Diversity and Generalization in Neural Network Ensembles

Andrés R. Masegosa

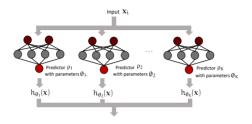
October 31, 2025

Aalborg University (Copenhagen Campus)

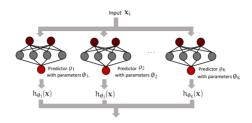
Ortega, L.A., Cabañas, R. and Masegosa, A. R., Diversity and Generalization in Neural Network Ensembles. AISTATS 2022.

Introduction and Motivation

Introduction: Ensembles of Neural Networks

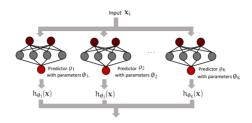


Introduction: Ensembles of Neural Networks



- Ensembles of NNs are recently getting a lot of attention.
 - Provide better uncertainty quantification.
 - More robust to **Out-Distribution-Data**.
 - Key properties in many real-world applications.

Introduction: Ensembles of Neural Networks



- Ensembles of NNs are recently getting a lot of attention.
 - Provide better uncertainty quantification.
 - More robust to Out-Distribution-Data.
 - Key properties in many real-world applications.
- Ongoing debate of why ensembles of NNs work so well:
 - Ensemble's diversity is widely used to justify ensemble performance.

- Ensemble's **diversity** is a broad concept:
 - Ensemble's performance jointly depends of the individual model's performance and the diversity among them.

- Ensemble's **diversity** is a broad concept:
 - Ensemble's performance jointly depends of the individual model's performance and the diversity among them.
 - An ensemble has null diversity if the predictions of individual models coincide on all the samples.

- Ensemble's **diversity** is a broad concept:
 - Ensemble's performance jointly depends of the individual model's performance and the diversity among them.
 - An ensemble has null diversity if the predictions of individual models coincide on all the samples.
 - No advantage of having an ensemble when diversity is null.

- Ensemble's **diversity** is a broad concept:
 - Ensemble's performance jointly depends of the individual model's performance and the diversity among them.
 - An ensemble has null diversity if the predictions of individual models coincide on all the samples.
 - No advantage of having an ensemble when diversity is null.
- Theoretically, diversity is **not** a **well-established concept**:
 - Different names: ambiguity, disagreement, etc.
 - Many different proposals to define diversity.
 - No theoretical analysis covering different different ensembles.

Our Contributions

- We built on previously published results:
 - (Krogh and Vedelsby, 1994): Ensemble of regression models.
 - (Masegosa, 2020): Bayesian model averaging.
 - (Masegosa et al., 2020): Weighted Majority Vote.

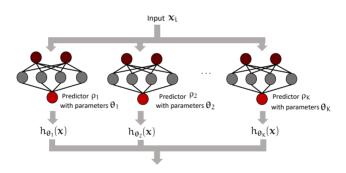
Our Contributions

- We built on previously published results:
 - (Krogh and Vedelsby, 1994): Ensemble of regression models.
 - (Masegosa, 2020): Bayesian model averaging.
 - (Masegosa et al., 2020): Weighted Majority Vote.
- We introduce a theoretical framework to answer these questions:
 - 1) How to measure the diversity of an ensemble?.
 - 2) How is diversity related to the ensemble's generalization performance?.
 - 3) How can diversity be promoted by ensemble learning algorithms?.
- We derive a common framework for different types of ensembles.

Previous Knowledge

Basics on NNs ensembles

An ensemble trained with $D = \{(x_1, y_1), \dots, (x_n, y_n)\}$ is the combination of different predictors.



Regression Ensemble: Multiple regression models.

• Weighted Model Averaging: $MA_{\rho}(\mathbf{x}) = \mathbb{E}_{\theta \sim \rho}[h_{\theta}(\mathbf{x})].$

Regression Ensemble: Multiple regression models.

