w|s) \triangleq \frac{\pi(w)e^{-\gamma L_E(w,s)}}{V(s,\gamma)}, \quad \gamma \geq 0$$

where

- ▶ inverse temperature γ , reduces to standard ERM if $\gamma \rightarrow \infty$
- ▶ $\pi(w)$ arbitrary prior distribution of W
- ▶ $V(s, \gamma) \triangleq \int \pi(w)e^{-\gamma L_E(w,s)} dw$ partition function

Practical Implementation of Gibbs algorithm

- ▶ Stochastic Gradient Langevin Dynamics (SGLD)
- ▶ Metropolis adjusted Langevin algorithm (MALA)

The SGLD can be viewed as the noisy version of SGD,

$$W_{k+1} = W_k - \eta_t \nabla L_E(W_k, s) + \sqrt{\frac{2\eta_t}{\gamma}} \zeta_k, \quad k = 0, 1, \dots,$$

where ζ_k standard Gaussian random vector; $\eta_t > 0$ step size.

- ▶ [Raginsky et al., 2017] shows that $P_{W_k|S}$ induced by SGLD converges to $(\gamma, \pi(W_0), L_E(w_k, s))$ -Gibbs distribution for sufficiently large k
- ▶ MALA is SGLD with Metropolis rejection, faster convergence

Table of Contents (Part II)

Generalization Error of Gibbs Algorithm

Gibbs algorithm

Exact Characterizations

Other Properties

Optimal Inverse Temperature

Asymptotic Behavior and Information Criteria for Model Selection

Generalization Error of Transfer learning

Conclusion

Expected Generalization Error

An **exact** characterization of generalization error for Gibbs algorithm

Theorem

For $(\gamma, \pi(w), L_E(w, s))$ -Gibbs algorithm,

$$P_{W|S}^\gamma(w|s) = \frac{\pi(w)e^{-\gamma L_E(w,s)}}{V(s, \gamma)}, \quad \gamma > 0,$$

the expected generalization error is

$$\overline{\text{gen}}(P_{W|S}^\gamma, P_S) = \frac{I_{\text{SKL}}(W; S)}{\gamma}.$$

- ▶ Highlights the fundamental role of $I_{\text{SKL}}(W; S)$ in learning theory
- ▶ Holds even for non-i.i.d training samples

G. Aminian*, Y. Bu*, L. Toni, M. R. Rodrigues, G. W. Wornell. "An Exact Characterization of the Generalization Error for the Gibbs Algorithm," in *Proc. Conference on Neural Information Processing Systems (NeurIPS)*, Dec. 2021.

Generalization Error of Gibbs Algorithm

Theorem

For Gibbs algorithm $P_{W|S}^\gamma(w|s) = \frac{\pi(w)e^{-\gamma L_E(w,s)}}{V(s,\gamma)}$,

$$\overline{\text{gen}}(P_{W|S}^\gamma, P_S) = \frac{I_{\text{SKL}}(W; S)}{\gamma}.$$

Sketch of Proof:

Symmetrized KL information can be written as

$$\begin{aligned} I_{\text{SKL}}(W; S) &= \mathbb{E}_{P_{W,S}} \left[\log \left(\frac{P_{W|S}^\gamma}{P_W} \right) \right] + \mathbb{E}_{P_W \otimes P_S} \left[\log \left(\frac{P_W}{P_{W|S}^\gamma} \right) \right] \\ &= \mathbb{E}_{P_{W,S}} \left[\log(P_{W|S}^\gamma) \right] - \mathbb{E}_{P_W \otimes P_S} \left[\log(P_{W|S}^\gamma) \right] \end{aligned}$$

Note that $P_{W,S}$ and $P_W \otimes P_S$ share the same marginal distribution,

$$\begin{aligned} I_{\text{SKL}}(W; S) &= \mathbb{E}_{P_{W,S}} [-\gamma L_E(W, S)] - \mathbb{E}_{P_W \otimes P_S} [-\gamma L_E(W, S)] \\ &= \gamma \overline{\text{gen}}(P_{W|S}^\gamma, P_S) \end{aligned}$$

□

Empirical risk of Gibbs algorithm

Theorem

$\log V(s, \gamma)$ is convex and differentiable infinitely many times with respect to γ . In particular,

$$\mathbb{E}_\gamma[L_E(W, s)] = -\frac{\partial \log V(s, \gamma)}{\partial \gamma},$$
$$\text{Var}_\gamma[L_E(W, s)] = \frac{\partial^2 \log V(s, \gamma)}{\partial \gamma^2},$$

where $\mathbb{E}_\gamma[\cdot] \triangleq \mathbb{E}_{P_{W|S=s}^\gamma}[\cdot]$, and $\text{Var}_\gamma[L_E(W, s)] \triangleq \mathbb{E}_\gamma[L_E(W, s)^2] - \mathbb{E}_\gamma[L_E(W, s)]^2$.

