

Machine Learning 1

Lecture 6 - Classification - Probabilistic
Generative Models - Maximum Likelihood -
Discrete Data - Discriminant Functions

Erik Bekkers



Decision theory: Misclassification Rate

\hat{x} : Decision boundary

x_0 : Optimal Decision boundary

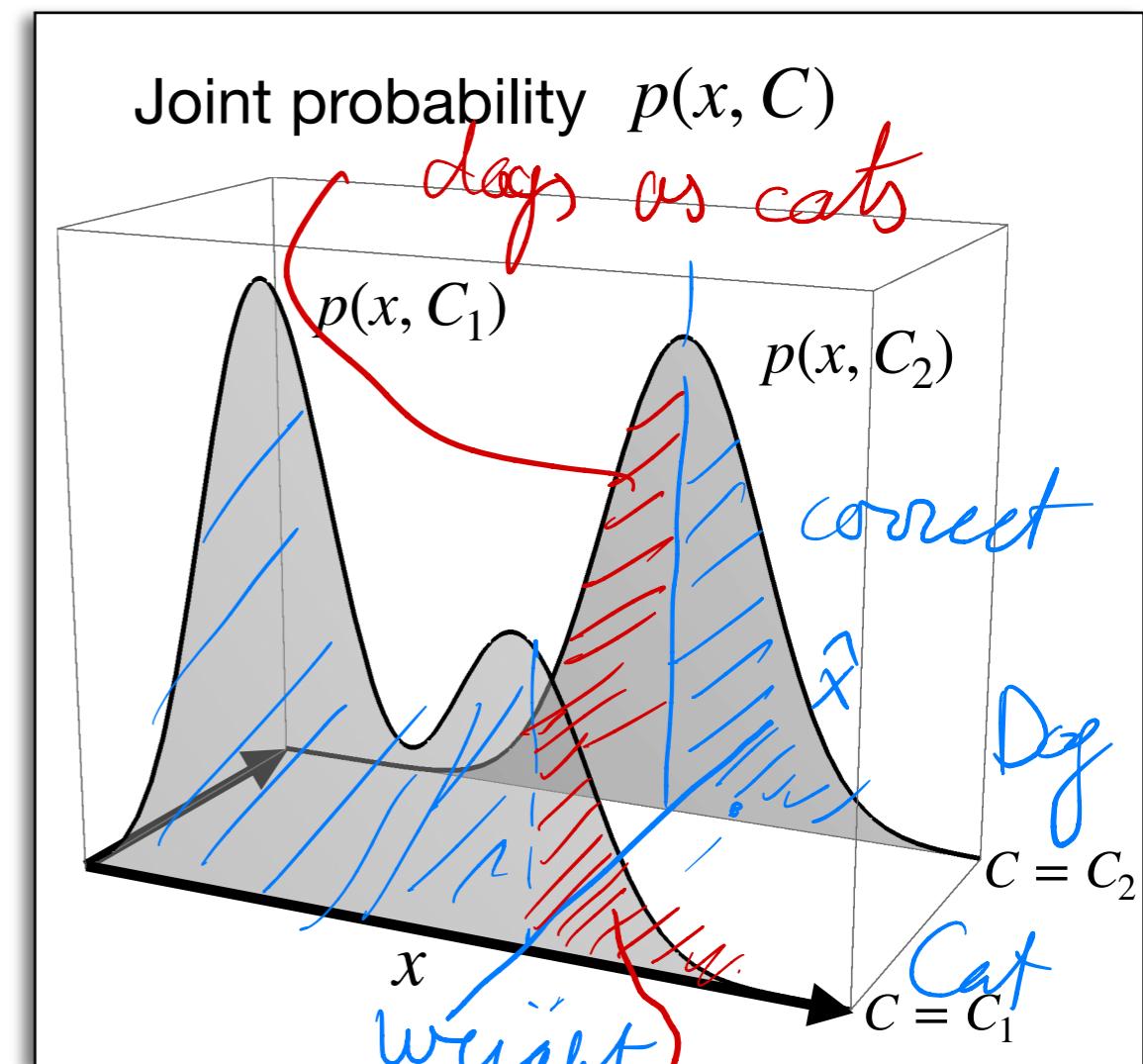
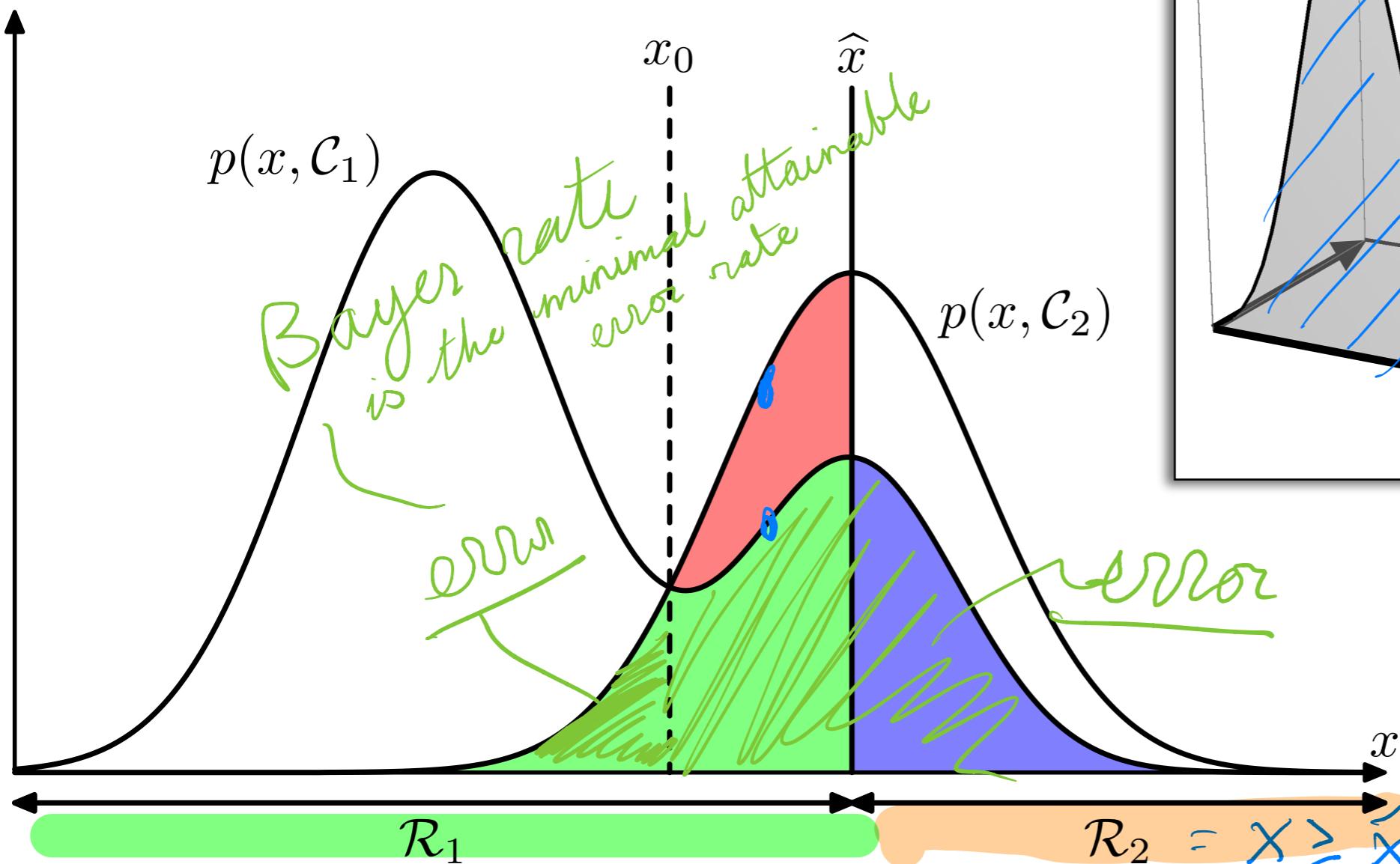


Figure: joint probability distributions and decision boundary (Bishop 1.24)

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:

$$p(x_1 | C_k) \cdot p(x, C_k)$$

Prior class probabilities:

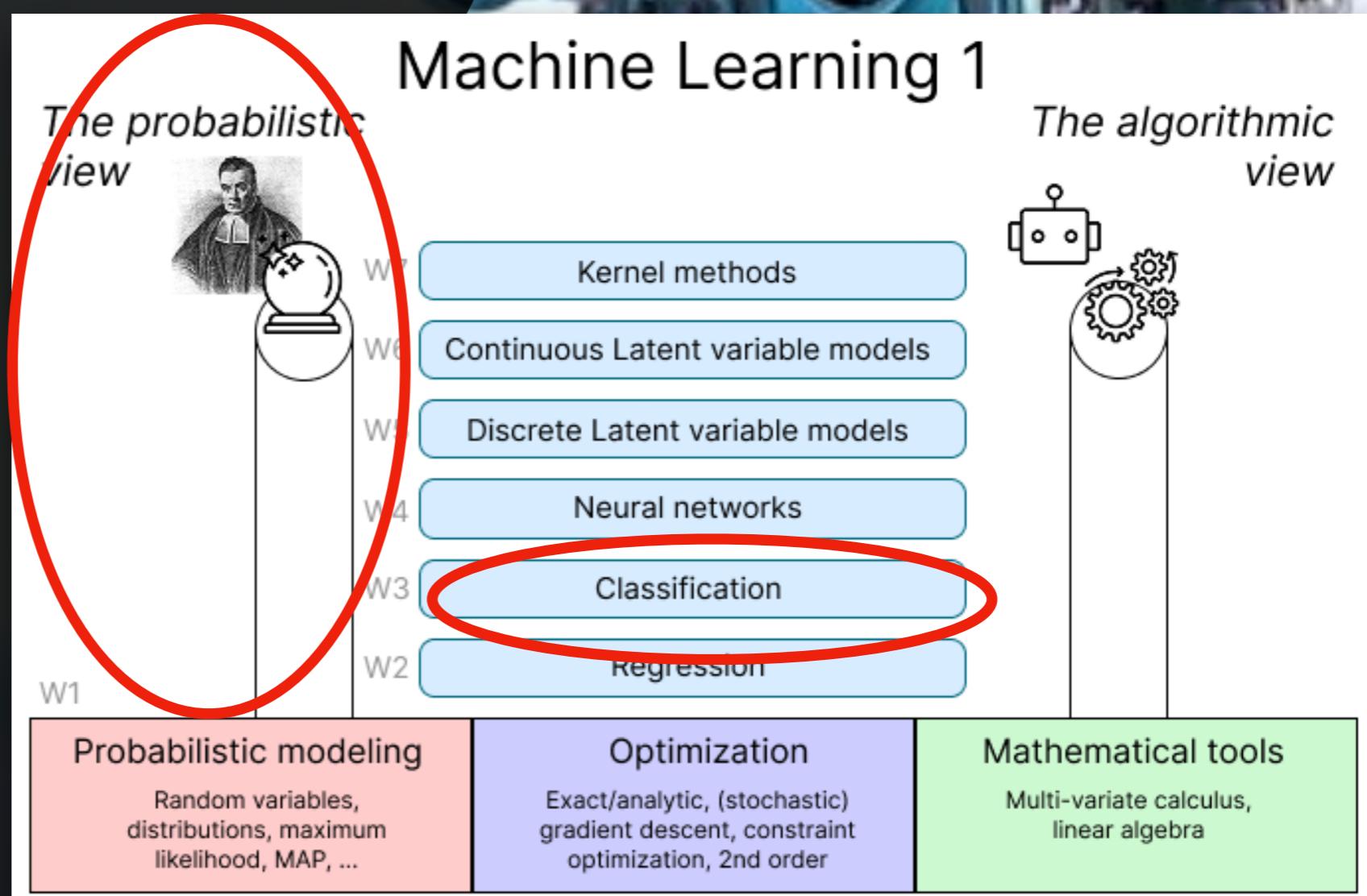
$$p(C_k) \cdot p(C_k | \underline{x})$$

Machine Learning 1

Lecture 5.6 - Supervised Learning
Classification - **Probabilistic Generative Models**

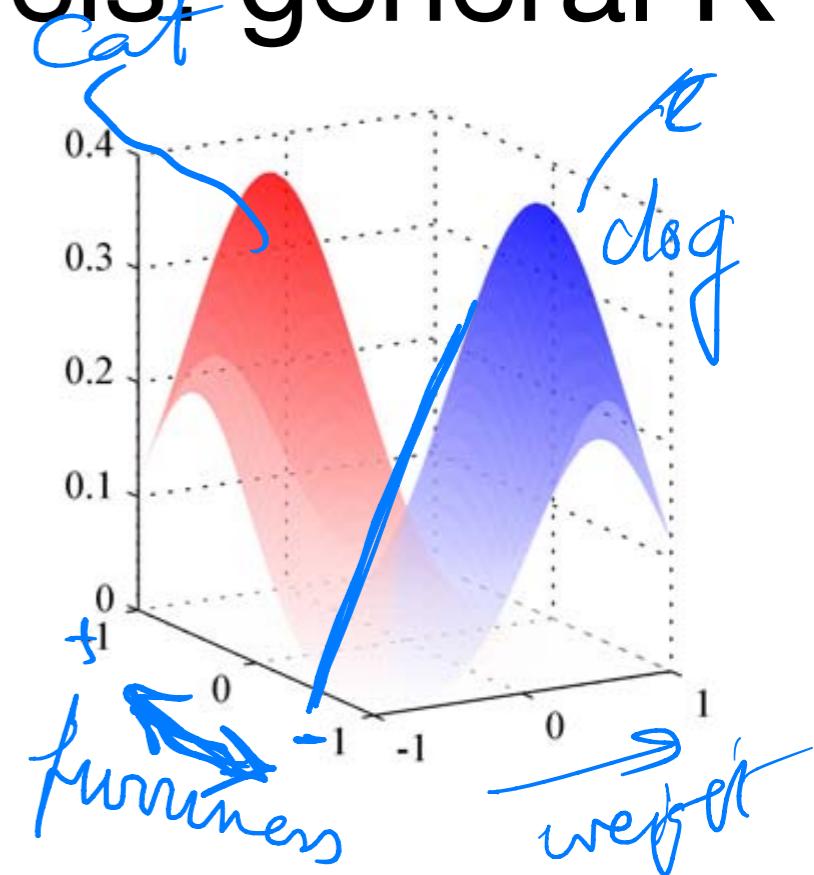
Erik Bekkers

(Bishop 1.5)



Probabilistic Generative Models: general K

- ▶ Model:
 - ▶ Class-conditional densities: $p(\mathbf{x} | C_k)$
 - ▶ Prior class probabilities: $p(C_k)$
- ▶ This gives access to joint and posterior:
 - ▶ $p(\mathbf{x}, C_k)$ and $p(C_k | \mathbf{x})$
- ▶ Why are we interested in joint or posterior?
 - ▶ Decision theory tells us that the best prediction for input \mathbf{x} , is to choose the class with highest joint $p(\mathbf{x}, C_k) = p(C_k | \mathbf{x}) p(\mathbf{x})$
 - ▶ Or equivalently: choose class with highest posterior $p(C_k | \mathbf{x})$
 - ▶ Decision boundary between C_k and C_j are at $p(C_k | \mathbf{x}) = p(C_j | \mathbf{x})$