- Weighted Model Averaging: $MA_{\rho}(\mathbf{x}) = \mathbb{E}_{\boldsymbol{\theta} \sim \rho}[h_{\boldsymbol{\theta}}(\mathbf{x})].$
- Squared loss:

$$L_{\mathsf{sq}}(\theta) = \mathbb{E}_{\nu}[(h_{\theta}(\mathbf{x}) - y)^2] \quad L_{\mathsf{sq}}(\rho) = \mathbb{E}_{\nu}[(MA_{\rho}(\mathbf{x}) - y)^2]$$

Regression Ensemble: Multiple regression models.

- Weighted Model Averaging: $MA_{\rho}(\mathbf{x}) = \mathbb{E}_{\boldsymbol{\theta} \sim \rho}[h_{\boldsymbol{\theta}}(\mathbf{x})].$
- Squared loss:

$$L_{sq}(\theta) = \mathbb{E}_{\nu}[(h_{\theta}(\mathbf{x}) - y)^2] \quad L_{sq}(\rho) = \mathbb{E}_{\nu}[(MA_{\rho}(\mathbf{x}) - y)^2]$$

Probabilistic Ensemble: Multiple probabilistic classification models.

Weighted Model Averaging

Regression Ensemble: Multiple regression models.

- Weighted Model Averaging: $MA_{\rho}(\mathbf{x}) = \mathbb{E}_{\boldsymbol{\theta} \sim \rho}[h_{\boldsymbol{\theta}}(\mathbf{x})].$
- Squared loss:

$$L_{\mathsf{sq}}(\boldsymbol{\theta}) = \mathbb{E}_{\nu}[(h_{\boldsymbol{\theta}}(\mathbf{x}) - y)^2] \quad L_{\mathsf{sq}}(\rho) = \mathbb{E}_{\nu}[(MA_{\rho}(\mathbf{x}) - y)^2]$$

Probabilistic Ensemble: Multiple probabilistic classification models.

- Weighted Model Averaging
- Cross-entropy loss

Regression Ensemble: Multiple regression models.

- Weighted Model Averaging: $MA_{\rho}(\mathbf{x}) = \mathbb{E}_{\boldsymbol{\theta} \sim \rho}[h_{\boldsymbol{\theta}}(\mathbf{x})].$
- Squared loss:

$$L_{\mathsf{sq}}(\boldsymbol{\theta}) = \mathbb{E}_{\nu}[(h_{\boldsymbol{\theta}}(\mathbf{x}) - y)^2] \quad L_{\mathsf{sq}}(\rho) = \mathbb{E}_{\nu}[(MA_{\rho}(\mathbf{x}) - y)^2]$$

Probabilistic Ensemble: Multiple probabilistic classification models.

- Weighted Model Averaging
- Cross-entropy loss

Majority Vote Ensemble: Multiple classification models.

Weighted Majority Vote

Regression Ensemble: Multiple regression models.

- Weighted Model Averaging: $MA_{\rho}(\mathbf{x}) = \mathbb{E}_{\boldsymbol{\theta} \sim \rho}[h_{\boldsymbol{\theta}}(\mathbf{x})].$
- Squared loss:

$$L_{\mathsf{sq}}(\boldsymbol{\theta}) = \mathbb{E}_{\nu}[(h_{\boldsymbol{\theta}}(\mathbf{x}) - y)^2] \quad L_{\mathsf{sq}}(\rho) = \mathbb{E}_{\nu}[(MA_{\rho}(\mathbf{x}) - y)^2]$$

Probabilistic Ensemble: Multiple probabilistic classification models.

- Weighted Model Averaging
- Cross-entropy loss

Majority Vote Ensemble: Multiple classification models.

- Weighted Majority Vote
- Zero-one loss

Diversity and Ensembles'

Performance

Diversity and Generalization

Theorem 1

General Upper-bound for all the ensembles considered in this work:

$$\underbrace{L(\rho)}_{\text{Ensemble's}} \leq \alpha \left(\underbrace{\mathbb{E}_{\rho}[L(\theta)]}_{\text{Enjected Loss}} - \underbrace{\mathbb{D}(\rho)}_{\text{Diversity}} \right)$$

where $\alpha = 4$ for the 0/1-loss, otherwise, $\alpha = 1$.

