

A PAC-BAYESIAN LINK BETWEEN GENERALISATION AND FLAT MINIMA

JOINT WORK WITH PAUL VIALARD, U MUT SIMSEKLI AND BENJAMIN GUEDJ

Maxime Haddouche

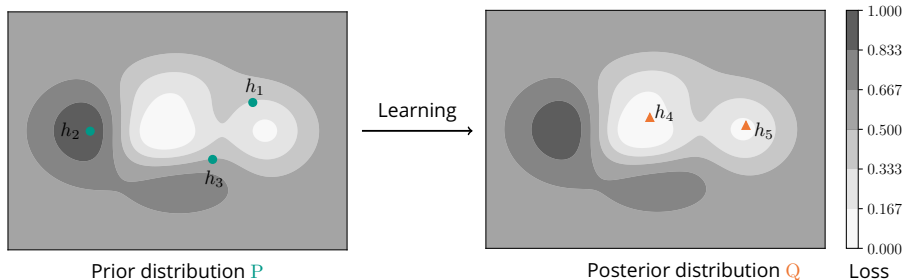
INRIA London
Université de Lille

Session MALIA JdS 2024

(SPECIAL CASE OF) PAC-BAYESIAN LEARNING

PAC-Bayesian learning

Learning a distribution Q over models from the data and a prior distribution P



PAC-Bayesian generalisation bounds in a nutshell

With probability at least $1 - \delta$

$$\text{performance gap}(Q) \leq \text{bound}\left(\text{complexity}(Q, P), \frac{1}{m}, \ln \frac{1}{\delta}\right)$$

Notations:

- Predictor/hypothesis $h \in \mathcal{H}$, Data space \mathcal{Z}
- Loss $\ell : \mathcal{H} \times \mathcal{Z} \rightarrow \mathbb{R}^+$, possibly heavy-tailed
- m -sized *i.i.d.* learning sample $\mathcal{S} \in \mathcal{Z}^m$, $\mathcal{S} := \{\mathbf{z}_i\}_{i=1}^m \sim \mathcal{D}^{\otimes m}$
- Population risk $R_{\mathcal{D}}(h) = \mathbb{E}_{\mathbf{z} \sim \mathcal{D}} \ell(h, \mathbf{z})$ and empirical risk $\hat{R}_{\mathcal{S}}(h) = \frac{1}{m} \sum_{i=1}^m \ell(h, \mathbf{z}_i)$
- Expected risks $R_{\mathcal{D}}(\mathbb{Q}) = \mathbb{E}_{h \sim \mathbb{Q}} [R_{\mathcal{D}}(h)]$, $\hat{R}_{\mathcal{S}}(\mathbb{Q}) = \mathbb{E}_{h \sim \mathbb{Q}} [\hat{R}_{\mathcal{S}}(h)]$
- Space of distributions over \mathcal{H} : $\mathcal{M}(\mathcal{H})$

Catoni's bound Alquier et al. (2016, Theorem 4.1) (σ -subgaussian losses)

For $\lambda > 0$, with probability $1 - \delta$ over $\mathcal{S} \sim \mathcal{D}^m$, for any $\mathbb{Q} \in \mathcal{M}(\mathcal{H})$,

$$R_{\mathcal{D}}(\mathbb{Q}) \leq \hat{R}_{\mathcal{S}}(\mathbb{Q}) + \frac{\text{KL}(\mathbb{Q}, \mathbb{P}) + \ln \frac{1}{\delta}}{\lambda} + \frac{\lambda \sigma^2}{2m}$$

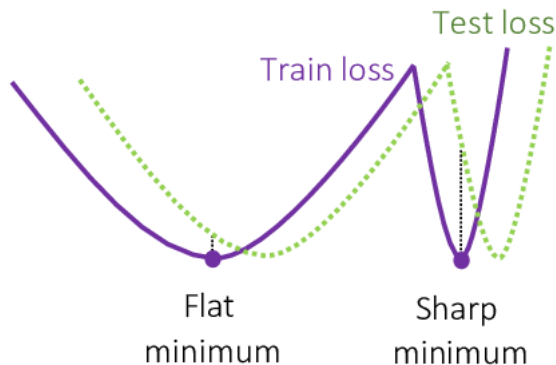
FLAT MINIMUM

What is a flat minimum?

FLAT MINIMUM

What is a flat minimum?

A minimum such that its neighbourhood nearly minimises the loss.



FLAT MINIMA AND GENERALISATION ARE CORRELATED!

