# Decision trees

- Constructing
- pruning (cost-complexity pruning)
  ↑
  depends on some $\alpha \geq 0$



A: 60% samples
B: 50% samples

How to choose a reasonable $\alpha$ for pruning?

- test on the testing data

  usually one sees the following behaviour:



on the training data

- use cross-correlation: we divide the training dataset into $K$ parts, $A_1, A_2, \ldots, A_k$
  $$A_i \cap A_j = \emptyset \text{ for } i \neq j, \quad \cup A_i = \{\text{training dataset}\} \quad (\text{randomly})$$
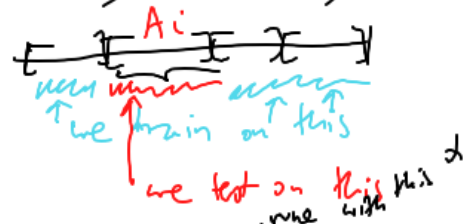  for each $i = 1, 2, \ldots, k$: train the tree on $\bigcup_{j \neq i} A_j$

  and test it on $A_i$
  check the performance of diff. $\alpha$'s

  find
  → $\alpha$ → train tree on the original training set and prune with this $\alpha$

$A_i$

we train on this

we test on this

# Confusion matrix

|  | actual − | actual + |
|---|---|---|
| predicted − | $a_{11}$ | $a_{12}$ |
| predicted + | $a_{21}$ | $a_{22}$ |

$a_{11} = \#\{$ samples in class − correctly predicted $\}$

$a_{12} = \#\{ \; -11 \underline{\quad\quad} + \;$ predicted as − $\}$

e.g.

|  | actual − | actual + |
|---|---|---|
| predicted − | 93 | 5 |
| + | 7 | 45 |

|  | actual − | actual + |
|---|---|---|
| predicted − | 96 | 11 |
| + | 4 | 39 |

cost $C_{12}$ = 9

cost $C_{21}$ = 1

Often errors of one type are more costly then those of the other type

Decision trees (and many other classifiers) can be easily converted to cost sensitive classifiers by changing the way we make final classification
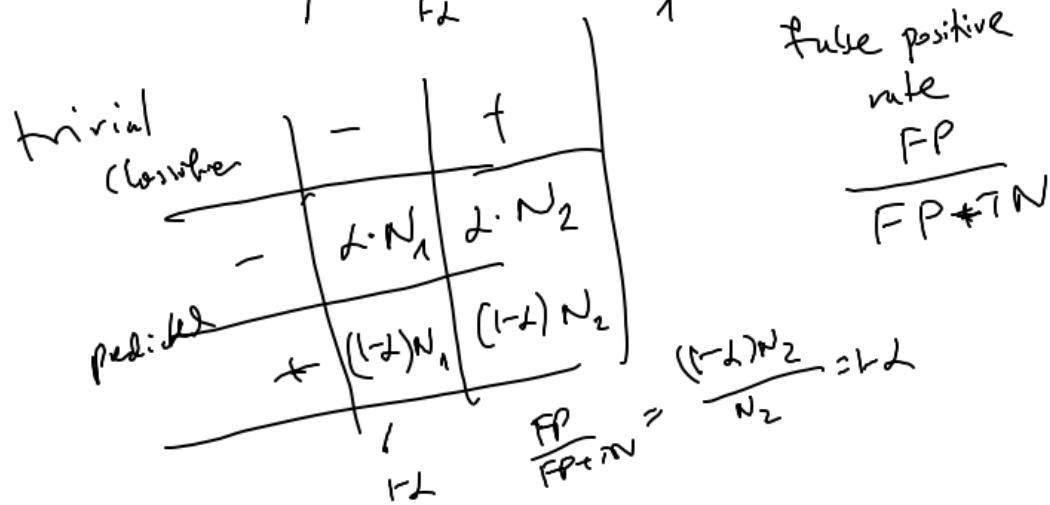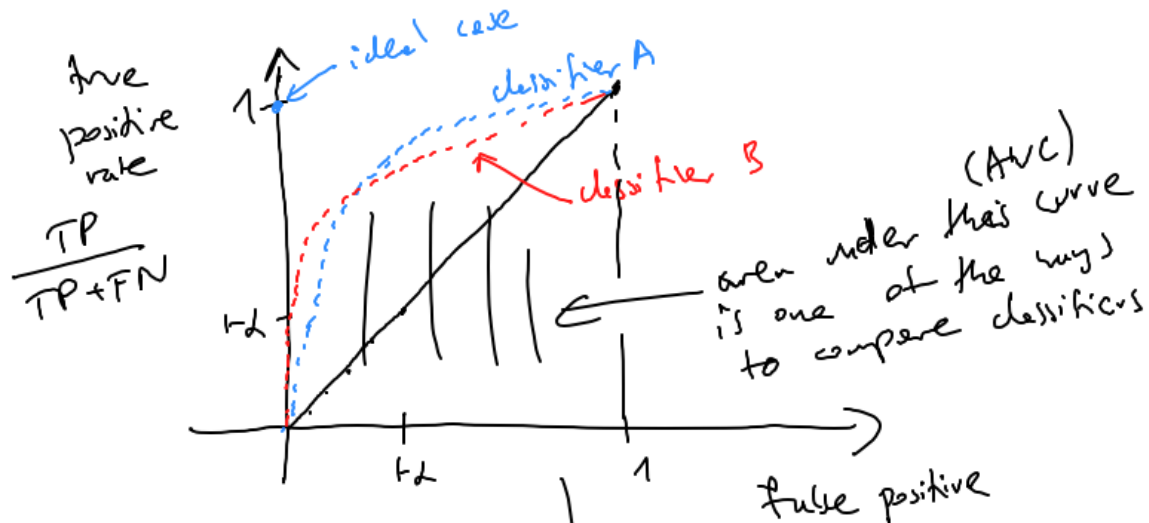
one may classify this as diff class "+"
if $\dfrac{\#\{\text{samples }+\}}{\#\{\text{samples}\}} > \dfrac{C_{12}}{C_{12}+C_{21}} = p^*$

