

k -Nearest Neighbors

Wiebke Köpp

Technische Universität München

Reading Material:

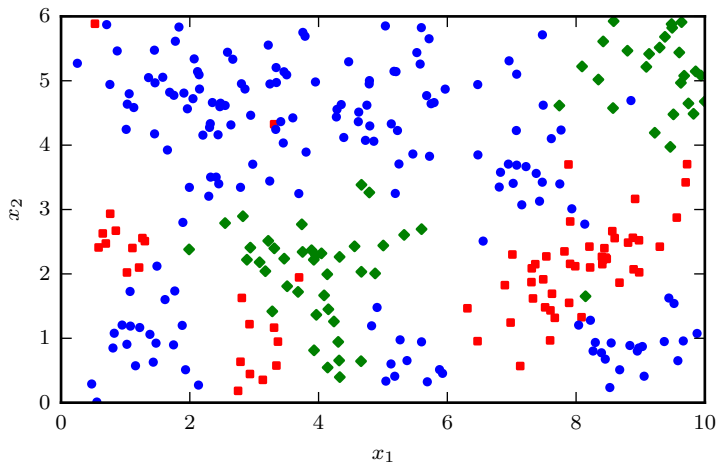
"Machine Learning: A Probabilistic Perspective" by Murphy [ch. 1.4.1 - 1.4.3]

Further extra reading:

"Bayesian Reasoning and Machine Learning" by Barber [ch. 14]

Note: These slides are adapted from slides originally by Daniała Korhammer

The data from last time



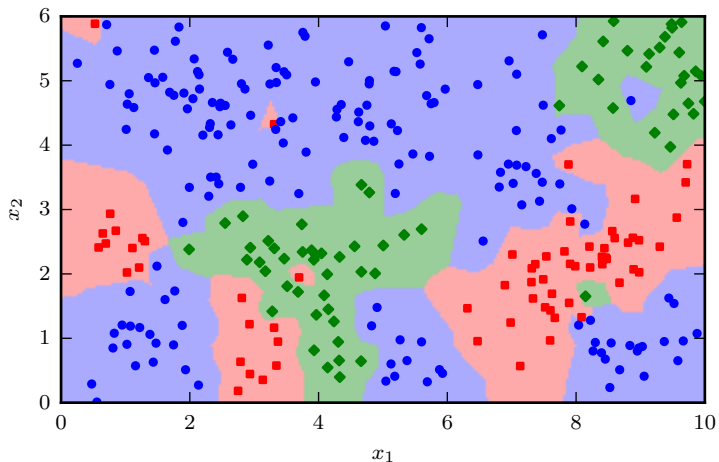
How do we intuitively label new samples by hand?
Look at the *surrounding* points.

1-NN Algorithm

- define a distance measure (e.g. Euclidean distance)
- compute the nearest neighbor for all new data points
- and label them with the label of their nearest neighbor

This works for both *classification* and *regression*.

1-NN



This corresponds to a Voronoi tessellation.
And results in overfitting...

k -Nearest Neighbor Classification

More *robust* against errors in the training set:

Look at a number of neighbors!

Let $\mathcal{N}_k(\mathbf{x})$ be the k nearest neighbors of a vector \mathbf{x} , then in classification tasks:

$$p(z = c \mid \mathbf{x}, k) = \frac{1}{k} \sum_{i \in \mathcal{N}_k(\mathbf{x})} \mathbb{I}(z_i = c),$$

$$y = \arg \max_c p(z = c \mid \mathbf{x}, k)$$

with the *indicator variable* $\mathbb{I}(e)$ is defined as:

$$\mathbb{I}(e) = \begin{cases} 1 & \text{if } e \text{ is true} \\ 0 & \text{if } e \text{ is false.} \end{cases}$$

i.e., the vector will be labeled by the mode of its neighbors' labels.

k -Nearest-Neighbor Regression

Regression is similar:

