

# Linear Classifiers II, and Tangent Space for $k - NN$

Tomáš Svoboda

Vision for Robots and Autonomous Systems, Center for Machine Perception  
Department of Cybernetics  
Faculty of Electrical Engineering, Czech Technical University in Prague

May 23, 2022

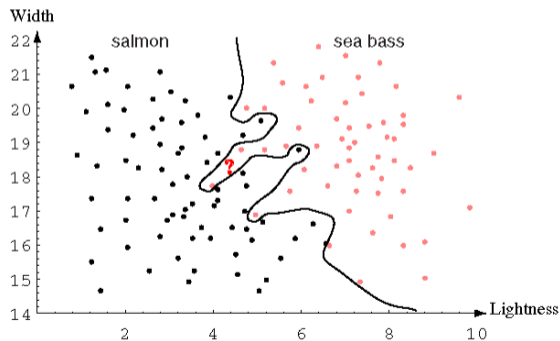
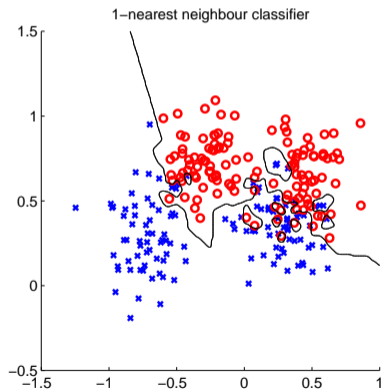
# Outline

- ▶  $k$  – NN, Tangent distance measure, invariance to rotation
- ▶ Better etalons by applying Fischer linear discriminator analysis.
- ▶ LSQ formulation of the learning task.

# $K$ -Nearest neighbors classification

For a query  $\mathbf{x}$ :

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



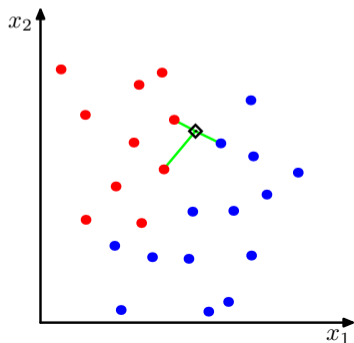
## $K$ - Nearest Neighbor and Bayes $j^* = \operatorname{argmax}_j P(s_j|\mathbf{x})$

Assume data:

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

We want to classify  $\mathbf{x}$ . Draw a sphere centered at  $\mathbf{x}$  containing  $K$  points irrespective of class.

$V$  is the volume of this sphere.  $P(s_j|\mathbf{x}) = ?$



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

$K_j$  is the number of points of class  $s_j$  among the  $K$  nearest neighbors.

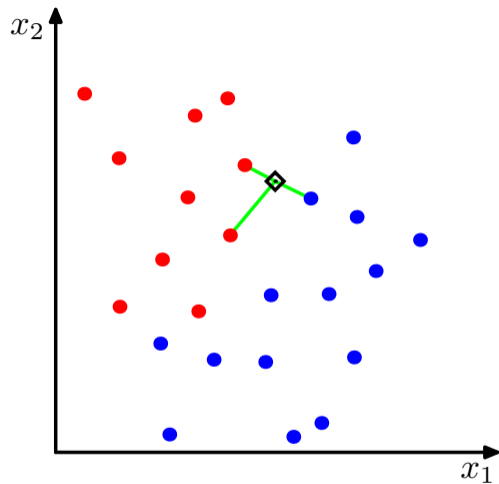
$$P(s_j) = \frac{N_j}{N}$$

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

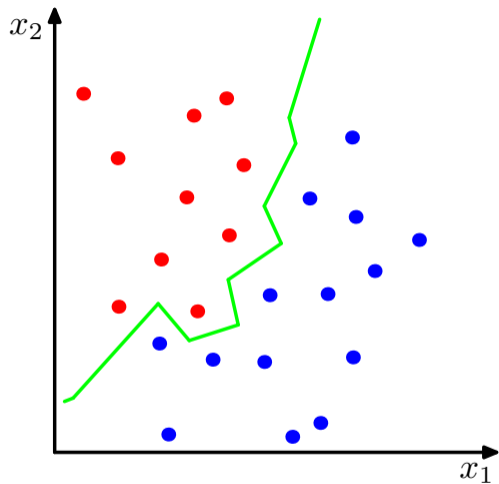
$$P(\mathbf{x}|s_j) = \frac{K_j}{N_j V}$$