$\dfrac{4}{15}$     $\dfrac{11}{10}$

So far we used $p^* = \dfrac{1}{2}$

"11 −" samples
"4 +" samples

# ROC curve
(receiver operating characteristic)



true positive rate
$$\frac{TP}{TP+FN}$$

ideal case

classifier A

classifier B

(AUC)
area under this curve is one of the ways to compare classifiers

false positive rate
$$\frac{FP}{FP+TN}$$

actual

|  | − | + |
|---|---|---|
| predicted − | TN | FN |
| + | FP | TP |

$$\frac{TP}{TP+FN}$$

trivial classifier:
choose some $\alpha$
predict a sample to be of
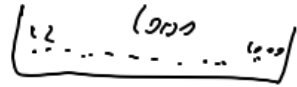class − with prob $\alpha$
       + with prob $1-\alpha$

| pr. | actual − | + |
|---|---|---|
| − | $\alpha \cdot N$ | |
| + | | $(1-\alpha)N$ |

trivial classifier

|  | − | + |
|---|---|---|
| predicted − | $\alpha \cdot N_1$ | $\alpha \cdot N_2$ |
| + | $(1-\alpha)N_1$ | $(1-\alpha)N_2$ |
|  | 1 | |
|  | $1-\alpha$ | |

$$\frac{FP}{FP+TN} = \frac{(1-\alpha)N_2}{N_2} = 1-\alpha$$
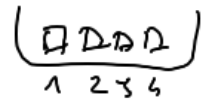
# Bagging (bootstrap aggregation)

A general method, but often used with trees.

1) take B training sets

2) build a separate model for each

3) average / take majority vote of resulting predictions

$\quad\quad\quad \uparrow \quad\quad\quad\quad \uparrow$

$\quad$ regression $\quad$ classification

Ad 1): Draw $\tilde{n}$ samples from the original set of $n$ observation with replacement

$\quad\quad\quad$ (usually $\tilde{n} = n$)

$\quad\quad$ Repeat B times to obtain B training sets

$\left\{\begin{array}{l} \text{If } n \text{ is large and we draw } \tilde{n} = n \text{ samples with replacement,} \\ \text{then we may expect} \approx 0.63 \cdot n \text{ of the original samples} \\ \quad\quad\quad\quad\quad\quad\quad\quad\quad \uparrow \quad\quad\quad\quad \text{unique samples} \\ \quad\quad\quad\quad\quad\quad\quad (1 - \frac{1}{e}) \end{array}\right.$

1: $(1,3,1,4)$ $\left.\begin{array}{l} \\ \\ \end{array}\right\}$ new training sets

2: $(3,3,2,1)$

3: $(4,2,3,1)$

For bagging one may use out-of-bag error for example to find the pruning threshold $\alpha$

- Each sample occurs in ca. 63% of the training datasets
  and does not occur in ca. 37% $-11$ _____

  So this sample $x$ may be used to check the performance of ca. 37%
  of the models that were trained on a training set not containing $x$

# Random forests

$x \in \mathbb{R}^n$ — $n$ predictors

I.e. they are like bagging, but with 2) replaced by

2') select some random subset of ~~predoton~~ $p$ predictors out of $n$
and build a model using just these $p$ predictors
(and a training set selected in step 1)

The advantage: different models will have smaller correlation

example: passenger of Titanic



reasonable choices for $p$:   $p = \lfloor \sqrt{n} \rfloor$ for classification
                               $p = \lfloor n/3 \rfloor$ for regression

(but $p$ is a hyperparameter that can be also subject to optimisation)

In step 2) or 2)' one may take a full tree without pruning
or also a small tree (heavily pruned)