

The  $k$ NN, or  $k$ -nearest neighbors algorithm, is a simple supervised learning algorithm, which works under the assumption that samples with similar features are more likely to share the same label. The algorithm is very straight-forward, although differing between trying to predict categoric and numeric labels: using a given distance metric (such as Euclidean, Manhattan, etc.), the algorithm finds the  $k$  closest samples to the sample we want to classify, and assigns the label that:

1. is most common (i.e. the mode) among the  $k$  neighbors, in the case of categoric labels;
2. is the average of the labels of the  $k$  neighbors, in the case of numeric labels.

1. Considering the following data set:

	$y_1$	$y_2$	$z_1$	$z_2$
$x_1$	1	1	A	1.4
$x_2$	2	1	B	0.5
$x_3$	2	3	B	2
$x_4$	3	3	B	2.2
$x_5$	2	2	A	0.7
$x_6$	1	2	A	1.2

Assuming  $k$ NN, with  $k = 3$  applied within a leave-one-out schema:

- (a) Considering an  $z_1$  categoric output variable and the Euclidean distance, provide the prediction for  $x_1$ .
- (b) Considering an  $z_2$  numeric output variable and the cosine similarities, provide the mean regression estimate for  $x_1$ .
- (c) Considering a weighted-distance  $k$ NN, with Manhattan distance, identify both the **weighted-mode** estimate of  $x_1$  for a  $z_1$  outcome and the **weighted-mean** estimate of  $x_1$  for a  $z_2$  outcome.

Here, working with a leave-one-out schema means that, for each sample, we can pick its  $k$  neighbors from the pool of all the other samples, excluding itself (which, in theory, is always its closest neighbor, with null distance).

- (a) The Euclidean distance is defined as the square root of the sum of the squared differences between the features of the two samples:

$$d(x_i, x_j) = \sqrt{\sum_{l=1}^p (x_i^{(l)} - x_j^{(l)})^2}$$

Since we're working with a leave-one-out schema, and trying to estimate the  $z_1$  label for  $x_1$ , we can pick the  $k$  neighbors from all samples except  $x_1$ . Below are illustrated the Euclidean distances between those samples and  $x_1$ , with the  $k$  closest neighbors highlighted in teal:

$$d(x_1, x_2) = \sqrt{(1 - 2)^2 + (1 - 1)^2} = 1$$

$$d(x_1, x_3) = \sqrt{5}, \quad d(x_1, x_4) = 2\sqrt{2}, \quad d(x_1, x_5) = \sqrt{2}, \quad d(x_1, x_6) = 1$$

Knowing the  $k$  closest neighbors, we can now estimate the  $z_1$  label for  $x_1$ , by picking the most common label among them. In this case, the estimated label will be **mode**( $B, A, A$ ) =  $A$ .

- (b) Here, instead of the usual Euclidean/Manhattan distance, we're using the cosine similarities, which are defined as the cosine of the angle between the two vectors of features. As we know since high school, the cosine of the angle between two vectors is equal to the dot product of the two vectors divided by the product of their norms:

$$\cos(\vec{a}, \vec{b}) = \frac{\vec{a} \cdot \vec{b}}{\|\vec{a}\| \|\vec{b}\|}$$

Here, each sample is essentially a vector of features, hence the usage of this distance metric. The higher the cosine similarity, the closer the two vectors are, and vice versa: the cosine is equal to 1 when the two vectors are identical, and 0 when they are orthogonal.

Let's now compute the cosine similarities between  $x_1$  and the other samples, picking the  $k$  closest neighbors (once again, in teal):

$$\cos(x_1, x_2) = \frac{1 \cdot 2 + 1 \cdot 1}{\sqrt{2}\sqrt{5}} = \frac{3}{\sqrt{10}} = 0.94868, \quad \cos(x_1, x_3) = \frac{1 \cdot 2 + 1 \cdot 3}{\sqrt{2}\sqrt{13}} = \frac{5}{\sqrt{26}} = 0.98058$$

$$\cos(x_1, x_4) = \frac{1 \cdot 3 + 1 \cdot 3}{\sqrt{2}\sqrt{18}} = \frac{6}{\sqrt{36}} = 1, \quad \cos(x_1, x_5) = \frac{1 \cdot 2 + 1 \cdot 2}{\sqrt{2}\sqrt{8}} = \frac{4}{\sqrt{16}} = 1$$

$$\cos(x_1, x_6) = \frac{1 \cdot 1 + 1 \cdot 2}{\sqrt{2}\sqrt{5}} = \frac{3}{\sqrt{10}} = 0.94868$$

Note, as mentioned above, that the closest neighbors here are the ones with a higher cosine similarity:  $x_3$ ,  $x_4$ , and  $x_5$ . Since we're working with numeric labels for  $z_2$ , we can now estimate the mean regression estimate for  $x_1$  by averaging the values of  $z_2$  for those samples. In this case, the estimated value will be **mean**(2, 2.2, 0.7) = 1.6(3).

