

# Classifiers, Learning

Tomáš Svoboda and Matěj Hoffmann

thanks to Daniel Novák and Filip Železný, Ondřej Drbohlav

Vision for Robots and Autonomous Systems, Center for Machine Perception

Department of Cybernetics

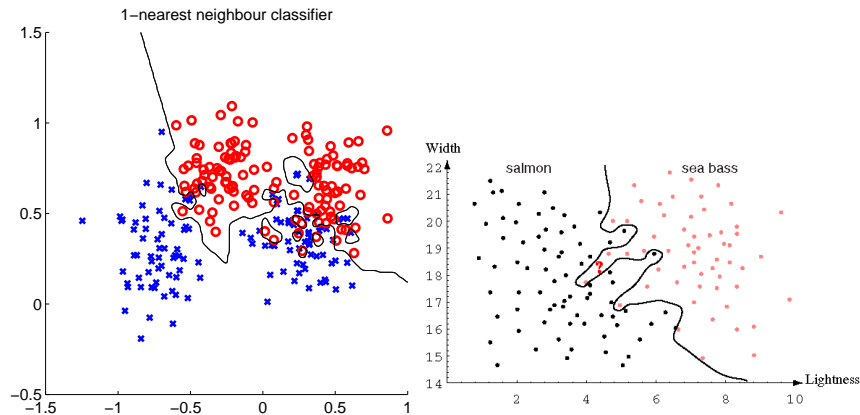
Faculty of Electrical Engineering, Czech Technical University in Prague

May 20, 2019

## K-Nearest neighbors classification

For a query  $\vec{x}$ :

- ▶ Find  $K$  nearest  $\vec{x}$  from the training (labeled) data.
- ▶ Classify to the class with the most exemplars in the set above.



Some properties:

- A *nonparametric* method – does not assume anything about the distribution (that it is Gaussian etc.)
- Can be used for classification or regression. Here: classification.
- Training: Only store feature vectors and their labels.
- Very simple and suboptimal. With unlimited nr. prototypes, error never worse than twice the Bayes rate (optimum).
- *instance-based* or *lazy* learning – function only approximated locally; computation only during inference.
- Limitations
  - Curse of dimensionality - for every additional dimension, one needs exponentially more points to cover the space.
  - Comp. complexity - has to look through all the samples all the time. Some speed-up is possible. E.g., storing data in a K-d tree.
  - Noise. Misclassified examples will remain in the database....

## $K$ — Nearest Neighbor and Bayes

Assume data:

- ▶  $N$  points  $\vec{x}$  in total.
- ▶  $N_j$  points in  $s_j$  class. Hence,  $\sum_j N_j = N$ .

We want classify  $\vec{x}$ . We draw a sphere centered at  $\vec{x}$  containing  $K$  points irrespective of class.  $V$  is the volume of this sphere.  $P(s_j|\vec{x}) = ?$

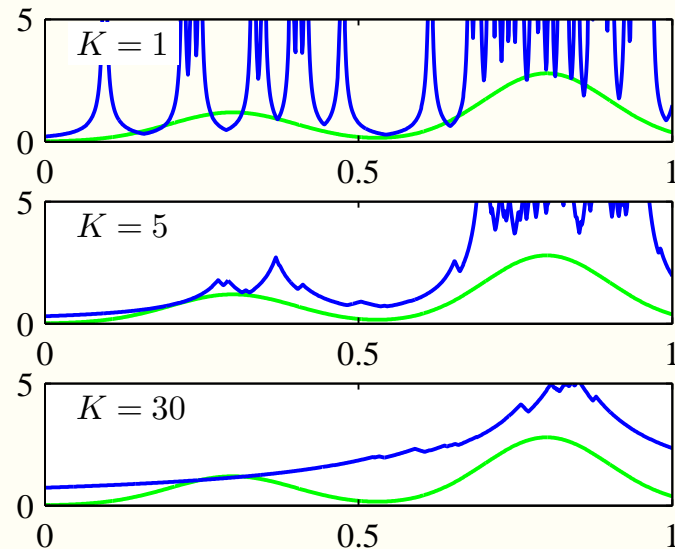
$$P(s_j|\vec{x}) = \frac{P(\vec{x}|s_j)P(s_j)}{P(\vec{x})}$$

$$P(s_j) = \frac{N_j}{N}$$

$$P(\vec{x}) = \frac{K}{NV}$$

$$P(\vec{x}|s_j) = \frac{K_j}{N_j V}$$

$$P(s_j|\vec{x}) = \frac{P(\vec{x}|s_j)P(s_j)}{P(\vec{x})} = \frac{K_j}{K}$$



A  $K$ —NN classifier can be understood as a non-parametric density estimator.

## $K$ — Nearest Neighbor and Bayes

Assume data:

- ▶  $N$  points  $\vec{x}$  in total.
- ▶  $N_j$  points in  $s_j$  class. Hence,  $\sum_j N_j = N$ .

We want classify  $\vec{x}$ . We draw a sphere centered at  $\vec{x}$  containing  $K$  points irrespective of class.  $V$  is the volume of this sphere.  $P(s_j|\vec{x}) = ?$

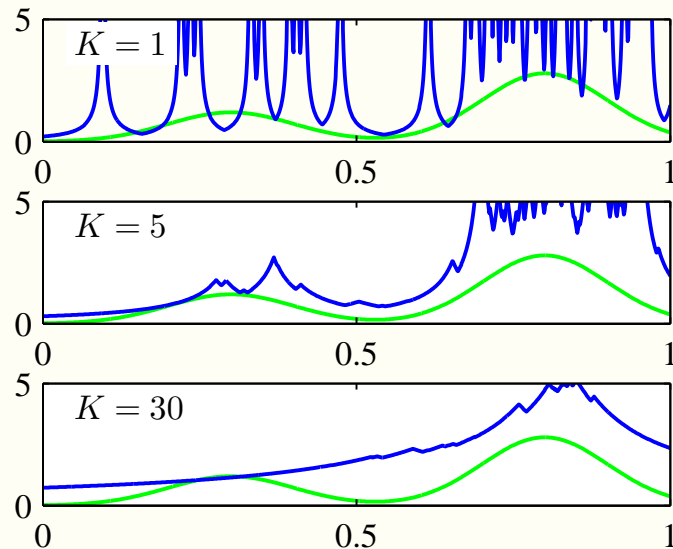
$$P(s_j|\vec{x}) = \frac{P(\vec{x}|s_j)P(s_j)}{P(\vec{x})}$$

$$P(s_j) = \frac{N_j}{N}$$

$$P(\vec{x}) = \frac{K}{NV}$$

$$P(\vec{x}|s_j) = \frac{K_j}{N_j V}$$

$$P(s_j|\vec{x}) = \frac{P(\vec{x}|s_j)P(s_j)}{P(\vec{x})} = \frac{K_j}{K}$$



A  $K$ —NN classifier can be understood as a non-parametric density estimator.

## $K$ — Nearest Neighbor and Bayes

Assume data:

- ▶  $N$  points  $\vec{x}$  in total.
- ▶  $N_j$  points in  $s_j$  class. Hence,  $\sum_j N_j = N$ .

We want classify  $\vec{x}$ . We draw a sphere centered at  $\vec{x}$  containing  $K$  points irrespective of class.  $V$  is the volume of this sphere.  $P(s_j|\vec{x}) = ?$

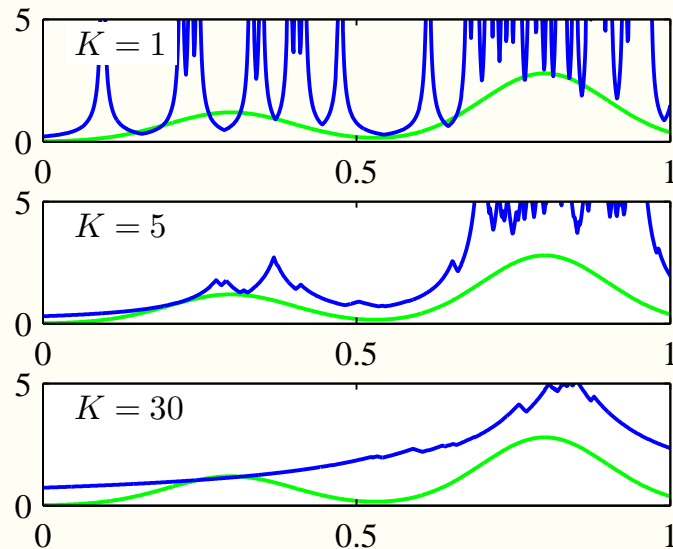
$$P(s_j|\vec{x}) = \frac{P(\vec{x}|s_j)P(s_j)}{P(\vec{x})}$$

$$P(s_j) = \frac{N_j}{N}$$

$$P(\vec{x}) = \frac{K}{NV}$$

$$P(\vec{x}|s_j) = \frac{K_j}{N_j V}$$

$$P(s_j|\vec{x}) = \frac{P(\vec{x}|s_j)P(s_j)}{P(\vec{x})} = \frac{K_j}{K}$$



A  $K$ —NN classifier can be understood as a non-parametric density estimator.

## $K$ — Nearest Neighbor and Bayes

Assume data:

- ▶  $N$  points  $\vec{x}$  in total.
- ▶  $N_j$  points in  $s_j$  class. Hence,  $\sum_j N_j = N$ .

