Learning under model misspecification Applications to Variational and Ensemble methods

Andrés R. Masegosa

University of Almería Spain

November 25, 2020

Introduction

Model Misspecification

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

Model Misspecification

- Virtually any model we use **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 misspecfication.
- 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 prefect).
- $p(\mathbf{x}|\boldsymbol{\theta})$ is our (parametric) probabilistic model class.

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

Assumption 3: Likelihood is Upper-Bounded

 $\bullet\,$ There exists a M>0

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

Assumption 3: Likelihood is Upper-Bounded

• There exists a M > 0

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

• 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 $\epsilon > 0$),.

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

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

- **Notation**: $\rho(\theta)$ is a probability distribution over the parameters of my model.
- The predictive posterior distribution for a given $\rho(\theta)$,

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

• The learning problem is defined as,

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

- **Notation**: $\rho(\theta)$ is a probability distribution over the parameters of my model.
- The predictive posterior distribution for a given $\rho(\theta)$,

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

• The learning problem is defined as,

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

• ... is equivalent to:

$$\rho^{\star} = \arg\min_{\rho} \underbrace{\mathbb{E}_{\nu(\mathbf{x})}[-\ln \mathbb{E}_{\rho(\theta)}[p(\mathbf{x}|\theta)]]}_{CE(\rho)}$$

- **Notation**: $\rho(\theta)$ is a probability distribution over the parameters of my model.
- The predictive posterior distribution for a given $\rho(\theta)$,

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

• The learning problem is defined as,

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

• ... is equivalent to:

$$\rho^{\star} = \arg\min_{\rho} \underbrace{\mathbb{E}_{\nu(\mathbf{x})}[-\ln \mathbb{E}_{\rho(\theta)}[p(\mathbf{x}|\theta)]]}_{CE(\rho)}$$

• $CE(\rho)$ measures the generalization error (or the predicitive risk) associated to ρ .

• How to solve this problem

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

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

• How to solve this problem

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

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

The learning strategy

• The solution is to employ upper-bounds:

$$CE(\rho) \leq Oracle-Bound(\rho) \leq mpirical-Bound(\rho, D, \xi)$$

Jensen inequality

• How to solve this problem

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

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

The learning strategy

• The solution is to employ upper-bounds:

$$CE(\rho) \leq Oracle-Bound(\rho) \leq Empirical-Bound(\rho, D, \xi)$$

• ... and **minimize** Empirical-Bound(ρ , D, ξ),

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

• How to solve this problem

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

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

The learning strategy

• The solution is to employ upper-bounds:

$$CE(\rho) \leq Oracle-Bound(\rho) \leq mpirical-Bound(\rho, D, \xi)$$

• ... and minimize Empirical-Bound(ρ , D, ξ),

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

• The quality of the solution is going to depend of the quality of the bounds.

First-order Jensen bounds and the Bayesian posterior



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



 $L(\boldsymbol{\theta})$ 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, $L(\boldsymbol{\theta}, D) = -\frac{1}{n} \ln p(D|\boldsymbol{\theta})$.

The Bayesian posterior (Germain et al. 2016)

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

$$\rho^{\star} = \arg\min_{\rho} \mathbb{E}_{\rho(\theta)}[L(\theta, D)] + \frac{KL(\rho, \pi)}{n} + cte$$

PAC-Bayes bound (Alquier et al. 2016)

The Bayesian posterior (Germain et al. 2016)

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

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

Evidence Lower Bound (ELBO)

The Bayesian posterior (Germain et al. 2016)

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

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

• ρ^{\star} is the **Bayesian posterior**,

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

The Bayesian posterior (Germain et al. 2016)

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

$$\rho^{\star} = \arg \min_{\rho} \underbrace{\mathbb{E}_{\rho(\theta)}[L(\theta, D)] + \frac{KL(\rho, \pi)}{n} + cte}_{\text{PAC-Bayes bound (Alquier et al. 2016)}}$$
$$= \arg \max_{\rho} \underbrace{\mathbb{E}_{\rho(\theta)}[\ln p(D|\theta)] - KL(\rho, \pi)}_{\text{Evidence I over Bound (FLBO)}}$$

• ρ^{\star} is the **Bayesian posterior**,

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

The Bayesian posterior is a proxy

$$p(\boldsymbol{\theta}|D) \approx \arg\min_{\boldsymbol{\rho}} KL(\boldsymbol{\nu}(\mathbf{x}), \mathbb{E}_{\boldsymbol{\rho}(\boldsymbol{\theta})}[\boldsymbol{p}(\mathbf{x}|\boldsymbol{\theta})])$$

Data distribution Predictive posterior

• The Bayesian learning strategy,



• The Bayesian learning strategy,



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

• The Bayesian learning strategy,



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

• The Bayesian learning strategy,



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

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

• The Bayesian learning strategy,



- The Bayesian posterior converges to the minimum of $\mathbb{E}_{\rho}[L(\theta)]$.
- The minimum of the first-order Jensen bound $\mathbb{E}_{\rho}[L(\theta)]$ is

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

Is the Bayesian approach an optimal learning strategy?

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

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

Data distribution Predictive posterior





Perfect Model Specification





Bayesian posterior not optimal under misspecification



Model Misspecification



Generalization Error

Second-order Jensen bound

Bayesian posterior not optimal under misspecification



Learning by Minimizing second-order PAC-Bayes bounds

• A variational-like method,

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

Second-order PAC-Bayes Bound

• A variational-like method,

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

Second-order PAC-Bayes Bound

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

• A variational-like method,

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

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).

