

# Machine Learning

## Lecture 3

### Simple Classifiers: Nearest Centroids and KNN

Felix Bießmann

Beuth University & Einstein Center for Digital Future

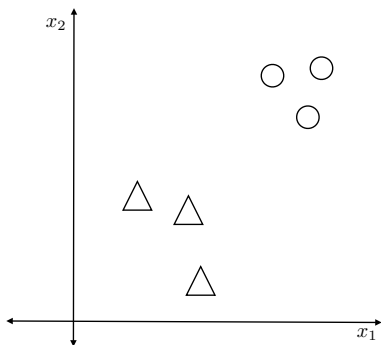


# Overview of today's lecture

- Today we will introduce two simple classifiers
  1. Nearest Centroid Classifier (NCC)
  2. K-Nearest Neighbor (KNN)
- These algorithms are extremely powerful
- Often they can compete with complex algorithms



# Prototypes: Psychological Models of Abstract Ideas



Psychologists postulated that we learn **prototypes** [Jaekel, 2007; Posner and Keele, 1968]

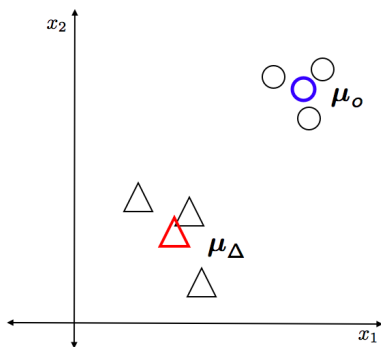
**Toy data example:**

Two dimensional input  $\mathbf{x} \in \mathbb{R}^2$

Two *classes* of data,  $\Delta$  and  $\circ$



# Prototypes: Psychological Models of Abstract Ideas



Prototypes  $\mu_{\Delta}$  and  $\mu_o$  can be the class means

$$\mu_{\Delta} = 1/N_{\Delta} \sum_n^{N_{\Delta}} \mathbf{x}_{\Delta,n}$$

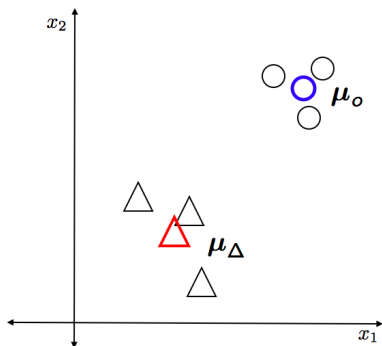
$$\mu_o = 1/N_o \sum_n^{N_o} \mathbf{x}_{o,n}$$

Distance from  $w_{\Delta}$  to new data  $\mathbf{x}$

$$\|\mu_{\Delta} - \mathbf{x}\|_2$$



# Prototypes: Psychological Models of Abstract Ideas



For new data  $x$  check:  
**Is  $x$  more similar to  $\mu_o$ ?**

$$\|\mu_{\Delta} - x\| > \|\mu_o - x\|$$

yes?  $\rightarrow x$  belongs to  $\mu_o$

no?  $\rightarrow x$  belongs to  $\mu_{\Delta}$

This is called a  
**nearest centroid classifier**



# Nearest Centroid Classification Algorithm (Batch Mode)

---

## Algorithm 1 Computation of Class-Centroids

---

**Require:** data  $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^D$ , labels  $y_1, \dots, y_N \in \{1, \dots, K\}$

**Ensure:** Class means  $\boldsymbol{\mu}_k$ ,  $k \in \{1, \dots, K\}$

- 1: # Initialize means and counters for each class
  - 2: # Computation of class means
  - 3: **for** Class  $k = 1, \dots, K$  **do**
  - 4:      $\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{i=1}^{N_k} \mathbf{x}_i$
  - 5: **end for**
- 



# Batch Computations vs. Streaming

Solutions for algorithms can be obtained

- In Batch Mode:
  - Use all available data at once
  - Requires to store all data in memory
- In Streaming Mode:
  - Use one data point at a time
  - Requires to store **only** centroids