We want classify  $\vec{x}$ . We draw a sphere centered at  $\vec{x}$  containing  $K$  points irrespective of class.  $V$  is the volume of this sphere.  $P(s_j|\vec{x}) = ?$

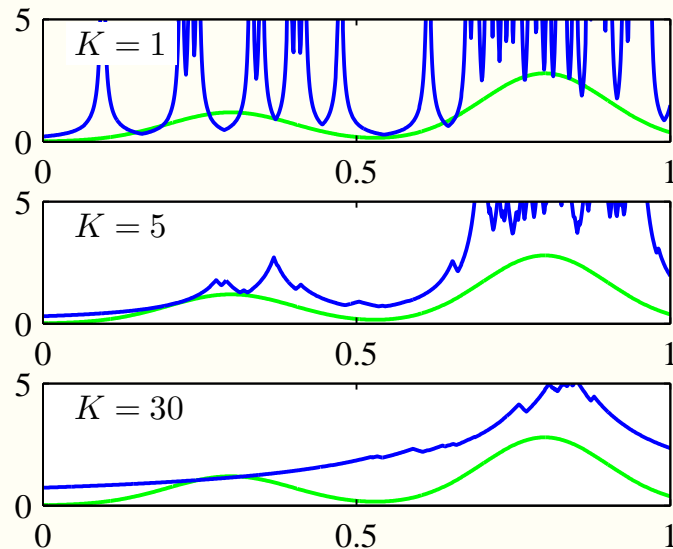
$$P(s_j|\vec{x}) = \frac{P(\vec{x}|s_j)P(s_j)}{P(\vec{x})}$$

$$P(s_j) = \frac{N_j}{N}$$

$$P(\vec{x}) = \frac{K}{NV}$$

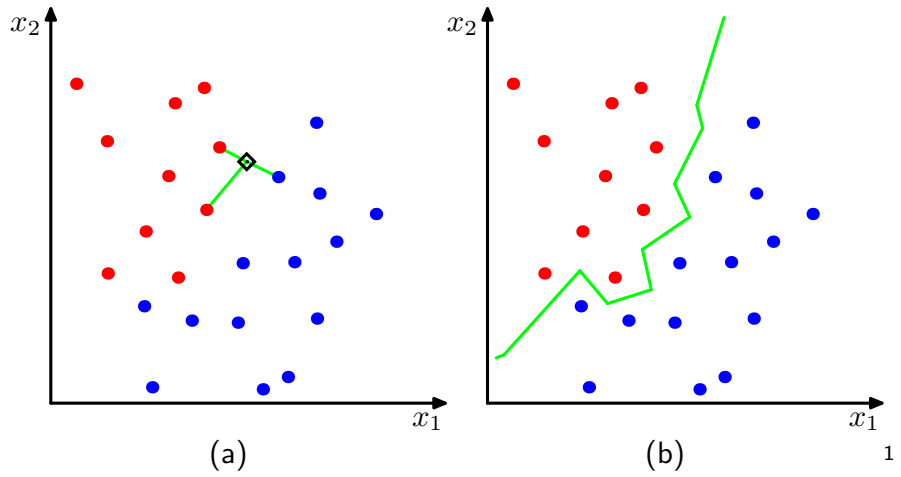
$$P(\vec{x}|s_j) = \frac{K_j}{N_j V}$$

$$P(s_j|\vec{x}) = \frac{P(\vec{x}|s_j)P(s_j)}{P(\vec{x})} = \frac{K_j}{K}$$



A  $K$ —NN classifier can be understood as a non-parametric density estimator.

# NN classification example

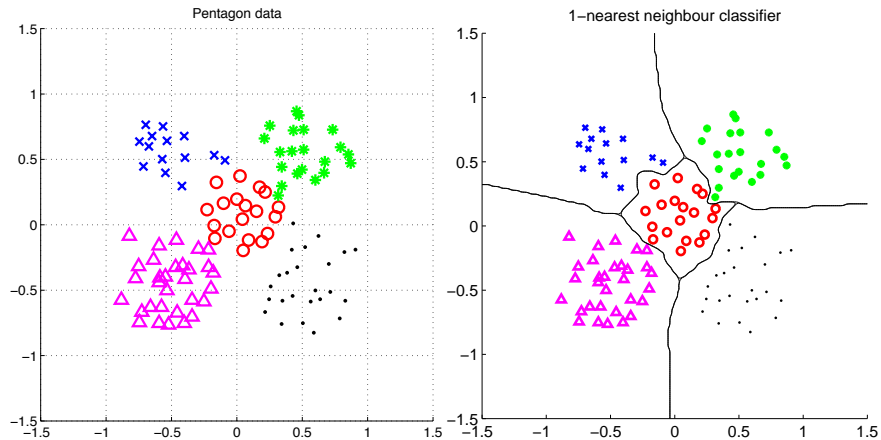


<sup>1</sup>Figs from [1]

## NN classification example

Fast on “learning”, very slow on decision.

There are ways for speeding it up, search for NN editing - making training data sparser, keeping only representative points.





## Metrics for NN classification

$$D(\mathbf{a}, \mathbf{b}) \geq 0$$

$$D(\mathbf{a}, \mathbf{b}) = 0 \text{ iff } \mathbf{a} = \mathbf{b}$$

$$D(\mathbf{a}, \mathbf{b}) = D(\mathbf{b}, \mathbf{a})$$

$$D(\mathbf{a}, \mathbf{b}) + D(\mathbf{b}, \mathbf{c}) \geq D(\mathbf{a}, \mathbf{c})$$

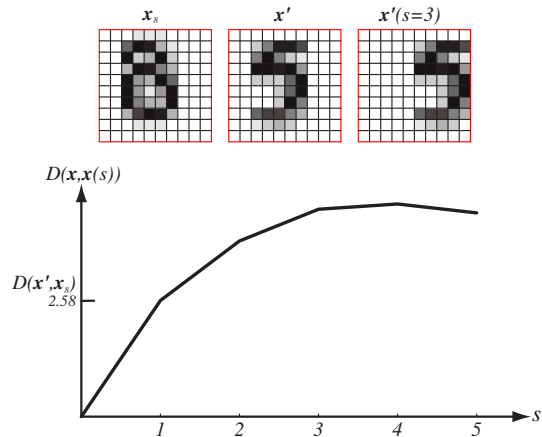
## Metrics for NN classification

$$D(\mathbf{a}, \mathbf{b}) \geq 0$$

$$D(\mathbf{a}, \mathbf{b}) = 0 \text{ iff } \mathbf{a} = \mathbf{b}$$

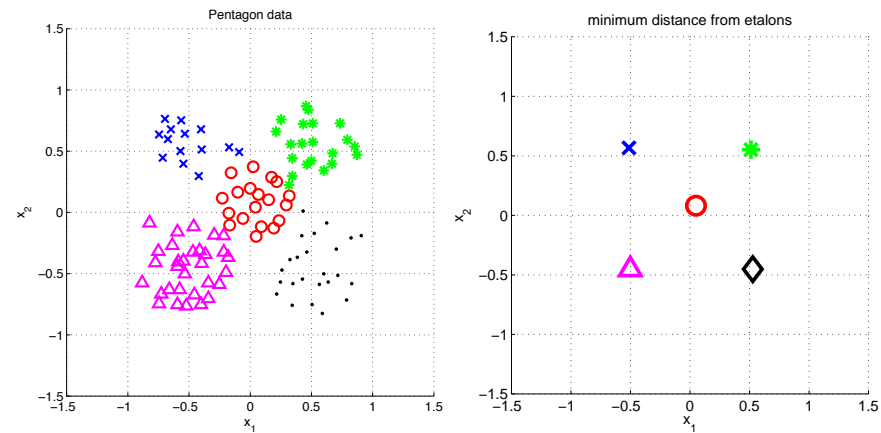
$$D(\mathbf{a}, \mathbf{b}) = D(\mathbf{b}, \mathbf{a})$$

$$D(\mathbf{a}, \mathbf{b}) + D(\mathbf{b}, \mathbf{c}) \geq D(\mathbf{a}, \mathbf{c})$$



Invariance to geometrical transformations?

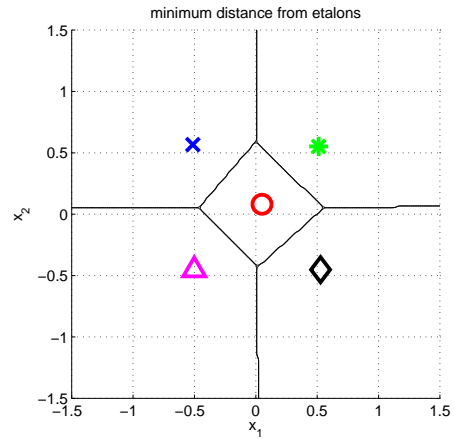
# Etalon based classification



Represent  $\vec{x}$  by **etalon** ,  $\vec{e}_s$  per each class  $s \in S$

## Separate etalons

$$s^* = \arg \min_{s \in S} (\|\vec{x} - \vec{e}_s\|^2 + o_s)$$

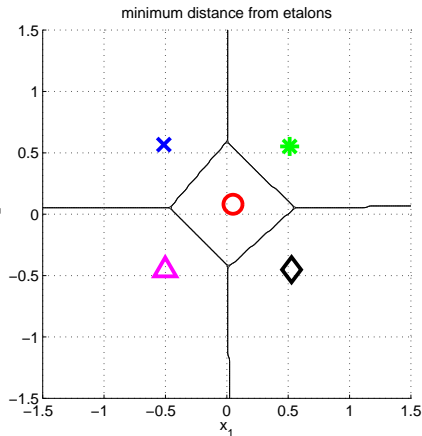


## What etalons?

If  $\mathcal{N}(\vec{x}|\vec{\mu}, \Sigma)$ ; all classes same covariance matrices, then

$$\vec{e}_s \stackrel{\text{def}}{=} \vec{\mu}_s = \frac{1}{|\mathcal{X}^s|} \sum_{i \in \mathcal{X}^s} \vec{x}_i^s$$

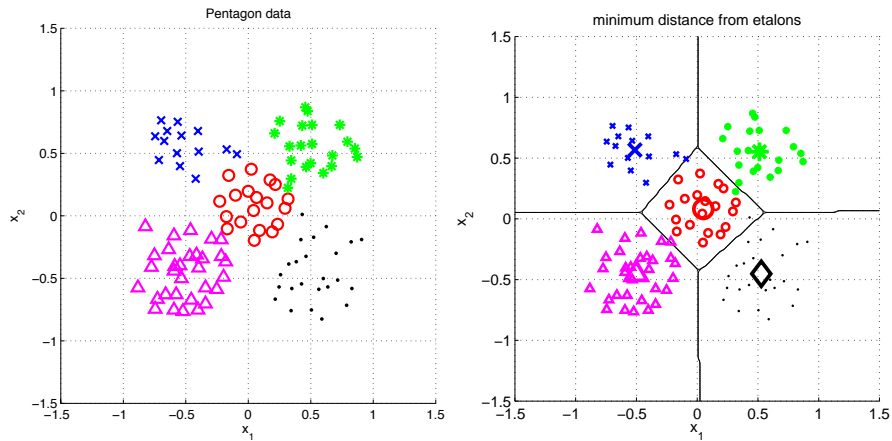
and separating hyperplanes halve distances between pairs.



