

The Nearest Neighbour algorithm

We now introduce a concrete learning algorithm for classification. This algorithm differs from ERM because it is not minimizing the training error in a given class of predictors. For now, we restrict our attention to binary classification tasks with numerical features, namely $\mathcal{X} = \mathbb{R}^d$ and $\mathcal{Y} = \{-1, 1\}$. Given a training set, the classifier generated by this algorithm is based on the following simple rule: predict every point in the training set with its own label, and predict any other point with the label of the point in the training set which is closest to it.

More formally, given a training set S of examples $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$, the **nearest neighbour** algorithm (NN) generates a classifier $h_{\text{NN}} : \mathbb{R}^d \rightarrow \{-1, 1\}$ defined by:

$$h_{\text{NN}}(\mathbf{x}) = \text{label } y_t \text{ of the point } \mathbf{x}_t \in S \text{ closest to } \mathbf{x}.$$

If there is more than one point in S with smallest distance to \mathbf{x} , then the algorithm predicts with the majority of the labels of these closest points. If there is an equal number of closest points with positive and negative labels, then the algorithm predicts a default value in $\{-1, 1\}$ (for instance, the most frequent label in the training set).

Note that $h_{\text{NN}}(\mathbf{x}_t) = y_t$ for every training example (\mathbf{x}_t, y_t) . The distance between $\mathbf{x} = (x_1, \dots, x_d)$ and $\mathbf{x}_t = (x_{t,1}, \dots, x_{t,d})$, denoted by $\|\mathbf{x} - \mathbf{x}_t\|$, is computed using the Euclidean distance,

$$\|\mathbf{x} - \mathbf{x}_t\| = \sqrt{\sum_{i=1}^d (x_i - x_{t,i})^2}.$$

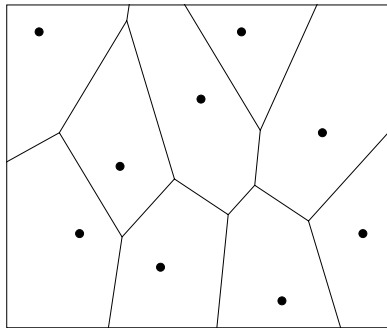


Figure 1: Voronoi diagram for a training set in \mathbb{R}^2 .

The classifier generated by NN induces a partition of \mathbb{R}^d in *Voronoi cells*, where each training instance \mathbf{x}_t (here called a “center”) is contained in a cell, and the border between two cells is the set of points in \mathbb{R}^d that have equal distance from the two cell centers (see Figure 1).

As NN typically stores the entire training set, the algorithm does not scale well with the number $|S| = m$ of training points. Moreover, given any test point \mathbf{x} , computing $h_{\text{NN}}(\mathbf{x})$ is costly, as it requires computing the distance between \mathbf{x} and every point of the training set, which in \mathbb{R}^d takes time $\Theta(dm)$ (shorter running times are possible when distances are approximated rather than being computed exactly). Finally, note that NN always generates a classifier h_{NN} such that $\widehat{\ell}(h_{\text{NN}}) = 0$. This is not surprising because, as we already said, NN stores the entire training set.

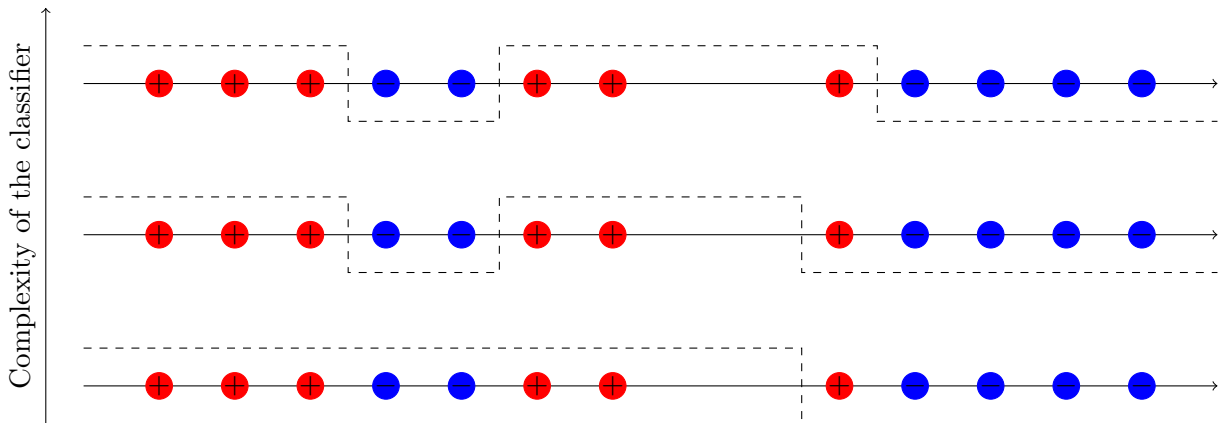


Figure 2: Plot of the $h_{k-\text{NN}}$ classifier for $k = 1, 3, 5$ on a 1-dimensional training set. As k increases, the classifier becomes simpler and the number of mistaken points in the training set increases.

