Bayesian Model Averaging is not Model Combination: A PAC-Bayesian Analysis of Deep Ensembles

Andrés R. Masegosa

University of Aalborg (Copenhagen Campus) Denmark

 Masegosa, A. R. Learning under model misspecification: Applications to variational and ensemble methods. NeurIPS 2020.
Ortega, L. A., Cabañas R., Masegosa, A. R. Diversity and Generalization in Neural Network Ensembles. arxiv/2110.13786. 2021

#### Bayesian Model Averaging vs Model Combination

#### Bayesian model averaging is not model combination Thomas P. Minka December 13, 2002

In a recent paper, Domingos (2000) compares Bayesian model averaging (BMA) to other model combination methods on some benchmark data sets, is surprised that BMA performs worst, and suggests that BMA may be flawed. These results are actually *not* surprising, especially in light of an earlier paper by Domingos (1997) where it was shown that model combination works by enriching the space of hypotheses, not by approximating a Bayesian model average. And the only flaw with BMA is the belief that it is an algorithm for model combination, when it is not.

#### Model Combination

- Model Combination works by enriching the model space.
- BMA represents the inability to distinguish the best single model when using limited data.

- Deep Ensembles provides SOTA performance in terms of:
  - Uncertainty Estimation.
  - Robustness Against Distributional Shifts.

- Deep Ensembles provides SOTA performance in terms of:
  - Uncertainty Estimation.
  - Robustness Against Distributional Shifts.

#### Open Question: Why do deep ensembles work so well?



operand a source adulting that there is able in interference to adult any advection of a source adulting to a there is a source adult of the source of the s

Coop creation are a popular gozzani hal work by reheating the areas musclession characteristic and a weight get a muscling models. When does assentiative van indexistenzis is the analysis and gezzanism and the presented in a social deviasion competitor traves and gezanism. Begins in thereas procedures, per closely were and all weisland indexistences. Your Sequence the gezanism competitor traves "Beginsian models, to the Assezzanism control and the analysis and the analysis and the second analysis and the second analysis and the increasion. The Market and Control and analysis and analysis and the second question at the service (Devian Harvet) in practice, There is and control and advid deag amendiates", the analysis and gezanism deep terming workshop, at Neuroffs. 2016.





Finite Sample Approximation

- Deep Ensembles provides SOTA performance in terms of:
  - Uncertainty Estimation.
  - Robustness Against Distributional Shifts.



- Benefits of deep ensembles are due to their Bayesian nature.
- Deep Ensembles do not perform model combination (Minka, 2002).

- Let us assume we have:
  - Model class:  $\{p(y|\mathbf{x}, \boldsymbol{\theta}) : \boldsymbol{\theta} \in \boldsymbol{\Theta}\}$ . (e.g. neural network fixed architecture).
  - Training Data Sample D.
- Bayesian posterior:



- Let us assume we have:
  - Model class:  $\{p(y|\mathbf{x}, \theta) : \theta \in \Theta\}$ . (e.g. neural network fixed architecture).
  - Training Data Sample D.
- Bayesian posterior:



• Bayesian model averaging (BMA):

$$\underbrace{\mathbb{E}_{p(\theta|D)}[p(y \mid \mathbf{x}, \boldsymbol{\theta})]}_{\text{Bayesian Predictive posterior}} = \int \underbrace{p(y \mid \mathbf{x}, \boldsymbol{\theta}) d\boldsymbol{\theta}}_{\text{Bayesian posterior}} \underbrace{p(\boldsymbol{\theta}|D)}_{\text{Bayesian posterior}} d\boldsymbol{\theta}$$

- $\rho(\theta)$  denotes a quasi-posterior:
  - Defines a probability distribution over  $\Theta$ .
  - May depends on the training data sample D.