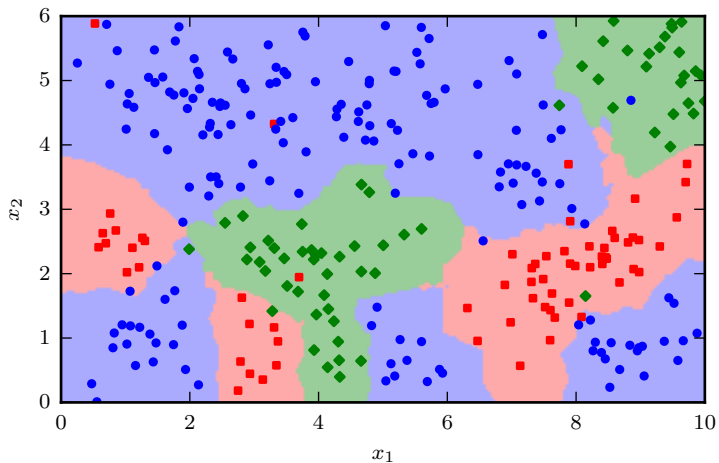
Let $\mathcal{N}_k(\mathbf{x})$ be the k nearest neighbors of a vector \mathbf{x} , then for regression:

$$y = \frac{1}{\mathcal{C}} \sum_{i \in \mathcal{N}_k(\mathbf{x})} \frac{1}{d(\mathbf{x}, \mathbf{x}_i)} z_i,$$

with $\mathcal{C} = \sum_{i \in \mathcal{N}_k(\mathbf{x})} \frac{1}{d(\mathbf{x}, \mathbf{x}_i)}$ the normalization constant and $d(\mathbf{x}, \mathbf{x}_i)$ being a distance measure between \mathbf{x} and \mathbf{x}_i ,

i.e., the vector will be labeled by a *weighted mean* of its neighbors' values.

3-NN

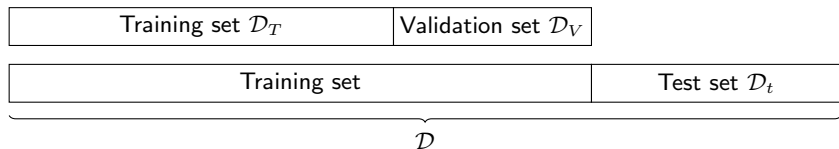


So, how many neighbors are best?

Determining k

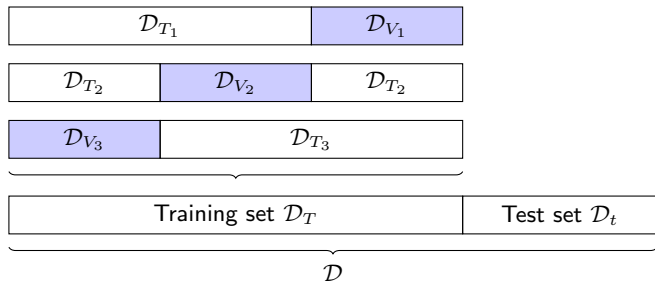
k is called a *hyper-parameter*.

Keep part of your training set separate and evaluate the algorithm with different k on this *validation set* to choose the best.



Better even: Do cross-validation!

K -fold Cross-Validation



- Split your training data into K folds (10-fold CV is common).
- Use $K - 1$ folds for training and the remaining for evaluation.
- Average over all folds to get an estimate of the error for a certain setting of your *hyper-parameters* (such as k in k -NN).
- Try different settings for your hyper-parameters.
- Use all your training data and the best hyper-parameters for final training (and testing) of your model.

The extreme case - LOOCV

In *leave-one-out-cross validation* (LOOCV) we train on all but one sample.