• A variational-like method,

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

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).

• A variational-like method,

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

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,

$$\arg\min_{\rho\in Q} \underbrace{\mathbb{E}_{\rho(\theta)}[L(\theta,D)] + \frac{KL(\rho,\pi)}{n} + cte}_{n}$$

First-order PAC-Bayes Bound

Ensemble Learning (Finite Mixture Models)

Mixture of Dirac-delta distribution

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

$$\boldsymbol{\rho}(\boldsymbol{\theta}) = \sum_{j=1}^{E} \frac{1}{E} \delta_{\boldsymbol{\theta}_{j}}(\boldsymbol{\theta})$$

where $\delta_{\boldsymbol{\theta}_{i}}$ is a **Dirac-delta distribution** centered around $\boldsymbol{\theta}_{j}$

Mixture of Dirac-delta distribution

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

$$\boldsymbol{\rho}(\boldsymbol{\theta}) = \sum_{j=1}^{E} \frac{1}{E} \delta_{\boldsymbol{\theta}_{j}}(\boldsymbol{\theta})$$

where $\delta_{\boldsymbol{\theta}_{i}}$ is a **Dirac-delta distribution** centered around $\boldsymbol{\theta}_{j}$

• The predictive posterior is defined as

$$p_E(\mathbf{x}) = \mathbb{E}_{
ho(heta)}[p(\mathbf{x}|oldsymbol{ heta})] = rac{1}{E}\sum_{j=1}^E p(\mathbf{x}|oldsymbol{ heta}_j)$$

Mixture of Dirac-delta distribution

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

$$\boldsymbol{\rho}(\boldsymbol{\theta}) = \sum_{j=1}^{E} \frac{1}{E} \delta_{\boldsymbol{\theta}_{j}}(\boldsymbol{\theta})$$

where $\delta_{\boldsymbol{\theta}_{j}}$ is a **Dirac-delta distribution** centered around $\boldsymbol{\theta}_{j}$

• The predictive posterior is defined as

$$p_E(\mathbf{x}) = \mathbb{E}_{\rho(\theta)}[p(\mathbf{x}|\boldsymbol{\theta})] = \frac{1}{E} \sum_{j=1}^{E} p(\mathbf{x}|\boldsymbol{\theta}_j)$$

• The learning problem,

$$\{\boldsymbol{\theta}_1^\star,\ldots,\boldsymbol{\theta}_E^\star\} = rg\min_{\{\theta_1,\ldots,\theta_E\}} KL(\nu(\mathbf{x}),p_E(\mathbf{x}))$$

Experimental Evaluation with Toy Data Sets

$$\begin{aligned}
\boldsymbol{\nu}(\boldsymbol{y}|\boldsymbol{x}) &= \mathcal{N}(\boldsymbol{\mu} = 1 + \boldsymbol{x}, \sigma^2 = 5) \\
\boldsymbol{p}(\boldsymbol{y}|\boldsymbol{x}, \boldsymbol{\theta}) &= \mathcal{N}(\boldsymbol{\mu} = \theta_0 + \theta_1 \boldsymbol{x}, \sigma^2 = 1) \\
\rho(\boldsymbol{\theta}) &= \mathcal{M}\mathcal{V}\mathcal{N}(\boldsymbol{\mu}, \Sigma)
\end{aligned}$$

Misspecified Linear Regression Model



Test Log-likelihood=-13.09

$$\begin{aligned} \boldsymbol{\nu}(\boldsymbol{y}|\boldsymbol{x}) &= \mathcal{N}(\boldsymbol{\mu} = 1 + \boldsymbol{x}, \sigma^2 = 5) \\ \boldsymbol{p}(\boldsymbol{y}|\boldsymbol{x}, \boldsymbol{\theta}) &= \mathcal{N}(\boldsymbol{\mu} = \theta_0 + \theta_1 \boldsymbol{x}, \sigma^2 = 1) \\ \boldsymbol{\rho}(\boldsymbol{\theta}) &= \mathcal{M}\mathcal{V}\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \end{aligned}$$

Misspecified Linear Regression Model







Test Log-likelihood=-7.89

$$\nu(y|x) = \mathcal{N}(\mu = 1 + x, \sigma^2 = 5)$$

$$p(y|x, \theta) = \mathcal{N}(\mu = \theta_0 + \theta_1 x, \sigma^2 = 1)$$

$$\rho(\theta) = \mathcal{MVN}(\mu, \Sigma)$$



Test Log-likelihood = -50.15

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



Test Log-likelihood = -50.15



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



Test Log-likelihood = -15.91

$$\nu(y|x) = \mathcal{N}(\mu = s(x), \sigma^2 = 10)$$

$$p(y|x, \theta) = \mathcal{N}(\mu = MLP_{20}(x; \theta), \sigma^2 = 1)$$

$$\rho(\theta) = \sum_{i=1}^{3} \frac{1}{E} \delta_{\theta_i}(\theta)$$

Experimental Evaluation on real data sets



Fahsion-Mnist



CIFAR 10



Task 1

• Supervised Classification: 10 classes.



Self-Supervised Classification

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



- 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.

• 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?

- 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?
 - 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?
 - 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
```