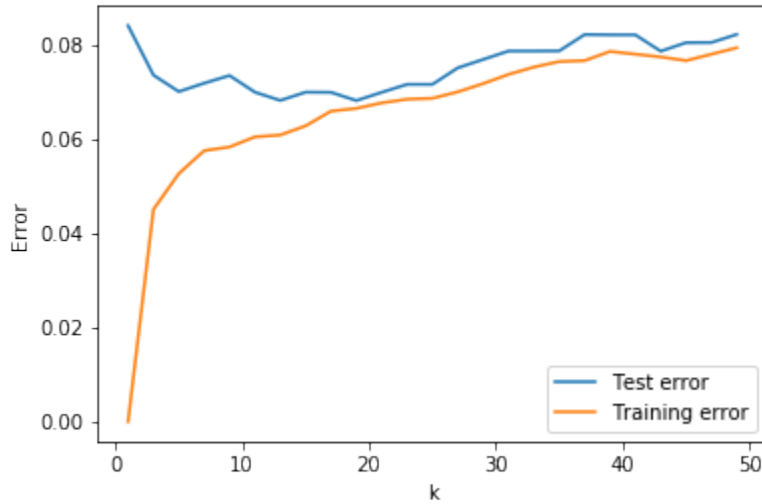
Starting from NN, we can obtain a family of algorithms denoted by k -NN for $k = 1, 3, 5, \dots$, where k cannot be taken larger than the size of the training set. These algorithms are defined as follows: given a training set $S = (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$, k -NN generates a classifier $h_{k-\text{NN}}$ such that $h_{k-\text{NN}}(\mathbf{x})$ is the label $y_t \in \{-1, 1\}$ appearing in the majority of the k points $\mathbf{x}_t \in S$ which are closest to \mathbf{x} .¹ Hence, in order to compute $h_{k-\text{NN}}(\mathbf{x})$, we perform the following operations:

1. Find the k training points $\mathbf{x}_{t_1}, \dots, \mathbf{x}_{t_k}$ closest to \mathbf{x} .¹ Let y_{t_1}, \dots, y_{t_k} be their labels.
2. If the majority of the labels y_{t_1}, \dots, y_{t_k} is $+1$, then $h_{k-\text{NN}}(\mathbf{x}) = +1$; if the majority is -1 , then $h_{k-\text{NN}}(\mathbf{x}) = -1$.

Note that, for each $k \geq 1$ and for each \mathbf{x}_t in the training set, \mathbf{x}_t is always included in the k points that are closest to \mathbf{x}_t .

It is important to note that, unlike 1-NN, in general we have that $\widehat{\ell}(h_{k-\text{NN}}) > 0$. Moreover, in Figure 2 we see that, as k grows, the classifiers generated by k -NN become simpler. In particular, when k is equal to the size of the training set, $h_{k-\text{NN}}$ becomes a constant classifier that always predicts the most common label in the training set.

¹Just like in the case of 1-NN, there could be training points at the same distance from \mathbf{x} such that more than k points are closest to \mathbf{x} . In this case we proceed by ranking the training points based on their distance from \mathbf{x} and then taking the k' closest points where k' is the smallest integer bigger or equal to k such that the $(k' + 1)$ -th point in the ranking has distance from \mathbf{x} strictly larger than the k' -th point. If no such k' exists, then we take all the points



The figure above shows the typical trend of training error (orange curve) and test error (blue curve) of the k -NN classifier for increasing values of the parameter k on a real dataset (Breast Cancer Wisconsin) for binary classification with zero-one loss. Note that the minimum of the test error is attained at a value corresponding to a $h_{k\text{-NN}}$ classifier with training error generally bigger than zero. Small values of k cause overfitting (small training error and large test error), while large values of k cause underfitting (training and test error both large).

In addition to binary classification, k -NN can be used to solve multiclass classification problems (where \mathcal{Y} contains more than two symbols) and also regression problems (where $\mathcal{Y} = \mathbb{R}$). In the first case, we operate like in the binary case and predict using the label corresponding to the majority of the labels of the k closest training points. In the second case, the prediction is the average of the labels of the k closest training points.

in the training set. If k' is strictly bigger than k , even, and there is an equal number of closest points with positive and negative labels, then the algorithm predicts a default value in $\{-1, 1\}$.