The diversity term depends on the considered **loss function**:

$$\mathbb{D}(
ho) = \mathbb{E}_{
u} \Big[\mathbb{V}_{
ho} \left(f(y, \boldsymbol{x}; \boldsymbol{\theta}) \right) \Big].$$

Diversity and Generalization

Theorem 1

General Upper-bound for all the ensembles considered in this work:

$$\underbrace{L(\rho)}_{\text{Ensemble's}} \leq \alpha \left(\underbrace{\mathbb{E}_{\rho}[L(\theta)]}_{\text{Enjected Loss}} - \underbrace{\mathbb{D}(\rho)}_{\text{Diversity}} \right)$$

where $\alpha = 4$ for the 0/1-loss, otherwise, $\alpha = 1$.

The diversity term depends on the considered **loss function**:

$$\mathbb{D}(
ho) = \mathbb{E}_{
u} \Big[\mathbb{V}_{
ho} \left(f(y, \boldsymbol{x}; \boldsymbol{\theta}) \right) \Big].$$

 Ensemble's performance depends on both individual models' performance and diversity.

Regression Ensemble: Multiple regression models.

$$\mathbb{D}_{\mathsf{sq}}(
ho) \ = \ \mathbb{E}_{
u} \Big[\mathbb{V}_{
ho}(h_{oldsymbol{ heta}}(oldsymbol{x})) \Big]$$

Regression Ensemble: Multiple regression models.

$$\mathbb{D}_{\sf sq}(
ho) \;\; = \;\; \mathbb{E}_{
u} \Big[\mathbb{V}_{
ho}(h_{oldsymbol{ heta}}(oldsymbol{x})) \Big]$$

Probabilistic Ensemble: Multiple probabilistic classification models.

$$\mathbb{D}_{\mathsf{ce}}(\rho) = \mathbb{E}_{\nu} \left[\mathbb{V}_{\rho} \left(\frac{p(y \mid \boldsymbol{x}, \boldsymbol{\theta})}{\sqrt{2} \max_{\boldsymbol{\theta}} p(y \mid \boldsymbol{x}, \boldsymbol{\theta})} \right) \right]$$

Regression Ensemble: Multiple regression models.

$$\mathbb{D}_{\sf sq}(
ho) \;\; = \;\; \mathbb{E}_{
u} \Big[\mathbb{V}_{
ho}(h_{oldsymbol{ heta}}(oldsymbol{x})) \Big]$$

Probabilistic Ensemble: Multiple probabilistic classification models.

$$\mathbb{D}_{\mathsf{ce}}(\rho) = \mathbb{E}_{\nu} \left[\mathbb{V}_{\rho} \left(\frac{p(y \mid \boldsymbol{x}, \boldsymbol{\theta})}{\sqrt{2} \max_{\boldsymbol{\theta}} p(y \mid \boldsymbol{x}, \boldsymbol{\theta})} \right) \right]$$

Majority Vote Ensemble: Multiple classification models.

$$\mathbb{D}_{0/1}(\rho) = \mathbb{E}_{\nu} \Big[\mathbb{V}_{\rho} \Big(\mathbb{1} \big(h_{\theta}(\mathbf{x}) \neq \mathbf{y} \big) \Big) \Big]$$

Is $\mathbb{D}(\rho)$ a diversity measure?

A General Measure of Diversity:

$$\mathbb{D}(\rho) = \mathbb{E}_{\nu}\Big[\mathbb{V}_{\rho}\left(f(y, \boldsymbol{x}; \boldsymbol{\theta})\right)\Big].$$

Lemma

- i) $\mathbb{D}(\rho) \geq 0$
- ii) If all individual models provide the same predictions, then $\mathbb{D}(
 ho)=0$.
- iii) $0 \leq \mathbb{D}(\rho) \leq \mathbb{E}_{\rho}[L(\boldsymbol{\theta})].$
- iv) $\mathbb{D}(\rho)$ is invariant to reparametrizations.

