

# AN INTRODUCTION TO PAC-BAYES LEARNING AND ITS LINKS TO FLAT MINIMA

SÉMINAIRE SO

**Maxime Haddouche**

Inria Paris

**Mardi 4 Février 2025**

1. Introduction to PAC-Bayes Learning
2. PAC-Bayes with Weak Statistical Assumptions
3. Involving Flat Minima in PAC-Bayes

### What is a learning theory problem?

A tuple  $(\mathcal{Z}, \mathcal{H}, \ell)$ : a data space  $\mathcal{Z}$ , a predictor space  $h \in \mathcal{H}$ , a mathematically well-defined problem  $\ell : \mathcal{H} \times \mathcal{Z} \rightarrow \mathbb{R}$

### What is a learning theory problem?

A tuple  $(\mathcal{Z}, \mathcal{H}, \ell)$ : a data space  $\mathcal{Z}$ , a predictor space  $h \in \mathcal{H}$ , a mathematically well-defined problem  $\ell : \mathcal{H} \times \mathcal{Z} \rightarrow \mathbb{R}$

### What is our goal?

We have access to a  $m$ -sized training set  $\mathcal{S}_m = \{\mathbf{z}_1, \dots, \mathbf{z}_m\}$ . We aim to learn the best  $h^* \in \mathcal{H}$  to answer  $\ell$  in a certain way

## What is a learning theory problem?

A tuple  $(\mathcal{Z}, \mathcal{H}, \ell)$ : a data space  $\mathcal{Z}$ , a predictor space  $h \in \mathcal{H}$ , a mathematically well-defined problem  $\ell : \mathcal{H} \times \mathcal{Z} \rightarrow \mathbb{R}$

## What is our goal?

We have access to a  $m$ -sized training set  $\mathcal{S}_m = \{\mathbf{z}_1, \dots, \mathbf{z}_m\}$ . We aim to learn the best  $h^* \in \mathcal{H}$  to answer  $\ell$  in a certain way

- **Optimisation:** minimise the empirical risk  
$$h^* \in \operatorname{argmin}_{h \in \mathcal{H}} \hat{R}_{\mathcal{S}_m}(h) := \frac{1}{m} \sum_{i=1}^m \ell(h, \mathbf{z}_i)$$
- **Generalisation:** if  $\mathcal{S}_m \sim \mathcal{D}^{\otimes m}$ , minimise the theoretical risk  
$$h^* \in \operatorname{argmin}_{h \in \mathcal{H}} R_{\mathcal{D}}(h) := \mathbb{E}_{\mathbf{z} \sim \mathcal{D}}[\ell(h, \mathbf{z})]$$

### Supervised learning with linear classifiers:

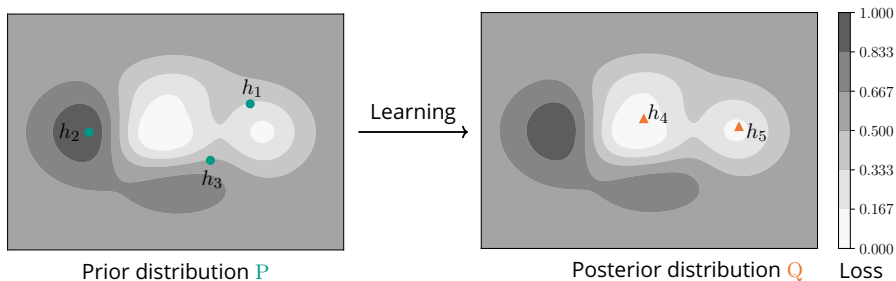
- $\mathcal{Z} = \mathbb{R}^k \times \mathcal{Y}$  with  $\mathcal{Y} = \{-1, 1\}$
- Loss  $\ell(h, (x, y)) = \mathbb{1}\{h(x) \neq y\}$
- Linear classifiers:  $\mathcal{H} := \{h_\theta(x) = \text{sgn}(\langle \theta, x \rangle)\}$ , where  $\text{sgn}(a)$  denotes the sign of  $a$ .

