# Lecture 20 : Classification Trees and RFs

$X_2$

in each rectangle, predict majority class
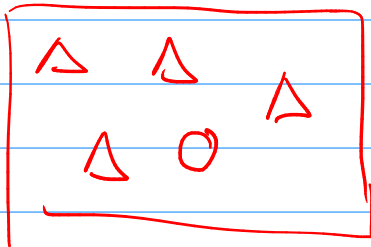
$X_1$

What makes a good split for class tree?

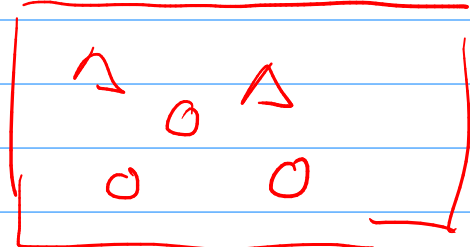Regression $\rightsquigarrow$ RSS

Classification $\rightsquigarrow$ reduce node impurity
increase node purity

Ex,

pure

impure

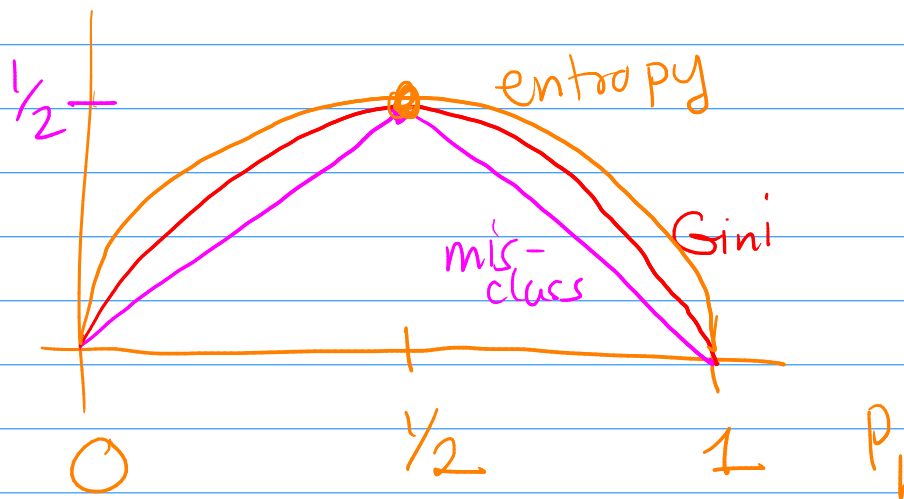Node impurity measures: $p_k$ = pct. of class $k$ in my rectangle

(1) mis-class rate : $1 - P_{\hat{k}}$ , $\hat{k} = \underset{k}{\text{argmax}} \; P_k$
= maj. class

② Gini - Index

$$\sum_k p_k(1-p_k)$$

③ Entropy: $\sum_k p_k \log p_k$

$\underline{K=2}$



— — — — — — — — — — — — — — —

## Categorical Vars

Splitting a cat var is just dividing cats into two groups

$$1, 2, 3, 4, 5$$



$$3, 5, 2 \qquad 1, 4$$

If I have $q$ levels then there are $2^{q-1}-1$ possible splits.

CARTs can deal w/ missing data very nicely.

cat. vars, just add a "missing" category

numeric vars. — Keep track of "surrogate" splits using other vars that divide the data similarly

Problem w/ CART they are really easy to over-fit

Tend to be low bias, high variance.

Recap of properties of means

If I have $X_n$ all w/ the same mean $\mu$, and same variance $\sigma^2$.

Let $\rho$ the correlation among them.

Consider $\bar{X} = \frac{1}{N} \sum_{n=1}^{N} X_n$

① $E\bar{X} = E\left[\frac{1}{N} \sum_n X_n\right] = \frac{1}{N} \sum_n E X_n = \frac{1}{N} N\mu = \mu$

② $\text{Var}(\bar{X}) = \text{Var}\left(\frac{1}{N} \sum_n X_n\right)$

$$= \frac{1}{N^2} \text{Var}\left(\sum_n X_n\right)$$

$$= \frac{1}{N^2}\left[\sum_n \underbrace{\text{Var}(X_n)}_{\sigma^2} + \sum_{i \neq j} \text{Cov}(X_i, X_j)\right]$$

$\text{Cov}(X_i, X_j)$
$$= \text{Cor}(X_i, X_j)$$
$$\cdot \sqrt{\text{Var}(X_i)} \sqrt{\text{Var}(X_2)}$$
$$= \rho\sigma^2$$

