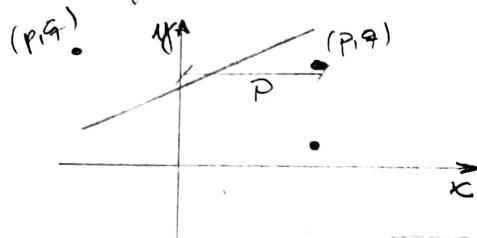


① Data Scientist (udacity)

- linear regression
 - Absolute Trick
 - square Trick

→ Absolute Trick

$$y = w_1 \cdot x + (w_2) \longrightarrow y = (w_1 + p) \cdot x + (w_2 + 1)$$



! we need small steps, let's add a learning rate (α)

$$y = (w_1 + \alpha p)x + (w_2 + \alpha)$$

what if the point is on the bottom of the line?

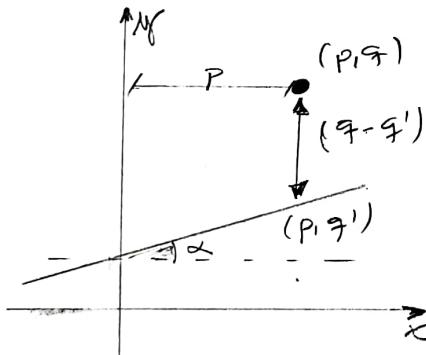
$$y = (w_1 - \alpha p)x + (w_2 - \alpha)$$

why we have p ?

→ IF the point is in the left, p is negative and will make the slope negative, turning clockwise

→ IF p is small, the line will rotate less, IF is big, will rotate more to get closer.

→ Square Trick



$$y = w_1 \cdot x + w_2$$

$$y = [w_1 + p\alpha(q - q')]x + [w_2 + \alpha(q - q')]$$

Consider the distance of horizontal (p) and vertical ($q - q'$).

Übung

$$\textcircled{1} \quad y = -0,6x + 4$$

$$(x_1, y_1) = (-5, 3)$$

$$\alpha = 0,1$$

Absolute trick

$$y = w_1 x + w_2$$

$$y = (w_1 + \alpha p)x + w_2 + \alpha$$

$$y = (w_1 - \alpha p)x + w_2 + \alpha$$

$$y = [-0,6 + 0,1 \cdot (-5)]x + 4 - 0,1$$

$$y = -0,6 + 0,5x + 3,9$$

$$y = -0,1x + 3,9$$

② Square trick

$$y = w_1 \alpha + w_2$$

$$y = -0,6x + 4 \Rightarrow y = 7$$

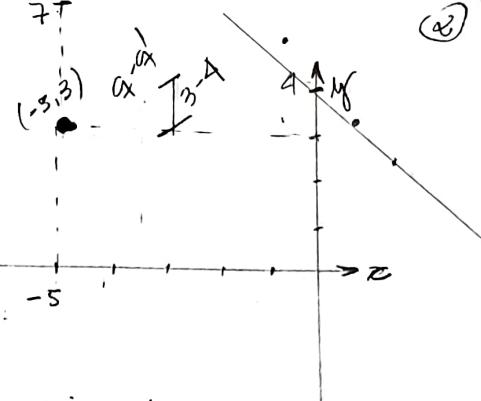
$$(x, y) = (-5, 3)$$

$$y = [w_1 + \alpha(q-q')]x + [w_2 + \alpha(1-q)]$$

$$y = [-0,6 + 5,01(3-7)]x + [4 + 0,01(3-7)]$$

$$y = (-0,6 + 9,2)x + 4 - 0,04$$

$$\underline{y = -0,4x + 3,96}$$

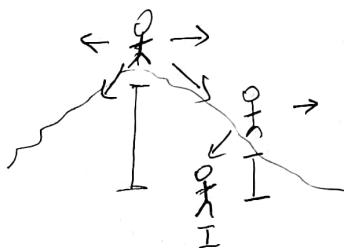
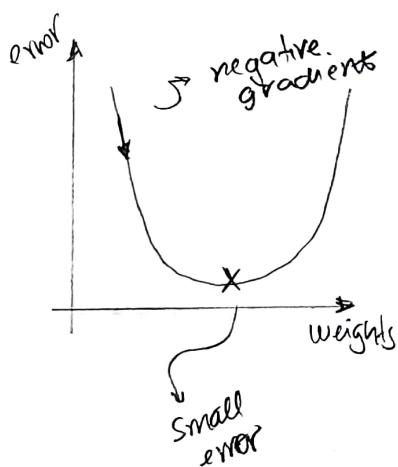


the magnitude by which the intercept and slope change is dependent on how large the error prediction is.

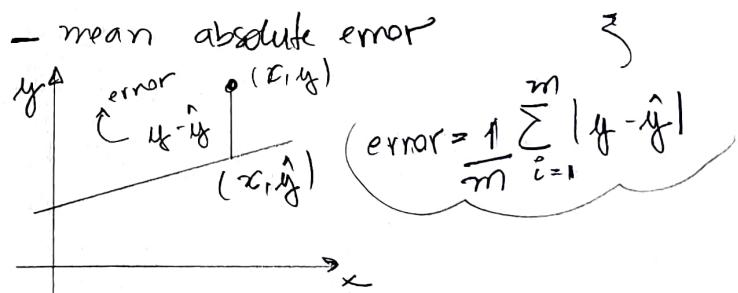
$$\text{error} \Rightarrow 7 - 3 = 4$$

\hat{y}

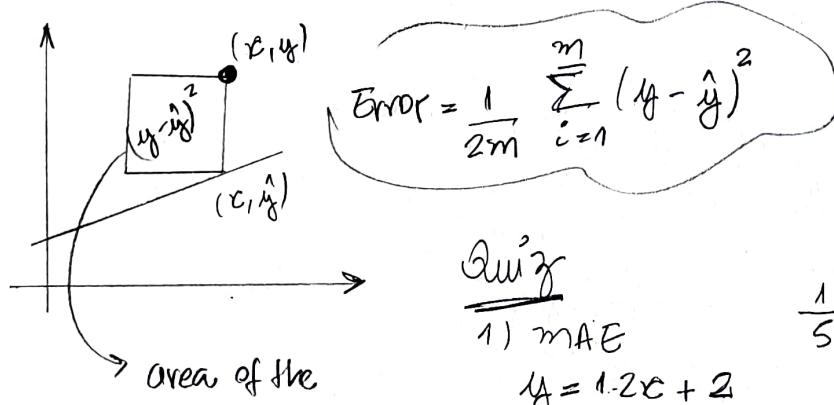
- linear regression
- gradient descent



always errors \hat{y}



- mean squared error



Quiz

1) MAE

$$y = 1,2x + 2$$

$$(2, -2) (8, 14)$$

$$(5, 6)$$

$$(-4, -4)$$

$$(-7, 1)$$

$$\frac{1}{5} [6,4 + 2 + 1,2 + 7,4 + 2,4] = 3,88$$

2) MSE

$$y = 1.2x + 2$$

$$(2, -2)(5, 6)(-4, -4)(-7, 1)(8, 14)$$

$$\begin{aligned} \text{MSE} &= \frac{1}{2m} \cdot \sum_{i=1}^m (y_i - \hat{y}_i)^2 \\ &= \frac{1}{10} \left[6.4^2 + 2^2 + 1.2^2 + 7.4^2 + 2.4^2 \right] \\ &= \underline{\underline{10.692}} \end{aligned}$$

