Bayesian Priors and Generalization in Probabilistic Machine Learning

 $\begin{array}{c} \text{Andr\'es R. Masegosa} \\ \text{(joint work with Luis A. Ortega (UPM, Spain))} \end{array}$

Aalborg University Copenhagen Campus Denmark

Danish Data Science 2022

Bayesian machine learning

Bayesian machine learning

- Bayesian methods are widely used in machine learning.
- They provide well founded approach for dealing with model uncertainty.
- Random variables + Probability Calculus.
- They automatically account for model complexity.
- They allow to combine data with prior knowledge.

Probabilistic Model

- Conditional generative models: $p(y|\mathbf{x}, \theta)$.
- Generative models: $p(\mathbf{x}|\boldsymbol{\theta})$.

Probabilistic Model

- Conditional generative models: $p(y|\mathbf{x}, \boldsymbol{\theta})$.
- Generative models: $p(\mathbf{x}|\boldsymbol{\theta})$.

Bayesian Posterior

$$\underbrace{p(\boldsymbol{\theta}|D)}_{\text{Bayesian posterior}} = \underbrace{\frac{\sum_{\text{Likelihood}}^{\text{Likelihood}}}_{\text{Normalization Constant}} \underbrace{\frac{p(\boldsymbol{\theta}|D)}{\pi(\boldsymbol{\theta})d\boldsymbol{\theta}}}_{\text{Prior}}$$

We have to resort to approximations to compute the integral.

Probabilistic Model

- Conditional generative models: $p(y|\mathbf{x}, \boldsymbol{\theta})$.
- Generative models: $p(\mathbf{x}|\boldsymbol{\theta})$.

Bayesian Posterior

$$\underbrace{p(\boldsymbol{\theta}|D)}_{\text{Bayesian posterior}} = \underbrace{\frac{\sum_{\text{Likelihood}}^{\text{Likelihood}}}_{\text{Normalization Constant}} \underbrace{\frac{p(\boldsymbol{\theta}|D)}{\pi(\boldsymbol{\theta})d\boldsymbol{\theta}}}_{\text{Prior}} \underbrace{\frac{p(\boldsymbol{\theta}|D)}{\pi(\boldsymbol{\theta})d\boldsymbol{\theta}}}_{\text{Prior}}$$

We have to resort to approximations to compute the integral.

Bayesian model averaging

$$\underbrace{p(y_{\text{test}} \mid \mathbf{x}_{\text{test}}, D)}_{\text{Predictive posterior}} = \int \underbrace{p(y_{\text{test}} \mid x_{\text{test}}, \boldsymbol{\theta})}_{\text{Bayesian posterior}} \underbrace{p(\boldsymbol{\theta} | D)}_{\text{Bayesian posterior}} d\boldsymbol{\theta}$$

Bayesian Priors in Bayesian Machine Learning (work in progress)

Bayesian Priors

Bayesian Priors in Bayesian Statistics

- (Weakly) Informative Priors
 - Priors providing information about the data generating process.
- (Non-informative) Reference Priors
 - Priors minimizing the impact they have in the Bayesian posterior.

Bayesian Priors

Bayesian Priors in Bayesian Statistics

- (Weakly) Informative Priors
 - Priors providing information about the data generating process.
- (Non-informative) Reference Priors
 - Priors minimizing the impact they have in the Bayesian posterior.

Bayesian Priors in Machine Learning

- Regularizing Priors (e.g., zero centered Gaussian distributions)
 - Promote small norm parameter that reduce overfitting.
 - Overwhelming empirical evidence.
 - Connections with L2, L1, etc. through MAP learning.
- What are Regularizing Priors?
 - Reference priors, (Weakly) Informative priors or something different.
- How a Bayesian prior should look like to guarantee generalization performance?

Main Conclusion (Take Away Message)

• Regularizing Priors introduce a biased against high-variance models.

Main Conclusion (Take Away Message)

- Regularizing Priors introduce a biased against high-variance models.
- For a fixed θ ,
 - We treat D as a random variable, $D \sim \nu(y, \mathbf{x})$ (and $D_{\mathsf{Test}} \sim \nu(y, \mathbf{x})$).
 - Training Loss: $\hat{L}(D, \theta) = -\frac{1}{n} \ln p(D|\theta)$.
 - Expected Loss: $L(\theta) = \mathbb{E}_D[-\frac{1}{n}\ln p(D|\theta)] = \mathbb{E}_{\nu}[-\ln p(y \mid \mathbf{x}, \theta)].$
 - Variance of the training Loss: $\mathbb{V}(\boldsymbol{\theta}) = \mathbb{V}_{\nu}(\frac{1}{n} \ln p(D|\boldsymbol{\theta}))$.