- $\rho(\theta)$  denotes a quasi-posterior:
  - Defines a probability distribution over  $\Theta$ .
  - May depends on the training data sample D.
- Model Combination defined by  $\rho(\theta)$ :

$$\underbrace{\mathbb{E}_{\rho}[p(y \mid \mathbf{x}, \boldsymbol{\theta})]}_{\text{Predictive posterior}} = \int \underbrace{p(y \mid \mathbf{x}, \boldsymbol{\theta})}_{\text{Quasi-posterior}} \underbrace{\rho(\boldsymbol{\theta})}_{\text{Quasi-posterior}} d\boldsymbol{\theta}$$

- $\rho(\theta)$  denotes a quasi-posterior:
  - Defines a probability distribution over  $\Theta$ .
  - May depends on the training data sample D.
- Model Combination defined by  $\rho(\theta)$ :

$$\underbrace{\mathbb{E}_{\rho}[p(y \mid \mathbf{x}, \boldsymbol{\theta})]}_{\text{Predictive posterior}} = \int \underbrace{p(y \mid \mathbf{x}, \boldsymbol{\theta})}_{\text{Quasi-posterior}} \underbrace{\rho(\boldsymbol{\theta})}_{\text{Quasi-posterior}} d\boldsymbol{\theta}$$

#### • Considerations:

- BMA is special case of model combination (i.e. when  $\rho(\theta) = p(\theta \mid D)$ ).
- We can choose any distribution  $\rho$  over  $\Theta$  using information from D.

- We assume there is an unknown data generating distribution  $\nu(y, \mathbf{x})$ .
  - The training/test data are generated from  $\nu$  (i.e  $D \sim \nu(y, \mathbf{x})$ ).

- We assume there is an unknown data generating distribution  $\nu(y, \mathbf{x})$ .
  - The training/test data are generated from  $\nu$  (i.e  $D \sim \nu(y, \mathbf{x})$ ).
- Generalization Error of a p-combined model using cross-entropy loss.

 $\underbrace{L_{ce}(\boldsymbol{\rho})}_{\boldsymbol{\mu}} = \mathbb{E}_{\boldsymbol{\nu}}[-\ln \underbrace{\mathbb{E}_{\boldsymbol{\rho}}[p(y \mid \mathbf{x}, \boldsymbol{\theta})]}_{\boldsymbol{\mu}}]$ Predictive posterior Gen. Error of

ρ-combined model

redictive poster of  $\rho(\theta)$ 

- We assume there is an unknown data generating distribution  $\nu(y, \mathbf{x})$ .
  - The training/test data are generated from  $\nu$  (i.e  $D \sim \nu(y, \mathbf{x})$ ).
- Generalization Error of a p-combined model using cross-entropy loss.



• Model combination works with an enriched model class parametrized by  $\rho$ :

 $\{p(y|\mathbf{x}, \boldsymbol{\theta}) : \boldsymbol{\theta} \in \boldsymbol{\Theta}\} \subseteq \{\mathbb{E}_{\rho}[p(y \mid \mathbf{x}, \boldsymbol{\theta})] : \rho \text{ is a quasi-posterior over } \boldsymbol{\Theta}\}$ 

Model class parametrized by  $\theta$ 

Model class parametrized by  $\rho$ 

- We assume there is an unknown data generating distribution  $\nu(y, \mathbf{x})$ .
  - The training/test data are generated from  $\nu$  (i.e  $D \sim \nu(y, \mathbf{x})$ ).
- Generalization Error of a *ρ*-combined model using cross-entropy loss.



Model combination works with an enriched model class parametrized by ρ:

 $\{\underline{p(y|\mathbf{x}, \boldsymbol{\theta}) : \boldsymbol{\theta} \in \boldsymbol{\Theta}}\} \subseteq \{\underline{\mathbb{E}_{\rho}[p(y \mid \mathbf{x}, \boldsymbol{\theta})] : \rho \text{ is a quasi-posterior over } \boldsymbol{\Theta}}\}$ 

Model class parametrized by  $\theta$ 

Model class parametrized by  $\rho$ 

In consequence



Jensen Inequality  $L_{ce}(\rho)$  $\mathbb{E}_{\boldsymbol{\rho}}[L_{ce}(\boldsymbol{\theta})]$ Gen. Error of Gibbs Error: p-average p-combined model individual models' error





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

PAC-Bayes bound (Germain et al. 2016)



