

(1)

# ML Week 2 - Supervised I

- (1) Linear Regression
  - (2) Decision trees
  - (3) k-nearest Neighbours
- } this week

## Decision trees

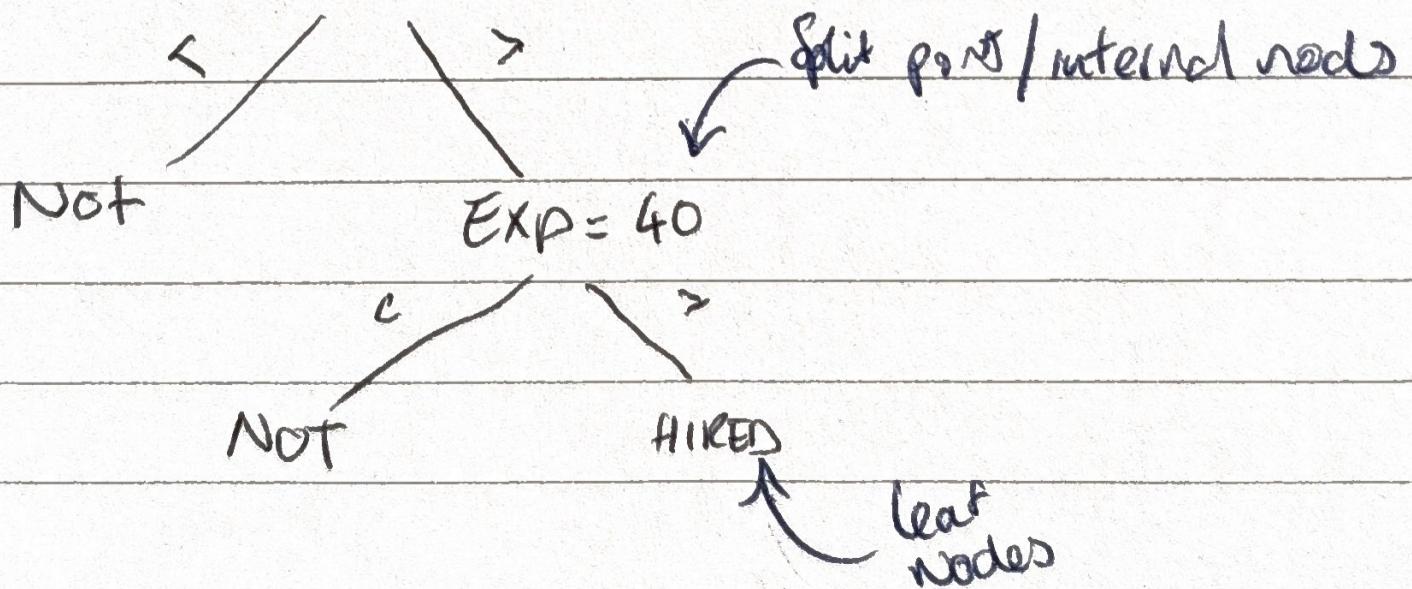
Recall, the goal of an ML is to find a function  $f(\cdot)$  that takes in  $\mathcal{X}_n$  & outputs  $y_n$ :

- $f(X) = \hat{y} = y$  TE set
- $f(x_m) = \hat{y}_m = y_m$  where  $x_m \in \mathcal{X}_n$

example:

feature: Exp, ML grade      output: Hired

$$ML = 50$$



(2)

- DTs work by partitioning the data recursively - binary tree structure
- Split point defined by a hyperplane that separates partitions
- leaf nodes = labels for data instances that are the same
- used for classif & regression

## Purity

- Purity helps decide when to split or not

$$\text{Purity}_m = \max_{C_k} \frac{n_{mk}}{n_m} \quad \text{where } n_m = \text{num data points in the partition region } m$$

$n_{mk}$  = num data points that belong to label class  $C_k$  in the partition region  $m$

$n_{mk}$  = num data points that belong to label class  $C_k$  in the partition region  $m$

Max purity = 1

e.g. ML = 51  $C_1 = \text{Not}$   $C_2 = \text{Hired}$

$$LSI = 3 \quad C_1 = \frac{3}{3} \quad C_2 = \frac{0}{3} \quad \max_{C_k} \left\{ \frac{3}{3} C_1, \frac{0}{3} C_2 \right\}$$

$$= \frac{3}{3} = 1 \text{ purity, not need to split further}$$

A purity of 1 means the whole cohort of a category/Split was partitioned

Others, > S1, make  $\sum_k \frac{2}{3} p_1, \frac{3}{3} p_2 = 0.67$   
less than 1 so split further

How to determine the split point?

