

Learning under model misspecification

Andrés R. Masegosa

University of Almería
Spain

Masegosa, Andrés R. Learning under model misspecification: Applications to variational and ensemble methods. NeurIPS 2020.

Introduction

Model Misspecification

- Virtually any model we use in ML **does not perfectly represent reality**.
- We mostly work in the **model misspecification** regime.

Model Misspecification

- Virtually any model we use in ML **does not perfectly represent reality**.
- We mostly work in the **model misspecification** regime.

Contributions

- **Generalization** analysis of Bayesian learning under model misspecification.
- Bayesian model averaging is **suboptimal** for generalization.
- **New learning framework** which explicitly addresses model misspecification.
- Empirical evaluations on **Bayesian deep learning** illustrate this approach.

Assumption 1: I.I.D. Data

- There exists an **underlying distribution** $\nu(\mathbf{x})$ generating the training/test data.
- The **training data sample**, $D = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$, is i.i.d. from $\nu(\mathbf{x})$.

Assumption 2: Model misspecification

- Our model class only **approximates** reality (not prefect).
- $p(\mathbf{x}|\boldsymbol{\theta})$ is our (parametric) **probabilistic model class**.

Assumption 2: Model misspecification

- Our model class only **approximates** reality (not perfect).
- $p(\mathbf{x}|\boldsymbol{\theta})$ is our (parametric) **probabilistic model class**.

$$\forall \boldsymbol{\theta} \in \Theta \quad \nu \neq p(\cdot|\boldsymbol{\theta})$$

Assumption 3: Likelihood is Upper-Bounded

- There exists a $M > 0$

$$\forall \mathbf{x} \in \mathcal{X}, \forall \boldsymbol{\theta} \in \Theta \quad p(\cdot | \boldsymbol{\theta}) \leq M,$$

Assumption 3: Likelihood is Upper-Bounded

- There exists a $M > 0$

$$\forall \mathbf{x} \in \mathcal{X}, \forall \boldsymbol{\theta} \in \Theta \quad p(\cdot | \boldsymbol{\theta}) \leq M,$$

- It holds in **supervised classification** (e.g. $M = 1$) and it may require to constrain the parameter space (e.g. the variance of the Gaussian higher than $\epsilon > 0$),.

Assumption 3: Likelihood is Upper-Bounded

- There exists a $M > 0$

$$\forall \mathbf{x} \in \mathcal{X}, \forall \boldsymbol{\theta} \in \Theta \quad p(\cdot | \boldsymbol{\theta}) \leq M,$$

- It holds in **supervised classification** (e.g. $M = 1$) and it may require to constrain the parameter space (e.g. the variance of the Gaussian higher than $\epsilon > 0$),.

This analysis also applies to a supervised settings!!

The Learning Problem

- **Notation:**

- $\rho(\theta)$ is a probability distribution over the parameters of my model.

- **Notation:**

- $\rho(\theta)$ is a probability distribution over the parameters of my model.
- $\rho(\theta)$ depends on data. It is a **quasi-posterior**.

- **Notation:**

- $\rho(\theta)$ is a probability distribution over the parameters of my model.
- $\rho(\theta)$ depends on data. It is a **quasi-posterior**.

- The **predictive posterior distribution** for a given $\rho(\theta)$,

$$p(\mathbf{x}) = \int p(\mathbf{x}|\theta)\rho(\theta)d\theta = \mathbb{E}_{\rho(\theta)}[p(\mathbf{x}|\theta)]$$

- **Notation:**

- $\rho(\theta)$ is a probability distribution over the parameters of my model.
- $\rho(\theta)$ depends on data. It is a **quasi-posterior**.

- The **predictive posterior distribution** for a given $\rho(\theta)$,

$$p(\mathbf{x}) = \int p(\mathbf{x}|\theta)\rho(\theta)d\theta = \mathbb{E}_{\rho(\theta)}[p(\mathbf{x}|\theta)]$$

- $\mathbb{E}_{\rho(\theta)}[p(\mathbf{x}|\theta)]$ is **Bayesian model averaging** when $\rho(\theta) = p(\theta|D)$,

- **Notation:**

- $\rho(\theta)$ is a probability distribution over the parameters of my model.
- $\rho(\theta)$ depends on data. It is a **quasi-posterior**.

- The **predictive posterior distribution** for a given $\rho(\theta)$,

$$p(\mathbf{x}) = \int p(\mathbf{x}|\theta)\rho(\theta)d\theta = \mathbb{E}_{\rho(\theta)}[p(\mathbf{x}|\theta)]$$

- $\mathbb{E}_{\rho(\theta)}[p(\mathbf{x}|\theta)]$ is **Bayesian model averaging** when $\rho(\theta) = p(\theta|D)$,

- The **learning problem** is defined as,

$$\rho^* = \arg \min_{\rho} KL(\underbrace{\nu(\mathbf{x})}_{\text{Data distribution}}, \underbrace{\mathbb{E}_{\rho(\theta)}[p(\mathbf{x}|\theta)]}_{p(\mathbf{x})})$$

- **Notation:**

- $\rho(\theta)$ is a probability distribution over the parameters of my model.
- $\rho(\theta)$ depends on data. It is a **quasi-posterior**.

- The **predictive posterior distribution** for a given $\rho(\theta)$,

$$p(\mathbf{x}) = \int p(\mathbf{x}|\theta)\rho(\theta)d\theta = \mathbb{E}_{\rho(\theta)}[p(\mathbf{x}|\theta)]$$

- $\mathbb{E}_{\rho(\theta)}[p(\mathbf{x}|\theta)]$ is **Bayesian model averaging** when $\rho(\theta) = p(\theta|D)$,

- The **learning problem** is defined as,

$$\rho^* = \arg \min_{\rho} KL(\underbrace{\nu(\mathbf{x})}_{\text{Data distribution}}, \underbrace{\mathbb{E}_{\rho(\theta)}[p(\mathbf{x}|\theta)]}_{p(\mathbf{x})}) = \arg \min_{\rho} \underbrace{\mathbb{E}_{\nu(\mathbf{x})}[-\ln \mathbb{E}_{\rho(\theta)}[p(\mathbf{x}|\theta)]]}_{CE(\rho)}$$

- **Notation:**

- $\rho(\theta)$ is a probability distribution over the parameters of my model.
- $\rho(\theta)$ depends on data. It is a **quasi-posterior**.

