

Decision trees tend to overfit

i.e., New training data usually lead to very different trees.

Solution:

Introduce randomised algorithms.

1. Bagging = Bootstrap aggregation

$$S = \{ \vec{x}^{(i)}, y^{(i)} \}_{i=1, \dots, N} \quad \text{Our sample}$$

(1) Bootstrapping: Produce B samples from S all of size N by random sampling with replacement: $S_b, b=1, \dots, B$

(2) Models: From each of the B samples obtain a DT:

$$\left\{ \hat{f}_b^{DT} \right\}_{b=1}^B$$

In each S_b there will be repeated samples and absent samples from the original S

(3) Aggregate the model:

$$(i) \hat{f}(\vec{x}^{in}) = \frac{1}{B} \sum_{b=1}^B \hat{f}_b^{DT}(\vec{x}^{in})$$

Regression

(ii) Classification:

$$\vec{x}^{in} \mapsto [\vec{\pi}_1^{DT}, \dots, \vec{\pi}_B^{DT}]$$

$$\vec{\pi}_b = \begin{pmatrix} \pi_{b,1} \\ \vdots \\ \pi_{b,q} \end{pmatrix}$$

q classes

$$\vec{\pi} = \frac{1}{B} \sum_{b=1}^B \vec{\pi}_b^{DT}$$

(aggregated probability)

$$\hat{y} = \arg \max_g \vec{\pi}$$

These outcomes ~~are~~ are still correlated usually because of dominance of some descriptors.

If the $\{\hat{f}_b^{DT}(\vec{x}^{in})\}$ for a given \vec{x}^{in} are all correlated, then averaging them does not reduce the variance beyond a limit.

Hence we need to uncorrelate ^{them} more radically.

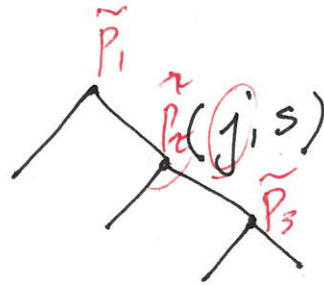
Random forests add an additional level of randomisation.

Steps:

(1) Generate a bootstrap sample

S_b

(2)



In DT

$j \in [1, \dots, P]$

maximise over all predictors

In RF:

At each split choose x_j from a random subset of descriptors: with \tilde{p} descriptors

$x_j \in \{x_i\}_{i \in RF}$ random subset of descriptors
 $|RF| = \tilde{p} < P$

Chosen at random among the $[1, \dots, P]$

\hat{f}^{RF}
 $[\tilde{P}_1, \tilde{P}_2, \tilde{P}_3, \dots]_S$

(3) Repeat S times
 $S = 1, \dots, S'$

The result is an ensemble of trees:

$$\left\{ \hat{f}_{[\tilde{p}_k]}^{RF} \right\}_{s=1}^{S'} = \underline{\text{random forest}}$$

k is defined on each of the splits.

The predictor is the aggregation of

$$\left\{ \hat{f}_{[\tilde{p}_k]}^{RF} \right\}_{s=1}^{S'}$$

- Regression: average of $\left\{ \hat{f}_{[\tilde{p}_k]}^{RF} \right\}_{s=1}^{S'}$
- Classification: $\left\{ \hat{T}_{[\tilde{p}_k]}^{RF} \right\}_{s=1}^{S'}$

Some comments:

- (1) Loss of interpretability:
with aggregation we lose the direct connection with the descriptors.

* Parameters (hyper-parameters):

• S = # of trees. Not very sensitive

• Minimum leaf size: $\begin{cases} \text{Classification} \rightarrow 1 \\ \text{Regression} \rightarrow 5 \end{cases}$

• $\tilde{p} < p$
of descriptors allowed at each split

$\begin{cases} \tilde{p} \sim \sqrt{p} & \text{prediction (regression)} \\ \tilde{p} \sim \frac{p}{3} & \text{classification.} \end{cases}$

$\begin{cases} \text{If } \tilde{p} = 1, \text{ trees are very uncorrelated.} \\ \text{If } \tilde{p} \approx p, \text{ trees are correlated} \end{cases}$

* Bootstrap buys us some extra knowledge:

For every sample S_b , there are some 'out-of-bag' samples:

out-of-bag err \equiv OOB error can be used for validation

• Example of ensemble methods

~~RF~~ is an example of an ensemble method.