Main Conclusion (Take Away Message)

- Regularizing Priors introduce a biased against high-variance models.
- For a fixed θ ,
 - We treat D as a random variable, $D \sim \nu(y, \mathbf{x})$ (and $D_{\mathsf{Test}} \sim \nu(y, \mathbf{x})$).
 - Training Loss: $\hat{L}(D, \theta) = -\frac{1}{n} \ln p(D|\theta)$.
 - Expected Loss: $L(\theta) = \mathbb{E}_D[-\frac{1}{n}\ln p(D|\theta)] = \mathbb{E}_{\nu}[-\ln p(y \mid \mathbf{x}, \theta)].$
 - Variance of the training Loss: $\mathbb{V}(\theta) = \mathbb{V}_{\nu}(\frac{1}{n} \ln p(D|\theta))$.
- If $\mathbb{V}(\theta)$ is large, then θ is an unreliable model.
 - Training loss can be small and Expected loss can be large.

Main Conclusion (Take Away Message)

- Regularizing Priors introduce a biased against high-variance models.
- For a fixed θ .
 - We treat D as a random variable, $D \sim \nu(y, \mathbf{x})$ (and $D_{\mathsf{Test}} \sim \nu(y, \mathbf{x})$).
 - Training Loss: $\hat{L}(D, \theta) = -\frac{1}{n} \ln p(D|\theta)$.
 - Expected Loss: $L(\theta) = \mathbb{E}_D[-\frac{1}{n}\ln p(D|\theta)] = \mathbb{E}_{\nu}[-\ln p(y \mid \mathbf{x}, \theta)].$
 - Variance of the training Loss: $\mathbb{V}(\theta) = \mathbb{V}_{\nu}(\frac{1}{n} \ln p(D|\theta))$.
- If $\mathbb{V}(\theta)$ is large, then θ is an unreliable model.
 - Training loss can be small and Expected loss can be large.
- Generalization requires penalizing high-variance models.
 - This should be encoded in the prior:

$$\pi_J(\boldsymbol{\theta}) \propto e^{-\mathbb{V}(\boldsymbol{\theta})}$$

• Connected to Gaussian zero-centered priors and other regularizing priors.

Generalization Analysis of Bayesian learning

PAC-Bayes Bounds

- ullet Frequentist analysis of Bayesian methods (D is **treated as a random quantity**).
- Theoretical tool to analyse the relationship between training loss and test loss:

 $\mathsf{Expected}\ \mathsf{Loss} \leq \mathsf{Train}\ \mathsf{Loss} + \mathsf{Complexity}$

Generalization Analysis of Bayesian learning

PAC-Bayes Bounds

- Frequentist analysis of Bayesian methods (D is treated as a random quantity).
- Theoretical tool to analyse the relationship between training loss and test loss:

Expected Loss
$$\leq$$
 Train Loss $+$ Complexity

• Expected Loss of Bayesian learning:

$$\underbrace{CE(p(\boldsymbol{\theta}|D))}_{\text{Bayesian Expected Loss}} = \mathbb{E}_{\nu(\mathbf{y},\mathbf{x})}[-\ln\underbrace{\mathbb{E}_{p(\boldsymbol{\theta}|D)}[p(\mathbf{y}|\mathbf{x},\boldsymbol{\theta})]]}_{\text{Bayesian Model Averaging}}$$

• $\nu(\mathbf{y}, \mathbf{x})$ is the data-generating distribution and $D \sim \nu(\mathbf{y}, \mathbf{x})$.

PAC-Bayesian Bound (Alquier et al. 2016, Germain et al. 2016, Masegosa 2020)

• For any prior $\pi(\theta)$ independent of D and any $\lambda > 0$,

$$\underbrace{CE(\frac{\mathbf{p}_{\pi}^{\lambda}}{p_{\pi}})}_{\text{Bayesian Expected Loss}} \underbrace{\lesssim}_{\text{VPAC-Bayes bound}} - \underbrace{\frac{\hat{LM}_{\lambda}(\pi,D)}{n\lambda} + \frac{R_{\lambda}(\pi)}{\lambda n} + \frac{\ln\frac{1}{\delta}}{\lambda n}}_{\text{PAC-Bayes bound}}$$

PAC-Bayesian Bound (Alquier et al. 2016, Germain et al. 2016, Masegosa 2020)

• For any prior $\pi(\theta)$ independent of D and any $\lambda > 0$,

$$\underbrace{CE(p_{\pi}^{\lambda})}_{\text{Bayesian Expected Loss}} \lesssim \underbrace{-\frac{\hat{LM}_{\lambda}(\pi,D)}{n\lambda} + \frac{R_{\lambda}(\pi)}{\lambda n} + \frac{\ln\frac{1}{\delta}}{\lambda n}}_{\text{PAC-Bayes bound}}$$

where p_{π}^{λ} denotes the **generalized Bayesian posterior**,

$$p_{\pi}^{\lambda}(\boldsymbol{\theta}|D) \propto p(D|\boldsymbol{\theta})^{\lambda}\pi(\boldsymbol{\theta})$$