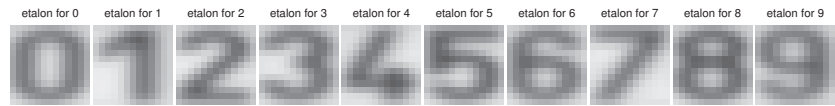
$$\mathcal{N}(\vec{x}|\vec{\mu}, \Sigma) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(\vec{x} - \vec{\mu})^\top \Sigma^{-1}(\vec{x} - \vec{\mu})\right\}$$

Etalon based classification,  $\vec{e}_s = \vec{\mu}_s$

Some wrongly classified samples. We like the simple idea. Are there better etalons? How to find them?

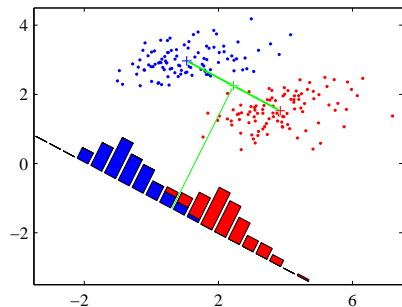


# Digit recognition - etalons $\vec{e}_s = \vec{\mu}_s$



Figures from [5]

## Better etalons – Fischer linear discriminant



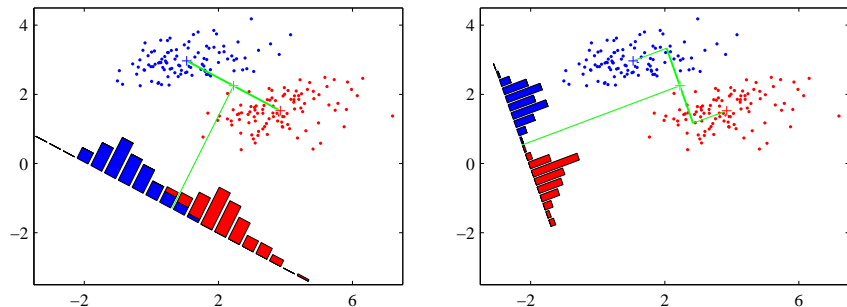
- ▶ Dimensionality reduction
- ▶ Maximize distance between means, ...
- ▶ ... and minimize within class variance. (minimize overlap)

Figures from [1]

Searching for a projection of the data to minimize intra-class variance and maximize inter-class variance.



## Better etalons – Fischer linear discriminant

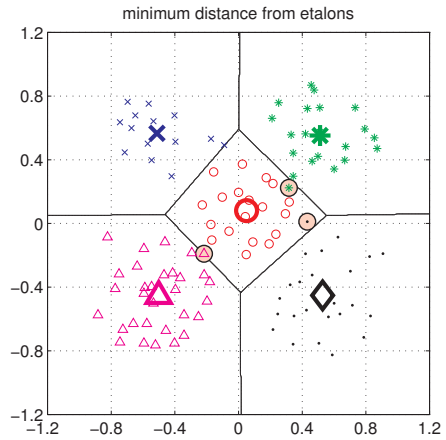


- ▶ Dimensionality reduction
- ▶ Maximize distance between means, ...
- ▶ ... and minimize within class variance. (minimize overlap)

Figures from [1]

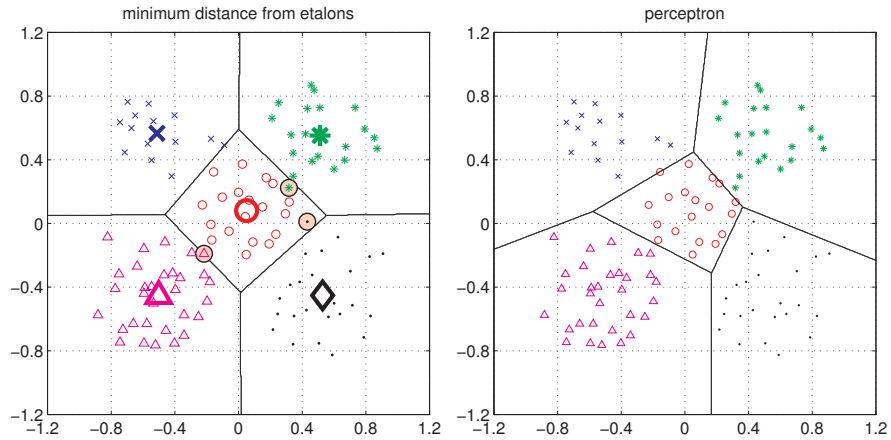
Searching for a projection of the data to minimize intra-class variance and maximize inter-class variance.

# Better etalons?



Figures from [5]

# Better etalons?



Figures from [5]

## Etalon classifier – Linear classifier

$$\begin{aligned}s^* &= \arg \min_{s \in S} (\|\vec{x} - \vec{e}_s\|^2 + o_s) = \arg \min_{s \in S} (\vec{x}^\top \vec{x} - 2 \vec{e}_s^\top \vec{x} + \vec{e}_s^\top \vec{e}_s + o_s) = \\&= \arg \min_{s \in S} \left( \vec{x}^\top \vec{x} - 2 \left( \vec{e}_s^\top \vec{x} - \frac{1}{2} (\vec{e}_s^\top \vec{e}_s + o_s) \right) \right) = \\&= \arg \min_{s \in S} (\vec{x}^\top \vec{x} - 2 (\vec{e}_s^\top \vec{x} + b_s)) = \\&= \boxed{\arg \max_{s \in S} (\vec{e}_s^\top \vec{x} + b_s)} = \arg \max_{s \in S} g_s(\vec{x}).\end{aligned}$$

$$b_s = -\frac{1}{2}(\vec{e}_s^\top \vec{e}_s + o_s)$$

Linear function (plus offset)

$$g_s(\mathbf{x}) = \mathbf{w}_s^\top \mathbf{x} + w_{s0}$$

The result is a *linear discriminant function* – hence etalon classifier is a linear classifier.

We classify into the class with highest value of the discriminant function.  $\mathbf{w}_s$  is a generalized etalon. How do we find it? Such that it is better than just the mean of the class members in the training set.

## Etalon classifier – Linear classifier

$$s^* = \arg \min_{s \in S} (\|\vec{x} - \vec{e}_s\|^2 + o_s) = \arg \min_{s \in S} (\vec{x}^\top \vec{x} - 2 \vec{e}_s^\top \vec{x} + \vec{e}_s^\top \vec{e}_s + o_s) =$$

$$= \arg \min_{s \in S} \left( \vec{x}^\top \vec{x} - 2 \left( \vec{e}_s^\top \vec{x} - \frac{1}{2} (\vec{e}_s^\top \vec{e}_s + o_s) \right) \right) =$$

$$= \arg \min_{s \in S} (\vec{x}^\top \vec{x} - 2 (\vec{e}_s^\top \vec{x} + b_s)) =$$

$$= \boxed{\arg \max_{s \in S} (\vec{e}_s^\top \vec{x} + b_s)} = \arg \max_{s \in S} g_s(\vec{x}).$$

$$b_s = -\frac{1}{2} (\vec{e}_s^\top \vec{e}_s + o_s)$$

Linear function (plus offset)

$$g_s(\mathbf{x}) = \mathbf{w}_s^\top \mathbf{x} + w_{s0}$$

The result is a *linear discriminant function* – hence etalon classifier is a linear classifier.

We classify into the class with highest value of the discriminant function.

$\mathbf{w}_s$  is a generalized etalon. How do we find it? Such that it is better than just the mean of the class members in the training set.

## Etalon classifier – Linear classifier

$$s^* = \arg \min_{s \in S} (\|\vec{x} - \vec{e}_s\|^2 + o_s) = \arg \min_{s \in S} (\vec{x}^\top \vec{x} - 2 \vec{e}_s^\top \vec{x} + \vec{e}_s^\top \vec{e}_s + o_s) =$$

$$= \arg \min_{s \in S} \left( \vec{x}^\top \vec{x} - 2 (\vec{e}_s^\top \vec{x} - \frac{1}{2} (\vec{e}_s^\top \vec{e}_s + o_s)) \right) =$$

$$= \arg \min_{s \in S} (\vec{x}^\top \vec{x} - 2 (\vec{e}_s^\top \vec{x} + b_s)) =$$

$$= \boxed{\arg \max_{s \in S} (\vec{e}_s^\top \vec{x} + b_s)} = \arg \max_{s \in S} g_s(\vec{x}).$$

$$b_s = -\frac{1}{2} (\vec{e}_s^\top \vec{e}_s + o_s)$$

Linear function (plus offset)

$$g_s(\mathbf{x}) = \mathbf{w}_s^\top \mathbf{x} + w_{s0}$$

The result is a *linear discriminant function* – hence etalon classifier is a linear classifier.

We classify into the class with highest value of the discriminant function.

$\mathbf{w}_s$  is a generalized etalon. How do we find it? Such that it is better than just the mean of the class members in the training set.

## Etalon classifier – Linear classifier

$$\begin{aligned}s^* &= \arg \min_{s \in S} (\|\vec{x} - \vec{e}_s\|^2 + o_s) = \arg \min_{s \in S} (\vec{x}^\top \vec{x} - 2 \vec{e}_s^\top \vec{x} + \vec{e}_s^\top \vec{e}_s + o_s) = \\&= \arg \min_{s \in S} \left( \vec{x}^\top \vec{x} - 2 \left( \vec{e}_s^\top \vec{x} - \frac{1}{2} (\vec{e}_s^\top \vec{e}_s + o_s) \right) \right) = \\&= \arg \min_{s \in S} (\vec{x}^\top \vec{x} - 2 (\vec{e}_s^\top \vec{x} + b_s)) = \\&= \boxed{\arg \max_{s \in S} (\vec{e}_s^\top \vec{x} + b_s)} = \arg \max_{s \in S} g_s(\vec{x}).\end{aligned}$$

Linear function (plus offset)

$$g_s(\mathbf{x}) = \mathbf{w}_s^\top \mathbf{x} + w_{s0}$$

The result is a *linear discriminant function* – hence etalon classifier is a linear classifier.

