

Abstract

- Most of machine learning models are misspecified.
- A novel PAC-Bayesian analysis shows that Bayesian model averaging is suboptimal for generalization under misspecification.
- A novel learning framework explicitly addressing misspecification is presented.

The learning problem

- Assumptions:**
 - $\nu(\mathbf{x})$ is the data generating distribution (**unknown**).
 - Model Misspecification:** $\forall \theta p(\cdot|\theta) \neq \nu$.
- The **predictive posterior distribution** for a given $\rho(\theta)$,

$$p(\mathbf{x}) = \int p(\mathbf{x}|\theta)\rho(\theta)d\theta = \mathbb{E}_{\rho}[p(\mathbf{x}|\theta)]$$
- The **learning problem** is defined as,

$$\rho^* = \arg \min_{\rho} \text{KL}(\nu(\mathbf{x}), \mathbb{E}_{\rho}[p(\mathbf{x}|\theta)]) = \arg \min_{\rho} \underbrace{\mathbb{E}_{\nu(\mathbf{x})}[-\ln \mathbb{E}_{\rho}[p(\mathbf{x}|\theta)]]}_{\text{CE}(\rho)}$$
- $\text{CE}(\rho)$ measures the **generalization error** associated to ρ .

First-order PAC-Bayes bounds and the Bayesian posterior

Germain et al. 2016 showed the Bayesian posterior minimize a (first-order) PAC-Bayesian bound:

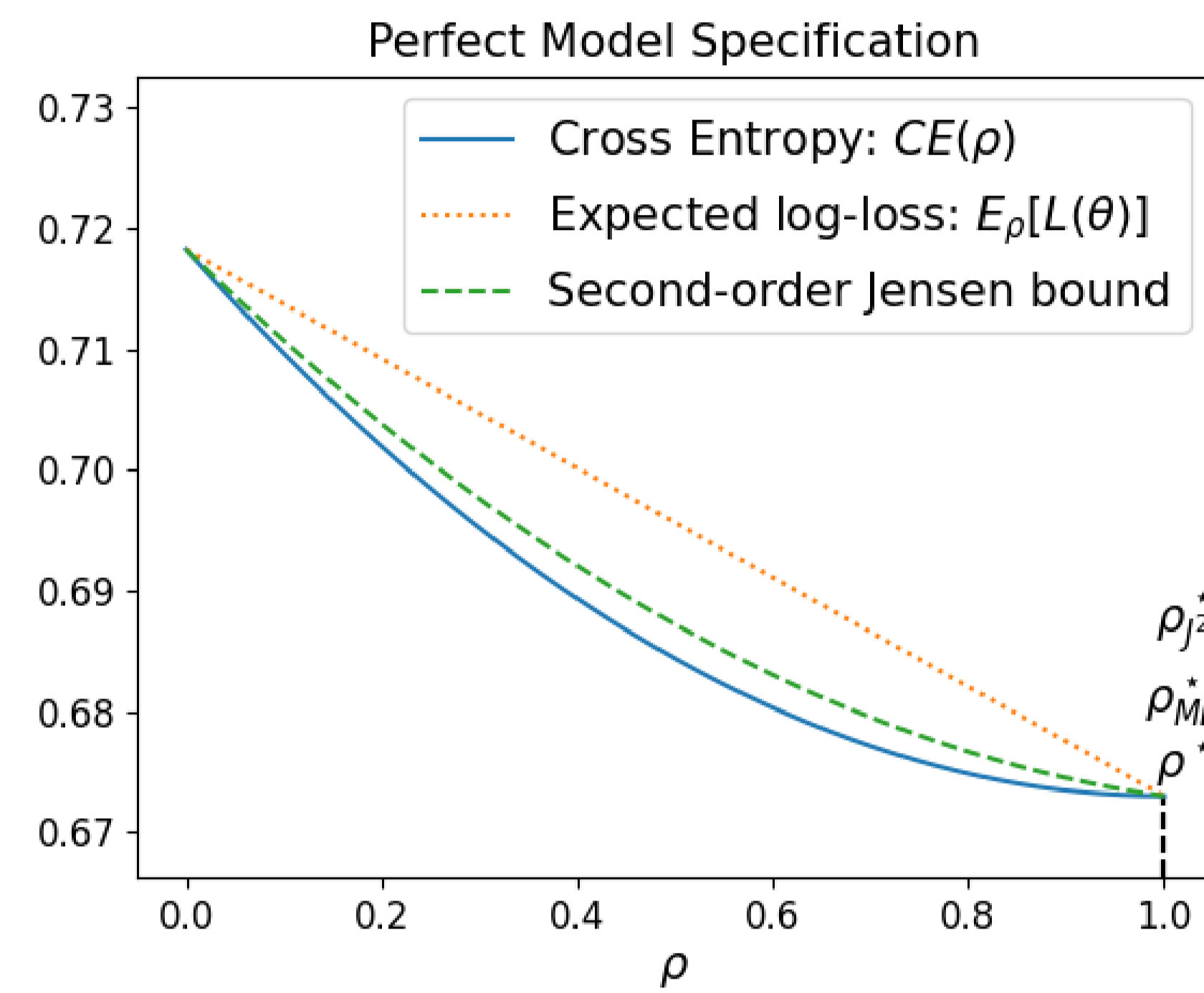
$$\underbrace{p(\theta|D)}_{\text{Bayesian Posterior}} = \arg \min_{\rho} \underbrace{\mathbb{E}_{\rho}[\hat{L}(\theta, D)]}_{\text{Empirical log-loss}} + \frac{\text{KL}(\rho, \pi)}{n} + \text{cte}$$