**It may be hard to find directly the best  $h$  for complex predictor classes (eg neural nets). What could we do?**

# PAC-BAYES LEARNING

**PAC-Bayes: Find the best distribution over  $\mathcal{H}$  !**

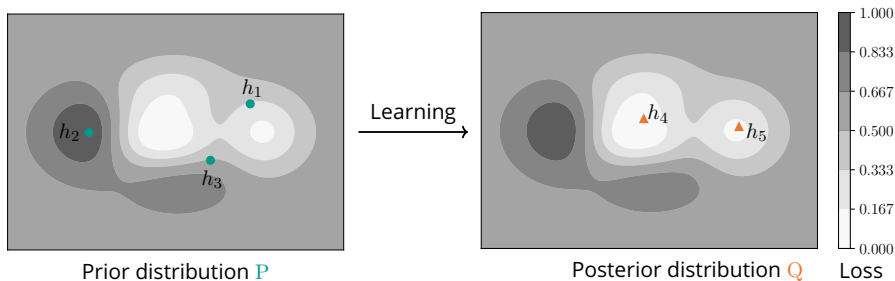
Learning a posterior  $Q$  over models from  $m$  data and a prior distribution  $P$



# PAC-BAYES LEARNING

## PAC-Bayes: Find the best distribution over $\mathcal{H}$ !

Learning a posterior  $\mathcal{Q}$  over models from  $m$  data and a prior distribution  $\mathcal{P}$



## PAC-Bayesian generalisation bounds in a nutshell

With probability at least  $1 - \delta$

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



**Notations:**

- Predictor/hypothesis  $h \in \mathcal{H}$ , Data space  $\mathcal{Z}$
- Loss  $\ell : \mathcal{H} \times \mathcal{Z} \rightarrow \mathbb{R}$ ,
- Countable learning sample  $\mathcal{S} = (\mathbf{z}_i)_{i \geq 1} \in \mathcal{Z}^{\mathbb{N}}$ , with distribution  $\mathcal{D}_{\mathcal{S}}$
- $\mathcal{S}_m$ : Restriction of  $\mathcal{S}$  to  $m$  first points with distribution  $\mathcal{D}_m$
- Space of distributions over  $\mathcal{H}$ :  $\mathcal{M}(\mathcal{H})$
- Posterior and prior distribution  $\mathbb{Q}, \mathbb{P} \in \mathcal{M}(\mathcal{H})^2$
- If  $\mathcal{S}_m \sim \mathcal{D}^m$  i.i.d., Risks:  $R_{\mathcal{D}}(h) = \mathbb{E}_{\mathbf{z} \sim \mathcal{D}} \ell(h, \mathbf{z})$ ,  $\hat{R}_{\mathcal{S}_m}(h) = \frac{1}{m} \sum_{i=1}^m \ell(h, \mathbf{z}_i)$
- Expected risks  $R_{\mathcal{D}}(\mathbb{Q}) = \mathbb{E}_{h \sim \mathbb{Q}} [R_{\mathcal{D}}(h)]$ ,  $\hat{R}_{\mathcal{S}_m}(\mathbb{Q}) = \mathbb{E}_{h \sim \mathbb{Q}} [\hat{R}_{\mathcal{S}_m}(h)]$

## A FUNDAMENTAL RESULT: MCALLESTER'S BOUND

**McAllester's bound (Maurer's improvement) Maurer (2004, Theorem 5)** ( $\ell \in [0, 1]$ )

For any  $\mathbf{P} \in \mathcal{M}(\mathcal{H})$ , with probability  $1-\delta$  over  $\mathcal{S}_m \sim \mathcal{D}^m$ , for any  $\mathbf{Q} \in \mathcal{M}(\mathcal{H})$ ,

$$R_{\mathcal{D}}(\mathbf{Q}) \leq \hat{R}_{\mathcal{S}_m}(\mathbf{Q}) + \sqrt{\frac{\text{KL}(\mathbf{Q}, \mathbf{P}) + \ln \frac{2\sqrt{m}}{\delta}}{2m}},$$

where  $\text{KL}(\mathbf{Q}, \mathbf{P}) = \mathbb{E}_{h \sim \mathbf{Q}} \left[ \frac{d\mathbf{Q}}{d\mathbf{P}}(h) \right]$ .