Others: $\begin{cases} (1) \text{ mixing different classifiers} \\ (2) \text{ generating ensemble of trees by } \underline{\text{boosting}} \end{cases}$

Support vector machines (SVM's)

- Vapnik

Classification as geometric separation.

$$\vec{x}^{(i)} = (x_1^{(i)}, \dots, x_p^{(i)})$$

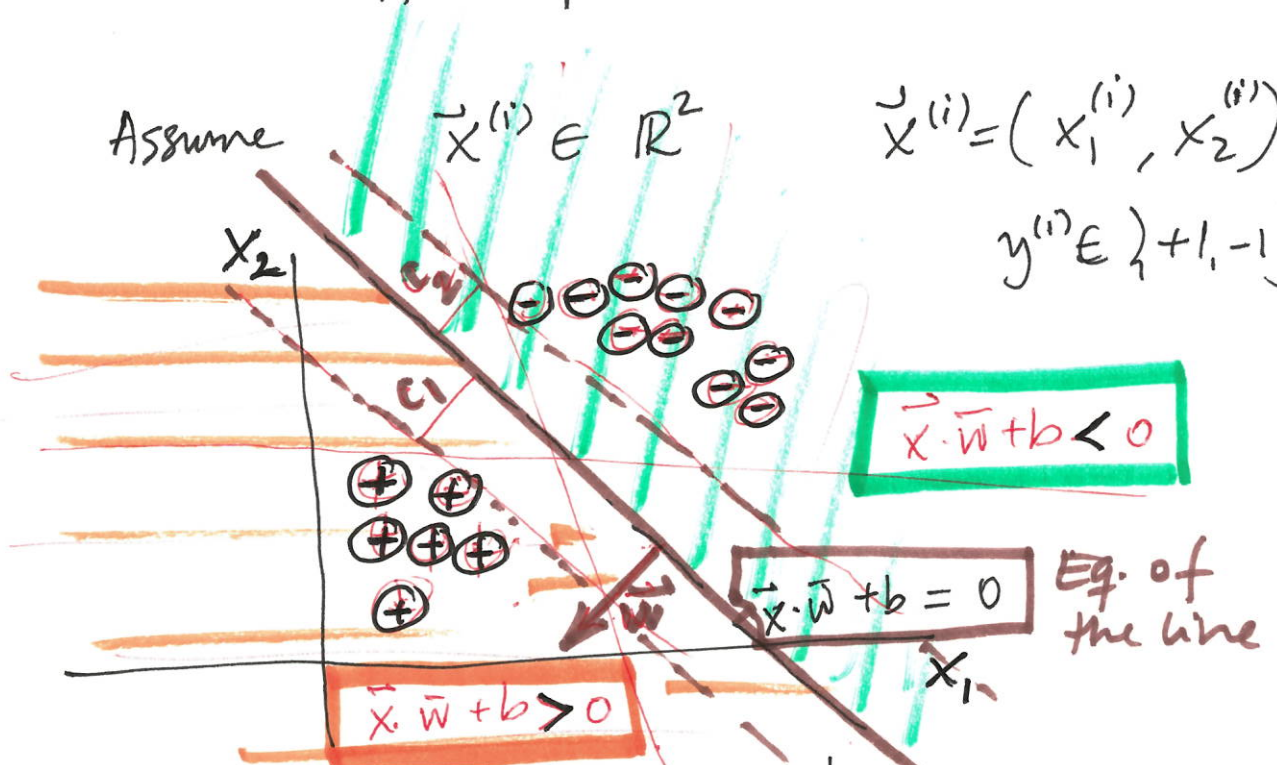
$$y^{(i)} \in \{-1, +1\}$$

Assume

$$\vec{x}^{(i)} \in \mathbb{R}^2$$

$$\vec{x}^{(i)} = (x_1^{(i)}, x_2^{(i)})$$

$$y^{(i)} \in \{+1, -1\}$$



In \mathbb{R}^2 : look for a line : $x_2 = \underline{w} x_1 + \underline{b}$

$$\vec{x} \cdot \vec{w} + b = 0$$

$$\vec{w} = (w, -1)$$

$$(w, -1) \cdot (x_1, x_2) + b = 0 \quad \checkmark$$

For $\vec{x} \in \mathbb{R}^p$

$$\left[\vec{x}^T \cdot \vec{\beta} \right] = \vec{x} \cdot \vec{w} + b = 0 \quad \text{Equation of a hyperplane}$$
$$\vec{x} = \begin{pmatrix} 1 \\ x_1 \\ \vdots \\ x_p \end{pmatrix} \quad \vec{\beta} = \begin{pmatrix} b \\ w_1 \\ \vdots \\ w_p \end{pmatrix} = \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_p \end{pmatrix} \quad \vec{w} = \begin{pmatrix} w_1 \\ \vdots \\ w_p \end{pmatrix}$$

Find the hyperplane

$$\vec{x} \cdot \vec{w} + b = 0$$

(\vec{w}, b) to be found.

$$c_1, c_2 \in \mathbb{R}^+ \quad \left\{ \begin{array}{l} \vec{x}^{(i)} \cdot \vec{w} + b > c_1 \quad \text{if } y^{(i)} = +1 \\ \vec{x}^{(i)} \cdot \vec{w} + b < -c_2 \quad \text{if } y^{(i)} = -1 \end{array} \right\}$$