Probabilistic Generative Models: $K = 2$

- Class-conditional densities: $p(x | C_k)$
- Prior class probabilities: $p(C_k)$
- Joint distribution: $p(x, C_k) = p(x | C_k) p(C_k)$
- Posterior distribution: $K = 2$

Bayes

$$p(C_1 | x) = \frac{p(x | C_1)p(C_1)}{p(x | C_1)p(C_1) + p(x | C_2)p(C_2)} = \frac{1}{1 + \frac{p(x | C_2)p(C_2)}{p(x | C_1)p(C_1)}}$$

"odds for class C_1 vs C_2 "

$$\text{odds} = \frac{p(C_1 | x)}{p(C_2 | x)} = \frac{p(x | C_1)p(C_1)}{p(x | C_2)p(C_2)} = \frac{\sigma}{1 - \sigma} = e^a$$

Log odds $a = \ln \frac{\sigma}{1 - \sigma}$

Logistic Sigmoid Function

$$\sigma(a) = \frac{1}{1 + \exp(-a)}$$

$$\sigma(-a) = 1 - \sigma(a)$$

$$\sigma'(a) = \sigma(a)(1 - \sigma(a))$$

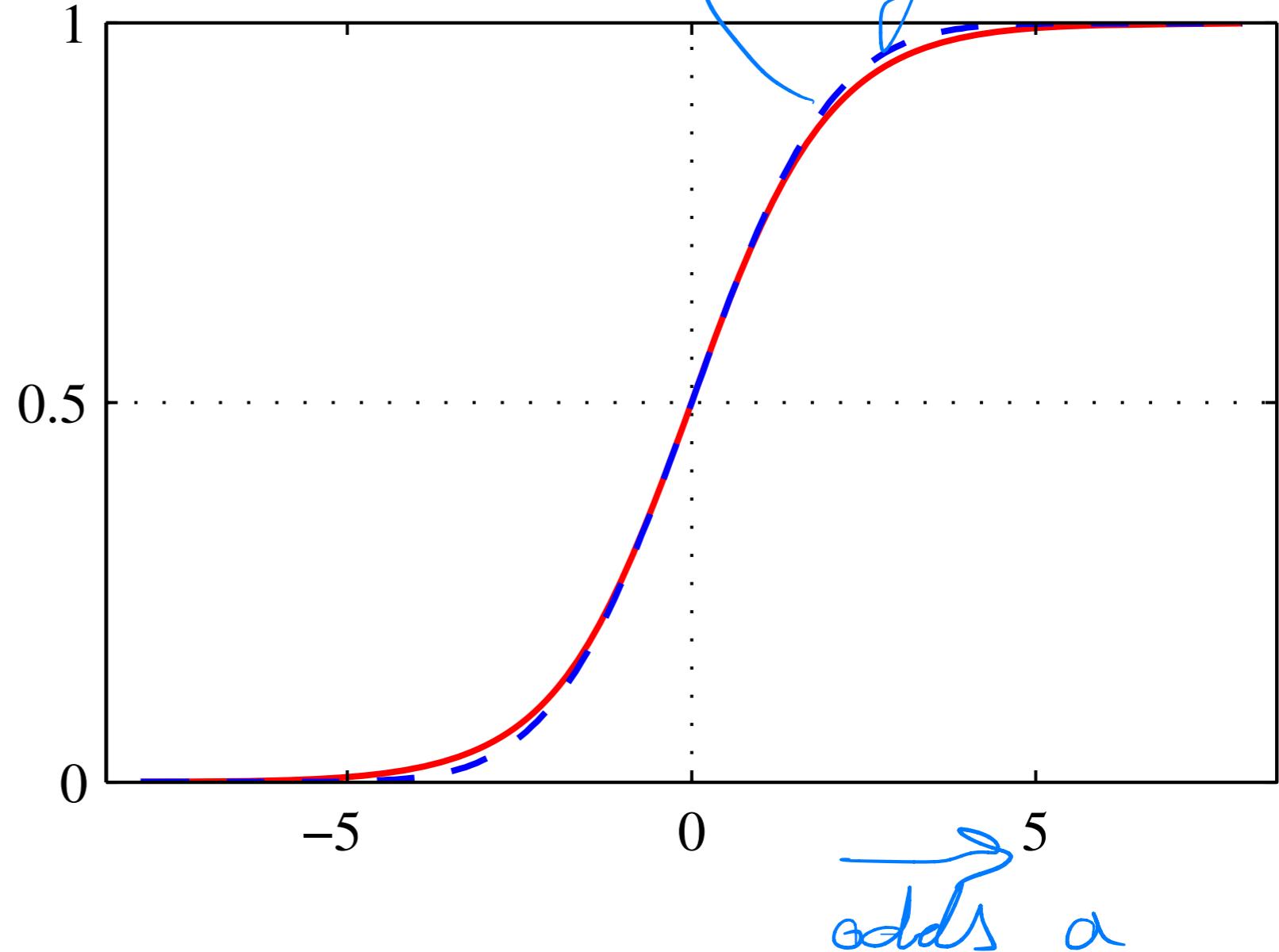


Figure: Logistic Sigmoid function (red) (Bishop 4.9)

Probabilistic Generative Models: general K

- For multiple classes (general K):

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

- Softmax: if $a_k \gg a_j$ for all $j \neq k$:

$$\begin{aligned} p(C_k | \mathbf{x}) &\approx 1 \\ p(C_j | \mathbf{x}) &\approx 0 \end{aligned}$$

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

called logits

- Note: for $K = 2$:

$$p(C_1 | \mathbf{x}) = \frac{p(\mathbf{x} | C_1)p(C_1)}{p(\mathbf{x} | C_1)p(C_1) + p(\mathbf{x} | C_2)p(C_2)} = \frac{1}{1 + \frac{p(\mathbf{x} | C_2)p(C_2)}{p(\mathbf{x} | C_1)p(C_1)}}$$

$$= \sigma(a), a = a_1 - a_2$$

$$a = \ln \frac{p(\mathbf{x} | C_2)p(C_2)}{p(\mathbf{x} | C_1)p(C_1)}$$

Class Conditional Densities: Continuous Inputs

Modeling Choice

- ▶ Gaussian Class-conditional densities:

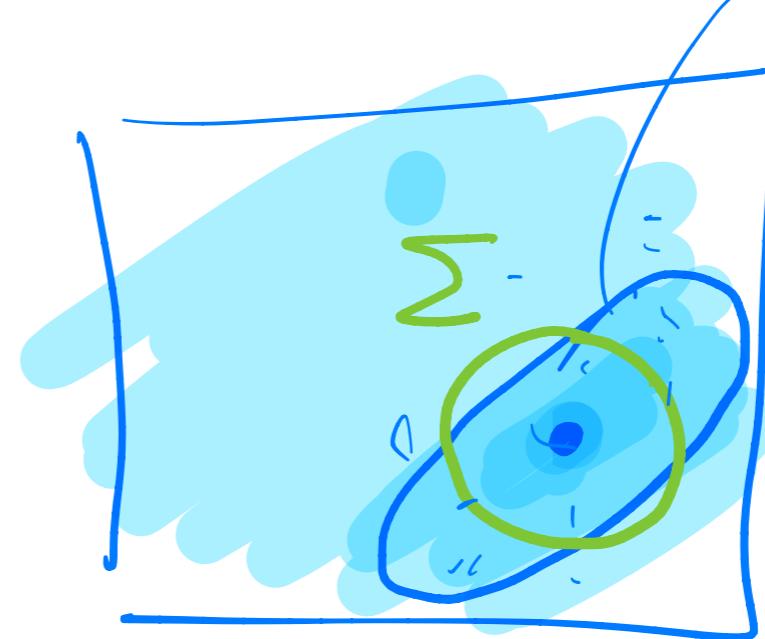
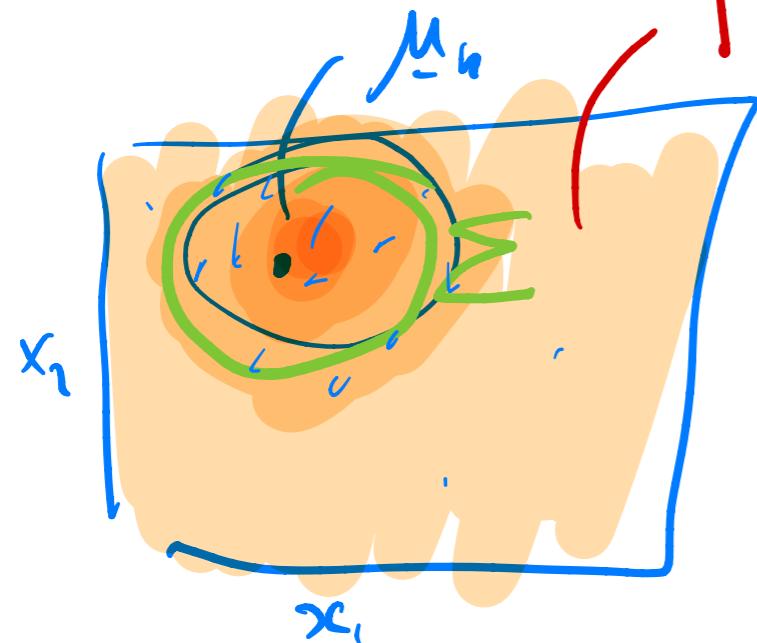
$$p(\mathbf{x} | C_k) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma_k|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k)\right\}$$

- ▶ Class posteriors simplify with the assumption of shared covariance matrix: $\Sigma_k = \Sigma$

- ▶ For $K = 2$ classes they are given as $p(C_1 | \mathbf{x}) = \frac{1}{1 + \exp(-a)} = \sigma(a)$ with

$$a = \ln \frac{p(\mathbf{x} | C_1)p(C_1)}{p(\mathbf{x} | C_2)p(C_2)} =$$

$p(C_1 | \mathbf{x})$ $p(C_2 | \mathbf{x})$



Class Conditional Densities: Continuous Inputs

- ▶ Gaussian Class-conditional densities:

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

- ▶ Class posteriors simplify with the assumption of shared covariance matrix: $\Sigma_k = \Sigma$

- ▶ For $K = 2$ classes they are given as $p(C_1 | \mathbf{x}) = \frac{1}{1 + \exp(-a)} = \sigma(a)$ with

$$a = \ln \frac{p(\mathbf{x} | C_1)p(C_1)}{p(\mathbf{x} | C_2)p(C_2)} = \ln \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_1, \Sigma) - \ln \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_2, \Sigma) + \ln \frac{p(C_1)}{p(C_2)}$$

Class Conditional Densities: Continuous Inputs

- ▶ Gaussian Class-conditional densities:

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

- ▶ Class posteriors simplify with the assumption of shared covariance matrix: $\Sigma_k = \Sigma$

- ▶ For $K = 2$ classes they are given as $p(C_1 | \mathbf{x}) = \frac{1}{1 + \exp(-a)} = \sigma(a)$ with

$$\begin{aligned} a &= \ln \frac{p(\mathbf{x} | C_1)p(C_1)}{p(\mathbf{x} | C_2)p(C_2)} = \ln \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_1, \Sigma) - \ln \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_2, \Sigma) + \ln \frac{p(C_1)}{p(C_2)} \\ &= -\frac{1}{2} \ln |\Sigma| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_1)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) + \frac{1}{2} \ln |\Sigma| + \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_2)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_2) + \ln \frac{p(C_1)}{p(C_2)} \end{aligned}$$

Handwritten annotations below the equation:

- Red circle: $-\frac{1}{2} \mathbf{x}^T \Sigma^{-1} \mathbf{x}$
- Orange circle: $+\mathbf{x}^T \Sigma \boldsymbol{\mu}_1$
- Red circle: $+\frac{1}{2} \boldsymbol{\mu}_1^T \Sigma^{-1} \boldsymbol{\mu}_1$
- Green circle: $-\frac{1}{2} \mathbf{x}^T \Sigma^{-1} \mathbf{x}$
- Blue circle: $+\mathbf{x}^T \Sigma \boldsymbol{\mu}_2$
- Green circle: $+\frac{1}{2} \boldsymbol{\mu}_2^T \Sigma^{-1} \boldsymbol{\mu}_2$

Class Conditional Densities: Continuous Inputs

- ▶ Gaussian Class-conditional densities:

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

- Class posteriors simplify with the assumption of shared covariance matrix: $\Sigma_k = \Sigma$

- For $K = 2$ classes they are given as $p(C_1 | \mathbf{x}) = \frac{1}{1 + \exp(-a)} = \sigma(a)$ with

$$a = \ln \frac{p(\mathbf{x} | C_1)p(C_1)}{p(\mathbf{x} | C_2)p(C_2)} = \ln \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_1, \Sigma) - \ln \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_2, \Sigma) + \ln \frac{p(C_1)}{p(C_2)}$$

$$= -\frac{1}{2} \ln |\Sigma| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_1)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) + \frac{1}{2} \ln |\Sigma| + \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_2)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_2) + \ln \frac{p(C_1)}{p(C_2)}$$

$$= (\mu_1 - \mu_2)^T \Sigma^{-1} \mathbf{x} - \frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 + \frac{1}{2} \mu_2^T \Sigma^{-1} \mu_2 + \ln \frac{p(C_1)}{p(C_2)}$$

||
 $\underline{\mathbf{w}}^T \mathbf{x}$
 +
 +
 w_0
 ||
 11

$a(x)$ is a linear model!