Expected **empirical risk** of the Gibbs algorithm is non-increasing w.r.t γ

- ▶ Monotonicity: $L_E(W, s)$ is non-increasing with γ
- ▶ Sub-Gaussianity: $L_E(W, s)$ is sub-Gaussian under Gibbs algorithm if $\text{Var}_\gamma[L_E(W, s)]$ is bounded

Perlaza, Samir M., Gaetan Bisson, Iñaki Esnaola, Alain Jean-Marie, and Stefano Rini. "Empirical risk minimization with relative entropy regularization," *IEEE Trans. Inf. Theory*, 2024.

Table of Contents (Part II)

Generalization Error of Gibbs Algorithm

Gibbs algorithm

Exact Characterizations

Other Properties

Optimal Inverse Temperature

Asymptotic Behavior and Information Criteria for Model Selection

Generalization Error of Transfer learning

Conclusion

Tighter Generalization Error Upper Bounds

Why do we care about upper bounds when we have exact characterization?

- ▶ Quantify how $\overline{\text{gen}}(P_{W|S}^\gamma, P_S)$ depends on number of i.i.d. samples n
- ▶ Useful when directly evaluating $I_{\text{SKL}}(W; S)$ is hard

Theorem

Suppose that

- ▶ $S = \{Z_i\}_{i=1}^n$ are i.i.d generated from the distribution P_Z
- ▶ $\ell(w, Z)$ is σ -sub-Gaussian
- ▶ $C_E \leq \frac{L(W; S)}{I(W; S)}$ for some $C_E \geq 0$,

$$\overline{\text{gen}}(P_{W|S}^\gamma, P_S) \leq \frac{2\sigma^2\gamma}{(1 + C_E)n}.$$

G. Aminian*, Y. Bu*, L. Toni, M. R. Rodrigues, G. W. Wornell. "An Exact Characterization of the Generalization Error for the Gibbs Algorithm," in *Proc. Conference on Neural Information Processing Systems (NeurIPS)*, Dec. 2021.

Sketch of Proof:

Recall the mutual information-based bound,

$$\begin{aligned}\sqrt{\frac{2\sigma^2}{n}I(S; W)} &\geq \overline{\text{gen}}(P_{W|S}^\gamma, P_S) = \frac{I(W; S) + L(W; S)}{\gamma} \\ &\geq \frac{(1 + C_E)}{\gamma} I(W; S)\end{aligned}$$

$$\overline{\text{gen}}(P_{W|S}^\gamma, P_S) \leq \sqrt{\frac{2\sigma^2}{n}I(S; W)} \leq \frac{2\sigma^2\gamma}{(1 + C_E)n}$$

□

[Choice of C_E]

- ▶ $C_E = 0$ is always valid, which gives $\overline{\text{gen}}(P_{W|S}^\gamma, P_S) \leq \frac{2\sigma^2\gamma}{n}$
- ▶ $C_E = 1$, $L(S; W) \geq I(S; W)$ holds for any Gaussian channel $P_{W|S}$

Example: Mean Estimation

- ▶ Learning mean $\boldsymbol{\mu} \in \mathbb{R}^d$ of Z using n i.i.d training samples $S = \{\mathbf{z}_i\}_{i=1}^n$
- ▶ Not necessary Gaussian, but covariance matrix $\Sigma_Z = \sigma_Z^2 \mathbf{I}_d$
- ▶ Mean-squared loss $\ell(\mathbf{w}, \mathbf{z}) = \|\mathbf{z} - \mathbf{w}\|_2^2$
- ▶ Gaussian prior $\pi(\mathbf{w}) = \mathcal{N}(\boldsymbol{\mu}_0, \sigma_0^2 \mathbf{I}_d)$
- ▶ Then, $(\gamma, \mathcal{N}(\boldsymbol{\mu}_0, \sigma_0^2 \mathbf{I}_d), L_E(\mathbf{w}, s))$ -Gibbs algorithm is given by the following Gaussian posterior

$$P_{W|S}^\gamma(\mathbf{w}|\mathbf{z}^n) \sim \mathcal{N}\left(\alpha \boldsymbol{\mu}_0 + (1 - \alpha) \bar{\mathbf{z}}, \alpha \sigma_0^2 \mathbf{I}_d\right),$$

with

$$\alpha \triangleq \frac{1}{2\sigma_0^2\gamma + 1}, \quad \bar{\mathbf{z}} \triangleq \frac{1}{n} \sum_{i=1}^n \mathbf{z}_i.$$

Example: Mean Estimation

Since $P_{W|S}^\gamma$ is Gaussian,

$$I(S; W) = \frac{d\sigma_0^2\sigma_Z^2\gamma}{(n\sigma_0^2 + \frac{1}{2\gamma})} - \text{KL}(P_W \| \mathcal{N}(\boldsymbol{\mu}_W, \sigma_1^2 I_d)),$$
$$L(S; W) = \frac{d\sigma_0^2\sigma_Z^2\gamma}{(n\sigma_0^2 + \frac{1}{2\gamma})} + \text{KL}(P_W \| \mathcal{N}(\boldsymbol{\mu}_W, \sigma_1^2 I_d)),$$

with $\boldsymbol{\mu}_W = \alpha\boldsymbol{\mu}_0 + (1 - \alpha)\boldsymbol{\mu}$.