01/08/19

2 algorithms
 ↗ tricks
 ↘ minimize errors
 } Are the same \checkmark

$$\text{MSE} = \frac{1}{2} (y - \hat{y})^2$$

$$\text{Gradient step} \rightarrow w_i \rightarrow w_i - \alpha \cdot \underbrace{\frac{\partial}{\partial w_i} \text{error}}_{-(y - \hat{y}) \cdot x}$$

$$\therefore w_i \rightarrow w_i + \underline{(y - \hat{y}) \cdot \alpha \cdot x}$$

Square trick

$$y = w_1 + p\alpha(\hat{y} - \hat{y})x + w_2 + \alpha(\hat{y} - \hat{y})$$

minimizing the mean square is equal to use the absolute or square trick

Exercise

Calculate the derivative of the error w.r.t w_2 and verify that it is precisely $-(y - \hat{y})$

$$\text{Error} = \frac{1}{2} (y - \hat{y})^2$$

$$\hat{y} = w_1 \cdot x + w_2$$

$$\frac{\partial \text{Error}}{\partial w_2} = \frac{\partial}{\partial w_2} \left(\frac{1}{2} (y - [w_1 \cdot x + w_2])^2 \right)$$

$$= \cancel{\alpha} \cdot \frac{1}{2} (y - w_1 \cdot x - w_2) \cdot (-1)$$

$$= w_1 \cdot x + w_2 - y$$

$$= \hat{y} - y \quad \text{or} \quad \underline{\underline{-(y - \hat{y})}}$$

→ Batch / stochastic gradient descent (SGD)

- calculate the absolute / square tick for all the points, add them, and update the weights with the sum of the values.

→ mini batch gradient descent

- Split the data into small batches, then use each batch to update the weights.

Cooling challenge

$$mse = \frac{1}{2} (y - \hat{y})^2 \quad \frac{\partial mse}{\partial w} = \frac{\partial}{\partial w} \frac{1}{2} (y - w \cdot x - b)^2$$

$$\hat{y} = w \cdot x + b \quad = (y - w \cdot x - b) \cdot (-x) \\ = (-y + \underbrace{w \cdot x + b}_{\hat{y}}) \cdot \underbrace{x}_{\hat{y}}$$

$$\frac{\partial mse}{\partial b} = \frac{\partial}{\partial b} \frac{1}{2} (y - w \cdot x - b)^2 \quad = (-y + \hat{y}) \cdot x \text{ or} \\ = (y - w \cdot x - b) \cdot (-1) \quad - (y - \hat{y}) \cdot x$$

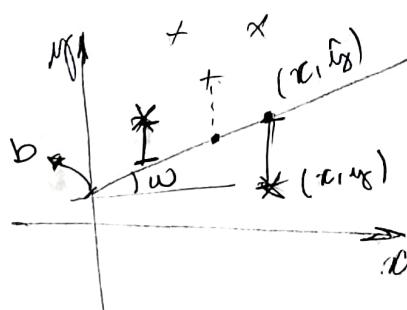
$$= -y + \underbrace{w \cdot x + b}_{\hat{y}} \\ = - (y - \hat{y})$$

$$\therefore \left| \begin{array}{l} \frac{\partial mse}{\partial w} = -(y - \hat{y}) \cdot x \\ \frac{\partial mse}{\partial b} = -(y - \hat{y}) \end{array} \right|$$

for the square tick:

$$y = w + \alpha p(\hat{y} - y) \cdot x + b + \alpha (y - \hat{y})$$

$$mse = \frac{1}{2m} \sum_{i=1}^m (y - \hat{y})^2 \\ = \frac{1}{2m} \sum_{i=1}^m (y_i - (w \cdot x_i + b))^2$$



$$y = w - \underbrace{\alpha x(y - \hat{y})}_{\text{new } w} + b - \underbrace{\alpha(y - \hat{y})}_{\text{new } b}$$

$$\hat{y} = \underbrace{w \cdot x + b}_{\text{points}} \\ \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \end{bmatrix} = \begin{bmatrix} w \cdot \begin{bmatrix} x_1 \\ x_2 \\ \vdots \end{bmatrix} + b \end{bmatrix}$$

$$\hat{y} = w \cdot x + b \Rightarrow \begin{bmatrix} \hat{y}_0 \\ \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{bmatrix}_{n \times 1} = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}_{n \times 1} \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_n \end{bmatrix}_{1 \times n} + \begin{bmatrix} b \\ b \\ \vdots \\ b \end{bmatrix}_{n \times 1}$$

$$\text{error} = y - \hat{y} \Rightarrow \begin{bmatrix} \text{error}_0 \\ \text{error}_1 \\ \vdots \\ \text{error}_n \end{bmatrix} = \begin{bmatrix} y_0 \\ y_1 \\ \vdots \\ y_n \end{bmatrix} - \begin{bmatrix} \hat{y}_0 \\ \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{bmatrix}$$

$$w = w + \alpha \cdot x \cdot \text{error} \Rightarrow \begin{bmatrix} \hat{w}_0 \\ \hat{w}_1 \\ \vdots \\ \hat{w}_n \end{bmatrix} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_n \end{bmatrix} + \begin{bmatrix} \alpha \\ \alpha \\ \vdots \\ \alpha \end{bmatrix} \begin{bmatrix} \text{error}_0 \\ \text{error}_1 \\ \vdots \\ \text{error}_n \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$$

