



# Machine Learning 1

Lecture 6.4 - Supervised Learning  
Classification - Discriminative Models - Least  
Squares Regression

*Erik Bekkers*

*(Bishop 4.1.3)*



# Least Squares for Classification

(K)

- ▶ Each class  $C_k$  has its own linear model:

$$y_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0}$$

- ▶ Shorter notation:  $\mathbf{y}(\mathbf{x}) = \widetilde{\mathbf{W}}^T \tilde{\mathbf{x}}$

- ▶ Matrix  $\widetilde{\mathbf{W}}$  : column  $k$  contains  $\tilde{\mathbf{w}}_k = (w_{k0}, \underline{w})^T \in \mathbb{R}^m$

- ▶ Vector  $\tilde{\mathbf{x}} = (1, \underline{x})^T \in \mathbb{R}^m$

- ▶ Vector  $\mathbf{y}(\mathbf{x}) = \begin{pmatrix} y_1(\underline{x}) \\ y_2(\underline{x}) \\ \vdots \\ y_k(\underline{x}) \end{pmatrix} = \widetilde{\mathbf{W}}^T \tilde{\mathbf{x}} \in \mathbb{R}^k$

- ▶ Assign  $\mathbf{x}$  to class  $C_k$  if

$$k = \underset{j}{\operatorname{argmax}} y_j(\underline{x})$$

# Least Squares for Classification (II)

- ▶ Data set:  $N \times (D+1)$  data matrix,  $N \times K$  target matrix

$$\tilde{\mathbf{X}} = \begin{pmatrix} -\tilde{x}_1^T - \\ \vdots \\ -\tilde{x}_N^T - \end{pmatrix}$$

$$\mathbf{T} = \begin{pmatrix} -t_1^T - \\ \vdots \\ -t_N^T - \end{pmatrix}$$

$$\underline{t}_n = (0, 0, 1, 0, 0)^T$$

$$(\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{T})_{nk} = \sum_m \tilde{x}_{nm} \tilde{w}_{mk} - T_{nk}$$

- ▶ Use sum-of-squares regression error function

$$E_D(\tilde{\mathbf{W}}) = \frac{1}{2} \text{Tr} \left[ (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{T})^T (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{T}) \right]$$

$$= \frac{1}{2} \sum_{k=1}^K \sum_{n=1}^N \left( \sum_{m=1}^M \tilde{x}_{nm} \tilde{w}_{mk} - T_{nk} \right)^2$$