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

The Bayesian posterior is a proxy

$$p(\boldsymbol{\theta}|D) \approx \arg\min_{\boldsymbol{\rho}} L_{ce}(\boldsymbol{\rho})$$

Bayesian Model Averaging is not Model Combination: A PAC-Bayesian Analysis

## Is the Bayesian approach an optimal distribution for model combination?

#### PAC-Bayesian upper-bounds (Alquier et al. 2016, Germain et al. 2016, Masegosa 2020)



**()** The Bayesian posterior converges (large sample limit) to  $\arg \min \mathbb{E}_{\rho}[L_{ce}(\theta)]$ .

# Is the Bayesian approach an optimal distribution for model combination?

#### PAC-Bayesian upper-bounds (Alquier et al. 2016, Germain et al. 2016, Masegosa 2020)



**()** The Bayesian posterior converges (large sample limit) to  $\arg \min \mathbb{E}_{\rho}[L_{ce}(\theta)]$ .

**(a)** The minimum of  $\mathbb{E}_{\rho}[L(\theta)]$  is a **Dirac-delta distribution** centered around the **best** possible single model  $\theta^*$ .

$$\underbrace{\delta(\boldsymbol{\theta} - \boldsymbol{\theta}^{\star})}_{\rho} = \arg\min_{\rho} \mathbb{E}_{\rho}[L_{ce}(\boldsymbol{\theta})]$$

Dirac-Delta distribution



• The Bayesian posterior converges (large sample limit) to  $\arg \min \mathbb{E}_{\rho}[L_{ce}(\theta)]$ .

Of The minimum of E<sub>ρ</sub>[L(θ)] is a Dirac-delta distribution centered around the best possible single model θ<sup>\*</sup>.

$$\underbrace{\delta(\boldsymbol{\theta} - \boldsymbol{\theta}^{\star})}_{\rho} = \arg\min_{\rho} \mathbb{E}_{\rho}[L_{ce}(\boldsymbol{\theta})]$$

Dirac-Delta distribution

Is this solution optimal for model combination? (Masegosa 2020)

$$\delta(\boldsymbol{\theta} - \boldsymbol{\theta}^{\star}) = \arg\min_{\rho} L_{ce}(\rho)$$

Bayesian Model Averaging is not Model Combination: A PAC-Bayesian Analysis



• The Bayesian posterior converges (large sample limit) to  $\arg \min \mathbb{E}_{\rho}[L_{ce}(\theta)]$ .

Of The minimum of E<sub>ρ</sub>[L(θ)] is a Dirac-delta distribution centered around the best possible single model θ<sup>\*</sup>.

$$\underbrace{\delta(\boldsymbol{\theta} - \boldsymbol{\theta}^{\star})}_{\rho} = \arg\min_{\rho} \mathbb{E}_{\rho}[L_{ce}(\boldsymbol{\theta})]$$

Dirac-Delta distribution

Is this solution optimal for model combination? (Masegosa 2020)

$$\delta(\boldsymbol{\theta} - \boldsymbol{\theta}^{\star}) = \arg\min_{\rho} L_{ce}(\rho)$$

Under perfect model specification (i.e.  $\nu(y|\mathbf{x})$  belongs to the model class).

# BMA optimal under perfect specification (Masegosa 2020)



# BMA optimal under perfect specification (Masegosa 2020)





# BMA optimal under perfect specification (Masegosa 2020)





# BMA optimal under model misspecification (Masegosa 2020)



# BMA optimal under model misspecification (Masegosa 2020)





# BMA optimal under model misspecification (Masegosa 2020)













Jensen Inequality <  $\mathbb{E}_{\rho}[L_{ce}(\boldsymbol{\theta})]$  $L_{ce}(\rho)$ Gen. Error of Gibbs Error: p-average

p-combined model

individual models' error





Bayesian Model Averaging is not Model Combination: A PAC-Bayesian Analysis

#### Second-Order PAC-Bayesian Bounds (Masegosa 2020)





Bayesian Model Averaging is not Model Combination: A PAC-Bayesian Analysis

Bayesian Model Averaging





 $L(\theta)$ <

 $L_{ce}(\rho)$ 

 $\underset{\theta^{\star} \text{model}}{\text{Gen. Error of}}$ 

Gen. Error of *p*-combined model

 $L(\boldsymbol{\theta}$ >

Gen. Error of  $\theta^*$  model

 $\underbrace{L_{ce}(\rho^{\star})}_{\text{Gen. Error of }}_{\rho^{\star}\text{-combined model}}$ 