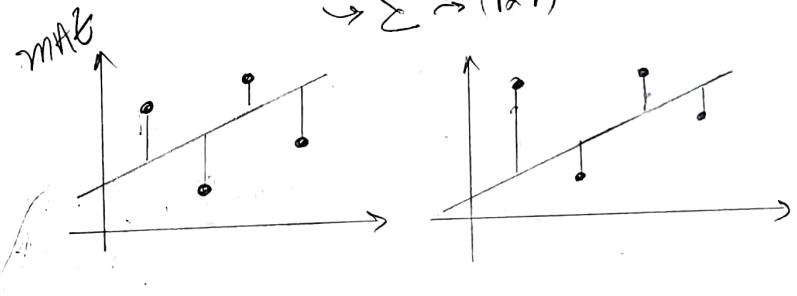
$$\hat{b} = b + \alpha \cdot \text{error}_V \Rightarrow \begin{bmatrix} \hat{b} \end{bmatrix} = \begin{bmatrix} b \end{bmatrix} + \alpha \cdot \begin{bmatrix} \text{error}_0 \\ \text{error}_1 \\ \vdots \\ \text{error}_n \end{bmatrix}_{1 \times 20} \quad \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}_{20 \times 1}$$

$$\sum \rightarrow (1 \times 1)$$

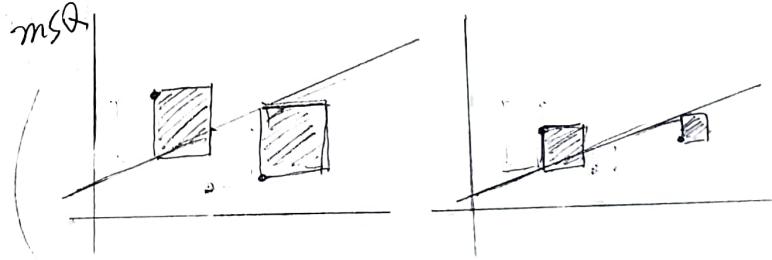
Dm3 2

$$\text{MAE} = \frac{1}{2m} \cdot \sum_{i=1}^m |y_i - \hat{y}_i|$$

$$\text{MSE} = \frac{1}{2m} \cdot \sum_{i=1}^m (y_i - \hat{y}_i)^2$$



Same MAE if the line goes up / down



→ MSQ the error is minimal when the line is in the middle.

Polynomial Regression

$$y = w_1 x^2 + w_2 x + w_0$$

→ Scikit learn

- Polynomial Features

▷ Predictor features

⇒ w and b
for linear regression

▷ outcome feature
⇒ result y

→ Regularization

• Adds model complexity into the error

Regularization L1. \rightarrow sum the absolute coefficient values.

Regularization L2 \rightarrow sum the squared coefficient values.

λ is used to give importance to the complexity error

$\uparrow \lambda \Rightarrow \uparrow$ complexity error

$$\text{error} = \frac{\text{model error}}{\text{model complexity}}$$

$\downarrow \lambda \Rightarrow \downarrow$ complexity error

Regularization (L1)
importance assigns zero to those coefficients with less

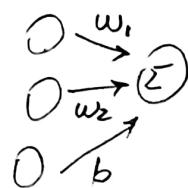
→ feature scaling

important for algorithms that use distance comparison
(k-nn, svm, etc)

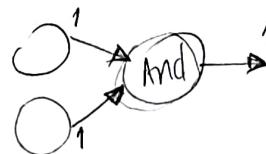
use Standard Scaler to fit and transform the coefficients

→ Classification

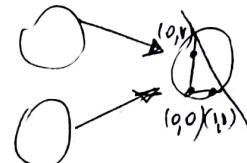
- Perceptron



\rightarrow logical operations.



IN	IN	OUT
1	1	1
1	0	0
0	1	0
0	0	0



for and

$$A \cdot w_0 + B \cdot w_1 = 0$$

$$1 \cdot 1 + 1 \cdot 1 - 2 = 0$$

$$1 \cdot 0 + 1 \cdot 1 - 2 = -1$$

$$1 \cdot 0 + 1 \cdot 0 - 2 = -2$$

$$\begin{aligned} w_0 &= 1 \\ w_1 &= 1 \\ b &= -2 \end{aligned}$$

→ for or

$$1 \cdot 1 + 1 \cdot 1 - 1 = 1$$

$$1 \cdot 0 + 1 \cdot 1 - 1 = 0$$

$$1 \cdot 0 + 1 \cdot 0 - 1 = -1$$

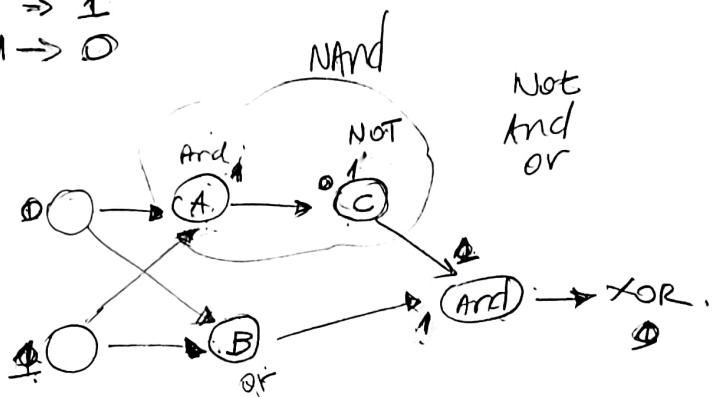
$\Rightarrow \text{not}$

$$A \cdot w_1 + B \cdot w_2 + b =$$

$$\begin{aligned} 0 \cdot (-1) + 0 &= -0 \Rightarrow 1 \\ 1 \cdot (-1) + 0 &= -1 \Rightarrow 0 \end{aligned}$$

$\Rightarrow \text{xor}$

1	1	0
1	0	1
0	1	1
0	0	0



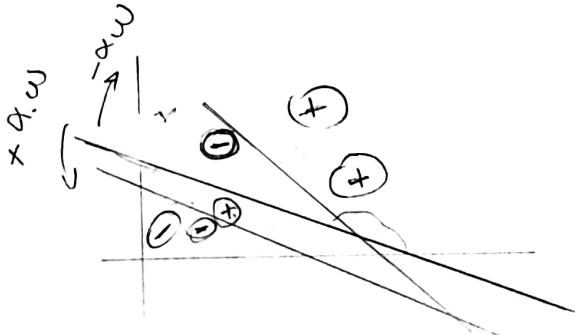
$$3x_1 + 4x_2 - 10 = 0$$

$$P = (1, 1)$$

$$\begin{array}{r}
 3 \quad 4 \quad - 10 \\
 + 1.0, 1 \quad 1.0, 1 \quad 1.0, 1 \\
 \hline
 3,1 \quad 4,1 \quad -9,9 \quad \geq 0 ? \text{ no} \\
 0,1 \quad 0,1 \quad 0,1 \\
 \hline
 3,2 \quad 4,2 \quad -9,8 \quad \geq 0 \quad \text{no}
 \end{array}$$

$$\begin{array}{r}
 4,2 \\
 3,2 \\
 \hline
 7,4 \\
 \hline
 2,1 \quad 3 \\
 \hline
 8
 \end{array}$$

$$\begin{array}{r}
 3 \quad 4 \quad - 10 \\
 + 1 \quad 1 \quad 1 \\
 \hline
 4 \quad 5 \quad -9 \quad \geq 0 ? \quad \text{sum} \\
 \hline
 \end{array}$$



$$\begin{aligned}
 \hat{y} &= x \cdot w + b \\
 \hat{y} &= x_1 w_1 + x_2 w_2 + b
 \end{aligned}$$

$$1 \times 1 = 1 \times 2 \cdot 2 + 1 + 1 \times 1$$

Decisions trees.