# Iterative Computation of the Mean

Given the mean  $\mu_{N-1}$  computed from  $N - 1$  samples we want to update  $\mu_{N-1}$  with the  $N$ th sample  $\mathbf{x}_N$  to obtain  $\mu_N$

$$\begin{aligned}\mu_N &= \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \\ &= \frac{1}{N} \sum_{n=1}^{N-1} \mathbf{x}_n + \frac{1}{N} \mathbf{x}_N \\ &= \frac{N-1}{N} \underbrace{\frac{1}{N-1} \sum_{n=1}^{N-1} \mathbf{x}_n}_{\mu_{N-1}} + \frac{1}{N} \mathbf{x}_N \\ &= \frac{N-1}{N} \mu_{N-1} + \frac{1}{N} \mathbf{x}_N\end{aligned}$$





# Nearest Centroid Classification Algorithm (Streaming)

---

## Algorithm 2 Iterative computation of Class-Centroids

---

**Require:** data  $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^D$ , labels  $y_1, \dots, y_N \in \{1, \dots, K\}$

**Ensure:** Class means  $\mu_k$ ,  $k \in \{1, \dots, K\}$

1: # Initialize means and counters for each class

2:  $\forall k : \mu_k = \mathbf{1} \cdot 0, N_k = 0$

3: # Iterative computation of class means

4: **for** Data point  $i = 1, \dots, N$  **do**

5:   # Update means and counters

6:    $k = y_i$

7:    $\mu_k = \frac{N_k}{N_k+1} \mu_k + \frac{1}{N_k+1} \mathbf{x}_i$

8:    $N_k = N_k + 1$

9: **end for**

---



# Nearest Centroid Classification

---

## Algorithm 3 Nearest Centroid Prediction

---

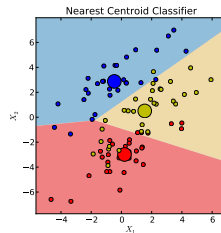
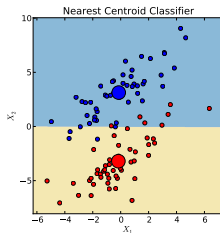
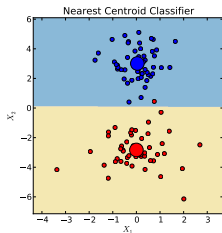
**Require:** Data point  $\mathbf{x} \in \mathbb{R}^D$ , class centroids  $\boldsymbol{\mu}_k$ ,  $k \in \{1, \dots, K\}$

**Ensure:** Class membership  $k^*$

- 1: # Compute nearest class centroid
  - 2:  $k^* = \operatorname{argmin}_k \|\boldsymbol{\mu}_k - \mathbf{x}\|_2$ .
- 



# Toy Data Example NCC



# From Prototypes to Linear Classification

$$\begin{aligned} \text{distance}(\mathbf{x}, \mu_{\Delta}) &> \text{distance}(\mathbf{x}, \mu_o) \\ \|\mathbf{x} - \mu_{\Delta}\| &> \|\mathbf{x} - \mu_o\| \end{aligned} \tag{1}$$



# From Prototypes to Linear Classification

$$\text{distance}(\mathbf{x}, \mu_{\Delta}) > \text{distance}(\mathbf{x}, \mu_o) \quad (1)$$

$$\|\mathbf{x} - \mu_{\Delta}\| > \|\mathbf{x} - \mu_o\|$$

$$\Leftrightarrow \|\mathbf{x} - \mu_{\Delta}\|^2 > \|\mathbf{x} - \mu_o\|^2$$

$$\Leftrightarrow \mathbf{x}^T \mathbf{x} - 2\mu_{\Delta}^T \mathbf{x} + \mu_{\Delta}^T \mu_{\Delta} > \mathbf{x}^T \mathbf{x} - 2\mu_o^T \mathbf{x} + \mu_o^T \mu_o$$

$$\Leftrightarrow \mu_{\Delta}^T \mathbf{x} - \mu_{\Delta}^2/2 < \mu_o^T \mathbf{x} - \mu_o^2/2$$

$$\Leftrightarrow 0 < \underbrace{(\mu_o - \mu_{\Delta})^T}_{\mathbf{w}} \mathbf{x} - 1/2 \underbrace{(\mu_o^T \mu_o - \mu_{\Delta}^T \mu_{\Delta})}_{\beta}$$