The generalization error can be computed exactly as:

$$\overline{\text{gen}}(P_{W|S}^\gamma, P_S) = \frac{I_{\text{SKL}}(W; S)}{\gamma} = \frac{2d\sigma_0^2\sigma_Z^2}{n(\sigma_0^2 + \frac{1}{2\gamma})}.$$

As a comparison, the ISMI-based bound gives a sub-optimal bound $\mathcal{O}(1/\sqrt{n})$, as $n \rightarrow \infty$.

Check Point

Generalization error or empirical risk is one part of the story

Our goal is to **design** (or guide the design) algorithms that minimize population risk.

There are three elements in $(\gamma, \pi(w), L_E(w, s))$ -Gibbs algorithm

- ▶ inverse temperature $\gamma \rightarrow$ Optimal hyper-parameter
- ▶ empirical risk $L_E(w, s)$, or model family \rightarrow Information criteria for model selection
- ▶ prior distribution $\pi(w) \rightarrow$ Transfer learning

Table of Contents (Part II)

Generalization Error of Gibbs Algorithm

Gibbs algorithm

Exact Characterizations

Other Properties

Optimal Inverse Temperature

Asymptotic Behavior and Information Criteria for Model Selection

Generalization Error of Transfer learning

Conclusion

Expected Test Loss

For fixed training data s and testing data s' , consider expected test loss

$$L_P(\gamma, s, s') \triangleq \mathbb{E}_\gamma[L_E(W, s')],$$

and expected generalization error

$$\overline{\text{gen}}(\gamma, s, s') \triangleq \mathbb{E}_\gamma[L_E(W, s') - L_E(W, s)].$$

Theorem

For $\gamma \geq 0$ such that $\log V(s, \gamma) < \infty$, the first order derivative of the expected test loss is given by

$$\frac{\partial}{\partial \gamma} L_P(\gamma, s, s') = -\text{Cov}_\gamma[L_E(W, s'), L_E(W, s)],$$

with

$$\text{Cov}_\gamma[L_E(W, s'), L_E(W, s)] \triangleq \mathbb{E}_\gamma[L_E(W, s)L_E(W, s')] - \mathbb{E}_\gamma[L_E(W, s)]\mathbb{E}_\gamma[L_E(W, s')].$$

$\text{Cov}_\gamma[L_E(W, s'), L_E(W, s)]$ can be positive/negative, no monotonicity

Expected Generalization Error

Corollary

For $\gamma \geq 0$ such that $\log V(s, \gamma) < \infty$, the first order derivative of the expected generalization error is given by

$$\frac{\partial}{\partial \gamma} \overline{\text{gen}}(\gamma, s, s') = \text{Var}_\gamma(L_E(W, s)) - \text{Cov}_\gamma[L_E(W, s'), L_E(W, s)].$$

- ▶ Cannot show that the $\overline{\text{gen}}$ is non-decreasing, Cauchy-Schwarz Inequality only guarantees that

$$|\text{Cov}_\gamma[L_E(W, s'), L_E(W, s)]| \leq \sqrt{\text{Var}_\gamma(L_E(W, s))\text{Var}_\gamma(L_E(W, s'))}.$$

- ▶ [Aminian et al., 2021] provides a bound of order $\mathcal{O}(\frac{\gamma}{n})$ by simply combining the I_{SKL} characterization with the MI bound, which may hint that $\overline{\text{gen}}$ is always increasing with γ .
- ▶ However, we will illustrate how $\overline{\text{gen}}$ rises from zero and then decreases as γ increases.

Example: Mean Estimation

$(\gamma, \mathcal{N}(\boldsymbol{\mu}_0, \sigma_0^2 \mathbf{I}_d), L_E(\mathbf{w}, s))$ -Gibbs algorithm is given by the following Gaussian posterior

$$P_{W|S}^\gamma(\mathbf{w}|\mathbf{z}^n) \sim \mathcal{N}(\alpha\boldsymbol{\mu}_0 + (1 - \alpha)\bar{\mathbf{z}}, \alpha\sigma_0^2 \mathbf{I}_d)$$

Population risk has the following exact characterization

$$\begin{aligned} L_P(P_{W|S}^\gamma, P_S) &= \underbrace{\frac{4d\sigma_0^2\sigma_Z^2\gamma}{n(1 + 2\sigma_0^2\gamma)}}_{\text{generalization error}} + \underbrace{\frac{\|\boldsymbol{\mu}_0 - \boldsymbol{\mu}\|_2^2 + d\sigma_Z^2/n}{(1 + 2\sigma_0^2\gamma)^2} + \frac{d\sigma_0^2}{1 + 2\sigma_0^2\gamma} + \frac{n-1}{n}d\sigma_Z^2}_{\text{empirical risk}}. \end{aligned}$$

To find optimal γ minimizes L_P

- ▶ Optimize over γ using the above equation directly
- ▶ Evaluate the derivative of $L_P(\gamma, s, s')$ by computing covariance