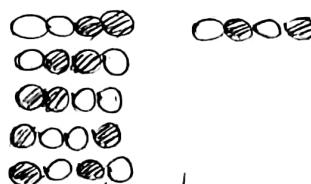
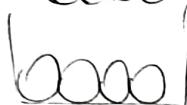
→ Entropy → How much freedom has to move around.

High → can move more.

Example



oooo



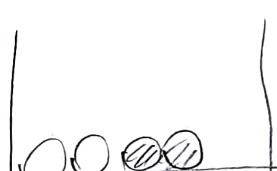
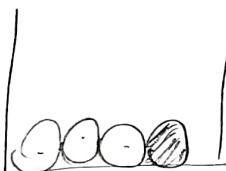
high knowledge / low

middle knowledge

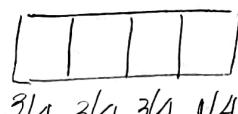
high knowledge

! As more homogenous less entropy

As more knowledge, less entropy !



$$1 \times 1 \times 1 \times 1 = 1$$



$$0,75 \times 0,75 \times 0,75 \times 0,25 = 0,105$$



$$4 \cdot 0,5 = 0,0625$$

Products are bad,
many elements, numbers
get too small !

Sum is
better

use logarithms

$$3 \cdot \log(0,75) + \log(0,25) = -3,245.$$

$$\log(a \cdot b) = \log a + \log b$$

$$\text{Entropy} = -\frac{m}{m+n} \cdot \log\left(\frac{m}{m+n}\right) - \frac{n}{m+n} \cdot \log\left(\frac{n}{m+n}\right)$$

Entropy

4 red balls

10 blue balls.

14 balls

$$\text{Entropy} = \frac{-4}{14} \cdot \log\left(\frac{4}{14}\right) - \frac{10}{14} \cdot \log\left(\frac{10}{14}\right) =$$

$$= 0,104$$

$$\text{entropy} = -p_1 \cdot \log(p_1) - p_2 \cdot \log(p_2)$$

Quiz

$$\text{entropy} = -p_1 \cdot \log(p_1) - p_2 \cdot \log(p_2) - p_3 \cdot \log(p_3)$$

3 blue

2 yellow

8 red.

$$\text{entropy} = -\frac{3}{13} \cdot \log\left(\frac{3}{13}\right) - \frac{2}{13} \cdot \log\left(\frac{2}{13}\right) - \frac{8}{13} \cdot \log\left(\frac{8}{13}\right)$$

$$\underbrace{-\frac{3}{13} \cdot (-2,1959)}_{0,2307} - \underbrace{\frac{2}{13} \cdot (-0,7008)}_{0,1538} - 0,6153 \cdot (-0,7006)$$

$$= \underline{\underline{1,334}}$$

> information gain

$$\text{information gain} = \text{parent's entropy} - 0,5 \cdot (\text{children 01} + \text{children 02})$$

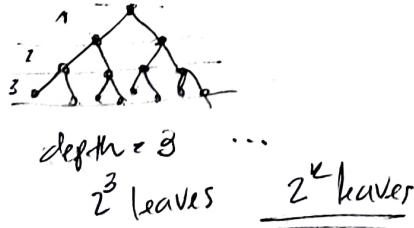
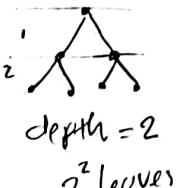
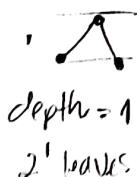
> decision trees.

→ hyperparameters. { maximum depth .

→ maximum depth

largest possible length between
the root to a leaf.

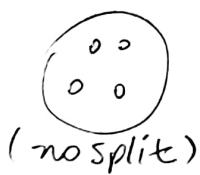
A tree of length k can have at most
 2^k leaves.



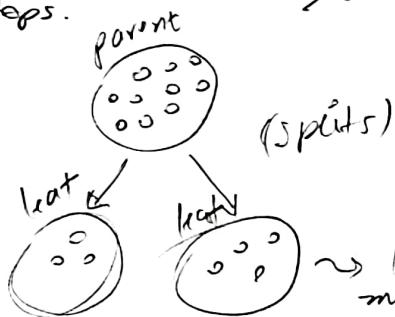
→ minimum number of samples to split

node \Rightarrow at least min-samples-split samples to be large enough to split.

⇒ If node has fewer than min-samples-split, the splitting process stops.



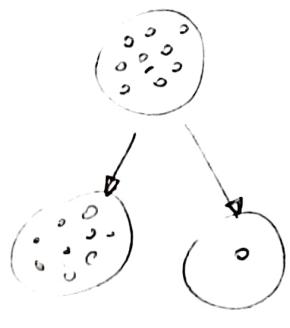
$$\text{min-samples-split} = 5$$



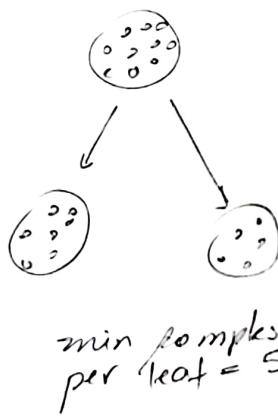
→ less than min-samples-split

no control of the minimum size of leaves.

→ minimal number of samples per leaf



$$\text{min samples per leaf} = 1$$



$$\text{min samples per leaf} = 5$$

▷ Avoid 9 samples in one leaf and 1 in another.

int : # of samples in a leaf

float : % of samples in a leaf

Guys

Underfitting / overfitting vs Feature size

Small max depth

→ underfitting

Large max depth

→ overfitting

Small min samples per split

→ overfitting

Large min samples per split

→ underfitting

too simple model.

too deep trees can memorize the data

complicated, highly branched tree, model memorized the data.

→ not enough flexibility to build the tree.

Sklearn

Decision tree classifier

Fitting the tree means finding the best tree that fits the data



→ titanic kaggle
using decision trees.