A General Measure of Diversity:

$$\mathbb{D}(
ho) = \mathbb{E}_{
u} \Big[\mathbb{V}_{
ho} \left(f(y, oldsymbol{x}; oldsymbol{ heta})
ight) \Big]$$

Theorem

The diversity term $\mathbb{D}(\rho)$ can be written as

$$\mathbb{D}(\rho) = \mathbb{V}_{\nu \times \rho} \Big(f(y, \boldsymbol{x}; \boldsymbol{\theta}) \Big) - \mathbb{E}_{\rho \times \rho} \Big[\mathsf{Cov}_{\nu} (f(y, \boldsymbol{x}; \boldsymbol{\theta}), f(y, \boldsymbol{x}; \boldsymbol{\theta}')) \Big]$$

where $Cov_{\nu}(\cdot,\cdot)$ is the co-variance between two models.

A General Measure of Diversity:

$$\mathbb{D}(
ho) = \mathbb{E}_{
u} \Big[\mathbb{V}_{
ho} \left(f(y, oldsymbol{x}; oldsymbol{ heta})
ight) \Big]$$

Theorem

The diversity term $\mathbb{D}(\rho)$ can be written as

$$\mathbb{D}(\rho) = \mathbb{V}_{\nu \times \rho}\Big(f(y, \boldsymbol{x}; \boldsymbol{\theta})\Big) - \mathbb{E}_{\rho \times \rho}\Big[\mathsf{Cov}_{\nu}(f(y, \boldsymbol{x}; \boldsymbol{\theta}), f(y, \boldsymbol{x}; \boldsymbol{\theta}'))\Big]$$

where $Cov_{\nu}(\cdot,\cdot)$ is the co-variance between two models.

 First term helps to explain why randomized models improve ensemble performance.

A General Measure of Diversity:

$$\mathbb{D}(
ho) = \mathbb{E}_{
u} \Big[\mathbb{V}_{
ho} \left(f(y, oldsymbol{x}; oldsymbol{ heta})
ight) \Big]$$

Theorem

The diversity term $\mathbb{D}(\rho)$ can be written as

$$\mathbb{D}(\rho) = \mathbb{V}_{\nu \times \rho}\Big(f(y, \boldsymbol{x}; \boldsymbol{\theta})\Big) - \mathbb{E}_{\rho \times \rho}\Big[\mathsf{Cov}_{\nu}(f(y, \boldsymbol{x}; \boldsymbol{\theta}), f(y, \boldsymbol{x}; \boldsymbol{\theta}'))\Big]$$

where $Cov_{\nu}(\cdot,\cdot)$ is the co-variance between two models.

- First term helps to explain why randomized models improve ensemble performance.
- Second term helps to explain why independent and anti-correlated models improve ensemble performance.

Diversity and Generalization

Theorem 1

General Upper-bound for all the ensembles considered in this work:

$$\underbrace{L(\rho)}_{\text{Ensemble's}} \leq \alpha \left(\underbrace{\mathbb{E}_{\rho}[L(\theta)]}_{\text{Individual Models'}} - \underbrace{\mathbb{D}(\rho)}_{\text{Diversity}}\right)$$

where $\alpha = 4$ for the 0/1-loss, otherwise, $\alpha = 1$.

 Ensemble's performance depends on both individual models' performance and diversity among them.

A General Measure of Diversity:

$$\mathbb{D}(
ho) = \mathbb{E}_{
u} \Big[\mathbb{V}_{
ho} \left(f(y, \boldsymbol{x}; \boldsymbol{\theta}) \right) \Big].$$

Question: How much do we gain by ensembling a set of models wrt randomly choosing them?

A General Measure of Diversity:

$$\mathbb{D}(
ho) = \mathbb{E}_{
u} \Big[\mathbb{V}_{
ho} \left(f(y, \boldsymbol{x}; \boldsymbol{\theta}) \right) \Big].$$

Question: How much do we gain by ensembling a set of models wrt randomly choosing them?

Corollary

Under these settings, we have that

$$\mathbb{D}(\rho) \leq \underbrace{\mathbb{E}_{\rho}[L(\theta)] - \frac{1}{\alpha}L(\rho)}_{\textit{Ensemble's Gap}}$$

Answer: Larger diversity induces larger gains.