**No explicit dependency in the dimension of the problem** (hidden in the KL term):  
positive phenomenon can be caught with the right priors (*e.g.* sparsity).

## Step 1: A key ingredient: change of measure inequality

For any function  $f$ , any  $Q \ll P$ :

$$\mathbb{E}_{h \sim Q} [f(h)] - \ln \left( \mathbb{E}_{h \sim P} [\exp \circ f(h)] \right) \leq \text{KL}(Q, P).$$

# A SIMPLE ROUTE OF PROOF

## Step 1: A key ingredient: change of measure inequality

For any function  $f$ , any  $Q \ll P$ :

$$\mathbb{E}_{h \sim Q} [f(h)] - \ln \left( \mathbb{E}_{h \sim P} [\exp \circ f(h)] \right) \leq \text{KL}(Q, P).$$

## Step 2: Markov's inequality

With probability at least  $1 - \delta$ :

$$\begin{aligned} \mathbb{E}_{h \sim P} [\exp \circ f(h)] &\leq \frac{1}{\delta} \mathbb{E}_{S_m} \left[ \mathbb{E}_{h \sim P} [\exp \circ f(h)] \right], \\ &= \frac{1}{\delta} \mathbb{E}_{h \sim P} \left[ \mathbb{E}_{S_m} [\exp \circ f(h)] \right]. \end{aligned} \quad (\text{P data-free + Fubini})$$

## A SIMPLE ROUTE OF PROOF (2)

### Step 3: Choosing the right $f$ .

Take  $f(h) = m \text{kl} \left( R_{\mathcal{D}}(h), \hat{R}_{\mathcal{S}_m}(h) \right)$  (kl= KL of Bernoullis).

Then Maurer (2004): for any  $h$ , loss in  $[0, 1]$ :

$$\mathbb{E}_{\mathcal{S}_m} [\exp \circ f(h)] \leq 2\sqrt{m}$$

To conclude:  $\text{kl}(p, q) \geq 2(p - q)^2$ .

## A SECOND KEY RESULT: CATONI'S BOUND

**Catoni's bound Alquier *et al.* (2016, Theorem 4.1)** ( $\ell$   $\sigma$ -subgaussian)

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

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

**Previous bounds:** both fully empirical  $\rightarrow$  optimisation in  $Q$  is feasible on  $\mathcal{C} \subseteq \mathcal{M}(\mathcal{H})$  !

$$\text{McAllester} \quad Q_M := \operatorname{argmin}_{Q \in \mathcal{C}} \hat{R}_{S_m}(Q) + \sqrt{\frac{\text{KL}(Q, P)}{2m}}.$$

For any  $\lambda > 0$ ,

$$\text{Catoni} \quad Q_C := \operatorname{argmin}_{Q \in \mathcal{C}} \hat{R}_{S_m}(Q) + \frac{\text{KL}(Q, P)}{\lambda}.$$

If  $\mathcal{C} = \mathcal{M}(\mathcal{H})$ , a *Gibbs posterior*  $P_{-\lambda \hat{R}_{S_m}}$  is the explicit minimiser of Catoni's bound:

$$dP_{-\lambda \hat{R}_{S_m}}(h) = \frac{\exp(-\lambda \hat{R}_{S_m}(h))}{\mathbb{E}_{h \sim P}[\exp(-\lambda \hat{R}_{S_m}(h))]} dP(h).$$

### Quick sum up

PAC-Bayes algorithms minimise theoretical bounds  $\rightarrow$  sound theoretical guarantees comes with our posterior.