(11)

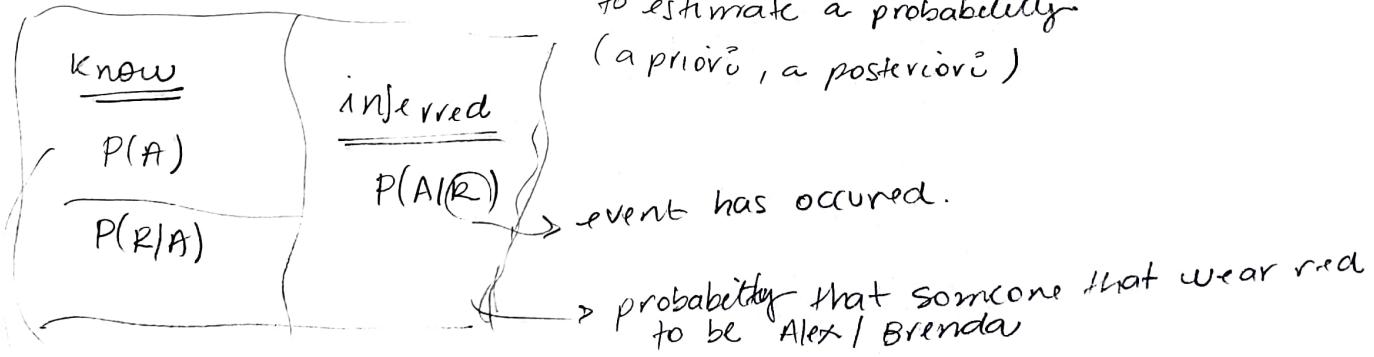
One-Hot-encoding → pandas.get_dummies (features).

! Achieving accuracy ↓ test accuracy = Overfitting

- reduce depth of the tree. ↘ → reduce memory and model complexity.
- increase min-sample-split ↑ → reduce number of branches.
- increase min-samples-leaf ↑ → reduce splitting, since it will stop splitting if the leaf has less than min-samples-leaf (will try different thresholds, though)

I Naive Bayes

→ Bayes theorem: use more information to estimate a probability (a priori, a posteriori)



→ probability Alex / Brenda wear red.

3/4

Alex		Brenda	
3 day/week		1 day/week	
w ₁	w ₂	w ₃	w ₄
A	A	A	B
A	A	A	B
A	A	A	B
A	A	B	B

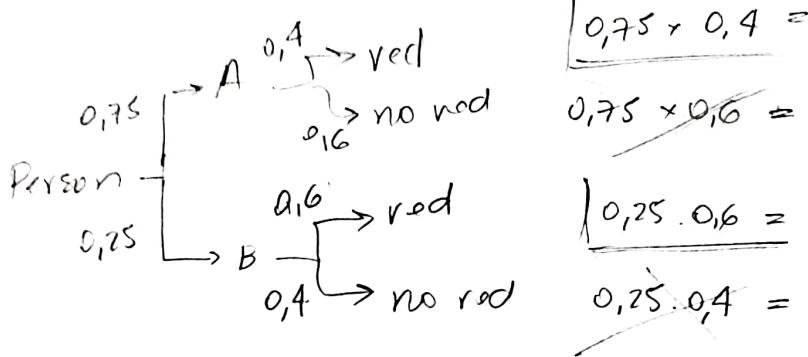
→ red sweater.

Prior: $P(\text{Alex}) = 0,75$
 $P(\text{Brenda}) = 0,25$
 $\rightarrow 1/4$

Person had a red sweater.
 Alex wears 2x a week.
 Brenda wears 3x a week.

Posterior: $P(\text{Alex}) = \frac{6}{9} = \frac{2}{3}$

$P(\text{Brenda}) = \frac{3}{9} = \frac{1}{3}$



$$P(A|R) = \frac{0,75 \cdot 0,4}{0,75 \cdot 0,4 + 0,25 \cdot 0,6} = 0,67 \quad \left. \begin{array}{l} \text{Bayes} \\ \text{normalize} \end{array} \right\}$$

$$P(B|R) = \frac{0,25 \cdot 0,6}{0,75 \cdot 0,4 + 0,25 \cdot 0,6} = 0,33 \quad \left. \begin{array}{l} \text{normalize} \end{array} \right\}$$

Event -

$P(A)$ → A	$\frac{P(R A)}{P(R^c A)}$ - R $P(R \cap A)$ $\frac{P(R^c A)}{P(R^c \cap A)}$ - R ^c $P(R^c \cap A)$ <small>c: complement</small>	$P(RNA) = P(A) \cdot P(R A)$
------------	--	------------------------------

$$P(RNA) = P(A) \cdot P(R|A)$$

$P(B)$ → B	$\frac{P(R B)}{P(R^c B)}$ - R $P(R \cap B)$ $\frac{P(R^c B)}{P(R^c \cap B)}$ - R ^c $P(R^c \cap B)$	$P(RnB) = P(B) \cdot P(R B)$
------------	--	------------------------------

$$P(RnB) = P(B) \cdot P(R|B)$$

Should
normalize
to sum 1

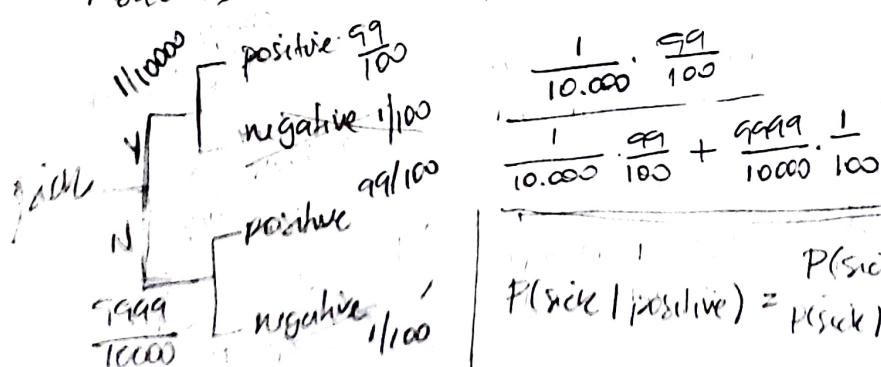
$$P(A|R) = \frac{P(A) \cdot P(R|A)}{P(A) \cdot P(R|A) + P(B) \cdot P(R|B)}$$

$$P(B|R) = \frac{P(B) \cdot P(R|B)}{P(B) \cdot P(R|B) + P(A) \cdot P(R|A)}$$

Quiz

99% accuracy

$$1 \text{ out of } 10.000 \Rightarrow P(\text{people sick}) = \frac{1}{10.000}$$



Posterior

$$P(\text{sick}) = 0,0001$$

$$P(\text{positive} | \text{sick}) = 0,99$$

$$P(\text{positive} | \text{not sick}) = 0,01$$

$$P(\text{not sick}) = 0,9999$$

$$P(\text{sick} | \text{positive}) = \frac{P(\text{sick}) \cdot P(\text{positive} | \text{sick})}{P(\text{sick}) \cdot P(\text{positive} | \text{sick}) + P(\text{not sick}) \cdot P(\text{positive} | \text{not sick})}$$