- (c) Note that, here, we're working with a weighted-distance  $k$ NN, which means that we're assigning different weights to each sample based on its distance from the sample we're trying to estimate - a closer neighbor will have a bigger impact on the final estimate than a further neighbor. Note, also, that we're now working with the Manhattan distance, which is defined as the sum of the absolute differences between the features of the two samples:

$$d(x_i, x_j) = \sum_{l=1}^p |x_i^{(l)} - x_j^{(l)}|$$

Let's now compute the Manhattan distances between  $x_1$  and the other samples, picking the  $k$  closest neighbors (once again, in teal):

$$d(x_1, x_2) = |1 - 2| + |1 - 1| = 1$$

$$d(x_1, x_3) = 3, \quad d(x_1, x_4) = 4, \quad d(x_1, x_5) = 2, \quad d(x_1, x_6) = 1$$

Now, the fun part begins: we have to correctly weigh each neighbor's label to estimate the label for  $x_1$ . We can do this by assigning a weight to each neighbor, based on its distance from  $x_1$ . The most common way to do this is by using the inverse of the distance. Regarding the weighted mode estimate of  $x_1$  for  $z_1$ :

$$\hat{z}_1 = \text{weighted\_mode}(\frac{1}{1}B, (\frac{1}{1} + \frac{1}{2})A) = A$$

The mean, of course, will take into account for the denominator the weights of each neighbor, instead of the amount of neighbors:

$$\hat{z}_2 = \text{weighted\_mean} = \frac{\frac{1}{1} \cdot 0.5 + \frac{1}{2} \cdot 0.7 + \frac{1}{1} \cdot 1.2}{\frac{1}{1} + \frac{1}{2} + \frac{1}{1}} = 0.82$$

2. Consider the following training data set:

	$y_1$	$y_2$	$z$
$x_1$	1	1	1.4
$x_2$	2	1	0.5
$x_3$	1	3	2
$x_4$	3	3	2.5

- Find the closed form solution for a linear regression, minimizing the sum of squared errors.
- Predict the target value for the query vector  $x_{new} = [2 \ 3]^T$ .
- Sketch the predicted three-dimensional hyperplane.
- Compute both the MSE and MAE produced by the linear regression.
- Are there biases on the residuals against any of the input variables?
- Compute the closed form solution, considering Ridge regularization term with  $\lambda = 0.2$ .
- Compare the hyperplanes obtained utilizing ordinary least squares and ridge regression.
- Why is the Lasso regression usually preferred over Ridge regression for data spaces with a larger number of features?

- 
- 
- 
- 
- 
-

(g)

(h)

3. Considering the following training data, with  $z$  as an ordinal variable:

	$y_1$	$y_2$	$z$
$x_1$	1	1	1
$x_2$	2	1	1
$x_3$	1	3	0
$x_4$	3	3	0

(a) Find a linear regression using the closed form solution.

(b) Assuming an output threshold  $\theta = 0.5$ , provide the predicted class for  $x_{new} = [2 \ 2.5]^T$ .

(a)

(b)

4. Considering the data below to learn the following model, compare:

$$z = w_1 y_1 + w_2 y_2 + \epsilon, \epsilon \sim \mathcal{N}(0, 0.1)$$

	$y_1$	$y_2$	$z$
$x_1$	3	-1	2
$x_2$	4	2	1
$x_3$	2	2	1

(a)  $w = [w_1 \ w_2]^T$ , using an MLE approach.

(b)  $w$  using the Bayesian approach, assuming  $p(w) = N(w \mid \mu = [0, 0], \Sigma = \begin{bmatrix} 0.2 & 0 \\ 0 & 0.2 \end{bmatrix})$ .

(a)

(b)

5. Identify a transformation to aid linearly modelling the data set above. Sketch the predicted surface.

6. Consider both logarithmic and quadratic transformations:

$$\phi_1(x_1) = \log(x_1), \quad \phi_2(x_2) = x_2^2$$

- (a) Plot both of the closed form regressions.
- (b) Which transformation minimizes the sum of squared errors on the original data?

(a)

(b)

7. Select the criteria promoting a smoother regression model:

- (a) Applying Ridge and Lasso regularizations to linear regression models.
- (b) Increasing the depth of a decision tree regressor.
- (c) Increasing the  $k$  parameter of a  $k$ NN regressor.
- (d) Parametrizing a  $k$ NN regressor with uniform weights, instead of the default distance-based weights.