PAC-Bayesian Bound (Alquier et al. 2016, Germain et al. 2016, Masegosa 2020)

• For any prior $\pi(\theta)$ independent of D and any $\lambda > 0$,

$$\underbrace{CE(p_{\pi}^{\lambda})}_{\text{Bayesian Expected Loss}} \underbrace{\lesssim \underbrace{-\frac{\hat{LM}_{\lambda}(\pi,D)}{n\lambda} + \frac{R_{\lambda}(\pi)}{\lambda n} + \frac{\ln\frac{1}{\delta}}{\lambda n}}_{\text{PAC-Bayes bound}}$$

where p_{π}^{λ} denotes the **generalized Bayesian posterior**,

$$p_{\pi}^{\lambda}(\boldsymbol{\theta}|D) \propto p(D|\boldsymbol{\theta})^{\lambda}\pi(\boldsymbol{\theta})$$

where $\hat{LM}_{\lambda}(\pi, D)$ denotes the **log-marginal**:

$$\hat{LM}_{\lambda}(\pi, D) = \ln \mathbb{E}_{\pi}[p(D|\boldsymbol{\theta})^{\lambda}]$$

PAC-Bayesian Bound (Alquier et al. 2016, Germain et al. 2016, Masegosa 2020)

• For any prior $\pi(\theta)$ independent of D and any $\lambda > 0$,

$$\underbrace{CE(p_{\pi}^{\lambda})}_{\text{Bayesian Expected Loss}} \underbrace{\lesssim \underbrace{-\frac{\hat{LM}_{\lambda}(\pi,D)}{n\lambda} + \frac{R_{\lambda}(\pi)}{\lambda n} + \frac{\ln\frac{1}{\delta}}{\lambda n}}_{\text{PAC-Bayes bound}}$$

where p_{π}^{λ} denotes the **generalized Bayesian posterior**,

$$p_{\pi}^{\lambda}(\boldsymbol{\theta}|\boldsymbol{D}) \propto p(\boldsymbol{D}|\boldsymbol{\theta})^{\lambda}\pi(\boldsymbol{\theta})$$

where $\hat{LM}_{\lambda}(\pi, D)$ denotes the **log-marginal**:

$$\hat{LM}_{\lambda}(\pi, D) = \ln \mathbb{E}_{\pi}[p(D|\boldsymbol{\theta})^{\lambda}]$$

where $R_{\lambda}(\pi)$ is a **cummulant generating function**, which can be expressed as:

$$R_{\lambda}(\pi) = \ln \mathbb{E}_{\pi D} \left[e^{\lambda n(L(\theta) - \hat{L}(\theta, D))} \right]$$

Upper Bounds

ullet PAC-Bayesian bound: For any prior $\pi(oldsymbol{ heta})$ independent of D and any $\lambda>0$, ,

$$\underbrace{CE(p_{\pi}^{\lambda})}_{\text{Bayesian Expected Loss}} \lesssim \underbrace{-\frac{\hat{LM}_{\lambda}(\pi,D)}{n\lambda} + \frac{R_{\lambda}(\pi)}{\lambda n} + \frac{\ln\frac{1}{\delta}}{\lambda n}}_{\text{PAC-Bayes bound}}$$

Upper Bounds

• PAC-Bayesian bound: For any prior $\pi(\theta)$ independent of D and any $\lambda>0$, ,

$$\underbrace{CE(p_{\pi}^{\lambda})}_{\text{Bayesian Expected Loss}} \underbrace{\lesssim}_{\text{W.p. } (1-\delta)} \underbrace{-\frac{\hat{LM}_{\lambda}(\pi,D)}{n\lambda} + \frac{R_{\lambda}(\pi)}{\lambda n} + \frac{\ln\frac{1}{\delta}}{\lambda n}}_{\text{PAC-Bayes bound}}$$

• Expectation bound: In expectation over different data samples D,

$$\underbrace{\mathbb{E}_D[CE(p_\pi^\lambda)]}_{\text{Bayesian Expected Loss}} \leq \underbrace{-\frac{\mathbb{E}_D[\hat{LM}_\lambda(\pi,D)]}{n\lambda} + \frac{R_\lambda(\pi)}{\lambda n}}_{\text{Deterministic bound}}$$

Upper Bounds

ullet PAC-Bayesian bound: For any prior $\pi(oldsymbol{ heta})$ independent of D and any $\lambda>0$, ,

$$\underbrace{CE(p_{\pi}^{\lambda})}_{\text{Bayesian Expected Loss}} \lesssim \underbrace{-\frac{\hat{LM}_{\lambda}(\pi,D)}{n\lambda} + \frac{R_{\lambda}(\pi)}{\lambda n} + \frac{\ln\frac{1}{\delta}}{\lambda n}}_{\text{PAC-Bayes bound}}$$