A General Measure of Diversity:

$$\mathbb{D}(
ho) = \mathbb{E}_{
u} \Big[\mathbb{V}_{
ho} \left(f(y, \boldsymbol{x}; \boldsymbol{\theta}) \right) \Big].$$

Question: When is an ensemble better that best individual model?

A General Measure of Diversity:

$$\mathbb{D}(
ho) = \mathbb{E}_{
u} \Big[\mathbb{V}_{
ho} \left(f(y, \boldsymbol{x}; \boldsymbol{\theta}) \right) \Big].$$

Question: When is an ensemble better that best individual model?

Corollary

$$\underbrace{\mathbb{E}_{\rho}[L(\theta)] - \frac{1}{\alpha}L(\theta^{\star})}_{Single\ Model's\ Error\ Gap} < \underbrace{\mathbb{D}(\rho)}_{Ensemble's\ Expected\ Loss} \underbrace{L(\rho)}_{Ensemble's\ Expected\ Loss} < \underbrace{L(\theta^{\star})}_{Single\ Model's\ Expected\ Loss}$$

Answer: If the diversity of the ensemble is large enough, then an ensemble outperforms the best single model.

How to Exploit Diversity to Learn Ensembles?

Diversity and Generalization

Theorem 1

General Upper-bound for all the ensembles considered in this work:

$$\underbrace{L(\rho)}_{\text{Ensemble's}} \leq \alpha \left(\underbrace{\mathbb{E}_{\rho}[L(\theta)]}_{\text{Individual Models'}} - \underbrace{\mathbb{D}(\rho)}_{\text{Diversity}}\right)$$

where $\alpha = 4$ for the 0/1-loss, otherwise, $\alpha = 1$.

• This inequality depends on the data generating distribution.

A PAC-Bayesian Bound

For distribution $\pi(\theta)$ independent of D, with probability at least $1-\delta$ over draws of training data $D \sim \nu^n(y, \mathbf{x})$ (i.e., i.i.d.), for all $\lambda > 0$, for all distribution ρ over Θ , simultaneously,

$$\underbrace{L(\rho)}_{\text{Ensemble's}} \leq \alpha \left(\underbrace{\mathbb{E}_{\rho}[\hat{L}(\boldsymbol{\theta}, D)]}_{\text{Averaged}} - \underbrace{\hat{\mathbb{D}}(\rho, D)}_{\text{Ensemble's}} + \underbrace{\frac{2\mathit{KL}(\rho \mid \pi)}{\lambda n}}_{\text{Regularization}} + \underbrace{\frac{\epsilon(\nu, \pi, \lambda, n, \delta)}{\lambda n}}_{\text{Regularization}}\right)$$

A PAC-Bayesian Bound

For distribution $\pi(\theta)$ independent of D, with probability at least $1-\delta$ over draws of training data $D \sim \nu^n(y, \mathbf{x})$ (i.e., i.i.d.), for all $\lambda > 0$, for all distribution ρ over Θ , simultaneously,

$$\underbrace{L(\rho)}_{\text{Ensemble's}} \leq \alpha \left(\underbrace{\mathbb{E}_{\rho}[\hat{L}(\boldsymbol{\theta}, D)]}_{\text{Averaged}} - \underbrace{\hat{\mathbb{D}}(\rho, D)}_{\text{Ensemble's}} + \underbrace{\frac{2\mathit{KL}(\rho \mid \pi)}{\lambda n}}_{\text{Regularization}} + \underbrace{\frac{\epsilon(\nu, \pi, \lambda, n, \delta)}{\lambda n}}_{\text{An}}\right)$$

 \bullet Find the ρ minimizing this PAC-Bayesian Bound.

A PAC-Bayesian Bound

For distribution $\pi(\theta)$ independent of D, with probability at least $1 - \delta$ over draws of training data $D \sim \nu^n(y, \mathbf{x})$ (i.e., i.i.d.), for all $\lambda > 0$, for all distribution ρ over Θ , simultaneously,

$$\underbrace{\mathcal{L}(\rho)}_{\text{Ensemble's}} \leq \alpha \left(\underbrace{\mathbb{E}_{\rho}[\hat{L}(\boldsymbol{\theta}, D)]}_{\text{Averaged}} - \underbrace{\hat{\mathbb{D}}(\rho, D)}_{\text{Ensemble's}} + \underbrace{\frac{2\mathit{KL}(\rho \mid \pi)}{\lambda n}}_{\text{Regularization}} + \underbrace{\frac{\epsilon(\nu, \pi, \lambda, n, \delta)}{\lambda n}}_{\text{Regularization}}\right)$$

- Find the ρ minimizing this PAC-Bayesian Bound.
- We move to a continuous hypothesis space.