First-Order PAC-Bayes bound

First-order PAC-Bayes upper bounds the generalization error $\text{CE}(\rho)$

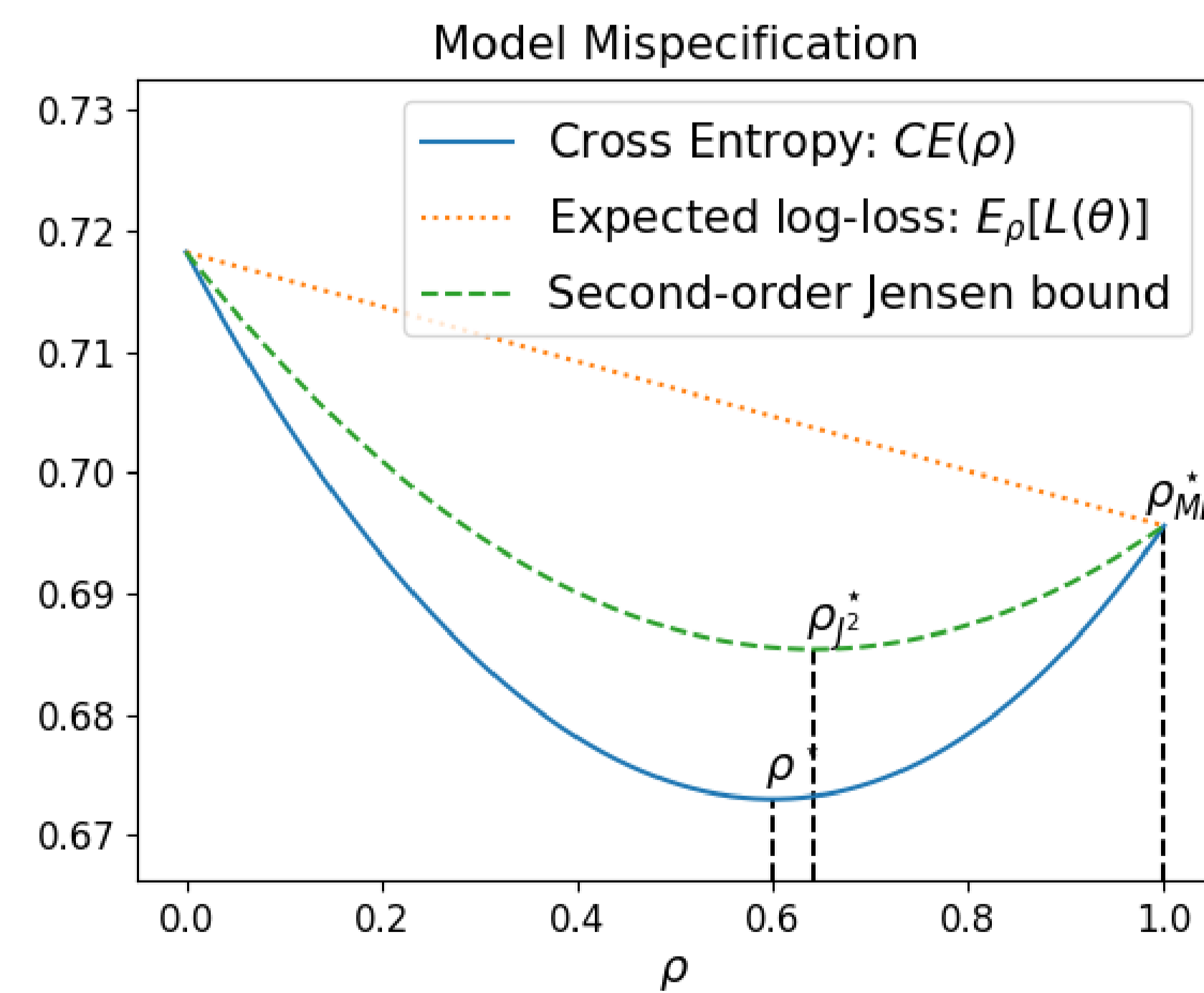
$$\underbrace{\text{CE}(\rho)}_{\text{Generalization Error}} \leq \underbrace{\mathbb{E}_{\rho}[L(\theta)]}_{\text{Jensen bound}} \leq \underbrace{\mathbb{E}_{\rho}[\hat{L}(\theta, D)]}_{\text{First-Order PAC-Bayes bound}} + \frac{\text{KL}(\rho, \pi)}{n} + \text{cte}$$