- The **predictive posterior distribution** for a given $\rho(\theta)$,

$$p(\mathbf{x}) = \int p(\mathbf{x}|\theta)\rho(\theta)d\theta = \mathbb{E}_{\rho(\theta)}[p(\mathbf{x}|\theta)]$$

- $\mathbb{E}_{\rho(\theta)}[p(\mathbf{x}|\theta)]$ is **Bayesian model averaging** when $\rho(\theta) = p(\theta|D)$,

- The **learning problem** is defined as,

$$\rho^* = \arg \min_{\rho} KL(\underbrace{\nu(\mathbf{x})}_{\text{Data distribution}}, \underbrace{\mathbb{E}_{\rho(\theta)}[p(\mathbf{x}|\theta)]}_{p(\mathbf{x})}) = \arg \min_{\rho} \underbrace{\mathbb{E}_{\nu(\mathbf{x})}[-\ln \mathbb{E}_{\rho(\theta)}[p(\mathbf{x}|\theta)]]}_{CE(\rho)}$$

- $CE(\rho)$ measures the **generalization error** (or the predictive risk) associated to ρ .

The learning strategy

- How to solve this problem

$$\rho^* = \arg \min_{\rho} \underbrace{CE(\rho)}_{\substack{\text{Generalization} \\ \text{Error}}}$$

if we do not have access to $\nu(\mathbf{x})$

The learning strategy

- How to solve this problem

$$\rho^* = \arg \min_{\rho} \underbrace{CE(\rho)}_{\text{Generalization Error}}$$

if we do not have access to $\nu(\mathbf{x})$

The learning strategy

- The solution is to **employ upper-bounds**:

$$CE(\rho) \underbrace{\leq}_{\text{Jensen inequality}} \text{Oracle-Bound}(\rho, \nu) \underbrace{\lesssim}_{\text{w.p. } (1-\xi)} \text{Empirical-Bound}(\rho, D, \xi)$$

The learning strategy

- How to solve this problem

$$\rho^* = \arg \min_{\rho} \underbrace{CE(\rho)}_{\text{Generalization Error}}$$

if we do not have access to $\nu(\mathbf{x})$

The learning strategy

- The solution is to **employ upper-bounds**:

$$CE(\rho) \underbrace{\leq}_{\text{Jensen inequality}} \text{Oracle-Bound}(\rho, \nu) \underbrace{\lesssim}_{\text{w.p. } (1-\xi)} \text{Empirical-Bound}(\rho, D, \xi)$$

- ... and **minimize** $\text{Empirical-Bound}(\rho, D, \xi)$,

$$\min_{\rho} \text{Empirical-Bound}(\rho, D, \xi)$$

The learning strategy

- How to solve this problem

$$\rho^* = \arg \min_{\rho} \underbrace{CE(\rho)}_{\text{Generalization Error}}$$

if we do not have access to $\nu(\mathbf{x})$

The learning strategy

- The solution is to **employ upper-bounds**:

$$CE(\rho) \underbrace{\leq}_{\text{Jensen inequality}} \text{Oracle-Bound}(\rho, \nu) \underbrace{\lesssim}_{\text{w.p. } (1-\xi)} \text{Empirical-Bound}(\rho, D, \xi)$$

- ... and **minimize** $\text{Empirical-Bound}(\rho, D, \xi)$,

$$\min_{\rho} \text{Empirical-Bound}(\rho, D, \xi)$$

- The quality of the solution is going to depend of the **quality of the bounds**.

First-order Jensen bounds and the Bayesian posterior

$$\underbrace{CE(\rho)}_{\text{Generalization Error}} \overset{\text{Jensen Inequality}}{\leq} \underbrace{\mathbb{E}_{\rho}[L(\boldsymbol{\theta})]}_{\text{Oracle bound}}$$

$L(\boldsymbol{\theta})$ is the **expected log-loss**, $L(\boldsymbol{\theta}) = -\mathbb{E}_{\nu(\mathbf{x})}[\ln p(\mathbf{x}|\boldsymbol{\theta})]$.

$$\underbrace{CE(\rho)}_{\text{Generalization Error}} \overset{\text{Jensen Inequality}}{\leq} \underbrace{\mathbb{E}_{\rho}[L(\boldsymbol{\theta})]}_{\text{Oracle bound}} \overset{\text{w.p. } (1-\xi)}{\lesssim} \underbrace{\mathbb{E}_{\rho}[\hat{L}(\boldsymbol{\theta}, D)] + \frac{KL(\rho, \pi)}{n} + \frac{cte}{n}}_{\text{PAC-Bayes bound (Alquier et al. 2016)}}$$

$L(\boldsymbol{\theta})$ is the **expected log-loss**, $L(\boldsymbol{\theta}) = -\mathbb{E}_{\nu(\mathbf{x})}[\ln p(\mathbf{x}|\boldsymbol{\theta})]$.

$\hat{L}(\boldsymbol{\theta}, D)$ is the **empirical log-loss**, $\hat{L}(\boldsymbol{\theta}, D) = -\frac{1}{n} \ln p(D|\boldsymbol{\theta})$.

$\pi(\boldsymbol{\theta})$ is a prior, which is independent of D .

The Bayesian posterior (Germain et al. 2016)

- The learning strategy is to **minimize the PAC-Bayes bound**,

$$\rho^* = \arg \min_{\rho} \underbrace{\mathbb{E}_{\rho(\theta)}[L(\theta, D)] + \frac{KL(\rho, \pi)}{n}}_{\text{PAC-Bayes bound (Alquier et al. 2016)}} + \frac{cte}{n}$$

The Bayesian posterior (Germain et al. 2016)

- The learning strategy is to **minimize the PAC-Bayes bound**,

$$\begin{aligned}\rho^* &= \arg \min_{\rho} \underbrace{\mathbb{E}_{\rho(\theta)}[L(\theta, D)] + \frac{KL(\rho, \pi)}{n} + \frac{cte}{n}}_{\text{PAC-Bayes bound (Alquier et al. 2016)}} \\ &= \arg \max_{\rho} \underbrace{\mathbb{E}_{\rho(\theta)}[\ln p(D|\theta)] - KL(\rho, \pi)}_{\text{Evidence Lower Bound (ELBO)}}\end{aligned}$$

The Bayesian posterior (Germain et al. 2016)