• Expectation bound: In expectation over different data samples D,

$$\underbrace{\mathbb{E}_D[CE(p_\pi^\lambda)]}_{\text{Bayesian Expected Loss}} \leq \underbrace{-\frac{\mathbb{E}_D[\hat{LM}_\lambda(\pi,D)]}{n\lambda} + \frac{R_\lambda(\pi)}{\lambda n}}_{\text{Deterministic bound}}$$

- According to these bounds, small predictive loss is attained if:
 - $-L\hat{M}_{\lambda}(\pi,D)$ and $R_{\lambda}(\pi)$ are both small.
 - Both depends on the prior $\pi(\theta)$.
 - Which priors $\pi(\theta)$ make these two terms small?

The Log-Marginal Likelihood

$$\hat{LM}_{\lambda}(\pi, D) = \ln \mathbb{E}_{\pi}[p(D|\boldsymbol{\theta})^{\lambda}]$$

- Widely used in **Bayesian model comparison**.
- Measures how well our model class explains the data.
- Depends on the prior $\pi(\theta)$.

The Log-Marginal Likelihood

$$\hat{LM}_{\lambda}(\pi, D) = \ln \mathbb{E}_{\pi}[p(D|\boldsymbol{\theta})^{\lambda}]$$

- Widely used in **Bayesian model comparison**.
- Measures how well our model class explains the data.
- Depends on the prior $\pi(\theta)$.

Theorem: Informative Priors improves the log-marginal likelihood

- Let $\pi_0(\boldsymbol{\theta})$ be a flat or reference prior.
- We build an informative prior using (expected) Bayesian updating:

$$\pi_I(\boldsymbol{\theta}) = \mathbb{E}_{D' \sim \nu^n} [p_{\pi_0}^{\lambda}(\boldsymbol{\theta}|D')]$$

The Log-Marginal Likelihood

$$\hat{LM}_{\lambda}(\pi, D) = \ln \mathbb{E}_{\pi}[p(D|\boldsymbol{\theta})^{\lambda}]$$

- Widely used in **Bayesian model comparison**.
- Measures how well our model class explains the data.
- Depends on the prior $\pi(\theta)$.

Theorem: Informative Priors improves the log-marginal likelihood

- Let $\pi_0(\boldsymbol{\theta})$ be a flat or reference prior.
- We build an informative prior using (expected) Bayesian updating:

$$\pi_I(\boldsymbol{\theta}) = \mathbb{E}_{D' \sim \nu^n} [p_{\pi_0}^{\lambda}(\boldsymbol{\theta}|D')]$$

Then, we have that

$$\mathbb{E}_{D \sim \nu^n} [-\hat{L} M_\lambda(\pi_I, D)] \le \mathbb{E}_{D \sim \nu^n} [-\hat{L} M_\lambda(\pi_0, D)]$$

• Informative priors reduce, in expectation, the negative log-marginal likelihood.

Upper Bounds

• PAC-Bayesian bound: For any prior $\pi(\theta)$ independent of D and any $\lambda > 0$, ,

$$\underbrace{CE(\frac{\mathbf{p}_{\pi}^{\lambda}}{\rho_{\pi}})}_{\text{Bayesian Expected Loss}} \underbrace{\sum_{\mathbf{p} \in \mathbb{Z}} \underbrace{-\frac{\hat{L}\hat{M}_{\lambda}(\pi,D)}{n\lambda} + \frac{R_{\lambda}(\pi)}{\lambda n} + \frac{\ln\frac{1}{\delta}}{\lambda n}}_{\mathbf{p} \in \mathbb{Z}}$$

• Expectation bound: In expectation over different data samples D,

$$\underbrace{\mathbb{E}_D[CE(\mathbf{p}_{\pi}^{\lambda})]}_{\text{Bayesian Expected Loss}} \leq \underbrace{-\frac{\mathbb{E}_D[\hat{LM}_{\lambda}(\pi,D)]}{n\lambda} + \frac{R_{\lambda}(\pi)}{\lambda n}}_{\text{Deterministic bound}}$$

Upper Bounds

• PAC-Bayesian bound: For any prior $\pi(\theta)$ independent of D and any $\lambda > 0$, ,

$$\underbrace{CE(p_{\pi}^{\lambda})}_{\text{Bayesian Expected Loss}} \underbrace{\lesssim}_{\text{W.p. } (1-\delta)} \underbrace{-\underbrace{\hat{LM}_{\lambda}(\pi,D)}_{n\lambda} + \underbrace{\frac{R_{\lambda}(\pi)}{\lambda n} + \frac{\ln\frac{1}{\delta}}{\lambda n}}_{\text{PAC-Bayes bound}}$$

• Expectation bound: In expectation over different data samples D,

$$\underbrace{\mathbb{E}_D[CE(\textbf{p}_{\pi}^{\lambda})]}_{\text{Bayesian Expected Loss}} \leq \underbrace{-\frac{\mathbb{E}_D[\hat{LM}_{\lambda}(\pi,D)]}{n\lambda} + \frac{R_{\lambda}(\pi)}{\lambda n}}_{\text{Deterministic bound}}$$

- Informative priors reduce the $\hat{LM}_{\lambda}(\pi,D)$ term:
 - But not enough to guarantee generalization performance.
 - Which priors reduce the $R_{\lambda}(\pi)$ term?