Guw3

$$P(\text{easy} | \text{spam}) = ?$$

$$P(\text{easy}) = \frac{2}{5}$$

$$P(\text{money}) = \frac{3}{5}$$

$$P(\text{easy}) \cdot P(\text{spam} | \text{easy})$$

$$\frac{P(\text{easy}) \cdot P(\text{spam} | \text{easy})}{P(\text{easy}) \cdot P(\text{spam} | \text{easy}) + P(\text{not spam}) \cdot P(\text{easy} | \text{not spam})}$$

▷ conditional probability

$$P(A \& B) = P(A \cap B) = P(A) \cdot P(B)$$

$$P(A | B) = P(A) \cdot P(B | A)$$

Naive =
ingenuous

$$P(\text{spam} | \text{easy, money})$$

$$= P(\text{spam}) \cdot P(\text{easy, money} | \text{spam})$$

we are
considering the
variables are
independent

$$P(\text{easy} | \text{spam}) \cdot P(\text{money} | \text{spam})$$

Hot and cold
are dependent

Naive assumption, since
it might be dependents.

$$P(\text{spam} | \text{easy, money}) = P(\text{spam}) \cdot P(\text{easy} | \text{spam}) \cdot P(\text{money} | \text{spam})$$

$$= \frac{3}{8} \cdot \frac{1}{3} \cdot \frac{2}{3} = \frac{3}{36} = \frac{1}{12}$$

$$P(\text{ham} | \text{easy, money}) = P(\text{ham}) \cdot P(\text{easy} | \text{ham}) \cdot P(\text{money} | \text{ham})$$

$$= \frac{5}{8} \cdot \frac{1}{8} \cdot \frac{1}{3} = \frac{1}{40}$$

not the probability,
are proportional because
it is not normalized.

For $1/12$:

$$\textcircled{1} = 1$$

$$1/40 :$$

$$\frac{\frac{1}{40}}{\frac{1}{12} + \frac{1}{40}} = \frac{\frac{1}{40}}{\frac{10+3}{120}} = \frac{\frac{1}{40}}{\frac{13}{120}} = \frac{1}{\frac{13}{120}} = \frac{120}{13} \Rightarrow \text{not spam}$$

$$\frac{\frac{1}{12}}{\frac{1}{12} + \frac{1}{40}} = \frac{\frac{1}{12}}{\frac{10+3}{120}} = \frac{\frac{1}{12}}{\frac{13}{120}} = \frac{1}{\frac{13}{120}} = \frac{120}{13} \Rightarrow \text{spam}$$

- what if we have many features?

- ① $P(\text{spam} | \text{easy, money, cheap}, \dots) \propto P(\text{spam}) \cdot P(\text{easy} | \text{spam})$
 - ② $P(\text{ham} | \text{easy, money, cheap}, \dots) \propto P(\text{ham}) \cdot P(\text{easy} | \text{ham})$
 - $P(\text{money} | \text{spam})$
 - $P(\text{cheap} | \text{spam})$
 - $P(\text{money} | \text{ham})$
 - $P(\text{cheap} | \text{ham})$
 - ③ Normalize
- $\frac{P_1}{P_1+P_2}$ and $\frac{P_2}{P_1+P_2}$

Ques

$$\begin{pmatrix} 1 & 2 & 3 \\ 1 & 5 & 6 \end{pmatrix}$$

3 Standard dice

$$\begin{pmatrix} 2 & 3 & 3 \\ 4 & 4 & 5 \end{pmatrix}$$

2 no standard dice

1 die \Rightarrow was $\alpha 3$.

$$P(\text{standard} | 3) = ?$$

$$P(\text{std} | 3) = P(3) \cdot P(3 | \text{std}) \quad \textcircled{1}$$

$$P(\text{not std} | 3) = P(3) \cdot P(3 | \text{not std}) \quad \textcircled{2}$$

$$\textcircled{1} \left(\underbrace{\frac{3 \cdot 1}{6} + 2 \cdot \frac{2}{6}}_{P(3 | \text{std})} \right) \cdot \underbrace{\left(3 \cdot \frac{1}{6} \right)}_{P(3 | \text{not std})} = \left(\frac{1}{2} + \frac{2}{3} \right) \cdot \frac{1}{2} = \left(\frac{3+4}{6} \right) \cdot \frac{1}{2} = \frac{7}{12}$$

$$\textcircled{2} \left(\frac{3 \cdot 1}{6} + 2 \cdot \frac{2}{6} \right) \cdot \underbrace{\left(2 \cdot \frac{1}{6} \right)}_{P(3 | \text{not std})} = \left(\frac{1}{2} + \frac{2}{3} \right) \cdot \frac{2}{3} = \frac{7}{6} \cdot \frac{2}{3} = \frac{7}{9}$$

$$\textcircled{3} P(\text{std} | 3) = \frac{\frac{7}{12}}{\frac{7}{12} + \frac{7}{9}} = \frac{7}{12} = \frac{1}{12} \cdot \frac{36}{49} = \frac{3}{7} = 0,428 //$$

1,2	3
4,3	3
1,1	2
2,1	2
1,1	1
1,1	1
2,3	6

→ hands-on

Naive Bayes

Example: $P(\text{threat}, \text{page}, \text{bag}, \text{nervous} \dots)$

considered independent

- Spam / not spam
based on joint probabilistic distributions of certain other events (presence of some words in this case).

→ ex: Age and nervousness

{ nervous kid does not seem to be a threat, but a man has a bigger likelihood to be a threat.

→ Interesting points

→ pandas reads tab spaced table
use `read_table`, params:

{ $\text{sep} = '\text{t}'$
 $\text{names} = [\text{'colA'}, \text{'colB'}]$.
↓
indicate cols names.

→ transform categorical vars to integers

using `map` function

`df['label'].map({'ham': 0, 'spam': 1})`.

assign a dict and transform the key to the value.

→ features using Bow

count tokenizer ⇒ separates the words into tokens (ids)

count the some IDs

{ transform to lowercase
remove punctuation (uses regex)
remove stop words (am, an, the..)

→ removing punctuation from string → maps punctuation to none.

`import string`

`string.translate(translator)` → `translator = str.maketrans("", "", string.punctuation)`

→ Python Counter from collections lib

`Counter(list) ⇒ {"hello": 1, "how": 1, "are": 1}`.

→ scikit learn Count Vectorizer

- `sklearn.feature-extraction.text import CountVectorizer`

→ `CountVectorizer.fit(documents)` → transform to a document-term matrix

`CountVectorizer.get_feature_names()`

→ feature names.

row	word1	word2	word3...	frequency of word.
doc #1			○	
doc #2				
doc #3				
doc #4				

for this, use `CountVectorizer.transform(documents)`
use `toarray()` to print

- `pd.DataFrame(doc_array, columns= get_feature_name)`

large dataset → Bow is not ideal

words like is, the, an, would appear a lot and mess with the features.