Class Conditional Densities: Continuous Inputs

- ▶ Gaussian Class-conditional densities:

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

- ▶ Class posteriors simplify with the assumption of shared covariance matrix: $\Sigma_k = \Sigma$

- ▶ For $K = 2$ classes they are given as $p(C_1 | \mathbf{x}) = \frac{1}{1 + \exp(-a)} = \sigma(a)$ with

$$\begin{aligned} a &= \ln \frac{p(\mathbf{x} | C_1)p(C_1)}{p(\mathbf{x} | C_2)p(C_2)} = \ln \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_1, \Sigma) - \ln \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_2, \Sigma) + \ln \frac{p(C_1)}{p(C_2)} \\ &= -\frac{1}{2} \ln |\Sigma| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_1)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) + \frac{1}{2} \ln |\Sigma| + \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_2)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_2) + \ln \frac{p(C_1)}{p(C_2)} \\ &= (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \Sigma^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_1^T \Sigma^{-1} \boldsymbol{\mu}_1 + \frac{1}{2} \boldsymbol{\mu}_2^T \Sigma^{-1} \boldsymbol{\mu}_2 + \ln \frac{p(C_1)}{p(C_2)} \end{aligned}$$

σ is non-linear

- ▶ Generalized Linear Model: $p(C_1 | \mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$

$$\mathbf{w} = \Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$$

$$w_0 = -\frac{1}{2} \boldsymbol{\mu}_1^T \Sigma^{-1} \boldsymbol{\mu}_1 + \frac{1}{2} \boldsymbol{\mu}_2^T \Sigma^{-1} \boldsymbol{\mu}_2 + \ln \frac{p(C_1)}{p(C_2)}$$

linear

Example: Linear Discriminant Analysis for K=2

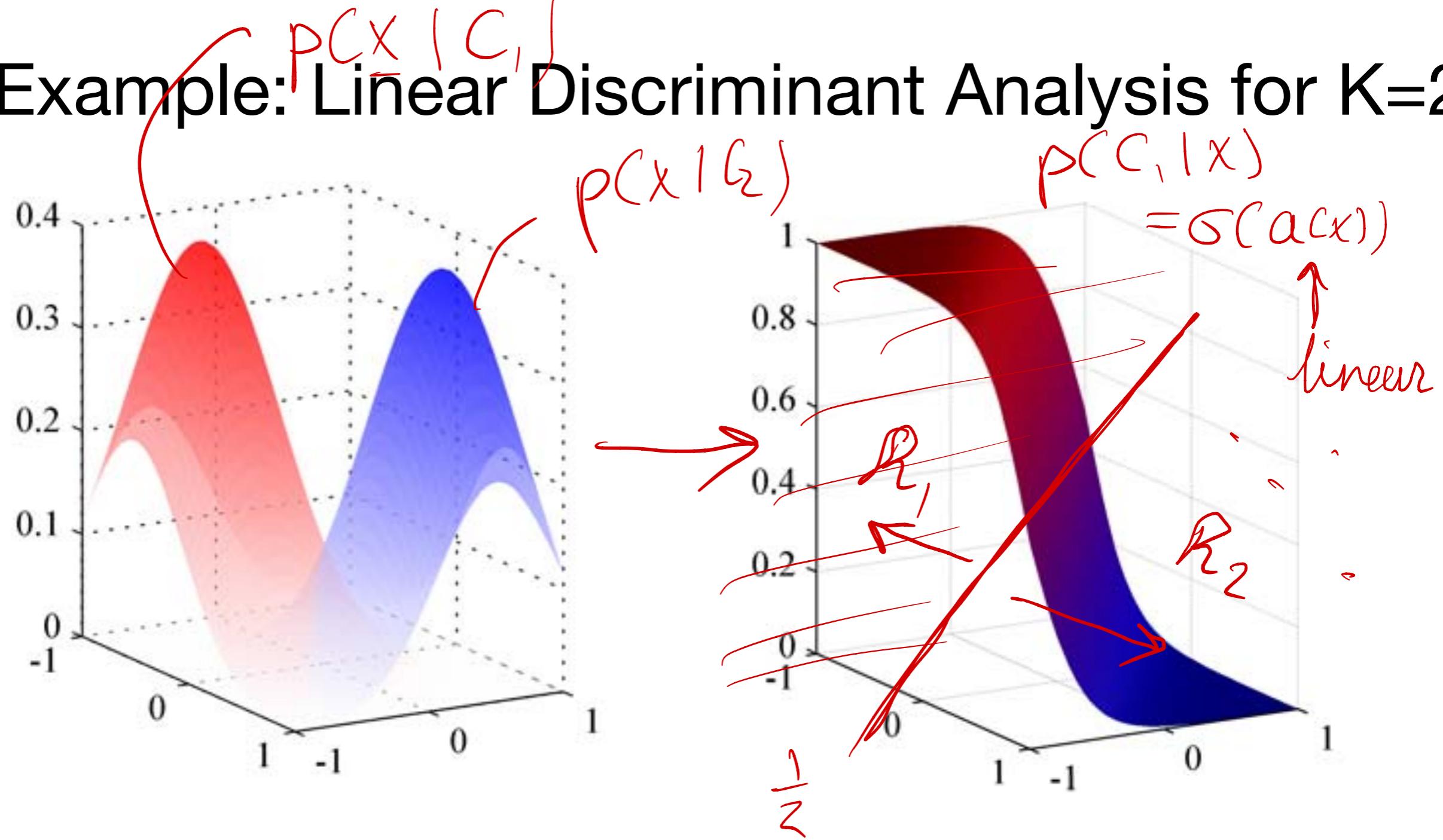


Figure: Left: class conditional densities $p(\mathbf{x} | C_k)$. Right: posterior $p(C_1 | \mathbf{x})$ as sigmoid of linear function of \mathbf{x} . (Bishop 4.9)

Linear Discriminant Analysis: General K

- Gaussian Class-conditional densities & fixed covariance:

$$\sum_k = \Sigma$$

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

- Posterior distributions:

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

one linear model
per class

has this
form

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

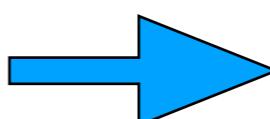
$$\mathbf{w}_k = \Sigma^{-1} \boldsymbol{\mu}_k$$

$$w_{k0} = -\frac{1}{2} \boldsymbol{\mu}_k^T \Sigma^{-1} \boldsymbol{\mu}_k + \ln p(C_k)$$

boundaries!

- Decision boundary:

$$p(C_k | \mathbf{x}) = p(C_j | \mathbf{x})$$



$$a_k(\underline{\mathbf{x}}) = a_j(\underline{\mathbf{x}})$$

so linear decision

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



non-linear d. b.

Example: LDA and QDA

Quadratic discriminant analysis

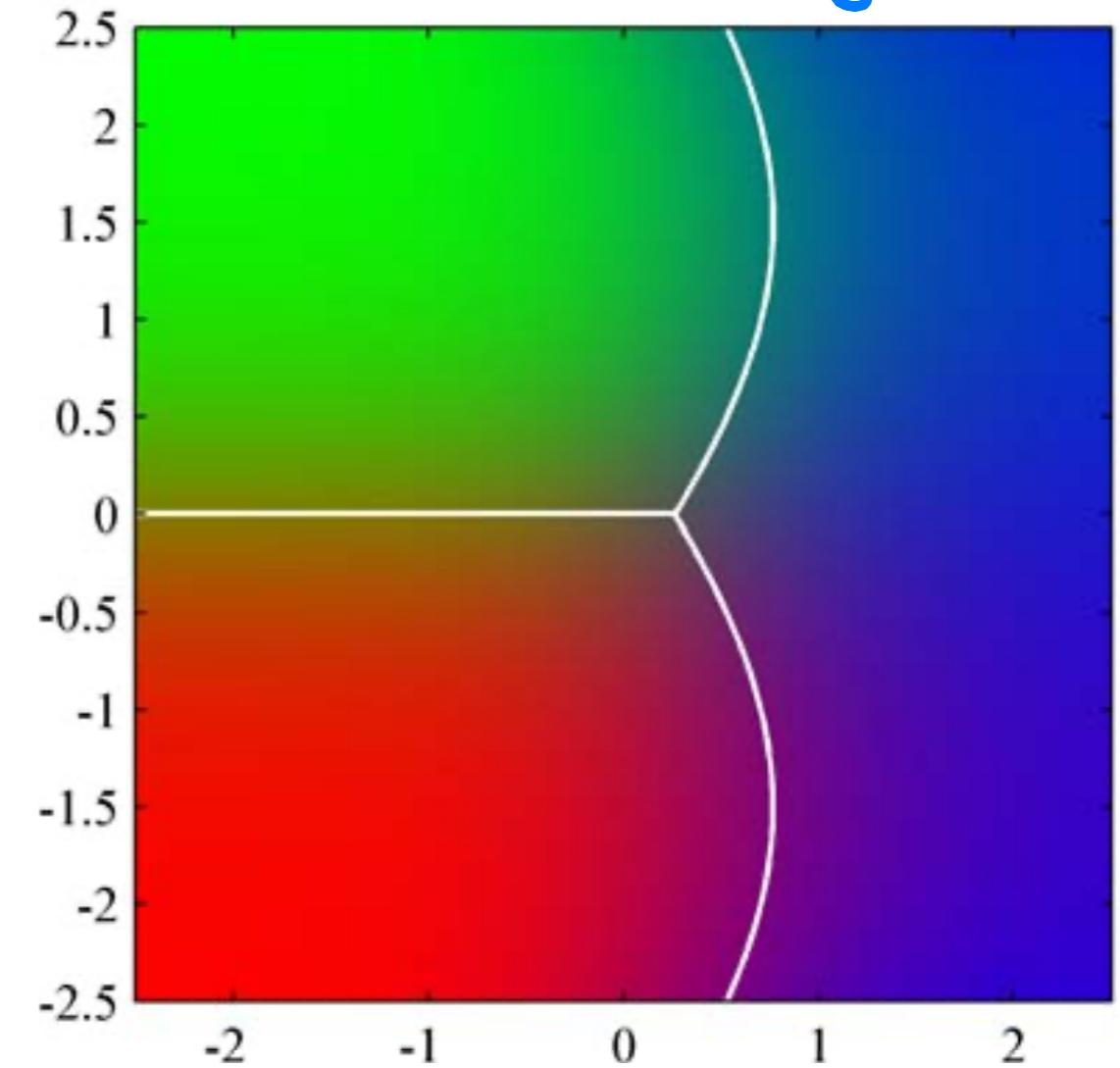
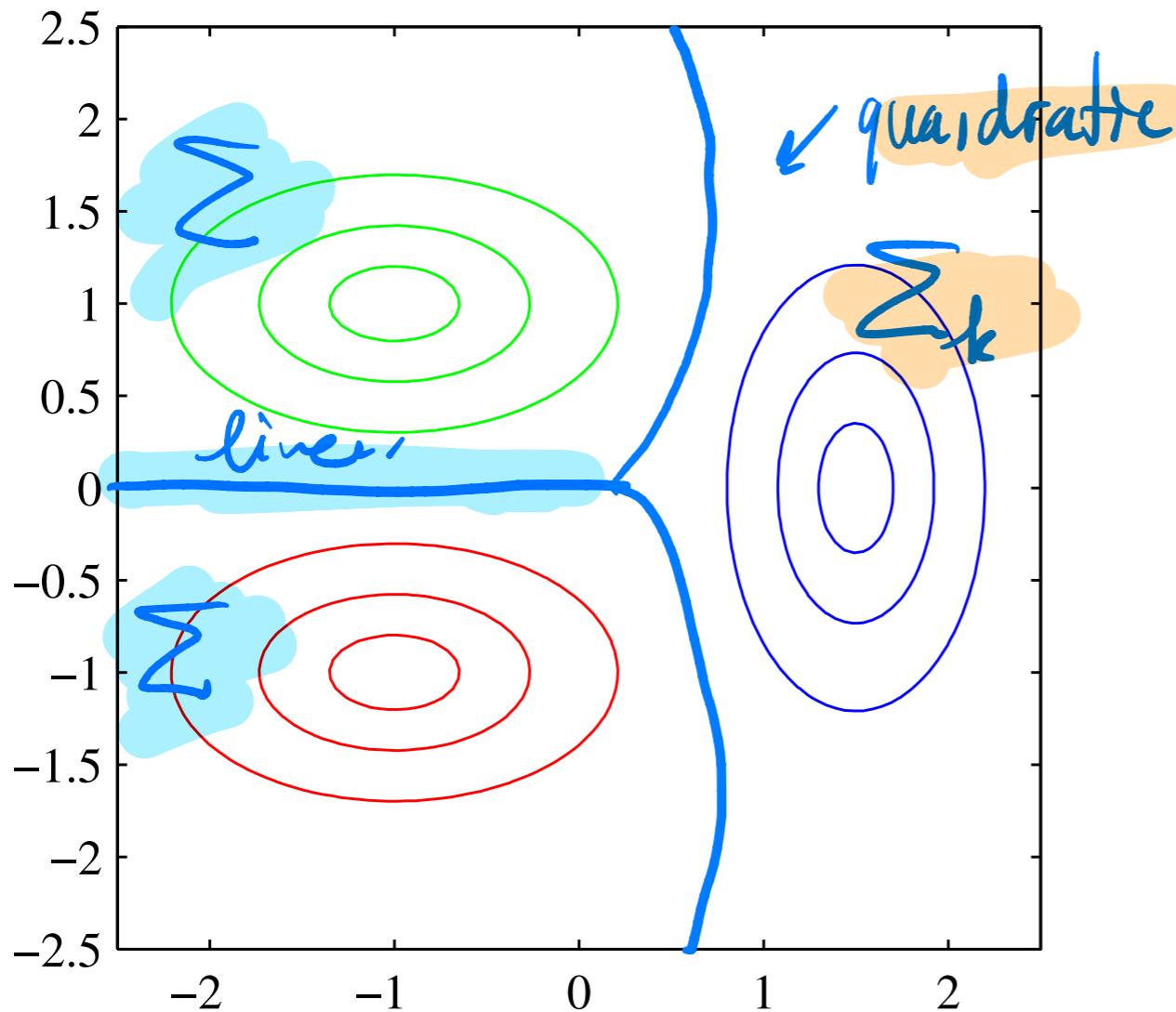


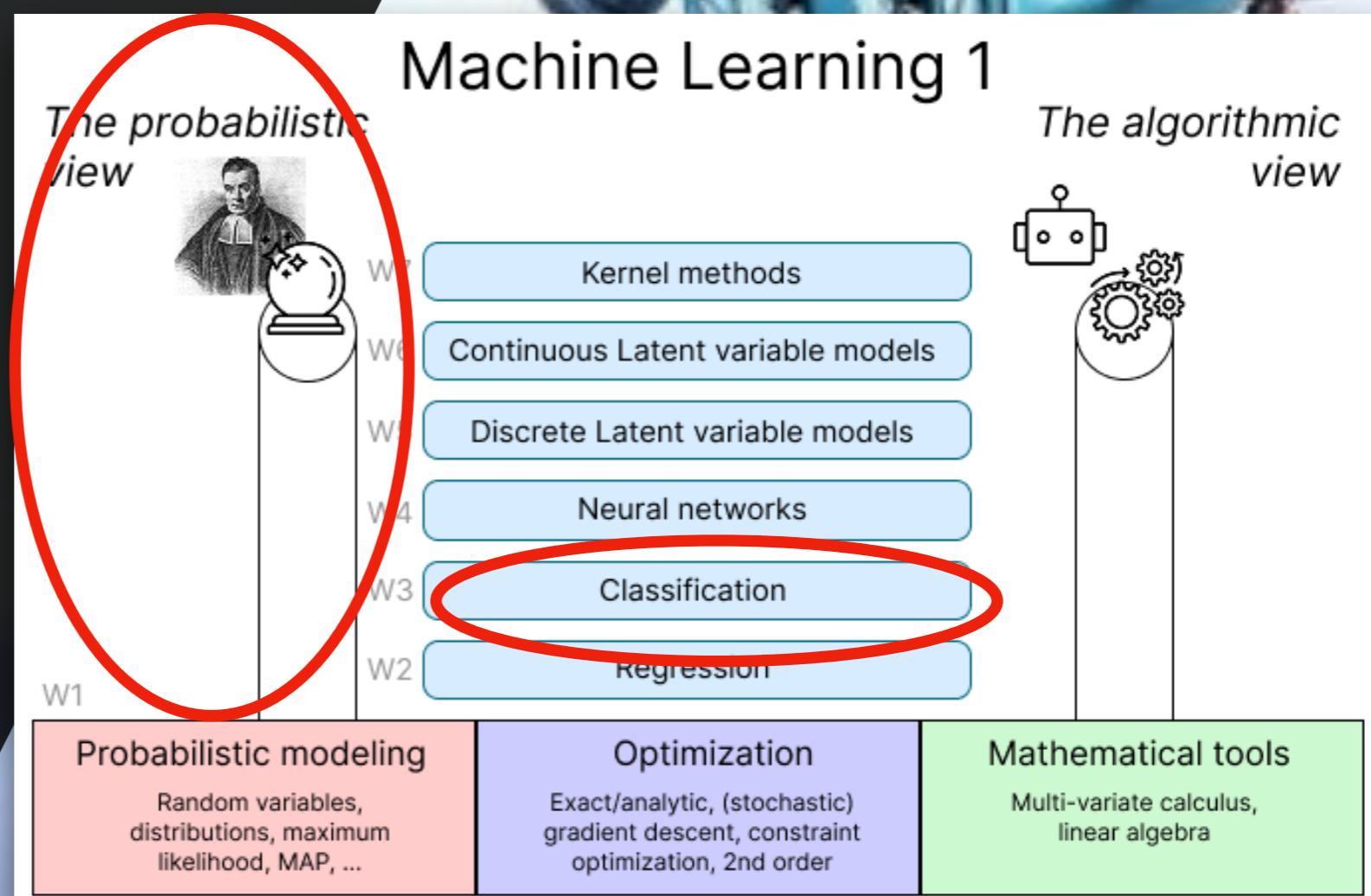
Figure: Left: Gaussian class conditional densities $p(\mathbf{x} | C_k)$, red and green have same covariance matrix. Right: posterior $p(C_k | \mathbf{x})$ distributions (RGB vectors) and decision boundaries. (Bishop 4.9)

