

18.3 knn

(1)

Zaki - Data Mining & ML, 2020

Knn is a non-parametric approach

This means it does not make any assumptions about the underlying distribution of the data

A parametric method assumes that the data follows a specific prob dist, i.e. Normal

These methods try to estimate the parameters of that distribution namely i.e. mean & standard deviation

Non-parametric does not do this

Uses the data itself to make predictions without trying to fit a predefined model

the knn works by locating the nearest neighbours of known/training data compared to an unseen data point

Then uses those labels/data to create pred

There is ~~not~~ no equation or function that describes new data

If you lose lose the data, you lose the ability to pred

Where as lin Reg can be performed w/  
just the equation & new data

This is also why knn is heavy on the cost to compute

The lack of assumption along fitting w/ complex data but exposes risk to overfitting

3

$D$  = train set

$n$  = points ,  $x_i \in \mathbb{R}^d$

$D_i$  = subset of  $D$  w/ class  $C_i$  w/  $n_i = |D_i|$

test point =  $x \in \mathbb{R}^d$

$K$  = number of neighbours

$r$  = dist of  $x$  test point to  $k^{th}$  neighbour in  $D$

$$B_d(x, r) = \{x_i \in D \mid \|x - x_i\| \leq r\}$$

$\downarrow$   
means  
Ball

↳ eq for a ball radius around  $x$ .

$$|B_d(x, r)| = K$$

$K_i$  = num of neighbours/points among the  $k$  nearest neighbours of  $x$  (test point) that are labelled w/ class  $C_i$

$$K_i = \{x_j \in B_d(x, r) \mid y_j = C_i\}$$

| = such that or where

4

The classes conditional probability & density @ point  $x$  can be estimated as the fraction of points from class  $C_i$  that lie within the ball/hyperball Divided by its volume

$$\hat{f}(x|C_i) = \frac{K_i/n_i}{V} = \frac{K_i}{n_i * V}$$

$K_i$  = Points of class,  $n_i$  = num of points in class  $i$ ,  $V$  = Volume of ball in ball

using posterior Probability  $P(C_i|x)$  eq:

$$P(C_i|x) = \frac{\hat{f}(x|C_i) \hat{P}(C_i)}{\sum \hat{f}(x|C_j) \hat{P}(C_j)}$$

This an application of Bayes theorem in the context of classification

Posterior Prob = Prob that point  $x$  belongs to class  $C$  given observed features of  $x$

(5)

$$\text{General Bayes} = P(A|B) = \frac{P(B|A) * P(A)}{P(B)}$$

Here, A = event that X belongs to  $c_i$

$$P(A|B) = P(c_i|x) = \text{Posterior Prob}$$

Prob that x belongs to  $c_i$  given feature of x

$$\hat{f}(x | c_i) = \text{Class conditional prob dens func}$$

Prob density of observing x feature given

we know it belongs to class  $x$ ,

i.e. how likely are these features for class  $c_i$

$$P(c_i) = \text{Prior prob of just being class } c_i \text{ in}$$

the training data

$\sum \hat{f}(x | c_j) \hat{P}(c_j)$  = This normalization or evidence term. Sums to the product of the likelihood for all classes. Ensures the posterior probs sum to 1 to form a valid prob distribution

6

In English, the equation calculates the posterior probability for class  $c_i$  by:

- (1) multiply the prob of observing  $x$  given class  $c_i$  by prior (overall) prob of class  $c_i$  in training data
- (2) Normalize by sum of similar products for all classes

However  $\hat{P}(c_i) = \frac{n_i}{n}$ , so:

$$\hat{s}(x | c_i) \hat{P}(c_i) \stackrel{(TOP)}{=} \frac{k_i}{n_i V} \cdot \frac{n_i}{n} = \frac{k_i}{n V}$$

thus posterior probability:

$$P(c_i | x) = \frac{\frac{k_i}{n V}}{\sum_{j=1}^k \frac{k_j}{n V}} = \frac{k_i}{K}$$

$k_i$  = neighbours of Class  $c_i$ ,  $k$  = Total neighbors

7

finally, the predicted class for  $x$ :

$$\hat{y} = \operatorname{argmax}_{c_i} \left\{ P(c_i | x) \right\} = \operatorname{argmax}_{c_i} \left\{ \frac{k_i}{K} \right\} = \operatorname{argmax}_{c_i} \left\{ k_i \right\}$$

$\{k_i\}/K$  represents the prop. of  $k$  nearest neighbours belong to  $c_i$

$\operatorname{argmax} \{k_i\}$  means the pred  $\hat{y}$  is whatever class  $c_i$  has the highest count