

Information-Theoretic Analysis for Generalization of Learning Algorithms

A Short Tutorial



uOttawa

Ziqiao Wang

University of Ottawa

School of Electrical Engineering and Computer Science

December 22, 2023

Preliminaries on Information Theory

Background on Information-Theoretic Generalization Bounds

Information-Theoretic Generalization Bounds for Black-Box Algorithms

Information-Theoretic Bounds in Stochastic Convex Optimization

Information-Theoretic Generalization Bounds for SGD

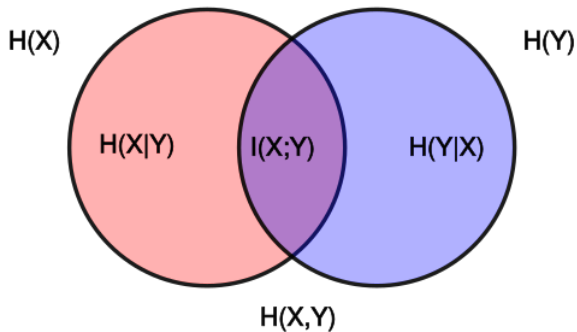
Information-Theoretic Analysis Beyond Supervised Learning

References

Preliminaries on Information Theory

- ▶ Entropy: $H(X) = \mathbb{E}_{P_X} \left[\log \frac{1}{P(X)} \right]$, $H(X, Y) = \mathbb{E}_{P_{X,Y}} \left[\log \frac{1}{P(X,Y)} \right]$,
 $H(X|Y) = \mathbb{E}_{P_{X,Y}} \left[\log \frac{1}{P(X|Y)} \right]$
 - ▶ For discrete X , $H(X) \geq 0$
 - ▶ $H(X, Y) = H(X|Y) + H(Y)$
 - ▶ Conditioning reduces entropy: $H(X|Y) \leq H(X)$
 - ▶ For discrete X , $H(X) \leq \log |\mathcal{X}|$
- ▶ Relative Entropy: $D_{\text{KL}}(Q||P) = \mathbb{E}_Q \left[\log \frac{Q(X)}{P(X)} \right]$
 - ▶ $D_{\text{KL}}(Q||P) \geq 0$ with equality holds iff $Q = P$.
 - ▶ Usually $D_{\text{KL}}(Q||P) \neq D_{\text{KL}}(P||Q)$

- ▶ Mutual Information: $I(X; Y) = \mathbb{E}_{P_{X,Y}} \left[\log \frac{P(X,Y)}{P(X)P(Y)} \right] = D_{\text{KL}}(P_{X,Y} || P_X P_Y)$.
 - ▶ $I(X; Y) \geq 0$ with equality holds iff $X \perp\!\!\!\perp Y$.
 - ▶ $I(X; Y) = H(X) - H(X|Y) = H(Y) - H(Y|X) = H(X) + H(Y) - H(X, Y)$.
 - ▶ $I(X; Y) = I(Y; X)$
 - ▶ $I(X; Y) = \mathbb{E}_{P_{X,Y}} \left[\log \frac{P(X|Y)}{P(X)} \right] = \mathbb{E}_{P_Y} [D_{\text{KL}}(P_{X|Y} || P_X)]$
- ▶ Conditional Mutual Information and Disintegrated Mutual Information:
$$I(X; Y|Z) = \mathbb{E}_{P_{X,Y,Z}} \left[\log \frac{P(X,Y|Z)}{P(X|Z)P(Y|Z)} \right] = H(X|Z) - H(X|Y, Z)$$
$$I^Z(X; Y) = \mathbb{E}_{P_{X,Y|Z=z}} \left[\log \frac{P(X,Y|Z=z)}{P(X|Z=z)P(Y|Z=z)} \right]$$
 - ▶ $\mathbb{E}_Z [I^Z(X; Y)] = I(X; Y|Z)$



Venn diagram. Credit: https://en.wikipedia.org/wiki/Mutual_information

► Chain-rule:

- $H(X_1, X_2, \dots, X_n) = \sum_{i=1}^n H(X_i | X_{i-1}, \dots, X_1)$
- $I(X_1, X_2, \dots, X_n; Y) = \sum_{i=1}^n I(X_i; Y | X_{i-1}, \dots, X_1)$
- $D_{\text{KL}}(Q_{X,Y} || P_{X,Y}) = D_{\text{KL}}(Q_X || P_X) + D_{\text{KL}}(Q_{Y|X} || P_{Y|X})$

► Data-processing inequality (DPI):

If $X - Y - Z$ forms a Markov chain (i.e. $P_{X,Z|Y} = P_{X|Y}P_{Z|Y}$), then

$$I(X; Y) \geq I(X; Z)$$

e.g., $(X, Y) - f(X, Y) - Z$ is a Markov chain : $I(X, Y; Z) \leq I(f(X, Y); Z)$

- Other useful stuff: Fano's inequality, KL divergence between two Gaussian, Gaussian distribution maximizes the entropy over all distributions with the same variance, ...
- Textbook for beginners: *Thomas M. Cover and Joy A. Thomas. Elements of Information Theory, Wiley-Interscience, 2006.*

Lemma 1 (Variational Representation of Mutual Information)

For two random variables X and Y , we have

$$I(X; Y) = \inf_Q \mathbb{E}_{P_Y} [\mathrm{D}_{\mathrm{KL}}(P_{X|Y} || Q)],$$

where the infimum is achieved at $Q = P_X$.

Note that $I(X; Y) = \mathbb{E}_{P_Y} [\mathrm{D}_{\mathrm{KL}}(P_{X|Y} || P_X)]$

Lemma 2 (Donsker and Varadhan's variational formula)

For any measurable function $f : \Theta \rightarrow \mathbb{R}$, we have

$$D_{\text{KL}}(Q||P) = \sup_f \mathbb{E}_{\theta \sim Q} [f(\theta)] - \log \mathbb{E}_{\theta \sim P} [\exp f(\theta)].$$

proof. Define the density of the Gibbs measure P_f : $P_f(\theta) = \frac{e^{f(\theta)}}{\mathbb{E}_{\theta \sim P}[e^{f(\theta)}]} P(\theta)$.

$$\begin{aligned} D_{\text{KL}}(Q||P_f) &= \mathbb{E}_Q \left[\log \frac{Q}{P_f} \right] = \mathbb{E}_Q [\log Q] - \mathbb{E}_Q \left[\log \frac{e^{f(\theta)}}{\mathbb{E}_P [e^{f(\theta)}]} P \right] \\ &= \mathbb{E}_Q [\log Q] - \mathbb{E}_Q [f(\theta)] - \mathbb{E}_Q [\log P] + \log \mathbb{E}_P [e^{f(\theta)}] \\ &= D_{\text{KL}}(Q||P) - \mathbb{E}_{\theta \sim Q} [f(\theta)] - \log \mathbb{E}_{\theta \sim P} [\exp f(\theta)] \\ &\geq 0 \end{aligned}$$

- *Polyanskiy, Y. and Wu, Y.. Information Theory: From Coding to Learning, Cambridge University Press, 2023 (book draft).*

Background on Information-Theoretic Generalization Bounds

- ▶ A learning algorithm $\mathcal{A} : S \rightarrow W$ i.e. mapping training sample S to a hypothesis W .
- ▶ Gen. err. = $\mathbb{E} [\text{Test err.} - \text{Train err.}] \leq \text{Gen. bound.}$

- ▶ A learning algorithm $\mathcal{A} : S \rightarrow W$ i.e. mapping training sample S to a hypothesis W .
- ▶ Gen. err. = $\mathbb{E} [\text{Test err.} - \text{Train err.}] \leq \text{Gen. bound.}$

Formal Notations:

- ▶ Training dataset: $S = \{Z_i\}_{i=1}^n \in \mathcal{Z}$, drawn i.i.d. from μ
- ▶ Hypothesis space: $\mathcal{W} \subseteq \mathbb{R}^d$
- ▶ Learning algorithm: $\mathcal{A} : \mathcal{Z}^n \rightarrow \mathcal{W}$ by $P_{W|S}$
- ▶ Loss: $\ell : \mathcal{W} \times \mathcal{Z} \rightarrow \mathbb{R}^+$
- ▶ We're interested in
 - ▶ Population risk: $L_\mu(w) \triangleq \mathbb{E}_{Z \sim \mu}[\ell(w, Z)]$
 - ▶ Empirical risk: $L_S(w) \triangleq \frac{1}{n} \sum_{i=1}^n \ell(w, Z_i)$
 - ▶ Expected generalization error: $\mathcal{E}_\mu(\mathcal{A}) \triangleq \mathbb{E}_{W, S}[L_\mu(W) - L_S(W)]$




Before Xu’s bound:

- ▶ *Russo, D. and Zou, J.. Controlling bias in adaptive data analysis using information theory. AISTATS 2016.*
Russo, D., and Zou, J. How much does your data exploration overfit? Controlling bias via information usage. TIT 2019.
- ▶ *Raginsky, M. et al. Information-theoretic analysis of stability and bias of learning algorithms. ITW 2016.*

Lemma 3 (Xu and Raginsky [2017])

Assume the loss $\ell(w, Z)$ is R -subgaussian¹ for any $w \in \mathcal{W}$. The generalization error of \mathcal{A} is bounded by

$$|\mathcal{E}| \leq \sqrt{\frac{2R^2}{n} I(W; S)}.$$

¹A random variable X is R -subgaussian if for any ρ , $\log \mathbb{E} \exp(\rho(X - \mathbb{E}X)) \leq \rho^2 R^2 / 2$.  uOttawa

► **Step 1: Finding the target.**

$$\begin{aligned}\mathcal{E} = \mathbb{E}_{S,W} [L_\mu(W) - L_S(W)] &= \mathbb{E}_{S,W} [\mathbb{E}_{S'} [L_{S'}(W)]] - \mathbb{E}_{S,W} [L_S(W)] \\ &= \mathbb{E}_{P_W \otimes P_{S'}} [L_{S'}(W)] - \mathbb{E}_{P_{W,S}} [L_S(W)]\end{aligned}$$

► **Step 2: Selecting the measurable function f .**

Recall DV Lemma:

$$\begin{aligned}I(W, S) &= D_{\text{KL}}(P_{W,S} \| P_W \otimes P_{S'}) \\ &\geq \sup_f \mathbb{E}_{(W,S) \sim P_{W,S}} [f(W, S)] - \log \mathbb{E}_{(W,S') \sim P_W \otimes P_{S'}} [\exp f(W, S')]\end{aligned}$$

Let $f(W, S) = tL_S(W)$ for some $t > 0$.



► **Step 3: Bounding the CGF.**

If $\ell(w, Z)$ is R -SubGaussian, $f(w, S') = L_{S'}(w)$ is R/\sqrt{n} -SubGaussian:

$$\log \mathbb{E}_{W, S'} [\exp \lambda (L_{S'} - \mathbb{E} [L_{S'}])] \leq t^2 R^2 / 2n$$

Thus, $\log \mathbb{E}_{W, S'} [\exp t L_{S'}(W)] \leq t \mathbb{E}_{W, S'} [L_{S'}(W)] + t^2 R^2 / 2n$.

► **Step 4: Optimizing the bound.**

$$I(W, S) \geq \sup_{t > 0} t \left(\mathbb{E}_{(W, S) \sim P_{W, S}} [L_S(W)] - \mathbb{E}_{(W, S') \sim P_W \otimes P_{S'}} [L_{S'}(W)] \right) - t^2 R^2 / 2n$$

$$\implies \mathbb{E}_{(W, S) \sim P_{W, S}} [L_S(W)] - \mathbb{E}_{(W, S') \sim P_W \otimes P_{S'}} [L_{S'}(W)] \leq \inf_t \frac{I(W, S)}{t} + \frac{t R^2}{2n} =$$

$$\sqrt{\frac{2R^2}{n} I(W, S)}$$

$$\implies |\mathcal{E}| \leq \sqrt{\frac{2R^2}{n} I(W, S)}$$

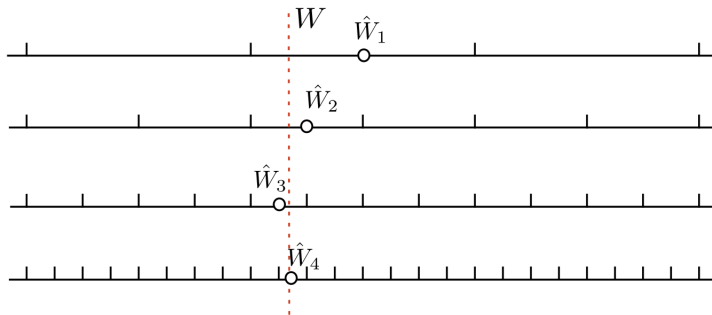
$I(W; S) \rightarrow \infty$ e.g., \mathcal{A} is deterministic $\implies I(W; S) = H(W) - H(W|S) = H(W)$.

Some previous efforts:

- ▶ **Chaining Method:** $3\sqrt{2} \sum_{k=k_1}^{\infty} 2^{-k} \sqrt{I([W]_k; Z_i)}$
Asadi, A. et al. Chaining mutual information and tightening generalization bounds. NeurIPS 2018.
- ▶ **Individual Technique/Sample-Wise Bound:** $\frac{1}{n} \sum_{i=1}^n \sqrt{I(W; Z_i)}$
Bu, Y. et al. Tightening Mutual Information Based Bounds on Generalization Error. ISIT 2019.
- ▶ **Random Subset Technique:** $\mathbb{E} \sqrt{\frac{1}{n-m} I^{S_J}(W; S_J^c)}$
Negrea, J. et al. Information-theoretic generalization bounds for SGLD via data-dependent estimates. NeurIPS 2019.
- ▶ **Solved by CMI:** $\sqrt{\frac{1}{n} I(W; U|\tilde{Z})} \leq \mathcal{O}(1)$
Steinke, T. and Zakyntinou, L.. Reasoning about generalization via conditional mutual information. COLT 2020.



Idea:



Quantization of W . Credit: Zhou R, et al. Stochastic Chaining and Strengthened Information-Theoretic Generalization Bounds ISIT 2022.

- **Step 1: Finding the target.** For any integers k_1 and k_0 such that $k_1 > k_0$, let $\mathcal{E}(W) = L_\mu(W) - L_S(W)$, we have

$$\mathcal{E}(W) = \mathcal{E}([W]_{k_0}) + \sum_{k=k_0+1}^{k_1} (\mathcal{E}([W]_k) - \mathcal{E}([W]_{k-1})) + \mathcal{E}(W) - \mathcal{E}([W]_{k_1}).$$

We require $\mathbb{E}[\mathcal{E}([W]_{k_0})] = 0$ and $\lim_{k_1 \rightarrow \infty} \mathcal{E}([W]_{k_1}) = \mathcal{E}(W)$.

Let $k_1 \rightarrow \infty$ and taking expectation over $(S, W) \sim P_{S,W}$ for both sides of the equation above, we have

$$\mathcal{E} = \sum_{k=k_0+1}^{\infty} \mathbb{E}_{S,[W]_k,[W]_{k-1}} [(\mathcal{E}([W]_k) - \mathcal{E}([W]_{k-1}))]. \quad (1)$$

- **Step 2: Selecting f , Q and P .**

$$f = t \cdot (\mathcal{E}([W]_k) - \mathcal{E}([W]_{k-1})), \quad Q = P_{S,[W]_k,[W]_{k-1}}, \quad P = P_S \otimes P_{[W]_k,[W]_{k-1}}$$

► **Step 3: Bounding the CGF.**

$\mathcal{E}([W]_k) - \mathcal{E}([W]_{k-1})$ is $d^2([W]_k, [W]_{k-1})$ -subGaussian:

$$\text{CGF} = \log \mathbb{E}_{S'} \left[\mathbb{E}_{[W]_k, [W]_{k-1}} \left[e^{t(\mathcal{E}([W]_k) - \mathcal{E}([W]_{k-1}))} \right] \right] \leq \frac{t^2 \mathbb{E} [d^2([W]_k, [W]_{k-1})]}{2}$$

► **Step 4: Optimizing the bound.**

$$\mathcal{E} \leq \sum_{k=k_0+1}^{\infty} \sqrt{2 \mathbb{E}_{[W]_k, [W]_{k-1}} [d^2([W]_k, [W]_{k-1}) I([W]_k, [W]_{k-1}; S)]}.$$

Notice that $S - W - [W]_k - [W]_{k-1}$ is a Markov chain, so

$$I([W]_k, [W]_{k-1}; S) = I([W]_k; S) + I([W]_k, [W]_{k-1}; S | [W]_k) = I([W]_k, [W]_{k-1}; S).$$

Special case: 2^{-k} -partition, $d([W]_k, W) \leq 2^{-k}$, then
 $d([W]_k, W) + d([W]_{k-1}, W) \leq 2^{-k} + 2^{-(k-1)} = 3 \times 2^{-k}.$



- **Step 1: Finding the target.**

$$\mathbb{E}_{W,S}[L_{\mu}(W) - L_S(W)] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{W,Z_i} [\mathbb{E}_{Z'} [\ell(W, Z')] - \ell(W, Z_i)]$$

- **Step 2: Selecting f , Q and P .**

$$f = t \cdot (\mathbb{E}_{Z'} [\ell(W, Z')] - \ell(W, Z_i)), \quad Q = P_{W,Z_i}, \quad P = P_{Z'} \otimes P_W$$