- The learning strategy is to **minimize the PAC-Bayes bound**,

$$\begin{aligned}\rho^* &= \arg \min_{\rho} \underbrace{\mathbb{E}_{\rho(\theta)}[L(\theta, D)] + \frac{KL(\rho, \pi)}{n} + \frac{cte}{n}}_{\text{PAC-Bayes bound (Alquier et al. 2016)}} \\ &= \arg \max_{\rho} \underbrace{\mathbb{E}_{\rho(\theta)}[\ln p(D|\theta)] - KL(\rho, \pi)}_{\text{Evidence Lower Bound (ELBO)}}\end{aligned}$$

- ρ^* is the **Bayesian posterior**,

$$\rho^* = p(\theta|D) = \frac{p(D|\theta)\pi(\theta)}{\int p(D|\theta)\pi(\theta)d\theta}$$

The Bayesian posterior (Germain et al. 2016)

- The learning strategy is to **minimize the PAC-Bayes bound**,

$$\begin{aligned}\rho^* &= \arg \min_{\rho} \underbrace{\mathbb{E}_{\rho(\theta)}[L(\theta, D)] + \frac{KL(\rho, \pi)}{n} + \frac{cte}{n}}_{\text{PAC-Bayes bound (Alquier et al. 2016)}} \\ &= \arg \max_{\rho} \underbrace{\mathbb{E}_{\rho(\theta)}[\ln p(D|\theta)] - KL(\rho, \pi)}_{\text{Evidence Lower Bound (ELBO)}}\end{aligned}$$

- ρ^* is the **Bayesian posterior**,

$$\rho^* = p(\theta|D) = \frac{p(D|\theta)\pi(\theta)}{\int p(D|\theta)\pi(\theta)d\theta}$$

The Bayesian posterior is a proxy

$$p(\theta|D) \approx \arg \min_{\rho} KL(\underbrace{\nu(\mathbf{x})}_{\text{Data distribution}}, \underbrace{\mathbb{E}_{\rho(\theta)}[p(\mathbf{x}|\theta)]}_{\text{Predictive posterior}})$$

Is the Bayesian approach an optimal learning strategy?

The **Bayesian learning strategy**,

$$\underbrace{CE(\rho)}_{\text{Generalization Error}} \overset{\text{Jensen Inequality}}{\leq} \underbrace{\mathbb{E}_{\rho}[L(\boldsymbol{\theta})]}_{\text{First-Order Jensen bound}} \overset{\text{w.p. } (1-\xi)}{\lesssim} \underbrace{\mathbb{E}_{\rho}[\hat{L}(\boldsymbol{\theta}, D)] + \frac{KL(\rho, \pi)}{n} + \frac{cte}{n}}_{\text{PAC-Bayes bound (Alquier et al. 2016)}}$$

Is the Bayesian approach an optimal learning strategy?

The **Bayesian learning strategy**,

$$\underbrace{CE(\rho)}_{\text{Generalization Error}} \overset{\text{Jensen Inequality}}{\leq} \underbrace{\mathbb{E}_{\rho}[L(\boldsymbol{\theta})]}_{\text{First-Order Jensen bound}} \overset{\text{w.p. } (1-\xi)}{\lesssim} \underbrace{\mathbb{E}_{\rho}[\hat{L}(\boldsymbol{\theta}, D)] + \frac{KL(\rho, \pi)}{n} + \frac{cte}{n}}_{\text{PAC-Bayes bound (Alquier et al. 2016)}}$$

- 1 The Bayesian posterior **converges** to the minimum of $\mathbb{E}_{\rho}[L(\boldsymbol{\theta})]$.

Is the Bayesian approach an optimal learning strategy?

The **Bayesian learning strategy**,

$$\underbrace{CE(\rho)}_{\text{Generalization Error}} \overset{\text{Jensen Inequality}}{\leq} \underbrace{\mathbb{E}_{\rho}[L(\theta)]}_{\text{First-Order Jensen bound}} \overset{\text{w.p. } (1-\xi)}{\lesssim} \underbrace{\mathbb{E}_{\rho}[\hat{L}(\theta, D)] + \frac{KL(\rho, \pi)}{n} + \frac{cte}{n}}_{\text{PAC-Bayes bound (Alquier et al. 2016)}}$$

- 1 The Bayesian posterior **converges** to the minimum of $\mathbb{E}_{\rho}[L(\theta)]$.
- 2 The minimum of $\mathbb{E}_{\rho}[L(\theta)]$ is

Is the Bayesian approach an optimal learning strategy?

The **Bayesian learning strategy**,

$$\underbrace{CE(\rho)}_{\text{Generalization Error}} \overset{\text{Jensen Inequality}}{\leq} \underbrace{\mathbb{E}_{\rho}[L(\boldsymbol{\theta})]}_{\text{First-Order Jensen bound}} \overset{\text{w.p. } (1-\xi)}{\lesssim} \underbrace{\mathbb{E}_{\rho}[\hat{L}(\boldsymbol{\theta}, D)] + \frac{KL(\rho, \pi)}{n} + \frac{cte}{n}}_{\text{PAC-Bayes bound (Alquier et al. 2016)}}$$

- 1 The Bayesian posterior **converges** to the minimum of $\mathbb{E}_{\rho}[L(\boldsymbol{\theta})]$.
- 2 The minimum of $\mathbb{E}_{\rho}[L(\boldsymbol{\theta})]$ is

A **Dirac-delta distribution** centered around $\boldsymbol{\theta}_J^* = \arg \min_{\boldsymbol{\theta}} KL(\nu(\mathbf{x}), p(\mathbf{x}|\boldsymbol{\theta}))$

Is the Bayesian approach an optimal learning strategy?

The **Bayesian learning strategy**,

$$\underbrace{CE(\rho)}_{\text{Generalization Error}} \stackrel{\text{Jensen Inequality}}{\leq} \underbrace{\mathbb{E}_{\rho}[L(\theta)]}_{\text{First-Order Jensen bound}} \stackrel{\text{w.p. } (1-\xi)}{\lesssim} \underbrace{\mathbb{E}_{\rho}[\hat{L}(\theta, D)] + \frac{KL(\rho, \pi)}{n} + \frac{cte}{n}}_{\text{PAC-Bayes bound (Alquier et al. 2016)}}$$

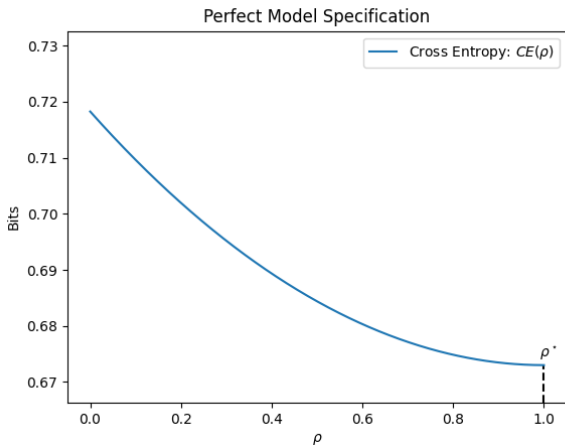
- 1 The Bayesian posterior **converges** to the minimum of $\mathbb{E}_{\rho}[L(\theta)]$.
- 2 The minimum of $\mathbb{E}_{\rho}[L(\theta)]$ is

A **Dirac-delta distribution** centered around $\theta_J^* = \arg \min_{\theta} KL(\nu(\mathbf{x}), p(\mathbf{x}|\theta))$

Is the Bayesian approach optimal for minimizing the generalization error?