# From Prototypes to Linear Classification

$$\text{distance}(\mathbf{x}, \mu_{\Delta}) > \text{distance}(\mathbf{x}, \mu_o) \quad (1)$$

$$\|\mathbf{x} - \mu_{\Delta}\| > \|\mathbf{x} - \mu_o\|$$

$$\Leftrightarrow \|\mathbf{x} - \mu_{\Delta}\|^2 > \|\mathbf{x} - \mu_o\|^2$$

$$\Leftrightarrow \mathbf{x}^T \mathbf{x} - 2\mu_{\Delta}^T \mathbf{x} + \mu_{\Delta}^T \mu_{\Delta} > \mathbf{x}^T \mathbf{x} - 2\mu_o^T \mathbf{x} + \mu_o^T \mu_o$$

$$\Leftrightarrow \mu_{\Delta}^T \mathbf{x} - \mu_{\Delta}^2/2 < \mu_o^T \mathbf{x} - \mu_o^2/2$$

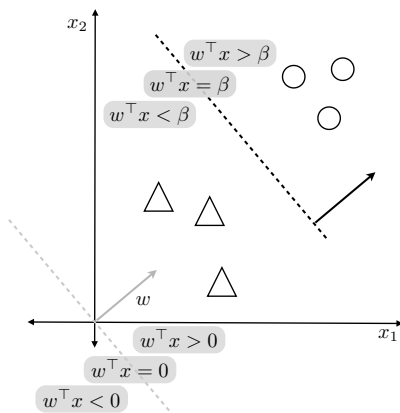
$$\Leftrightarrow 0 < \underbrace{(\mu_o - \mu_{\Delta})^T \mathbf{x}}_{\mathbf{w}} - 1/2 \underbrace{(\mu_o^T \mu_o - \mu_{\Delta}^T \mu_{\Delta})}_{\beta}$$

## Linear Classification

$$\mathbf{w}^T \mathbf{x} - \beta = \begin{cases} > 0 & \text{if } \mathbf{x} \text{ belongs to class } o \\ < 0 & \text{if } \mathbf{x} \text{ belongs to class } \Delta \end{cases} \quad (2)$$



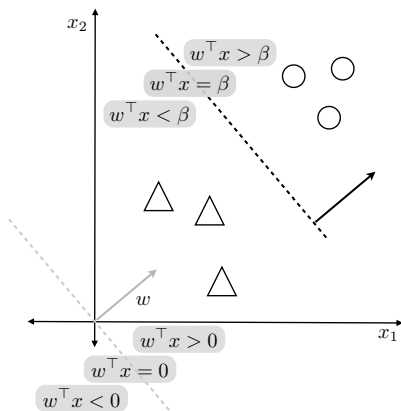
# Linear Classification



$$\mathbf{w}^\top \mathbf{x} - \beta = \begin{cases} > 0 & \text{if } \mathbf{x} \text{ belongs to } o \\ < 0 & \text{if } \mathbf{x} \text{ belongs to } \Delta \end{cases}$$



# Linear Classification



$$\mathbf{w}^T \mathbf{x} - \beta = \begin{cases} > 0 & \text{if } \mathbf{x} \text{ belongs to } o \\ < 0 & \text{if } \mathbf{x} \text{ belongs to } \Delta \end{cases}$$

The *offset*  $\beta$  can be included in  $\mathbf{w}$

$$\tilde{\mathbf{x}} \leftarrow \begin{bmatrix} 1 \\ \mathbf{x} \end{bmatrix} \quad \tilde{\mathbf{w}} \leftarrow \begin{bmatrix} -\beta \\ \mathbf{w} \end{bmatrix}$$