$$= \frac{1}{N^2}\left[N\sigma^2 + N(N-1)\rho\sigma^2\right]$$

$$= \frac{\sigma^2}{N} + \frac{(N-1)}{N}\rho\sigma^2$$

$$= \cdots$$

$$= \sigma^2\rho + \frac{\sigma^2}{N}(1-\rho)$$

If $\rho = 0$ then $\text{Var}(\bar{X}) = \sigma^2/N$

If $\rho = 1$ then $\text{Var}(\bar{x}) = \sigma^2$

# Bagging: Ensemble Method

↑ *Bootstrap Aggregating*

↑ *combining many learning methods*

**(1) Draw a series of Bootstrap samples**

Assume training: $\{(x_n, y_n)\}_{n=1}^{N}$

Sample $B$ bootstrap samples

For $b = 1, \ldots, B$

  draw a sample of $N$ training pts w/ replacement

Call these samples $S_1, \ldots, S_B$

**(2) Train a series of methods on each resample**

For $b = 1, \ldots, B$

  $\hat{f}_b = $ method fit to $S_b$

**(3)** combine these to make a ensemble $\hat{f}$

**(i)** Regression : $\hat{f}(x) = \frac{1}{B} \sum_{b=1}^{B} \hat{f}_b(x)$

**(ii)** Classification! $\hat{f}(x) = $ most common predicted class among $\hat{f}_b(x)$

(plurality)

**Binary:**

If $\hat{f}_b(x) \in \{\pm 1\}$

$$\hat{f}(x) = \text{Sign}\left(\sum_{b=1}^{B} \hat{f}_b(x)\right)$$

## Why does this work?

For regression

$$MSE(\hat{f}) = \text{Bias}(\hat{f})^2 + Var(\hat{f})$$

$$\text{bias}(\hat{f}) = \mathbb{E}\hat{f} - f$$

For **bagging** in regression

$$\text{bias}(\hat{f}) = \mathbb{E}\hat{f} - f = \mathbb{E}\left[\frac{1}{B}\sum_{b=1}^{B} \hat{f}_b(x)\right] - f(x)$$

$$= \mathbb{E}\left[\hat{f}_b(x)\right] - f(x)$$

$$= \text{bias}(\hat{f}_b)$$

So bagging doesn't change bias.

However,

$$\text{Var}(\hat{f}) = \rho \sigma^2 + (1-\rho) \frac{\sigma^2}{B}$$

$$\left[ \text{Var}(\hat{f_b}) = \sigma^2, \quad \text{Cor}(\hat{f_b}, \hat{f_{b'}}) = \rho \right]$$

and so if $\rho \approx 0$ then

$$\text{Var}(\hat{f}) \approx \text{Var}(\hat{f_b}) / B.$$

So bagging <u>reduces</u> variance.

So to make bagging effective we need

① to build $\hat{f_b}$ that are $\approx$ uncorrelated

② want to bag $\hat{f_b}$ that are high var and low bias (e.g. trees!)

<u>Random Forest</u> : basically bagged set of trees
(bagged randomized trees)

# RF: algo

(1) Fit B trees (randomized trees)

For $b = 1, \ldots, B$

<span style="color:red">help make $\rho$ between trees smaller</span>

(i) bootstrap sample from training data

(ii) Fit a randomized tree to this bootstrap sample — at each point where I split I only consider a random subset of variables

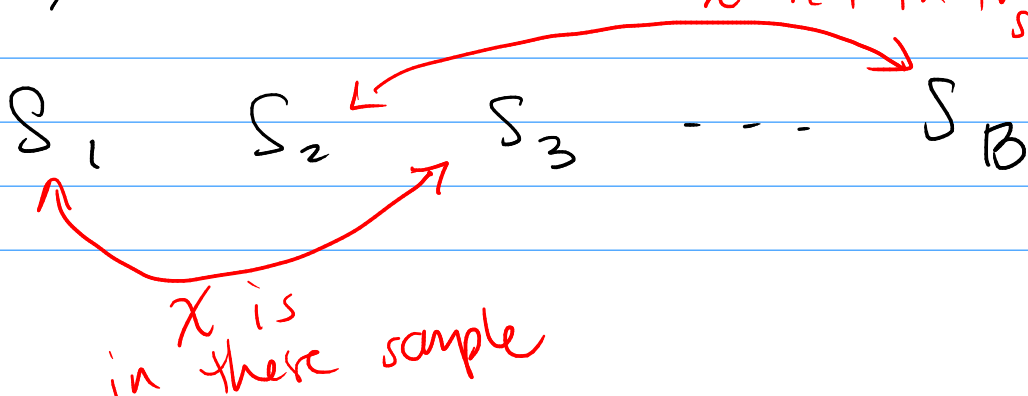<span style="color:orange">↰ fit a huge tree</span>

(2) bag the trees

---

## Out-of-bag Error (OOB)

<span style="color:red">↰ estimate of my test error</span>

When I bootstrap sample, for any point $x$ it will end up in same of my bootstrap samples, and not others <span style="color:red">$x$ not in these samples</span>

$$S_1 \quad S_2 \quad S_3 \quad \cdots \quad S_B$$

<span style="color:red">$x$ is in there sample</span>

Consider bagging only those trees trained on $S_b$ not including $x$ : $\hat{f}_{-x}$

As far as $\hat{f}_{-x}$ is concerned, $x$ is a validation point and so

$$\hat{y}_{OOB} = \hat{f}_{-x}(x)$$

then $y - \hat{y}_{OOB}$ is an est of my test/val error for $y$

So if I do this for all pts, I can est test err using these OOB predictions.