Correlations with generalisation recently emerged:

- Flat minima of \hat{R}_S .
PAC-Bayes based correlation measure : works for many datasets (Neyshabur *et al.*, 2017; Dziugaite *et al.*, 2020; Jiang *et al.*, 2020)
- Flat minima of the adversarial loss in the context of adversarially robust learning. (Stutz *et al.*, 2021)
- Flat minima implies generalisation for 2-layers nets (Wen *et al.*, 2023).

FLAT MINIMA AND GENERALISATION ARE CORRELATED!

Correlations with generalisation recently emerged:

- Flat minima of \hat{R}_S .
PAC-Bayes based correlation measure : works for many datasets (Neyshabur *et al.*, 2017; Dziugaite *et al.*, 2020; Jiang *et al.*, 2020)
- Flat minima of the adversarial loss in the context of adversarially robust learning. (Stutz *et al.*, 2021)
- Flat minima implies generalisation for 2-layers nets (Wen *et al.*, 2023).

Can we go beyond correlation or 2-layers net and obtain sound generalisation bounds involving directly flat minima?

ESSENTIAL TOOLS: POINCARÉ AND LOG-SOBOLEV INEQUALITIES

Notation: for any Q , $H^1(Q) := \{f \in L^2(Q) \cap D_1(\mathbb{R}^d) \mid \|\nabla f\| \in L^2(Q)\}$

Poincaré

Q is Poinc(c_P) if for all $f \in H^1(Q)$:

$$\text{Var}(f) \leq c_P(Q) \mathbb{E}_{h \sim Q} [\|\nabla f(h)\|^2],$$

Log-Sobolev

Q is L-Sob(c_{LS}) if for all function $f \in H^1(Q)$:

$$\mathbb{E}_{h \sim Q} \left[f^2(h) \log \left(\frac{f^2(h)}{\mathbb{E}_{h \sim Q} [f^2(h)]} \right) \right] \leq c_{LS}(Q) \mathbb{E}_{h \sim Q} [\|\nabla f(h)\|^2],$$

ESSENTIAL TOOLS: POINCARÉ AND LOG-SOBOLEV INEQUALITIES

Notation: for any Q , $H^1(Q) := \{f \in L^2(Q) \cap D_1(\mathbb{R}^d) \mid \|\nabla f\| \in L^2(Q)\}$

Poincaré

Q is $\text{Poinc}(c_P)$ if for all $f \in H^1(Q)$:

$$\text{Var}(f) \leq c_P(Q) \mathbb{E}_{h \sim Q} [\|\nabla f(h)\|^2],$$

Log-Sobolev

Q is $L\text{-Sob}(c_{LS})$ if for all function $f \in H^1(Q)$:

$$\mathbb{E}_{h \sim Q} \left[f^2(h) \log \left(\frac{f^2(h)}{\mathbb{E}_{h \sim Q} [f^2(h)]} \right) \right] \leq c_{LS}(Q) \mathbb{E}_{h \sim Q} [\|\nabla f(h)\|^2],$$

Gaussian distributions and Gibbs posteriors are Poinc and L-Sob!

FAST-RATE GENERALISATION BOUNDS FOR FLAT MINIMA (1)

Notation: $\text{Err}(\ell, \mathbb{Q}, \mathbf{z}) := \mathbb{E}_{h \sim \mathbb{Q}}[\ell(h, \mathbf{z})]$

Assumption

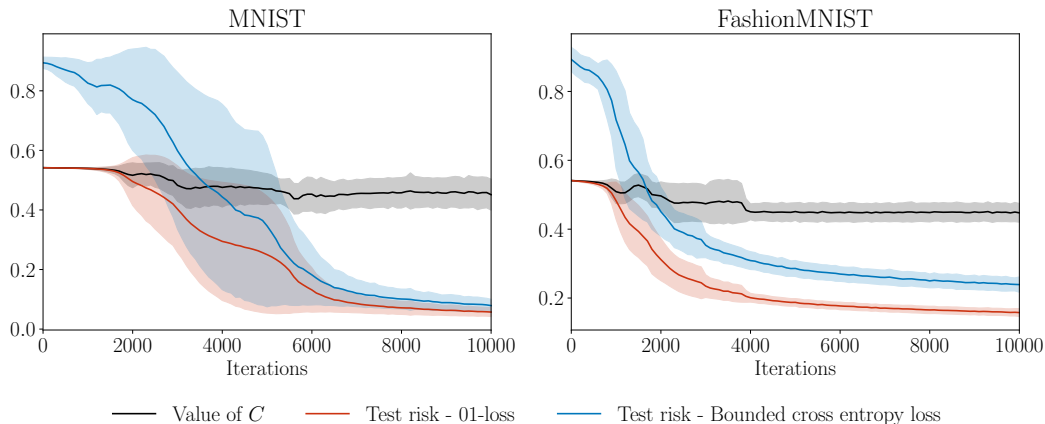
$\mathbb{Q} \in \mathcal{M}(\mathcal{H})$ is *quadratically self-bounded w.r.t. ℓ and $C > 0$* (namely $\text{QSB}(\ell, C)$) if