→ use tf idf instead.

▷ Bayes implementation

priors → probabilities we are aware / was given

postiors → want to compute using priors

$P(D)$ = probability of having diabetes = 0,01

$P(\text{pos})$ = probability of getting \oplus test result

$P(\text{neg})$ = " " \ominus test result

$P(\text{pos}/D)$ = probability of getting \oplus on test, given you have diabetes = 0,9

$P(\text{neg}/\text{noD})$ = " " \ominus on test given you don't have diabetes = 0,9

Sensitivity
True positive
rate

diabetes
= 0,9

Specificity
true neg rate

Bayes

$$\frac{P(A|B)}{P(B)} = \frac{P(A) \cdot P(B/A)}{P(B)} \rightarrow \text{given } P(B/A) = P(\text{Pos}/D)$$

(posterior)

$$P(D/\text{Pos}) \quad \left. \begin{array}{l} P(A) = P(D) \\ P(B) = P(\text{Pos}) \end{array} \right\} \rightarrow \text{given}$$

$$\rightarrow P(D/\text{Pos}) = P(D) \cdot \frac{P(\text{Pos}/D)}{P(\text{Pos})} \rightarrow \text{sensitivity}$$

$$P(\text{pos}) = [P(D) \cdot \text{sensitivity}] + [P(\sim D) \cdot (1 - \text{specificity})]$$

$$\downarrow \qquad \qquad \qquad \downarrow$$

$$P(\text{pos}/D) \qquad \qquad \qquad P(\text{neg}/\sim D)$$

$$\rightarrow P(\sim D/\text{pos}) = P(\sim D) \cdot \frac{P(\text{pos}/\sim D)}{P(\text{pos})} \rightarrow 1 - \frac{P(\text{neg}/\sim D)}{\text{specificity}}$$

▷ Naive Bayes in scikit-learn.

- multinomial Naive Bayes → good for classifying discrete values
- gaussian naive Bayes → better for continuous data as it assumes input data has a Gaussian (normal) distribution.

→ Evaluation model.

- Accuracy : how often the classifier makes the correct prediction

$$Ac = \frac{\# \text{correct}}{\# \text{total pred.}}$$

- Precision : what proportion of messages were classified as spam, actually were spam.

$$Pr = \frac{\# \text{true pos} \rightarrow \text{spam.}}{(\# \text{true pos} + \# \text{false pos})}$$

all classified as spam.

- Recall (positivity): proportion of messages that actually were spam were classified by us as spam

$$\text{Recall} = \frac{\# \text{ true positives}}{\# \text{ true positives} + \# \text{ false negatives}}$$

→ predicted as spam.
all words that are spam.

* Accuracy is not a good metric when the data is highly unbalanced or skewed.

! Naive Bayes is very good on handling a large number of features

→ Support Vector machines (SVM)

- Increase boundary and reduce error.

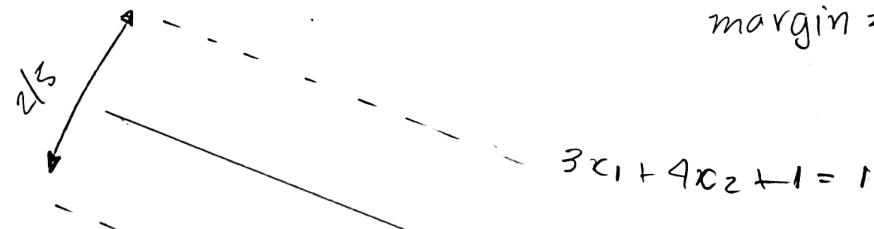
$$\text{margin} = \frac{2}{\|w\|} \quad \|w\| = \sqrt{x^2 + y^2}$$

$$\text{Error} = \|w\|^2 = x^2 + y^2$$

Ex: $w = (3, 4)$ $w_1x_1 + w_2x_2 + b \geq 0$
 $b = 1$ $3x_1 + 4x_2 + 1 \geq 0$

$$\|w\|^2 = 3^2 + 4^2 = 25 = \text{Error}$$

$$\text{margin} = \frac{2}{\sqrt{25}} = \frac{2}{5}$$



$$\text{ERROR} = 25$$

$$3x_1 + 4x_2 + 1 = 0$$

$$3x_1 + 4x_2 + 1 = -1$$

what if weights are :

$$w = (6, 8)$$

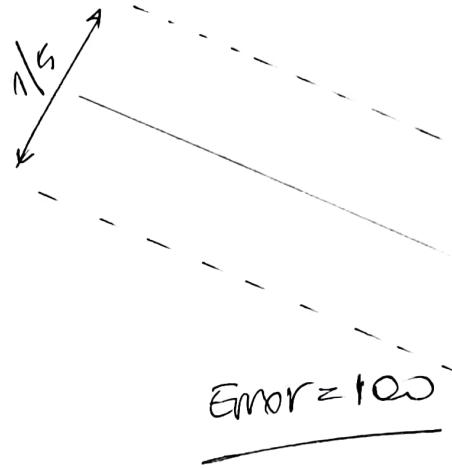
$$b = 2$$

$$w_1 x_1 + w_2 x_2 + b = 0$$

$$6x_1 + 8x_2 + 2 = 0$$

$$\text{Error} = \|w\|^2 = 6^2 + 8^2 = 100$$

$$\text{Margin} = \frac{2}{\sqrt{100}} = \frac{1}{5}$$



Summary:

$$w = (3, 4)$$

$$b = 1 \Rightarrow$$

$$\text{Error} = 25$$

$$\text{Margin} = 2/5$$

$$\left. \begin{array}{ll} w = (6, 8) & \text{Error} = 100 \\ b = 2 & \text{Margin} = 1/5 \end{array} \right\}$$

→ margin error is a norm of the vector w^2

> Error function

$$\boxed{\text{Error} = \text{Classification error} + \text{margin Error}}$$

How to minimize? → Gradient descent?

↑ C ⇒ large margin

↓ Classification Errors

↑ C ⇒ small margin

classifies well.