Proposition: $R_{\lambda}(\pi)$ is a prior regularizer

Over joint draws of $\theta \sim \pi(\theta)$ and $D \sim \nu^n(\mathbf{x}, \mathbf{y})$, we have that

$$\underbrace{L(\boldsymbol{\theta}) - \hat{L}(\boldsymbol{\theta}, D)}_{\text{Overfitting}} \lesssim \frac{1}{\lambda n} R_{\lambda}(\pi) + \frac{1}{\lambda n} \ln \frac{1}{\delta}. \tag{1}$$

• If $R_{\lambda}(\pi)$ is small, then $\pi(\theta)$ prefers models with small overfitting.

Proposition: $R_{\lambda}(\pi)$ is a prior regularizer

Over joint draws of $\theta \sim \pi(\theta)$ and $D \sim \nu^n(\mathbf{x}, \mathbf{y})$, we have that

$$\underbrace{L(\boldsymbol{\theta}) - \hat{L}(\boldsymbol{\theta}, D)}_{\text{Overfitting}} \lesssim \frac{1}{\lambda n} R_{\lambda}(\pi) + \frac{1}{\lambda n} \ln \frac{1}{\delta}. \tag{1}$$

• If $R_{\lambda}(\pi)$ is small, then $\pi(\theta)$ prefers models with small overfitting.

Proposition: $R_{\lambda}(\pi)$ is a prior regularizer

- $R_{\lambda}(\pi) \geq 0$ for any prior $\pi(\theta)$ and any $\lambda \geq 0$.
- $R_{\lambda}(\pi) = 0$ iif $\pi(\theta)$ is Dirac-Delta distribution around θ_0 ,

$$\underbrace{L(\boldsymbol{\theta}_0) - \hat{L}(\boldsymbol{\theta}_0, D)}_{\text{Overfitting}} = 0$$

• E.g., A neural network with all the weights set to zero.

The Information-Regularization Trade-off

• PAC-Bayesian bound: For any prior $\pi(\theta)$ independent of D and any $\lambda > 0$, ,

$$\underbrace{CE(p_{\pi}^{\lambda})}_{\text{Bayesian Expected Loss}} \underbrace{\lesssim}_{\text{W.p.}} \underbrace{\frac{(1-\delta)}{\sum_{j=1}^{N} \frac{1}{j}} \underbrace{-\frac{\hat{LM}_{\lambda}(\pi,D)}{n\lambda} + \frac{R_{\lambda}(\pi)}{\lambda n} + \frac{\ln\frac{1}{\delta}}{\lambda n}}_{\text{PAC-Bayes bound}}$$

• Expectation bound: In expectation over different data samples D,

$$\underbrace{\mathbb{E}_D[CE(p_\pi^\lambda)]}_{\text{Bayesian Expected Loss}} \leq \underbrace{-\frac{\mathbb{E}_D[\hat{LM}_\lambda(\pi,D)]}{n\lambda} + \frac{R_\lambda(\pi)}{\lambda n}}_{\text{Deterministic bound}}$$

- Priors minimizing these upper-bounds face a trade-off:
 - Informative priors reduce $\hat{LM}_{\lambda}(\pi, D)$.
 - Regularizing priors reduce $R_{\lambda}(\pi)$.

Theorem: Optimal Priors

- If $\pi_0(\boldsymbol{\theta})$ is a flat or reference prior.
- We define a new priors as:

$$\pi_1(m{ heta}) \propto \underbrace{\pi_I(m{ heta})}_{ ext{Informative Prior}} \underbrace{e^{-nJ_
u(m{ heta},\lambda)}}_{ ext{Regularizing Prior}}$$

where $J_{\nu}(\theta,\lambda)$ is the so-called **Jensen-Gap function**, defined as:

$$J_{\nu}(\theta, \lambda) = \ln \mathbb{E}_{\nu}[p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})] - \mathbb{E}_{\nu}[\ln p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})]$$