- **Step 3: Bounding the CGF.**

$$\ell(W, Z'_i) \text{ is } R\text{-subGaussian: } \log \mathbb{E}_{Z'} \left[\mathbb{E}_W \left[e^{t(\mathbb{E}_{Z'} [\ell(W, Z')] - \ell(W, Z_i))} \right] \right] \leq \frac{t^2 R^2}{2}$$

- **Step 4: Optimizing the bound.**

$$\mathcal{E} \preceq \frac{1}{n} \sum_{i=1}^n \sqrt{I(W; Z_i)} \leq \sqrt{\frac{I(W; S)}{n}}$$



► **Step 1: Finding the target.**

Let $J \subseteq [n]$, $|J| = m$,

$$\begin{aligned}\mathbb{E}_{W,S}[L_\mu(W) - L_S(W)] &= \mathbb{E}_{W,S} \left[\frac{1}{n} \sum_{i=1}^n (\mathbb{E}_{Z'} [\ell(W, Z')] - \ell(W, Z_i)) \right] \\ &= \mathbb{E}_J \left[\mathbb{E}_{W,S_J} \left[\frac{1}{m} \sum_{i=1}^m (\mathbb{E}_{Z'} [\ell(W, Z')] - \ell(W, S_{J_i}^c)) \right] \right]\end{aligned}$$

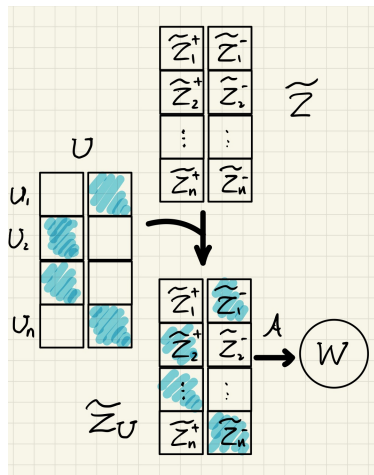
► **Step 2: Selecting f , Q and P .**

$$f = t \cdot \left(\frac{1}{m} \sum_{i=1}^m (\mathbb{E}_{Z'} [\ell(W, Z')] - \ell(W, S_{J_i}^c)) \right), \quad Q = P_{W, S_J^c | S, J}, \quad P = P_{S_J^c} \otimes P_{W' | S_J}$$

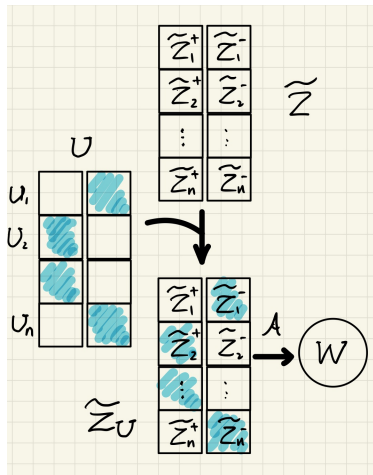
\implies Data-Dependent Prior of W

- $\mathcal{E} \lesssim \mathbb{E} \sqrt{\frac{I^{S_J}(W; S_J^c)}{n-m}}$; Individual Technique is a special case for $m = n - 1$.





- ▶ Let \tilde{Z} drawn i.i.d. from μ
- ▶ Let $U = (U_1, U_2, \dots, U_n)^T \sim \text{Unif}(\{0, 1\}^n)$.
- ▶ Learning algorithm $\mathcal{A} : \mathcal{Z}^n \rightarrow \mathcal{W}$
- ▶ $\mathcal{E} = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{W, U_i, \tilde{Z}} \left[(-1)^{U_i} \left(\ell(W, \tilde{Z}_{i,1}) - \ell(W, \tilde{Z}_{i,0}) \right) \right]$



- ▶ Let \tilde{Z} drawn i.i.d. from μ
- ▶ Let $U = (U_1, U_2, \dots, U_n)^T \sim \text{Unif}(\{0, 1\}^n)$.
- ▶ Learning algorithm $\mathcal{A} : \mathcal{Z}^n \rightarrow \mathcal{W}$
- ▶ $\mathcal{E} = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{W, U_i, \tilde{Z}} \left[(-1)^{U_i} \left(\ell(W, \tilde{Z}_{i,1}) - \ell(W, \tilde{Z}_{i,0}) \right) \right]$

Lemma 4 (Steinke and Zakynthinou [2020])

Assume the loss is bounded between $[0, 1]$, we have

$$|\mathcal{E}| \leq \sqrt{\frac{2I(W; U | \tilde{Z})}{n}}.$$

- **Step 1: Finding the target.**

$$\mathcal{E} = \mathbb{E}_{W, U, \tilde{Z}} \left[\frac{1}{n} \sum_{i=1}^n (-1)^{U_i} \left(\ell(W, \tilde{Z}_i^-) - \ell(W, \tilde{Z}_i^+) \right) \right]$$

- **Step 2: Selecting f , Q and P .**

$$f = t \cdot \frac{1}{n} \sum_{i=1}^n (-1)^{U_i} \left(\ell(W, \tilde{z}_i^-) - \ell(W, \tilde{z}_i^+) \right), \quad Q = P_{W, U|\tilde{z}}, \quad P = P_{U'} \otimes P_{W|\tilde{z}}$$

- **Step 3: Bounding the CGF.**

$(-1)^{U_i} \left(\ell(w, \tilde{z}_i^-) - \ell(w, \tilde{z}_i^+) \right)$ is $|\ell(w, \tilde{z}_i^-) - \ell(w, \tilde{z}_i^+)|^2$ -subGaussian:

$$\log \mathbb{E}_{W|\tilde{z}} \left[\mathbb{E}_{U'} \left[e^{t \frac{1}{n} \sum_{i=1}^n (-1)^{U_i} \left(\ell(W, \tilde{z}_i^-) - \ell(W, \tilde{z}_i^+) \right)} \right] \right] \leq \frac{t^2}{2n}$$

- **Step 4: Optimizing the bound.**

$$\mathcal{E} \preceq \sqrt{\frac{I(W; U|\tilde{Z})}{n}} \leq \sqrt{\frac{I(W; S)}{n}}.$$

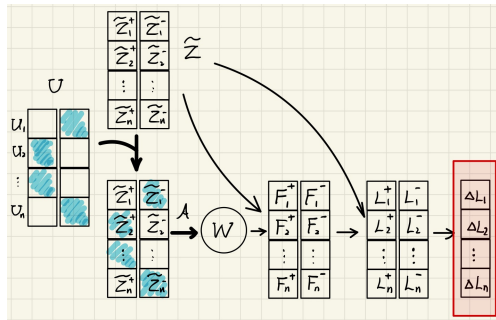


- ▶ **Random Subset CMI:** Haghifam, M. et al. *Sharpened generalization bounds based on conditional mutual information and an application to noisy, iterative algorithms. NeurIPS 2020.*
- ▶ **Individual CMI:** Rodríguez-Gálvez, B. et al. *On random subset generalization error bounds and the stochastic gradient Langevin dynamics algorithm. ITW 2020.* Zhou R, et al. *Individually conditional individual mutual information bound on generalization error. TIT 2022.*
- ▶ **Stochastic Chaining IOMI/CMI:** Zhou R, et al. *Stochastic Chaining and Strengthened Information-Theoretic Generalization Bounds ISIT 2022.*
- ▶ **Leave-One-Out CMI:** Haghifam, M. et al. *Understanding Generalization via Leave-One-Out Conditional Mutual Information. ISIT 2022.* Rammal, M. R. et al. *On leave-one-out conditional mutual information for generalization. NeurIPS 2022.*



Information-Theoretic Generalization Bounds for Black-Box Algorithms

- ▶ Wang, Z., and Mao, Y.. *Tighter Information-Theoretic Generalization Bounds from Supersamples. ICML 2023.*
- ▶ Main Contribution: New **Conditional Mutual Information (CMI)** bounds which are **either theoretically or empirically tighter** than previous CMI bounds for the **same supersample** setting.



$$\blacktriangleright F_i^+ := f_W(\tilde{X}_i^+), F_i^- := f_W(\tilde{X}_i^-), \\ F_i := (F_i^+, F_i^-)$$

$$\Rightarrow \textbf{f-CMI Bound: } |\mathcal{E}| \leq \frac{1}{n} \sum_{i=1}^n \sqrt{I(F_i; U_i | \tilde{Z})}$$

[Harutyunyan et al., 2021]

$$\blacktriangleright L_i^+ := \ell(W, \tilde{Z}_i^+), L_i^- := \ell(W, \tilde{Z}_i^-), \\ L_i := (L_i^+, L_i^-)$$

$$\Rightarrow \textbf{e-CMI Bound: } |\mathcal{E}| \leq \frac{1}{n} \sum_{i=1}^n \sqrt{I(L_i; U_i | \tilde{Z})}$$

[Hellström and Durisi, 2022]

$$\blacktriangleright \text{This paper: } \Delta L_i := L_i^- - L_i^+$$

$$\Rightarrow \textbf{ld-CMI: } I(\Delta L_i; U_i | \tilde{Z})$$

► Step 1: Finding the target.

