PAC-Chernoff Bounds: Understanding Generalization in the Interpolation Regime

Andrés R. Masegosa and Luis A. Ortega

ECAI, November 2025

https://www.jair.org/index.php/jair/article/view/17036

Aalborg University and Autonomous University of Madrid

Motivations

Generalization bounds that solely depend on the training data are **provably vacuous** for overparameterized model classes; unable to explain generalization.

$$L(\theta) \leq \hat{L}(D,\theta) + \mathcal{O}(\sqrt{\frac{p}{n}})$$

Why current machine learning techniques find overparameterized interpolators with strong generalization performance is an open question.

Contributions

A perfectly tight distribution-dependent PAC-Chernoff bound for interpolators, even in over-parameterized models.

$$L(\theta) \leq \hat{L}(D,\theta) + C_{\nu,\theta}(\frac{p}{n})$$

where ν is the data-generating distribution.

Contributions

A perfectly tight distribution-dependent PAC-Chernoff bound for interpolators, even in over-parameterized models.

$$L(\theta) \leq \hat{L}(D, \theta) + C_{\nu, \theta}(\frac{p}{n})$$

where ν is the data-generating distribution.

A theoretical framework that explains why some interpolators generalize well, while others do not, based on a novel characterization of smoothness.

Contributions

A perfectly tight distribution-dependent PAC-Chernoff bound for interpolators, even in over-parameterized models.

$$L(\theta) \leq \hat{L}(D, \theta) + C_{\nu, \theta}(\frac{p}{n})$$

where ν is the data-generating distribution.

A **theoretical framework** that explains why some interpolators generalize well, while others do not, based on a novel **characterization of smoothness**.

We explain why regularization, data augmentation, invariant architectures, and over-parameterization produce smoother interpolators with superior generalization.

The Rate Function

Chernoff's Theorem. For any fixed $\theta \in \Theta$ and a > 0, it satisfies

$$\mathbb{P}_{D \sim \nu^n} \Big(L(\theta) - \hat{L}(D, \theta) \ge a \Big) \le e^{-n\mathcal{I}_{\theta}(a)}$$
.

with

$$\mathcal{I}_{m{ heta}}(a) := \sup_{\lambda > 0} \ \lambda a - J_{m{ heta}}(\lambda) \quad ext{ and } \quad J_{m{ heta}}(\lambda) := \ln \mathbb{E}_{
u} \Big[e^{\lambda(L(m{ heta}) - \ell(m{y}, m{x}, m{ heta}))} \Big] \,,$$

- Cramér's Theorem: For large *n*, the bound is tight
- **Proposition 3.4**: When *a* is large, the bound is tight.

Smoother Interpolators Generalize Better

Smoothness: A model $\theta \in \Theta$ is β -smoother than a model $\theta' \in \Theta'$

$$\forall a \in (0, \beta] \quad \mathcal{I}_{\theta}(a) \geq \mathcal{I}_{\theta'}(a)$$
.

Theorem 4.6. For any $\epsilon \geq 0$ with h.p. $1-\delta$, for all $\theta \in \Theta$, $\theta' \in \Theta'$, simultaneously,

$$\hat{L}(D, m{ heta}) \leq \epsilon$$
 and $m{ heta}$ is eta -smoother than $m{ heta}'$ \Downarrow
$$L(m{ heta}) \leq L(m{ heta}') + \epsilon$$

PAC-Chernoff Bound

Theorem 4.1. With h.p. $1 - \delta$, for all $\theta \in \Theta$, simultaneously,

$$L(\theta) \leq \hat{L}(D,\theta) + \mathcal{I}_{\theta}^{-1}(\frac{1}{n}\ln\frac{k^p}{\delta})$$
.

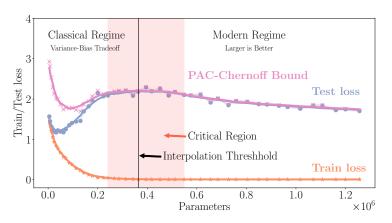
with

$$\mathcal{I}_{ heta}^{-1}(s) := \inf_{\lambda>0} rac{J_{oldsymbol{ heta}}(\lambda) + s}{\lambda} \quad orall s \geq 0 \, .$$

Tightness

Proposition 4.4. The bound is **perfectly tight** for interpolators.

$$L(\theta) \leq \hat{L}(D,\theta) + \mathcal{I}_{\theta}^{-1}(\frac{1}{n}\ln\frac{k^p}{\delta}) \leq L(\theta) + \hat{L}(D,\theta)$$
.



Optimal Regularization

• The inverse rate is an **regularizer** towards **smoother** models.

$$\boldsymbol{\theta}_{\epsilon}^{\times} = \mathop{\arg\min}_{\boldsymbol{\theta} \,:\, \hat{L}(D,\boldsymbol{\theta}) \,\leq\, \epsilon} \, \hat{L}(D,\boldsymbol{\theta}) + \underbrace{\mathcal{I}_{\boldsymbol{\theta}}^{-1}\big(\frac{1}{n}\ln\frac{k^p}{\delta}\big)}_{\text{Regularizer}},$$

$$oldsymbol{ heta}^\star_\epsilon = \mathop{\mathrm{arg\,min}}_{oldsymbol{ heta}: \, \hat{L}(D, oldsymbol{ heta}) \, \leq \, \epsilon} L(oldsymbol{ heta}) \, .$$

• How close is $\theta_{\epsilon}^{\times}$ from the best possible interpolator $\theta_{\epsilon}^{\star}$.