Example: Mean Estimation

γ^* depends on other parameters of the problem in a non-trivial manner

$$\gamma^* = \begin{cases} +\infty, & \text{if } \frac{\sigma_Z^2}{n} \in [0, \frac{\sigma_0^2}{2}), \text{ (high-SNR)} \\ \frac{\|\boldsymbol{\mu} - \boldsymbol{\mu}_0\|^2 + d\sigma_0^2/2}{d(2\sigma_Z^2/n - \sigma_0^2)\sigma_0^2}, & \text{if } \frac{\sigma_Z^2}{n} \in [\frac{\sigma_0^2}{2}, \infty). \text{ (low-SNR)} \end{cases}$$

- ▶ $\frac{\sigma_Z^2}{n}$ only depends on S , can be interpreted as normalized noise
- ▶ σ_0^2 and $\|\boldsymbol{\mu} - \boldsymbol{\mu}_0\|^2$ captures the confidence and bias of prior knowledge
- ▶ high-SNR regime, high-quality training samples, discarding prior distribution and employing standard ERM
- ▶ low-SNR regime, where we should incorporate knowledge from both training samples and prior, optimal γ depends on everything
- ▶ If $\boldsymbol{\mu}_0 = \boldsymbol{\mu}$ and $\sigma_0^2 = 0$, $\gamma^* = 0$

Example: Linear Regression

- ▶ Training data $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$, with $\mathcal{X} = \mathbb{R}^d$ and $\mathcal{Y} = \mathbb{R}$
- ▶ Data is generated using true weights $W^* \in \mathbb{R}^d$ with additive noise,

$$Y_i = X_i \cdot W^* + \varepsilon_i, \quad \varepsilon \sim \mathcal{N}(0, \sigma_\varepsilon^2).$$

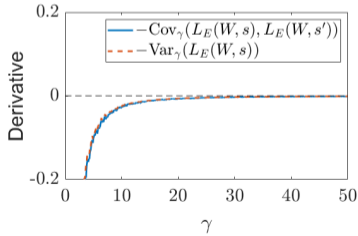
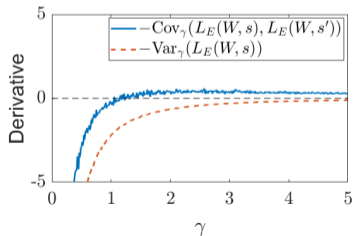
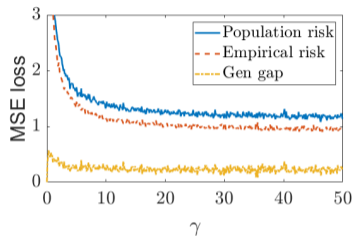
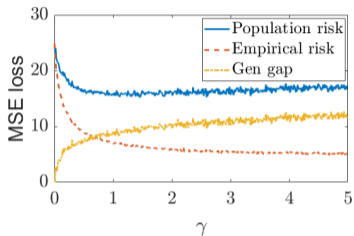
- ▶ Mean-squared loss $\ell(\mathbf{w}, \mathbf{z}) = (y - \mathbf{x} \cdot \mathbf{w})^2$
- ▶ Gaussian prior $\pi(\mathbf{w}) = \mathcal{N}(0, \sigma_0^2 \mathbf{I}_d)$
- ▶ $(\gamma, \mathcal{N}(0, \sigma_0^2 \mathbf{I}_d), L_E(\mathbf{w}, s))$ -Gibbs algorithm is Gaussian

$$P_{W|S}^\gamma(\mathbf{w}|S) \sim \mathcal{N}\left(\Sigma^{-1} \mathbf{X}^\top \mathbf{Y}, \frac{n}{2\gamma} \Sigma^{-1}\right),$$

with $\Sigma \triangleq \mathbf{X}^\top \mathbf{X} + \frac{n}{2\sigma_0^2\gamma} \mathbf{I}_d$, and $\mathbf{X} \in \mathbb{R}^{n \times d}$, $\mathbf{Y} \in \mathbb{R}^n$ are the matrix form of the training data.

Simulation of Linear Regression

Low SNR regime, $n = 10$ and $\sigma_\varepsilon^2 = 3$; high SNR regime, $n = 100$ and $\sigma_\varepsilon^2 = 1$.



Low-SNR

High-SNR

Table of Contents (Part II)

Generalization Error of Gibbs Algorithm

Gibbs algorithm

Exact Characterizations

Other Properties

Optimal Inverse Temperature

Asymptotic Behavior and Information Criteria for Model Selection

Generalization Error of Transfer learning

Conclusion

Asymptotic Behavior of Generalization Error

- ▶ Can we show something for ERM by letting $\gamma \rightarrow \infty$?