$$\begin{aligned}\mathcal{E} &= \mathbb{E}_{W, U, \tilde{Z}} \left[\frac{1}{n} \sum_{i=1}^n (-1)^{U_i} \left(\ell(W, \tilde{Z}_i^-) - \ell(W, \tilde{Z}_i^+) \right) \right] \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\Delta L_i, U_i, \tilde{Z}} \left[(-1)^{U_i} \Delta L_i \right]\end{aligned}$$

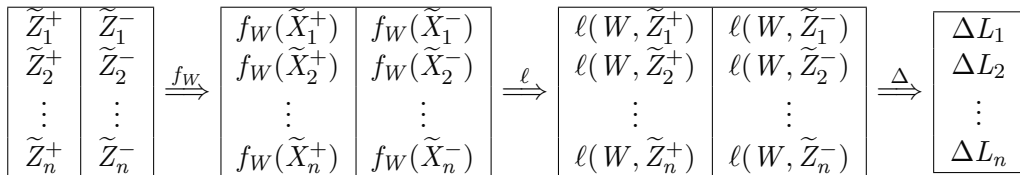
► Step 2: Selecting f , Q and P .

$$\begin{aligned}f &= t \cdot \frac{1}{n} \sum_{i=1}^n (-1)^{U_i} \left(\ell(W, \tilde{z}_i^-) - \ell(W, \tilde{z}_i^+) \right) \\ &= (-1)^{U_i} \Delta L_i\end{aligned}$$

$$Q = P_{W, U | \tilde{z}} = P_{\Delta L_i, U_i | \tilde{z}} \text{ or } P_{\Delta L_i, U_i}$$

$$P = P_{U'} \otimes P_{W | \tilde{z}} = P_{U'} \otimes P_{\Delta L_i | \tilde{z}} \text{ or } P_{U'} \otimes P_{\Delta L_i}$$





$$\underbrace{I(W; U_i | \tilde{Z})}_{\text{CMI}} \geq \underbrace{I(f_W(\tilde{Z}_i); U_i | \tilde{Z})}_{f\text{-CMI [Harutyunyan et al., 2021]}} \geq \underbrace{I(L_i; U_i | \tilde{Z})}_{\text{e-CMI [Hellström and Durisi, 2022]}} \geq \underbrace{I(\Delta L_i; U_i | \tilde{Z})}_{\text{Id-CMI (Ours)}}$$

Theorem 1

Assume the loss is bounded between $[0, 1]$, we have

$$|\mathcal{E}| \leq \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\tilde{Z}} \sqrt{2I^{\tilde{Z}}(\Delta L_i; U_i)} \leq \frac{1}{n} \sum_{i=1}^n \sqrt{2I(\Delta L_i; U_i | \tilde{Z})}, \quad (2)$$

$$|\mathcal{E}| \leq \frac{1}{n} \sum_{i=1}^n \sqrt{2I(\Delta L_i; U_i)}. \quad (3)$$

Theorem 1

Assume the loss is bounded between $[0, 1]$, we have

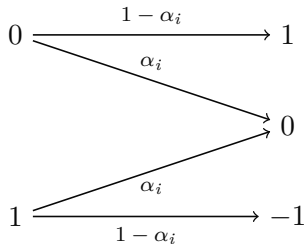
$$|\mathcal{E}| \leq \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\tilde{Z}} \sqrt{2I^{\tilde{Z}}(\Delta L_i; U_i)} \leq \frac{1}{n} \sum_{i=1}^n \sqrt{2I(\Delta L_i; U_i | \tilde{Z})}, \quad (2)$$

$$|\mathcal{E}| \leq \frac{1}{n} \sum_{i=1}^n \sqrt{2I(\Delta L_i; U_i)}. \quad (3)$$

Estimating $I(W; U_i | \tilde{Z}_i)$ vs $I(\Delta L_i; U_i)$:

- ▶ W is a high-dimensional R.V.
- ▶ ΔL_i is an one-dimensional R.V. \implies Easy-to-Compute!





Channel from U_i to ΔL_i . Zero-one loss assumed.

Theorem 2

Under zero-one loss and for any interpolating algorithm \mathcal{A} , $I(\Delta L_i; U_i) = (1 - \alpha_i) \ln 2$ nats for each i , and $|\mathcal{E}| = L_\mu = \sum_{i=1}^n \frac{I(\Delta L_i; U_i)}{n \ln 2}$.

\Rightarrow Generalization error is exactly determined by the communication rate over the channel in the figure averaged over all such channels.

Key observation:

$$\mathcal{E} = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{W, U_i, \tilde{Z}} \left[(-1)^{U_i} \left(\ell(W, \tilde{Z}_i^+) - \ell(W, \tilde{Z}_i^-) \right) \right] = \frac{2}{n} \sum_{i=1}^n \mathbb{E}_{L_i^+, \varepsilon_i} [\varepsilon_i L_i^+],$$

where $\varepsilon_i = (-1)^{\overline{U}_i}$ is the Rademacher variable.

Key observation:

$$\mathcal{E} = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{W, U_i, \tilde{Z}} \left[(-1)^{U_i} \left(\ell(W, \tilde{Z}_i^+) - \ell(W, \tilde{Z}_i^-) \right) \right] = \frac{2}{n} \sum_{i=1}^n \mathbb{E}_{L_i^+, \varepsilon_i} [\varepsilon_i L_i^+],$$

where $\varepsilon_i = (-1)^{\overline{U}_i}$ is the Rademacher variable.

Lemma 5

Consider the weighted generalization error, $\mathcal{E}_{C_1} \triangleq L_\mu - (1 + C_1)L_n$. We have

$$\mathcal{E}_{C_1} = \frac{2 + C_1}{n} \sum_{i=1}^n \mathbb{E}_{L_i^+, \tilde{\varepsilon}_i} [\tilde{\varepsilon}_i L_i^+],$$

where $\tilde{\varepsilon}_i = (-1)^{\overline{U}_i} - \frac{C_1}{C_1+2}$ is a shifted Rademacher variable with mean $-\frac{C_1}{C_1+2}$.

Theorem 3

Let $\ell(\cdot, \cdot) \in [0, 1]$. There exist $C_1, C_2 > 0$ such that

$$L_\mu \leq (1 + C_1)L_n + \sum_{i=1}^n \frac{I(L_i^+; U_i)}{C_2 n}, \quad (4)$$

$$L_\mu \leq L_n + \sum_{i=1}^n \frac{4I(L_i^+; U_i)}{n} + 4\sqrt{\sum_{i=1}^n \frac{L_n I(L_i^+; U_i)}{n}}. \quad (5)$$

Theorem 3

Let $\ell(\cdot, \cdot) \in [0, 1]$. There exist $C_1, C_2 > 0$ such that

$$L_\mu \leq (1 + C_1)L_n + \sum_{i=1}^n \frac{I(L_i^+; U_i)}{C_2 n}, \quad (4)$$

$$L_\mu \leq L_n + \sum_{i=1}^n \frac{4I(L_i^+; U_i)}{n} + 4\sqrt{\sum_{i=1}^n \frac{L_n I(L_i^+; U_i)}{n}}. \quad (5)$$

If $L_n \rightarrow 0$, then (3)(4) vanish with a faster rate.

Theorem 4

For any $\lambda \in (0, 1)$, the “ λ -sharpness” at position i of the training set is defined as

$$F_i(\lambda) \triangleq \mathbb{E}_{W, Z_i} [\ell(W, Z_i) - (1 + \lambda)\mathbb{E}_{W|Z_i} \ell(W, Z_i)]^2.$$

Let $F(\lambda) = \frac{1}{n} \sum_{i=1}^n F_i(\lambda)$. Assume $\ell(\cdot, \cdot) \in \{0, 1\}$, $\lambda \in (0, 1)$. Then, there exist $C_1, C_2 > 0$ such that

$$\mathcal{E} \leq C_1 F(\lambda) + \sum_{i=1}^n \frac{I(L_i^+; U_i)}{C_2 n}. \quad (6)$$

Theorem 4

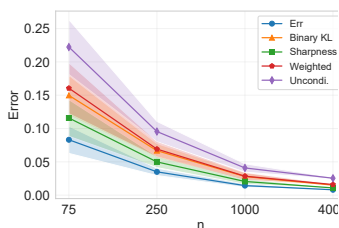
For any $\lambda \in (0, 1)$, the “ λ -sharpness” at position i of the training set is defined as

$$F_i(\lambda) \triangleq \mathbb{E}_{W, Z_i} [\ell(W, Z_i) - (1 + \lambda) \mathbb{E}_{W|Z_i} \ell(W, Z_i)]^2.$$

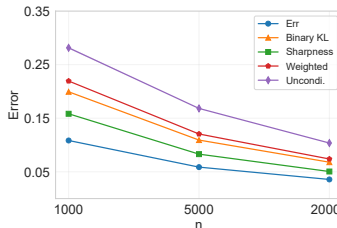
Let $F(\lambda) = \frac{1}{n} \sum_{i=1}^n F_i(\lambda)$. Assume $\ell(\cdot, \cdot) \in \{0, 1\}$, $\lambda \in (0, 1)$. Then, there exist $C_1, C_2 > 0$ such that

$$\mathcal{E} \leq C_1 F(\lambda) + \sum_{i=1}^n \frac{I(L_i^+; U_i)}{C_2 n}. \quad (6)$$

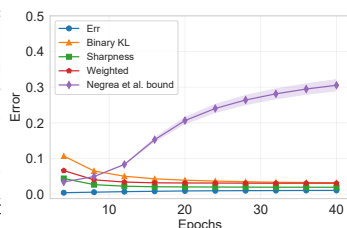
- ▶ $L_n = 0 \rightarrow F(\lambda) = 0$, but $L_n = 0 \nleftrightarrow F(\lambda) = 0$;
- ▶ For any fixed C_1 and C_2 , Eq. (6) is tighter than Eq. (4).