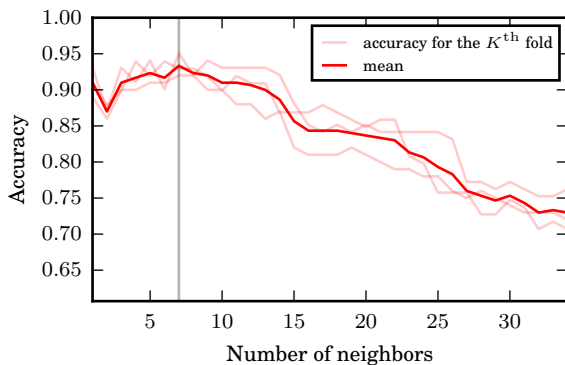
If we have N samples, this is the same as N -fold cross-validation.

LOOCV is interesting if we do not have a lot of data and we want to use as much of it for training as possible but still get a good estimate of model performance.

But it also means that we need to train our model N times...

If we have sufficiently large amounts of data and training our model is computationally expensive, we better stick to lower numbers of K or a single validation set.

Cross-validation for determining k

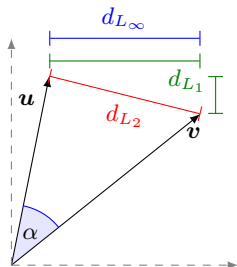


We choose $k = 7$.

Distance Measures

- Euclidean distance (L_2 norm): $\sqrt{\sum_i (u_i - v_i)^2}$
- L_1 norm: $\sum_i |u_i - v_i|$
- L_∞ norm: $\max_i |u_i - v_i|$
- Angle:

$$\cos \alpha = \frac{\mathbf{u}^T \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|}$$

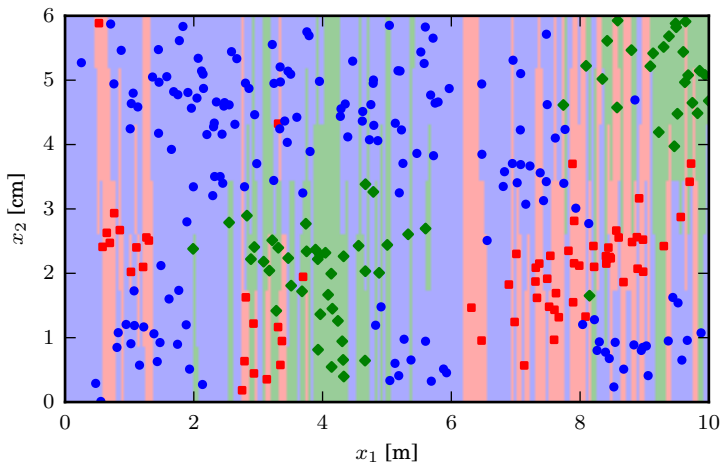


- Mahalanobis distance (Σ is positive (semi) definite and symmetric):

$$\sqrt{(\mathbf{u} - \mathbf{v})^T \Sigma^{-1} (\mathbf{u} - \mathbf{v})}$$

- Hamming distance, Edit distance, ...

Scaling Issues



The same old example but one of our features is in the order of meters, the other in the order of centimeters. ($k = 1$)

Circumventing Scaling Issues

- Data *standardization*
Scale each feature to zero mean and unit variance.

$$x_{i,\text{std}} = \frac{x_i - \mu_i}{\sigma_i}$$

(This is a standard procedure in machine learning. Many models are sensitive to differences in scale.)