Theorem: Optimal Priors

- If $\pi_0(\boldsymbol{\theta})$ is a flat or reference prior.
- We define a new priors as:

$$\pi_1(oldsymbol{ heta}) \propto \underbrace{\pi_I(oldsymbol{ heta})}_{ ext{Informative Prior}} \underbrace{e^{-nJ_
u(heta,\lambda)}}_{ ext{Regularizing Prior}}$$

where $J_{\nu}(\theta,\lambda)$ is the so-called **Jensen-Gap function**, defined as:

$$J_{\nu}(\theta, \lambda) = \ln \mathbb{E}_{\nu}[p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})] - \mathbb{E}_{\nu}[\ln p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})]$$

Then, we have that

$$\underbrace{\mathbb{E}_D[CE(p_{\pi_1}^{\lambda})]}_{\text{Bayesian Expected Loss}} \leq \underbrace{-\frac{\mathbb{E}_D[\hat{LM}_{\lambda}(\pi_1, D)]}{n\lambda} + \frac{R_{\lambda}(\pi_1)}{\lambda n}}_{\text{Upper bound for } \pi_1(\theta)} \leq \underbrace{-\frac{\mathbb{E}_D[\hat{LM}_{\lambda}(\pi_0, D)]}{n\lambda} + \frac{R_{\lambda}(\pi_0)}{\lambda n}}_{\text{Upper bound for } \pi_0(\theta)}$$

PAC-Bayesian Analysis of Regularizing Priors

Regularizing Prior

• We define an Jensen-Gap prior:

$$\pi_J(\boldsymbol{\theta}) \propto e^{-nJ_{\nu}(\boldsymbol{\theta},\lambda)}$$

- Naturally emerges when minimizing a (PAC-Bayes) upper-bound over the Bayesian Expected Loss.
- **Proposition:** For any $\theta \in \Theta$, over random draws of $D \sim \nu^n(\mathbf{x}, \mathbf{y})$, we have that

$$\underbrace{L(\boldsymbol{\theta}) - \hat{L}(\boldsymbol{\theta}, D)}_{\text{Overfitting}} \lesssim \frac{1}{\lambda n} J_{\nu}(\boldsymbol{\theta}, \lambda) + \frac{1}{\lambda n} \ln \frac{1}{\delta}. \tag{2}$$

- $\pi_J(\theta)$ assigns low probability to models with high risk of overfitting.
- $\pi_J(\theta)$ addresses overfitting (i.e., a regularizing prior).
 - It is a functional prior.
 - The prior depends on $p(y \mid \mathbf{x}, \boldsymbol{\theta})$ and the $\nu(y, \mathbf{x})$.

$$\pi_J(oldsymbol{ heta})$$
 and existing regularization methods.

MAP estimate using $\pi_J(\boldsymbol{\theta})$

$$= \arg\min_{\boldsymbol{\theta}} \hat{L}(\boldsymbol{\theta}, D) - \underbrace{\frac{\ln \pi_{J}(\boldsymbol{\theta})}{\lambda n}}_{\text{log Prior}}$$

$$= \arg\min_{\boldsymbol{\theta}} \hat{L}(\boldsymbol{\theta}, D) + \underbrace{\frac{J_{\nu}(\boldsymbol{\theta}, \lambda)}{\lambda}}_{\text{Regularizer}}$$

 $\theta_{\mathsf{MAP}} = \arg\max_{a} p_{\pi_J}^{\lambda}(\boldsymbol{\theta}|D)$

 $\pi_J(\boldsymbol{\theta})$ and existing regularization methods.

MAP estimate using $\pi_J(\boldsymbol{\theta})$

$$\begin{split} \theta_{\mathsf{MAP}} &= \arg\max_{\boldsymbol{\theta}} p_{\pi_J}^{\lambda}(\boldsymbol{\theta}|D) \\ &= \arg\min_{\boldsymbol{\theta}} \hat{L}(\boldsymbol{\theta},D) - \underbrace{\frac{\ln\pi_J(\boldsymbol{\theta})}{\lambda n}}_{\text{log Prior}} \\ &= \arg\min_{\boldsymbol{\theta}} \hat{L}(\boldsymbol{\theta},D) + \underbrace{\frac{J_{\nu}(\boldsymbol{\theta},\lambda)}{\lambda}}_{\text{Regularizer}} \end{split}$$

$\pi_J(\boldsymbol{\theta})$ and existing regularization methods.

$\pi_J(oldsymbol{ heta})$ and frequentist estimation theory

Proposition: Under a 2nd-order Taylor approximation of $J_{\nu}(\theta, \lambda)$ wrt λ :

$$J_{\nu}(\boldsymbol{\theta}, \lambda) \approx \frac{\lambda^2}{2} \mathbb{V}_{D \sim \nu^n} \left(\hat{L}(\boldsymbol{\theta}, D) \right)$$

- Connection with frequentist estimation theory:
 - $\hat{L}(\theta, D)$ is an unbiased estimator of $L(\theta)$.
 - $\bullet \ \mathbb{V}_{D \sim \nu^n} \left(\hat{L}(oldsymbol{ heta}, D) \right)$ is the variance of the estimator.
- Regularization means preferring models with low variance.
 - ullet For low variance models, $\hat{L}(m{ heta},D)$ is a better estimator of $L(m{ heta}).$
- Existing literature: (Namkoong et al. 2017), (Xie et al., 2021), etc.

