



# Machine Learning 1

Lecture 6.1 - Supervised Learning  
Classification - Probabilistic Generative  
Models - Maximum Likelihood Solution

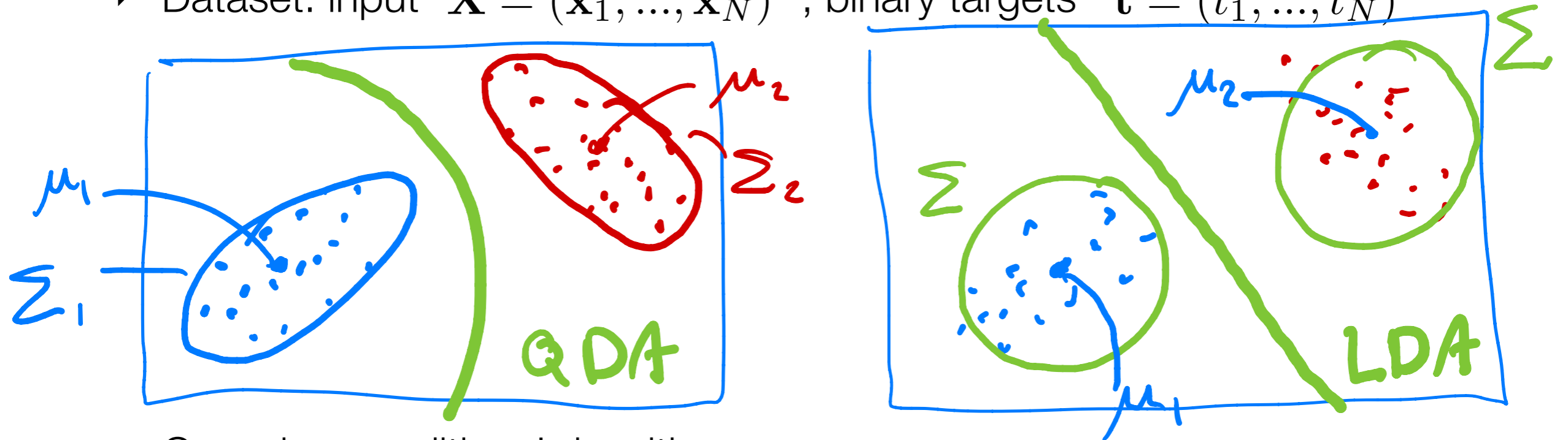
*Erik Bekkers*

*(Bishop 4.2.2)*



# LDA: Maximum Likelihood for K=2

- Dataset: input  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T$ , binary targets  $\mathbf{t} = (t_1, \dots, t_N)^T$



- Gaussian conditional densities

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

- Use Maximum likelihood to estimate  $\boldsymbol{\mu}_k$ ,  $\Sigma$  and priors  $p(C_k)$
- Denote  $p(C_1) = q$  and  $p(C_2) = 1 - q$

- For  $\mathbf{x}_n$  with  $t_n = 1$ :  $p(\mathbf{x}_n, C_1) = p(\mathbf{x}_n|C_1)p(C_1) = q \mathcal{N}(\underline{x}_n | \underline{\mu}_1, \Sigma)$
- For  $\mathbf{x}_n$  with  $t_n = 0$ :  $p(\mathbf{x}_n, C_2) = p(\mathbf{x}_n|C_2)p(C_2) = (1-q) \mathcal{N}(\underline{x}_n | \underline{\mu}_2, \Sigma)$

# LDA: Maximum Likelihood for K=2

- ▶ Gaussian conditional densities

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

- ▶ Use Maximum likelihood to estimate  $\boldsymbol{\mu}_k$ ,  $\boldsymbol{\Sigma}$  and priors  $p(C_k)$