We classify into the class with highest value of the discriminant function.  $\mathbf{w}_s$  is a generalized etalon. How do we find it? Such that it is better than just the mean of the class members in the training set.

$$b_s = -\frac{1}{2}(\vec{e}_s^\top \vec{e}_s + o_s)$$

## Etalon classifier – Linear classifier

$$\begin{aligned}s^* &= \arg \min_{s \in S} (\|\vec{x} - \vec{e}_s\|^2 + o_s) = \arg \min_{s \in S} (\vec{x}^\top \vec{x} - 2 \vec{e}_s^\top \vec{x} + \vec{e}_s^\top \vec{e}_s + o_s) = \\&= \arg \min_{s \in S} \left( \vec{x}^\top \vec{x} - 2 \left( \vec{e}_s^\top \vec{x} - \frac{1}{2} (\vec{e}_s^\top \vec{e}_s + o_s) \right) \right) = \\&= \arg \min_{s \in S} (\vec{x}^\top \vec{x} - 2 (\vec{e}_s^\top \vec{x} + b_s)) = \\&= \boxed{\arg \max_{s \in S} (\vec{e}_s^\top \vec{x} + b_s)} = \arg \max_{s \in S} g_s(\vec{x}).\end{aligned}$$

$$b_s = -\frac{1}{2}(\vec{e}_s^\top \vec{e}_s + o_s)$$

Linear function (plus offset)

$$g_s(\mathbf{x}) = \mathbf{w}_s^\top \mathbf{x} + w_{s0}$$

The result is a *linear discriminant function* – hence etalon classifier is a linear classifier.

We classify into the class with highest value of the discriminant function.  $\mathbf{w}_s$  is a generalized etalon. How do we find it? Such that it is better than just the mean of the class members in the training set.



## Etalon classifier – Linear classifier

$$\begin{aligned}s^* &= \arg \min_{s \in S} (\|\vec{x} - \vec{e}_s\|^2 + o_s) = \arg \min_{s \in S} (\vec{x}^\top \vec{x} - 2 \vec{e}_s^\top \vec{x} + \vec{e}_s^\top \vec{e}_s + o_s) = \\&= \arg \min_{s \in S} \left( \vec{x}^\top \vec{x} - 2 (\vec{e}_s^\top \vec{x} - \frac{1}{2} (\vec{e}_s^\top \vec{e}_s + o_s)) \right) = \\&= \arg \min_{s \in S} (\vec{x}^\top \vec{x} - 2 (\vec{e}_s^\top \vec{x} + b_s)) = \\&= \boxed{\arg \max_{s \in S} (\vec{e}_s^\top \vec{x} + b_s)} = \arg \max_{s \in S} g_s(\vec{x}).\end{aligned}$$

$$b_s = -\frac{1}{2}(\vec{e}_s^\top \vec{e}_s + o_s)$$

Linear function (plus offset)

$$g_s(\mathbf{x}) = \mathbf{w}_s^\top \mathbf{x} + w_{s0}$$

The result is a *linear discriminant function* – hence etalon classifier is a linear classifier.

We classify into the class with highest value of the discriminant function.  $\mathbf{w}_s$  is a generalized etalon. How do we find it? Such that it is better than just the mean of the class members in the training set.

# Learning and decision

**Learning** stage - learning models/function/parameters from data.

**Decision** stage - decide about a query  $\vec{x}$ .

What to learn?

- ▶ **Generative model** : Learn  $P(\vec{x}, s)$ . Decide by computing  $P(s|\vec{x})$ .
- ▶ **Discriminative model** : Learn  $P(s|\vec{x})$
- ▶ **Discriminant function** : Learn  $g(\vec{x})$  which maps  $\vec{x}$  directly into class labels.

Generative models because by sampling from them it is possible to generate synthetic data points  $\vec{x}$ .

For the discriminative model one can consider, e.g. logistic function:

$$f(x) = \frac{1}{1 + e^{-k(x-x_0)}}$$

## Linear discriminant function - two class case

$$g(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + w_0$$

Decide  $s_1$  if  $g(\mathbf{x}) > 0$  and  $s_2$  if  $g(\mathbf{x}) < 0$

Think about 1-D, 2-D case, how does it generalize to  $d$ -D case?

$g(\mathbf{x}) = 0$  is the *separating hyperplane*. Its dimension is one less than that of the input space – for 2D space, it is a line. (This is a bit counterintuitive - “hyper” normally means above, more...)

What is the geometric meaning of the weight vector  $\mathbf{w}$ ?

## Linear discriminant function - two class case

$$g(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + w_0$$

Decide  $s_1$  if  $g(\mathbf{x}) > 0$  and  $s_2$  if  $g(\mathbf{x}) < 0$

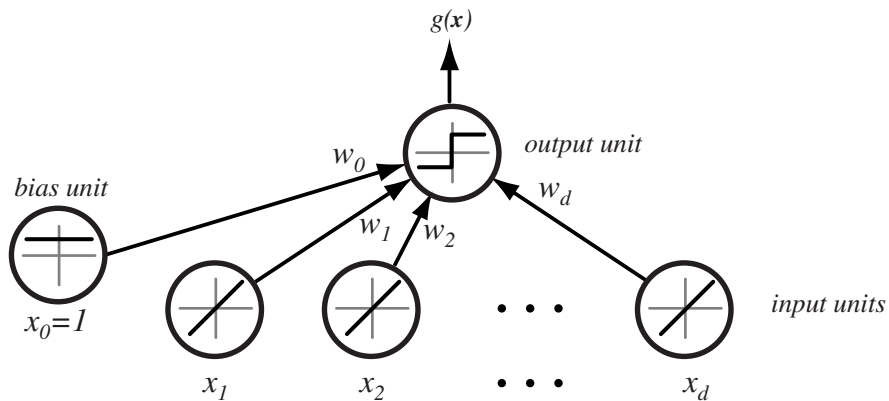


Figure from [2]

Think about 1-D, 2-D case, how does it generalize to  $d$ -D case?  
 $g(\mathbf{x}) = 0$  is the *separating hyperplane*. Its dimension is one less than that of the input space – for 2D space, it is a line. (This is a bit counterintuitive - “hyper” normally means above, more...) What is the geometric meaning of the weight vector  $\mathbf{w}$ ?

## Separating hyperplane

$$\mathbf{w}^\top \mathbf{x}_1 + w_0 = \mathbf{w}^\top \mathbf{x}_2 + w_0$$

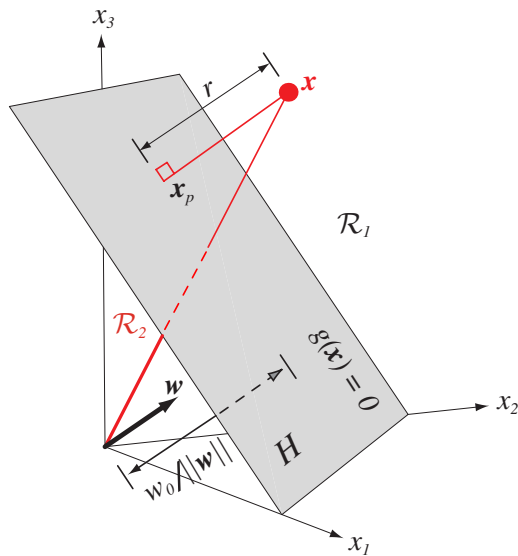
$$\mathbf{w}^\top (\mathbf{x}_1 - \mathbf{x}_2) = 0$$

$g(\mathbf{x})$  gives an algebraic measure of the distance from  $\mathbf{x}$  to the hyperplane.

$$\mathbf{x} = \mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}$$

as  $g(\mathbf{x}_p) = 0$ ,  
and  $g(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + w_0$ , then:

$$g(\mathbf{x}) = r \|\mathbf{w}\|$$



(any) vector  $(\mathbf{x}_1 - \mathbf{x}_2)$  lies on the separating hyperplane,  $\mathbf{w}$  is perpendicular to it

Summary: A linear discriminant function divides the feature space by a hyperplane decision surface.

- The orientation of the surface is determined by the normal vector  $\mathbf{w}$ .
- The location of the surface is determined by the bias term  $w_0$ .

## Separating hyperplane

$$\mathbf{w}^\top \mathbf{x}_1 + w_0 = \mathbf{w}^\top \mathbf{x}_2 + w_0$$

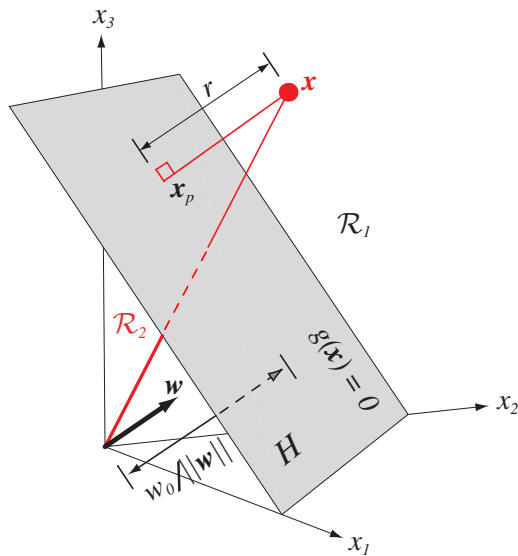
$$\mathbf{w}^\top (\mathbf{x}_1 - \mathbf{x}_2) = 0$$

$g(\mathbf{x})$  gives an algebraic measure of the distance from  $\mathbf{x}$  to the hyperplane.

$$\mathbf{x} = \mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}$$

as  $g(\mathbf{x}_p) = 0$ ,  
and  $g(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + w_0$ , then:

$$g(\mathbf{x}) = r \|\mathbf{w}\|$$



(any) vector  $(\mathbf{x}_1 - \mathbf{x}_2)$  lies on the separating hyperplane,  $\mathbf{w}$  is perpendicular to it

Summary: A linear discriminant function divides the feature space by a hyperplane decision surface.