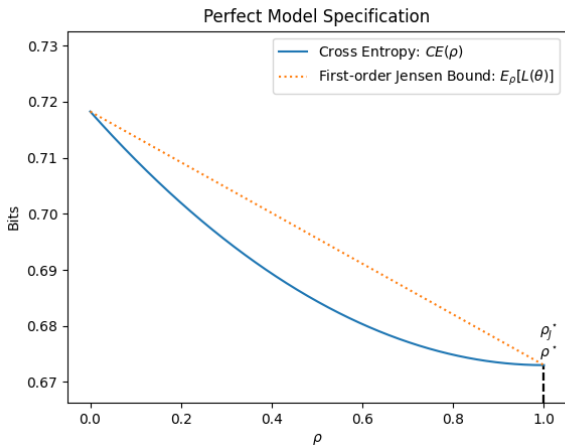
- Is this **Dirac-delta distribution** centered around θ_J^* a good proxy of ρ^* ?

$$\rho^* = \arg \min_{\rho} KL(\underbrace{\nu(\mathbf{x})}_{\text{Data distribution}}, \underbrace{\mathbb{E}_{\rho(\theta)}[p(\mathbf{x}|\theta)]}_{\text{Predictive posterior}})$$



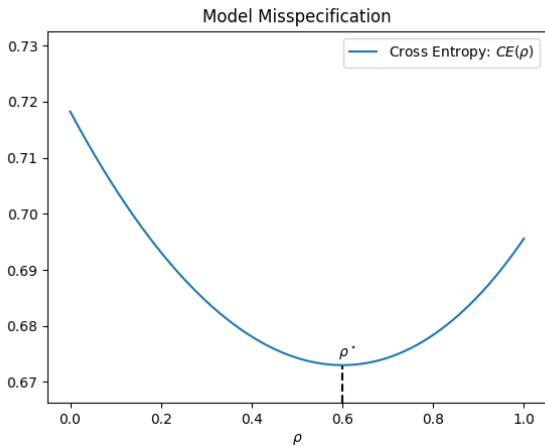
$$\arg \min_{\rho} \underbrace{CE(\rho)}_{\text{Generalization Error}} = \underbrace{\delta_{\theta_J}(\theta)}_{\text{Dirac-Delta distribution}}$$

Bayesian posterior **optimal** under perfect specification



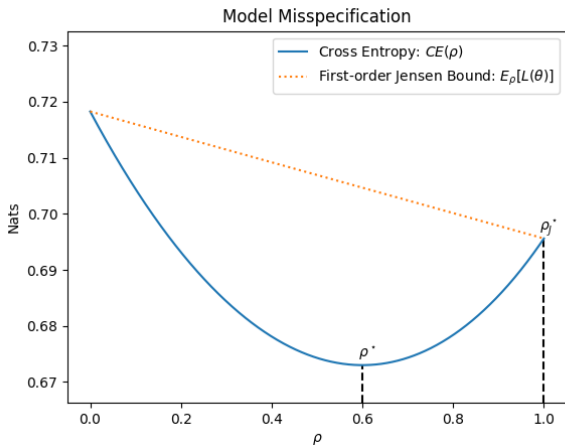
$$\arg \min_{\rho} \underbrace{CE(\rho)}_{\text{Generalization Error}} = \underbrace{\delta_{\theta_J}(\theta)}_{\text{Dirac-Delta distribution}} = \arg \min_{\rho} \underbrace{\mathbb{E}_{\rho}[L(\theta)]}_{\text{First-Order Jensen Bound}}$$

Bayesian posterior **not optimal** under misspecification



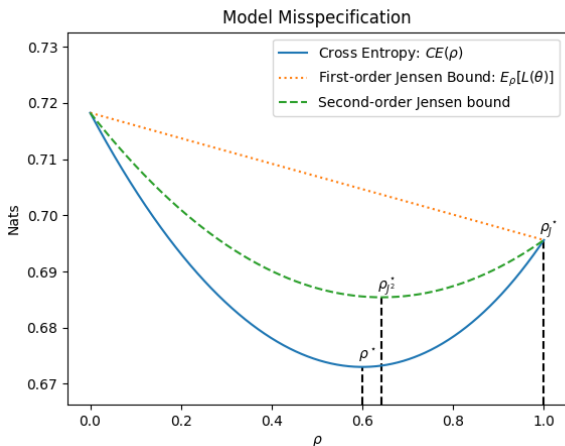
$$\arg \min_{\rho} \underbrace{CE(\rho)}_{\text{Generalization Error}} \neq \underbrace{\delta_{\theta_J}(\theta)}_{\text{Dirac-Delta distribution}}$$

Bayesian posterior **not optimal** under misspecification



$$\arg \min_{\rho} \underbrace{CE(\rho)}_{\text{Generalization Error}} \neq \underbrace{\delta_{\theta_J}(\theta)}_{\text{Dirac-Delta distribution}} = \arg \min_{\rho} \underbrace{\mathbb{E}_{\rho}[L(\theta)]}_{\text{First-Order Jensen Bound}}$$

Bayesian posterior **not optimal** under misspecification



$$\underbrace{CE(\rho)}_{\text{Generalization Error}} \leq \underbrace{E_{\rho}[L(\theta)]}_{\text{Second-order Jensen bound (Liao et al. 2019)}} - \underbrace{\text{Variance}}_{\text{Variance}}$$

$$CE(\rho) \leq \underbrace{\mathbb{E}_{\rho}[L(\boldsymbol{\theta})]}_{\text{Second-order Jensen bound (Liao et al. 2019)}} - \overbrace{\mathbb{V}(\rho)}^{\text{Variance}}$$

$$CE(\rho) \leq \underbrace{\mathbb{E}_{\rho}[L(\boldsymbol{\theta})]}_{\text{Second-order Jensen bound (Liao et al. 2019)}} - \overbrace{\mathbb{V}(\rho)}^{\text{Variance}}$$

