

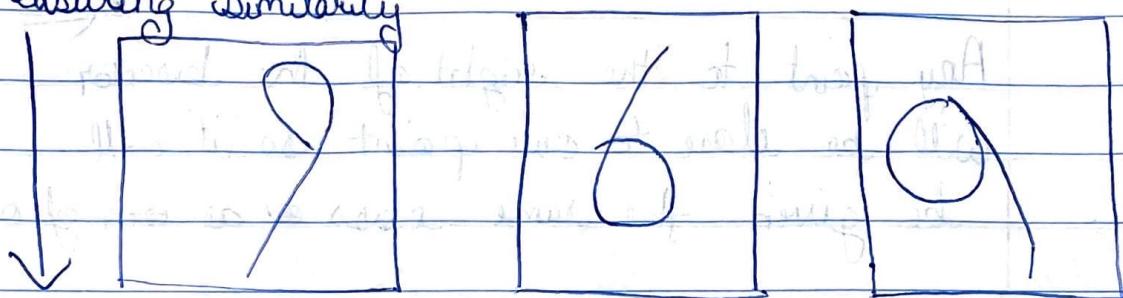
Lecture - 16

vRPnuc

Topic: Similarity

* Similarity

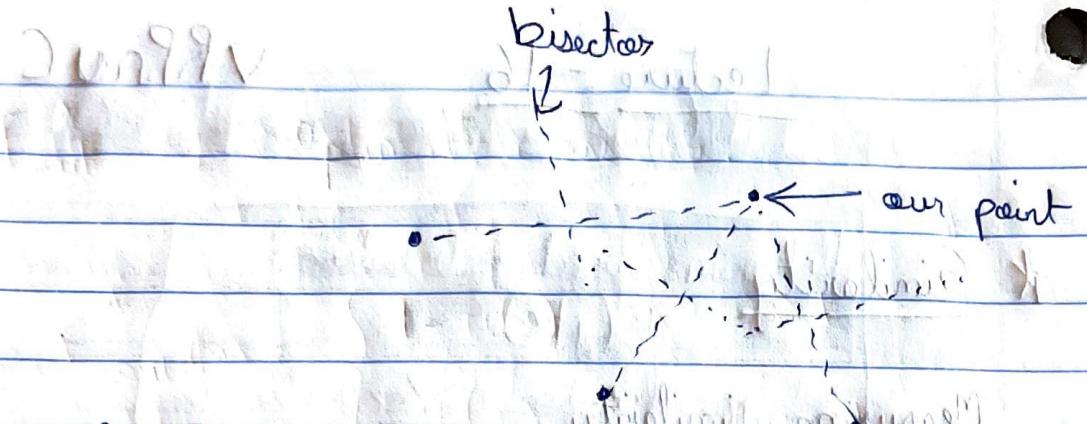
Measuring Similarity



Compute features.

* The k nearest neighbours Algorithm (KNN).

Given a test point x ;classify using k of its nearest neighbour. $x_{[1]}$ → Closest to our test point. $x_{[2]}$ → A bit far $x_{[N]}$ → too far from our test pointwhere $d(x_{[1]}, x) \leq d(x_{[2]}, x) \dots d(x_{[N]}, x)$.Output = $\text{sign}(y_{[1]})$.



Any point to the right of the bisector will be closer to our point so it will be given the same class as ~~as~~ of our point.

Any point to the left will be closer to that point.

Now we ~~can't~~ ^{don't} need to make tessellation (region of the data points)

Just calculate the distance, and find the nearest data point. Assigning them +ve or -ve depending on the distance.

* KNN algorithm ~~it's~~ will give you 0 in sample error. Just by considering itself as the closest data point.

KNN will always have 0 in sample error ($E_{in} = 0$)

though we cannot calculate E_{in} & relation.
We can calculate E_{out} .

$$E_{out(g)} \leq 2 E_{out}^* \rightarrow \text{Noisy target func'}$$

\rightarrow KNN algo.

$$E_{out(g)} \leq 2 E_{out}^*$$

as $N \rightarrow \infty$ (for high data sample),
the $E_{out(g)} \approx 2 E_{out}^*$.

i.e. $E_{out(g)}$ will be small.

Noisy

$$\pi(x) = P[y = +1 | x]$$

(given x what is the probability of
to give $y = +1$)

Unknown Target.

$$1 - \pi(x) = P[y = -1 | x].$$

* Optimal classifier:

$$\pi(x) \geq \frac{1}{2} \rightarrow +1$$

$$\pi(x) < \frac{1}{2} \rightarrow -1.$$

E_{out}^*

(Optimal classifier)

$$E_{\text{out}}^*(x) = \begin{cases} 1 - \pi(x) & \pi(x) \geq \frac{1}{2} \rightarrow -1 \\ \pi(x) & \pi(x) < \frac{1}{2} \rightarrow +1. \end{cases}$$

$$= \min \{(\pi(x), 1 - \pi(x))\}$$

$$= \pi(x)$$

$$E_{\text{out}}^* = \int dx P(x) E_{\text{out}}^*(x).$$

optimal

* A few assumptions:

* [Property 1]: $\pi(x)$ is continuous.