Bayesian posterior is optimal under perfect specification



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

Bayesian posterior is not optimal under model misspecification



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

A new learning framework

Minimizing second-order PAC-Bayes bounds

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

where Q is a tractable family of densities.

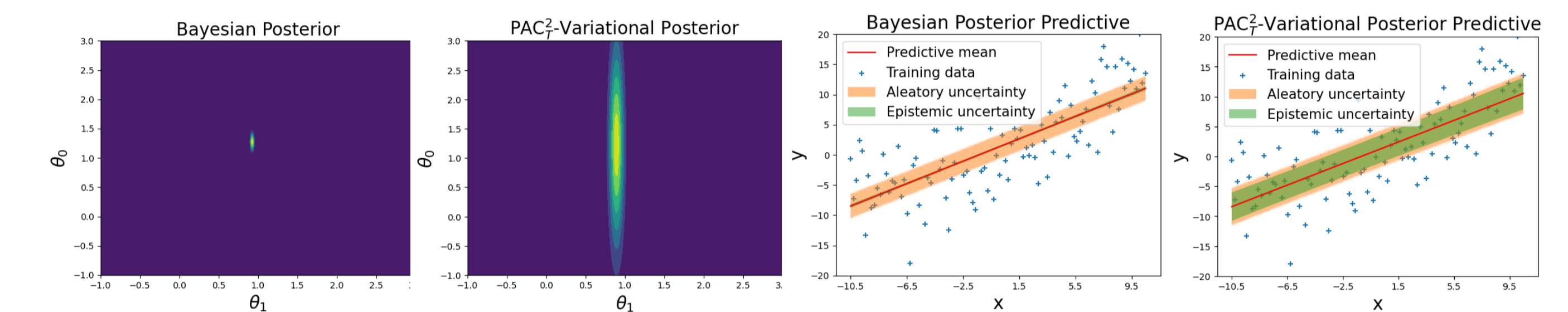


Figure 1: Bayesian Linear Regression.

