



# Machine Learning 1

Lecture 7.2 - Supervised Learning  
Classification - Probabilistic Discriminative  
Models - Logistic Regression

*Erik Bekkers*

*(Bishop 4.3.2)*



# Classification Strategies

- ▶ Discriminant functions

Direct mapping of input to target

$$t = y(\underline{x}, \underline{w})$$

- ▶ Probabilistic discriminative models

Posterior class probabilities:

$$p(C_k | \underline{x})$$

- ▶ Probabilistic generative models

Class-conditional densities:

$$p(\underline{x} | C_k)$$

Prior class probabilities:

$$\left. \begin{array}{l} p(\underline{x} | C_k) \\ p(C_k) \end{array} \right\} \xRightarrow{\text{Bayes}} p(C_k | \underline{x})$$

# Logistic Regression for Two Classes

- ▶ Given: Dataset  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T$  with binary targets  $\mathbf{t} = (t_1, \dots, t_N)^T$  with  $t_n \in \{\mathcal{C}_1, \mathcal{C}_2\} = \{1, 0\}$
- ▶ Basis functions  $\phi = \phi(\mathbf{x}) = (\phi_0(\mathbf{x}), \dots, \phi_{m-1}(\mathbf{x}))^T$   
 $\phi_0(\mathbf{x}) = 1$
- ▶ **Probabilistic Discriminative Linear Models:** posteriors  $p(\mathcal{C}_k | \phi)$  are nonlinear functions with a linear function of  $\phi$  as input.

$$p(\mathcal{C}_k | \phi, \mathbf{w}) = f(\mathbf{w}^T \phi)$$

non-linear

linear

Logistic sigmoid

Generative setting

$$p(\mathcal{C}_1 | \mathbf{x}) = \sigma\left(\ln \frac{p(\mathbf{x} | \mathcal{C}_1) p(\mathcal{C}_1)}{p(\mathbf{x} | \mathcal{C}_2) p(\mathcal{C}_2)}\right)$$

- ▶ Logistic regression (K=2)

$$p(\mathcal{C}_1 | \phi, \mathbf{w}) = y(\phi) = \sigma(\mathbf{w}^T \phi)$$

$$p(\mathcal{C}_2 | \phi, \mathbf{w}) = 1 - y(\phi) = 1 - \sigma(\mathbf{w}^T \phi)$$

$$p(t | \phi, \mathbf{w}) = y(\phi)^t (1 - y(\phi))^{1-t}$$

# Logistic Regression for Two Classes

$$\underline{w}, \underline{\phi} \in \mathbb{R}^M$$

- ▶ Logistic Regression:

$$p(\mathcal{C}_1 | \phi, \mathbf{w}) = \sigma(\mathbf{w}^T \phi) \quad p(\mathcal{C}_2 | \phi, \mathbf{w}) = 1 - \sigma(\mathbf{w}^T \phi)$$

# parameters:  $\underline{w}: M$  linearly with  $M$

- ▶ Gaussian conditional densities:

$$p(\mathbf{x} | \mathcal{C}_k) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}_k|^{1/2}} \exp\left\{\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_k)\right\}$$

class priors  $p(\mathcal{C}_k)$

# parameters:  $\underline{\mu}_k \in \mathbb{R}^D$   $M$   
 $\boldsymbol{\Sigma} \sim M^2$  quadratically with  $M$

# Logistic Regression for Two Classes

- ▶ Given: Dataset  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T$  with binary targets  
 $\mathbf{t} = (t_1, \dots, t_N)^T$  with  $t_n \in \{\mathcal{C}_1, \mathcal{C}_2\} = \{1, 0\}$

- ▶ Conditional likelihood function:

$$p(\mathbf{t}|\mathbf{X}, \mathbf{w}) = \prod_{n=1}^N p(t_n | \mathbf{x}_n, \mathbf{w}) = \prod_{n=1}^N y_n^{t_n} (1-y_n)^{1-t_n}$$

$$y_n = p(\mathcal{C}_1 | \phi_n) = \sigma(\mathbf{w}^T \phi_n) \quad \phi_n = \phi(\mathbf{x}_n)$$