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

# NN classification example



(a)

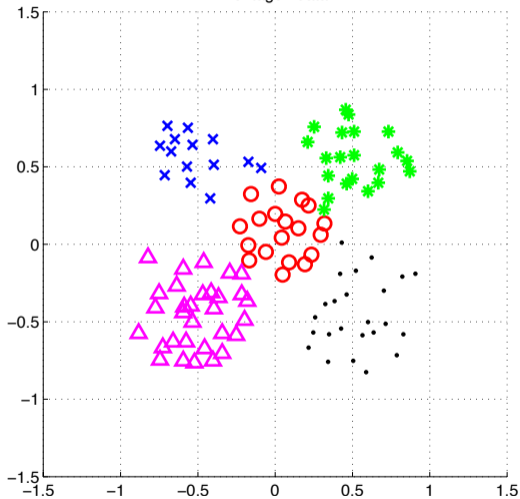


(b)

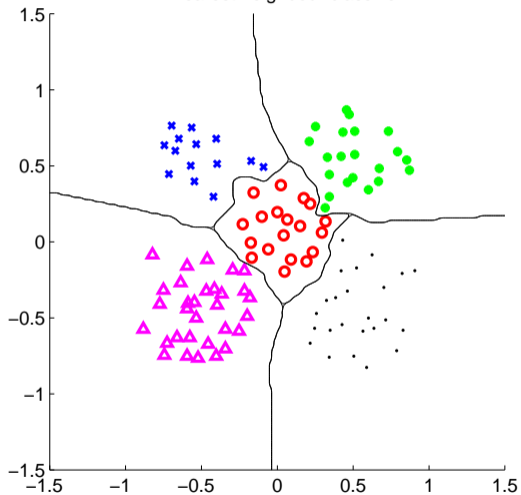
<sup>1</sup>Figs from [1]

# NN classification example

Pentagon data



1-nearest neighbour classifier



## What is *nearest*? Metrics for NN classification ...

**Metrics** : a function  $D$  which is

- ▶ nonnegative,
- ▶ reflexive,
- ▶ symmetrical,
- ▶ satisfying triangle inequality:

$$D(\mathbf{x}_1, \mathbf{x}_2) \geq 0$$

$$D(\mathbf{x}_1, \mathbf{x}_2) = 0 \text{ iff } \mathbf{x}_1 = \mathbf{x}_2$$

$$D(\mathbf{x}_1, \mathbf{x}_2) = D(\mathbf{x}_2, \mathbf{x}_1)$$

$$D(\mathbf{x}_1, \mathbf{x}_2) + D(\mathbf{x}_2, \mathbf{x}_3) \geq D(\mathbf{x}_1, \mathbf{x}_3)$$

$D(\mathbf{x}_1, \mathbf{x}_2) = \|\mathbf{x}_1 - \mathbf{x}_2\|$  just fine, but

...

## What is *nearest*? Metrics for NN classification ...

**Metrics** : a function  $D$  which is

- ▶ nonnegative,
- ▶ reflexive,
- ▶ symmetrical,
- ▶ satisfying triangle inequality:

$$D(\mathbf{x}_1, \mathbf{x}_2) \geq 0$$

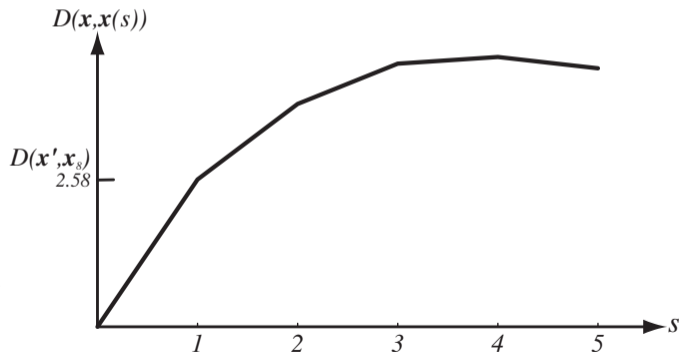
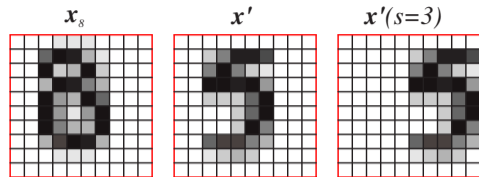
$$D(\mathbf{x}_1, \mathbf{x}_2) = 0 \text{ iff } \mathbf{x}_1 = \mathbf{x}_2$$

$$D(\mathbf{x}_1, \mathbf{x}_2) = D(\mathbf{x}_2, \mathbf{x}_1)$$

$$D(\mathbf{x}_1, \mathbf{x}_2) + D(\mathbf{x}_2, \mathbf{x}_3) \geq D(\mathbf{x}_1, \mathbf{x}_3)$$

$$D(\mathbf{x}_1, \mathbf{x}_2) = \|\mathbf{x}_1 - \mathbf{x}_2\| \text{ just fine, but}$$

...



Invariance to geometrical transformations? (figure from [2]) 7 / 21



# Tangent space

Consider continuous transformation:  
e.g. rotation or translation not mirror reflection.

$\mathbf{x} = [x_1, x_2]^T$  move along manifold  $\mathcal{M}$   
 $\alpha$  is a transformation parameter (e.g. angle)  
Tangent vector  $\boldsymbol{\tau}$  is a linearization

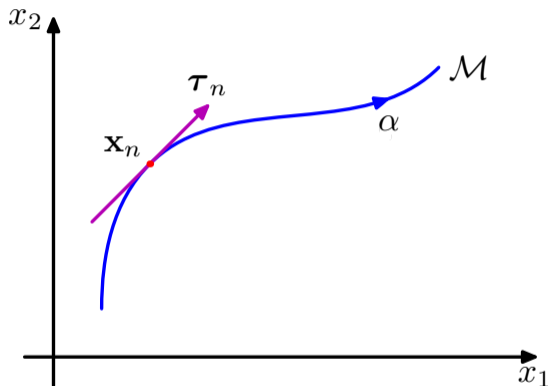
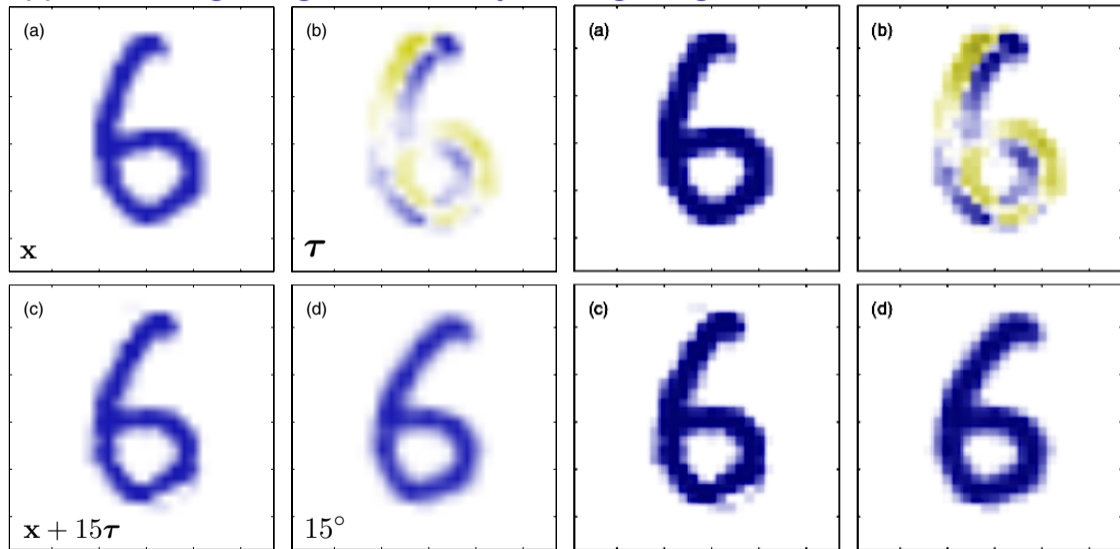


Figure from [1], slightly adapted

## Approximating image rotation by adding tangent vector



Figures from [1], slightly adapted.