Ensemble Learning algorithm as a mixture model

$$\rho(\boldsymbol{\theta}|\boldsymbol{\theta}_1,\ldots,\boldsymbol{\theta}_K,\sigma^2) = \frac{1}{K} \sum_{k=1}^K \mathcal{N}(\boldsymbol{\theta}; \; \boldsymbol{\theta}_k,\sigma^2 \boldsymbol{I}).$$

Ensemble Learning algorithm as a mixture model

$$\rho(\boldsymbol{\theta}|\boldsymbol{\theta}_1,\ldots,\boldsymbol{\theta}_K,\sigma^2) = \frac{1}{K} \sum_{k=1}^K \mathcal{N}(\boldsymbol{\theta}; \; \boldsymbol{\theta}_k,\sigma^2 \boldsymbol{I}).$$

Learning Objective (P2B-Ensemble)

$$\min_{\theta_1, \dots, \theta_K} \underbrace{\mathbb{E}_{\rho}[\hat{L}(\boldsymbol{\theta}, D)]}_{\text{Averaged}} - \underbrace{\hat{\mathbb{D}}(\rho, D)}_{\text{Ensemble's}} - \underbrace{\frac{2\mathbb{E}_{\rho}[\ln \pi(\boldsymbol{\theta})]}{\lambda n}}_{\text{Regularization}}$$

Empirical validation: Experimental Settings

Tasks

- Regression Task: Wine-Quality dataset.
- Classification Task: Cifar10 and Cifar100 data sets.

Empirical validation: Experimental Settings

Tasks

- Regression Task: Wine-Quality dataset.
- Classification Task: Cifar10 and Cifar100 data sets.

Models

- **Regression Task**: MLP with 50 hidden units.
- Classification Task: LeNet5 and ResNet20 convolutional networks.

Empirical validation: Experimental Settings

Tasks

- Regression Task: Wine-Quality dataset.
- Classification Task: Cifar10 and Cifar100 data sets.

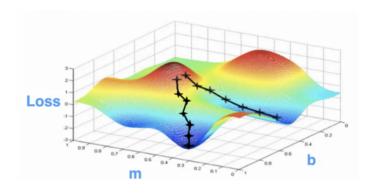
Models

- **Regression Task**: MLP with 50 hidden units.
- Classification Task: LeNet5 and ResNet20 convolutional networks.

Learning Algorithms

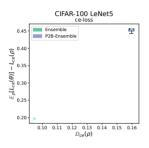
- P2B-Ensemble: K models jointly learned promoting diversity.
- **Ensemble**: *K* models independently learned.

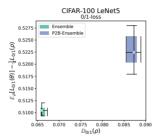
Ensemble Learning

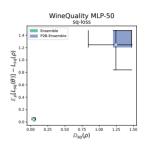


(Manish Kumar, 2018)

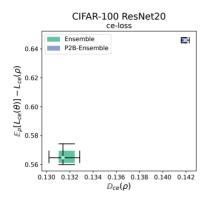
• Ensemble composed by K different local minima.

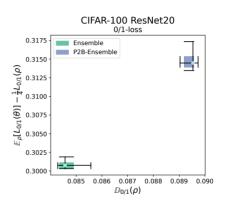




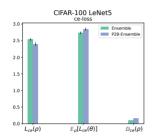


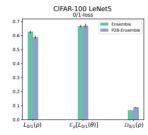
- Higher diversity correlates with higher gains by ensembling.
- Standard ensemble methods implicitly promote diversity.
- P2B-Ensemble finds ensembles with higher diversity.