$$\mathbb{E}_{\mathbf{z} \sim \mathcal{D}} [\text{Err}(\ell, \mathbb{Q}, \mathbf{z})^2] \leq C R_{\mathcal{D}}(\mathbb{Q}) (= C \mathbb{E}_{\mathbf{z} \sim \mathcal{D}} [\text{Err}(\ell, \mathbb{Q}, \mathbf{z})])$$

- QSB intricates $\mathcal{D} \in \mathcal{M}(\mathcal{Z})$ with $\mathbb{Q} \in \mathcal{M}(\mathcal{H})$
- Satisfied if $\ell \in [0, K]$ with $C = K$.
- QSB quantifies the 'flatness' of the post-training minima reached by \mathbb{Q} .

IS THE QSB ASSUMPTION VERIFIED IN PRACTICE?

QSB holds for 3-layer neural nets trained on MNIST (black curve)!



FAST-RATE GENERALISATION BOUNDS VIA FLAT MINIMA (2)

Theorem

For any $C > 0$, data-free prior \mathbf{P} , with probability at least $1 - \delta$ for any $m > 0$, and \mathbf{Q} being $\text{Poinc}(c_P)$, $\text{QSB}(\ell, C)$,

$$R_{\mathcal{D}}(\mathbf{Q}) \leq 2\hat{R}_{\mathcal{S}}(\mathbf{Q}) + 2C \frac{KL(\mathbf{Q}, \mathbf{P}) + \log(1/\delta)}{m} + \frac{1}{C} c_P(\mathbf{Q}) \mathbb{E}_{\mathbf{z} \sim \mathcal{D}} \left[\mathbb{E}_{h \sim \mathbf{Q}} (\|\nabla_h \ell(h, \mathbf{z})\|^2) \right].$$

FAST-RATE GENERALISATION BOUNDS VIA FLAT MINIMA (2)

Theorem

For any $C > 0$, data-free prior \mathbf{P} , with probability at least $1 - \delta$ for any $m > 0$, and \mathbf{Q} being $\text{Poinc}(c_P)$, $\text{QSB}(\ell, C)$,

$$R_{\mathcal{D}}(\mathbf{Q}) \leq 2\hat{R}_S(\mathbf{Q}) + 2C \frac{KL(\mathbf{Q}, \mathbf{P}) + \log(1/\delta)}{m} + \frac{1}{C} c_P(\mathbf{Q}) \mathbb{E}_{\mathbf{z} \sim \mathcal{D}} \left[\mathbb{E}_{h \sim \mathbf{Q}} (\|\nabla_h \ell(h, \mathbf{z})\|^2) \right].$$

If \mathcal{D} is also Poinc :

With more minor technical assumptions, for any \mathbf{Q} being $\text{Poinc}(c_P)$ with $R_{\mathcal{D}}(\mathbf{Q}) \leq C$:

$$R_{\mathcal{D}}(\mathbf{Q}) \leq 2\hat{R}_S(\mathbf{Q}) + 2C \frac{KL(\mathbf{Q}, \mathbf{P}) + \log(1/\delta)}{m} + \frac{1}{C} \left(c_P(\mathbf{Q}) \mathbb{E}_{\mathbf{z} \sim \mathcal{D}} \left[\mathbb{E}_{h \sim \mathbf{Q}} (\|\nabla_h \ell(h, \mathbf{z})\|^2) \right] + c_P(\mathcal{D}) \mathbb{E}_{\mathbf{z} \sim \mathcal{D}} \left(\left\| \mathbb{E}_{h \sim \mathbf{Q}} [\nabla_z \ell(h, \mathbf{z})] \right\|^2 \right) \right).$$

FULLY EMPIRICAL FAST RATE

Current drawback: bounds are not empirical.

Current drawback: bounds are not empirical.

Solution: \mathcal{C}^2 gradient-lipschitz losses!

