

Reasoning about Generalization via Conditional Mutual Information



Thomas Steinke*
Google Brain

Lydia Zakynthinou*
Northeastern University

Generalization in Machine Learning

- We sample an i.i.d dataset Z of size n from an unknown distribution D over \mathcal{Z} : $Z \leftarrow D^n$.
- Loss function ℓ indicates the quality of a model w .

Empirical loss

$$\ell(w, Z)$$

||

$$\frac{1}{n} \sum_{i=1}^n \ell(w, Z_i)$$

True loss

$$\ell(w, D)$$

||

$$\mathbb{E}_{Z \leftarrow D} [\ell(w, Z)]$$

- How can we ensure:

$$\ell(A(Z), Z) \approx \ell(A(Z), D)$$

Conditional Mutual Information (CMI)

- Draw $2n$ samples $\tilde{Z} \leftarrow D^{2n}$.
- Selector function $S \in \{0,1\}^n$ uniformly random defines partition of real samples \tilde{Z}_S and “ghost” samples $\tilde{Z}_{\bar{S}}$.

$$\tilde{Z}_S = (\tilde{Z}_{1,S_1}, \dots, \tilde{Z}_{n,S_n})$$

- Run $A(\tilde{Z}_S)$.
- CMI of A with respect to D :

$$I(A(\tilde{Z}_S); S | \tilde{Z})$$

$\tilde{Z}_{1,0}$	$\tilde{Z}_{1,1}$
$\tilde{Z}_{2,0}$	$\tilde{Z}_{2,1}$
\vdots	\vdots
$\tilde{Z}_{n,0}$	$\tilde{Z}_{n,1}$

- ✓ Post-processing and (non-adaptive) composition
- ✓ CMI is finite (“normalized” MI)

Low CMI implies Generalization

For loss in $[0,1]$.

- $|\mathbb{E}[\ell(A(Z), Z) - \ell(A(Z), D)]| \leq \sqrt{\frac{2 \cdot CMI_D(A)}{n}}$
- $\mathbb{E}[(\ell(A(Z), Z) - \ell(A(Z), D))^2] \leq \frac{3 \cdot CMI_D(A) + 2}{n}$
- $\mathbb{E}[\ell(A(Z), D)] \leq 2 \cdot \mathbb{E}[\ell(A(Z), Z)] + \frac{3 \cdot CMI_D(A)}{n}$

Extension to “unbounded” loss (e.g. hinge loss) and to non-linear losses (e.g. Area Under the ROC Curve).

Methods used to prove generalization

- Uniform Convergence: VC dimension [VC71]
- Distributional Stability: Differential Privacy [DMNS06]
- Uniform Stability [BE02]
- Local Statistical Stability [LS19]
- Mutual Information Methods**
 - ✓ Low MI \Rightarrow Generalization [RZ16/RX17]
 - ✓ ϵ -DP $\Rightarrow I(A(Z); Z) \leq \frac{1}{2} \epsilon^2 n$
 - ✓ $I(A(Z); Z) \leq \log |\mathcal{W}|$
 - ! (ϵ, δ) -DP $\not\Rightarrow$ bound on MI [D12/MMPRTV10]
 - ! For 1-D thresholds, any consistent learner has $I(A(Z); Z) = \infty$ for some D . [BMNSY18]

Bounding CMI

- ✓ $\sqrt{\epsilon}$ -Differential Privacy
- ✓ ϵ -Mutual Information Stability
- ✓ ϵ -KL Stability
- ✓ ϵ -Average-Leave-One-Out KL Stability
- ✓ δ -TV Stability (i.e., $(0, \delta)$ -DP) $\Rightarrow CMI_D(A) \leq \delta n$
- ✓ Compression schemes of size $k \Rightarrow CMI_D(A) \leq k \log(n)$
- ✓ Hypothesis class with VC dimension $d \Rightarrow \exists$ ERM such that $CMI_D(A) \leq d \log(n)$

$$CMI_D(A) \leq \epsilon n$$

Conclusion

CMI is a new framework for reasoning about generalization, which:

- Unifies existing frameworks
- Provides a variety of forms of generalization guarantees.