**Drawbacks Often hard to optimise on  $\mathcal{M}(\mathcal{H})$ , and Gibbs posterior implementation is time-consuming.**



### Quick sum up

PAC-Bayes algorithms minimise theoretical bounds  $\rightarrow$  sound theoretical guarantees comes with our posterior.

**Drawbacks Often hard to optimise on  $\mathcal{M}(\mathcal{H})$ , and Gibbs posterior implementation is time-consuming.**

### Questions:

- **How are those algorithms instantiated in practice?**
- **Are these algorithms efficient and do they come with non-vacuous theoretical guarantees?**

## Instantiation

- Use of multiple data-free priors (grid + union bounds)
- Sacrifice some part of the data to train the prior.
- $\mathcal{C}$  is often a set of Gaussians (closed form of the KL)

## Efficiency

- Non-vacuous generalisation guarantees attainable for small deep nets (Dziugaite *et al.*, 2017 and following works)
- Faster convergence rates via small variance (Tolstikhin *et al.*, 2013)
- When vacuous, use of PAC-Bayes bounds as correlation measures for generalisation (Neyshabur *et al.*, 2017)

In 20+ years of development:

- Inspiration from the Bayesian paradigm
- Little attention on statistical assumptions (*i.i.d.* data, subgaussian losses) except few works *e.g.* (Seldin *et al.*, 2012; Kuzborskij *et al.*, 2019).
- Priors and posteriors are designed *w.r.t.* to the KL divergence: either Gibbs (closed form) or Gaussian (computation)

**Can we weaken the statistical assumption on the loss?**

In 20+ years of development:

- Inspiration from the Bayesian paradigm
- Little attention on statistical assumptions (*i.i.d.* data, subgaussian losses) except few works *e.g.* (Seldin *et al.*, 2012; Kuzborskij *et al.*, 2019).
- Priors and posteriors are designed *w.r.t.* to the KL divergence: either Gibbs (closed form) or Gaussian (computation)

**Can we weaken the statistical assumption on the loss?**

**Yes: we can extend Catoni's bound for any countable dataset  $\mathcal{S}$  and finite variance assumption**

# A PAC-BAYESIAN BOUND FOR UNBOUNDED MARTINGALES

## Theorem

For any data-free prior  $\mathbf{P} \in \mathcal{M}(\mathcal{H})$ , any  $\lambda > 0$ , any collection of martingales  $(M_m(h))_{m \geq 1}$  indexed by  $h \in \mathcal{H}$ , the following holds with probability  $1 - \delta$  over the sample  $\mathcal{S} = (\mathbf{z}_i)_{i \in \mathbb{N}}$ , for all  $m \in \mathbb{N}/\{0\}$ ,  $\mathbf{Q} \in \mathcal{M}(\mathcal{H})$ :

$$|M_m(\mathbf{Q})| \leq \frac{\text{KL}(\mathbf{Q}, \mathbf{P}) + \log(2/\delta)}{\lambda} + \frac{\lambda}{2} ([M]_m(\mathbf{Q}) + \langle M \rangle_m(\mathbf{Q})).$$

# A PAC-BAYESIAN BOUND FOR UNBOUNDED MARTINGALES

## Theorem

For any data-free prior  $\mathbf{P} \in \mathcal{M}(\mathcal{H})$ , any  $\lambda > 0$ , any collection of martingales  $(M_m(h))_{m \geq 1}$  indexed by  $h \in \mathcal{H}$ , the following holds with probability  $1 - \delta$  over the sample  $\mathcal{S} = (\mathbf{z}_i)_{i \in \mathbb{N}}$ , for all  $m \in \mathbb{N}/\{0\}$ ,  $\mathbf{Q} \in \mathcal{M}(\mathcal{H})$ :

$$|M_m(\mathbf{Q})| \leq \frac{\text{KL}(\mathbf{Q}, \mathbf{P}) + \log(2/\delta)}{\lambda} + \frac{\lambda}{2} ([M]_m(\mathbf{Q}) + \langle M \rangle_m(\mathbf{Q})).$$

