

SUPERVISED LEARNING

BINARY CLASSIFICATION

$$\left\{ \begin{array}{l} x_1, \dots, x_n \\ y_1, \dots, y_n \end{array} \right\} \quad \begin{array}{l} x_i \in \mathbb{R}^d \\ y_i \in \{0, 1\} / \{-1, +1\} \end{array}$$

Goal:

$$h: \mathbb{R}^d \rightarrow \{0, 1\}$$

LOSS/ERROR

$$\text{Loss}(h) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}(h(x_i) \neq y_i)$$

0-1-loss

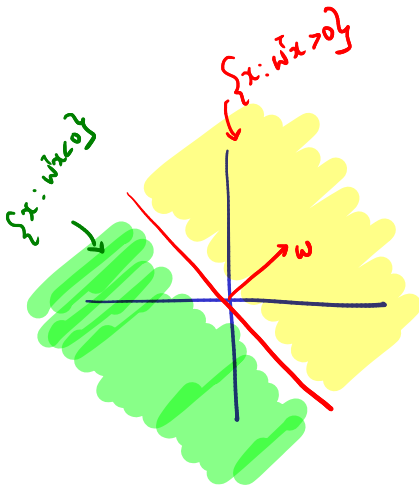
$$\mathbb{1}(z) = \begin{cases} 1 & \text{if } z \text{ is true} \\ 0 & \text{o/w} \end{cases}$$

$$\min_{h \in \mathcal{H}_{\text{linear}}} \sum_{i=1}^n \mathbb{1}(h(x_i) \neq y_i)$$

NP-hard problem in general

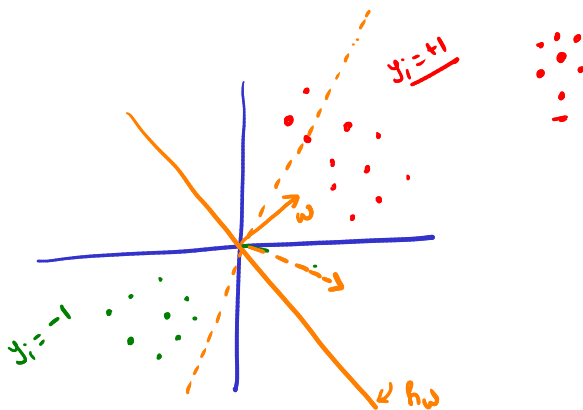
$$\mathcal{H}_{\text{linear}} = \left\{ h_w: \begin{array}{l} h_w(x) = \text{Sign}(w^T x) \end{array} \right\}$$

$$\text{Sign}(z) = \begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{o/w} \end{cases}$$



- Can we use linear regression to solve classification problem?

$$\{(x_1, y_1), \dots, (x_n, y_n)\} \rightarrow \boxed{\text{Lin Reg}} \rightarrow w \in \mathbb{R}^d \rightarrow h_w: \mathbb{R}^d \rightarrow \{0, 1\}$$



Conclusion

Regression is SENSITIVE TO location of the data points and not just the "Side" on which the data lies wrt separator.

SIMPLE ALGORITHMS FOR CLASSIFICATION

- Given a test point $x_{\text{test}} \in \mathbb{R}^d$,
find the closest point x^* to x_{test} in the training set
- Predict $y_{\text{test}} = y^*$

ISSUE: can get affected by outliers

Fix: Ask more neighbours

K-NN (k-Nearest Neighbours)

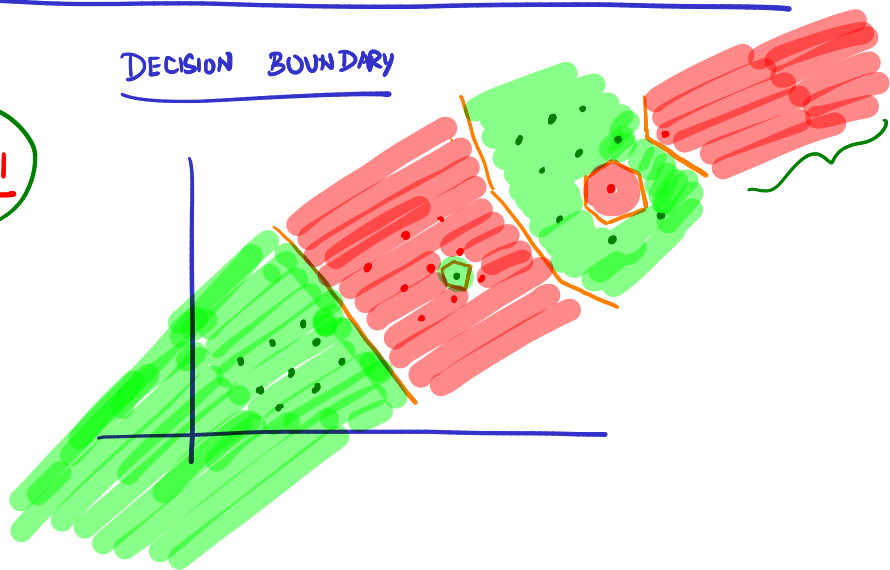
- Given x_{test} , find the k-closest points in the training set - $(x_1^*, x_2^*, \dots, x_k^*)$