- ▶ Denote  $p(C_1) = q$  and  $p(C_2) = 1 - q$

- ▶ Likelihood

$$\begin{aligned} p(\mathbf{t}, \mathbf{X} | q, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}) &= \prod_{n=1}^N p(\mathbf{x}_n, t_n) = \prod_{n=1}^N p(\mathbf{x}_n | t_n) p(t_n) \\ &= \prod_{n=1}^N [p(\mathbf{x}_n | C_1) p(C_1)]^{t_n} [p(\mathbf{x}_n | C_2) p(C_2)]^{1-t_n} \\ &= \prod_{n=1}^N [q \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \boldsymbol{\Sigma})]^{t_n} [(1 - q) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_2, \boldsymbol{\Sigma})]^{1-t_n} \end{aligned}$$

# LDA: Maximum Likelihood for K=2

- ▶ Likelihood

$$p(\mathbf{t}, \mathbf{X} | q, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}) = \prod_{n=1}^N \underbrace{[q \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \boldsymbol{\Sigma})]^{t_n}}_{\text{blue underline}} \underbrace{[(1-q) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_2, \boldsymbol{\Sigma})]^{1-t_n}}_{\text{red underline}}$$

- ▶ Log likelihood

$$\ln p(\mathbf{t}, \mathbf{X} | q, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}) = \sum_{n=1}^N \underbrace{t_n \ln q + t_n \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \boldsymbol{\Sigma})}_{\text{blue underline}} + \underbrace{(1-t_n) \ln(1-q) + (1-t_n) \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_2, \boldsymbol{\Sigma})}_{\text{red underline}}$$

- ▶ Estimate for q:

$$\frac{\partial}{\partial q} \ln p(\mathbf{t}, \mathbf{X} | q, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}) = \sum_{n=1}^N \frac{t_n}{q} + \frac{1-t_n}{1-q} \cdot (-1) = \sum_{n=1}^N \frac{t_n - (1-t_n)}{q(1-q)}$$

$\frac{t_n - (1-t_n)}{q(1-q)} = \frac{t_n - 1 + t_n}{q(1-q)} = \frac{2t_n - 1}{q(1-q)}$

$$= \sum_{n=1}^N \frac{t_n - q}{q(1-q)} = 0$$

$(q \neq 0)$   
 $(q \neq 1) \Rightarrow \sum_{n=1}^N q = \sum_{n=1}^N t_n \Rightarrow$

$$q_{ML} = \frac{1}{N} \sum_{n=1}^N t_n = \frac{N_1}{N}$$

# LDA: Maximum Likelihood for K=2

- ▶ log likelihood:

$$\ln p(\mathbf{t}, \mathbf{X} | q, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}) = \sum_{n=1}^N t_n \ln q + t_n \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}) + (1 - t_n) \ln(1 - q) + (1 - t_n) \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_2, \boldsymbol{\Sigma})$$

- ▶ Estimate for  $\boldsymbol{\mu}_1$

$$\frac{\partial}{\partial \boldsymbol{\mu}_1} \ln p(\mathbf{t}, \mathbf{X} | q, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}) = \frac{\partial}{\partial \boldsymbol{\mu}_1} \sum_{n=1}^N t_n \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \boldsymbol{\Sigma})$$

$$= -\frac{1}{2} \frac{\partial}{\partial \boldsymbol{\mu}_1} \sum_{n=1}^N t_n (\mathbf{x}_n - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_1)$$

( $\boldsymbol{\Sigma}$  is symmetric)

$$\sum_{n=1}^N t_n (\mathbf{x}_n - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1} = \mathbf{0}^T$$

( $\boldsymbol{\Sigma}$  pos. def.)

$$\sum_{n=1}^N t_n \boldsymbol{\mu}_1 = \sum_{n=1}^N t_n \mathbf{x}_n$$

$$\Rightarrow N_1 \boldsymbol{\mu}_1 = \sum_{n=1}^N t_n \mathbf{x}_n$$

$$\boldsymbol{\mu}_{1,ML} = \frac{1}{N_1} \sum_{n=1}^N t_n \mathbf{x}_n, \quad \boldsymbol{\mu}_{2,ML} = \frac{1}{N_2} \sum_{n=1}^N (1 - t_n) \mathbf{x}_n$$

# LDA: Maximum Likelihood for K=2

- ▶ log likelihood:

$$\ln p(\mathbf{t}, \mathbf{X} | q, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}) = \sum_{n=1}^N t_n \ln q + t_n \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}) + (1 - t_n) \ln(1 - q) + (1 - t_n) \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_2, \boldsymbol{\Sigma})$$