Required: finiteness of  $([M]_m, \langle M \rangle_m)_{m \geq 1}$  (variance terms)

Toolbox: Ville's inequality and supermartingales

## Corollary

For any data-free prior  $P \in \mathcal{M}(\mathcal{H})$ , any  $\lambda > 0$ , the following holds with probability  $1 - \delta$  over the sample  $\mathcal{S} = (\mathbf{z}_i)_{i \in \mathbb{N}}$ , for all  $m \in \mathbb{N}/\{0\}$ ,  $Q \in \mathcal{M}(\mathcal{H})$

$$\mathbb{E}_{h \sim Q}[\mathbf{R}(h)] \leq \mathbb{E}_{h \sim Q} \left[ \hat{\mathbf{R}}_{\mathcal{S}_m}(h) + \frac{\lambda}{2m} \sum_{i=1}^m (\ell(h, z_i) - \mathbf{R}_{\mathcal{D}}(h))^2 \right] \\ + \frac{\text{KL}(Q, P) + \log(2/\delta)}{\lambda m} + \frac{\lambda}{2} \mathbb{E}_{h \sim Q} [\text{Var}_{\mathcal{D}}(h)],$$

where  $\text{Var}_{\mathcal{D}}(h)$  is the variance of  $\ell(h, \cdot)$ .

**Interesting property:** time-uniform bound.

## A TIGHTER BOUND FOR NON-NEGATIVE LOSSES

In Chugg *et al.* (2023): tighter bound for nonnegative loss:

### Corollary

For  $\ell \geq 0$ , any data-free prior  $P \in \mathcal{M}(\mathcal{H})$ , any  $\lambda > 0$ , the following holds with probability  $1 - \delta$  over the sample  $\mathcal{S} = (\mathbf{z}_i)_{i \in \mathbb{N}}$ , for all  $m \in \mathbb{N}/\{0\}$ ,  $Q \in \mathcal{M}(\mathcal{H})$

$$\mathbb{E}_{h \sim Q}[\mathbf{R}(h)] \leq \mathbb{E}_{h \sim Q}[\hat{\mathbf{R}}_{\mathcal{S}_m}(h)] + \frac{\text{KL}(Q, P) + \log(2/\delta)}{\lambda m} + \frac{\lambda}{2} \mathbb{E}_{h \sim Q}[\ell(h, \mathbf{z})^2],$$



**We can recover Catoni's bound at the sole price of uniformly bounded variance. We also reached time-uniform generalisation bounds**

We can recover Catoni's bound at the sole price of uniformly bounded variance. We also reached time-uniform generalisation bounds

**Drawback:** those bounds holds for any  $\mathcal{Q}$  and  $\mathcal{D}$  simultaneously. Gastpar *et al.* (2023) showed such bounds were limited in the overparametrised setting.

**We can recover Catoni's bound at the sole price of uniformly bounded variance. We also reached time-uniform generalisation bounds**

**Drawback: those bounds holds for any  $\mathcal{Q}$  and  $\mathcal{D}$  simultaneously. Gastpar *et al.* (2023) showed such bounds were limited in the overparametrised setting.**

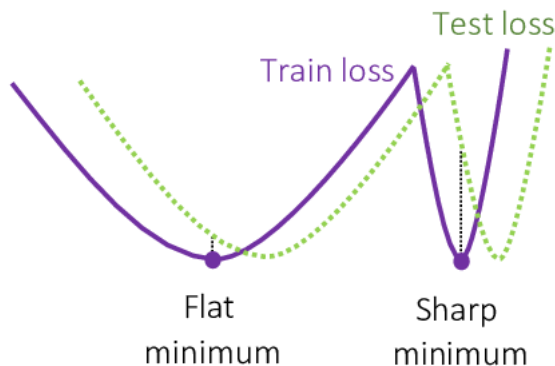
**Question: can we incorporate some benefits of a successful learning process in such bounds?**

# FLAT MINIMUM

## Yes for flat minima !!

### 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 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  $f \in H^1(Q)$ :