Theorem

For any $C_1, C_2, c > 0$, with probability at least $1 - \delta$, for any $m > 0$, \mathbb{Q} being Poinc(c_P) with constant c , $\text{QSB}(\ell, C_1)$, $\text{QSB}(\|\nabla_h \ell\|^2, C_2)$,

$$R_{\mathcal{D}}(\mathbb{Q}) \leq 2\hat{R}_{\mathcal{S}}(\mathbb{Q}) + \mathcal{O}\left(\mathbb{E}_{h \sim \mathbb{Q}} \left[\frac{1}{m} \sum_{i=1}^m \|\nabla_h \ell(h, \mathbf{z}_i)\|^2 \right] + \frac{\text{KL}(\mathbb{Q}, \mathbb{P}) + \log(1/\delta)}{m}\right).$$

If Q satisfies either

1 Flat minima for \hat{R}_S and R_D ,

2 if ℓ gradient-lipschitz, flat minima for \hat{R}_S and empirical gradient norms,

then Q generalises well!

Current limitation: with Poincaré posteriors, KL is uncontrolled.

Current limitation: with Poincaré posteriors, KL is uncontrolled.

Solution: consider Gibbs posterior with log-Sobolev priors!

Definition

$\mathbb{P}_{-\gamma\hat{R}_S}$ is the Gibbs posterior *w.r.t.* prior \mathbb{P} with *inverse temperature* $\gamma > 0$ if

$$d\mathbb{P}_{-\gamma\hat{R}_S}(h) \propto \exp\left(-\gamma\hat{R}_S(h)\right) d\mathbb{P}(h)$$

.

Why focus on those?

- Minimise Catoni's bound (Alquier *et al.*, 2016, Theorem 4.1)
- if \mathbb{P} L-Sob(+ technical assumptions) and $\ell = \ell_1 + \ell_2$ (ℓ_1 convex, twice differentiable, ℓ_2 bounded) then $\mathbb{P}_{-\gamma\hat{R}_S}$ is L-Sob.

UNDERSTANDING GIBBS POSTERIOBS THROUGH FLAT MINIMA

Theorem

For any $C > 0$, any $\gamma > 0$, any prior \mathbf{P} L-Sob(c_{LS}) (+ technical assumptions), if $\ell = \ell_1 + \ell_2$ (as above), then with probability at least $1 - \delta$, for any $m > 0$, \mathbf{Q} being QSB(ℓ, C):

$$\mathbf{R}_{\mathcal{D}}(\mathbf{P}_{-\gamma\hat{\mathbf{R}}_S}) \leq 2\hat{\mathbf{R}}_S(\mathbf{P}_{-\gamma\hat{\mathbf{R}}_S}) + \mathcal{O}\left(C \frac{\gamma^2 \mathbb{E}_{h \sim \mathbf{P}_{-\gamma\hat{\mathbf{R}}_S}} [\|\nabla_h \hat{\mathbf{R}}_S(h)\|^2]}{m} + \log(1/\delta) + \frac{1}{C} \mathbb{E}_{\mathbf{z} \sim \mathcal{D}} \left[\mathbb{E}_{h \sim \mathbf{P}_{-\gamma\hat{\mathbf{R}}_S}} (\|\nabla_h \ell(h, \mathbf{z})\|^2) \right]\right).$$

UNDERSTANDING GIBBS POSTERIOBS THROUGH FLAT MINIMA

Theorem

For any $C > 0$, any $\gamma > 0$, any prior \mathbb{P} L-Sob(c_{LS}) (+ technical assumptions), if $\ell = \ell_1 + \ell_2$ (as above), then with probability at least $1 - \delta$, for any $m > 0$, \mathbb{Q} being QSB(ℓ, C):

$$R_{\mathcal{D}}(\mathbb{P}_{-\gamma\hat{R}_S}) \leq 2\hat{R}_S(\mathbb{P}_{-\gamma\hat{R}_S}) + \mathcal{O}\left(C \frac{\gamma^2 \mathbb{E}_{h \sim \mathbb{P}_{-\gamma\hat{R}_S}} [\|\nabla_h \hat{R}_S(h)\|^2]}{m} + \log(1/\delta) + \frac{1}{C} \mathbb{E}_{\mathbf{z} \sim \mathcal{D}} \left[\mathbb{E}_{h \sim \mathbb{P}_{-\gamma\hat{R}_S}} (\|\nabla_h \ell(h, \mathbf{z})\|^2) \right]\right).$$

KL small if a flat minima on \hat{R}_S is reached:

→ **Flat minima fully explain generalisation here!**

- 1 Gibbs posterior generalises well if they reach a flat minima on both $\hat{R}_{\mathcal{S}}$ and $R_{\mathcal{D}}$.
- 2 Flatness of the minimum on $\hat{R}_{\mathcal{S}}$ controls the expansion of KL.

