Second Order PAC-Bayesian Bounds for the Weighted Majority Vote

Andrés R. Masegosa¹ Stephan S. Lorenzen² Christian Igel² Yevgeny Seldin²

¹University of Almería

²University of Copenhagen

NeurIPS, December 2020

Weighted Majority Vote

- Fundamental technique for combining predictions of multiple classifiers
- ▶ Used in Bagging, Boosting, etc.
- Wins most ML competitions

Weighted Majority Vote

- Fundamental technique for combining predictions of multiple classifiers
- Used in Bagging, Boosting, etc.
- Wins most ML competitions

Ensemble's Key Power

Cancellation of errors effect
 If the errors are independent, they average out

Weighted Majority Vote

- Fundamental technique for combining predictions of multiple classifiers
- Used in Bagging, Boosting, etc.
- Wins most ML competitions

Ensemble's Key Power

Cancellation of errors effect
 If the errors are independent, they average out

Our contributions

- Second order PAC-Bayesian generalization bound for the weighted majority vote
- Minimization of the bound guides weighting of ensemble members and does not deteriorate the test error

If $\rho\text{-weighted}$ majority vote makes an error, then at least a $\rho\text{-weighted}$ half of the classifiers make an error

If ρ -weighted majority vote makes an error, then at least a ρ -weighted half of the classifiers make an error

$$\underbrace{\mathcal{L}(\mathsf{MV}_{\rho})}_{\substack{\mathsf{Expected loss of}\\ \rho\text{-weighted majority vote}}} \leq \mathbb{P}(\underbrace{\mathbb{E}_{\rho}[\mathbbm{1}(h(X) \neq Y)]}_{\rho\text{-weighted mass}} \geq 0.5)$$

If ρ -weighted majority vote makes an error, then at least a ρ -weighted half of the classifiers make an error



First order bound: use Markov's inequality $(\mathbb{P}(X \ge \varepsilon) \le \frac{1}{\varepsilon}\mathbb{E}[X])$

$$\leq$$
 2 $\mathbb{E}_D[\mathbb{E}_
ho[\mathbb{1}(h(X)
eq Y)]]$

$$=$$
 2 $\mathbb{E}_{\rho}[L(h)]$

Expected loss of ρ -weighted randomized classifier

If ρ -weighted majority vote makes an error, then at least a ρ -weighted half of the classifiers make an error



First order bound: use Markov's inequality $(\mathbb{P}(X \ge \varepsilon) \le \frac{1}{\varepsilon}\mathbb{E}[X])$

$$\leq 2 \mathbb{E}_{D}[\mathbb{E}_{\rho}[\mathbb{1}(h(X) \neq Y)]]$$

= 2 $\mathbb{E}_{\rho}[L(h)]$
Expected loss of ρ -weighted

Expected loss of ho-weighted randomized classifier

Issues

- Ignores correlation of errors (the key power)
- Minimization of the corresponding PAC-Bayes bound degrades the test error (Lorenzen et al., 2019)

Prior second order analysis

The C-bounds (Lacasse et al., 2007, Germain et al., 2015, Laviolette et al., 2017)

Based on Chebyshev-Cantelli inequality

$$\mathbb{P}(X \ge \varepsilon) \le \frac{\mathbb{E}\left[X^2\right] - \mathbb{E}\left[X\right]^2}{\mathbb{E}\left[X^2\right] - \mathbb{E}\left[X\right] + \varepsilon^2}$$

Prior second order analysis

The C-bounds (Lacasse et al., 2007, Germain et al., 2015, Laviolette et al., 2017)

Based on Chebyshev-Cantelli inequality

$$\mathbb{P}(X \ge \varepsilon) \le \frac{\mathbb{E}\left[X^2\right] - \mathbb{E}\left[X\right]^2}{\mathbb{E}\left[X^2\right] - \mathbb{E}\left[X\right] + \varepsilon^2}$$

Issues

- $\mathbb{E}[X^2]$ and $\mathbb{E}[X]$ in the denominator make empirical estimation hard
- Empirically weaker than the first order bound (Lorenzen et al., 2019)
- Impossible to optimize the weighting except in very restrictive cases

$$\mathsf{V}_{
ho}) \leq \mathbb{P}(\ \mathbb{E}_{
ho}[\mathbbm{1}(h(X)
eq Y)] \geq 0.5$$

 $\substack{\rho\text{-weighted mass}\\ \text{of errors}}$

)

Expected loss of ρ -weighted majority vote

L(M)

$$\underbrace{\mathcal{L}(\mathsf{MV}_{\rho})}_{\substack{\mathsf{Expected loss of } \\ \rho \text{-weighted majority vote}}} \leq \mathbb{P}(\underbrace{\mathbb{E}_{\rho}[\mathbbm{1}(h(X) \neq Y)]}_{\rho \text{-weighted mass}} \geq 0.5)$$