- $\mathbb{V}(\rho)$ measures the variance of the predictive posterior:

$$\mathbb{V}(\rho) = \mathbb{E}_{\nu(\mathbf{x})} \left[\frac{1}{2 \max_{\theta} p(\mathbf{x}|\boldsymbol{\theta})^2} \mathbb{E}_{\rho(\theta)} [(p(\mathbf{x}|\boldsymbol{\theta}) - p(\mathbf{x}))^2] \right] \geq 0$$

$$CE(\rho) \leq \underbrace{\mathbb{E}_{\rho}[L(\boldsymbol{\theta})]}_{\text{Second-order Jensen bound (Liao et al. 2019)}} - \overbrace{\mathbb{V}(\rho)}^{\text{Variance}}$$

- $\mathbb{V}(\rho)$ measures the variance of the predictive posterior:

$$\mathbb{V}(\rho) = \mathbb{E}_{\nu(\mathbf{x})} \left[\frac{1}{2 \max_{\theta} p(\mathbf{x}|\theta)^2} \mathbb{E}_{\rho(\theta)} [(p(\mathbf{x}|\theta) - p(\mathbf{x}))^2] \right] \geq 0$$

- $\mathbb{V}(\rho)$ accounts for *model diversity*:

$$\mathbb{V}(\rho) = 0 \text{ if } \forall \theta \neq \theta' \quad p(\mathbf{x}|\theta) = p(\mathbf{x}|\theta')$$

$$\underbrace{CE(\rho)}_{\text{Generalization Error}} \leq \underbrace{\mathbb{E}_{\rho}[L(\boldsymbol{\theta})] - \mathbb{V}(\rho)}_{\text{Second-order Jensen bound}} \lesssim \underbrace{\mathbb{E}_{\rho}[\hat{L}(\boldsymbol{\theta}, D)] - \overbrace{\hat{\mathbb{V}}(\rho, D)}^{\text{Empirical Variance}} + \frac{KL(\rho, \pi)}{n} + \frac{cte}{n}}_{\text{Second-order PAC-Bayes bound}}$$

$$\underbrace{CE(\rho)}_{\text{Generalization Error}} \leq \underbrace{\mathbb{E}_\rho[L(\boldsymbol{\theta})] - \mathbb{V}(\rho)}_{\text{Second-order Jensen bound}} \lesssim \underbrace{\mathbb{E}_\rho[\hat{L}(\boldsymbol{\theta}, D)] - \overbrace{\hat{\mathbb{V}}(\rho, D)}^{\text{Empirical Variance}}}_{\text{Second-order PAC-Bayes bound}} + \frac{KL(\rho, \pi)}{n} + \frac{cte}{n}$$

- $\mathbb{E}_\rho[\hat{L}(\boldsymbol{\theta}, D)]$ encourages to place ρ around individual models with small error.

$$\underbrace{CE(\rho)}_{\text{Generalization Error}} \leq \underbrace{\mathbb{E}_\rho[L(\boldsymbol{\theta})] - \mathbb{V}(\rho)}_{\text{Second-order Jensen bound}} \lesssim \underbrace{\mathbb{E}_\rho[\hat{L}(\boldsymbol{\theta}, D)] - \overbrace{\hat{\mathbb{V}}(\rho, D)}^{\text{Empirical Variance}} + \frac{KL(\rho, \pi)}{n} + \frac{cte}{n}}_{\text{Second-order PAC-Bayes bound}}$$

- $\mathbb{E}_\rho[\hat{L}(\boldsymbol{\theta}, D)]$ encourages to place ρ around individual models with small error.
- $\hat{\mathbb{V}}(\rho, D)$ encourages *diversity* among models.

$$\underbrace{CE(\rho)}_{\text{Generalization Error}} \leq \underbrace{\mathbb{E}_\rho[L(\theta)] - \mathbb{V}(\rho)}_{\text{Second-order Jensen bound}} \lesssim \underbrace{\mathbb{E}_\rho[\hat{L}(\theta, D)] - \hat{\mathbb{V}}(\rho, D)}_{\text{Second-order PAC-Bayes bound}} + \frac{KL(\rho, \pi)}{n} + \frac{cte}{n}$$

Empirical Variance
 $\hat{\mathbb{V}}(\rho, D)$

- $\mathbb{E}_\rho[\hat{L}(\theta, D)]$ encourages to place ρ around individual models with small error.
- $\hat{\mathbb{V}}(\rho, D)$ encourages *diversity* among models.
 - Key factor when learning under model misspecification.

$$\underbrace{CE(\rho)}_{\text{Generalization Error}} \leq \underbrace{\mathbb{E}_\rho[L(\theta)] - \mathbb{V}(\rho)}_{\text{Second-order Jensen bound}} \lesssim \underbrace{\mathbb{E}_\rho[\hat{L}(\theta, D)] - \overbrace{\hat{\mathbb{V}}(\rho, D)}^{\text{Empirical Variance}}}_{\text{Second-order PAC-Bayes bound}} + \frac{KL(\rho, \pi)}{n} + \frac{cte}{n}$$

- $\mathbb{E}_\rho[\hat{L}(\theta, D)]$ encourages to place ρ around individual models with small error.
- $\hat{\mathbb{V}}(\rho, D)$ encourages *diversity* among models.
 - Key factor when learning under model misspecification.
- $\frac{KL(\rho, \pi)}{n}$ encourages ρ to be close to π (i.e. acts as a regularizer).

Learning by Minimizing second-order PAC-Bayes bounds

PAC²-Bayesian Learning

- A **variational-like method**,

$$\arg \min_{\rho \in \mathcal{Q}} \underbrace{\mathbb{E}_{\rho(\theta)}[L(\theta, D)] - \hat{\mathbb{V}}(\rho, D) + \frac{KL(\rho, \pi)}{n} + \frac{cte}{n}}_{\text{Second-order PAC-Bayes Bound}}$$

PAC²-Bayesian Learning

- A **variational-like method**,

$$\arg \min_{\rho \in Q} \underbrace{\mathbb{E}_{\rho(\theta)}[L(\theta, D)] - \hat{V}(\rho, D)}_{\text{Second-order PAC-Bayes Bound}} + \frac{KL(\rho, \pi)}{n} + \frac{cte}{n}$$

where Q is a tractable family of densities (i.e. fully factorized Gaussian distribution).