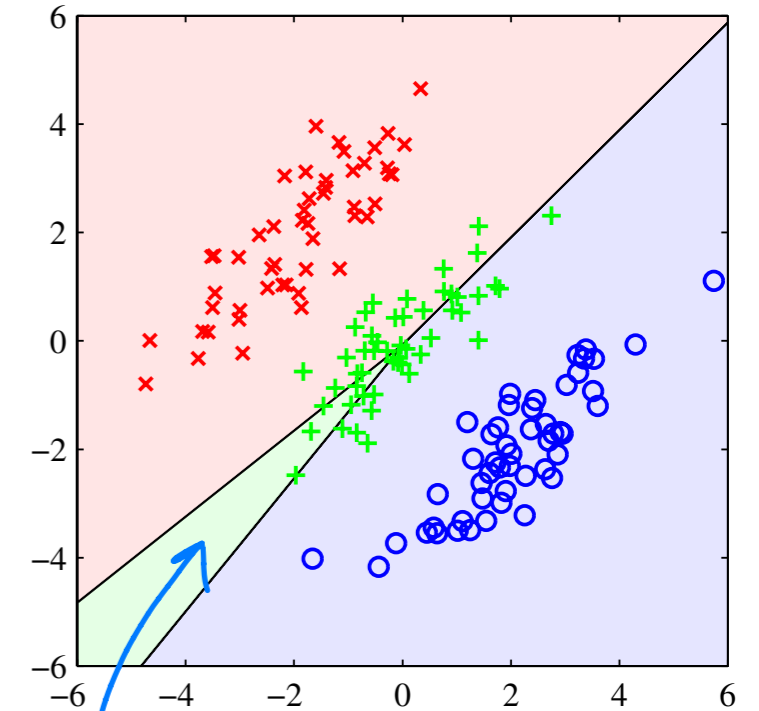
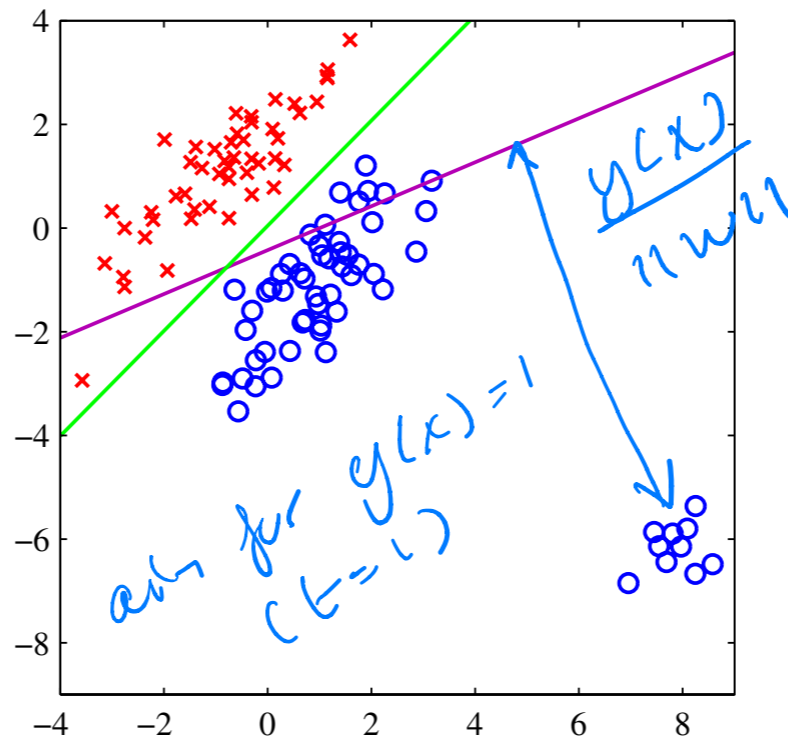
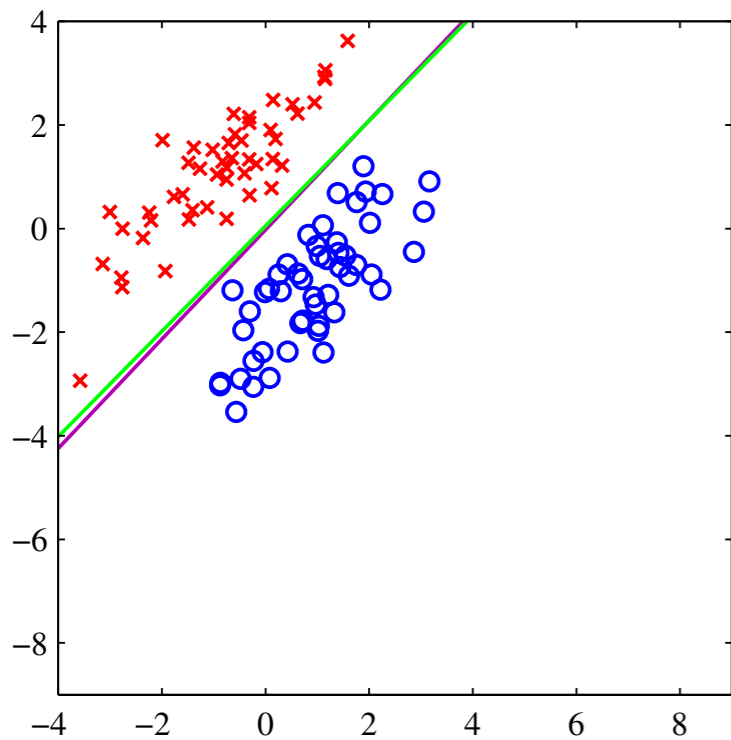
- ▶ Minimize  $E_D(\tilde{\mathbf{W}})$  as a function of  $\tilde{\mathbf{W}}$ :

$$\frac{\partial E_D(\tilde{\mathbf{W}})}{\partial \tilde{\mathbf{W}}} = 0$$

- ▶ Solution:  $\tilde{\mathbf{W}}_{\text{LS}} = (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \mathbf{T} = \tilde{\mathbf{X}}^\dagger \mathbf{T}$

- ▶ Discriminant function:  $\mathbf{y}_{\text{LS}}(\mathbf{x}) = \tilde{\mathbf{W}}_{\text{LS}}^T \tilde{\mathbf{x}}$

# Least Squares for Classification: Problems



**Figure:** least squares is very sensitive to outliers (Bishop 4.4)

**Figure:** masking for least squares for  $K > 2$  (Bishop 4.5)

1. The decision boundaries are very sensitive to outliers
2. For  $K > 2$  some decision regions can become very small or are even completely ignored
3. The components of  $\mathbf{y}_{LS}(\mathbf{x})$  are not real probabilities!

▸  $y_k(\mathbf{x})$

▸ if  $\sum_{k=1}^K t_k = 1 \rightarrow \sum_k y_k = 1$