## Combining more transformations

Approximate derivative by difference.

For all exemplars  $\mathbf{x}'$

and all  $r$  transformations  $\mathcal{F}_i$

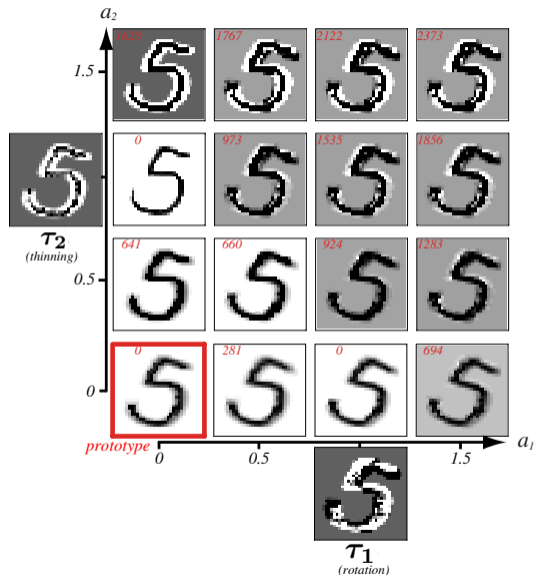
$$\blacktriangleright \tau_i = \mathcal{F}_i(\mathbf{x}', \alpha_i) - \mathbf{x}'$$

For each exemplar we have  $d \times r$  matrix  $\mathbf{T}$

$$\mathbf{T} = [\tau_1, \tau_2, \dots, \tau_r]$$

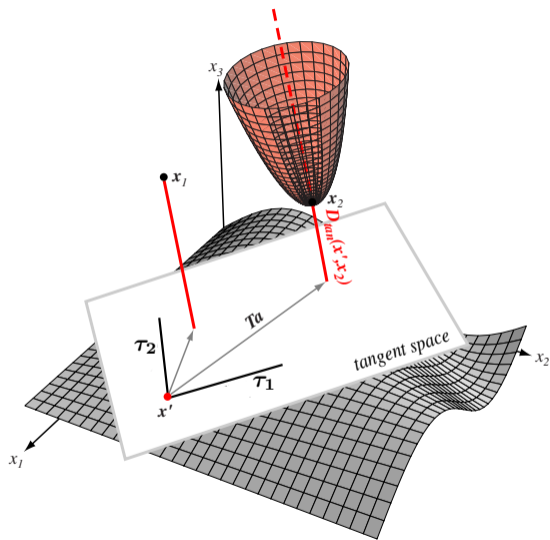
Grouping coefficients  $\mathbf{a} = [a_1, a_2]^T$

Right image visualizes  $\mathbf{x}' + \mathbf{T}\mathbf{a}$



Figures from [2], slightly adapted.

## Minimizing distance to tangent space



$$D_{tan}(\mathbf{x}', \mathbf{x}) = \min_{\mathbf{a}} \|(\mathbf{x}' + T\mathbf{a}) - \mathbf{x}\|$$

Gradient descent will do.

Figures from [2], slightly adapted.