PAC²-Bayesian Learning

- A **variational-like method**,

$$\arg \min_{\rho \in Q} \underbrace{\mathbb{E}_{\rho(\theta)}[L(\theta, D)] - \hat{V}(\rho, D) + \frac{KL(\rho, \pi)}{n} + \frac{cte}{n}}_{\text{Second-order PAC-Bayes Bound}}$$

where Q is a tractable family of densities (i.e. fully factorized Gaussian distribution).

- This is a **generalized variational inference** method (Knoblauch et al. 2019).

PAC²-Bayesian Learning

- A **variational-like method**,

$$\arg \min_{\rho \in Q} \underbrace{\mathbb{E}_{\rho(\theta)}[L(\theta, D)] - \hat{V}(\rho, D)}_{\text{Second-order PAC-Bayes Bound}} + \frac{KL(\rho, \pi)}{n} + \frac{cte}{n}$$

where Q is a tractable family of densities (i.e. fully factorized Gaussian distribution).

- This is a **generalized variational inference** method (Knoblauch et al. 2019).
- Different **solvers** are available in the literature (Wang et al. 2017).

PAC²-Bayesian Learning

- A **variational-like method**,

$$\arg \min_{\rho \in Q} \underbrace{\mathbb{E}_{\rho(\theta)}[L(\theta, D)] - \hat{V}(\rho, D) + \frac{KL(\rho, \pi)}{n} + \frac{cte}{n}}_{\text{Second-order PAC-Bayes Bound}}$$

where Q is a tractable family of densities (i.e. fully factorized Gaussian distribution).

- This is a **generalized variational inference** method (Knoblauch et al. 2019).
- Different **solvers** are available in the literature (Wang et al. 2017).

Variational Inference

- **Standard Variational methods** tries to minimize the first-order PAC-Bayes bound,

$$\arg \min_{\rho \in Q} \underbrace{\mathbb{E}_{\rho(\theta)}[L(\theta, D)] + \frac{KL(\rho, \pi)}{n} + \frac{cte}{n}}_{\text{First-order PAC-Bayes Bound}}$$

Ensembles through Mixture of Dirac-delta distributions

- ρ defined as a **mixture of Dirac-delta distributions** parametrized by $\{\theta_1, \dots, \theta_E\}$,

$$\rho_E(\theta) = \sum_{j=1}^E \frac{1}{E} \delta_{\theta_j}(\theta)$$

where δ_{θ_j} is a **Dirac-delta distribution** centered around θ_j

Ensembles through Mixture of Dirac-delta distributions

- ρ defined as a **mixture of Dirac-delta distributions** parametrized by $\{\theta_1, \dots, \theta_E\}$,

$$\rho_E(\theta) = \sum_{j=1}^E \frac{1}{E} \delta_{\theta_j}(\theta)$$

where δ_{θ_j} is a **Dirac-delta distribution** centered around θ_j

- **New ensemble learning framework:**

$$\arg \min_{\{\theta_1, \dots, \theta_E\}} \underbrace{\frac{1}{E} \sum_{j=1}^E \hat{L}(\theta_j, D)}_{\text{Individual Model Errors}} - \underbrace{\hat{\mathbb{V}}(\rho_E, D)}_{\text{Ensemble Diversity}} - \underbrace{\frac{1}{E} \sum_{j=1}^E \frac{\ln \pi(\theta_j)}{n}}_{\text{Regularizer}}$$

Ensembles through Mixture of Dirac-delta distributions

- ρ defined as a **mixture of Dirac-delta distributions** parametrized by $\{\theta_1, \dots, \theta_E\}$,

$$\rho_E(\theta) = \sum_{j=1}^E \frac{1}{E} \delta_{\theta_j}(\theta)$$

where δ_{θ_j} is a **Dirac-delta distribution** centered around θ_j

- **New ensemble learning framework:**

$$\arg \min_{\{\theta_1, \dots, \theta_E\}} \underbrace{\frac{1}{E} \sum_{j=1}^E \hat{L}(\theta_j, D)}_{\text{Individual Model Errors}} - \underbrace{\hat{\mathbb{V}}(\rho_E, D)}_{\text{Ensemble Diversity}} - \underbrace{\frac{1}{E} \sum_{j=1}^E \frac{\ln \pi(\theta_j)}{n}}_{\text{Regularizer}}$$

- **Diversity** has been widely recognized as a key factor in ensemble methods:

Ensembles through Mixture of Dirac-delta distributions

- ρ defined as a **mixture of Dirac-delta distributions** parametrized by $\{\theta_1, \dots, \theta_E\}$,

$$\rho_E(\theta) = \sum_{j=1}^E \frac{1}{E} \delta_{\theta_j}(\theta)$$

where δ_{θ_j} is a **Dirac-delta distribution** centered around θ_j

- **New ensemble learning framework:**

$$\arg \min_{\{\theta_1, \dots, \theta_E\}} \underbrace{\frac{1}{E} \sum_{j=1}^E \hat{L}(\theta_j, D)}_{\text{Individual Model Errors}} - \underbrace{\hat{\mathbb{V}}(\rho_E, D)}_{\text{Ensemble Diversity}} - \underbrace{\frac{1}{E} \sum_{j=1}^E \frac{\ln \pi(\theta_j)}{n}}_{\text{Regularizer}}$$

- **Diversity** has been widely recognized as a key factor in ensemble methods:
 - $\mathbb{V}(\rho_E, D)$ is a well-founded diversity measure.

Ensembles through Mixture of Dirac-delta distributions

- ρ defined as a **mixture of Dirac-delta distributions** parametrized by $\{\theta_1, \dots, \theta_E\}$,

$$\rho_E(\theta) = \sum_{j=1}^E \frac{1}{E} \delta_{\theta_j}(\theta)$$

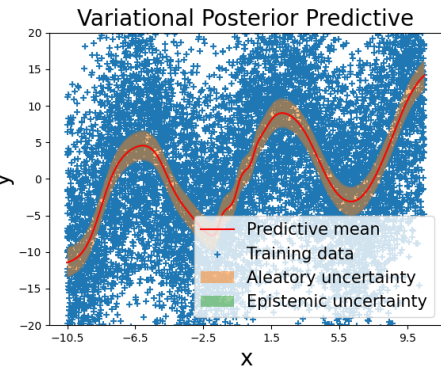
where δ_{θ_j} is a **Dirac-delta distribution** centered around θ_j

- **New ensemble learning framework:**

$$\arg \min_{\{\theta_1, \dots, \theta_E\}} \underbrace{\frac{1}{E} \sum_{j=1}^E \hat{L}(\theta_j, D)}_{\text{Individual Model Errors}} - \underbrace{\hat{\mathbb{V}}(\rho_E, D)}_{\text{Ensemble Diversity}} - \underbrace{\frac{1}{E} \sum_{j=1}^E \frac{\ln \pi(\theta_j)}{n}}_{\text{Regularizer}}$$

- **Diversity** has been widely recognized as a key factor in ensemble methods:
 - $\mathbb{V}(\rho_E, D)$ is a well-founded diversity measure.
 - Help to explain **why diversity is key for generalization in ensembles**.