Obs: $\text{margin} = \frac{2}{\|w\|} \therefore \text{margin} = \frac{2}{\sqrt{\text{Error}}}$

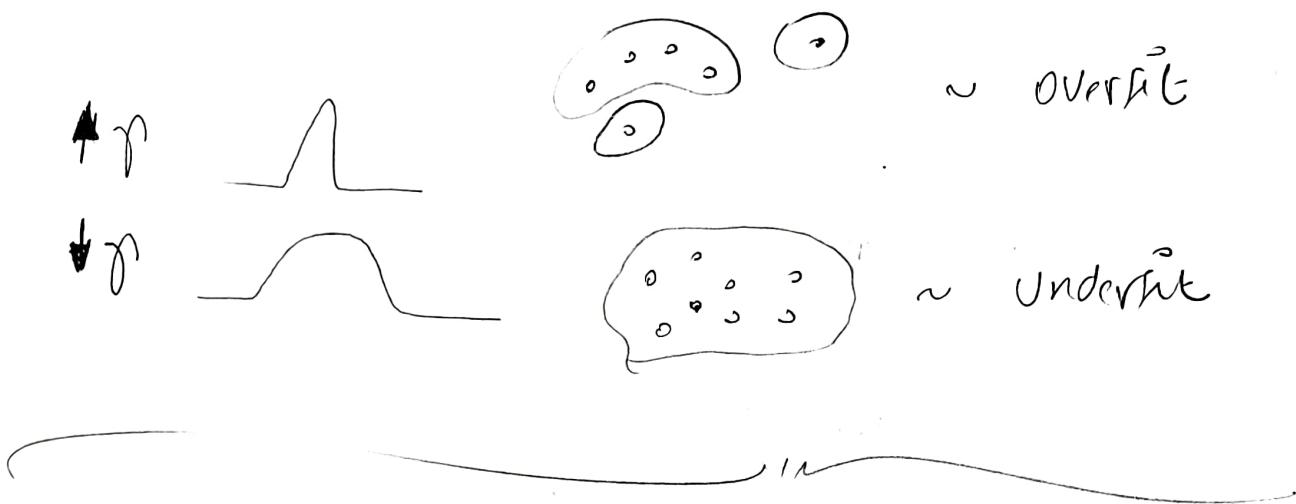
$$\text{Error} = \|w\|^2$$

C = hyperparameter,
use grid search for
optimal value

↑ error ↓ margin
↓ error ↑ margin

▷ Using kernels.

⇒ Radial Basis Function (RBF)



▷ Ensembles

1. Bagging → Get all friends and hear their opinion and combine (Average/ Voting)

2. Boosting → Select best friend for each topic. (explore friend's strengths)

! Combine models for better prediction

Weak and learner (friends) Strong learner (super genius)

Bias and Variance are important.

→ High bias ⇒ does not do a good job of bending to the data.

curve
↓
linear

→ High variance → model changes drastically to meet the needs of every point.

good bias + good variance = better model

- use central limit theorem to reduce bias / variance
- introduce randomness into ensembles

↳ Avoid overfit

1. Sample data with replacement

2. Use features only from a subset

→ Random forest → is a ensemble algorithm

⚠️ Decisions tree → memorize tends to overfit the data a lot

⚠️ Random forest is a set of decision tree models and a voting

→ Bagging

1. Split the data
2. Train in different models
3. Combine them

→ Boosting (AdaBoost)

$$y = \ln\left(\frac{x}{1-x}\right) \quad x = \text{accuracy}$$

↳ weight

+	-
•	• 0
•	0

Acc = 1
(Truthful)

$$w = \ln\left(\frac{1}{1-1}\right)$$

$$w = \infty$$

+	-
•	• 0 0 0
•	0 0 0

Acc = 0,5
(Random)

$$w = \ln\left(\frac{0,5}{0,5}\right) =$$

$$w \rightarrow 0$$

-	+
•	• 0
•	0 0

Acc = 0
(liar)

$$w = \ln\left(\frac{0}{1}\right)$$

$$w = 0$$

Quiz
 $\begin{array}{c} \oplus \\ \ominus \end{array}$

+	-
b	
•	•
•	•
•	•

+	-
•	•
•	•
•	•
•	•

-	+
•	•
•	•
•	•
•	•

$$Acc = \frac{7}{8}$$

$$Acc = \frac{4}{8} = \frac{1}{2}$$

$$Acc = \frac{2}{8} = \frac{1}{4}$$

$$\omega = \ln\left(\frac{x}{1-x}\right)$$

$$= \ln\left(\frac{7/8}{1-7/8}\right)$$

$$= 1.95$$

$$\omega = \ln\left(\frac{x}{1-x}\right)$$

$$= \ln\left(\frac{0,5}{0,5}\right)$$

$$= 0$$

$$\omega = \ln\left(\frac{x}{1-x}\right)$$

$$= \ln\left(\frac{0,25}{0,75}\right)$$

$$= -1.11$$

What about this?

+	-
•	•
•	•
•	•
•	•

-	+
•	•
•	•
•	•
•	•

$$Acc = \frac{8}{8} = 1$$

$$\omega = \ln\left(\frac{1}{1-1}\right) =$$

$$\ln\left(\frac{1}{0}\right) \sim \infty$$

$$Acc = 0$$

$$\omega = \ln\left(\frac{0}{1}\right)$$

$$\ln\left(\frac{0}{1}\right) \sim -\infty$$

do the complete opposite

→ Project (01) - finding donors.

I chose the following classifiers:

- Logistic Regression
- Random Forest
- AdaBoost

) Ensembles

logman's terms

- describe complex things with words to make everybody understand.

classifier	Strengths	weakness
Logistic Regression	<ul style="list-style-type: none"> • linear model • light • easy to expand data using SVC 	<ul style="list-style-type: none"> • hard to capture non-linear data relationship
Random Forest	<ul style="list-style-type: none"> • can model non-linear boundaries • parallel training each decision tree • reduce overfitting 	<ul style="list-style-type: none"> • hard to interpretate
AdaBoost	<ul style="list-style-type: none"> • use different classifiers as weak classifiers • simple to implement • resistant to overfitting 	<ul style="list-style-type: none"> • training time is consuming • imbalanced class problems. • needs a termination condition

→ STEPS to train a classifier

① Data exploration (analyse features, categoricals, numericals)

② Preparing data

→ Skewed features → Logarithmic transformation
(reduce distance between max and min values)

→ Normalisation → MinMaxScaler make feature values between 0-1. make the features treated equally by supervised learners. (keep the distribution shape).

→ Treat non-numeric features. → pandas.get_dummies to transform categorical data into numerical data using one-hot encoding

→ Shuffle / split data → scikit train-test-split helps to split the data.

(3) Naïve Predictor Performance

→ check a Naïve Predictor performance → consider the model always predicts 1 calculates Accuracy and Fscore.
(No 0, so FP=0, FN=0)

(4) Using / selecting Supervised models.

→ create a training and prediction pipeline → generic function to receive a classifier and the train/test data. fit the model, predict and calculate accuracy, f-score, training time, inference time

(5) Analyse classifier results

→ understand the metrics for each classifier and test/train step.

→ If train accuracy is ↑ and test accuracy is ↓ = OVERFITTING

Check the best Fscore for the specific problem.

(Some need high recall, others high precision)

(6) Model tuning

→ Find best parameters for the model

→ Use grid search to try the model with different parameters to find the best.

Use fbeta score as grid search optimization target