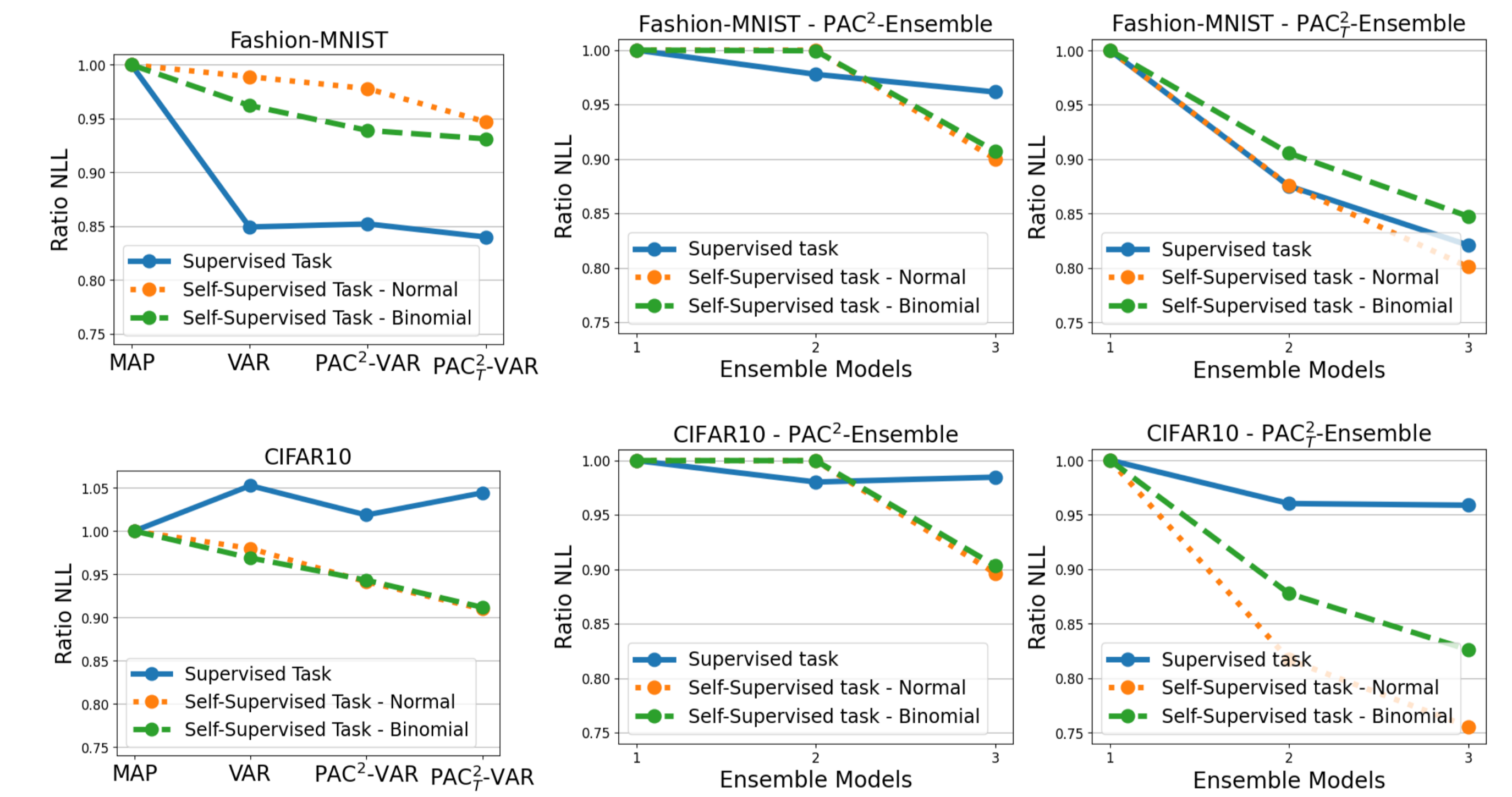


Figure 2: Bayesian Neural Networks.

Summary

- Bayesian methods are suboptimal for learning predictive models when the model family is misspecified.
- Second order PAC-Bayes bounds directly address model misspecification when learning.