Second-order Markov's inequality $\mathbb{P}(X \ge \varepsilon) \le \frac{1}{\varepsilon^2} \mathbb{E}[X^2]$:

$$\underbrace{L(\mathsf{MV}_{\rho})}_{\text{Expected loss of everythet majority vote}} \leq \mathbb{P}(\underbrace{\mathbb{E}_{\rho}[\mathbbm{1}(h(X) \neq Y)]}_{\rho \text{-weighted mass of errors}} \geq 0.5)$$

 ρ

Second-order Markov's inequality $\mathbb{P}(X \ge \varepsilon) \le \frac{1}{\varepsilon^2} \mathbb{E}[X^2]$: $\le 4 \mathbb{E}_D[\mathbb{E}_{\rho}[\mathbb{1}(h(X) \neq Y)]^2]$

$$\underbrace{\mathcal{L}(\mathsf{MV}_{\rho})}_{\substack{\mathsf{Expected loss of } \\ \rho \text{-weighted majority vote}}} \leq \mathbb{P}(\underbrace{\mathbb{E}_{\rho}[\mathbbm{1}(h(X) \neq Y)]}_{\rho \text{-weighted mass } } \geq 0.5)$$

Second-order Markov's inequality $\mathbb{P}(X \ge \varepsilon) \le \frac{1}{\varepsilon^2} \mathbb{E}[X^2]$:

$$\leq 4 \mathbb{E}_D[\mathbb{E}_\rho[\mathbb{1}(h(X) \neq Y)]^2] \\ = 4 \mathbb{E}_{\rho^2}[\underbrace{\mathbb{E}_D[\mathbb{1}(h(X) \neq Y \land h'(X) \neq Y)]]}_{=}$$

Expected Tandem Loss: L(h,h')

$$\underbrace{\mathcal{L}(\mathsf{MV}_{\rho})}_{\substack{\mathsf{Expected loss of } \\ \rho \text{-weighted majority vote}}} \leq \mathbb{P}(\underbrace{\mathbb{E}_{\rho}[\mathbbm{1}(h(X) \neq Y)]}_{\rho \text{-weighted mass } } \geq 0.5)$$

Second-order Markov's inequality $\mathbb{P}(X \ge \varepsilon) \le \frac{1}{\varepsilon^2} \mathbb{E}[X^2]$:

$$\leq 4 \mathbb{E}_{D}[\mathbb{E}_{\rho}[\mathbb{1}(h(X) \neq Y)]^{2}]$$

$$= 4 \mathbb{E}_{\rho^{2}}[\underbrace{\mathbb{E}_{D}[\mathbb{1}(h(X) \neq Y \land h'(X) \neq Y)]}_{\text{Expected Tandem Loss: } L(h,h')}]$$

$$= 4 \mathbb{E}_{\rho^{2}}[(h,h')]$$

 $= 4 \mathbb{E}_{\rho^2}[L(h, h')]$

$$\underbrace{L(\mathsf{MV}_{\rho})}_{\text{Expected loss of everythet majority vote}} \leq \mathbb{P}(\underbrace{\mathbb{E}_{\rho}[\mathbbm{1}(h(X) \neq Y)]}_{\rho \text{-weighted mass of errors}} \geq 0.5)$$

 ρ

Second-order Markov's inequality $\mathbb{P}(X \ge \varepsilon) \le \frac{1}{\varepsilon^2} \mathbb{E}[X^2]$:

$$\leq 4 \mathbb{E}_{D}[\mathbb{E}_{\rho}[\mathbb{1}(h(X) \neq Y)]^{2}]$$

$$= 4 \mathbb{E}_{\rho^{2}}[\mathbb{E}_{D}[\mathbb{1}(h(X) \neq Y \land h'(X) \neq Y)]]$$
Expected Tandem Loss: $L(h,h')$

$$= 4 \mathbb{E}_{\rho^{2}}[L(h,h')]$$

Tandem loss counts an error if both *h* and *h'* err on a sample

$$\underbrace{L(\mathsf{MV}_{\rho})}_{\text{Expected loss of everythet majority vote}} \leq \mathbb{P}(\underbrace{\mathbb{E}_{\rho}[\mathbbm{1}(h(X) \neq Y)]}_{\rho \text{-weighted mass of errors}} \geq 0.5)$$

 ρ

Second-order Markov's inequality $\mathbb{P}(X \ge \varepsilon) \le \frac{1}{\varepsilon^2} \mathbb{E}[X^2]$:

$$\leq 4 \mathbb{E}_{D}[\mathbb{E}_{\rho}[\mathbb{1}(h(X) \neq Y)]^{2}]$$

$$= 4 \mathbb{E}_{\rho^{2}}[\underbrace{\mathbb{E}_{D}[\mathbb{1}(h(X) \neq Y \land h'(X) \neq Y)]}_{\text{Expected Tandem Loss: } L(h,h')}]$$

$$= 4 \mathbb{E}_{\rho^{2}}[L(h, h')]$$

- **Tandem loss** counts an error if both *h* and *h'* err on a sample
- Second order oracle bound: $L(MV_{\rho}) \leq 4\mathbb{E}_{\rho^2}[L(h, h')]$

A specialized oracle bound for binary classification In binary classification **tandem loss** L(h, h') satisfies