(a) CNN on MNIST



(b) ResNet on CIFAR10



(c) SGLD (MNIST)

$$\text{Uncondi.}: \frac{1}{n} \sum_{i=1}^n \sqrt{2I(\Delta L_i; U_i)}; \quad \text{Binary KL: Hellström and Durisi [2022]}; \quad \text{Weighted:} \\ \sum_{i=1}^n \frac{4I(L_i^+; U_i)}{n} + 4\sqrt{\sum_{i=1}^n \frac{L_n I(L_i^+; U_i)}{n}}; \quad \text{Sharpness: } C_1 F(\lambda) + \sum_{i=1}^n \frac{I(L_i^+; U_i)}{C_2 n}.$$

Information-Theoretic Bounds in Stochastic Convex Optimization

Limitations of Information-Theoretic (IT) bounds:

- ▶ Original input-output mutual information (IOMI) (e.g., $I(W; S)$ [Xu and Raginsky, 2017]) based bound can $\rightarrow \infty$
 \implies solved by conditional mutual information (CMI) $I(W; U|\tilde{Z})$ [Steinke and Zakynthinou, 2020]
- ▶ Slow convergence rate, e.g., $\mathcal{O}(1/\sqrt{n})$
 \implies mitigated by [Haghifam et al., 2021, Hellström and Durisi, 2021, 2022, Wang and Mao, 2023, Wu et al., 2023, Zhou et al., 2023]

Limitations of Information-Theoretic (IT) bounds:

- ▶ Original input-output mutual information (IOMI) (e.g., $I(W; S)$ [Xu and Raginsky, 2017]) based bound can $\rightarrow \infty$
 \implies solved by conditional mutual information (CMI) $I(W; U|\tilde{Z})$ [Steinke and Zakynthinou, 2020]
- ▶ Slow convergence rate, e.g., $\mathcal{O}(1/\sqrt{n})$
 \implies mitigated by [Haghifam et al., 2021, Hellström and Durisi, 2021, 2022, Wang and Mao, 2023, Wu et al., 2023, Zhou et al., 2023]
- ▶ Non-vanishing in Stochastic Convex Optimization (SCO) problems for (nearly) all previous IT bounds! [Haghifam et al., 2023]

Limitations of Information-Theoretic (IT) bounds:

- ▶ Original input-output mutual information (IOMI) (e.g., $I(W; S)$ [Xu and Raginsky, 2017]) based bound can $\rightarrow \infty$
 \implies solved by conditional mutual information (CMI) $I(W; U|\tilde{Z})$ [Steinke and Zakynthinou, 2020]
- ▶ Slow convergence rate, e.g., $\mathcal{O}(1/\sqrt{n})$
 \implies mitigated by [Haghifam et al., 2021, Hellström and Durisi, 2021, 2022, Wang and Mao, 2023, Wu et al., 2023, Zhou et al., 2023]
- ▶ Non-vanishing in Stochastic Convex Optimization (SCO) problems for (nearly) all previous IT bounds! [Haghifam et al., 2023]
Wang, Z. and Mao, Y.. Sample-Conditioned Hypothesis Stability Sharpens Information-Theoretic Generalization Bounds. NeurIPS 2023.
Our contribution: **Incorporating stability-based analysis into IT framework which improves both stability-based bounds and IT bounds.**

- ▶ Z is one-hot vector in \mathbb{R}^d
- ▶ Loss: $-\langle w, z \rangle$; ERM solution $W = \frac{1}{n} \sum_{i=1}^n Z_i$
- ▶ Birthday Paradox Problem: For a large d , the probability that no pair of instances in \tilde{Z} sharing the same non-zero coordinate (referred to as event E_0) is smaller than a constant probability (independent of n).
- ▶ If $d \geq \frac{2n-1}{1-c^{1/(2n-1)}}$, then $P(E_0) \geq c \geq \left(1 - \frac{2n-1}{d}\right)^{2n-1}$, e.g., $d = 2n^2 \implies c \geq 0.1$.
- ▶ Let $d = 2n^2$, $I(W; U_i | \tilde{Z}_i) = \log 2 - H(U_i | W, \tilde{Z}_i) \geq 0.1 \cdot \log 2$.
- ▶ **CMI bound is non-vanishing but $\mathcal{E} \leq \mathcal{O}(1/\sqrt{n})$.**

$$\begin{array}{l} Z_1, \dots, \left| \begin{array}{c} \textcolor{red}{Z}_i, \\ \end{array} \right| \dots, Z_n \xrightarrow{\mathcal{A}} \left| \begin{array}{c} \textcolor{red}{W} \\ \end{array} \right| \Rightarrow \ell(W, Z) \\ Z_1, \dots, \left| \begin{array}{c} \textcolor{red}{Z}'_i, \\ \end{array} \right| \dots, Z_n \xrightarrow{\mathcal{A}} \left| \begin{array}{c} \textcolor{red}{W}^{-i} \\ \end{array} \right| \Rightarrow \ell(W^{-i}, Z) \end{array}$$

\mathcal{A} is Stable \iff Loss of (W^{-i}, Z) is close to Loss of (W, Z) .

- ▶ Uniform Stability [Bousquet and Elisseeff, 2002]:
 $\sup_{W, W^{-i}, Z} |\ell(W, Z) - \ell(W^{-i}, Z)| \leq \text{Unif. Stability Param.}$
- ▶ Sample-Conditioned Hypothesis (SCH) Stability in our paper
 $\mathbb{E}_{W, W^{-i}} [\sup_Z |\ell(W, Z) - \ell(W^{-i}, Z)|] \leq \text{SCH Stability Param.},$
where Z can be either Z_i or Z'_i .

By DV lemma: $\mathcal{E} \leq \inf_{t>0} \frac{\text{IOMI or CMI+CGF}}{t}$.

where

$$\text{CGF} = \log \mathbb{E} [\exp (t \cdot f_{\text{DV}})].$$

By DV lemma: $\mathcal{E} \leq \inf_{t>0} \frac{\text{IOMI or CMI+CGF}}{t}$.

where

$$\text{CGF} = \log \mathbb{E} [\exp (t \cdot f_{\text{DV}})].$$

► Previous works:

$f_{\text{DV}} = \ell(W, Z')$ e.g., [Bu et al., 2019]

$f_{\text{DV}} = \ell(W, Z') - \mathbb{E}_{Z'} [\ell(W, Z')]$ e.g., [Wu et al., 2023]

$f_{\text{DV}} = (-1)^U (\ell(W, Z_1) - \ell(W, Z_2))$ e.g., [Steinke and Zakynthinou, 2020]

By DV lemma: $\mathcal{E} \leq \inf_{t>0} \frac{\text{IOMI or CMI+CGF}}{t}$.

where

$$\text{CGF} = \log \mathbb{E} [\exp (t \cdot f_{\text{DV}})].$$

- Previous works:

$$f_{\text{DV}} = \ell(W, Z') \text{ e.g., [Bu et al., 2019]}$$

$$f_{\text{DV}} = \ell(W, Z') - \mathbb{E}_{Z'} [\ell(W, Z')] \text{ e.g., [Wu et al., 2023]}$$

$$f_{\text{DV}} = (-1)^U (\ell(W, Z_1) - \ell(W, Z_2)) \text{ e.g., [Steinke and Zakynthinou, 2020]}$$