- The orientation of the surface is determined by the normal vector  $\mathbf{w}$ .
- The location of the surface is determined by the bias term  $w_0$ .

## Separating hyperplane

$$\mathbf{w}^\top \mathbf{x}_1 + w_0 = \mathbf{w}^\top \mathbf{x}_2 + w_0$$

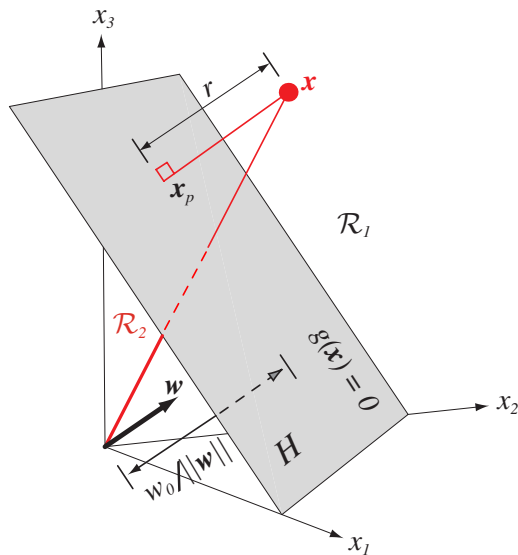
$$\mathbf{w}^\top (\mathbf{x}_1 - \mathbf{x}_2) = 0$$

$g(\mathbf{x})$  gives an algebraic measure of the distance from  $\mathbf{x}$  to the hyperplane.

$$\mathbf{x} = \mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}$$

as  $g(\mathbf{x}_p) = 0$ ,  
and  $g(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + w_0$ , then:

$$g(\mathbf{x}) = r \|\mathbf{w}\|$$



(any) vector  $(\mathbf{x}_1 - \mathbf{x}_2)$  lies on the separating hyperplane,  $\mathbf{w}$  is perpendicular to it

Summary: A linear discriminant function divides the feature space by a hyperplane decision surface.

- The orientation of the surface is determined by the normal vector  $\mathbf{w}$ .
- The location of the surface is determined by the bias term  $w_0$ .

## Multiclass case

each class has its own discriminant function

$$g_s(\mathbf{x}) = \mathbf{w}_s^\top \mathbf{x} + w_{s0}$$

and the classification  $s^*$  is along the max.

Show alternatives ([2], pg. 218, Section 5.2.2.)

- $c$  two-class problems, where the  $i$ th problem is solved by a linear discriminant function separating points assigned to  $w_i$  from the rest
- $c(c - 1)/2$  lin. discriminant functions – one for every pair of classes

In both cases, there are ambiguous regions.

How do the separating hyperplanes look like? Derive the equations?

What about the regions? Are they convex?



## Two classes set-up

$|S| = 2$ , i.e. two states (typically also classes)

$$g(\mathbf{x}) = \begin{cases} s = 1, & \text{if } \mathbf{w}^\top \mathbf{x} + w_0 > 0, \\ s = -1, & \text{if } \mathbf{w}^\top \mathbf{x} + w_0 < 0. \end{cases}$$

$$\mathbf{x}'_j = s_j \begin{bmatrix} 1 \\ \mathbf{x}_j \end{bmatrix}, \mathbf{w}' = \begin{bmatrix} w_0 \\ \mathbf{w} \end{bmatrix}$$

for all  $\mathbf{x}'$

$$\mathbf{w}'^\top \mathbf{x}' > 0$$

drop the dashes to avoid notation clutter.

There are two steps here:

1. Transformation to homogenous notation with augmented feature vector and augmented weight vector.
2. “Normalization” that simplifies treatment of the two-class case: labels can be ignored. Just look for a weight vector  $\mathbf{w}$  such that  $\mathbf{w}^\top \mathbf{x} > 0$

It means, the sign of  $\mathbf{x}$  depends on the class it belongs to! Keep in mind.

## Two classes set-up

$|S| = 2$ , i.e. two states (typically also classes)

$$g(\mathbf{x}) = \begin{cases} s = 1, & \text{if } \mathbf{w}^\top \mathbf{x} + w_0 > 0, \\ s = -1, & \text{if } \mathbf{w}^\top \mathbf{x} + w_0 < 0. \end{cases}$$

$$\mathbf{x}'_j = s_j \begin{bmatrix} 1 \\ \mathbf{x}_j \end{bmatrix}, \mathbf{w}' = \begin{bmatrix} w_0 \\ \mathbf{w} \end{bmatrix}$$

for all  $\mathbf{x}'$

$$\mathbf{w}'^\top \mathbf{x}' > 0$$

drop the dashes to avoid notation clutter.

There are two steps here:

1. Transformation to homogenous notation with augmented feature vector and augmented weight vector.
2. “Normalization” that simplifies treatment of the two-class case: labels can be ignored. Just look for a weight vector  $\mathbf{w}$  such that  $\mathbf{w}^\top \mathbf{x} > 0$

It means, the sign of  $\mathbf{x}$  depends on the class it belongs to! Keep in mind.

## Two classes set-up

$|S| = 2$ , i.e. two states (typically also classes)

$$g(\mathbf{x}) = \begin{cases} s = 1, & \text{if } \mathbf{w}^\top \mathbf{x} + w_0 > 0, \\ s = -1, & \text{if } \mathbf{w}^\top \mathbf{x} + w_0 < 0. \end{cases}$$

$$\mathbf{x}'_j = s_j \begin{bmatrix} 1 \\ \mathbf{x}_j \end{bmatrix}, \mathbf{w}' = \begin{bmatrix} w_0 \\ \mathbf{w} \end{bmatrix}$$

for all  $\mathbf{x}'$

$$\mathbf{w}'^\top \mathbf{x}' > 0$$

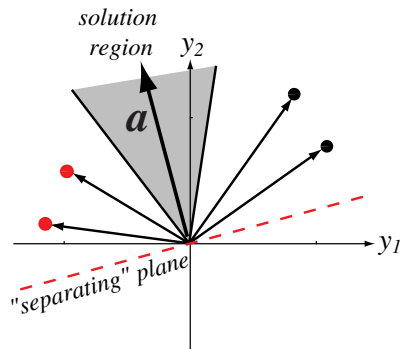
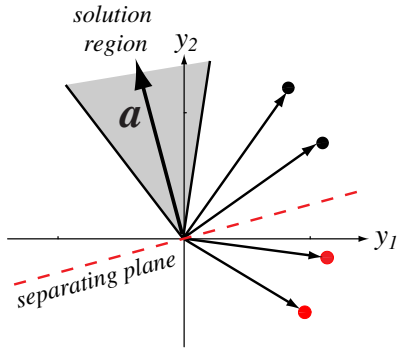
drop the dashes to avoid notation clutter.

There are two steps here:

1. Transformation to homogenous notation with augmented feature vector and augmented weight vector.
2. “Normalization” that simplifies treatment of the two-class case: labels can be ignored. Just look for a weight vector  $\mathbf{w}$  such that  $\mathbf{w}^\top \mathbf{x} > 0$

It means, the sign of  $\mathbf{x}$  depends on the class it belongs to! Keep in mind.

## Solution (graphically)



Four training samples (black for class/category  $w_1$ , red for  $w_2$ ). Left: Raw data Right: "Normalized data". Class  $w_2$  member replaced by their negatives... Simplifies the situation: labels can be ignored. Just look for a weight vector  $\mathbf{w}$  such that  $\mathbf{w}^\top \mathbf{x} > 0$

Before: defining the linear discriminant function.  
Now: How can we obtain it from (labeled) data?

Dif

notation in the book: substitute  $\mathbf{a} \leftarrow \mathbf{w}$  and  $y_1, y_2 \leftarrow x_1, x_2$

Figure from [2]

## Learning $\mathbf{w}$ , gradient descent

We're looking into *error-based classification* methods: missclassified examples are used to tune the classifier...

We already discussed (stochastic) Gradient descent when talking about  $Q$ -function learning

A criterion to be minimized  $J(\mathbf{w})$

Initialize  $\mathbf{w}$ , threshold  $\theta$ , learning rate  $\alpha$

$k \leftarrow 0$

**repeat**

$k \leftarrow k + 1$

$\mathbf{w} \leftarrow \mathbf{w} - \alpha(k) \nabla J(\mathbf{w})$

**until**  $|\alpha(k) \nabla J(\mathbf{w})| < \theta$

return  $\mathbf{w}$

## Learning $\mathbf{w}$ - Perceptron criterion

**Goal:** Find a weight vector  $\mathbf{w} \in \mathbb{R}^{D+1}$  (original feature space dimensionality is  $D$ ) such that:

$$\mathbf{w}^\top \mathbf{x}_j > 0 \quad (\forall j \in \{1, 2, \dots, m\})$$

(Perceptron) Criterion to be minimized:

$$J(\mathbf{w}) = \sum_{\mathbf{x} \in \mathcal{X}} -\mathbf{w}^\top \mathbf{x}$$

where  $\mathcal{X}$  is a set of misclassified  $\mathbf{x}$ .

$$\nabla J(\mathbf{w}) = \sum_{\mathbf{x} \in \mathcal{X}} -\mathbf{x}$$

What are the possible choices for  $J(\mathbf{w})$ ? First choice: number of misclassified examples. Problem: this function is piecewise constant.

Better choice: perceptron criterion function.

Mind that  $\mathbf{w}^\top \mathbf{x}_j \leq 0$  for  $\mathbf{x} \in \mathcal{X}$

Geometrically:  $J(\mathbf{w}) \propto$  sum of the distance of the misclassified samples to the decision boundary.

What is  $\nabla J(\mathbf{w})$  equal to?