- ▶ Estimate for  $\boldsymbol{\Sigma}$

$$\frac{\partial}{\partial \boldsymbol{\Sigma}} \ln p(\mathbf{t}, \mathbf{X} | q, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}) = \frac{\partial}{\partial \boldsymbol{\Sigma}} \left[ -\frac{N}{2} \ln |\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{n=1}^N t_n (\mathbf{x}_n - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_1) - \frac{1}{2} \sum_{n=1}^N (1 - t_n) (\mathbf{x}_n - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_2) \right] = 0$$

- ▶ ML solution:

$$\boldsymbol{\Sigma}_{\text{ML}} = \frac{N_1}{N} \left[ \frac{1}{N_1} \sum_{n=1}^N t_n (\mathbf{x}_n - \boldsymbol{\mu}_{1,\text{ML}}) (\mathbf{x}_n - \boldsymbol{\mu}_{1,\text{ML}})^T \right] + \frac{N_2}{N} \left[ \frac{1}{N_2} \sum_{n=1}^N (1 - t_n) (\mathbf{x}_n - \boldsymbol{\mu}_{2,\text{ML}}) (\mathbf{x}_n - \boldsymbol{\mu}_{2,\text{ML}})^T \right]$$

*(Bishop 4.2.2 x notes on Canvas)*

# LDA: Maximum Likelihood for K=2

The ML solutions:

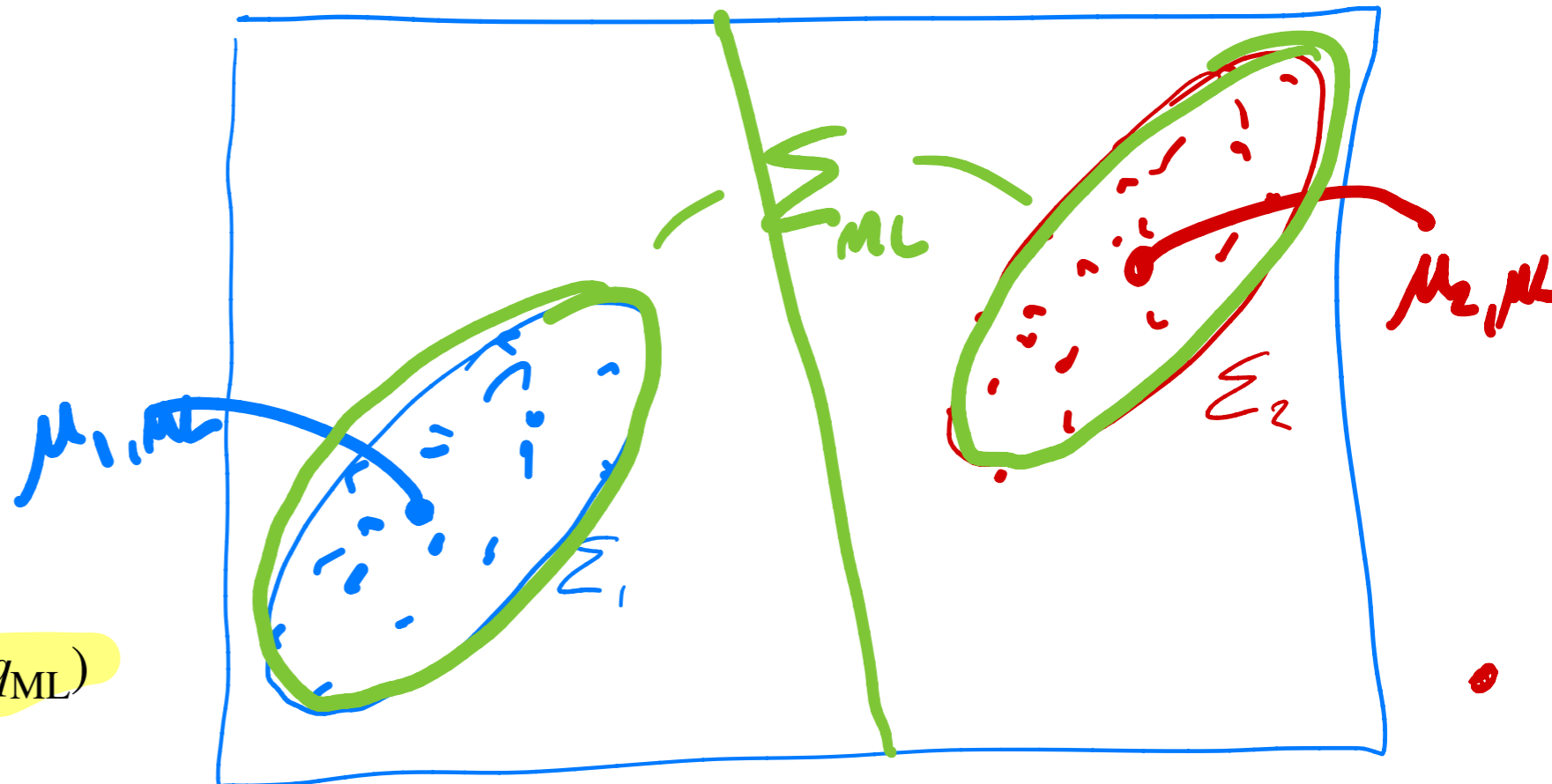
$$\boldsymbol{\mu}_{1,ML} = \frac{1}{N_1} \sum_{n=1}^N t_n \mathbf{x}_n \quad \boldsymbol{\mu}_{2,ML} = \frac{1}{N_2} \sum_{n=1}^N (1 - t_n) \mathbf{x}_n \quad q_{ML} = \frac{1}{N} \sum_{n=1}^N t_n = \frac{N_1}{N}$$

