

Machine Learning In Action — Chapter 2

Example 2.1— Explaining kNN

In this example, the parameters are

$$\begin{aligned} k &= 3, & in_X &= \begin{pmatrix} 0.6 & 0.6 \end{pmatrix}, \\ data_set &= \begin{pmatrix} 1.0 & 1.1 \\ 1.0 & 1.0 \\ 0.0 & 0.0 \\ 0.0 & 0.1 \end{pmatrix}, & labels &= \begin{pmatrix} A \\ A \\ B \\ B \end{pmatrix} \end{aligned}$$

Using the above arguments,

$$\begin{aligned} diff_matrix &= \mathbf{tile}(in_X, (4, 1)) - data_set \\ &= \begin{pmatrix} 0.6 & 0.6 \\ 0.6 & 0.6 \\ 0.6 & 0.6 \\ 0.6 & 0.6 \end{pmatrix} - \begin{pmatrix} 1.0 & 1.1 \\ 1.0 & 1.0 \\ 0.0 & 0.0 \\ 0.0 & 0.1 \end{pmatrix} \\ &= \begin{pmatrix} -0.4 & -0.5 \\ -0.4 & -0.4 \\ 0.6 & 0.6 \\ 0.6 & 0.5 \end{pmatrix} \end{aligned}$$

and *diff_matrix* represent the straight line distance the vector is from the other data points in the set. To calculate the Euclidean Distance, take the square of every element in *diff_matrix*, sum them across the rows, then take the square root.

$$sq_diff_matrix = \begin{pmatrix} 0.16 & 0.25 \\ 0.16 & 0.16 \\ 0.36 & 0.36 \\ 0.36 & 0.25 \end{pmatrix}, sq_distances = \begin{pmatrix} 0.41 \\ 0.32 \\ 0.72 \\ 0.61 \end{pmatrix}$$

sq_distances represents the sum of all the distances from the vector is from the data point. Finally, taking square root of the distances and then sorting them,

$$distances = \begin{pmatrix} 0.64031 \\ 0.56568 \\ 0.84852 \\ 0.78102 \end{pmatrix}, sorted_dist_indices = \begin{pmatrix} 1 \\ 0 \\ 3 \\ 2 \end{pmatrix}$$

Taking the first 3 labels, we can see that there are 2 data points labelled A and 1 data point labelled B . Hence, by majority, the data point in_X is classified in group A .