

UNIVERSITY OF AMSTERDAM Informatics Institute



# Machine Learning 1

Lecture 5.6 - Supervised Learning Classification - Probabilistic Generative Models

Erik Bekkers

(Bishop 1.5)

Slide credits: Patrick Forré and Rianne van den Berg

Image credit: Kirillm | Getty Images

### Probabilistic Generative Models: K=2

- Prior class probabilities:  $p(C_k)$
- Joint distribution:  $p(X, C_k) = p(X|C_k) p(C_k)$
- Posterior distribution: K=2

$$p(C_{1}|\mathbf{x}) = \frac{p(\mathbf{x}|C_{1})p(C_{1})}{p(\mathbf{x}|C_{1})p(C_{1}) + p(\mathbf{x}|C_{2})p(C_{2})} = p(\mathbf{x})$$

$$= \frac{1}{1 + \frac{p(\mathbf{x}|C_{2})p(C_{2})}{p(\mathbf{x}|C_{2})p(C_{2})}} = \frac{1}{1 + e^{-\alpha}}$$

$$a = \ln \frac{\sigma}{1 - \sigma} = \ln \frac{p(\mathbf{x}|C_{1})p(C_{2})}{p(\mathbf{x}|C_{2})p(C_{2})}$$

$$\log odds \mathcal{N}$$



Figure: Logistic Sigmoid function (red) (Bishop 4.9)

#### Probabilistic Generative Models: general K

• For multiple classes (general K):

$$p(C_k | \mathbf{x}) = \frac{p(\mathbf{x} | C_k) p(C_k)}{\sum_{j=1}^{K} p(\mathbf{x} | C_j) p(C_j)}$$

• 
$$a_k = \ln(p(\mathbf{x}|C_k)p(C_k))$$

$$\frac{e \kappa p(a_k)}{\sum_{j=1}^{k} e \kappa p(a_j)}$$

• Softmax: if 
$$a_k >> a_j$$
 for all  $j \neq k$ 

$$p(G_{k}|x) \approx 1$$
$$p(G_{j}|x) \approx 0$$

 $p(\mathbf{x})$ 

• Note: for K=2:

$$p(C_{1}|\mathbf{x}) = \frac{p(\mathbf{x}|C_{1})p(C_{1})}{p(\mathbf{x}|C_{1})p(C_{1}) + p(\mathbf{x}|C_{2})p(C_{2})} = \frac{1}{1 + \frac{p(\mathbf{x}|C_{2})p(C_{2})}{p(\mathbf{x}|C_{1})p(C_{1})}}$$
$$= O(a), \quad a = a, \quad -a, \quad a = \ln \frac{p(\mathbf{x}|C_{1})p(C_{1})}{p(\mathbf{x}|C_{1})p(C_{1})}$$

#### Class Conditional Densities: Continuous Inputs

Gaussian Class-conditional densities:

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

• Assume shared covariance matrix:  $\Sigma_k = \sum$ 

• K=2 classes: 
$$p(C_1|\mathbf{x}) = \frac{1}{1 + \exp(-a)} = \sigma(a)$$

 $a = \ln \frac{p(\mathbf{x}|C_{1})p(C_{1})}{p(\mathbf{x}|C_{2})p(C_{2})} = \ln \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_{1},\boldsymbol{\Sigma}) - \ln \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_{2},\boldsymbol{\Sigma}) + \ln \frac{p(C_{1})}{p(C_{2})}$  $= -\frac{1}{2}\ln|\boldsymbol{\Sigma}| - \frac{1}{2}(\mathbf{x}-\boldsymbol{\mu}_{1})^{T}\boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu}_{1}) + \frac{1}{2}\ln|\boldsymbol{\Sigma}| + \frac{1}{2}(\mathbf{x}-\boldsymbol{\mu}_{2})^{T}\boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu}_{2}) + \ln \frac{p(C_{1})}{p(C_{2})}$  $= (\mathcal{M}_{1} - \mathcal{M}_{2})^{T}\boldsymbol{\Sigma}^{-1}\boldsymbol{\Sigma} - \frac{1}{2}\mathcal{M}_{1}^{T}\boldsymbol{\Sigma}^{-1}\mathcal{M}_{2} + \frac{1}{2}\mathcal{M}_{2}^{T}\boldsymbol{\Sigma}^{-1}\mathcal{M}_{2} + \ln \frac{p(C_{2})}{p(C_{2})}$ 

• Generalized Linear Model:  $p(C_1|\mathbf{x}) = \sigma(\mathbf{w}^T\mathbf{x} + w_0)$   $\mathbf{w} = \mathbf{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$   $w_0 = -\frac{1}{2}\boldsymbol{\mu}_1^T \mathbf{\Sigma}^{-1} \boldsymbol{\mu}_1 + \frac{1}{2}\boldsymbol{\mu}_2^T \mathbf{\Sigma}^{-1} \boldsymbol{\mu}_2 + \ln \frac{p(C_1)}{p(C_2)}$  $Q = \alpha_z$  ( $\alpha = \infty$ )

Machine Learning 1



**Figure:** Left: class conditional densities  $p(x | C_k)$ . Right: posterior  $P(C_1|x)$  as sigmoid of linear function of x. (Bishop 4.9)

## Linear Discriminant Analysis: General K

• Gaussian Class-conditional densities & fixed covariance:

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

Posterior distributions:

$$p(C_k | \mathbf{x}) = \frac{\exp(a_k(\mathbf{x}))}{\sum_{j=1}^{K} \exp(a_j(\mathbf{x}))}$$

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

Decision boundary:

$$\mathbf{w}_{k} = \sum^{-1} \mathcal{M}_{k}$$

$$w_{k0} = -\frac{1}{2} \mathcal{M}_{k} \sum^{-1} \mathcal{M}_{k} + \ln p(\mathbf{C}_{k})$$

$$p(C_k|\mathbf{x}) = p(C_j|\mathbf{x}) \longrightarrow \mathcal{O}_k(\mathbf{x}) = \mathbf{a}_j(\mathbf{x})$$

• If all covariance matrices are different  $\Sigma_k \neq \Sigma_j$  then  $a_k(\mathbf{x})$  will also contain quadratic terms in  $\mathbf{x}$ 

rl.

#### Example: LDA and QDA



**Figure:** Left: Gaussian class conditional densities p(x | C<sub>k</sub>), red and green have same covariance matrix. Right: posterior P(C<sub>k</sub> |x) distributions (RGB vectors) and decision boundaries. (Bishop 4.9)