- ▶ Previous upper bound has order $\mathcal{O}(\frac{\gamma}{n})$

- ▶ Asymptotic normality of Gibbs algorithm

- ▶ **Single-well** case: there exists a unique $W^*(S)$

$$W^*(S) = \arg \min_{w \in \mathcal{W}} L_E(w, S).$$

- ▶ If $H^*(S) \triangleq \nabla_w^2 L_E(w, S)|_{w=W^*(S)}$ is invertible [Hwang, 1980],

$$P_{W|S}^\gamma \rightarrow \mathcal{N}(W^*(S), \frac{1}{\gamma} H^*(S)^{-1})$$

G. Aminian*, Y. Bu*, L. Toni, M. R. Rodrigues, G. W. Wornell. "An Exact Characterization of the Generalization Error for the Gibbs Algorithm," in *Proc. Conference on Neural Information Processing Systems (NeurIPS)*, Dec. 2021.

Asymptotic Behavior of MLE

Maximum likelihood estimates (MLE) in the asymptotic regime $n \rightarrow \infty$.

- ▶ n i.i.d. training samples generated from distribution P_Z
- ▶ Fit training data with distribution family $f(z_i|\theta)$, $\theta \in \mathbb{R}^p$
- ▶ $P_Z = f(\cdot|\theta^*)$ for $\theta^* \in \mathcal{W}$
- ▶ log-loss $\ell(\mathbf{w}, z) = -\log f(z|\mathbf{w})$

As $n \rightarrow \infty$, Gibbs algorithm converges to ERM algorithm (MLE),

$$\hat{W}_{\text{ML}} \triangleq \arg \max_{\theta \in \mathcal{W}} \sum_{i=1}^n \log f(Z_i|\theta).$$

Compute $I_{\text{SKL}}(W; S)$ using Gaussian approximation

$$\overline{\text{gen}}(P_{W|S}^\infty, P_S) = \frac{d}{n}.$$

Connection to Model Selection

- ▶ K candidate models M_1, M_2, \dots, M_K
- ▶ Each model M_k is characterized by parametric probabilistic model $P_k(\mathbf{z}|\boldsymbol{\theta}_k)$ and prior $\pi_k(\boldsymbol{\theta}_k)$
- ▶ log likelihood as the loss function $\ell_{\log}(w, \mathbf{z}) \triangleq -\log P(\mathbf{z}|w)$

How to select the *optimal* model?

- ▶ Information Criteria for Model Selection
 - ▶ Akaike Information Criterion (AIC)
 - ▶ Bayesian Information Criterion (BIC)

Akaike Information Criterion (AIC)

AIC selects the model that minimizes **population risk**:

$$\arg \min_k \text{KL}(P_Z \| P_k(\mathbf{z} | \hat{\theta}_{\text{ML}}^{(k)})) = \arg \min_k \mathbb{E}_{P_Z} [-\log P_k(Z | \hat{\theta}_{\text{ML}}^{(k)})].$$

AIC approximates it using **empirical risk** and **generalization error**

$$\text{AIC} = \arg \min_k L_E(\hat{\theta}_{\text{ML}}^{(k)}, S) + \overline{\text{gen}}(\hat{\theta}_{\text{ML}}^{(k)}, P_Z).$$

In classic regime where $n \rightarrow \infty$, and certain regularization conditions

$$\text{AIC} = \arg \min_k L_E(\hat{\theta}_{\text{ML}}^{(k)}, S) + \frac{p}{n}.$$

Bayesian Information Criterion (BIC)

BIC selects the model that maximizes **marginal likelihood**:

$$m_k(\mathbf{z}^n) \triangleq \int P_k(\mathbf{z}^n | \boldsymbol{\theta}_k) \pi_k(\boldsymbol{\theta}_k) d\boldsymbol{\theta}_k,$$

which is equivalent to maximizing posterior probability $P(M_k | \mathbf{z}^n)$.

$$\begin{aligned} \text{BIC} &= \arg \min_k -\frac{1}{n} \log m_k(\mathbf{z}^n) \\ &= \arg \min_k L_E(\hat{\boldsymbol{\theta}}_{\text{ML}}^{(k)}, \mathcal{S}) + \frac{p_k \log n}{2n}, \end{aligned}$$

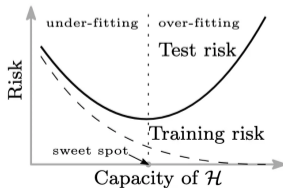
where Laplace approximation is applied as $n \rightarrow \infty$.

Comparison between AIC and BIC

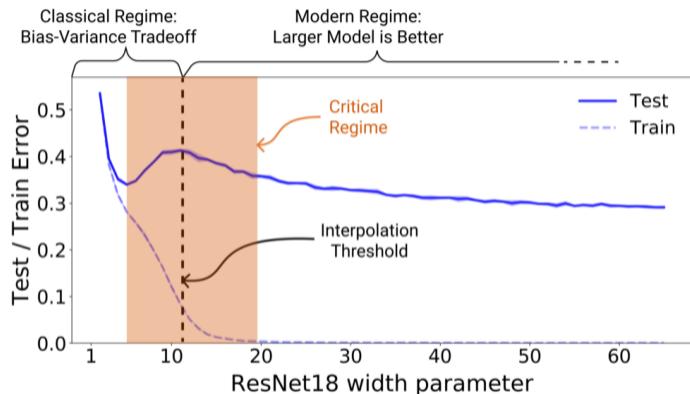
$$\text{AIC} = \arg \min L_E(\hat{\theta}_{\text{ML}}, S) + \frac{p}{n}$$

$$\text{BIC} = \arg \min L_E(\hat{\theta}_{\text{ML}}, S) + \frac{p \log n}{2n}.$$

- ▶ AIC minimizes **population risk** (optimal prediction performance)
- ▶ BIC maximizes the **marginal likelihood** (identifying the true model)
- ▶ BIC imposing a **larger** penalty for more complex models.



