

KNN

K - Nearest Neighbours

① Proximity:

Two observations: $\begin{array}{cc} x_1 & x_2 \\ \hline (x_{11}, x_{21}) & \\ (x_{12}, x_{22}) & \end{array}$

Euclidean distance (in vectorized notation)

let $\underline{x}_1 = (x_{11}, x_{21})^T$ and $\underline{x}_2 = (x_{12}, x_{22})^T$

be two vectors.

Euclidean distance (in general)

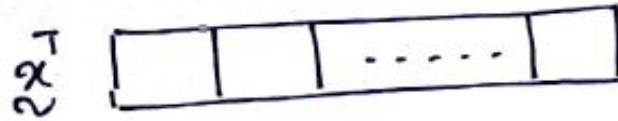
	x_1	x_2	...	x_p
\tilde{x}_1^T	x_{11}	x_{21}	...	x_{p1}
\tilde{x}_2^T	x_{12}	x_{22}	...	x_{p2}
...

The euclidean distance between \tilde{x}_1 and \tilde{x}_2 can be written as:

$$\text{dist} = \sqrt{(x_{11} - x_{12})^2 + (x_{21} - x_{22})^2 + \dots + (x_{p1} - x_{p2})^2}$$

$$= \sqrt{(\tilde{x}_1 - \tilde{x}_2)^T (\tilde{x}_1 - \tilde{x}_2)}$$

② Finding the nearest neighbour :



	x_1	x_2	...	x_p
x_1^T				
x_2^T				
x_3^T				
\vdots		\vdots	\cdot	\vdots
x_{i-1}^T				
x_i^T				
x_{i+1}^T				

1-NN

Step 1: calculate the dist. of \underline{x} from \underline{x}_i , for all $i = 1(1)n$

Step 2: $NN =$ the obs. that corresponds to the shortest distance.

Simplest approach (1-NN):

1: min-dist = Inf

2: for i in ~~range(1, len(train))~~ From 1 to len(train)

3: d = dist. of the test obs. from the i th obs. in train-data

4: if $d < \text{min-dist}$:

5: min-dist = d #update min-dist

6: nearest-obs-index = i

7: end if

8: end for

9: nearest_obs = train [i ,]

③ 1-NN algorithm for classification :

Training data :

	x_1	x_2	...	x_p	y
x_1^T	x_{11}	x_{21}	...	x_{p1}	y_1
x_2^T	x_{12}	x_{22}	...	x_{p2}	y_2
\vdots	\vdots	\vdots		\vdots	\vdots
$x_{n_1}^T$	x_{1n_1}	x_{2n_1}	...	x_{pn_1}	y_{n_1}

Test data :

	x_1	x_2	...	x_p
x_1^{*T}	x_{11}^*	x_{21}^*	...	x_{p1}^*
x_2^{*T}	x_{12}^*	x_{22}^*	...	x_{p2}^*
\vdots	\vdots	\vdots		\vdots
$x_{n_2}^{*T}$	$x_{1n_2}^*$	$x_{2n_2}^*$...	$x_{pn_2}^*$

n_1 = no. of obs. in training data

n_2 = no. of obs. in test data

□ function 1: function to calculate the euclidean distance.

$\text{distance}(\underline{x}, \underline{y}) \rightarrow$ returns euclidean dist. between the vectors \underline{x} and \underline{y} .

□ function 2: function to calculate the (or get the) NN for one test observation.

$\text{NN}(\text{train_x}, \text{train_y}, x^*) \rightarrow$ finds the nearest obs. to x^* in the training data and get the value of y for that obs.

□ function 3:

$1NN(\text{train-X}, \text{test}, \text{train y}) \rightarrow$ returns the values of y for all the obs. in the test data corresp. to the nearest neighbors in the train data.

(4)

