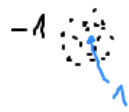


Ada Boost for binary classification



Data set (x_i, y_i) $i=1, \dots, N$ Two classes: $-1, 1$.
 \uparrow \uparrow
 \mathbb{R}^d $\{-1, 1\}$

• $C_0(x) = 0$

• $C_m(x_i) = C_{m-1}(x_i) + \underbrace{d_m}_{\text{constant}} \underbrace{k_m}_{\text{classifier with values in } \{-1, 1\}}(x_i)$

How to choose d_m, k_m ? To minimise the total error E of C_m :

$$E = \sum_{i=1}^N e^{-y_i C_m(x_i)} = \sum_{i=1}^N \underbrace{e^{-y_i C_{m-1}(x_i)}}_{W_i^{(m)}} e^{-y_i d_m k_m(x_i)} = \sum_{i=1}^N W_i^{(m)} e^{-y_i d_m k_m(x_i)}$$

$$= \sum_{i: y_i = k_m(x_i)} W_i^{(m)} e^{-d_m} + \sum_{i: y_i \neq k_m(x_i)} W_i^{(m)} e^{d_m} = \underbrace{\left(\sum_{i=1}^N W_i^{(m)} \right)}_{\text{does not depend on } k_m} e^{-d_m} + \sum_{i: y_i \neq k_m(x_i)} W_i^{(m)} \underbrace{\left(e^{d_m} - e^{-d_m} \right)}_{\text{does not depend on } i \text{ nor on } k_m}$$

So k_m should be chosen to minimise $\sum_{i: y_i \neq k_m(x_i)} W_i^{(m)}$, i.e., weighted error.

S.5 - grade

find some interesting project connected with ML, ~~share~~ share with me, if accepted, carry out and get S.5

Then, given k_m , choose d_m to minimize E :

$$L_m = \frac{1}{2} \ln \frac{\sum_{i: y_i = k_m(x_i)} w_i^{(m)}}{\sum_{i: y_i \neq k_m(x_i)} w_i^{(m)}}$$

After repeating for $m=1, 2, \dots, M$ we obtain some C_m .

Classification is done by $\text{sgn}(C_m(x))$

Drawbacks:

excessive weight is assigned to outliers
(there are some other versions)

Handling missing data

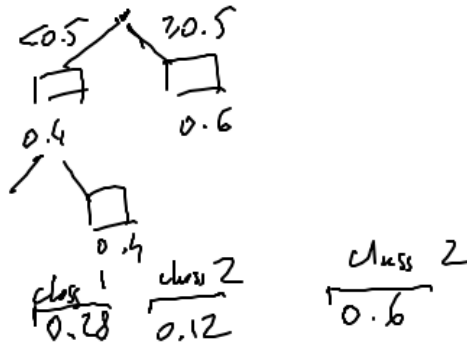
(x_1, \dots, x_n, y)

- ignore records with missing values (it could result in a bias)
- 'guess' the missing attribute value
- treat missing values as a special category
- split the sample to "subsamples" (used in e.g. C4.5)

Ross Quinlan



for the classification:



there are different methods to treat missing data in decision tree, but which is scipy

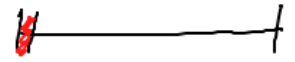
weka

Imbalanced data

classification problem, in which one class is substantially smaller than ^{one of} the others

e.g. testing for a (rare) disease

one method to tackle this problem: use bootstrap to
make the classes balanced



bootstrap



Confusion matrix

(binary classification)

	actual -	actual +
predicted (0) -	60%	10% $C(0,1)$
predicted (1) +	10% $C(1,0)$	20%

$$p^* = \frac{1}{9+1} = \frac{1}{10}$$

cost of misclassifying 1

$$accuracy = 60\% + 20\% = 80\%$$

often, one wants to avoid one of the types of the error, usually that one

* Cost-sensitive classification:

introduce a threshold $p^* = \frac{C(1,0)}{C(0,1) + C(1,0)}$



classify x as of class 1 (i.e., +), if

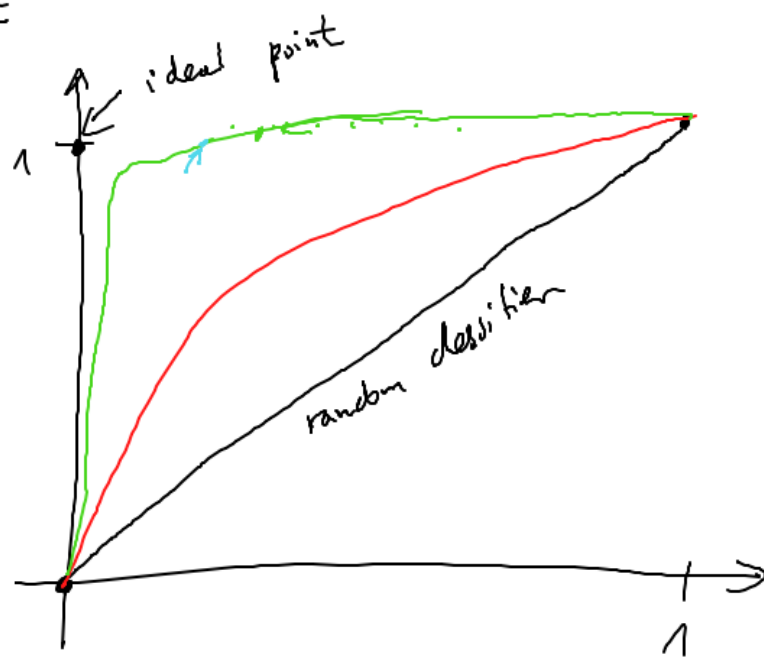
$$P(1|x) \geq p^*$$

↑
estimated prob. of being in class 1

ROC curve (receiver operating characteristic)

true positive rate

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



Confusion matrix

	actual -	actual +	
predicted -	TN	FN	avoid that error
predicted +	FP	TP	

	-	+
-	1- α	α
+	0	0

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

	-	+
-	0	0
+	1- α	α

random classifier:

randomly classify each sample as being in class 1 (with prob. β) or -1 (with prob. $1-\beta$)

$$\frac{TP}{TP + FN} = \frac{\beta \alpha}{\alpha} = \beta$$

$$\frac{FP}{FP + TN} = \frac{\beta(1-\alpha)}{1-\alpha} = \beta$$

1- α	α
$\beta(1-\alpha)$	$\beta\alpha$

There is some tradeoff between FPR and TPR to be made

Logistic regression

binary classification problem

X - predictors ($\in \mathbb{R}^n$), $Y \in \{0, 1\}$
class

model for $p(x) = \Pr(Y=1|x)$:

$$p(x) = \varphi(\underbrace{\beta_0 + \beta_1 x}_{\text{linear regression}}) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

$(\beta_0 \in \mathbb{R}, \beta_1 \in \mathbb{R}^n)$

we take logistic function:

$$\varphi(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{1 + e^x}$$



$$(1 + e^{\beta_0 + \beta_1 x}) p(x) = e^{\beta_0 + \beta_1 x}$$

$$p(x) = e^{\beta_0 + \beta_1 x} (1 - p(x))$$

$$e^{\beta_0 + \beta_1 x} = \frac{p(x)}{1 - p(x)}$$

fitting: maximising likelihood function:

$$\beta_0 + \beta_1 x = \log\left(\frac{p(x)}{1 - p(x)}\right)$$
$$L(\beta_0, \beta_1) = \prod_{i: y_i=1} p(x_i) \cdot \prod_{i: y_i=0} (1 - p(x_i))$$