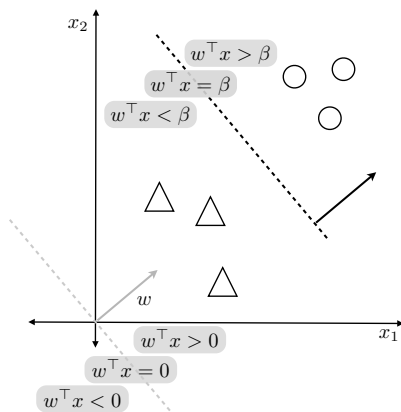
such that

$$\tilde{\mathbf{w}}^T \tilde{\mathbf{x}} = \mathbf{w}^T \mathbf{x} - \beta.$$





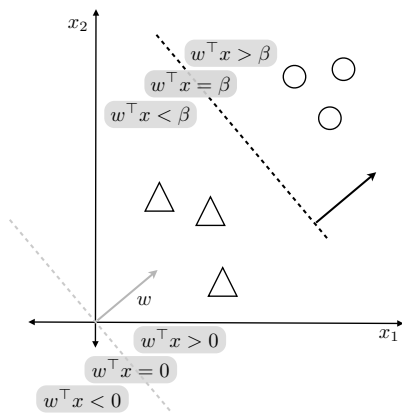
# Linear Classification



- What is a good  $w$ ?
- Some proposals:
  - Logistic Regression
  - Perceptrons
  - Support Vector Machines
  - Ridge Regression
- Linear methods have different ways of defining what a good  $w$  is



# Linear Classification



Nearest Centroid Classification is a simple linear classifier

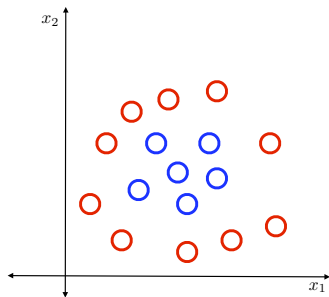
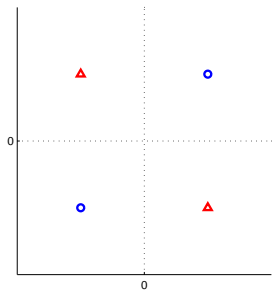
$$\mathbf{w} = \mu_o - \mu_{\Delta} \quad (3)$$

$$\beta = -1/2(\mu_o^\top \mu_o - \mu_{\Delta}^\top \mu_{\Delta})$$



# Problems with Linear Classification

Linear Classifiers fail on non-linear problems:

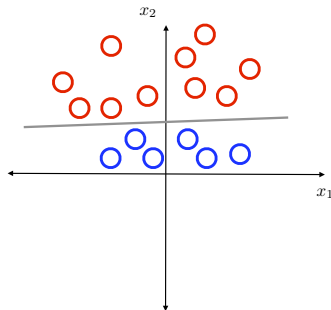
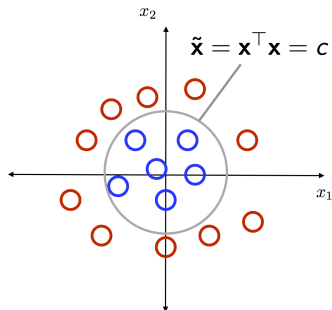


Try to separate these two-class data sets with **a single line**



# Problems with Perceptrons

If we can model the non-linearity,  
we can create new features  $\tilde{\mathbf{x}}$  which are linearly separable



But what if we do not know the non-linearity?



# K-Nearest Neighbor Classifier

- Nearest Centroid Classifiers require estimation of centroids
- K-Nearest Neighbor is simpler
- Idea:
  1. Find the  $k$  closest neighbors for a new data point  $\mathbf{x}$
  2. Look up labels of  $k$  closest neighbors
  3. Assign majority vote label for  $\mathbf{x}$



# K-Nearest Neighbor Classifier

We do not need to train a model for KNN –  
**the data is the model**

But (as with NCC) we need a distance function

Usually the euclidean distance is chosen:

$$d(\mathbf{x}_i, \mathbf{x}_j) = \|\mathbf{x}_i - \mathbf{x}_j\|_2 \quad j \in [1, \dots, N] \quad (4)$$

where  $N$  is the number of data points



# K-Nearest Neighbor Classifier: Pseudocode

---

## Algorithm 4 K-Nearest Neighbour Prediction

---

**Require:** Test data point  $\mathbf{x}_i \in \mathbb{R}^D$ , training data  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$ , corresponding labels  $\mathbf{y} = [y_1, \dots, y_N]^\top \in \mathbb{R}^N$ , number of neighbours  $K$