$\pi_J(oldsymbol{ heta})$ and L2 regularization (i.e., zero-centered Gaussian priors)

Proposition: For a logistic regression model and under a 2nd-order Taylor approximation of $J_{\nu}(\theta, \lambda)$ wrt θ :

$$J_{\nu}(\boldsymbol{\theta}, \lambda) \approx 0.25 \lambda^2 \boldsymbol{\theta}^T \mathsf{Cov}_{\nu}(y\mathbf{x}) \boldsymbol{\theta}$$

$\pi_J(oldsymbol{ heta})$ and L2 regularization (i.e., zero-centered Gaussian priors)

Proposition: For a logistic regression model and under a 2nd-order Taylor approximation of $J_{\nu}(\theta, \lambda)$ wrt θ :

$$J_{\nu}(\boldsymbol{\theta}, \lambda) \approx 0.25 \lambda^2 \boldsymbol{\theta}^T \mathsf{Cov}_{\nu}(y\mathbf{x}) \boldsymbol{\theta}$$

• $\pi_J(\theta)$ would be a multivariate normal distribution:

$$\pi_J(\boldsymbol{\theta}) \propto e^{-n0.25\lambda^2 \theta^T \mathsf{Cov}_{\nu}(y\mathbf{x})\theta}$$

$\pi_J(oldsymbol{ heta})$ and L2 regularization (i.e., zero-centered Gaussian priors)

Proposition: For a logistic regression model and under a 2nd-order Taylor approximation of $J_{\nu}(\theta, \lambda)$ wrt θ :

$$J_{\nu}(\boldsymbol{\theta}, \lambda) \approx 0.25 \lambda^2 \boldsymbol{\theta}^T \mathsf{Cov}_{\nu}(y\mathbf{x}) \boldsymbol{\theta}$$

• $\pi_J(\theta)$ would be a multivariate normal distribution:

$$\pi_J(\boldsymbol{\theta}) \propto e^{-n0.25\lambda^2 \theta^T \text{Cov}_{\nu}(y\mathbf{x})\theta}$$

 If the data is normalized and features are conditionally independent, it is equal to L2-regularization ,

$$\boldsymbol{\theta}^T \mathsf{Cov}_{\nu}(y\mathbf{x})\boldsymbol{\theta} = \boldsymbol{\theta}^T kI\boldsymbol{\theta} = k||\boldsymbol{\theta}||^2$$

$\pi_J(oldsymbol{ heta})$ and L2 regularization (i.e., zero-centered Gaussian priors)

Proposition: For a logistic regression model and under a 2nd-order Taylor approximation of $J_{\nu}(\theta,\lambda)$ wrt θ :

$$J_{\nu}(\boldsymbol{\theta}, \lambda) \approx 0.25 \lambda^2 \boldsymbol{\theta}^T \mathsf{Cov}_{\nu}(y\mathbf{x}) \boldsymbol{\theta}$$

• $\pi_J(\theta)$ would be a multivariate normal distribution:

$$\pi_J(\boldsymbol{\theta}) \propto e^{-n0.25\lambda^2 \theta^T \mathsf{Cov}_{\nu}(y\mathbf{x})\theta}$$

 If the data is normalized and features are conditionally independent, it is equal to L2-regularization ,

$$\boldsymbol{\theta}^T \mathsf{Cov}_{\nu}(y\mathbf{x})\boldsymbol{\theta} = \boldsymbol{\theta}^T kI\boldsymbol{\theta} = k||\boldsymbol{\theta}||^2$$

- Explains why L2-regularization improves generalization:
 - Small-norm models tends to have lower variance.
 - Lower variance implies better estimators $\hat{L}(D, \theta)$.
 - Better estimators leads to less overfitting.

$\pi_J(oldsymbol{ heta})$ and L2 regularization (i.e., zero-centered Gaussian priors)

Proposition: For a logistic regression model and under a 2nd-order Taylor approximation of $J_{\nu}(\theta, \lambda)$ wrt θ :

$$J_{\nu}(\boldsymbol{\theta}, \lambda) \approx 0.25 \lambda^2 \boldsymbol{\theta}^T \mathsf{Cov}_{\nu}(y\mathbf{x}) \boldsymbol{\theta}$$

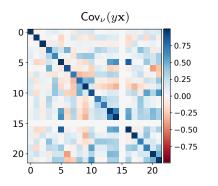
• $\pi_J(\theta)$ would be a multivariate normal distribution:

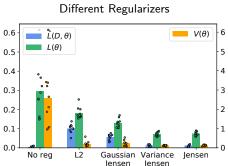
$$\pi_J(\boldsymbol{\theta}) \propto e^{-n0.25\lambda^2 \theta^T \text{Cov}_{\nu}(y\mathbf{x})\theta}$$

 If the data is normalized and features are conditionally independent, it is equal to L2-regularization ,

$$\boldsymbol{\theta}^T \mathsf{Cov}_{\nu}(y\mathbf{x})\boldsymbol{\theta} = \boldsymbol{\theta}^T kI\boldsymbol{\theta} = k||\boldsymbol{\theta}||^2$$

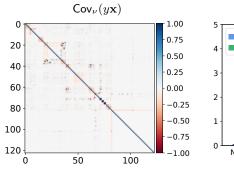
- Explains why L2-regularization improves generalization:
 - Small-norm models tends to have lower variance.
 - Lower variance implies better estimators $\hat{L}(D, \theta)$.
 - Better estimators leads to less overfitting.
- Also explains the **limitations** of L2-regularization:
 - L2-regularization does not take into account parameter correlations.

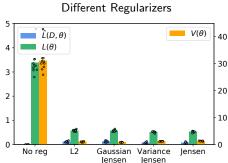




Mushroom Dataset:

- Attributes are highly conditionally (un)correlated.
- $Cov_{\nu}(y\mathbf{x})$ very different from a identity matrix.
- L2 performs poorly.





Adult Dataset:

- Attributes are not conditionally correlated.
- $Cov_{\nu}(y\mathbf{x})$ very similar to identity matrix.
- L2 performs well.

More connections with existing regularizations

- ullet For linear regression models, $\pi_J(m{ heta})$ is directly related to **g-priors** (Zellner, 1986).
- $\pi_J(\theta)$ is directly related to input gradient-normalization (Drucker et al., 1992, Varga et al., 2017).
- Working with more connections with other regularization techniques.

Conclusions and Future Works

- PAC-Bayesian bounds and the **generalization performance** of Bayesian methods.
 - Generalization is a key property in machine learning.
- PAC-Bayesian bounds allow to better understand **Bayesian priors**.
 - Open problem in Bayesian statistics.
 - We can explain the role of regularizing and informative priors.
 - Explain why (some) regularization methods work.
- PAC-Bayesian bounds allow to identify and correct weaknesses of Bayesian methods.
 - When learning under model misspecification, Bayesian posterior is not optimal (Masegosa, 2020).
 - We can get better performance for the same price.

- PAC-Bayesian bounds and the **generalization performance** of Bayesian methods.
 - Generalization is a key property in machine learning.
- PAC-Bayesian bounds allow to better understand **Bayesian priors**.
 - Open problem in Bayesian statistics.
 - We can explain the role of regularizing and informative priors.
 - Explain why (some) regularization methods work.
- PAC-Bayesian bounds allow to identify and correct weaknesses of Bayesian methods.
 - When learning under model misspecification, Bayesian posterior is not optimal (Masegosa, 2020).
 - We can get better performance for the same price.
- PAC-Bayesian bounds allow to better understand ensembles.
 Ortega et al. Diversity and Generalization in Neural Network Ensembles. AISTATS 2022.

- PAC-Bayesian bounds and the **generalization performance** of Bayesian methods.
 - Generalization is a key property in machine learning.
- PAC-Bayesian bounds allow to better understand **Bayesian priors**.
 - Open problem in Bayesian statistics.
 - We can explain the role of regularizing and informative priors.
 - Explain why (some) regularization methods work.
- PAC-Bayesian bounds allow to identify and correct weaknesses of Bayesian methods.
 - When learning under model misspecification, Bayesian posterior is not optimal (Masegosa, 2020).
 - We can get better performance for the same price.
- PAC-Bayesian bounds allow to better understand ensembles.
 Ortega et al. Diversity and Generalization in Neural Network Ensembles. AISTATS 2022.
- Future/Ongoing works:
 - Explain the Cold Posterior Effect (Wenzel et al., 2020).

- PAC-Bayesian bounds and the **generalization performance** of Bayesian methods.
 - Generalization is a key property in machine learning.
- PAC-Bayesian bounds allow to better understand **Bayesian priors**.
 - Open problem in Bayesian statistics.
 - We can explain the role of regularizing and informative priors.
 - Explain why (some) regularization methods work.
- PAC-Bayesian bounds allow to identify and correct weaknesses of Bayesian methods.
 - When learning under model misspecification, Bayesian posterior is not optimal (Masegosa, 2020).
 - We can get better performance for the same price.
- PAC-Bayesian bounds allow to better understand ensembles.
 Ortega et al. Diversity and Generalization in Neural Network Ensembles. AISTATS 2022.
- Future/Ongoing works:
 - Explain the Cold Posterior Effect (Wenzel et al., 2020).