 $\mathbb{E}_{
ho^2}[L(h,h')] =$

Expected loss of p-weighted randomized classifier

 $\mathbb{E}_{\rho}[L(h)] \quad - \quad \frac{1}{2} \mathbb{E}_{\rho^2}[\mathbb{E}_D[\mathbb{1}(h(X) \neq h'(X))]].$

Expected Disagreement $\mathbb{D}(h,h')$

A specialized oracle bound for binary classification In binary classification tandem loss L(h, h') satisfies

$$\mathbb{E}_{
ho^2}[L(h,h')] =$$

$$\underbrace{\mathbb{E}_{\rho}[L(h)]}$$

Expected loss of ρ -weighted randomized classifier

<

$$\frac{1}{2} \mathbb{E}_{\rho^2} [\underbrace{\mathbb{E}_D[\mathbb{1}(h(X) \neq h'(X))]}_{\rho^2}].$$

Expected Disagreement $\mathbb{D}(h,h')$

Specialized oracle bound for binary classification

$$\underbrace{L(MV_{\rho})}$$

Expected loss of ρ -weighted majority vote

$$4 \underbrace{\mathbb{E}_{\rho}[L(h)]}_{\rho}$$

Expected loss of *p*-weighted randomized classifier

 $2 \mathbb{E}_{\rho^2}[\mathbb{D}(h, h')]$

Expected Disagreement of ρ -weighted rand. classifier

A specialized oracle bound for binary classification In binary classification tandem loss L(h, h') satisfies

$$\mathbb{E}_{\rho^2}[L(h,h')] =$$

$$\underbrace{\mathbb{E}_{\rho}[L(h)]}$$

Expected loss of ρ -weighted randomized classifier

<

$$\frac{1}{2} \mathbb{E}_{\rho^2} [\underbrace{\mathbb{E}_D[\mathbb{1}(h(X) \neq h'(X))]}_{\rho^2}].$$

Expected Disagreement $\mathbb{D}(h,h')$

Specialized oracle bound for binary classification

$$\underbrace{L(\mathsf{MV}_{\rho})}$$

Expected loss of ρ -weighted majority vote

$$4 \underbrace{\mathbb{E}_{\rho}[L(h)]}_{\rho}$$

Expected loss of *p*-weighted randomized classifier

 $2 \mathbb{E}_{\rho^2}[\mathbb{D}(h, h')]$

Expected Disagreement of ρ -weighted rand. classifier

$\mathbb{D}(h, h')$ only depends on **unlabeled data**!!

From oracle to empirical bounds

PAC-Bayes- λ (Thiemann et al., 2017):

For π independent of S, with probability at least $1-\delta$ for all ρ and $\lambda\in(0,2)$

From oracle to empirical bounds

PAC-Bayes- λ (Thiemann et al., 2017): For π independent of S, with probability at least $1 - \delta$ for all ρ and $\lambda \in (0, 2)$



From oracle to empirical bounds

PAC-Bayes- λ (Thiemann et al., 2017):

For π independent of S, with probability at least $1-\delta$ for all ρ and $\lambda \in (0,2)$ and $\gamma > 0$



Second-order PAC-Bayesian bound

$$L(\mathsf{MV}_{\rho}) \leq 4 \quad \underbrace{\mathbb{E}_{\rho^2}[L(h, h')]}^{\mathsf{Expected}}$$

Second-order PAC-Bayesian bound



Second-order PAC-Bayesian bound



Advantages

- Takes correlation of errors into account
- Easy to minimize and tight
- Minimization of the bound does not degrade the test error

Expected loss of majority vote

$$\widetilde{L(MV_{\rho})} \leq$$

 \leq 4 $\mathbb{E}_{\rho}[L(h)]$

Expected Disagreement

$$(h)] - 2 \mathbb{E}_{\rho^2}[\mathbb{D}(h, h')]$$



PAC-Bayes **upper** bound on $\mathbb{E}_{\rho^2}[L(h, h')]$





It can exploit unlabeled data

Empirical evaluation

 Test error of optimized majority vote over uniformly weighted baseline for first order [FO] and new second order [TND] bound (the lower the better)



Empirical evaluation

• The optimized weights ρ^* generated by the first order [FO] and the new second order [TND] bound.



State-of-the-art

Minimization of existing first-order bound significantly deteriorates the test error

State-of-the-art

- Minimization of existing first-order bound significantly deteriorates the test error
- Existing second-order bounds are looser and can not be optimized

State-of-the-art

- Minimization of existing first-order bound significantly deteriorates the test error
- Existing second-order bounds are looser and can not be optimized

Contributions

- Novel second order oracle bound for the weighted majority vote based on second order Markov's inequality
- Novel second order PAC-Bayesian bound for the weighted majority vote

State-of-the-art

- Minimization of existing first-order bound significantly deteriorates the test error
- Existing second-order bounds are looser and can not be optimized

Contributions

- Novel second order oracle bound for the weighted majority vote based on second order Markov's inequality
- Novel second order PAC-Bayesian bound for the weighted majority vote
- Minimization of the bound guides weighting of ensemble members and does not deteriorate the test error