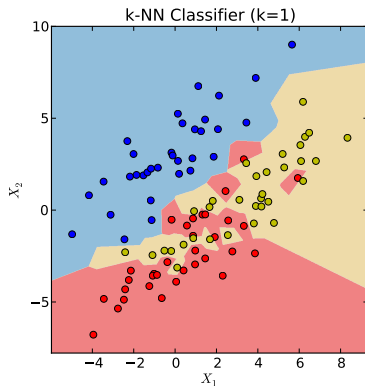
**Ensure:** Predicted label  $\mathbf{y}_i$

- 1: **for**  $\mathbf{x}_j$  in  $\mathbf{X}$  **do**
  - 2:   # Compute distance between test data  $\mathbf{x}_i$  and training data  $\mathbf{x}_j$
  - 3: **end for**
  - 4: Initialize a  $K$ -dimensional zero vector  $\gamma$
  - 5: **for**  $k$  in  $\mathbf{do}$
  - 6:   Count label of the  $k$ -nearest neighbor
  - 7:   #  $\gamma_{y_k} \leftarrow \gamma_{y_k} + 1$
  - 8: **end for**
  - 9: Predicted label is that with most votes (ties are broken at random)
  - 10:  $\mathbf{y}_i = \operatorname{argmax}_k(\gamma)$
- 

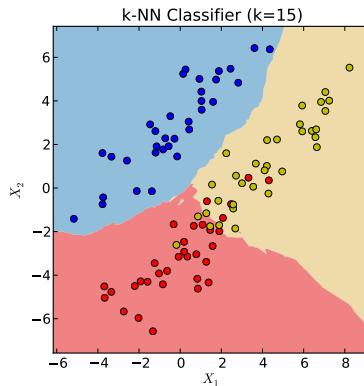


# K-Nearest Neighbor Classifier

Toy data problem: Linear classification



$k = 1$



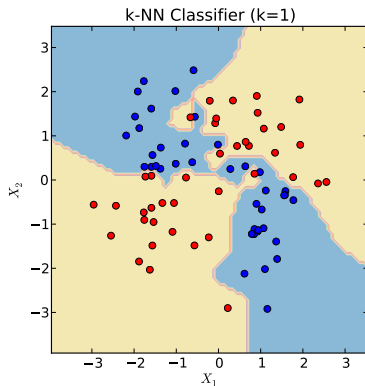
$k = 15$



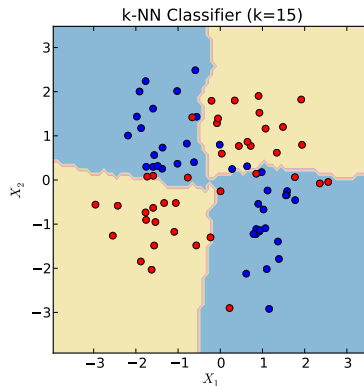


# K-Nearest Neighbor Classifier

Toy data problem: Nonlinear classification



$k = 1$



$k = 15$



# Problems with KNN

- Hyperparameter  $k$  needs to be set appropriately:
  - $k$  small: complex decision boundaries
  - $k$  large: smooth/simple decision boundaries
- Consider  $N$  data points  $\mathbf{x} \in \mathbb{R}^D$
- Find K Neighbors requires  $O(NND)$  operations
- For large data sets this is too costly
- Speedups can be gained by:
  - Trees for distance computations
  - Locality Sensitive Hashing for finding neighbors



# Summary

## Psychologists postulated we learn **Prototypes**

- Prototypes can be the class means

- Prototype theory is closely related to linear classification

## Nearest Centroid Classification

- New data is assigned to class with closest centroid

- Memory efficient: only requires to store  $D$ -dimensional centroids

- Iterative/Streaming version can scale to large data sets

## K-Nearest Neighbor Classifier

- Simple nonlinear classifier

- State-of-the-art prediction performance

- No training required

- Needs to evaluate pairwise distances of **all** data points

→ **Slow**



# References

- F. Jaekel. *Some Theoretical Aspects of Human Categorization Behaviour: Similarity and Generalization*. PhD thesis, 2007.
- M. I. Posner and S. W. Keele. On the genesis of abstract ideas. *Journal of Experimental Psychology*, 77(3):353–363, 1968.