- PREDICT $y_{\text{test}} = \text{majority}(y_1^*, y_2^*, \dots, y_k^*)$

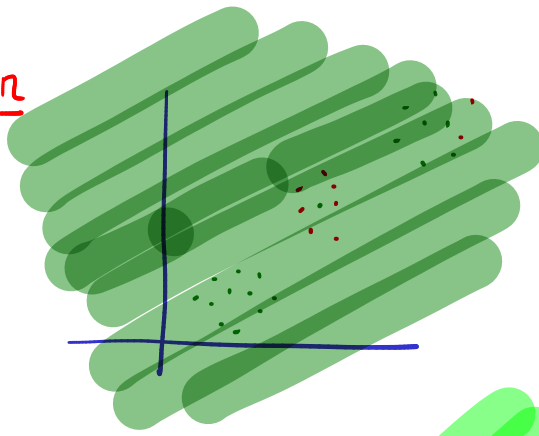
SENSITIVE TO
OUTLIERS!

K=1

DECISION BOUNDARY

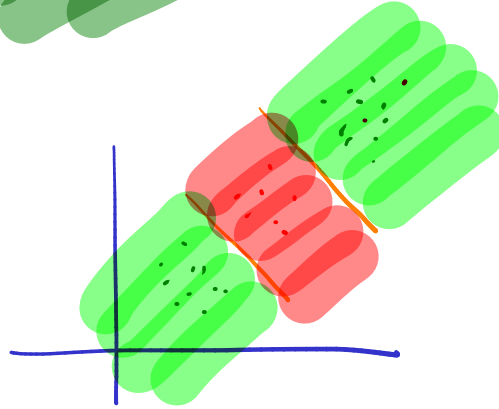


$k=n$



Every point is predicted $y = -1$ (green).

Too Smooth!



k^*

A good choice.

Choosing k

- Can treat as a HYPER PARAMETER
- Smaller the k , complicated the decision boundary
- Soln: CROSS-VALIDATE for k

ISSUES with K-NN

- Choosing a distance function
- PREDICTION is COMPUTATIONALLY EXPENSIVE.
- NO MODEL IS LEARNED.
 - Cannot throw away data after "LEARNING"

DECISION TREES

INPUT:

Dataset $\{ (x_1, y_1), \dots, (x_n, y_n) \}$

$x_i \in \mathbb{R}^d$
 $y_i \in \{+1, 0\}$

OUTPUT:

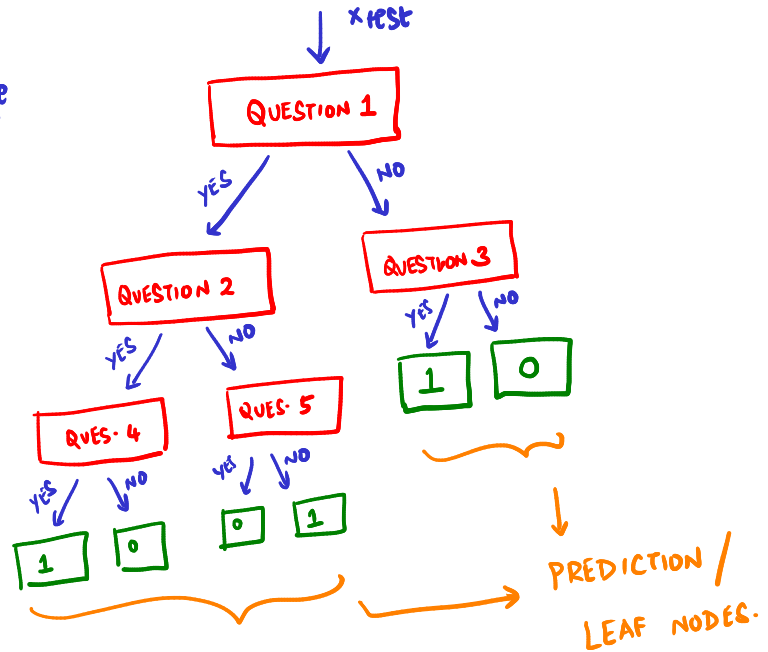
DECISION TREE

DECISION tree

PREDICTION : Given x_{test} ,

traverse through the tree to reach a leaf node.

y_{test} = value in leaf node.



QUESTION :

A question is a (feature, value) pair.

Eg: $\text{height} \leq 180\text{cm} ?$
(f_3) θ

How to measure "goodness" of a question?

DATASET
 $D = \{ (x_1, y_1), \dots, (x_n, y_n) \}$

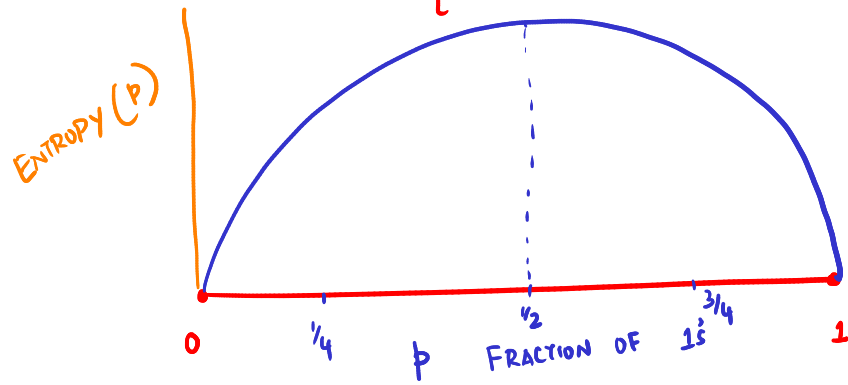
$f_k \leq \theta ?$

YES / NO

$D_{\text{yes}} = \{ (x_1, y_1), (x_{10}, y_{10}), \dots \}$

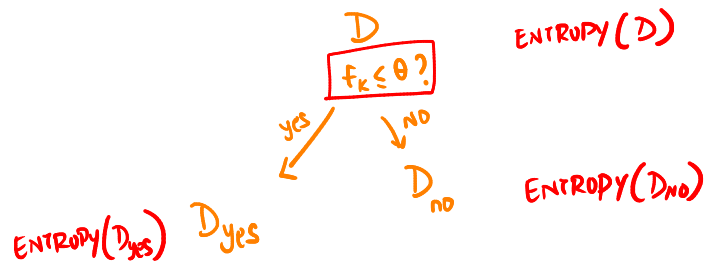
$\{ (x_2, y_2), (x_3, y_3), \dots \} D_{\text{no}}$

- Need is a measure of "Impurity" for a set of labels $\{y_1, \dots, y_n\}$



$$\begin{aligned} \text{ENTROPY}(\{y_1, \dots, y_n\}) &= \text{ENTROPY}(p) \\ &= -(p \log p + (1-p) \log (1-p)) \end{aligned}$$

[convention $\log(0)=0$]



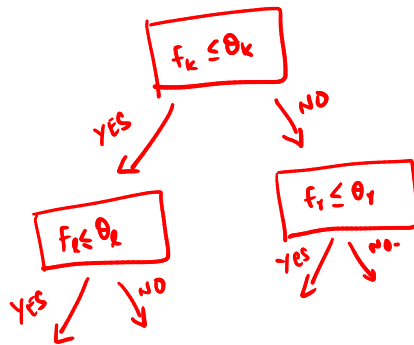
$$\text{INFORMATION GAIN}(\text{feature, value}) =$$

$$\text{ENTROPY}(D) - \left[\gamma \text{ENTROPY}(D_{\text{yes}}) + (1-\gamma) \text{ENTROPY}(D_{\text{no}}) \right]$$