- This paper: let W^{-i} be obtained by replacing one data in S ,

$$f_{\text{DV}} = \ell(W, Z') - \mathbb{E}_{W^{-i}|W} [\ell(W^{-i}, Z')] \implies \text{IOMI}$$

$$f_{\text{DV}} = (-1)^U (\ell(W, Z) - \ell(W^{-i}, Z)) \implies \text{New CMI}$$

$$\begin{array}{|c|c|} \hline \tilde{Z}_{1,0} & \blacksquare \\ \hline \tilde{Z}_{2,0} & \blacksquare \\ \hline \vdots & \vdots \\ \hline \tilde{Z}_{n,0} & \blacksquare \\ \hline \end{array} \xrightarrow{A} W, \quad \begin{array}{|c|c|} \hline \blacksquare & \tilde{Z}_{1,1} \\ \hline \tilde{Z}_{2,0} & \blacksquare \\ \hline \vdots & \vdots \\ \hline \tilde{Z}_{n,0} & \blacksquare \\ \hline \end{array} \xrightarrow{A} \underbrace{W^{-1}, \dots, W^{-n}}_{n \text{ new neighboring hypotheses}} \Rightarrow \begin{array}{|c|c|} \hline W & W^{-1} \\ \hline W & W^{-2} \\ \hline \vdots & \vdots \\ \hline W & W^{-n} \\ \hline \end{array}$$

The generalization game:

$$\begin{array}{|c|c|} \hline \tilde{Z}_{1,0} & \tilde{Z}_{1,1} \\ \hline \tilde{Z}_{2,0} & \tilde{Z}_{2,1} \\ \hline \vdots & \vdots \\ \hline \tilde{Z}_{n,0} & \tilde{Z}_{n,1} \\ \hline \end{array} + \begin{array}{|c|c|} \hline \square & \blacksquare \\ \hline \blacksquare & \square \\ \hline \vdots & \vdots \\ \hline \square & \blacksquare \\ \hline \end{array} \Rightarrow \begin{array}{|c|c|} \hline \tilde{Z}_{1,0} & \blacksquare \\ \hline \blacksquare & \tilde{Z}_{2,1} \\ \hline \vdots & \vdots \\ \hline \tilde{Z}_{n,0} & \blacksquare \\ \hline \end{array} + \begin{array}{|c|c|} \hline W & W^{-1} \\ \hline W & W^{-2} \\ \hline \vdots & \vdots \\ \hline W & W^{-n} \\ \hline \end{array} \Rightarrow \begin{array}{|c|c|} \hline \text{Tr. Err}_1 & \text{Ts. Err}_1 \\ \hline \text{Ts. Err}_2 & \text{Tr. Err}_2 \\ \hline \vdots & \vdots \\ \hline \text{Ts. Err}_n & \text{Tr. Err}_n \\ \hline \end{array}$$

$$\Rightarrow \mathcal{E} = \frac{1}{n} \sum_{i=1}^n \mathbb{E} [\text{Ts. Err}_i - \text{Tr. Err}_i] (\leq \text{Stability Param.})$$

Given $\begin{array}{|c|c|} \hline W & W^{-1} \\ \hline W & W^{-2} \\ \hline \vdots & \vdots \\ \hline W & W^{-n} \\ \hline \end{array}$ and $\begin{array}{|c|} \hline \tilde{Z}_{1,0} \\ \hline \tilde{Z}_{2,1} \\ \hline \vdots \\ \hline \tilde{Z}_{n,0} \\ \hline \end{array}$, can you infer $\begin{array}{|c|c|} \hline \square & \blacksquare \\ \hline \blacksquare & \square \\ \hline \vdots & \vdots \\ \hline \square & \blacksquare \\ \hline \end{array}$ 🤔?

Theorem 5 (Informal.)

If \mathcal{A} is β -stable, we have $\mathcal{E} \lesssim \beta \sqrt{I(Z_U; U | W, W^{-i})} \leq \beta \sqrt{I(W; Z_i)}$

In SCO counterexamples given by Haghifam et al. [2023]:

$$\mathcal{E} \leq \mathcal{O}(1/\sqrt{n}).$$

- ▶ Previous IOMI or CMI bound in these examples:
SubGaussian param. $R = \mathcal{O}(1)$ (=Lip. Para. \times Diam. of hypo. space)
and $\text{IOMI} \geq \text{CMI} = \mathcal{O}(1)$.
 $\implies \text{IOMI bound} \geq \text{CMI bound} \in \mathcal{O}(1) \implies$ fail to explain the learnability.
- ▶ New CMI bound in these examples:
 $\beta = \mathcal{O}(1/\sqrt{n})$
and New CMI = $\mathcal{O}(1)$.
 \implies **New CMI bound** $\in \mathcal{O}(1/\sqrt{n}) \implies$ can explain the learnability.
- ▶ More bounds, e.g., fast-rate bounds and second-moment bounds.
- ▶ More examples, e.g., our bounds can also improve stability-based bounds.



CMI and VC-dim:

Theorem 6

Let $\mathcal{Z} = \mathcal{X} \times \{0, 1\}$, and let $\mathcal{F} = \{f_w : \mathcal{X} \rightarrow \{0, 1\} | w \in \mathcal{W}\}$ be a functional hypothesis class with finite VC dimension d . Let $n > d + 1$, for any algorithm \mathcal{A} ,

$$\frac{1}{n} \sum_{i=1}^n \sqrt{I(F_i^+, F_i^-; U_i | \tilde{Z}_i)} \leq \mathcal{O} \left(\sqrt{\frac{d}{n} \log \left(\frac{n}{d} \right)} \right).$$

Proof Sketch.

For a given \tilde{Z} , the number of distinct values of their predictions, denoted by k , by Sauer-Shelah lemma for $n > d + 1$, $k \leq \sum_{i=1}^d \binom{n}{i} \leq \left(\frac{en}{d}\right)^d$.

$$I(F^+, F^-; U | \tilde{Z}) \leq H(F^+, F^- | \tilde{Z}) \leq H(F^+ | \tilde{Z}) + H(F^- | \tilde{Z}) \leq 2d \log \left(\frac{en}{d} \right).$$

- ▶ *Steinke, T., and Zakynthinou, L.. Open problem: Information complexity of vc learning. COLT 2020.*
- ▶ *Hafez-Kolahi, H. et al. Conditioning and processing: Techniques to improve information-theoretic generalization bounds. NeurIPS 2020.*
- ▶ *Haghifam, M. et al. Towards a unified information-theoretic framework for generalization. NeurIPS 2021.*
- ▶ *LOO CMI: Haghifam, M. et al. Understanding Generalization via Leave-One-Out Conditional Mutual Information. ISIT 2022.*
- ▶ *f -CMI and e -CMI: Harutyunyan, H. et al. Information-theoretic generalization bounds for black-box learning algorithms. NeurIPS 2021.*
Hellström, F. and Durisi, G.. A new family of generalization bounds using samplewise evaluated CMI. NeurIPS 2022.
- ▶ *Bassily, R. et al. Learners that use little information. ALT 2018.*
- ▶ *Livni, R.. Information Theoretic Lower Bounds for Information Theoretic Upper Bounds. NeurIPS 2023.*



Information-Theoretic Generalization Bounds for SGD


Lemma 6 (Xu and Raginsky [2017])

Assume the loss $\ell(w, Z)$ is R -subgaussian² for any $w \in \mathcal{W}$. The generalization error of \mathcal{A} is bounded by

$$|\mathcal{E}| \leq \sqrt{\frac{2R^2}{n} I(W; S)},$$

Mutual information $I(W; S) \triangleq D_{\text{KL}}(P_{W,S} \| P_W \otimes P_S)$.

\implies Distribution-dependent and Algorithm-dependent

²A random variable X is R -subgaussian if for any ρ , $\log \mathbb{E} \exp(\rho(X - \mathbb{E}X)) \leq \rho^2 R^2 / 2$.  uOttawa

SGLD updates:

$$W_t \triangleq W_{t-1} - \lambda_t g(W_{t-1}, B_t) + N_t,$$

where

$$g(w, B_t) \triangleq \frac{1}{b} \sum_{z \in B_t} \nabla_w \ell(w, z),$$

- ▶ λ_t : learning rate
- ▶ b : batch size
- ▶ B_t denotes the batch used for the t^{th} update
- ▶ $N_t \sim \mathcal{N}(0, \sigma_t^2 \mathbf{I}_d)$

Assume SGLD outputs W_T as the learned model parameter.



$$\begin{aligned} I(W_T; S) &= I(W_{T-1} - \lambda_T g(W_{T-1}, B_T) + N_T; S) \\ &\leq I(W_{T-1}, -\lambda_T g(W_{T-1}, B_T) + N_T; S) \end{aligned} \quad (7)$$

$$= I(W_{T-1}; S) + I(-\lambda_T g(W_{T-1}, B_T) + N_T; S | W_{T-1}) \quad (8)$$

$$\vdots$$

$$\leq \sum_{t=1}^T I(-\lambda_t g(W_{t-1}, B_t) + N_t; S | W_{t-1})$$

$$\begin{aligned} & I(-\lambda_t g(W_{t-1}, B_t) + N_t; S | W_{t-1}) \\ &= \mathbb{E}_{S, W_{t-1}} \left[\text{D}_{\text{KL}} \left(Q_{-\lambda_t g(W_{t-1}, B_t) + N_t | S, W_{t-1}} \parallel P_{-\lambda_t g(W_{t-1}, B'_t) + N_t | W_{t-1}} \right) \right] \\ &\leq \frac{d}{2} \mathbb{E}_{W_{t-1}} \log \left(\frac{\lambda_t^2 \mathbb{E}_S^{W_{t-1}} \|g - \mathbb{E}[g]\|_2^2}{d\sigma_t^2} + 1 \right). \end{aligned}$$