## Learning $\mathbf{w}$ - Perceptron criterion

**Goal:** Find a weight vector  $\mathbf{w} \in \mathbb{R}^{D+1}$  (original feature space dimensionality is  $D$ ) such that:

$$\mathbf{w}^\top \mathbf{x}_j > 0 \quad (\forall j \in \{1, 2, \dots, m\})$$

(Perceptron) Criterion to be minimized:

$$J(\mathbf{w}) = \sum_{\mathbf{x} \in \mathcal{X}} -\mathbf{w}^\top \mathbf{x}$$

where  $\mathcal{X}$  is a set of misclassified  $\mathbf{x}$ .

$$\nabla J(\mathbf{w}) = \sum_{\mathbf{x} \in \mathcal{X}} -\mathbf{x}$$

What are the possible choices for  $J(\mathbf{w})$ ? First choice: number of misclassified examples. Problem: this function is piecewise constant.

Better choice: perceptron criterion function.

Mind that  $\mathbf{w}^\top \mathbf{x}_j \leq 0$  for  $\mathbf{x} \in \mathcal{X}$

Geometrically:  $J(\mathbf{w}) \propto$  sum of the distance of the misclassified samples to the decision boundary.

What is  $\nabla J(\mathbf{w})$  equal to?

## Learning $\mathbf{w}$ - Perceptron criterion

**Goal:** Find a weight vector  $\mathbf{w} \in \Re^{D+1}$  (original feature space dimensionality is  $D$ ) such that:

$$\mathbf{w}^\top \mathbf{x}_j > 0 \quad (\forall j \in \{1, 2, \dots, m\})$$

(Perceptron) Criterion to be minimized:

$$J(\mathbf{w}) = \sum_{\mathbf{x} \in \mathcal{X}} -\mathbf{w}^\top \mathbf{x}$$

where  $\mathcal{X}$  is a set of misclassified  $\mathbf{x}$ .

$$\nabla J(\mathbf{w}) = \sum_{\mathbf{x} \in \mathcal{X}} -\mathbf{x}$$

What are the possible choices for  $J(\mathbf{w})$ ? First choice: number of misclassified examples. Problem: this function is piecewise constant. Better choice: perceptron criterion function.

Mind that  $\mathbf{w}^\top \mathbf{x}_j \leq 0$  for  $\mathbf{x} \in \mathcal{X}$

Geometrically:  $J(\mathbf{w}) \propto$  sum of the distance of the misclassified samples to the decision boundary.

What is  $\nabla J(\mathbf{w})$  equal to?



## (Batch) Perceptron algorithm

Initialize  $\mathbf{w}$ , threshold  $\theta$ , learning rate  $\alpha$

$k \leftarrow 0$

**repeat**

$k \leftarrow k + 1$

$\mathbf{w} \leftarrow \mathbf{w} + \alpha(k) \sum_{\mathbf{x} \in \mathcal{X}(k)} \mathbf{x}$

**until**  $|\alpha(k) \sum_{\mathbf{x} \in \mathcal{X}(k)} \mathbf{x}| < \theta$

return  $\mathbf{w}$

Next weight vector  $\sim$  adding some multiple of the sum of the misclassified samples to the present weight vector.

## Fixed-increment single-sample Perceptron

As we are looping over all patterns repeatedly, it is not an on-line algorithm

$n$  patterns/samples, we are looping over all patterns repeatedly

Initialize  $\mathbf{w}$

$k \leftarrow 0$

**repeat**

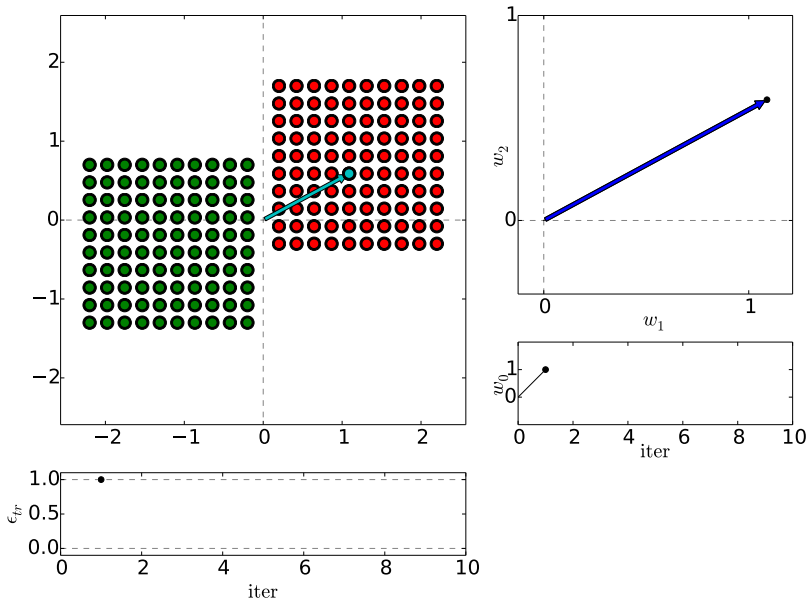
$k \leftarrow (k + 1) \bmod n$

**if**  $\mathbf{x}^k$  missclassified, **then**  $\mathbf{w} \leftarrow \mathbf{w} + \mathbf{x}^k$

**until** all  $\mathbf{x}$  correctly classified

return  $\mathbf{w}$

## Perceptron iterations/loops



Keep in mind the  $\pm$  normalization of  $\mathbf{x}$ .

$$g(\mathbf{x}) = \begin{cases} s = 1, & \text{if } \mathbf{w}^\top \mathbf{x} + w_0 > 0, \\ s = -1, & \text{if } \mathbf{w}^\top \mathbf{x} + w_0 < 0. \end{cases}$$

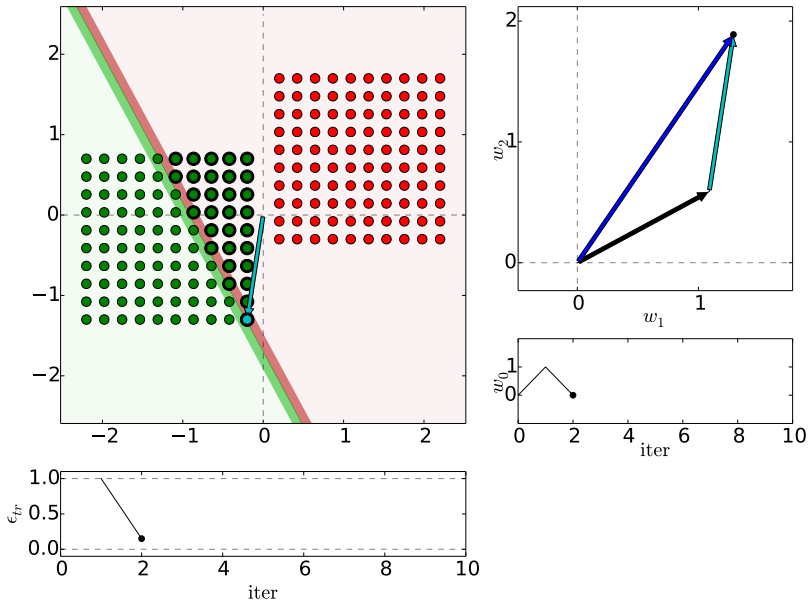
$$\mathbf{x}'_j = s_j \begin{bmatrix} 1 \\ \mathbf{x}_j \end{bmatrix}, \mathbf{w}' = \begin{bmatrix} w_0 \\ \mathbf{w} \end{bmatrix}$$

(as discussed few slides ago)

Red  $\mathbf{x}$  are  $+$ , green are  $-$

Track the iteration steps. After each update  $\mathbf{x}$ , draw a separating line for the next and verify.

## Perceptron iterations/loops



Keep in mind the  $\pm$  normalization of  $\mathbf{x}$ .

$$g(\mathbf{x}) = \begin{cases} s = 1, & \text{if } \mathbf{w}^\top \mathbf{x} + w_0 > 0, \\ s = -1, & \text{if } \mathbf{w}^\top \mathbf{x} + w_0 < 0. \end{cases}$$

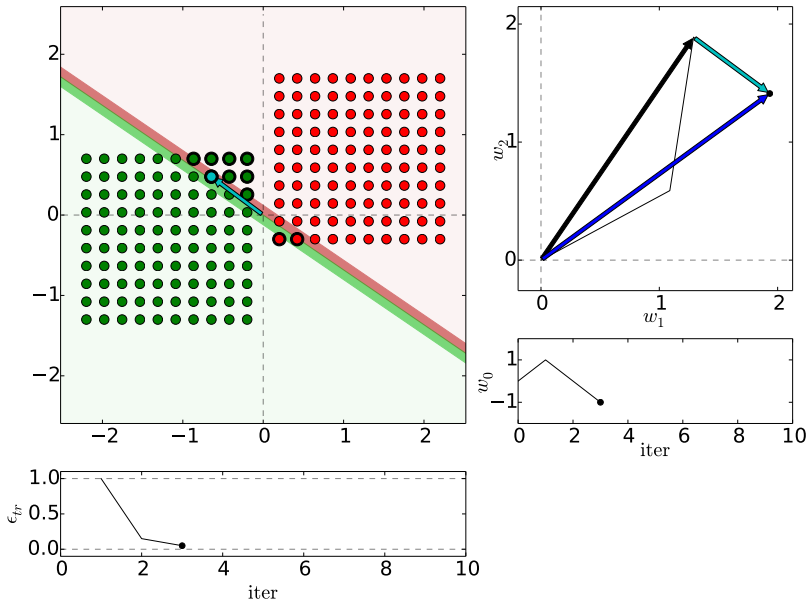
$$\mathbf{x}'_j = s_j \begin{bmatrix} 1 \\ \mathbf{x}_j \end{bmatrix}, \mathbf{w}' = \begin{bmatrix} w_0 \\ \mathbf{w} \end{bmatrix}$$

(as discussed few slides ago)

Red  $\mathbf{x}$  are +, green are -

Track the iteration steps. After each update  $\mathbf{x}$ , draw a separating line for the next and verify.