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

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



# 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

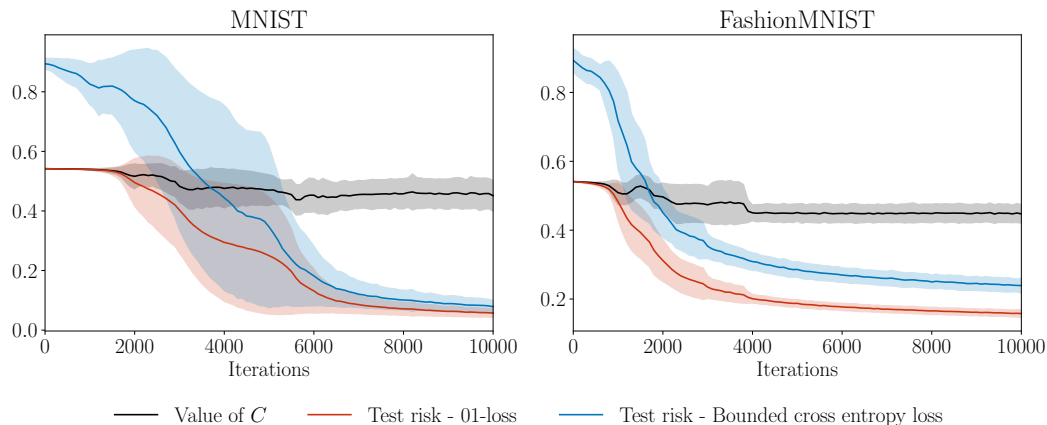
$\mathbb{Q} \in \mathcal{M}(\mathcal{H})$  is *quadratically self-bounded w.r.t.  $\ell$  and  $C > 0$*  (namely  $\text{QSB}(\ell, C)$ ) if

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

- QSB intricates  $\mathcal{D} \in \mathcal{M}(\mathcal{Z})$  with  $\mathbb{Q} \in \mathcal{M}(\mathcal{H})$
- Satisfied if  $\ell \in [0, K]$  with  $C = K$ .
- Also satisfied for unbounded lipschitz losses in a certain setting.

# IS THE QSB ASSUMPTION VERIFIED IN PRACTICE?

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



## 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  $\text{Poinc}(c_P)$ ,  $\text{QSB}(\ell, C)$ ,

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

## 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 $\text{Poinc}$ :

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

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

**Drawback: bounds are not empirical.**

**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  $\text{Poinc}(c_P)$  with constant  $c$ ,  $\text{QSB}(\ell, C_1)$ ,  $\text{QSB}(\|\nabla_h \ell\|^2, C_2)$ ,

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

If  $Q$  satisfies either

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

**2** if  $\ell$  gradient-lipschitz, flat minima for  $\hat{R}_S$  and small gradient norms on each training data,

**then  $Q$  generalises well!**

**Drawback: with Poincaré posteriors, KL is uncontrolled.**



**Drawback: with Poincaré posteriors, KL is uncontrolled.**

**Solution: 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
- if  $\mathbb{P}$  L-Sob(+ technical assumptions) and  $\ell = \ell_1 + \ell_2$  ( $\ell_1$  convex, twice differentiable,  $\ell_2$  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( $\ell, C$ ):

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

# 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  $\text{QSB}(\ell, C)$ :

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

**KL small if a flat minima on  $\hat{\mathbb{R}}_S$  is reached**

- 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)**

# A NEW CHANGE OF MEASURE INEQUALITY

## 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

- A relaxation of 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  $\Delta_{S_m}$  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!
- 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>