- Use the Mahalanobis distance.

$$\text{mahalanobis}(\mathbf{x}_1, \mathbf{x}_2) = \sqrt{(\mathbf{x}_1 - \mathbf{x}_2)^T \mathbf{\Sigma}^{-1} (\mathbf{x}_1 - \mathbf{x}_2)}$$

If $\mathbf{\Sigma} = \begin{bmatrix} \sigma_1^2 & 0 & 0 \\ 0 & \cdots & 0 \\ 0 & 0 & \sigma_n^2 \end{bmatrix}$, this is equal to *normalized Euclidean distance*

The curse of dimensionality

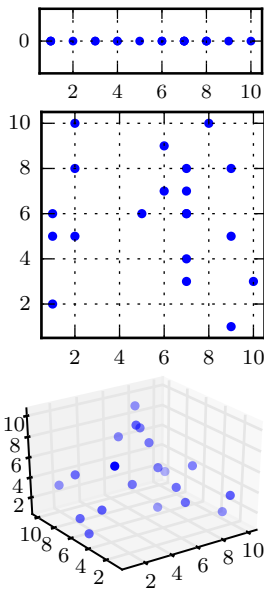
Say, you have a discrete one-dimensional input space $x \in \{1, 2, \dots, 10\}$. Your data is uniformly distributed over this space.

You have $N = 20$ samples and your data covers 100% of the input space.

Add a second dimension (now $x \in \{1, \dots, 10\}^2$) and your data only covers 18% of the input space.

Once you add a third dimension you only cover 2%. The nearest neighbor will now be pretty far away..

N has to grow exponentially with the number of features. Consider this when using k -NN on high-dimensional data.



A probabilistic interpretation of k -NN

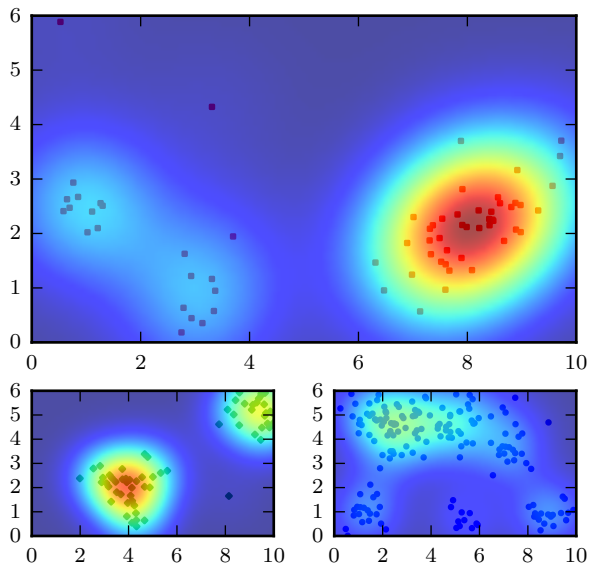
Assumption: Data points are distributed normally.

A *generative* model for data from class c :

$$\begin{aligned} p(\mathbf{x} \mid z = c) &= \frac{1}{N_c} \sum_{\mathbf{x}_n \in \text{class } c} \mathcal{N}(\mathbf{x} \mid \mathbf{x}_n, \sigma^2 \mathbf{I}) \\ &= \frac{1}{N_c} \frac{1}{(2\pi\sigma^2)^{D/2}} \sum_{\mathbf{x}_n \in \text{class } c} e^{-\frac{\|\mathbf{x} - \mathbf{x}_n\|^2}{(2\sigma^2)}} , \end{aligned}$$

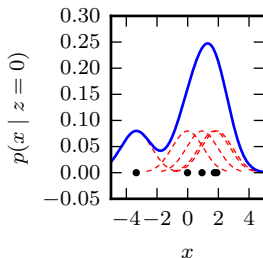
where D is the dimension of the data points, N_c is the number of training points in class c .

Generative models for classes red, blue and green



A probabilistic interpretation of k -NN

A *Parzen density estimator*: Uniform weighted sum of Gaussian distributions centered on the training points.