## Perceptron iterations/loops



Keep in mind the  $\pm$  normalization of  $\mathbf{x}$ .

$$g(\mathbf{x}) = \begin{cases} s = 1, & \text{if } \mathbf{w}^\top \mathbf{x} + w_0 > 0, \\ s = -1, & \text{if } \mathbf{w}^\top \mathbf{x} + w_0 < 0. \end{cases}$$

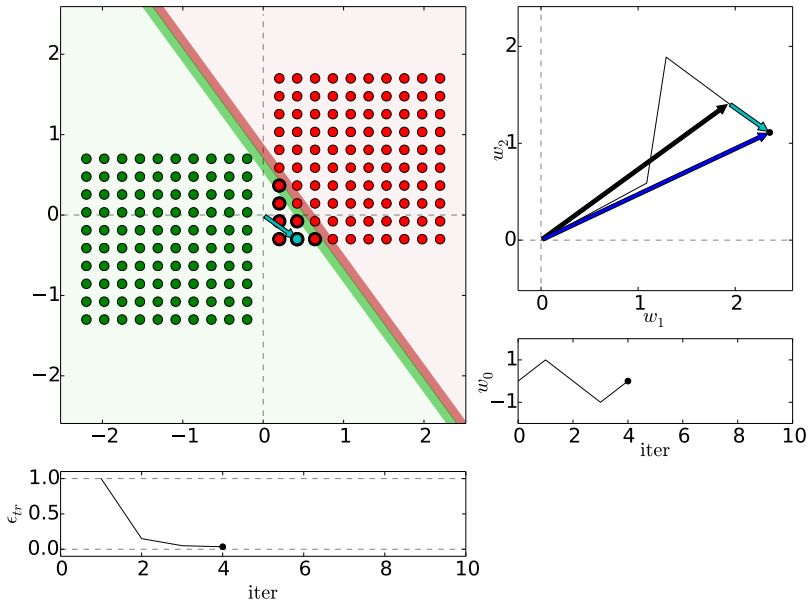
$$\mathbf{x}'_j = s_j \begin{bmatrix} 1 \\ \mathbf{x}_j \end{bmatrix}, \mathbf{w}' = \begin{bmatrix} w_0 \\ \mathbf{w} \end{bmatrix}$$

(as discussed few slides ago)

Red  $\mathbf{x}$  are  $+$ , green are  $-$

Track the iteration steps. After each update  $\mathbf{x}$ , draw a separating line for the next and verify.

## Perceptron iterations/loops



Keep in mind the  $\pm$  normalization of  $\mathbf{x}$ .

$$g(\mathbf{x}) = \begin{cases} s = 1, & \text{if } \mathbf{w}^\top \mathbf{x} + w_0 > 0, \\ s = -1, & \text{if } \mathbf{w}^\top \mathbf{x} + w_0 < 0. \end{cases}$$

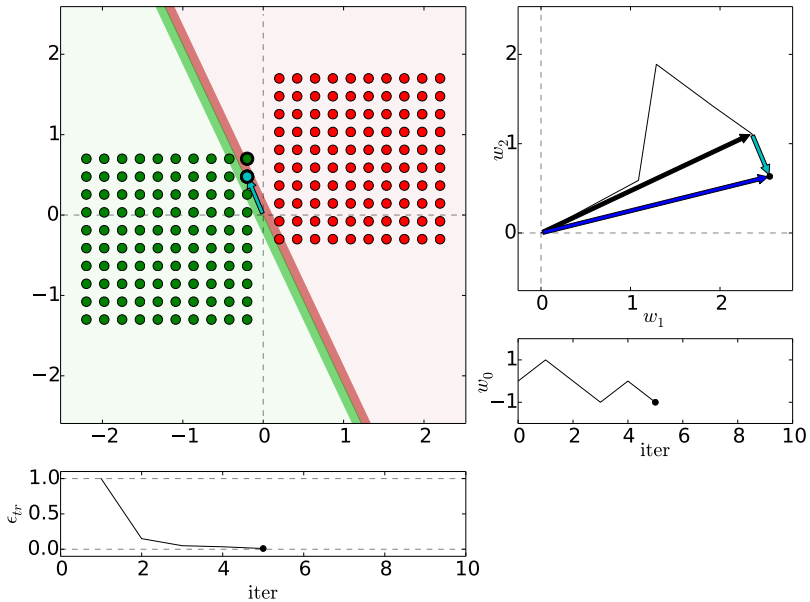
$$\mathbf{x}'_j = s_j \begin{bmatrix} 1 \\ \mathbf{x}_j \end{bmatrix}, \mathbf{w}' = \begin{bmatrix} w_0 \\ \mathbf{w} \end{bmatrix}$$

(as discussed few slides ago)

Red  $\mathbf{x}$  are  $+$ , green are  $-$

Track the iteration steps. After each update  $\mathbf{x}$ , draw a separating line for the next and verify.

## Perceptron iterations/loops



Keep in mind the  $\pm$  normalization of  $\mathbf{x}$ .

$$g(\mathbf{x}) = \begin{cases} s = 1, & \text{if } \mathbf{w}^\top \mathbf{x} + w_0 > 0, \\ s = -1, & \text{if } \mathbf{w}^\top \mathbf{x} + w_0 < 0. \end{cases}$$

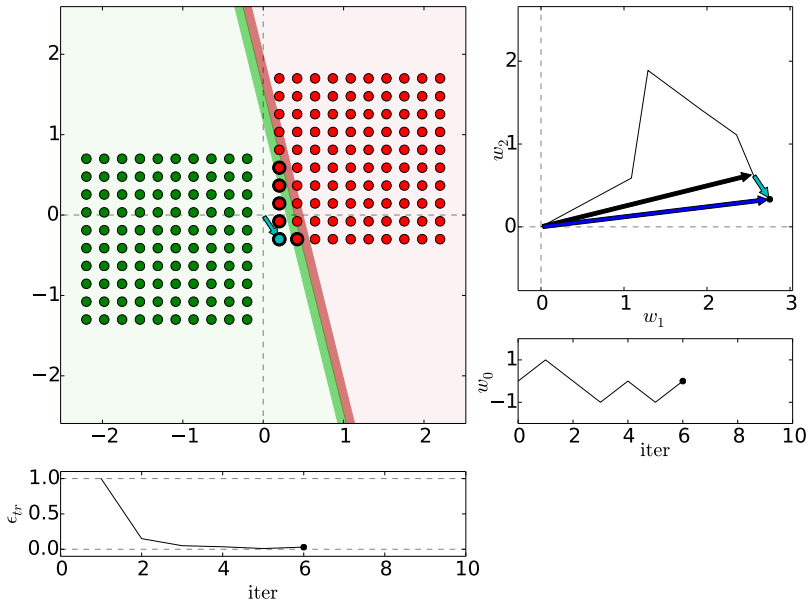
$$\mathbf{x}'_j = s_j \begin{bmatrix} 1 \\ \mathbf{x}_j \end{bmatrix}, \mathbf{w}' = \begin{bmatrix} w_0 \\ \mathbf{w} \end{bmatrix}$$

(as discussed few slides ago)

Red  $\mathbf{x}$  are +, green are -

Track the iteration steps. After each update  $\mathbf{x}$ , draw a separating line for the next and verify.

## Perceptron iterations/loops



Keep in mind the  $\pm$  normalization of  $\mathbf{x}$ .

$$g(\mathbf{x}) = \begin{cases} s = 1, & \text{if } \mathbf{w}^\top \mathbf{x} + w_0 > 0, \\ s = -1, & \text{if } \mathbf{w}^\top \mathbf{x} + w_0 < 0. \end{cases}$$

$$\mathbf{x}'_j = s_j \begin{bmatrix} 1 \\ \mathbf{x}_j \end{bmatrix}, \mathbf{w}' = \begin{bmatrix} w_0 \\ \mathbf{w} \end{bmatrix}$$

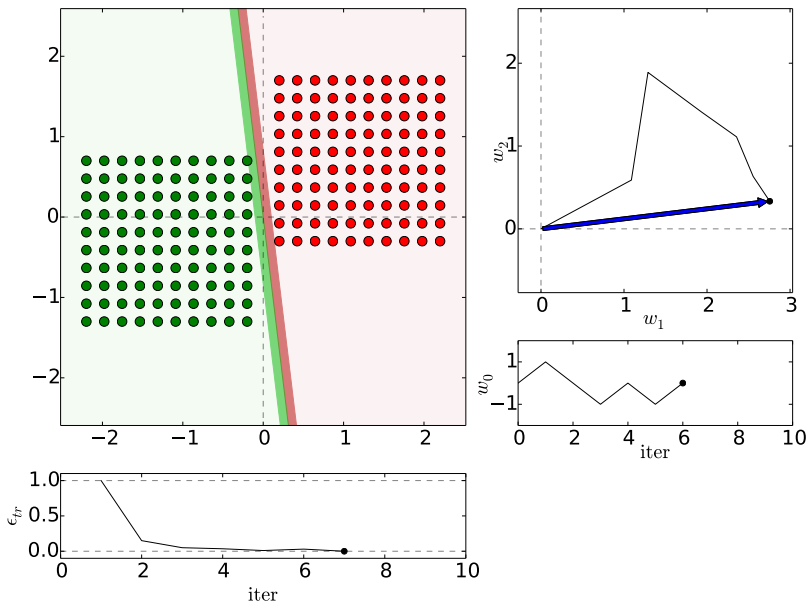
(as discussed few slides ago)

Red  $\mathbf{x}$  are  $+$ , green are  $-$

Track the iteration steps. After each update  $\mathbf{x}$ , draw a separating line for the next and verify.