[Property 2]: for all $N \rightarrow \infty$

$$x_{[i]} \rightarrow x.$$

$E_{\text{out}}(x) = \text{prob. of error}(x).$

$$= P[Y_{[i]} \neq y]$$

Nearest
label

True
label.

$$E_{out}(n) = P\left[\underbrace{y_{E,I} = +1, y = -1}_{\pi(y_{E,I})}\right] + P\left[y_{E,I} = -1, y = +1\right]$$

$$= \pi(n_{E,I}) (1 - \pi(n)) + (1 - \pi(n_{E,I})) (\pi(n))$$

* As $N \rightarrow \infty$

As $n_{E,I}$ converges to n [Continuous]
 $\pi(n_{E,I})$ also converges to $\pi(n)$.

~~$E_{out}(n) \approx 2\pi(n)(1-\pi(n))$~~

$$E_{out}(n) = 2\pi(n)(1 - \pi(n)) \\ \leq 2 \min \{ \pi(n), 1 - \pi(n) \}$$

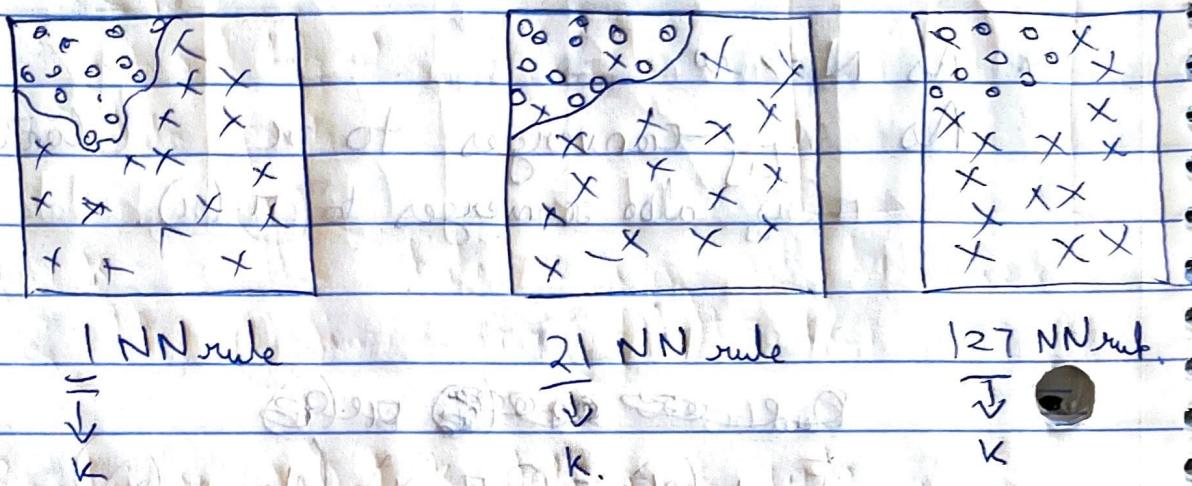
$$E_{out} \approx 2 E^*$$

So if E^* out is small, E_{out} will be small and we will have tight bound.

If E^* out is large, E_{out} will be large and we will have loose bound.

with $N \rightarrow \infty$ the model will become complex.
 We can tackle this by increasing $K \uparrow$

But if K is very ~~too~~ large we will underfit.



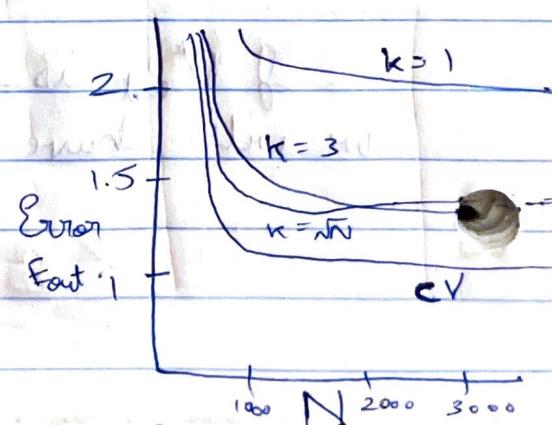
Will eliminate outliers.

As $N \uparrow$, try to $\uparrow K$

Rule of thumb $K = \lceil \sqrt{N} \rceil$

Choosing K .

$CV \Rightarrow$ Cross Validation.



NN

- Nb parameters
- Expressive, flexible
- $g(x)$ defined by the data.
- generic, can model any target

Linear

- $d+1$ parameters.
- Rigid, always linear
- $g(x)$ defined by weights.
- specialized.

* KNN :- simple

No training

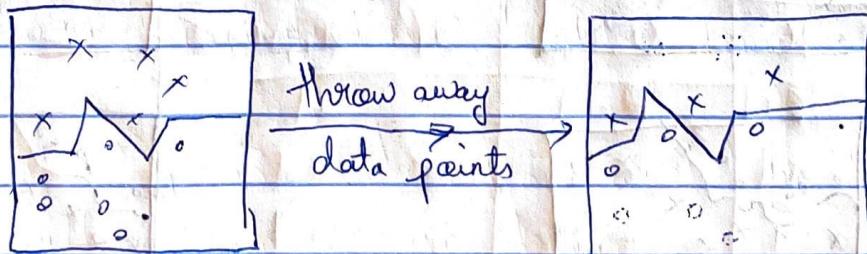
Near Optimal out of sample error

Can be used for regression or logistic regression

Problem:- computationally demanding

Do it:- Condensed Nearest Neighbour

↳ choosing a subset of the dataset



but still maintaining the boundary