

## ENSEMBLE METHOD

= approach that combines many simple "building block" models in order to obtain a single model.

weak learners  
(since they usually lead to mediocre results on their own)

### 1. BAGGING / BOOTSTRAP AGGREGATION

procedure for reducing the variance of a statistical learning method (frequently used in the context of decision trees).

- Recall: averaging a set of observations reduces variance (if each obs  $Z_1, \dots, Z_n$  has variance  $\sigma^2$ , the variance of the mean  $\bar{Z}$  is given by  $\sigma^2/n$ )

→ BOOTSTRAPPING uses ↑ as intuition and takes repeated samples (with replacement) from the training data set. It then trains a model on the  $b$ th bootstrapped training set in order to get  $\hat{f}^*{}^b(x)$ . Finally it averages all the predictions to get

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

$$\text{Classification: } \hat{f}_{\text{bag}}(x) = \underset{1 \leq k \leq K}{\operatorname{argmax}} \sum_{b=1}^B \text{Ind}(\hat{f}_b(x) = k)$$

record class predicted by each of the  $B$  trees and take a majority vote: the overall prediction is the most commonly occurring class among the  $B$  predictions.

Note: - Bagging is only useful for high variance/low bias models.

- Using a very large value of  $B$  will not lead to overfitting ( $B = 100$  is usually sufficiently large that the error has settled down).
- On average, each bagged tree makes use of around two thirds of the observations

→ Why? if we draw items (with replacement) from an original set of size  $n$ , the probability of choosing any one item on the first draw is

$$P(x_i) = \frac{1}{n} \text{ and so } P(\overline{x}) = 1 - \frac{1}{n}$$

As  $n$  gets larger

$$\lim_{n \rightarrow \infty} \left(1 - \frac{1}{n}\right)^n = \frac{1}{e}$$

$$\approx 0.368$$

The probability of never choosing this item on any of the draws is

$$P = \left(1 - \frac{1}{n}\right)^n$$

- the remaining  $\frac{1}{3}$  of the observations not used to fit a given bagged tree are referred to as the out-of-bag (OOB) observations.

It is possible to predict the response of the  $i^{\text{th}}$  observation using each of the trees in which that observation was OOB

- ↳ around  $\frac{2}{3}$  predictions for the  $i^{\text{th}}$  obs
- ↳ if we average these predictions / take a majority vote we get a single OOB prediction for the  $i^{\text{th}}$  observation.
- if we do this for all  $m$  observations we can use these predictions to compute the MSE or classification error.  
(similar to leave-one-out cross validation)

## 2. RANDOM FORESTS

- ↳ Similar to bagging but forces each split to consider only a subset of the predictors

$m$  = predictor subset size

On average  $\frac{(P-m)}{P}$  of the splits will not even consider the strong predictor  
(the trees are "decorrelated")

Note: - when  $m=p$ , random forest amounts simply to bagging

## 3. BOOSTING

designed for binary classification  
 $y \in \{-1, 1\}$

Each tree is grown using information from previously grown trees. (the trees are grown sequentially).

Boosting does not involve bootstrap sampling. Instead, each tree is fit on a modified version of the original dataset.

Each tree attempts to capture a signal that is not yet accounted for by the current set of trees.

Note:-

- Boosting can overfit if  $B$  is too large!
- The shrinkage parameter controls the rate at which boosting learns.

how it works: (ADABOOST)

GIVEN: labeled training data with labels  $y_j \in \{-1, 1\}$

and a set / class  $H$  of  $T$  possible weak classifiers.

**Initialize:** We introduce a weight  $w_j^{(1)}$  for each training sample and set

$$w_j^{(1)} = \frac{1}{m} \text{ for each } j$$

For  $t = 1, 2, \dots, T$

- 1) We train a weak classifier  $h_t \in H$  using the data and weights associated with them. We then choose the classifier with the lowest error

$$\epsilon_t = \sum_{j=1}^m w_j^{(t)} \text{Ind}(y_j \neq h_t(x_j))$$

- 2) Update the weights (they increase for samples that are misclassified and decrease for samples that are correctly classified)

$$w_j^{(t+1)} = w_j^{(t)} \exp(-\alpha_t y_j h_t(x_j))$$

- 3) Normalize weights so that they sum to 1.

$$\alpha_t = \log \left( \frac{1 - \epsilon_t}{\epsilon_t} \right)$$

( $\alpha$  always > 0)

note the sign!

sign of output  $h_t(x_j)$  is the same as table  $\rightarrow +$   
otherwise  $\rightarrow -$

correct answer  $\rightarrow y_j h_t(x_j) > 0$

$$\rightarrow \exp(-\alpha_t y_j h_t(x_j)) < 0$$

$\rightarrow$  weight less important

#### 4. DROPOUT

Regularization method for Artificial Neural Networks.

$\rightarrow$  Can be seen as an ensemble method.