Machine Learning 1

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

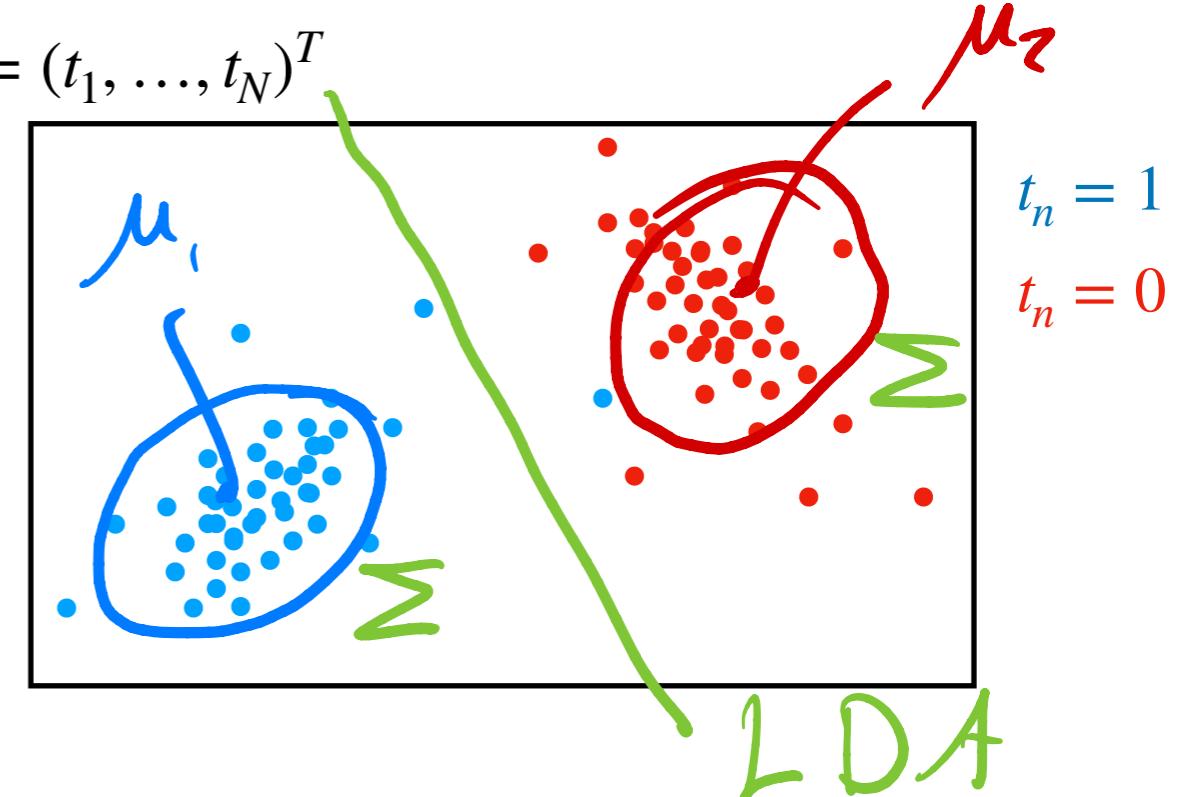
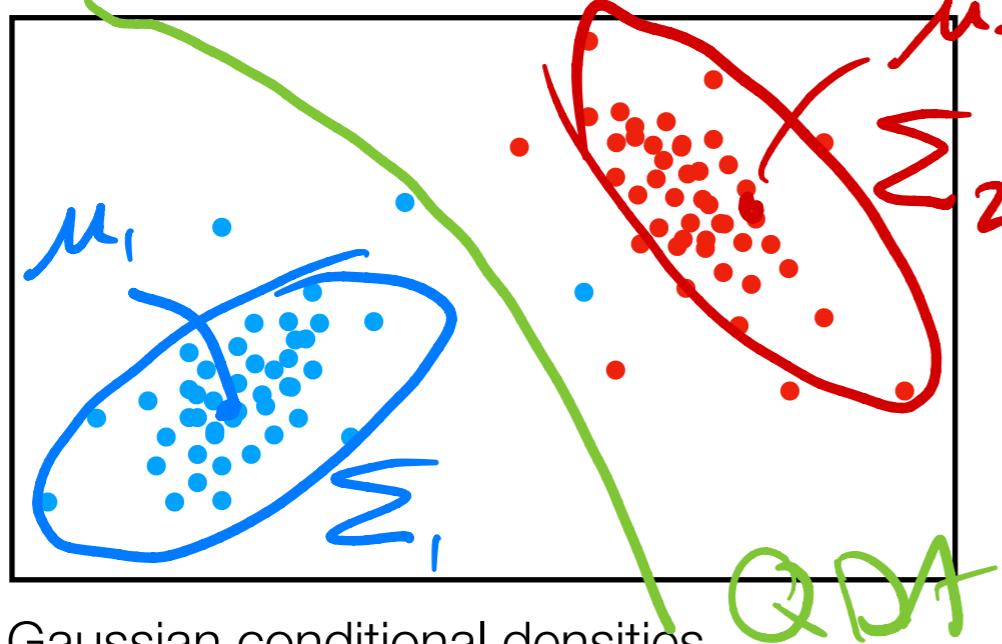
Erik Bekkers

(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

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

- Strategy: (the parameters ...)

use maximum likelihood to estimate μ_k , Σ and priors $p(C_k)$

let q parametrize the prior probabilities with q : $p(C_1) = q$ and $p(C_2) = 1 - q$

How to obtain? ML !!

$\mathcal{E}[0, 1]$

Bernoulli

For \mathbf{x}_n with $t_n = 1$:

$$p(\mathbf{x}_n, C_1) = p(\mathbf{x}_n | C_1)p(C_1) = q \mathcal{N}(\underline{x}_n | \mu_1, \Sigma)$$

For \mathbf{x}_n with $t_n = 0$:

$$p(\mathbf{x}_n, C_2) = p(\mathbf{x}_n | C_2)p(C_2) = (1-q) \mathcal{N}(\underline{x}_n | \mu_2, \Sigma)$$

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

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

- Strategy:

- use maximum likelihood to estimate $\boldsymbol{\mu}_k$, Σ and priors $p(C_k)$
- let q parametrize the prior probabilities with q : $p(C_1) = q$ and $p(C_2) = 1 - q$

- Likelihood:

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

$C_1 : t = 1$
 $C_2 : t = 0$

$p(t_n = 1) = q$
 $p(t_n = 0) = 1 - q$

if / else statement

LDA: Maximum Likelihood for K=2

- ▶ Likelihood

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

check!

- ▶ Log likelihood

$$\begin{aligned} \ln p(\mathbf{t}, \mathbf{X} | q, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \Sigma) &= \sum_{n=1}^N t_n \ln q + t_n \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \Sigma) \\ &\quad + (1 - t_n) \ln(1 - q) + (1 - t_n) \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_2, \Sigma) \end{aligned} \quad \left. \right\}$$

prev prob

- ▶ Let's estimate q by solving $\frac{\partial}{\partial q} \ln p(\mathbf{t}, \mathbf{X} | q, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \Sigma) = 0$

$$\frac{\partial}{\partial q} \ln p(\mathbf{t}, \mathbf{X} | q, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \Sigma) =$$

$$\ln a^t = t \ln a$$

$$\ln a \cdot b = \ln a + \ln b$$

LDA: Maximum Likelihood for K=2

- › Likelihood

$$p(\mathbf{t}, \mathbf{X} | q, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \Sigma) = \prod_{n=1}^N \left[q \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \Sigma) \right]^{t_n} \left[(1 - q) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_2, \Sigma) \right]^{1-t_n}$$

- › Log likelihood

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

- › Let's estimate q by solving $\frac{\partial}{\partial q} \ln p(\mathbf{t}, \mathbf{X} | q, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \Sigma) = 0$

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

LDA: Maximum Likelihood for K=2

- › Likelihood

$$p(\mathbf{t}, \mathbf{X} | q, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \Sigma) = \prod_{n=1}^N \left[q \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \Sigma) \right]^{t_n} \left[(1 - q) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_2, \Sigma) \right]^{1-t_n}$$

- › Log likelihood

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

- › Let's estimate q by solving $\frac{\partial}{\partial q} \ln p(\mathbf{t}, \mathbf{X} | q, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \Sigma) = 0$

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

- › Solve for q :

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

LDA: Maximum Likelihood for K=2

- › Likelihood

$$p(\mathbf{t}, \mathbf{X} | q, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \Sigma) = \prod_{n=1}^N \left[q \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \Sigma) \right]^{t_n} \left[(1 - q) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_2, \Sigma) \right]^{1-t_n}$$

- › Log likelihood

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

- › Let's estimate q by solving $\frac{\partial}{\partial q} \ln p(\mathbf{t}, \mathbf{X} | q, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \Sigma) = 0$

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

- › Solve for q :

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

$$\sum_{n=1}^N q = \sum_{n=1}^N t_n$$

N · q *N*, *is # timer* *t_n = 1*

LDA: Maximum Likelihood for K=2

- › Likelihood

$$p(\mathbf{t}, \mathbf{X} | q, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \Sigma) = \prod_{n=1}^N \left[q \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \Sigma) \right]^{t_n} \left[(1 - q) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_2, \Sigma) \right]^{1-t_n}$$

- › Log likelihood

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

- › Let's estimate q by solving $\frac{\partial}{\partial q} \ln p(\mathbf{t}, \mathbf{X} | q, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \Sigma) = 0$

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

- › Solve for q :

$$\sum_{n=1}^N \frac{t_n - q}{q(1 - q)} = 0 \stackrel{(q \neq 0, q \neq 1)}{\Rightarrow} \sum_{n=1}^N q = \sum_{n=1}^N t_n \Leftrightarrow q = \frac{1}{N} \sum_{n=1}^N t_n$$

ML solution to $p(C_1) = q$ is the fraction
of points with label $t_n = 1$

LDA: Maximum Likelihood for K=2

- Log likelihood

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

- Let's estimate $\boldsymbol{\mu}_1$ by solving $\frac{\partial}{\partial \boldsymbol{\mu}_1} \ln p(\mathbf{t}, \mathbf{X} | q, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \Sigma) = \mathbf{0}^T$

$$\begin{aligned}\frac{\partial}{\partial \boldsymbol{\mu}_1} \ln p(\mathbf{t}, \mathbf{X} | q, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \Sigma) &= \frac{\partial}{\partial \boldsymbol{\mu}_1} \sum_{n=1}^N t_n \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \Sigma) \\ &= -\frac{1}{2} \frac{\partial}{\partial \boldsymbol{\mu}_1} \sum_{n=1}^N t_n (\mathbf{x}_n - \boldsymbol{\mu}_1)^T \Sigma^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_1)\end{aligned}$$

LDA: Maximum Likelihood for K=2

- Log likelihood

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

- Let's estimate $\boldsymbol{\mu}_1$ by solving $\frac{\partial}{\partial \boldsymbol{\mu}_1} \ln p(\mathbf{t}, \mathbf{X} | q, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \Sigma) = \mathbf{0}^T$

$$\begin{aligned}\frac{\partial}{\partial \boldsymbol{\mu}_1} \ln p(\mathbf{t}, \mathbf{X} | q, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \Sigma) &= \frac{\partial}{\partial \boldsymbol{\mu}_1} \sum_{n=1}^N t_n \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \Sigma) \\ &= -\frac{1}{2} \frac{\partial}{\partial \boldsymbol{\mu}_1} \sum_{n=1}^N t_n (\mathbf{x}_n - \boldsymbol{\mu}_1)^T \Sigma^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_1) \quad (\Sigma \text{ is symmetric}) \\ &\quad \sum_{n=1}^N t_n (\mathbf{x}_n - \boldsymbol{\mu}_1)^T \Sigma^{-1}\end{aligned}$$

- Solve for $\boldsymbol{\mu}_1$:

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

(Σ is pos. def.)