Optimal Regularization

The inverse rate is an regularizer towards smoother models.

$$\boldsymbol{\theta}_{\epsilon}^{\times} = \underset{\boldsymbol{\theta} \,:\, \hat{\boldsymbol{L}}(\boldsymbol{D},\boldsymbol{\theta}) \,\leq\, \epsilon}{\arg\min} \, \, \hat{\boldsymbol{L}}(\boldsymbol{D},\boldsymbol{\theta}) + \underbrace{\mathcal{I}_{\boldsymbol{\theta}}^{-1}\big(\frac{1}{n}\ln\frac{k^p}{\delta}\big)}_{\text{Regularizer}},$$

$$oldsymbol{ heta}^\star_\epsilon = \mathop{\mathrm{arg\,min}}_{oldsymbol{ heta}: \, \hat{L}(D,oldsymbol{ heta}) \, \leq \, \epsilon} L(oldsymbol{ heta}) \, .$$

- ullet How close is $m{ heta}_{\epsilon}^{ imes}$ from the best possible interpolator $m{ heta}_{\epsilon}^{\star}$.
- Very close!! Theorem 6.1 For any $\epsilon>0$, with h.p. $1-\delta$ over $D\sim \nu^n$

$$|L(\boldsymbol{\theta}_{\epsilon}^{\star}) - L(\boldsymbol{\theta}_{\epsilon}^{\times})| \leq \epsilon.$$

The inverse rate is an optimal regularizer.

Understanding Existing Regularizers

Many common regularization techniques are **approximations** to the **optimal regularizer**:

• Distance from initialization and ℓ_2 -norm:

$$\mathcal{I}_{m{ heta}}^{-1}ig(rac{1}{n}\lnrac{k^p}{\delta}ig) \leq \sqrt{2Ma} \; \|m{ heta}\|_2 \, ,$$

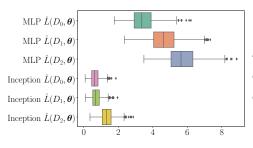
• Input-gradient norm:

$$\mathcal{I}_{\boldsymbol{\theta}}^{-1}\big(\tfrac{1}{n}\ln\tfrac{k^p}{\delta}\big) \leq \sqrt{\tfrac{1}{n}\ln\tfrac{k^p}{\delta}}\sqrt{M\mathbb{E}_{\nu}\big[\big\|\nabla_{\mathbf{x}}\ell(\boldsymbol{y},\boldsymbol{x},\boldsymbol{\theta})\big\|_2^2\big]}\,.$$

Transformed Input Data

Input data in many machine learning problems undergo **transformations**, often due to the measuring process, such as sensor noise or image distortions.

Transformed input-data makes the expected loss $L(\theta)$ higher and the distribution of $\hat{L}(D,\theta)$ with $D \sim \nu^n$ less concentrated.

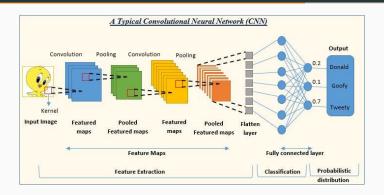


 D_0 - CIFAR-10's test set.

 D_1 - random translations.

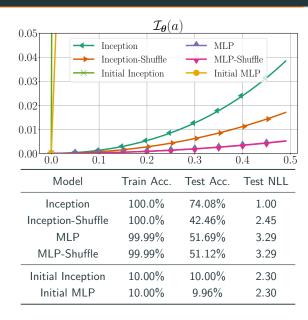
 D_2 - random rotations.

Invariant Architectures (I/II)



- PAC-Chernoff bounds explain why interpolating with invariant architecture leads to better generalization performance.
- The $\hat{L}(D,\theta)$ of invariant architectures is more concentrated under transformed inputs.

Invariant Architectures (II/II)



Over-parameterization

Modern machine learning models are highly overparametrized.

Previous works have established links between overparametrization and generalization, but under very limited settings.

The distribution-dependent PAC-Chernoff Bound can be used to obtain **bounds over the number of parameters** of **interpolators**:

Theorem 8.1. For any $\epsilon \in (0, L^*)$ and any $\delta \in (0, 1)$, with high probability $1 - \delta$ over $D \sim \nu^n$, for all $\theta \in \Theta$, simultaneously,

if
$$\hat{L}(D, \theta) \le \epsilon$$
 then $p \ge \frac{n\mathcal{I}_{\theta}(L^{\star} - \epsilon) + \ln \delta}{\ln k}$.

where $L^* = \arg\min_{\theta} L(\theta)$.

Conclusions and Limitations (I/II)

- Traditional bounds relying solely on training data are unable to explain generalization of over-parameterized interpolators.
- Distribution-dependent PAC-Chernoff bounds are a promising tool able to explain a wide range of learning techniques.
- Smoother interpolators generalize better.

Conclusions and Limitations (II/II)

- Connected to a wide range of regularization methods.
- Explain why invariant architectures and data-augmentation works under transformed input-data.
- Over-parameterization is a neccessary condition for smooth interpolation.
- Limitation: Assumption of a finite model class. It can be addressed by using PAC-Bayes Chernoff bounds (Casado et al. 2024).