Theorem 7

Gen. err. of SGLD is upper bounded by

$$\mathcal{E} \lesssim \sqrt{\frac{d}{n} \sum_{t=1}^T \mathbb{E} \log \left(\frac{\lambda_t^2 \mathbb{E} \|g - \mathbb{E}[g]\|_2^2}{d\sigma_t^2} + 1 \right)}.$$

- ▶ Bu, Y. et al. *Tightening Mutual Information Based Bounds on Generalization Error. ISIT 2019.*
Negrea, J. et al. *Information-theoretic generalization bounds for SGLD via data-dependent estimates. NeurIPS 2019.*
Haghifam, M. et al. *Sharpened generalization bounds based on conditional mutual information and an application to noisy, iterative algorithms. NeurIPS 2020.*
Rodríguez-Gálvez, B. et al. *On random subset generalization error bounds and the stochastic gradient Langevin dynamics algorithm. ITW 2020.*
- ▶ Wang, Hao et al. *Analyzing the generalization capability of sgld using properties of gaussian channels. NeurIPS 2021.*
- ▶ Li, J. et al. *On generalization error bounds of noisy gradient methods for non-convex learning. ICLR 2020.*

- ▶ *Mou, W.. Generalization bounds of sgld for non-convex learning: Two theoretical viewpoints. COLT 2018.*
- ▶ *Banerjee, A. et al. Stability based generalization bounds for exponential family langevin dynamics. ICML 2022.*
- ▶ *Futami, F., and Fujisawa, M.. Time-Independent Information-Theoretic Generalization Bounds for SGLD. NeurIPS 2023.*

SGD updates:

$$W_t \triangleq W_{t-1} - \lambda_t g(W_{t-1}, B_t),$$

where

$$g(w, B_t) \triangleq \frac{1}{b} \sum_{z \in B_t} \nabla_w \ell(w, z),$$

- ▶ λ_t : learning rate
- ▶ b : batch size
- ▶ B_t denotes the batch used for the t^{th} update.

Assume SGD outputs W_T as the learned model parameter.



SGD updates:

$$W_t \triangleq W_{t-1} - \lambda_t g(W_{t-1}, B_t),$$

where

$$g(w, B_t) \triangleq \frac{1}{b} \sum_{z \in B_t} \nabla_w \ell(w, z),$$

- ▶ λ_t : learning rate
- ▶ b : batch size
- ▶ B_t denotes the batch used for the t^{th} update.

Assume SGD outputs W_T as the learned model parameter.

Difficulty of using MI based bound: $I(W_T; S) \rightarrow \text{too large for SGD}$



Follow up the work of Neu et al. [2021], let $\{\sigma_t\}_{t=1}^T$ be a sequence of positive real numbers.

Define $\widetilde{W}_0 \triangleq W_0$, and $\widetilde{W}_t \triangleq \widetilde{W}_{t-1} - \lambda_t g(W_{t-1}, B_t) + N_t$, for $t > 0$, where $N_t \sim \mathcal{N}(0, \sigma_t^2 \mathbf{I}_d)$ is a Gaussian noise.

$$\begin{array}{cccccccc}
 & & N_1 & & N_2 & & \cdots & & N_{T-1} & & N_T \\
 & & \downarrow & & \downarrow & & & & \downarrow & & \downarrow \\
 \widetilde{W}_0 & \rightarrow & \widetilde{W}_1 & \rightarrow & \widetilde{W}_2 & \rightarrow & \cdots & \rightarrow & \widetilde{W}_{T-1} & \rightarrow & \widetilde{W}_T \\
 \parallel & \nearrow & & \nearrow & & \nearrow & & & & \nearrow & \\
 W_0 & \rightarrow & W_1 & \rightarrow & W_2 & \rightarrow & \cdots & \rightarrow & W_{T-1} & \rightarrow & W_T
 \end{array}$$

Let $\Delta_t = \sum_{\tau=1}^t N_\tau$. Notice that $\widetilde{W}_t = W_t + \Delta_t$.

Denote this auxiliary weight process by \mathcal{A}_{AWP} . Let \mathcal{A}_{SGD} be the original algorithm of SGD,

$$\begin{aligned}\mathcal{E}_\mu(\mathcal{A}_{SGD}) &= \mathcal{E}_\mu(\mathcal{A}_{SGD}) + \mathcal{E}_\mu(\mathcal{A}_{AWP}) - \mathcal{E}_\mu(\mathcal{A}_{AWP}) \\ &\leq \underbrace{\mathcal{O}\left(\sqrt{\frac{I(\widetilde{W}_T; S)}{n}}\right)}_{\text{Lemma 3}} + \underbrace{|\mathcal{E}_\mu(\mathcal{A}_{SGD}) - \mathcal{E}_\mu(\mathcal{A}_{AWP})|}_{\text{residual term}}\end{aligned}\tag{9}$$

\asymp Trajectories of Gradient Variance/Dispersion + Sharpness.

Theorem 8 (Wang and Mao [2022])

The generalization error of SGD is upper bounded by

$$\mathcal{E} \lesssim \sqrt[3]{\sum_{t=1}^T \frac{\mathbb{E} [\mathbb{V}_t(W_{t-1})] \mathbb{E} [\text{Tr}(H_{W_T}(Z))]}{n}} \quad (10)$$

► Gradient Dispersion: $\mathbb{V}_t(w) \triangleq \mathbb{E}_S [\|g(w, B_t) - \mathbb{E}_{W,Z} [\nabla_w \ell(W, Z)]\|_2^2]$

SDE updates: $W_t \triangleq W_{t-1} - \eta g(W_{t-1}, S) + \eta C_t^{1/2} N_t$, where

$$C_t \triangleq \frac{n-b}{b(n-1)} \left(\frac{1}{n} \sum_{i=1}^n \nabla \ell_i \nabla \ell_i^T - G_t G_t^T \right)$$

is the *gradient noise covariance* matrix.

Denote SDE approximation as \mathcal{A}_{SDE} ,

$$\begin{aligned} \mathcal{E}_\mu(\mathcal{A}_{SGD}) &= \mathcal{E}_\mu(\mathcal{A}_{SGD}) + \mathcal{E}_\mu(\mathcal{A}_{SDE}) - \mathcal{E}_\mu(\mathcal{A}_{SDE}) \\ &\leq \underbrace{\mathcal{O} \left(\sqrt{\frac{I(W_{SDE}; S)}{n}} \right)}_{\text{Lemma 3}} + \underbrace{|\mathcal{E}_\mu(\mathcal{A}_{SGD}) - \mathcal{E}_\mu(\mathcal{A}_{SDE})|}_{\text{residual term}}, \end{aligned} \quad (11)$$

where W_{SDE} is the output hypothesis by \mathcal{A}_{SDE} .

$$\mathcal{E}_\mu(\mathcal{A}_{SGD}) \leq \underbrace{\mathcal{O}\left(\sqrt{\frac{I(W_{SDE}; S)}{n}}\right)}_{\text{Lemma 3}} + \underbrace{|\mathcal{E}_\mu(\mathcal{A}_{SGD}) - \mathcal{E}_\mu(\mathcal{A}_{SDE})|}_{\text{residual term}}.$$

Empirical evidence from [Wu et al., 2020, Li et al., 2021] suggests that the residual term is small.

⇒ It is safe to investigate the generalization of SGD using the IT bounds of SDE directly.

Information-Theoretic Analysis Beyond Supervised Learning

Wang, Z. and Mao Y.. Information-theoretic analysis of unsupervised domain adaptation. ICLR 2023.