## Linear classifiers II

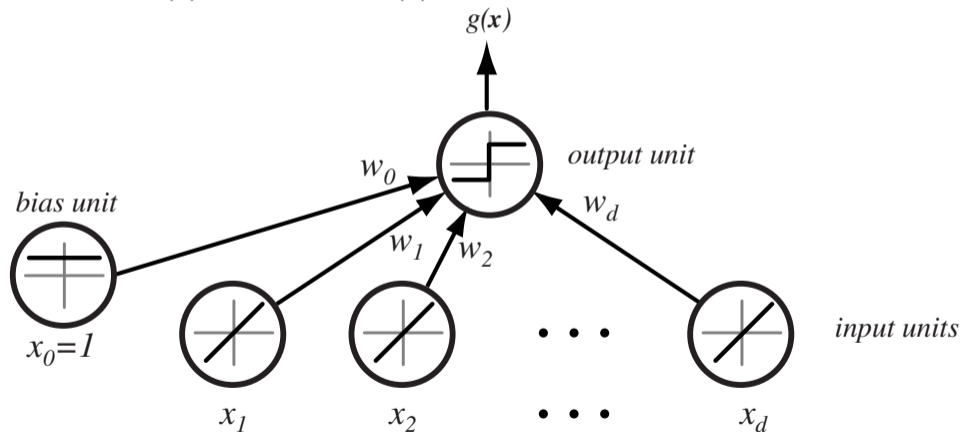
$$g(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + w_0$$

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

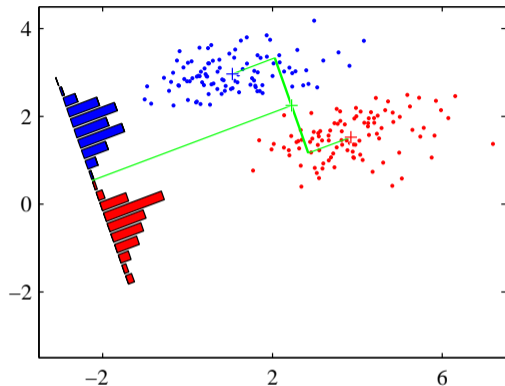
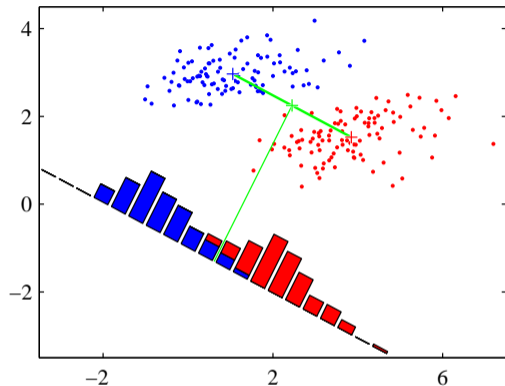
## Linear classifiers II

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

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



## Fischer linear discriminant



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

Figures from [1]

# Projections to lower dimensions $y = \mathbf{w}^\top \mathbf{x}$

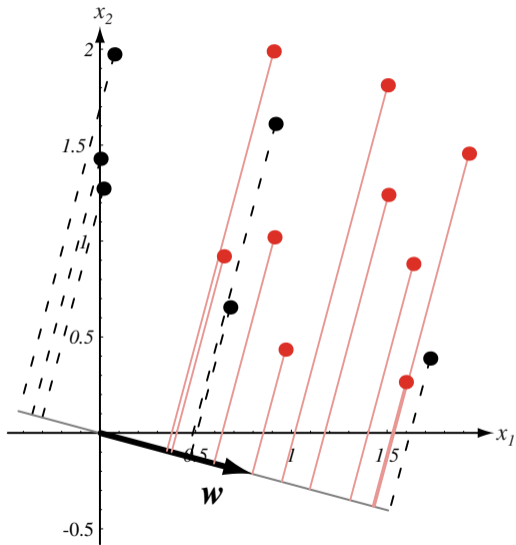
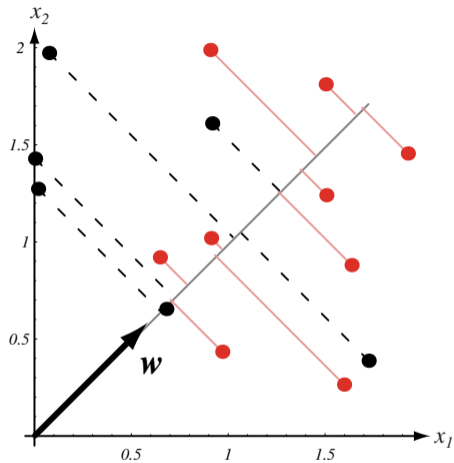