Model Combination Bayesian Model Averaging  $\nu(y \mid x)$  $\mathbb{E}_{\delta(\theta-\theta^*)}[p(y|x,\theta)]$  $\nu(y \mid x)$  $\mathbb{E}_{\delta(\theta=\theta^*)}[p(y|x,\theta)]$ Θ  $\mathbb{E}_{e^{s}}[p(y|x, \theta$  $\nu(y|x)$  $L_{ce}(\rho^{\star})$ Gen. Error of  $\theta^*$  model Gen. Error of p\*-combined model Gen. Error of Gen. Error of  $\theta^*$  model p-combined model

# Which is the regime of Neural Networks Ensembles?

Two Regimes (Masegosa 2020, Ortega et al. 2021)

# If Deep Ensembles work in the BMA regime

![](_page_36_Figure_2.jpeg)

# If Deep Ensembles work in the Model Combination regime

![](_page_37_Figure_2.jpeg)

# If Deep Ensembles work in the Model Combination regime

![](_page_38_Figure_2.jpeg)

## Empirical Evidence (Ortega et al. 2021)

![](_page_39_Picture_1.jpeg)

CIFAR 10

![](_page_39_Picture_3.jpeg)

CIFAR 100

#### Experimental Settings

- Ensembles of four ResNet20/leNet5 Neural Networks.
- Each model was independently run until convergence (random initialization).
- Five repetitions.

### Empirical Evidence (Ortega et al. 2021)

![](_page_40_Figure_1.jpeg)

- Model Combination Regime above the black-line.
- Bayesian Model Averaging Regime below the black-line.

Two Regimes (Masegosa 2020, Ortega et al. 2021)

# If Deep Ensembles work in the Model Combination regime

Each Local Optima is a Different Model

![](_page_41_Figure_3.jpeg)

### Corollary 5 of (Ortega et al. 2021)

If the diversity of the ensemble is large enough:

![](_page_41_Figure_6.jpeg)

Two Regimes (Masegosa 2020, Ortega et al. 2021)

# If Deep Ensembles work in the Model Combination regime

Each Local Optima is a Different Model

![](_page_42_Figure_3.jpeg)

#### Corollary 5 of (Ortega et al. 2021)

If the diversity of the ensemble is large enough:

![](_page_42_Figure_6.jpeg)

then, the generalization of the ensemble is better:

![](_page_42_Figure_8.jpeg)

Gen. Error of model  $\theta^*$ 

Gen. Error of  $\rho$ -combined model

# Empirical Evidence (Ortega et al. 2021)

![](_page_43_Figure_1.jpeg)

- Model Combination Regime above the black-line.
- Bayesian Model Averaging Regime below the black-line.

- This analysis corroborates the notions of Thomas Minka:
  - BMA only tries to identify the best single model.
  - Model Combination exploit an enrich model space.

- This analysis corroborates the notions of Thomas Minka:
  - BMA only tries to identify the best single model.
  - Model Combination exploit an enrich model space.
- We have two regimes:
  - BMA regime is optimal when your model class is not wrong.
  - MC regime is optimal when your model class is wrong.

- This analysis corroborates the notions of Thomas Minka:
  - BMA only tries to identify the best single model.
  - Model Combination exploit an enrich model space.
- We have two regimes:
  - BMA regime is optimal when your model class is not wrong.
  - MC regime is optimal when your model class is wrong.
- Evidence supporting that Deep Ensembles works on the model combination regime.

- This analysis corroborates the notions of Thomas Minka:
  - BMA only tries to identify the best single model.
  - Model Combination exploit an enrich model space.
- We have two regimes:
  - BMA regime is optimal when your model class is not wrong.
  - MC regime is optimal when your model class is wrong.
- Evidence supporting that Deep Ensembles works on the model combination regime.

Ensembles are great when Your model is Wrong!!