Double-descent in Over-parameterized Regime



- ▶ When $p \leq n$, the classical U-shaped curve is valid.
- ▶ When $p \geq n$, test loss can decrease again.

Challenges in Over-parameterized regime

Asymptotic normality (AIC) and Laplace Approximation (BIC) do not hold in this new regime!

There are some efforts to extend these information criteria:

- ▶ Akaike's Information Corrected Criterion (AICC), fixed p , small n
- ▶ Widely applicable BIC (WBIC), singular Hessian matrix

More recent work trying to demystify double-descent

- ▶ Neural Tangent Kernel (NTK), lazy training
- ▶ Random feature model
- ▶ Mean-field approach

Marginal likelihood of Gibbs algorithm

Recall the **information risk minimization** for motivating the Gibbs algorithm.

$$\min_{P_{W|S}} \mathbb{E}_{P_{W|S=s}} [L_E(W, s)] + \frac{1}{\gamma} \text{KL}(P_{W|S=s} \| \pi) = -\frac{1}{\gamma} \log V(s, \gamma).$$

If we adopt log-loss function $\ell(w, z) = -\log P(z|w)$, and set $\gamma = n$

$$\begin{aligned} -\frac{1}{\gamma} \log V(s, \gamma) &= -\frac{1}{n} \log \int \pi(w) e^{-nL_E(w, s)} dw \\ &= -\frac{1}{n} \log \int \pi(w) P(z^n|w) dw \\ &= -\frac{1}{n} \log m(z^n) \end{aligned}$$

Gibbs based Information Criteria

Gibbs-based AIC:

$$\text{AIC}^+ \triangleq L_E(\hat{W}_{\text{Gibbs}}, \mathbf{z}^n) + \frac{1}{n} \text{I}_{\text{SKL}}(P_{\hat{W}|S}^*, P_S).$$

Gibbs-based BIC:

$$\text{BIC}^+ \triangleq L_E(\hat{W}_{\text{Gibbs}}, \mathbf{z}^n) + \frac{1}{n} \text{KL}(P_{\hat{W}|S=\mathbf{z}^n}^* \parallel \pi),$$

$$\text{BIC}^- \triangleq \mathbb{E}_\pi [L_E(W, \mathbf{z}^n)] - \frac{1}{n} \text{KL}(\pi \parallel P_{\hat{W}|S=\mathbf{z}^n}^*).$$

We can show that in the classic regime where p is fixed and $n \rightarrow \infty$, they all reduce back to their classical forms.

H. Chen, Y. Bu, G. W. Wornell, "Gibbs-Based Information Criteria and the Over-Parameterized Regime," in *Proc. International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2024.

Random Feature Model

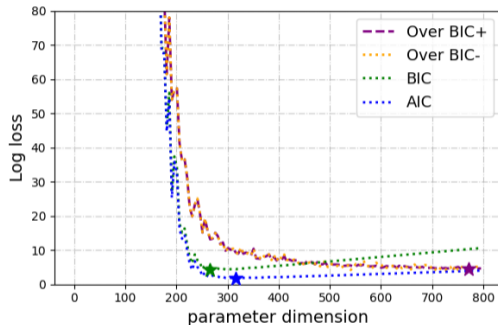
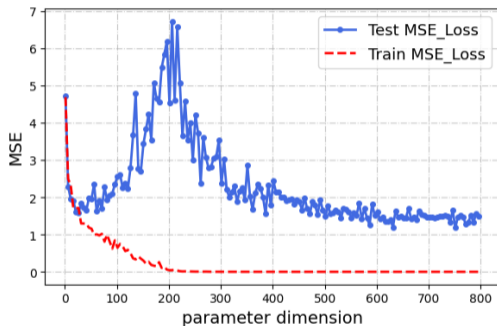
The output of **Random Feature (RF) model** with input data $\mathbf{x} \in \mathbb{R}^d$ is

$$g(\mathbf{x}) \triangleq \sum_{j=1}^p f\left(\frac{\langle \mathbf{x}, \mathbf{F}_j \rangle}{\sqrt{d}}\right) \mathbf{w}_j = f\left(\frac{\mathbf{x}^\top \mathbf{F}}{\sqrt{d}}\right) \mathbf{w},$$

- ▶ Two-layer neural network with i.i.d Gaussian weights $\mathbf{F} \in \mathbb{R}^{d \times p}$ in the first layer, only the second layer is trainable
- ▶ $f()$ is the non-linear activation function
- ▶ The dimensionality of input data d is not entangled with number of parameters p

Experiment

Evaluating the BIC^+ and BIC^- using $n = 200$ samples in RF models



- ▶ We observe **Double-descent** in population risk for RF model
- ▶ Our Gibbs-based BICs prefer **over-parameterized models**

Check point

- ▶ Provide information criteria for the Gibbs algorithm, with different **information measures** as the penalty terms.
- ▶ Generalize our information-theoretic analysis to **over-parameterized** random feature.
- ▶ The **mismatch** between **marginal likelihood (BIC)** and **generalization error (AIC)** in the over-parameterized setting, which highly depends on the prior distributions.