Bayes' rule helps with a new data point \mathbf{x}^* :

$$p(z = c | \mathbf{x}^*) = \frac{p(\mathbf{x}^* | z = c) p(z = c)}{\sum_{i \in C} p(\mathbf{x}^* | z = i) p(z = i)}$$

$$p(z = c | \mathbf{x}^*) \propto p(\mathbf{x}^* | z = c) p(z = c)$$

Neighborhood Component Analysis - NCA

Instead of defining the distance metric for k -NN, we can also learn it!
We want to find a *good* Q for the (squared) Mahalanobis distance:

$$\text{mahalanobis}(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{x}_1 - \mathbf{x}_2)^T Q (\mathbf{x}_1 - \mathbf{x}_2)$$

For *any* such Q we can find A , such that $Q = A^T A$, so we can rewrite the distance as:

$$\text{mahalanobis}(\mathbf{x}_1, \mathbf{x}_2) = (A\mathbf{x}_1 - A\mathbf{x}_2)^T (A\mathbf{x}_1 - A\mathbf{x}_2)$$

NCA - a good metric

How do we define *good*?

A should be chosen in such a way that when we take one exemplar from the training set and try to classify it based on the remaining training set, it should be classified into its *known* category. And this should hold for *all* exemplars in the training set (*Leave One Out Cross Validation—LOOCV*).

But an *infinitesimal* change in A can effect the classification performance by a *finite* amount.

NCA - the probabilistic perspective

A better-behaving performance measure (still based on LOOCV):
Don't assign classes directly but use probabilities instead.

$$p_{ij} = \frac{\exp(-\|\mathbf{A} \mathbf{x}_i - \mathbf{A} \mathbf{x}_j\|^2)}{\sum_{k \neq i} \exp(-\|\mathbf{A} \mathbf{x}_i - \mathbf{A} \mathbf{x}_k\|^2)}, \quad p_{ii} = 0$$

Each training point i selects another point j as its neighbor with some probability p_{ij} .

We compute the probability p_i that point i will be correctly classified:

$$p_i = \sum_{j \in C_i} p_{ij}, \quad C_i = \{j \mid c_i = c_j\}$$

NCA - Finding \mathbf{A}

The goal is now to find some \mathbf{A} that maximizes the *expected number of correctly classified points*:

$$f(\mathbf{A}) = \sum_i \sum_{j \in C_i} p_{ij} = \sum_i p_i$$

A function (*objective*) that depends solely on \mathbf{A} .

How to *maximize* it? Follow the *gradient* until the top (*gradient ascent*).

A gradient with respect to a matrix!

$$\frac{\partial f}{\partial \mathbf{A}} = 2\mathbf{A} \sum_i \left(p_i \sum_k p_{ik} \mathbf{x}_{ik} \mathbf{x}_{ik}^T - \sum_{j \in C_i} p_{ij} \mathbf{x}_{ij} \mathbf{x}_{ij}^T \right)$$

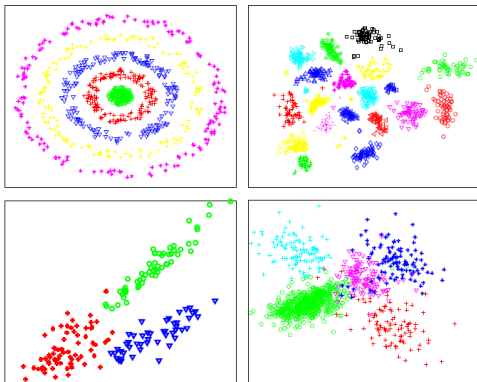
with $\mathbf{x}_{ij} = \mathbf{x}_i - \mathbf{x}_j$.

NCA - Impact of A

The algorithm also learns automatically a real-valued estimate of the optimal number of neighbors: by scaling up A uniformly, the method consults fewer neighbors and vice versa.

By restricting A to be a non-square matrix of size $d \times D$, NCA can also do *linear dimensionality reduction* (thus, the metric will be of low rank if $d \leq D$).

NCA for dimensionality reduction



NCA on data that was originally 3-D (top left), 13-D (lower left), 560-D (top right, face data) and 256-D (lower right, digits data)

What we learned

- k -NN Algorithm
- Cross-Validation
- Distance Metrics
- Curse of Dimensionality
- Probabilistic k -NN and NCA