Figure from [2]



Projection to lower dimension  $\mathbf{y} = \mathbf{W}^T \mathbf{x}$

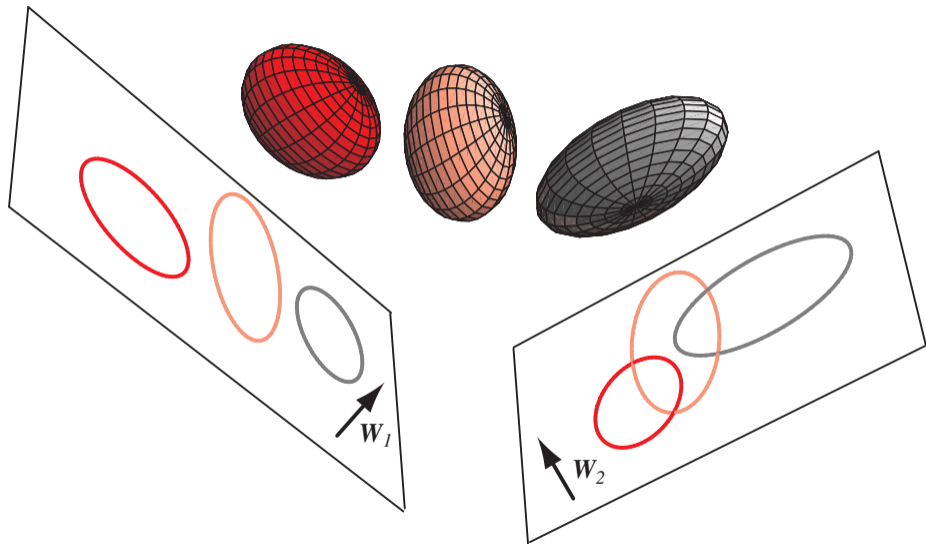
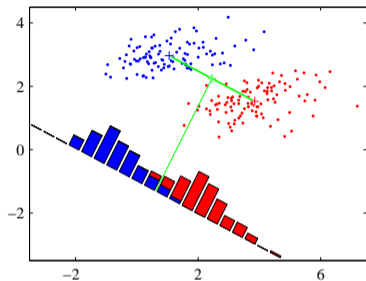
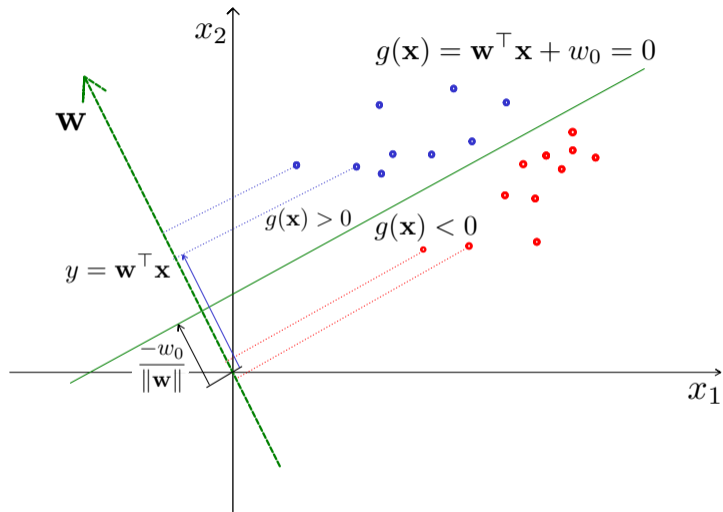


Figure from [2]

Finding the best projection  $y = \mathbf{w}^\top \mathbf{x}$ ,  $y \geq -w_0 \Rightarrow C_1$ , otherwise  $C_2$



Finding the best projection  $y = \mathbf{w}^\top \mathbf{x}$ ,  $y \geq -w_0 \Rightarrow C_1$ , otherwise  $C_2$