Table of Contents (Part II)

Generalization Error of Gibbs Algorithm

Gibbs algorithm

Exact Characterizations

Other Properties

Optimal Inverse Temperature

Asymptotic Behavior and Information Criteria for Model Selection

Generalization Error of Transfer learning

Conclusion

Generalization of Transfer Learning

- ▶ Source data set $D^s = \{Z_i^s\}_{i=1}^m$, generated from P_{D^s}
- ▶ Target data set $D^t = \{Z_j^t\}_{j=1}^n$, generated from P_{D^t}
- ▶ The empirical risk of source and target task

$$L_E(w, d^s) \triangleq \frac{1}{m} \sum_{j=1}^m \ell(w, z_j^s), \quad L_E(w, d^t) \triangleq \frac{1}{n} \sum_{j=1}^n \ell(w, z_j^t).$$

- ▶ The population risk of the target task

$$L_P(w, P_{D^t}) \triangleq \mathbb{E}_{P_{D^t}} [L_E(w, D^t)].$$

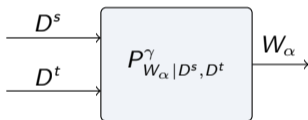
- ▶ Expected Transfer Generalization Error

$$\overline{\text{gen}}(P_{W|D^s, D^t}, P_{D^s}, P_{D^t}) \triangleq \mathbb{E}_{P_{W, D^s, D^t}} [L_P(W, P_{D^t}) - L_E(W, D^t)].$$

Transfer Learning: α -Weighted ERM

- ▶ Output hypothesis w_α is trained by minimizing a convex combination of the source and target task empirical risks [Ben-David et al., 2010], for $\alpha \in [0, 1]$

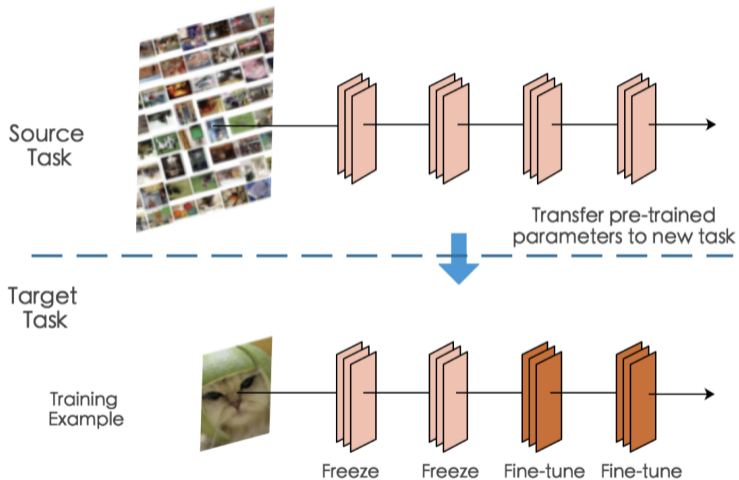
$$L_E(w_\alpha, d^s, d^t) = (1 - \alpha)L_E(w_\alpha, d^s) + \alpha L_E(w_\alpha, d^t)$$



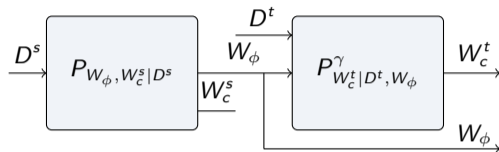
- ▶ **α -weighted Gibbs algorithm** generalizes the α -weighted-ERM by considering the $(\gamma, \pi(w_\alpha), L_E(w_\alpha, d^s, d^t))$ -Gibbs algorithm

$$P_{W_\alpha | D^s, D^t}^\gamma(w_\alpha | d^s, d^t) = \frac{\pi(w_\alpha) e^{-\gamma L_E(w_\alpha, d^s, d^t)}}{V_\alpha(d^s, d^t, \gamma)}.$$

Transfer Learning: Two-stage ERM



Two-stage-ERM Transfer Learning



- **First Stage:** Learn shared feature extractor $w_\phi \in \mathcal{W}_\phi$

$$[W_\phi, W_c^s] = \arg \min_w L_E^{S1}(w, d^s).$$

- **Second Stage:** Freeze W_ϕ , and learn target-specific hypothesis w_c^t

$$W_c^t = \arg \min_{w_c} L_E^{S2}([W_\phi, w_c], d^t)$$

Expected Transfer Generalization Error

Theorem

The expected transfer generalization error of the α -weighted Gibbs algorithm is given by

$$\overline{\text{gen}}_{\alpha}(P_{D^s}, P_{D^t}) = \frac{I_{\text{SKL}}(W_{\alpha}; D^t | D^s)}{\alpha\gamma}.$$