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

①

LDA: Maximum Likelihood for K=2

- Log likelihood

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

- Let's estimate $\boldsymbol{\mu}_1$ by solving $\frac{\partial}{\partial \boldsymbol{\mu}_1} \ln p(\mathbf{t}, \mathbf{X} | q, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \Sigma) = \mathbf{0}^T$

$$\begin{aligned}\frac{\partial}{\partial \boldsymbol{\mu}_1} \ln p(\mathbf{t}, \mathbf{X} | q, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \Sigma) &= \frac{\partial}{\partial \boldsymbol{\mu}_1} \sum_{n=1}^N t_n \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \Sigma) \\ &= -\frac{1}{2} \frac{\partial}{\partial \boldsymbol{\mu}_1} \sum_{n=1}^N t_n (\mathbf{x}_n - \boldsymbol{\mu}_1)^T \Sigma^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_1) \stackrel{(\Sigma \text{ is symmetric})}{=} \sum_{n=1}^N t_n (\mathbf{x}_n - \boldsymbol{\mu}_1)^T \Sigma^{-1}\end{aligned}$$

- Solve for $\boldsymbol{\mu}_1$:

$$\sum_{n=1}^N t_n (\mathbf{x}_n - \boldsymbol{\mu}_1)^T \Sigma^{-1} = \mathbf{0}^T \quad (\Sigma \text{ is pos. def.} \Rightarrow)$$

$\boldsymbol{\mu}_1, \text{ML} = \frac{1}{N_1} \sum_{n=1}^{N_1} t_n \mathbf{x}_n$

average value for \mathbf{x}_n for the cases $t_n = 1$

$$\boldsymbol{\mu}_2, \text{ML} = \frac{1}{N_2} \sum_{n=1}^{N_2} (1 - t_n) \mathbf{x}_n$$

LDA: Maximum Likelihood for K=2

- Log likelihood

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

- Let's estimate Σ by solving $\frac{\partial}{\partial \Sigma} \ln p(\mathbf{t}, \mathbf{X} | q, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \Sigma) = \mathbf{O}$

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

- Solve $\frac{\partial}{\partial \Sigma} \ln p(\mathbf{t}, \mathbf{X} | q, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \Sigma) = \mathbf{O}$ for $\boldsymbol{\mu}_1$ (Bishop 4.2.2): *+ Notes on Canvas*

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

Σ_1 Σ_2
sample cov of class 1 *.. class 2*
 Σ_1 *Σ_2*

LDA: Maximum Likelihood for K=2

The ML solutions:

$$\boldsymbol{\mu}_{1,\text{ML}} = \frac{1}{N_1} \sum_{n=1}^N t_n \mathbf{x}_n$$

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

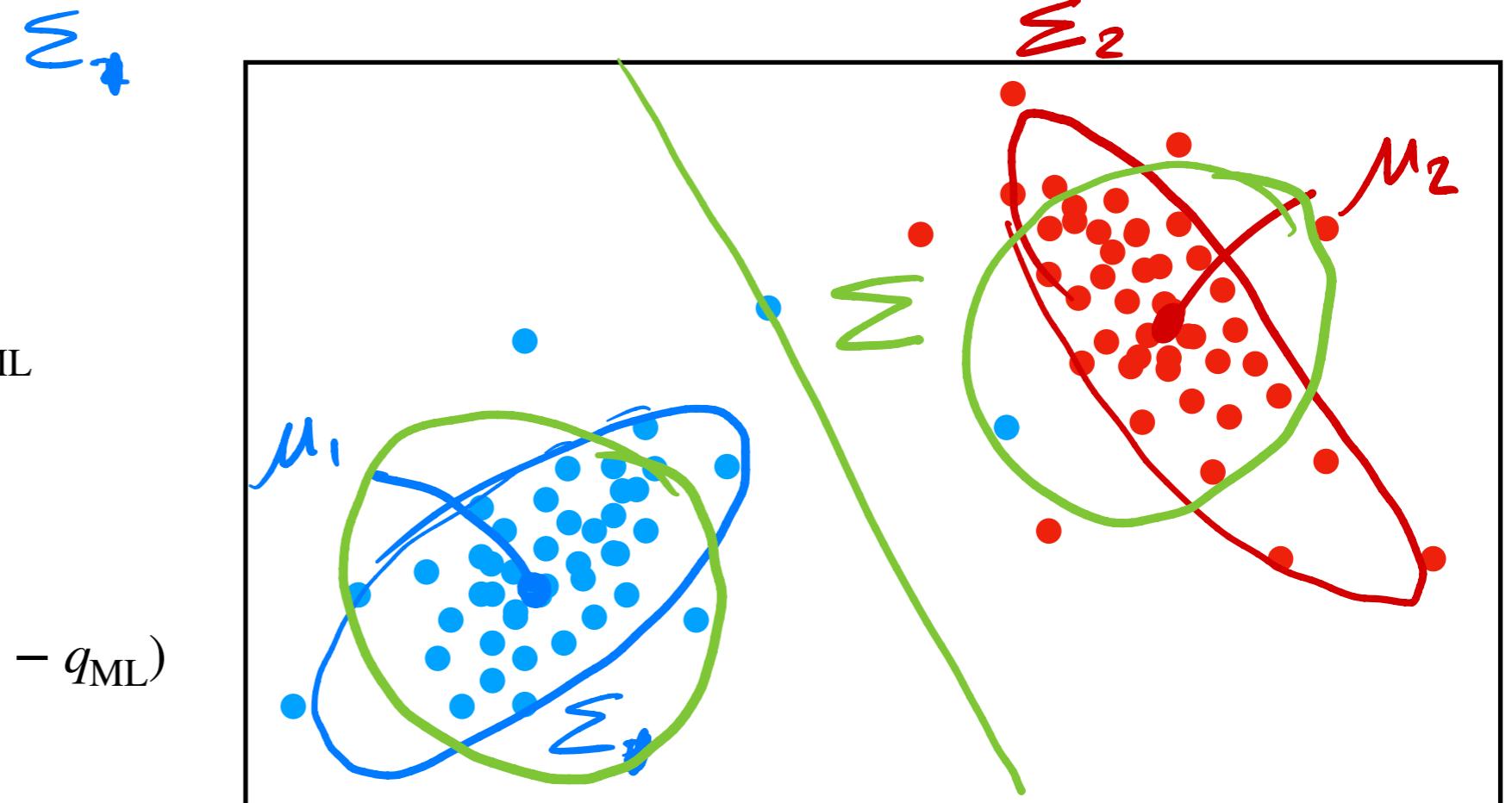
$$q_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N t_n = \frac{N_1}{N}$$

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

For the joint probabilities:

$$\begin{aligned} p(\mathbf{x}, C_1) &= p(\mathbf{x} | C_1)p(C_1) \\ &= \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_{1,\text{ML}}, \Sigma) q_{\text{ML}} \end{aligned}$$

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



LDA: prediction for K=2

- For new datapoint \mathbf{x}' :

derived before

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

lin model given by ML

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

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

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

- Disadvantage of LDA:

► Gaussian distribution is sensitive to outliers

(see averaging)

► Linearity/handcrafted features restrict application

X or use $\phi(\mathbf{x})$

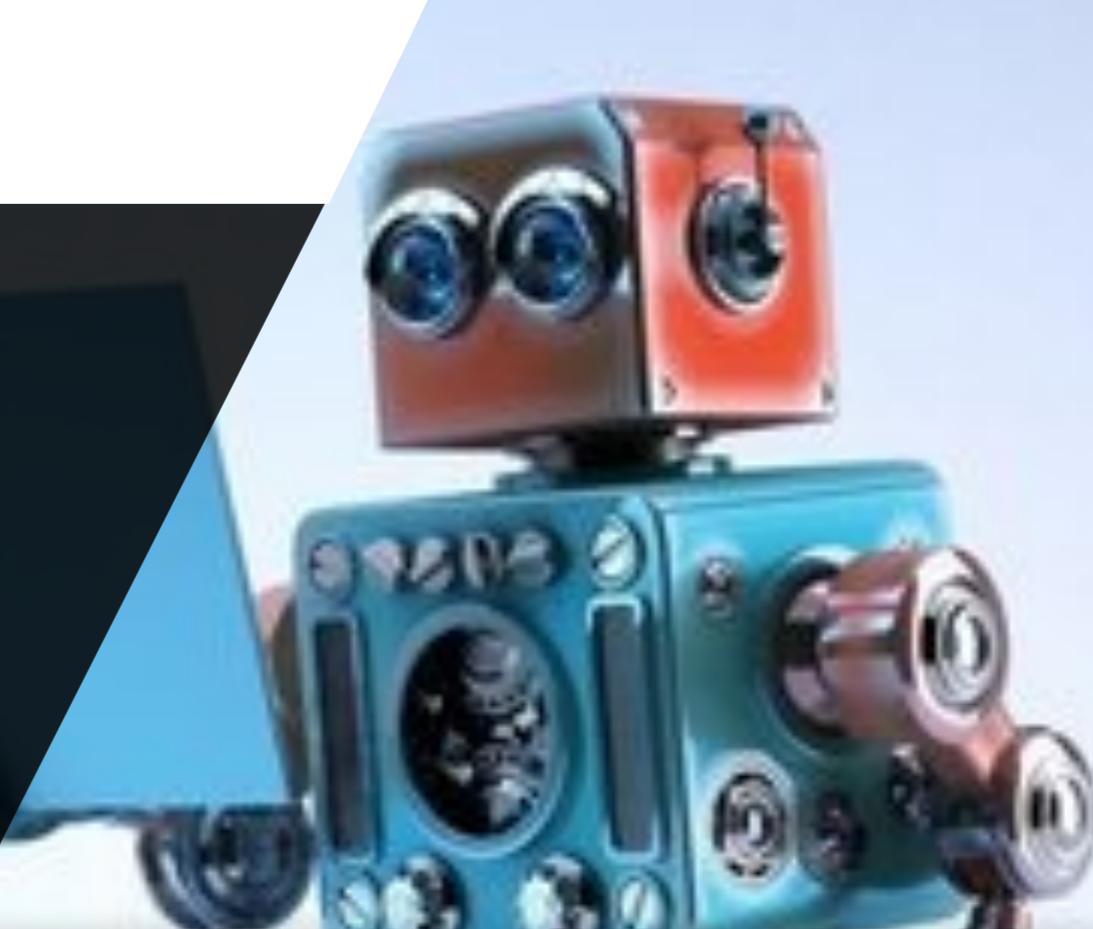
► Maximum likelihood is prone to overfitting

no regularization

Machine Learning 1

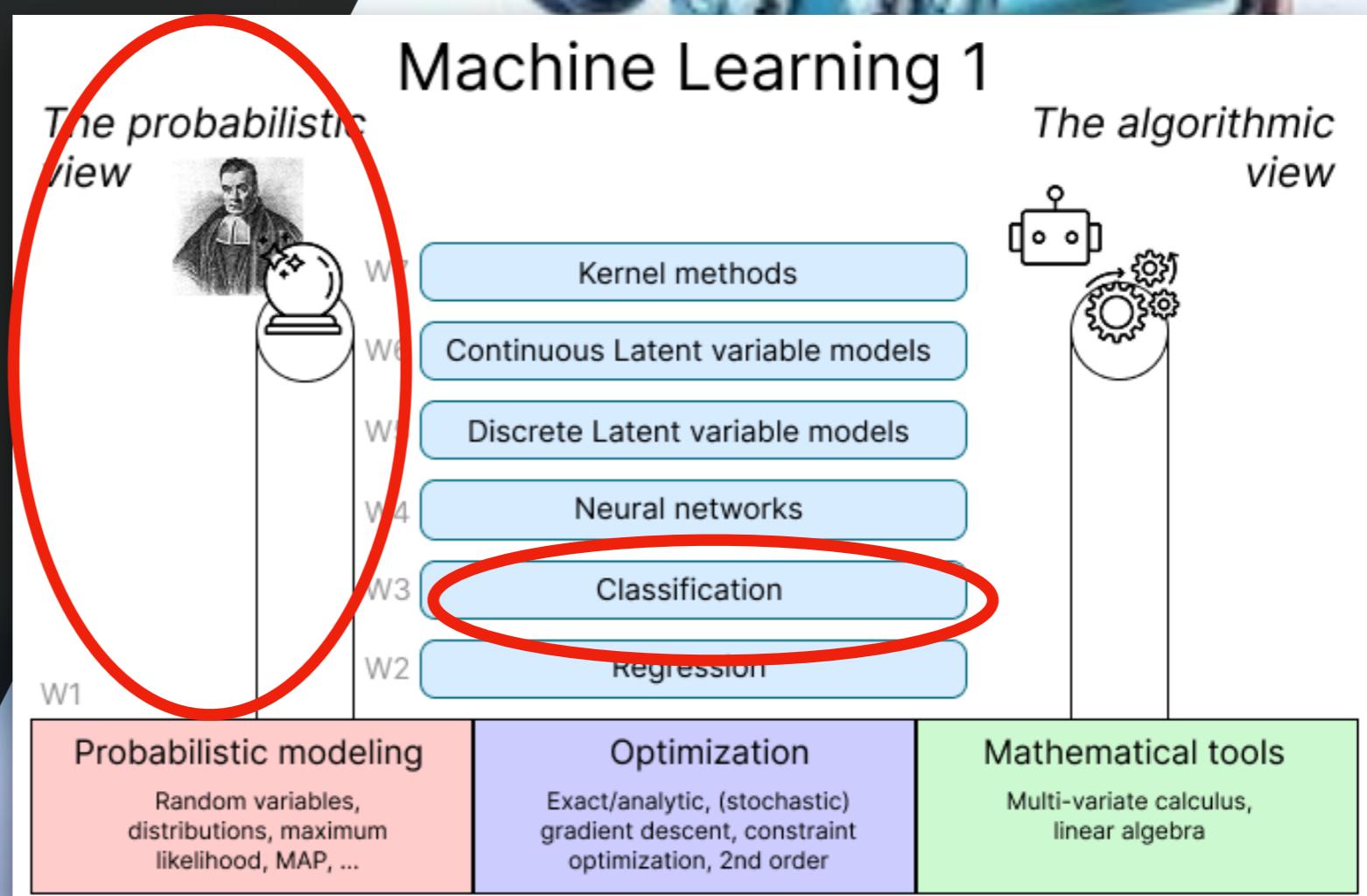
Lecture 6.2 - Supervised Learning

Classification - Probabilistic Generative Models - Discrete Variables (Naive Bayes)



Erik Bekkers

(Bishop 4.2.3)



*Slide credits: Patrick Forré and
Rianne van den Berg*

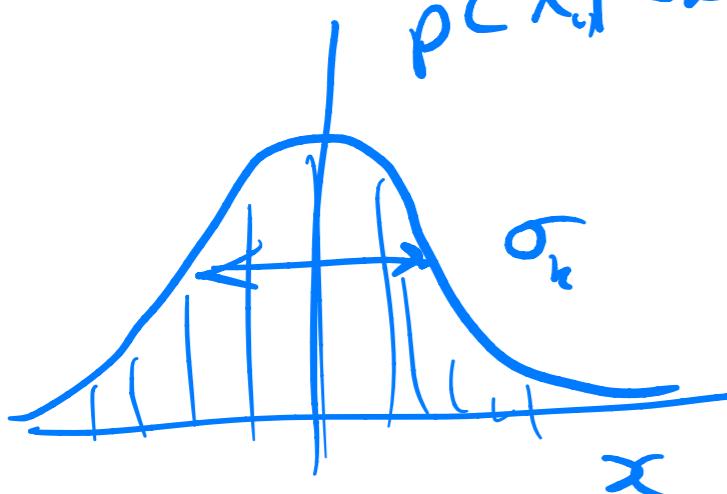
Image credit: Kirillm | Getty Images

Probabilistic Generative Models: Discrete ($K = 2$)

- Input: discrete feature vectors $\mathbf{x}_n = (x_1, \dots, x_D)^T = (0, 1, 1, 0, \dots)^T$
 - $x_i \in \{0, 1\}$
- For D-dimensional input # of parameters for $p(\mathbf{x} | C_k)$ per class: $\underline{\mathbf{x}}_n \in \{0, 1\}^D$ ask after drawing? $2^D - 1$

(D=1)

