

On The Robustness of Decision Tree Learning Under Label Noise

<https://arxiv.org/abs/1605.06296> and Hastie 2009

Goal: theoretical investigation of decision trees + label noise

to show for large samples, split rule under noisy data same for noise free

$S = (X \times Y)^n \sim$ ideal noise free iid

SL where $\tilde{y}_i = y_i$ w/ prob $(1 - \eta x_i)$

$\tilde{y}_i = -y_i$ w/ prob ηx_i

↑ noise ↑ true
symmetric/uniform if $\eta x = \eta$ $\forall x$
- η unknown to learning algo

Split rule & impurity

$$C(f) = G(v) - (a G(v_L) + (1-a) G(v_R))$$

say $p = \text{pos}$, $q = (1-p) = \text{frac negative}$

then Gini = $2pq$

$$C_{\text{entropy}} = -p \log p - q \log q$$

Noise Tolerance:

$$\arg \min_{f \in F} C(f) = \arg \min_{f \in F} C^\eta(f)$$

Under symmetric label noise,

$$p^\eta = p(1-\eta) + \eta n = p(1-2\eta) + \eta$$

Theorem 3: splitting criterion based on gini impurity is noise tolerant if $\eta < 0.5$

Lemma 7: if a leaf node v has n samples, under symmetric label noise w/ $\eta < 0.5$, majority voting will not fail w/ probability at least $1-\delta$ when $n \geq \frac{2}{p^2(1-2\eta)^2} \cdot \ln(\frac{1}{\delta})$ where p is the difference between the fraction of pos + negative samples in the noise free case.

Sample Complexity of RF

For a classifier:

$$\text{error}_{\text{gen}} = \text{error}_{\text{bias}} + \text{error}_{\text{var}} + \sigma_{\text{noise}}^2$$

For a classifier:

$$\text{error}_{\text{gen}} = \text{error}_{\text{bias}} + \text{error}_{\text{var}} + \sigma_{\text{noise}}^2$$

Under symmetric label noise, $\text{error}_{\text{bias}}$ is same for a single decision tree & random forest.
Claim: $\text{error}_{\text{gen}}$ controlled by $\text{error}_{\text{var}}$.

If ρ = pairwise correlation and $\text{Var } \sigma^2$ for each tree, then random forest w/ N trees has variance:

$$\text{error}_{\text{var}} = \rho \sigma^2 + \frac{1-\rho}{N} \sigma^2$$

Could one say that the different weighting of kernel polytopes decreases the correlation between trees and decreases the variance?

Hastie 2009

Notation:

trees $b=1$ to B

Classification: $\hat{C}_{\text{rf}}^B(x) = \text{maj vote } \{ \hat{C}_b(x) \}_{b=1}^B$

Regression: $\hat{f}_{\text{rf}}^B(x) = \frac{1}{B} \sum_{b=1}^B T_b(x)$

Notes

- expectation of B iid bagged trees is the same for any one
- bias of bagged trees is the same for any. so improvement must be through variance reduction

An avg of B iid vars all w/ variance σ^2 has variance $\frac{1}{B} \sigma^2$. If vars are identically distributed and have pairwise correlation ρ , then variance of the avg is:

$$\rho \sigma^2 + \frac{1-\rho}{B} \sigma^2$$

as $B \rightarrow \infty$, second term disappears

\therefore reducing the correlation btw trees w/o increasing the variance too much