$$\gamma = \frac{|D_{\text{yes}}|}{|D|}$$

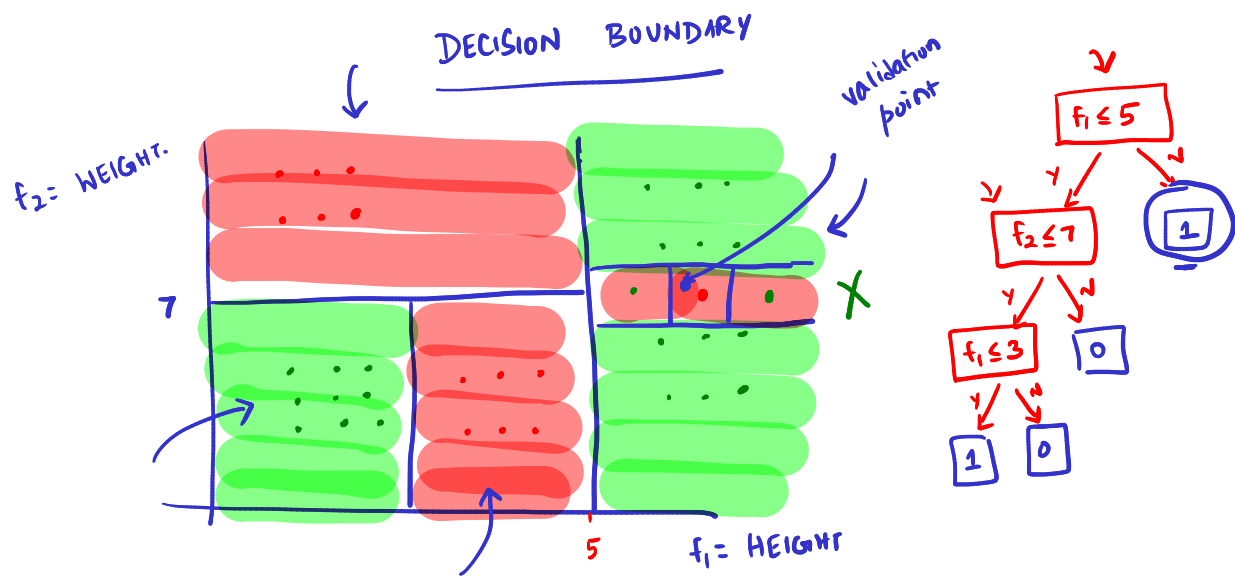
ALGORITHM - DECISION TREE

- DISCRETIZE each feature in $[\min, \max]$ range
 - Pick the Question that has the largest Information gain.
 - Repeat the procedure for D_{yes}, D_{no}
-

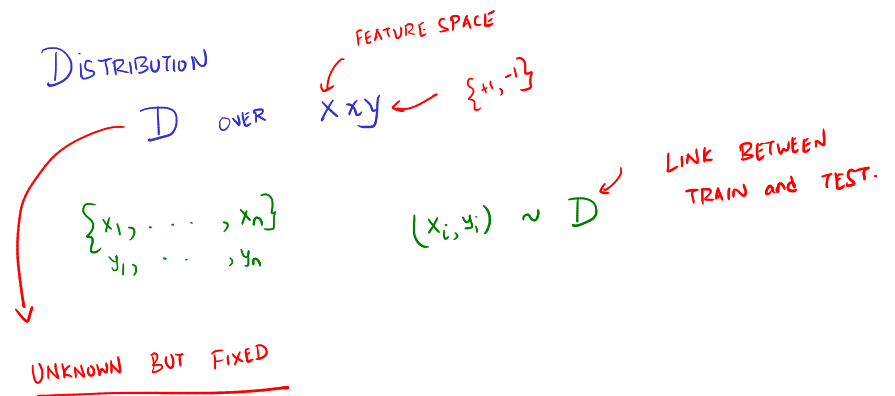


POINTS

- Can stop growing a tree if a node becomes "SUFFICIENTLY" pure.
 - DEPTH of the tree is a hyperparameter
 - There are alternate measures for "goodness" of a question
→ GINI INDEX
-



TYPES OF MODELING



CLASSIFICATION

- GENERATIVE MODEL
- DISCRIMINATIVE MODEL

GENERATIVE MODEL

- MODEL
- $P(x, y)$

↳ NEXT.

DISCRIMINATIVE MODEL

MODEL

- $P(y|x)$

Eg: K-NN
DECISION-TREES

↳ $P(y=1|x) = 1$ if
decision tree for x
says 1.

$P(y=1|x) = 1$ if majority of
neighbours say 1
= 0 otherwise