Before $\mathbf{x} \in \mathbb{R}$

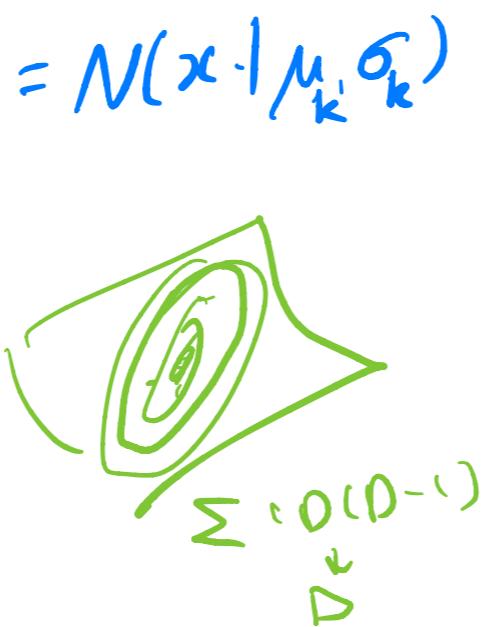


$$\mu_k$$

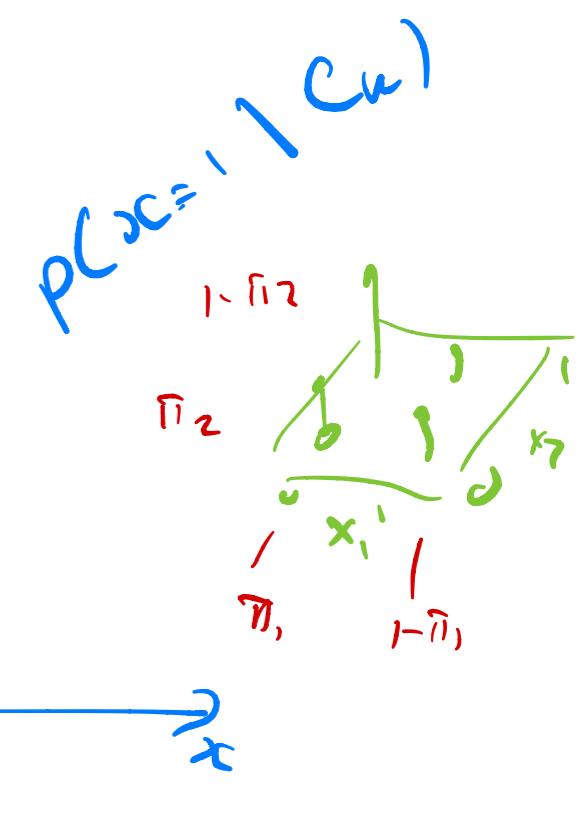
A Gaussian $N(x | \mu, \sigma)$ assigns prob to all possible $x \in \mathbb{R}$

(D=1)

Now



$$p(x_{>0} | C_k)$$



assign probs to all possible cases separately

Probabilistic Generative Models: Discrete ($K = 2$)

- Input: discrete feature vectors $\mathbf{x}_n = (x_1, \dots, x_D)^T$
 - $x_i \in \{0,1\}$
- For D-dimensional input # of parameters for $p(\mathbf{x} | C_k)$ per class: $2^D - 1$
- Naive Bayes assumption: feature values are treated as independent when conditioned on class C_k !

$$p(\mathbf{x} | C_k) = \prod_{i=1}^D p(x_i | C_k) = \prod_{i=1}^D (\pi_{k,i})^{x_i} (1 - \pi_{k,i})^{1-x_i}$$
$$p(x_1, x_2, \dots, x_D | C_k) = p(x_1 | C_k) p(x_2 | C_k) \cdots p(x_D | C_k)$$

with class conditional parameters π_{ki} .

params
= D

Bernoulli: $p(x_i = 1 | C_k) = \bar{\pi}_{k,i}$

Probabilistic Generative Models: Discrete ($K = 2$)

- Input: discrete feature vectors $\mathbf{x}_n = (x_1, \dots, x_D)^T$
 - $x_i \in \{0,1\}$
- For D-dimensional input # of parameters for $p(\mathbf{x} | C_k)$ per class: $2^D - 1$
- Naive Bayes assumption: feature values are treated as independent when conditioned on class C_k !

$$p(\mathbf{x} | C_k) = \prod_{i=1}^D p(x_i | C_k) = \prod_{i=1}^D \pi_{ki}^{x_i} (1 - \pi_{ki})^{1-x_i}$$

with class conditional parameters π_{ki} .

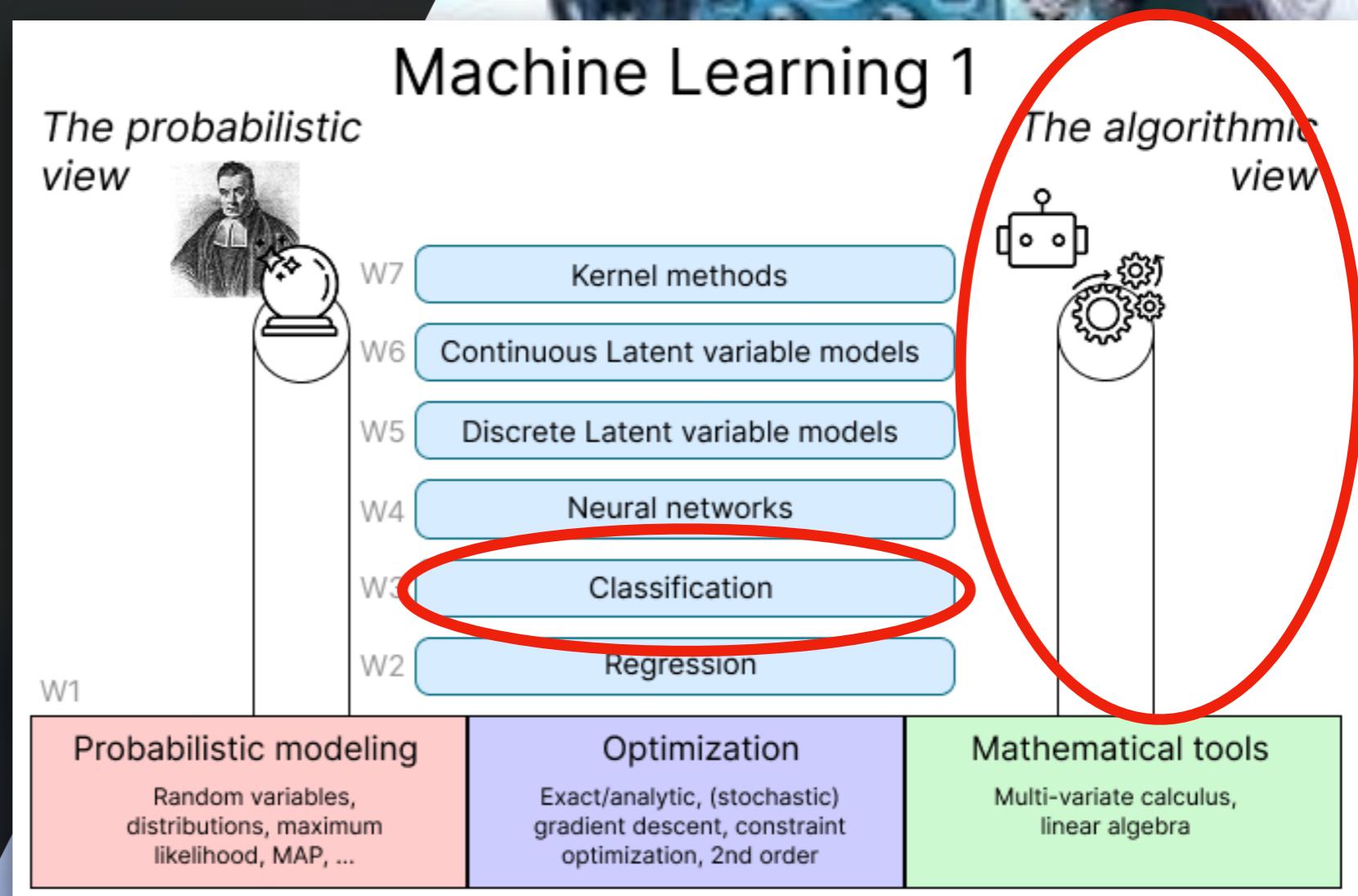
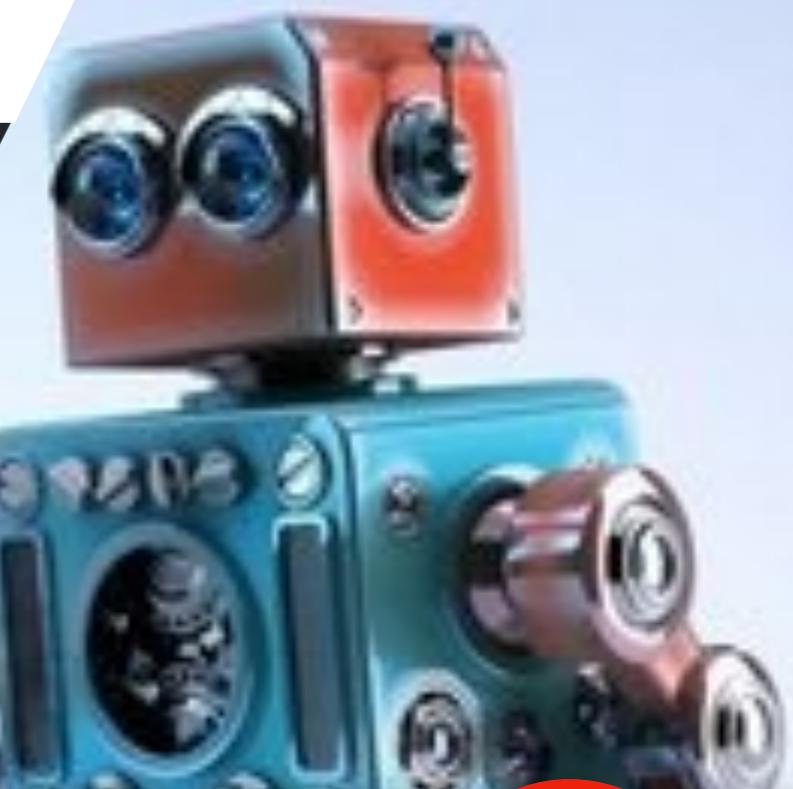
- Class posterior: $p(C_k | \mathbf{x}) = \frac{\exp(a_k(\mathbf{x}))}{\sum_{j=1}^K \exp(a_j(\mathbf{x}))}$ and prior model parameters $p(C_k) = q_k$
 - $a_k = \ln p(\mathbf{x} | C_k) p(C_k) = \ln p(\mathbf{x} | C_k) + \ln p(C_k)$
- then simplifies to:
- $$= \sum_{i=1}^D x_i \ln \pi_{ki} + (1 - x_i) \ln (1 - \pi_{ki}) + \ln q_k$$

Machine Learning 1

Lecture 6.3 - Supervised Learning
Classification - **Discriminative Models**

Erik Bekkers

(Bishop 4.0, 4.1.1, 4.1.2)



Classification Strategies

- ▶ Discriminant functions

Direct mapping of input to target:

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

- ▶ Probabilistic discriminative models

Posterior class probabilities:

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

- ▶ Probabilistic generative models

Class-conditional densities:

means we can sample from the learned data distribution

Prior class probabilities:

$$p(C_k)$$

LDA

Discriminant Functions: Two Classes

- ▶ Input: $\mathbf{x} \in \mathbb{R}^D$
- ▶ Targets: $t \in \{C_1, C_2\}$ = $t=0$, $t=1$

Discriminant functions:

- ▶ Direct mapping of input to target

We now consider generalized Linear Models (GLM)

- ▶ $y(\mathbf{x}, \tilde{\mathbf{w}}) = f(\tilde{\mathbf{w}}^T \boldsymbol{\phi}(\mathbf{x}))$
- ▶ with $\boldsymbol{\phi}(\mathbf{x}) = (1, \phi_1(\mathbf{x}), \dots, \phi_{M-1}(\mathbf{x}))^T$
- ▶ Decision boundary: $y(\mathbf{x}, \tilde{\mathbf{w}}) = \text{const}$

"so, why not directly model such a y ?"

Discriminant Functions: Two Classes

- Simplest discriminant function ($f(\mathbf{x}) = x$, canonical basis $\phi_0(\mathbf{x}) = 1, \phi_j(\mathbf{x}) = x_j, j = 1, \dots, D$)

$$y(\mathbf{x}, \tilde{\mathbf{w}}) = \mathbf{w}^T \mathbf{x} + w_0$$

- Decision boundary: $y(\mathbf{x}, \tilde{\mathbf{w}}) = 0$

$$\underline{\mathbf{w}}^T \underline{\mathbf{x}} + w_0 = 0$$

- Consider 2 datapoints \mathbf{x}_A and \mathbf{x}_B on decision boundary

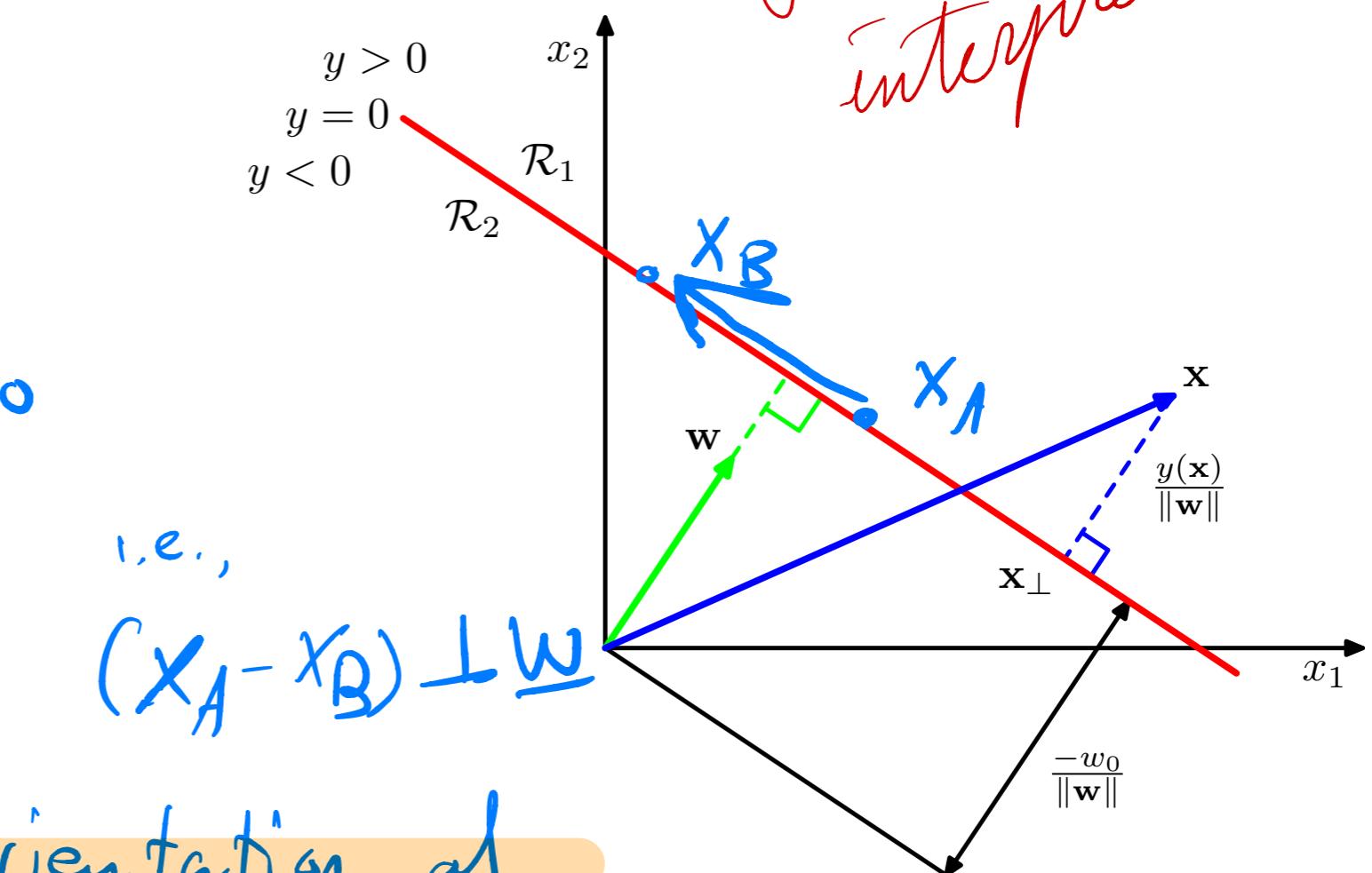
$$\underline{\mathbf{w}}^T \underline{\mathbf{x}}_A + w_0 = \underline{\mathbf{w}}^T \underline{\mathbf{x}}_B + w_0$$

$$\underline{y(\mathbf{x}_A) = y(\mathbf{x}_B) = 0}$$

$$\underline{\mathbf{w}}^T (\underline{\mathbf{x}}_B - \underline{\mathbf{x}}_A) = 0$$

Still useful to establish some intuition.

