

STAT406- Methods of Statistical Learning Lecture 17

Matias Salibian-Barrera

UBC - Sep / Dec 2017

Failing hard and failing often



Using Packages You Don't Understand

Hey, it runs!

O RLY?

Joe Hilgard

@JoeHilgard

Random forests

(1) `for (b in 1:B)`

- (a) Draw a bootstrap sample from the training data
- (b) Grow a “random forest tree” as follows: for each terminal node:
 - (i) Randomly select m features
 - (ii) Pick the best split among these
 - (iii) Split the node into two children
- (c) Repeat (b) to grow a (very very) large tree

(2) Return the ensemble of trees $(T_b)_{1 \leq b \leq B}$

Random forests

- Given a new point \mathbf{x} , for regression we use

$$\hat{f}(\mathbf{x}) = \frac{1}{B} \sum_{b=1}^B T_b(\mathbf{x})$$

- For classification:

$$\hat{f}(\mathbf{x}) = \text{majority vote among } \left\{ T_b(\mathbf{x}), \right. \\ \left. 1 \leq b \leq B \right\}$$

Q: why not average conditional prob's?

Out-of-bag error estimates

- Each bagged tree is trained on a bootstrap sample
- Predict the observations not in the bootstrap sample with that tree
- One will have “about” $B/3$ predictions for each point in the training set
- These can be used to estimate the prediction error (classification error rate) without having to use CV

Out-of-bag error estimates

- For each training observation (y_i, \mathbf{x}_i) , obtain a prediction using only those trees in which (y_i, \mathbf{x}_i) was **NOT** used
- In other words, let \mathcal{I}_i the set of trees (bootstrap samples) where (y_i, \mathbf{x}_i) does not appear, then

$$\hat{y}_i = \frac{1}{|\mathcal{I}_i|} \sum_{j \in \mathcal{I}_i} T_j(\mathbf{x}_i)$$

Random forests

- This error estimate can be computed at the same time as the trees are being built
- When this error estimate is stabilized we can stop adding trees to the ensemble

Example

OOB example

Random forests

- Feature ranking - relative importance of each variable
- Given a single tree T , at each node t split we can compute the sum of reductions in sum of squares (or gini or deviance measures) m_t^2
- We assign this squared measure m_t^2 to the variable (feature) used in the split

Random forests

- To each feature, we assign the sum of “squared gains” attributed to it
- For the i -th variable X_i we have

$$\mathcal{J}_i^2(T) = \begin{cases} m_t^2 & \text{if split involved } X_i \\ 0 & \text{otherwise} \end{cases}$$

Random forests

- For a random forest we use

$$\mathcal{J}_i^2 = \sum_T \mathcal{J}_i^2(T)$$

- In other words, we sum (or average) the importance of the variable across the trees in the forest

Random Forest - plot

Zip codes