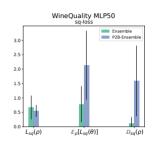




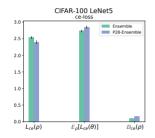
- Higher diversity correlates with higher gains by ensembling.
- Standard ensemble methods implicitly promote diversity.
- P2B-Ensemble finds ensembles with higher diversity.

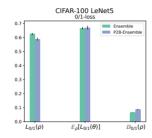


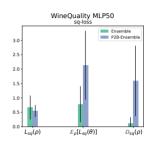




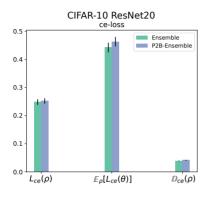
 Explicitly promoting diversity (ie. P2B-Ensemble) gives rise to better ensembles.

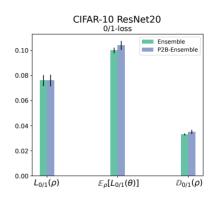




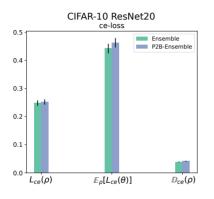


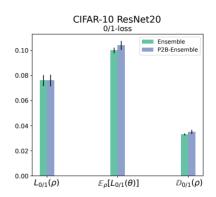
- Explicitly promoting diversity (ie. P2B-Ensemble) gives rise to better ensembles.
- But not always...





• P2B-Ensemble is not able to learn better ensembles.





- P2B-Ensemble is not able to learn better ensembles.
- Why?
 - Because big neural networks works in the interpolation regime.

Learning Objective

$$\min_{\theta_1, \dots, \theta_K} \underbrace{\mathbb{E}_{\rho}[\hat{L}(\boldsymbol{\theta}, D)]}_{\text{Averaged}} - \underbrace{\hat{\mathbb{D}}(\rho, D)}_{\text{Ensemble's}} - \underbrace{\frac{2\mathbb{E}_{\rho}[\ln \pi(\boldsymbol{\theta})]}{\lambda n}}_{\text{Regularization}}$$

Learning Objective

$$\min_{\theta_1, \dots, \theta_K} \underbrace{\mathbb{E}_{\rho}[\hat{L}(\theta, D)]}_{\text{Averaged}} - \underbrace{\hat{\mathbb{D}}(\rho, D)}_{\text{Ensemble's}} - \underbrace{\frac{2\mathbb{E}_{\rho}[\ln \pi(\theta)]}{\lambda n}}_{\text{Regularization}}$$

Inequality

$$0 \leq \hat{\mathbb{D}}(\rho, D) \leq \mathbb{E}_{\rho}[\hat{\mathcal{L}}(\boldsymbol{\theta}, D)]$$

Learning Objective

$$\min_{\theta_1, \dots, \theta_K} \underbrace{\mathbb{E}_{\rho}[\hat{L}(\theta, D)]}_{\substack{\text{Averaged} \\ \text{Empirical Loss}}} - \underbrace{\hat{\mathbb{D}}(\rho, D)}_{\substack{\text{Ensemble's} \\ \text{Empirical Diversity}}} \underbrace{\frac{2\mathbb{E}_{\rho}[\ln \pi(\theta)]}{\lambda n}}_{\substack{\text{Regularization}}}$$

Inequality

$$0 \leq \hat{\mathbb{D}}(\rho, D) \leq \mathbb{E}_{\rho}[\hat{\mathcal{L}}(\boldsymbol{\theta}, D)]$$

In the interpolation regime

$$\mathbb{E}_{
ho}[\hat{L}(oldsymbol{ heta},D)]pprox 0\Rightarrow\hat{\mathbb{D}}(
ho,D)pprox 0$$

• The empirical diversity does not provide any signal to the gradient.

Conclusions and Future Work

Conclusions

- We can formally speak about ensemble's diversity.
- Applies to very different ensemble methods.
- Useful to understand and derive learning algorithms.

Conclusions and Future Work

Conclusions

- We can formally speak about ensemble's diversity.
- Applies to very different ensemble methods.
- Useful to understand and derive learning algorithms.