Geometric interpretation



$\underline{\mathbf{w}}$ determines orientation of
"normal vector of the hyperplane" decision boundary

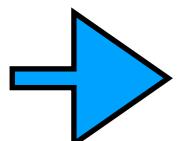
Discriminant Functions: Two Classes

- Take \mathbf{x}' a point on a decision surface: $y(\mathbf{x}') = 0$

$$\mathbf{w}^T \mathbf{x}' + w_0 = 0 \Leftrightarrow \mathbf{w}^T \mathbf{x}' = -w_0$$

- normal vector*
- Normal (signed) distance d from origin to decision surface

$$d = \frac{\mathbf{w}^T \mathbf{x}'}{\|\mathbf{w}\|} = \frac{-w_0}{\|\mathbf{w}\|}$$



- Normal (signed) distances r from general \mathbf{x} to decision surface:

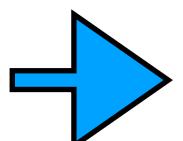
Let $\mathbf{x} = \mathbf{x}_\perp + r \frac{\mathbf{w}}{\|\mathbf{w}\|}$

Then, $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 = \cancel{\mathbf{w}^T \mathbf{x}_\perp} + r \frac{\mathbf{w}^T \mathbf{w}}{\|\mathbf{w}\|} + w_0 = r \cdot \|\mathbf{w}\|^2 + w_0$

$$y(\mathbf{x}_\perp) = \mathbf{w}^T \mathbf{x}_\perp + w_0 = 0$$

$$r \cdot \|\mathbf{w}\|^2 + w_0 = r \cdot \|\mathbf{w}\|^2$$

$$y(\mathbf{x}_\perp) = 0$$



Distance r of a point to the d.b. is $r = \frac{w_0}{\|\mathbf{w}\|}$

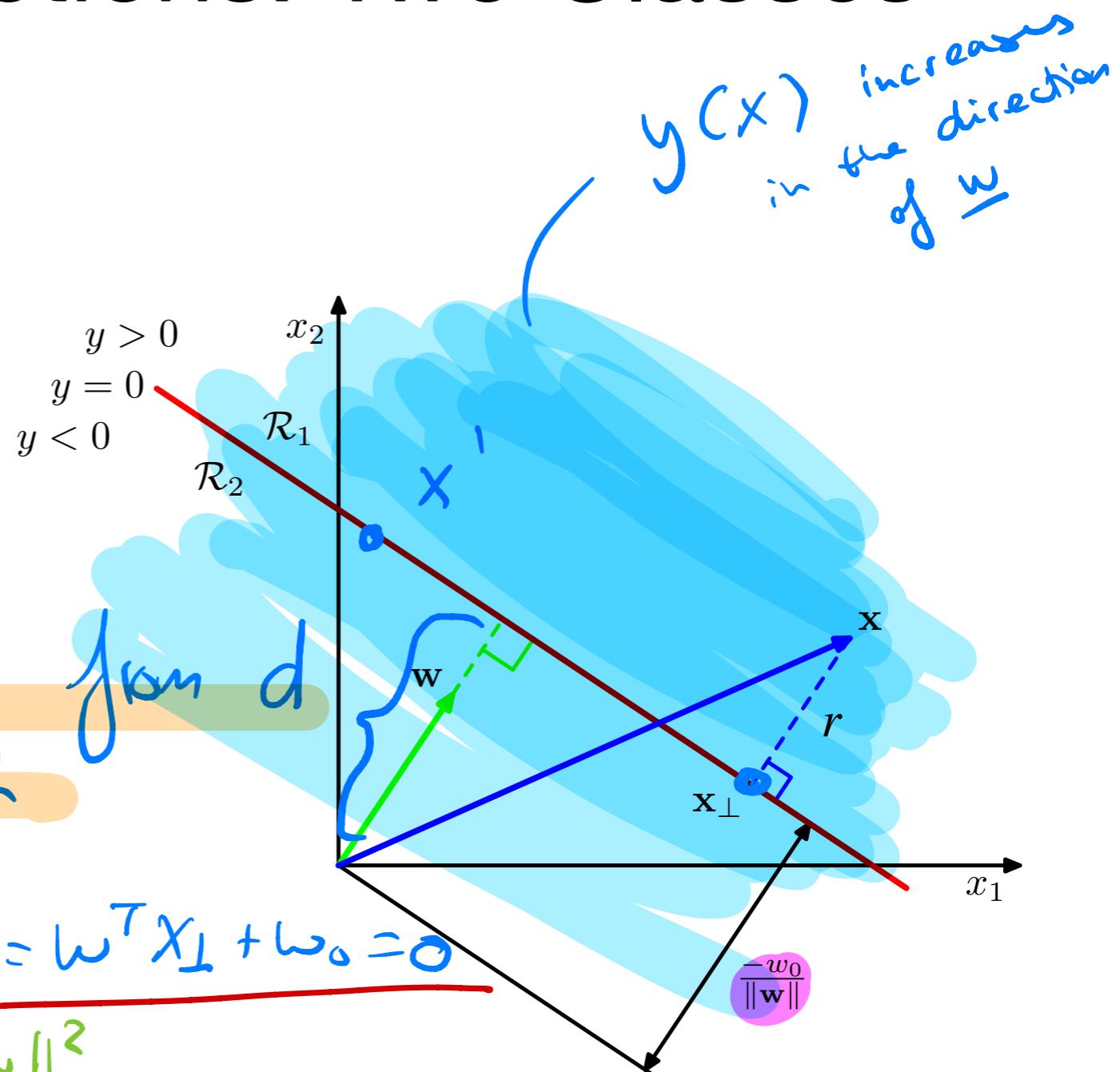


Figure: 2-class decision boundary (Bishop 4.1) 6

So the distance
of a point \underline{x} to
the decision boundary
is proportional to $y(\underline{x})$

$$r = \frac{y(\underline{x})}{\|\underline{w}\|}$$

Discriminant Functions: Multiple Classes

- K -class discriminant

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

like assign to max prob

- Assign \mathbf{x} to C_k if $\forall j \neq k : y_k(\mathbf{x}) > y_j(\mathbf{x})$

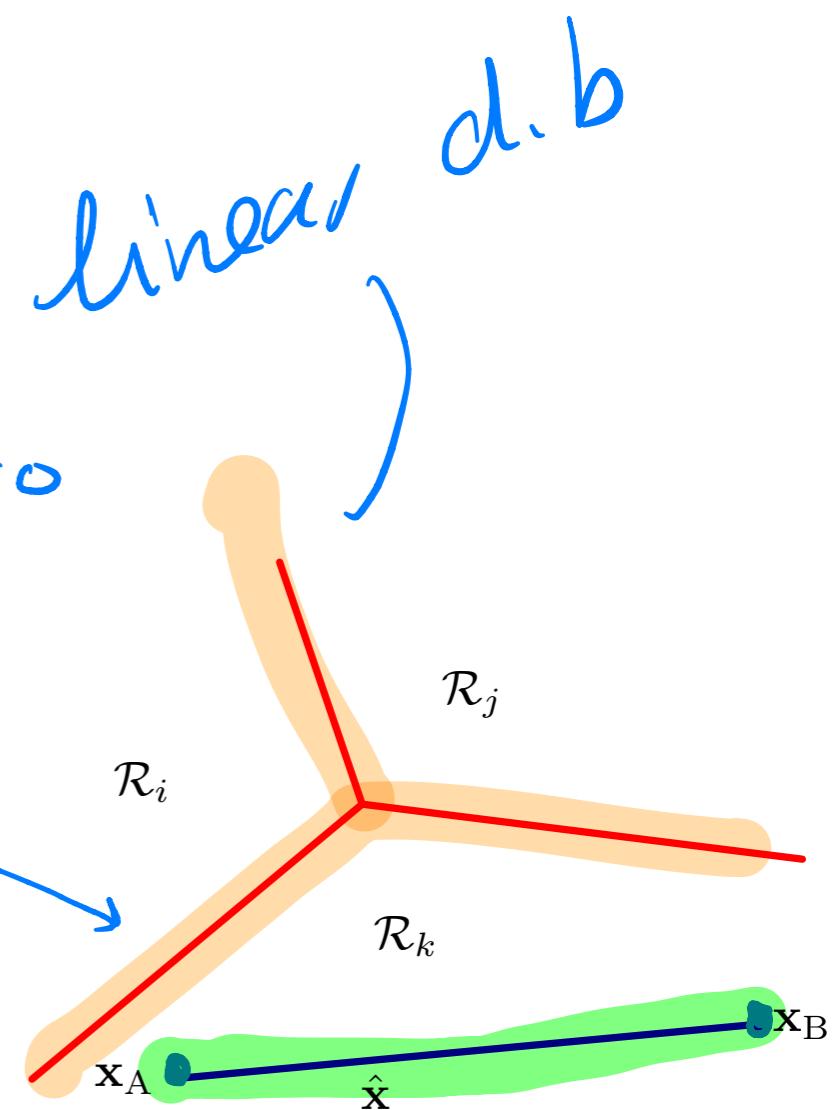
- Decision boundary between R_k and R_j :

$$y_k(\mathbf{x}) = y_j(\mathbf{x}) \iff y_k(\mathbf{x}) - y_j(\mathbf{x}) = 0$$

$$(\underline{w}_k - \underline{w}_j) \mathbf{x} + (w_{k0} - w_{j0}) = 0$$

thus linear

- Decision regions (for GLM) are convex



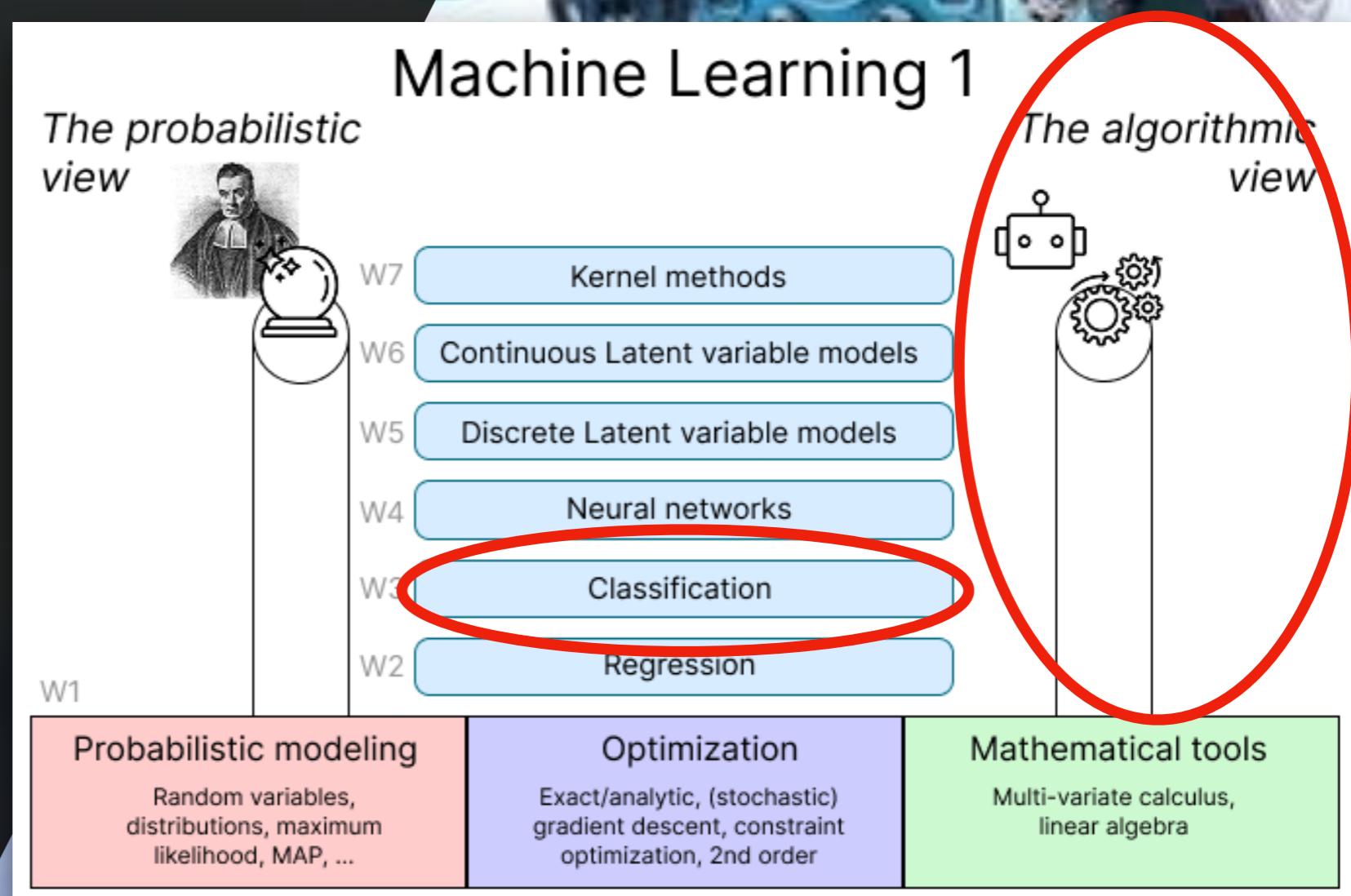
Machine Learning 1

Lecture 6.4 - Supervised Learning

Classification? - Discriminative Models - Least Squares Regression

Erik Bekkers

(Bishop 4.1.3)



Least Squares for Classification

- Each class C_k has its own linear model:

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

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

- Matrix $\tilde{\mathbf{W}}$, column k contains: $\tilde{w}_k = \text{Concat}(w_{k0}, \mathbf{w}_k) \in \mathbb{R}^M$

- Input vector: $\tilde{\mathbf{x}} = \text{Concat}(1, \mathbf{x}) \in \mathbb{R}^M$

- Output/prediction vector: $\mathbf{y}(\mathbf{x}) = \begin{pmatrix} y_1(\mathbf{x}) \\ y_2(\mathbf{x}) \\ \vdots \\ y_K(\mathbf{x}) \end{pmatrix} = \tilde{\mathbf{W}}^T \tilde{\mathbf{x}} \in \mathbb{R}^K$

- Classification:** Assign \mathbf{x} to class C_K if $k = \operatorname{argmax}_j y_j(\mathbf{x})$

Least Squares for Classification (II)

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

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

- Use sum-of-squares regression error function

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

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

- 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: $\mathbf{y}_{\text{LS}}(\mathbf{x}) =$