Theorem

The expected transfer generalization error of the two-stage Gibbs algorithm is given by

$$\overline{\text{gen}}_{\beta}(P_{D^s}, P_{D^t}) = \frac{I_{\text{SKL}}(W_c^t; D^t | W_{\phi})}{\gamma}.$$

Y. Bu*, G. Aminian*, L. Toni, M. R. Rodrigues, G. W. Wornell. "Characterizing and Understanding the Generalization Error of Transfer Learning with Gibbs Algorithm," in *Proc. International Conference on Artificial Intelligence and Statistics (AISTATS)* 2022.

Asymptotic Behavior of MLE

Maximum likelihood estimates

- ▶ n i.i.d. target samples, m i.i.d. source samples
- ▶ Fit training data with distribution family $f(z|\mathbf{w})$, $\mathbf{w} = (\mathbf{w}_\phi, \mathbf{w}_c) \in \mathbb{R}^d$, $\mathbf{w}_c \in \mathbb{R}^{d_c}$
- ▶ $P_{Z^t} = f(\cdot|\mathbf{w}^*)$ for $\mathbf{w}^* \in \mathcal{W}$
- ▶ log-loss $\ell(\mathbf{w}, z) = -\log f(z|\mathbf{w})$

	Standard target ERM	α -weighted ERM	Two-stage ERM
$\overline{\text{gen}}$	$\mathcal{O}(\frac{d}{n})$	$\mathcal{O}(\frac{d}{m+n})$	$\mathcal{O}(\frac{d_c}{n})$

Table of Contents (Part II)

Generalization Error of Gibbs Algorithm

Gibbs algorithm

Exact Characterizations

Other Properties

Optimal Inverse Temperature

Asymptotic Behavior and Information Criteria for Model Selection

Generalization Error of Transfer learning




Conclusion

Conclusion

- ▶ Connect **operational quantity** in learning theory (generalization error, marginal likelihood) with different **information measures** for Gibbs algorithm
- ▶ Demonstrate the **versatility** of our approach in multiple applications
 - ▶ Optimal Inverse temperature
 - ▶ Gibbs-based BIC for over-parameterized model selection
 - ▶ Gibbs based-**transfer** learning
- ▶ Our Gibbs-based analysis provides an information-theoretic **framework** for understanding **generalization** behavior in modern machine learning, still a lot to be explored!

Thank you for your attention!

References I

-  Aminian, G., Arjmandi, H., Gohari, A., Nasiri-Kenari, M., and Mitra, U. (2015).
Capacity of diffusion-based molecular communication networks over lti-poisson channels.
IEEE Transactions on Molecular, Biological and Multi-Scale Communications, 1(2):188–201.
-  Aminian, G., Bu, Y., Toni, L., Rodrigues, M., and Wornell, G. (2021).
An exact characterization of the generalization error for the gibbs algorithm.
Advances in Neural Information Processing Systems, 34:8106–8118.
-  Ben-David, S., Blitzer, J., Crammer, K., Kulesza, A., Pereira, F., and Vaughan, J. W. (2010).
A theory of learning from different domains.
Machine learning, 79(1):151–175.

References II



Bu, Y., Zou, S., and Veeravalli, V. V. (2020).

Tightening mutual information-based bounds on generalization error.

IEEE Journal on Selected Areas in Information Theory, 1(1):121–130.



Harutyunyan, H., Raginsky, M., Ver Steeg, G., and Galstyan, A. (2021).

Information-theoretic generalization bounds for black-box learning algorithms.

Advances in Neural Information Processing Systems, 34:24670–24682.






Hwang, C.-R. (1980).

Laplace's method revisited: weak convergence of probability measures.

The Annals of Probability, pages 1177–1182.

References III

-  Palomar, D. P. and Verdú, S. (2008).
Lautum information.
IEEE transactions on information theory, 54(3):964–975.
-  Raginsky, M., Rakhlin, A., and Telgarsky, M. (2017).
Non-convex learning via stochastic gradient langevin dynamics: a nonasymptotic analysis.
In *Conference on Learning Theory*, pages 1674–1703. PMLR.
-  Steinke, T. and Zakyntinou, L. (2020).
Reasoning about generalization via conditional mutual information.
arXiv preprint arXiv:2001.09122.

References IV



Wang, Z. and Mao, Y. (2023).

Tighter information-theoretic generalization bounds from supersamples.

In *International Conference on Machine Learning*, pages 36111–36137. PMLR.



Xu, A. and Raginsky, M. (2017).

Information-theoretic analysis of generalization capability of learning algorithms.

In *Advances in Neural Information Processing Systems*, pages 2524–2533.



Zhang, T. (2006).

Information-theoretic upper and lower bounds for statistical estimation.

IEEE Transactions on Information Theory, 52(4):1307–1321.