Limitations

- Diversity's linear dependency: not accurate in all cases (Germain et al. 2015, Wu et al. 2021)
- Only second-order interactions.
- Learning in the interpolation-regime.

Conclusions and Future Work

Conclusions

- We can formally speak about ensemble's diversity.
- Applies to very different ensemble methods.
- Useful to understand and derive learning algorithms.

Limitations

- Diversity's linear dependency: not accurate in all cases (Germain et al. 2015, Wu et al. 2021)
- Only second-order interactions.
- Learning in the interpolation-regime.

Future Works

Promote diversity using a external (non-labelled) dataset.

Questions?

Andrés, L.A.O., Cabañas, R. and Masegosa, A. R.,
Diversity and Generalization in Neural Network
Ensembles. AISTATS 2022.

Measuring Diversity

• For regression ensembles, (Krogh and Vedelsby, 1994) showed that:

$$\underbrace{L_{\mathsf{sq}}(\rho)}_{\text{Ensemble's}} = \mathbb{E}_{\rho} \underbrace{\left[\mathbb{E}_{\nu} [(y - h_{\theta}(\mathbf{x}))^{2}]\right]}_{\text{Individual Models'}} - \mathbb{E}_{\nu} \underbrace{\left[\mathbb{E}_{\rho} [(h_{\theta}(\mathbf{x}) - \mathbb{E}_{\rho} [h_{\theta}(\mathbf{x})])^{2}]\right]}_{\text{Variance among individual models}}$$

Measuring Diversity

• For regression ensembles, (Krogh and Vedelsby, 1994) showed that:

$$\underbrace{L_{\mathsf{sq}}(\rho)}_{\mathsf{Ensemble's}} = \mathbb{E}_{\rho} \underbrace{ \mathbb{E}_{\nu} [(y - h_{\theta}(\mathbf{x}))^2]]}_{\mathsf{Individual Models'}} - \mathbb{E}_{\nu} \underbrace{ \mathbb{E}_{\rho} [(h_{\theta}(\mathbf{x}) - \mathbb{E}_{\rho}[h_{\theta}(\mathbf{x})])^2]]}_{\mathsf{Variance among individual models}}$$

• **Strong ensembles** require strong and diverse individual models: small individual error and high variance.

Measuring Diversity

• For regression ensembles, (Krogh and Vedelsby, 1994) showed that:

$$\underbrace{L_{\mathsf{sq}}(\rho)}_{\text{Ensemble's}} = \mathbb{E}_{\rho} \underbrace{\left[\mathbb{E}_{\nu} [(y - h_{\theta}(\mathbf{x}))^2]\right]}_{\text{Individual Models'}} - \mathbb{E}_{\nu} \underbrace{\left[\mathbb{E}_{\rho} [(h_{\theta}(\mathbf{x}) - \mathbb{E}_{\rho} [h_{\theta}(\mathbf{x})])^2]\right]}_{\text{Variance among individual models}}$$

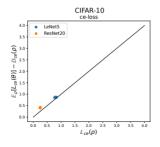
• **Strong ensembles** require strong and diverse individual models: small individual error and high variance.

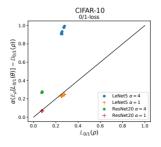
 Existing literature only contains ad-hoc decompositions for other kind of ensembles, but there is not a general decomposition.

Theorem 1

$$\underbrace{L(\rho)}_{\text{Ensemble's}} \leq \alpha \left(\underbrace{\mathbb{E}_{\rho}[L(\theta)]}_{\text{Individual Models'}} - \underbrace{\mathbb{D}(\rho)}_{\text{Diversity}}\right)$$

where $\alpha = 4$ for the 0/1-loss, otherwise, $\alpha = 1$.





Theorem 1

$$\underbrace{L(\rho)}_{\text{Ensemble's}} \leq \alpha \left(\underbrace{\mathbb{E}_{\rho}[L(\theta)]}_{\text{Individual Models'}} - \underbrace{\mathbb{D}(\rho)}_{\text{Diversity}}\right)$$

where $\alpha = 4$ for the 0/1-loss, otherwise, $\alpha = 1$.