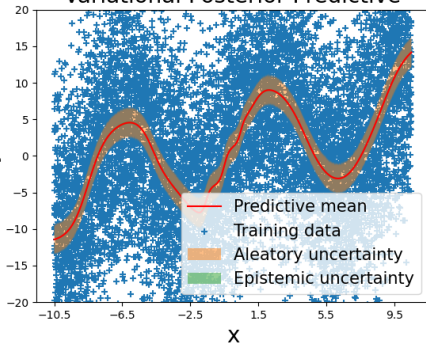
Experimental Evaluation with Toy Data Sets



Test Log-likelihood = -50.15

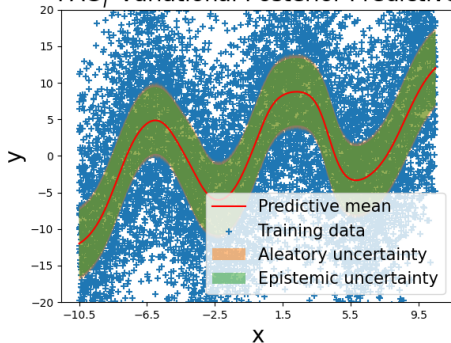
$$\begin{aligned}\nu(y|x) &= \mathcal{N}(\mu = s(x), \sigma^2 = 10) \\ p(y|x, \theta) &= \mathcal{N}(\mu = MLP_{20}(x; \theta), \sigma^2 = 1) \\ \rho(\theta) &= \prod_i \mathcal{N}(\mu_i, \sigma_i)\end{aligned}$$

Variational Posterior Predictive



Test Log-likelihood = -50.15

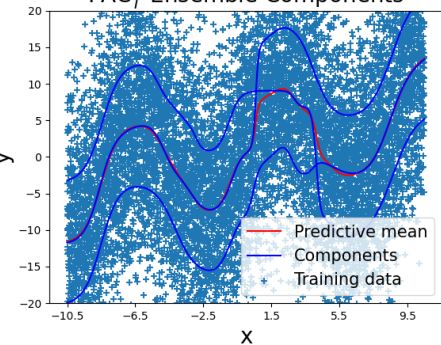
PAC_T²-Variational Posterior Predictive



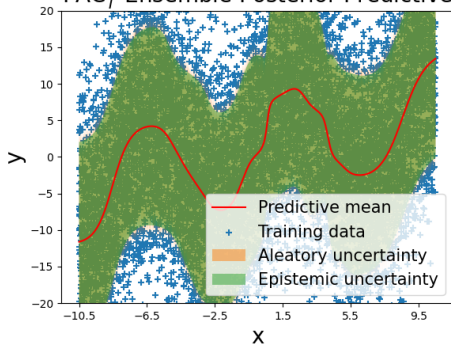
Test Log-likelihood = -25.23

$$\begin{aligned}
 \nu(y|x) &= \mathcal{N}(\mu = s(x), \sigma^2 = 10) \\
 p(y|x, \theta) &= \mathcal{N}(\mu = MLP_{20}(x; \theta), \sigma^2 = 1) \\
 \rho(\theta) &= \prod_i \mathcal{N}(\mu_i, \sigma_i)
 \end{aligned}$$

PAC_T²-Ensemble Components



PAC_T²-Ensemble Posterior Predictive



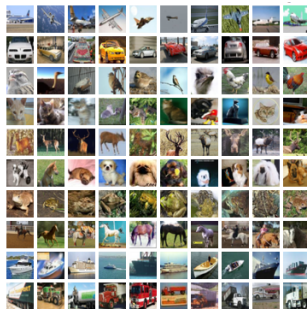
Test Log-likelihood = -15.91

$$\begin{aligned}\nu(y|x) &= \mathcal{N}(\mu = s(x), \sigma^2 = 10) \\ p(y|x, \theta) &= \mathcal{N}(\mu = MLP_{20}(x; \theta), \sigma^2 = 1) \\ \rho(\theta) &= \sum_{j=1}^3 \frac{1}{E} \delta_{\theta_j}(\theta)\end{aligned}$$

Experimental Evaluation on real data sets



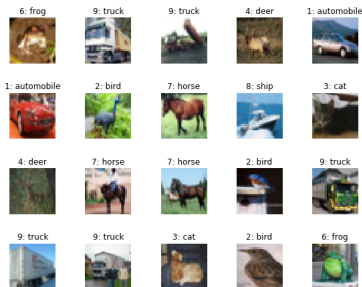
Fashion-Mnist



CIFAR 10



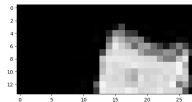
Fashion-Mnist



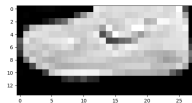
CIFAR 10

Task 1

- **Supervised Classification: 10 classes.**



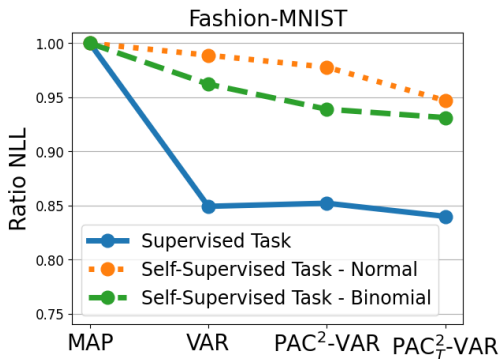
x



y

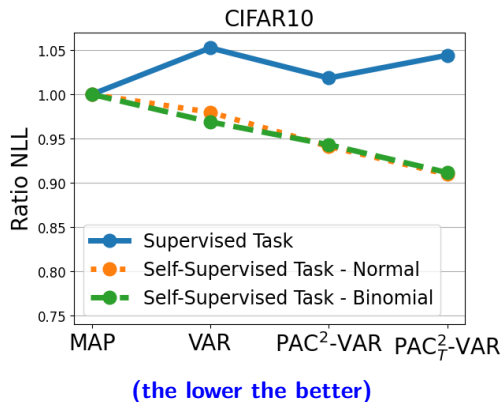
Self-Supervised Classification

- **Task 2** as a regression/Normal data model.
- **Task 3** as a Binomial data model.