- ▶ Maximizing the conditional likelihood/minimizing the cross-entropy

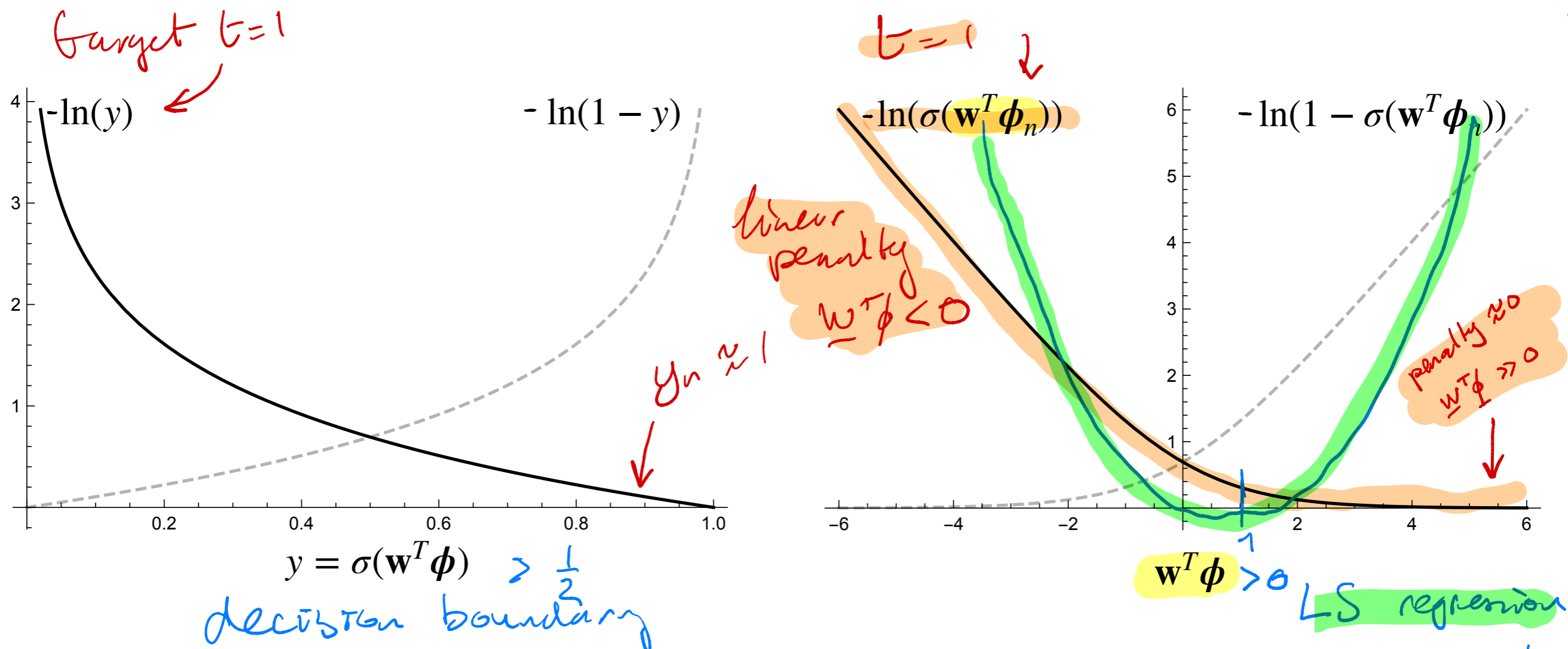
$$E(\mathbf{w}) = -\ln p(\mathbf{t}|\mathbf{X}, \mathbf{w}) = -\sum_{n=1}^N t_n \ln y_n + (1-t_n) \ln(1-y_n)$$

- ▶  $E(\mathbf{w})$ : convex, but no closed form solution!