## Perceptron iterations/loops



Keep in mind the  $\pm$  normalization of  $\mathbf{x}$ .

$$g(\mathbf{x}) = \begin{cases} s = 1, & \text{if } \mathbf{w}^\top \mathbf{x} + w_0 > 0, \\ s = -1, & \text{if } \mathbf{w}^\top \mathbf{x} + w_0 < 0. \end{cases}$$

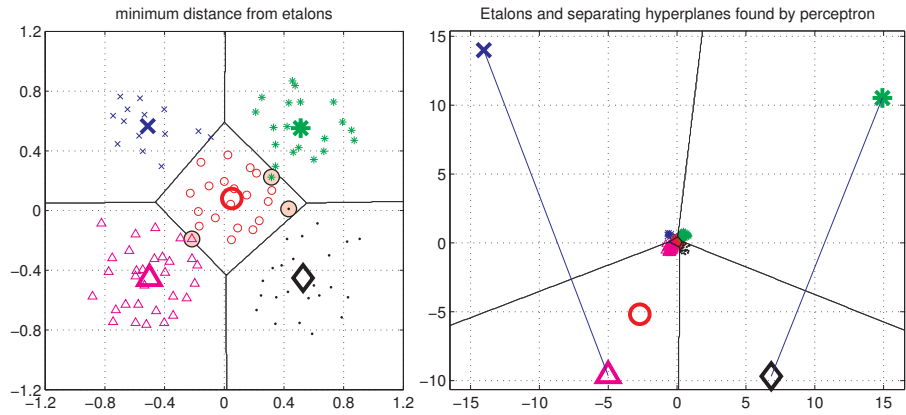
$$\mathbf{x}'_j = s_j \begin{bmatrix} 1 \\ \mathbf{x}_j \end{bmatrix}, \mathbf{w}' = \begin{bmatrix} w_0 \\ \mathbf{w} \end{bmatrix}$$

(as discussed few slides ago)

Red  $\mathbf{x}$  are  $+$ , green are  $-$

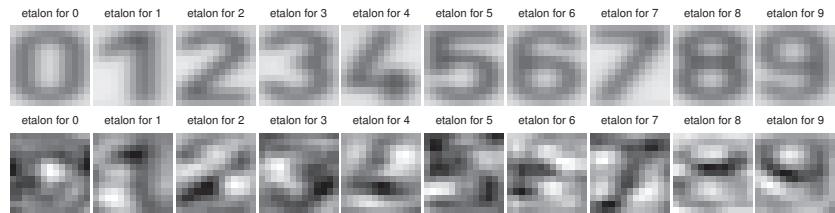
Track the iteration steps. After each update  $\mathbf{x}$ , draw a separating line for the next and verify.

# Etalons: means vs. found by perceptron



Figures from [5]

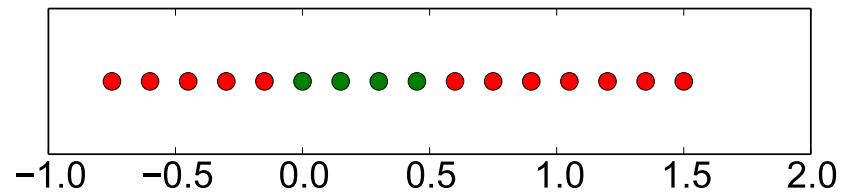
# Digit recognition - etalons means vs. perceptron



“Prototypes” resulting from the perceptron algorithm are harder to interpret because they are not means – instead, they are optimized for separating the classes.

Figures from [5]

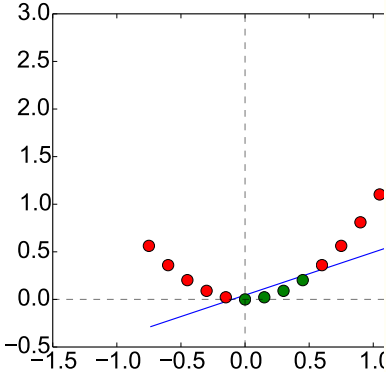
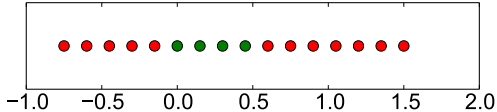
What if not lin separable?



Dimension lifting

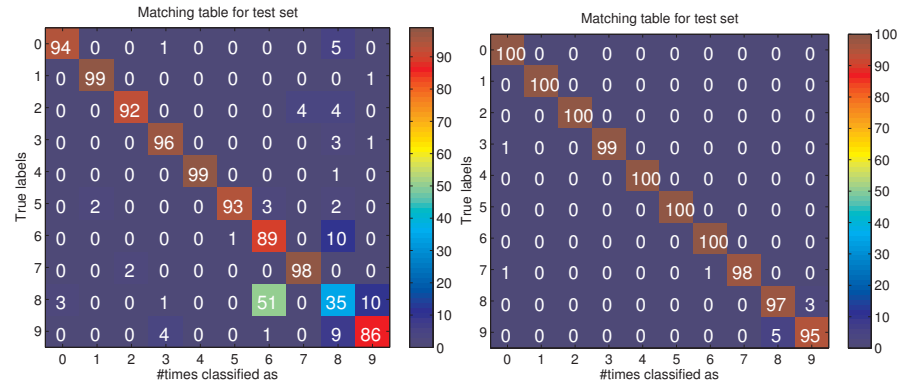
$$\mathbf{x} = [x, x^2]^T$$

Dimension lifting,  $\mathbf{x} = [x, x^2]^\top$



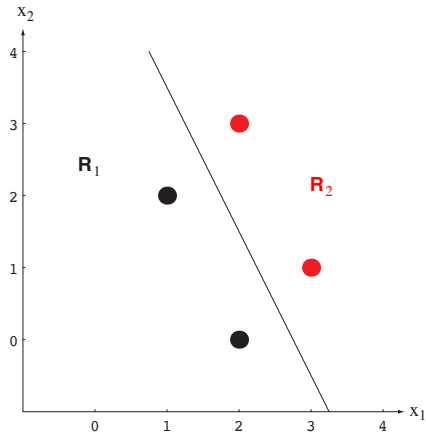
# Performance comparison, parameters fixed

Why there some errors in perceptron results? We said zero error on training set.



## LSQ approach to linear classification

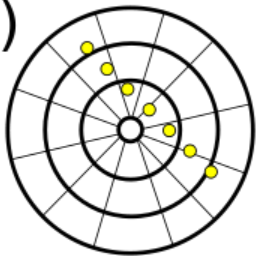
$$X\mathbf{w} = \mathbf{b}$$



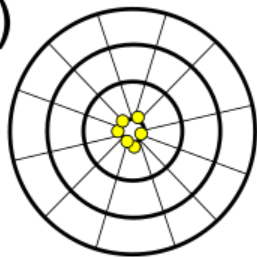
Linear least squares not guaranteed to correctly classify everything on the training set. It's objective function not perfect for classification.  
Outliers can shift the decision boundary.

## Accuracy vs precision

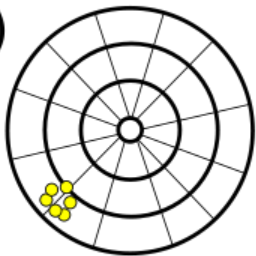
(a)



(b)



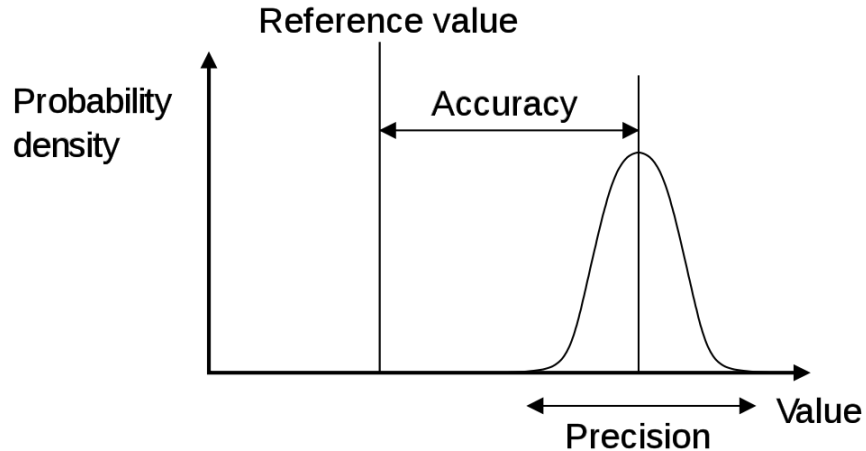
(c)



Accuracy: how close (is your model) to the truth. Precision: how consistent/stable

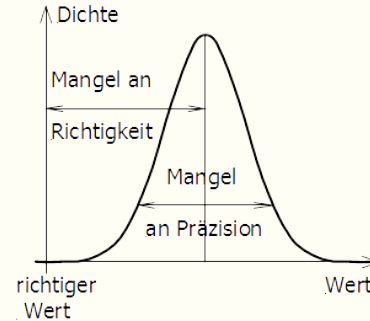


## Accuracy vs precision



Accuracy: how close (is your model) to the truth. Precision: how consistent/stable.

Think about terms *bias* and *error*. I



[https://en.wikipedia.org/wiki/Accuracy\\_and\\_precision](https://en.wikipedia.org/wiki/Accuracy_and_precision)

## References I

Further reading: Chapter 18 of [4], or chapter 4 of [1], or chapter 5 of [2].  
Many Matlab figures created with the help of [3]. You may also play with  
demo functions from [5].

[1] Christopher M. Bishop.

*Pattern Recognition and Machine Learning.*

Springer Science+Bussiness Media, New York, NY, 2006.

PDF freely downloadable.

[2] Richard O. Duda, Peter E. Hart, and David G. Stork.

*Pattern Classification.*

John Wiley & Sons, 2nd edition, 2001.

[3] Votjěch Franc and Václav Hlaváč.

Statistical pattern recognition toolbox.

<http://cmp.felk.cvut.cz/cmp/software/stprtool/index.html>.

## References II

- [4] Stuart Russell and Peter Norvig.

*Artificial Intelligence: A Modern Approach.*

Prentice Hall, 3rd edition, 2010.

<http://aima.cs.berkeley.edu/>.

- [5] Tomáš Svoboda, Jan Kybic, and Hlaváč Václav.

*Image Processing, Analysis and Machine Vision — A MATLAB Companion.*

Thomson, Toronto, Canada, 1<sup>st</sup> edition, September 2007.

<http://visionbook.felk.cvut.cz/>.