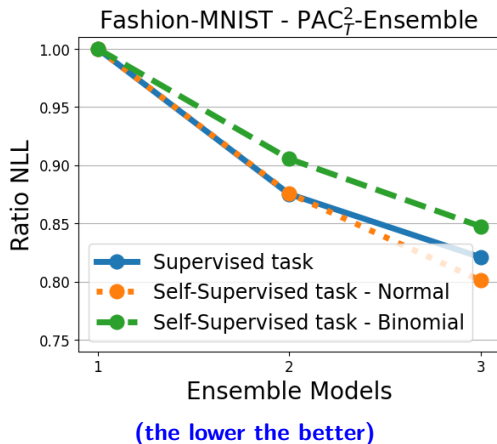


(the lower the better)

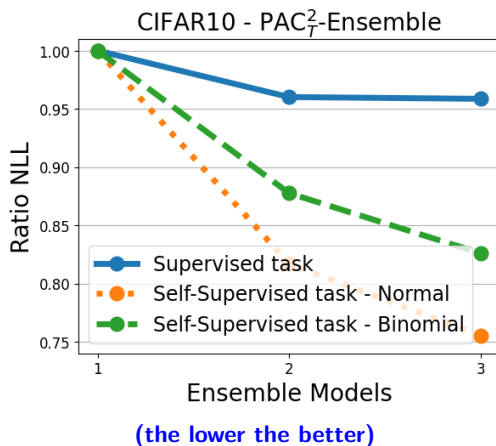
- MLP model with 20 hidden units, Relu activation.
- 100 data batches, 100 epochs, AdamOptimizer default learning rate.



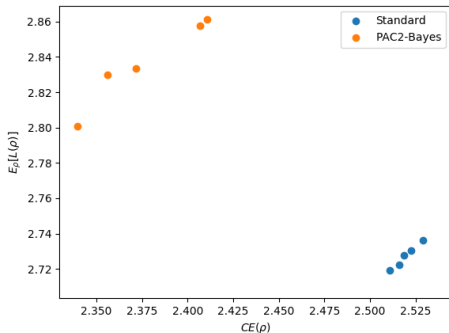
- MLP model with 20 hidden units, Relu activation.
- 100 data batches, 100 epochs, AdamOptimizer default learning rate.



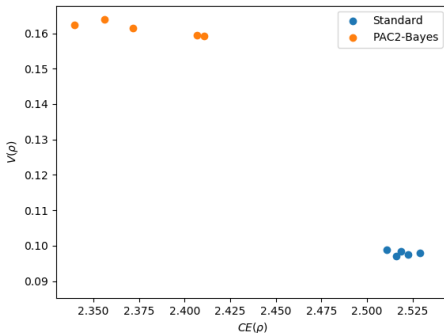
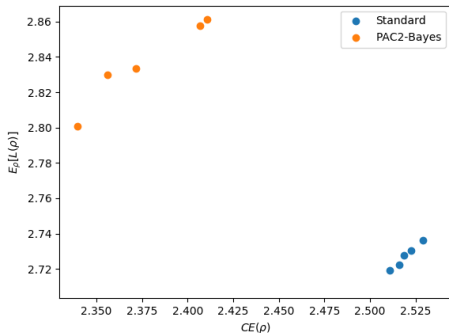
- **Models initialized with the same parameters.**
- MLP model with 20 hidden units, Relu activation.
- 100 data batches, 100 epochs, AdamOptimizer default learning rate.



- **Models initialized with the same parameters.**
- MLP model with 20 hidden units, Relu activation.
- 100 data batches, 100 epochs, AdamOptimizer default learning rate.



- Ensembles of four LeNet5 neural network on CIFAR-100



- Ensembles of four LeNet5 neural network on CIFAR-100

Conclusions and Future Works

- The Bayesian approach seems to be **not optimal strategy for learning**.

- The Bayesian approach seems to be **not optimal strategy for learning**.
- **Second-order PAC-Bayesian bounds** directly address mode misspecification.

- The Bayesian approach seems to be **not optimal strategy for learning**.
- **Second-order PAC-Bayesian bounds** directly address mode misspecification.
- Novel **variational and ensemble learning algorithms**.

- The Bayesian approach seems to be **not optimal strategy for learning**.
- **Second-order PAC-Bayesian bounds** directly address mode misspecification.
- Novel **variational and ensemble learning algorithms**.
- **Future works:**
 - Extensive empirical evaluation (new SOTA results in Bayesian deep learning?).

- The Bayesian approach seems to be **not optimal strategy for learning**.
- **Second-order PAC-Bayesian bounds** directly address mode misspecification.
- Novel **variational and ensemble learning algorithms**.
- **Future works:**
 - Extensive empirical evaluation (new SOTA results in Bayesian deep learning?).
 - What happens at the **interpolation regime**? In this case, $\mathbb{V}(\rho, D) = 0$

- The Bayesian approach seems to be **not optimal strategy for learning**.
- **Second-order PAC-Bayesian bounds** directly address mode misspecification.
- Novel **variational and ensemble learning algorithms**.
- **Future works:**
 - Extensive empirical evaluation (new SOTA results in Bayesian deep learning?).
 - What happens at the **interpolation regime**? In this case, $\mathbb{V}(\rho, D) = 0$
- Related work on Majority Voting:

Masegosa, A. R., Lorenzen, S. S., Igel, C., & Seldin, Y. Second order PAC-Bayesian bounds for the weighted majority vote. NeurIPS 2020.

- The Bayesian approach seems to be **not optimal strategy for learning**.
- **Second-order PAC-Bayesian bounds** directly address mode misspecification.
- Novel **variational and ensemble learning algorithms**.
- **Future works:**
 - Extensive empirical evaluation (new SOTA results in Bayesian deep learning?).
 - What happens at the **interpolation regime**? In this case, $\mathbb{V}(\rho, D) = 0$
- Related work on Majority Voting:

Masegosa, A. R., Lorenzen, S. S., Igel, C., & Seldin, Y. Second order PAC-Bayesian bounds for the weighted majority vote. NeurIPS 2020.

<https://github.com/PGM-Lab/PAC2BAYES>