- optimal point determined by a split criterion
- two types: info gain & gini index

Information Gain: Entropy

- Entropy = measure of uncertainty
- First compute entropy then info gain
- Choose split w/ highest info gain

Gini Index

- Compute formulae
- Compute Gini for both sides of split
- Lowest index is chosen

(4)

Example: best split point

- Order the data w/ a parameter
- Calc split potential point inbetween all existing points - min points instances &
- Count split outcomes  $O = 1 = n_{m_{<}} , n_{m_{>}}$
- Plug into Gini or Info formula

less than side of split      more than

Gini Example

$$G_{m_{<}} / m_{>} = \frac{n_{m_{<}}}{n_m} \left( 1 - \sum_{k=1}^K \left( \frac{n_{m_{k<}}}{n_{m_{<}}} \right)^2 \right) + \frac{n_{m_{>}}}{n_m} \left( 1 - \sum_{k=1}^K \left( \frac{n_{m_{k>}}}{n_{m_{>}}} \right)^2 \right)$$

27.5 Split

$$\frac{1}{8} \left( 1 - \left( \left( \frac{0^2}{2} + \frac{1^2}{2} \right) \right) \right) + \frac{7}{8} \left( 1 - \left( \left( \frac{3^2}{7} + \frac{4^2}{7} \right) \right) \right) = 0.43$$

- 8 = Total , 1, 7 = Splits , inner = ~~split total~~ ~~split total~~ & ~~total~~  $\times n_{m_{<}}$

- Calc for all mid points & select the best metric

(5)

## Categorical split points

- Each ~~point has~~ value / category is a potential split point
- Otherwise same process

## Building An Optimised Model

- Purity threshold
- Max num of leaf nodes
- Initial region
- Purity of region

}

parameters

region = the  
Training data  
itself

if purity  $\geq$  threshold

Stop splitting

ELSE

- for each feature find best split point
- Select best feature
- Split using feature & split
- Repeat for each new region created

(6)

## Overfitting w/ tree models:

- = large tree with many nodes
- address overfit by containing nodes

## Pruning

- reduce overfit by removing nodes
- Prune criterion used to determine if given node should be pruned

$$V_T = \sum_{T_i} E_T + \alpha |T|$$

$|T|$  = total leaf nodes in Tree

$\alpha$  = Prune hyperpara

$E_T$  = error measure for  $T$  (gini or entropy)

$\epsilon$  = contribution impact

$|T|$  = size factor

(7)

## Ensemble learning

- Bagging - Bootstrap aggregating  
 Bootstrap = each model trained  
~~parametrized~~ on random subset  
 Aggregate = average pred of  
 Sub models

Pros - Capture complexity of data,  
 average out errors,  
 Robust to overfitting as sub  
 models don't see full data

## Boosting

- Cumulative training, multiple  
 Models built from a base model
- Data weighting - worse predicted  
 data is weighted for closer pred  
 on next model
- Prediction weighting - pred of  
 Community model weighted over  
 the multiple

# K-Nearest Neighbors

With KNN you use the classes of the K nearest neighbours of  $x_i$  to determine its output  $\hat{y}_i$

KNN does not produce a model  $f(\cdot)$

Step 1: Decide distance metric

~~Euclidean distance~~ for each value of  $x_i$  find its K nearest neighbour based on this metric

Cosine Similarity

Manhattan Distance  
2 Dimension measure

Step 2 Determine  $\hat{y}_i$  of  $x_i$  based on the labels of neighbours  
Voting/average system

- for class = majority regression = average  
 $k \rightarrow$  need to choose num of neighbors

$$\hat{y}_i = \frac{1}{k} \sum y_j$$

- can distance-weight to prioritize closer points

(9)

## time complexity for kNN

$$O(N^2 D) \quad N = \text{data instances}$$

$$D = \text{total number features}$$

kNN is computationally heavy /  
many points, lots of calculations

each data point M is compared  
to each of N data points &  
each data point has D  
number of dimensions to be compared

$$\text{Dist}(x_i, x_j) = \sqrt{\sum_{d=1}^D (x_{id} - x_{jd})^2}, \forall i, i=1, 2, \dots, N, i \neq j, \forall j, j=1, 2, 3, \dots, M$$

• Euclidean

## Minimizing Complexity

- { indexing
  - kd tree - indexing tree
  - recursive partition the data region
- Similar to decision tree

- Reduce Dimensionality of week 6