Least Squares for Classification: Problems

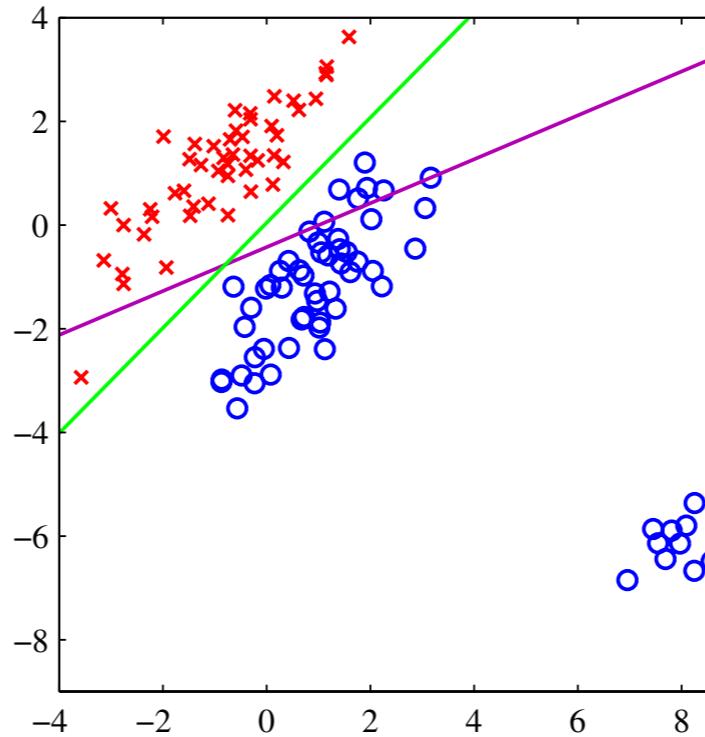
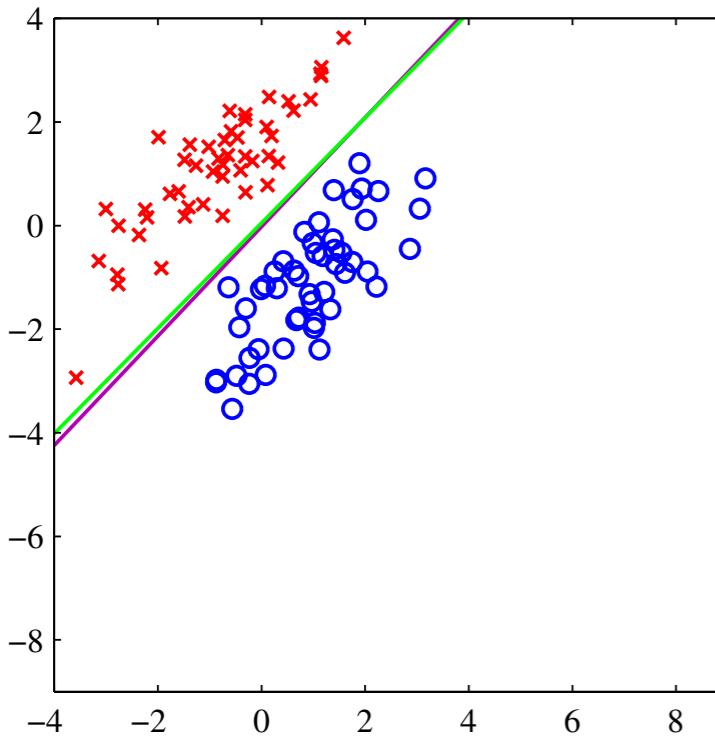


Figure: least squares is very sensitive to outliers (Bishop 4.4)

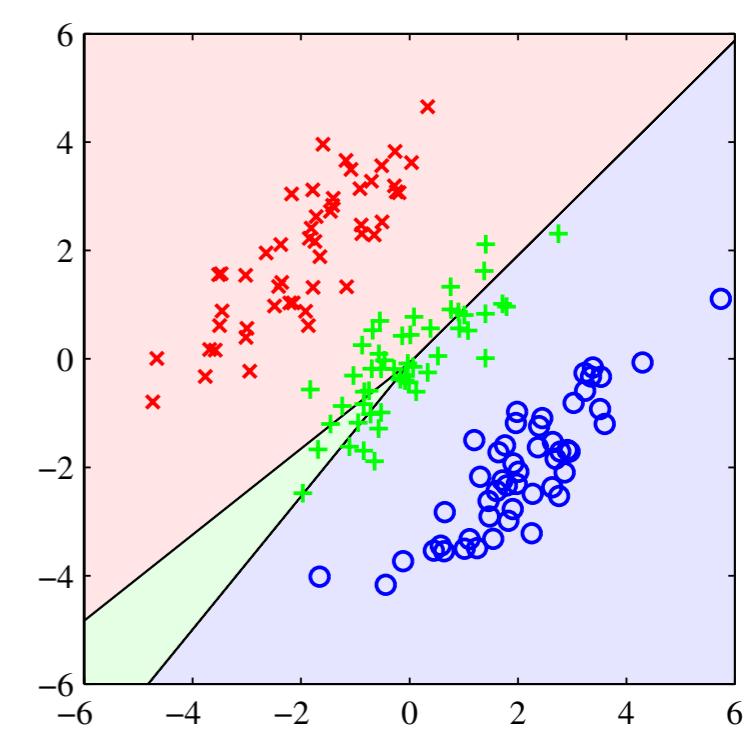


Figure: masking for least squares for $K > 2$ (Bishop 4.5)

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 $\mathbf{y}_{\text{LS}}(\mathbf{x})$ are not probabilities!

- $y_k(\mathbf{x})$

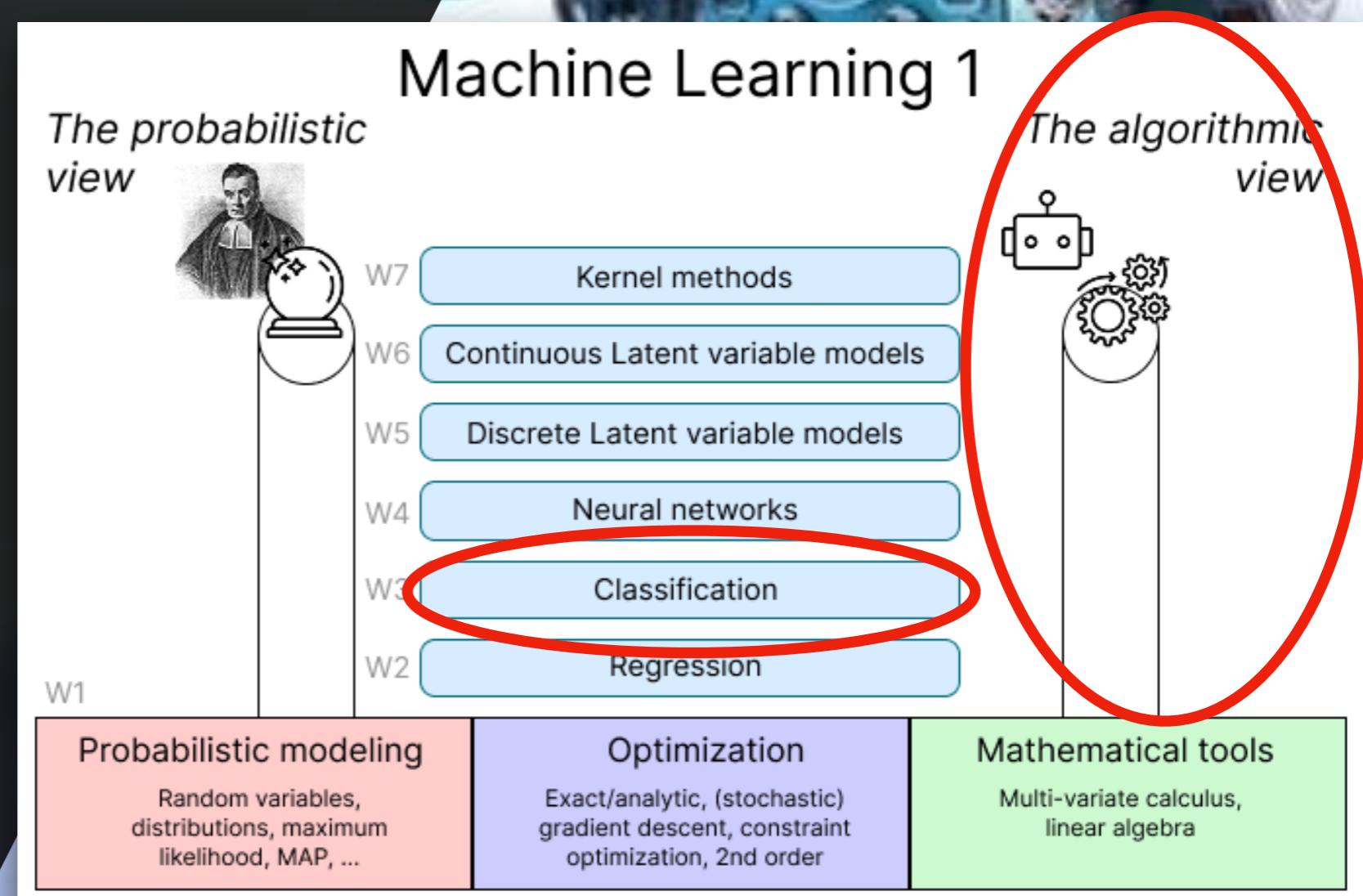
- if $\sum_{k=1}^K t_k = 1$

Machine Learning 1

Lecture 6.5 - Supervised Learning
Classification - Discriminative Models -
The Perceptron

Erik Bekkers

(Bishop 4.1.7)



The Perceptron Algorithm

- ▶ Input: $\mathbf{x} \in \mathbb{R}^D$
- ▶ Targets: $t \in \{C_1, C_2\}$
- ▶ Prediction: $y(\mathbf{x}) = f(\mathbf{w}^T \phi(\mathbf{x}))$, with $f(a) = \begin{cases} 1 & , a \geq 0 \\ -1 & , a < 0 \end{cases}$
- ▶ Class decisions: assign \mathbf{x} to class C_1 if
(and to C_2 if)
- ▶ For correct classification: find \mathbf{w} such that for all (\mathbf{x}_n, t_n) :
- ▶ Perceptron criterion: $E_P(\mathbf{w}) = - \sum_{n \in \mathcal{M}} \mathbf{w}^T \phi(\mathbf{x}_n) t_n$

with $\mathcal{M} =$

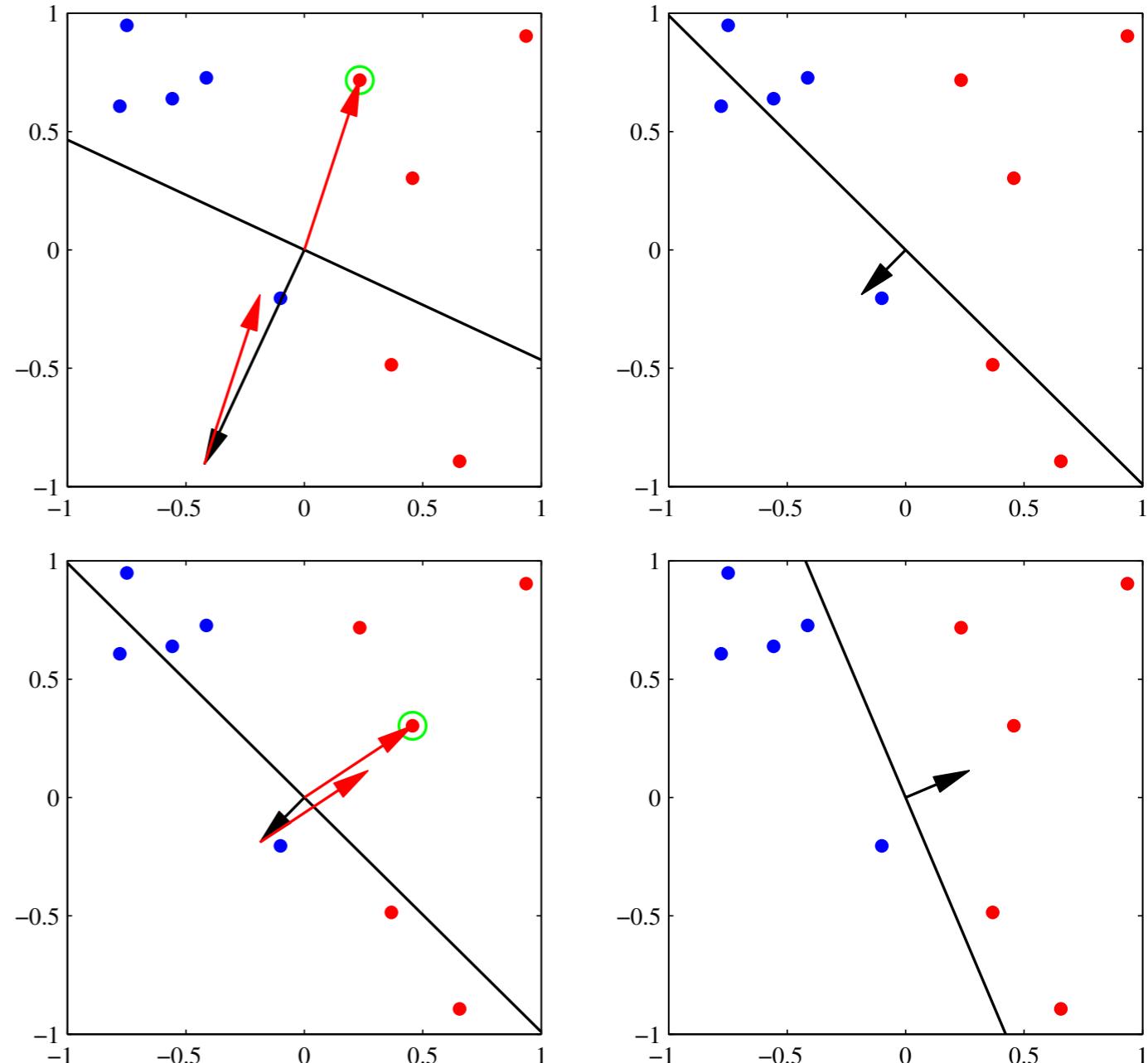
Perceptron: Stochastic Gradient Descent

- $$E_P(\mathbf{w}) = - \sum_{n \in \mathcal{M}} \mathbf{w}^T \phi(\mathbf{x}_n) t_n$$

$$= - \sum_{n \in \mathcal{M}} E_n(\mathbf{w})$$
- Stochastic Gradient Descent (SGD). For each misclassified \mathbf{x}_n :

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_n(\mathbf{w})^T$$

=



- If \mathbf{X} is linearly separable, then perceptron SGD will converge

Figure: SGD for perceptron criterion (Bishop 4.7)
for \mathbf{x}_n in C_1 : add $\phi(\mathbf{x}_n)$ to \mathbf{w}
for \mathbf{x}_n in C_2 : subtract $\phi(\mathbf{x}_n)$ from \mathbf{w} .

Problems: Perceptron

- Perceptron only works for 2 classes
- There might be many solutions depending on the initialization of \mathbf{w} and on the order in which data is presented in SGD
- If dataset is not linearly separable, the perceptron algorithm will not converge.
- Based on linear combination of fixed basis functions.