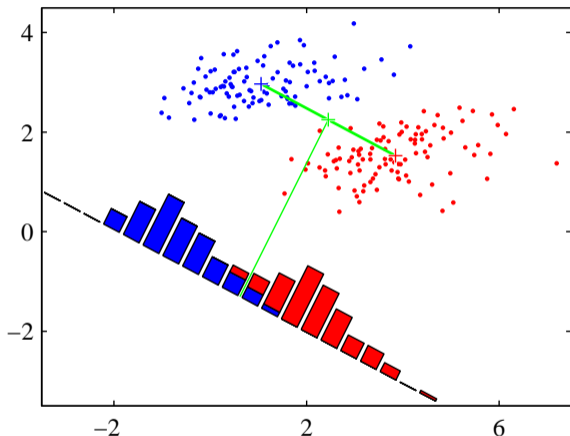
$$m_2 - m_1 = \mathbf{w}^\top (\mathbf{m}_2 - \mathbf{m}_1)$$

Within class scatter of projected samples

$$s_i^2 = \sum_{y \in C_i} (y - m_i)^2$$

Fischer criterion:

$$J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2}$$



Finding the best projection  $y = \mathbf{w}^\top \mathbf{x}$ ,  $y \geq -w_0 \Rightarrow C_1$ , otherwise  $C_2$

$$m_2 - m_1 = \mathbf{w}^\top (\mathbf{m}_2 - \mathbf{m}_1)$$

$$s_i^2 = \sum_{y \in C_i} (y - m_i)^2$$

$$J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2}$$

$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = 0$$

$$S_i = \sum_{\mathbf{x} \in C_i} (\mathbf{x} - \mathbf{m}_i)(\mathbf{x} - \mathbf{m}_i)^\top$$

$$S_W = S_1 + S_2$$

$$S_B = (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^\top$$

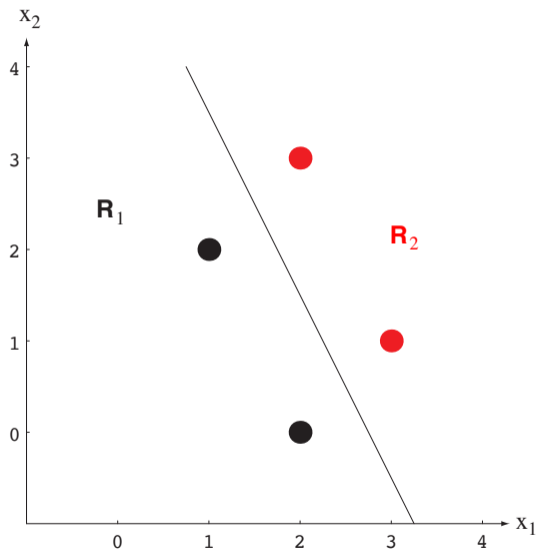
$$J(\mathbf{w}) = \frac{\mathbf{w}^\top S_B \mathbf{w}}{\mathbf{w}^\top S_W \mathbf{w}}$$

# LSQ approach to linear classification

$$\mathbf{w} = \begin{bmatrix} w_0 \\ \mathbf{w} \end{bmatrix}$$

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

$$J(\mathbf{w}) = \|X\mathbf{w} - \mathbf{b}\|^2$$



LSQ approach, better margins **b**?

$$\mathbf{X} = \begin{bmatrix} 1_1 & X_1 \\ -1_2 & -X_2 \end{bmatrix}$$

$$\mathbf{b} = \begin{bmatrix} \frac{n}{n_1} 1_1 \\ \frac{n}{n_2} 1_2 \end{bmatrix}$$

## References I

Further reading: Chapter 4 of [1], or chapter 3 and 5 of [2].

[1] Christopher M. Bishop.

*Pattern Recognition and Machine Learning.*

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

<https://www.microsoft.com/en-us/research/uploads/prod/2006/01/Bishop-Pattern-Recognition-and-Machine-Learning-2006.pdf>.

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

*Pattern Classification.*

John Wiley & Sons, 2nd edition, 2001.