**Thank you for your attention!**





# REFERENCES I

-  Andreas Maurer. A note on the PAC-Bayesian theorem. *arXiv. cs/0411099*. (2004).
-  Yevgeny Seldin, François Laviolette, Nicolò Cesa-Bianchi, John Shawe-Taylor, and Peter Auer. PAC-Bayesian Inequalities for Martingales. *IEEE Transactions on Information Theory*. (2012).
-  Ilya O. Tolstikhin and Yevgeny Seldin. PAC-Bayes-Empirical-Bernstein Inequality. *Advances in Neural Information Processing Systems (NeurIPS)*. (2013).
-  Pierre Alquier, James Ridgway, and Nicolas Chopin. On the properties of variational approximations of Gibbs posteriors. *Journal of Machine Learning Research*. (2016).
-  Gintare Karolina Dziugaite and Daniel Roy. Computing Nonvacuous Generalization Bounds for Deep (Stochastic) Neural Networks with Many More Parameters than Training Data. *Conference on Uncertainty in Artificial Intelligence (UAI)*. (2017).






## REFERENCES II

-  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>.
-  Ilja Kuzborskij and Csaba Szepesvári. Efron-Stein PAC-Bayesian Inequalities. *arXiv. abs/1909.01931*. (2019).
-  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/86d7c8a08b4aaa1bc7c599473f5d4dda-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 III

-  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>.
-  Ben Chugg, Hongjian Wang, and Aaditya Ramdas. A Unified Recipe for Deriving (Time-Uniform) PAC-Bayes Bounds. *Journal of Machine Learning Research.* (2023). url: <http://jmlr.org/papers/v24/23-0401.html>.
-  Michael Gastpar, Ido Nachum, Jonathan Shafer, and Thomas Weinberger. Fantastic Generalization Measures are Nowhere to be Found. (2023).
-  Kaiyue Wen, Zhiyuan Li, and Tengyu Ma. Sharpness Minimization Algorithms Do Not Only Minimize Sharpness To Achieve Better Generalization. *Thirty-seventh Conference on Neural Information Processing Systems.* (2023). url: <https://openreview.net/forum?id=Dkmpa6wCIx>.





# REFERENCES I

-  Andreas Maurer. A note on the PAC-Bayesian theorem. *arXiv. cs/0411099*. (2004).
-  Yevgeny Seldin, François Laviolette, Nicolò Cesa-Bianchi, John Shawe-Taylor, and Peter Auer. PAC-Bayesian Inequalities for Martingales. *IEEE Transactions on Information Theory*. (2012).
-  Ilya O. Tolstikhin and Yevgeny Seldin. PAC-Bayes-Empirical-Bernstein Inequality. *Advances in Neural Information Processing Systems (NeurIPS)*. (2013).
-  Pierre Alquier, James Ridgway, and Nicolas Chopin. On the properties of variational approximations of Gibbs posteriors. *Journal of Machine Learning Research*. (2016).
-  Gintare Karolina Dziugaite and Daniel Roy. Computing Nonvacuous Generalization Bounds for Deep (Stochastic) Neural Networks with Many More Parameters than Training Data. *Conference on Uncertainty in Artificial Intelligence (UAI)*. (2017).

## REFERENCES II

-  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>.
-  Ilja Kuzborskij and Csaba Szepesvári. Efron-Stein PAC-Bayesian Inequalities. *arXiv. abs/1909.01931*. (2019).
-  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/86d7c8a08b4aaa1bc7c599473f5d4dda-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 III

-  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>.
-  Ben Chugg, Hongjian Wang, and Aaditya Ramdas. A Unified Recipe for Deriving (Time-Uniform) PAC-Bayes Bounds. *Journal of Machine Learning Research.* (2023). url: <http://jmlr.org/papers/v24/23-0401.html>.
-  Michael Gastpar, Ido Nachum, Jonathan Shafer, and Thomas Weinberger. Fantastic Generalization Measures are Nowhere to be Found. (2023).
-  Kaiyue Wen, Zhiyuan Li, and Tengyu Ma. Sharpness Minimization Algorithms Do Not Only Minimize Sharpness To Achieve Better Generalization. *Thirty-seventh Conference on Neural Information Processing Systems.* (2023). url: <https://openreview.net/forum?id=Dkmpa6wCIx>.



**Thank you for your attention!**