$$\boldsymbol{\Sigma}_{ML} = \frac{N_1}{N} \left[ \frac{1}{N_1} \sum_{n=1}^N t^n (\mathbf{x}_n - \boldsymbol{\mu}_{1,ML}) (\mathbf{x}_n - \boldsymbol{\mu}_{1,ML})^T \right] + \frac{N_2}{N} \left[ \frac{1}{N_2} \sum_{n=1}^N (1 - t^n) (\mathbf{x}_n - \boldsymbol{\mu}_{2,ML}) (\mathbf{x}_n - \boldsymbol{\mu}_{2,ML})^T \right]$$

For the joint probabilities:

$$p(\mathbf{x}, C_1) = p(\mathbf{x} | C_1) p(C_1) \\ = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_{1,ML}, \boldsymbol{\Sigma}_1) q_{ML}$$

$$p(\mathbf{x}, C_2) = p(\mathbf{x} | C_2) p(C_2) \\ = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_{2,ML}, \boldsymbol{\Sigma}_2) (1 - q_{ML})$$





# LDA: prediction for K=2

- ▶ For new datapoint  $\mathbf{x}'$ :

$$p(C_1 | \mathbf{x}') = \sigma(\underline{\mathbf{w}_{ML}^T \mathbf{x}' + w_{0,ML}})$$

$$\mathbf{w}_{ML} = \Sigma_{ML}^{-1} (\boldsymbol{\mu}_{1,ML} - \boldsymbol{\mu}_{2,ML})$$

$$w_{0,ML} = -\frac{1}{2} \boldsymbol{\mu}_{1,ML}^T \Sigma_{ML}^{-1} \boldsymbol{\mu}_{1,ML} + \frac{1}{2} \boldsymbol{\mu}_{2,ML}^T \Sigma_{ML}^{-1} \boldsymbol{\mu}_{2,ML} + \ln \frac{q_{ML}}{1 - q_{ML}}$$

- ▶ Assign  $\mathbf{x}'$  to  $C_1$  if  $p(C_1 | \mathbf{x}') \geq \frac{1}{2}$

$$\left( \mathbf{w}_{ML}^T \mathbf{x}' + w_{0,ML} \geq 0 \right)$$

- ▶ Disadvantage of LDA:

- ▶ Gaussian distribution is sensitive to outliers
- ▶ Linearity/handcrafted features restrict application
- ▶ Maximum likelihood is prone to overfitting