Drawback: results hold for probabilistic predictors

Drawback: results hold for probabilistic predictors

Answer: Exploit the 2-Wasserstein distance to obtain guarantees valid for deterministic predictors (Diracs)

CONVERGENCE GUARANTEES FOR NON-CONVEX SGD

Key tool: a novel change of measure inequality

For any f gradient lipschitz, any P, Q :

$$\mathbb{E}_{h \sim Q}[f(h)] \leq \frac{G}{2} W_2^2(Q, P) + \mathbb{E}_{h \sim P}[f(h)] + D \mathbb{E}_{h \sim Q}[\|\nabla f(h)\|].$$

NB: a variant of this formula with a KL is attainable if $Q \ll P$ and P is L-Sob !

Assumption

- Gradient-lipschitz loss.
- $P \propto \exp(-V(h))dh$

Theorem

Let $\delta \in (0, 1)$ and $P \in \mathcal{M}(\mathcal{H})$ a data-free prior. Assume \mathcal{H} has a finite diameter $D > 0$, $\ell \geq 0$ and that for any m , the generalisation gap Δ_{S_m} is G gradient-Lipschitz. Assume that $\mathbb{E}_{h \sim P} \mathbb{E}_{z \sim \mathcal{D}} [\ell(h, z)^2] \leq \sigma^2$, then the following holds with probability at least $1 - \delta$, for any $m > 0$ and any \mathbb{Q} :

$$R_D(\mathbb{Q}) \leq \hat{R}_{S_m}(\mathbb{Q}) + \frac{G}{2} W_2^2(\mathbb{Q}, P) + \sqrt{\frac{2\sigma^2 \log(\frac{1}{\delta})}{m}} + D \mathbb{E}_{h \sim \mathbb{Q}} \left(\left\| \nabla_h R_{\mathcal{D}}(h) - \nabla_h \hat{R}_{S_m}(h) \right\| \right)$$

CONCLUSION

- We mathematically quantify the impact of flat minima in generalisation: momentum in Catoni's bound!
- The QSB condition is verified on basic neural nets (classification) with constant C sharper than 1!
- A crucial future lead: understanding why optimisation procedures on deep nets lead to flat minima: **here we are only able to explain why flat minima generalise well, not how we reach them.**

Full paper available at <https://arxiv.org/abs/2402.08508>

REFERENCES I

-  Pierre Alquier, James Ridgway, and Nicolas Chopin. On the properties of variational approximations of Gibbs posteriors. *Journal of Machine Learning Research*. (2016).
-  Behnam Neyshabur, Srinadh Bhojanapalli, David McAllester, and Nati Srebro. Exploring Generalization in Deep Learning. *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*. (2017). URL: <https://proceedings.neurips.cc/paper/2017/hash/10ce03a1ed01077e3e289f3e53c72813-Abstract.html>.
-  Gintare Karolina Dziugaite, Alexandre Drouin, Brady Neal, Nitarshan Rajkumar, Ethan Caballero, Linbo Wang, Ioannis Mitliagkas, and Daniel M. Roy. In search of robust measures of generalization. *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*. (2020). URL: <https://proceedings.neurips.cc/paper/2020/hash/86d7c8a08b4aaa1bc7c599473f5dddda-Abstract.html>.
-  Yiding Jiang, Behnam Neyshabur, Hossein Mobahi, Dilip Krishnan, and Samy Bengio. Fantastic Generalization Measures and Where to Find Them. *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. (2020). URL: <https://openreview.net/forum?id=SJgIPJBFvH>.

REFERENCES II



Lucas Liebenwein, Ramin Hasani, Alexander Amini, and Daniela Rus. Sparse flows: Pruning continuous-depth models. *Advances in Neural Information Processing Systems*. 34. (2021).



David Stutz, Matthias Hein, and Bernt Schiele. Relating Adversarially Robust Generalization to Flat Minima. *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021. IEEE*. (2021). DOI: 10.1109/ICCV48922.2021.00771. URL: <https://doi.org/10.1109/ICCV48922.2021.00771>.



Kaiyue Wen, Zhiyuan Li, and Tengyu Ma. Sharpness Minimization Algorithms Do Not Only Minimize Sharpness To Achieve Better Generalization. *CoRR*. abs/2307.11007. (2023). DOI: 10.48550/ARXIV.2307.11007. arXiv: 2307.11007. URL: <https://doi.org/10.48550/arXiv.2307.11007>.

Thank you for your attention!