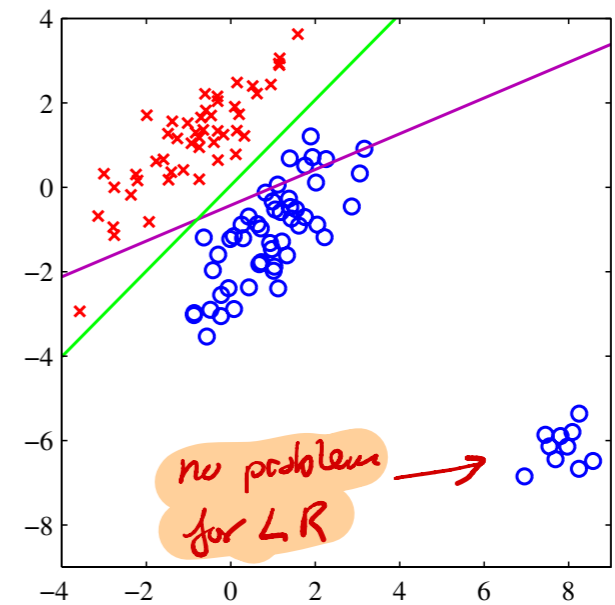
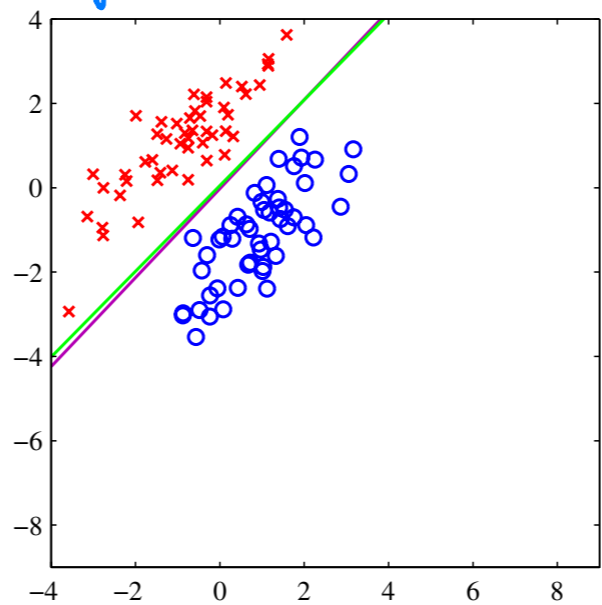
$$y_n = \sigma(\mathbf{w}^T \phi_n) \text{ is nonlinear in } \mathbf{w}$$

(Bishop 4.3.3)  
(Kleinman 2.0)

# The cross-entropy loss $E_n(\mathbf{w}) = -\left\{t_n \ln(y_n) + (1 - t_n)\ln(1 - y_n)\right\}$



Least squares outlier problem does not take place with logistic regression



$y = \mathbf{w}^T \boldsymbol{\phi}$

loss

$(y - t)^2$

no problem for LR

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

# Classification with Logistic Regression

- ▶ Dataset:  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T$  with targets  
 $\mathbf{t} = (t_1, \dots, t_N)^T$  with  $t_n \in \{\mathcal{C}_1, \mathcal{C}_2\} = \{1, 0\}$
- ▶ Basis functions  $\phi(\mathbf{x}) = (\phi_0(\mathbf{x}), \dots, \phi_{M-1}(\mathbf{x}))^T$
- ▶ Posterior distributions:  $p(\mathcal{C}_1 | \mathbf{x}, \mathbf{w}) = \sigma(\mathbf{w}^T \phi(\mathbf{x}))$
- ▶ Minimize  $E(\mathbf{w}) = -\ln p(\mathbf{t} | \mathbf{X}, \mathbf{w}) = -\sum_{n=1}^N t_n \ln y_n + (1 - t_n) \ln(1 - y_n)$   
with stochastic gradient descent or iterative reweighted least squares, to find  $\mathbf{w}^*$  (next videos)
- ▶ New datapoint  $\mathbf{x}'$  is assigned to class  $\mathcal{C}_1$  if  
 $p(\mathcal{C}_1 | \mathbf{x}', \mathbf{w}^*) = \sigma((\mathbf{w}^*)^T \phi(\mathbf{x}')) > \frac{1}{2}$   $\underline{\mathbf{w}}^* \cdot \underline{\phi}(\underline{x}) > 0$
- ▶ Decision boundaries:  
 $\underline{\mathbf{w}}^* \cdot \underline{\phi}(\underline{x}) = 0$