- ▶ Novel upper bounds for generalization error of UDA.
- ▶ Simple regularization technique for improving generalization of UDA

- ▶ Source data $Z = (X, Y) \sim \mu$ and target data $Z' = (X', Y') \sim \mu'$
- ▶ Labeled source sample: $S = \{Z_i\}_{i=1}^n \stackrel{\text{i.i.d}}{\sim} \mu^{\otimes n}$; Unlabelled target sample $S'_{X'} = \{X'_j\}_{j=1}^m \stackrel{\text{i.i.d}}{\sim} P_{X'}^{\otimes m}$
- ▶ *Generalization error = testing error of target domain - training error of source domain:*

$$\begin{aligned}\mathcal{E} &\triangleq \mathbb{E}_{W, S, S'_{X'}} [R_{\mu'}(W) - R_S(W)] \\ &= \mathbb{E}_{W, S, S'_{X'}} [L_{\mu'}(W) - \cancel{L_{\mu}(W)} + \cancel{L_{\mu}(W)} - L_S(W)]\end{aligned}$$

Theorem 9

Assume $\ell(f_w(X'), Y')$ is R -subgaussian. Then

$$|\mathcal{E}| \leq \frac{1}{nm} \sum_{j=1}^m \sum_{i=1}^n \mathbb{E}_{X'_j} \sqrt{2R^2 I^{X'_j}(W; Z_i)} + \sqrt{2R^2 D_{\text{KL}}(\mu || \mu')}.$$

Consider SGLD. At each time step t ,

- ▶ labelled source mini-batch: Z_{B_t} ; unlabelled target mini-batch: X'_{B_t}
- ▶ gradient: $G_t = g(W_{t-1}, Z_{B_t}, X'_{B_t})$
- ▶ updating rule: $W_t = W_{t-1} - \eta_t G_t + N_t$ where $N_t \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_d)$.

Theorem 10

Under the assumption of Theorem 9. Let the total iteration number be T , then

$$|\mathcal{E}| \leq \sqrt{\frac{R^2}{n} \sum_{t=1}^T \frac{\eta_t^2}{\sigma_t^2} \mathbb{E}_{S'_{X'}, W_{t-1}, S} \left[\left\| G_t - \mathbb{E}_{Z_{B_t}} [G_t] \right\|^2 \right]} + \sqrt{2R^2 \text{D}_{\text{KL}}(\mu || \mu')}.$$

restrict the gradient norm \implies reduce $|\mathcal{E}|$.

RotatedMNIST is built based on the MNIST dataset and consists of six domains, which are rotated MNIST images with rotation angle 0° , 15° , 30° , 45° , 60° and 75° .

RotatedMNIST.

Method	RotatedMNIST (0° as source domain)					Ave
	15°	30°	45°	60°	75°	
ERM	97.5 ± 0.2	84.1 ± 0.8	53.9 ± 0.7	34.2 ± 0.4	22.3 ± 0.5	58.4
DANN	97.3 ± 0.4	90.6 ± 1.1	68.7 ± 4.2	30.8 ± 0.6	19.0 ± 0.6	61.3
MMD	97.5 ± 0.1	95.3 ± 0.4	73.6 ± 2.1	44.2 ± 1.8	32.1 ± 2.1	68.6
CORAL	97.1 ± 0.3	82.3 ± 0.3	56.0 ± 2.4	30.8 ± 0.2	27.1 ± 1.7	58.7
WD	96.7 ± 0.3	93.1 ± 1.2	64.1 ± 3.3	41.4 ± 7.6	27.6 ± 2.0	64.6
KL	97.8 ± 0.1	97.1 ± 0.2	93.4 ± 0.8	75.5 ± 2.4	68.1 ± 1.8	86.4
ERM-GP	97.5 ± 0.1	86.2 ± 0.5	62.0 ± 1.9	34.8 ± 2.1	26.1 ± 1.2	61.2
KL-GP	98.2 ± 0.2	96.9 ± 0.1	95.0 ± 0.6	88.0 ± 8.1	78.1 ± 2.5	91.2

- ▶ **Semi-supervised learning:** He, H. et al. *Information-theoretic characterization of the generalization error for iterative semisupervised learning. JMLR 2022.*
Aminian, G. et al. *An information-theoretical approach to semi-supervised learning under covariate-shift. AISTATS 2022.*
- ▶ **Transfer Learning:** Wu, X. et al. *Informationtheoretic analysis for transfer learning. ISIT 2020.*
Bu, Y. et al. *Characterizing and understanding the generalization error of transfer learning with gibbs algorithm. AISTATS 2022.*
- ▶ **Meta Learning:** Jose, S. T. and Simeone, O. *Information-theoretic generalization bounds for meta-learning and applications. Entropy 2021.*
Chen, Q. et al. *Generalization bounds for meta-learning: An information-theoretic analysis. NeurIPS 2021*
Hellström, F. and Durisi, G. *Evaluated CMI bounds for meta learning: Tightness and expressiveness. NeurIPS 2022.*

- ▶ **Federated Learning:** *Yagli, S. et al. Information-theoretic bounds on the generalization error and privacy leakage in federated learning. SPAWC 2020.*
- ▶ **Transductive Learning:** *Tang, H. and Liu, Y. Information-Theoretic Generalization Bounds for Transductive Learning and its Applications. arXiv preprint arXiv:2311.04561, 2023.*
- ▶ **Quantum Learning:** *Caro, M. et al. Information-theoretic generalization bounds for learning from quantum data. arXiv preprint arXiv:2311.05529, 2023.*

- ▶ Wang, H. et al. *An information-theoretic view of generalization via Wasserstein distance*. ISIT 2019.
- ▶ Rodríguez-Gálvez, B. et al. *Tighter expected generalization error bounds via Wasserstein distance*. NeurIPS 2021.
- ▶ Lugosi, G., and Neu, G.. *Generalization bounds via convex analysis*. COLT 2022.
- ▶ Lugosi, G., and Neu, G.. *Online-to-PAC Conversions: Generalization Bounds via Regret Analysis*. arXiv preprint arXiv:2305.19674, 2023.
- ▶ Chu, Y., and Raginsky, M.. *A unified framework for information-theoretic generalization bounds*. NeurIPS 2023.

- Aolin Xu and Maxim Raginsky. Information-theoretic analysis of generalization capability of learning algorithms. *Advances in Neural Information Processing Systems*, 2017.
- Thomas Steinke and Lydia Zakynthinou. Reasoning about generalization via conditional mutual information. In *Conference on Learning Theory*. PMLR, 2020.
- Hrayr Harutyunyan, Maxim Raginsky, Greg Ver Steeg, and Aram Galstyan. Information-theoretic generalization bounds for black-box learning algorithms. In *Advances in Neural Information Processing Systems*, 2021.
- Fredrik Hellström and Giuseppe Durisi. A new family of generalization bounds using samplewise evaluated CMI. In *Advances in Neural Information Processing Systems*, 2022.



- Mahdi Haghifam, Gintare Karolina Dziugaite, Shay Moran, and Dan Roy. Towards a unified information-theoretic framework for generalization. *Advances in Neural Information Processing Systems*, 34:26370–26381, 2021.
- Fredrik Hellström and Giuseppe Durisi. Fast-rate loss bounds via conditional information measures with applications to neural networks. In *2021 IEEE International Symposium on Information Theory (ISIT)*, pages 952–957. IEEE, 2021.
- Ziqiao Wang and Yongyi Mao. Tighter information-theoretic generalization bounds from supersamples. In *International Conference on Machine Learning*. PMLR, 2023.
- Xuetong Wu, Jonathan H Manton, Uwe Aickelin, and Jingge Zhu. On the tightness of information-theoretic bounds on generalization error of learning algorithms. *arXiv preprint arXiv:2303.14658*, 2023.



- Ruida Zhou, Chao Tian, and Tie Liu. Exactly tight information-theoretic generalization error bound for the quadratic gaussian problem. *arXiv preprint arXiv:2305.00876*, 2023.
- Mahdi Haghifam, Borja Rodríguez-Gálvez, Ragnar Thobaben, Mikael Skoglund, Daniel M Roy, and Gintare Karolina Dziugaite. Limitations of information-theoretic generalization bounds for gradient descent methods in stochastic convex optimization. In *International Conference on Algorithmic Learning Theory*, pages 663–706. PMLR, 2023.
- Olivier Bousquet and André Elisseeff. Stability and generalization. *The Journal of Machine Learning Research*, 2:499–526, 2002.
- Yuheng Bu, Shaofeng Zou, and Venugopal V Veeravalli. Tightening mutual information based bounds on generalization error. In *2019 IEEE International Symposium on Information Theory (ISIT)*, pages 587–591. IEEE, 2019.



- Gergely Neu, Gintare Karolina Dziugaite, Mahdi Haghifam, and Daniel M Roy. Information-theoretic generalization bounds for stochastic gradient descent. In *Conference on Learning Theory*. PMLR, 2021.
- Ziqiao Wang and Yongyi Mao. On the generalization of models trained with SGD: Information-theoretic bounds and implications. In *International Conference on Learning Representations*, 2022.
- Jingfeng Wu, Wenqing Hu, Haoyi Xiong, Jun Huan, Vladimir Braverman, and Zhanxing Zhu. On the noisy gradient descent that generalizes as sgd. In *International Conference on Machine Learning*. PMLR, 2020.
- Zhiyuan Li, Sadhika Malladi, and Sanjeev Arora. On the validity of modeling sgd with stochastic differential equations (sdes). *Advances in Neural Information Processing Systems*, 2021.



Thank You!

