

# Advanced Probabilistic Machine Learning SSY316

## Linear Models for Classification

Alexandre Graell i Amat

[alexandre.graell@chalmers.se](mailto:alexandre.graell@chalmers.se)

<https://sites.google.com/site/agraellamat>

November 9 and 10, 2023



**CHALMERS**

# Classification

**Goal:** Given a training set  $\mathcal{D} = \{(\mathbf{x}_1, t_1), \dots, (\mathbf{x}_N, t_N)\}$ , with  $t \in \{C_1, \dots, C_K\}$ , assign (classify) a new example  $\mathbf{x}$  to one of the classes  $C_i$ .

- For convenience:  $t \in \{1, \dots, K\}$

Classification divides input space into **decision regions**  $\mathcal{R}_i$ ,  $i \in [K]$  (points in  $\mathcal{R}_i$  assigned to class  $C_i$ ).

**Linear models:** Decision boundaries (surfaces) are linear functions of  $\mathbf{x}$  ( $(d-1)$ -dimensional hyperplanes).

## Classification: Examples

**Data set:** annual income, credit card balance, whether or not the person has defaulted

**Aim:** Predict whether a person will default or not

**Data set:** Blood tests, X-rays, MRI scans

**Aim:** Predict whether a person has cancer or not

# Classification as a supervised learning problem

## Binary classification:

$$\mathbf{x} \in C_1 \quad \text{if } t = 1$$

$$\mathbf{x} \in C_2 \quad \text{if } t = 0 \text{ (or } t = -1)$$

e.g.,  $C_1 \equiv \text{cancer}$ ,  $C_2 \equiv \text{no cancer}$

## Multi-class classification ( $K > 2$ ):

Label a point  $\mathbf{x} \in C_k$  with a  $K \times 1$  one-hot vector  $\mathbf{t}$  with a single one at position  $k$  and zeroes elsewhere

$K = 5$ , class  $C_2$ :

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

# Three distinct modeling approaches

**Discriminative deterministic models:** Model directly the **deterministic** mapping between  $\mathbf{x}$  and  $t$  via a **parametrized function**  $\mu(\mathbf{x}, \mathbf{w})$  (**discriminant function**).

**More powerful approach:** model  $p(C_k|\mathbf{x})$  in an inference stage, then use this distribution to make optimal decisions.

**Discriminative probabilistic models:** Model  $p(C_k|\mathbf{x})$  via a parametrized model.

**Generative probabilistic models:** Model  $p(\mathbf{x}, t)$  specifying  $p(t)$  ( $p(C_k)$ ) and  $p(\mathbf{x}|t)$  ( $p(\mathbf{x}|C_k)$ ). Then,

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

# Why a probabilistic approach?

## Compensating for class priors (Cancer):

- Assume prevalence 1/1000
- Naïve classifier that assigns every new case result to “no cancer” class: 99.9% accuracy!
- An unbalanced data set will lead to learning model that likely won't generalize well
- A balanced data set will lead to a more accurate model ... but we need to compensate for modification to training data!

## Combining models (Cancer):

- 3 models based on blood tests, X-rays, MRIs
- If models provide posterior probabilities, they can be properly combined,

$$\begin{aligned} p(C_k | \mathbf{x}_b, \mathbf{x}_x, \mathbf{x}_m) &\propto p(\mathbf{x}_b, \mathbf{x}_x, \mathbf{x}_m | C_k) \\ &= p(\mathbf{x}_b | C_k) p(\mathbf{x}_x | C_k) p(\mathbf{x}_m | C_k) p(C_k) \end{aligned}$$

# Activation function

Linear regression:  $\mu(\mathbf{x}, \mathbf{w})$  a linear function of  $\mathbf{w}$

- Simplest case:  $\mu(\mathbf{x}, \mathbf{w}) = \mathbf{w}^T \mathbf{x} + w_0$

Classification problem: want to predict discrete class labels, or posterior probabilities.

Generalized linear models:

- Generalization of the linear regression model:  $\mu(\mathbf{x}, \mathbf{w}) = a(\mathbf{w}^T \mathbf{x} + w_0)$
- Decision boundaries: constant  $\mu(\mathbf{x}) \rightarrow$  linear functions of  $\mathbf{x}$ !

# Discriminative deterministic models

**Discriminative deterministic models:** Model directly the **deterministic** mapping between  $\mathbf{x}$  and  $t$  via a **parametrized function**  $\mu(\mathbf{x}, \mathbf{w})$  (**discriminant function**).

Binary classification:

- $\hat{t}(\mathbf{x})$  obtained by applying a **threshold rule** on  $\mu(\mathbf{x}, \mathbf{w})$

**Simplest form:**

$$\begin{aligned}\mu(\mathbf{x}, \mathbf{w}) &= a(\mathbf{w}^\top \mathbf{x} + w_0) \\ &= \mathbf{w}^\top \mathbf{x} + w_0\end{aligned}$$

and

$$\begin{aligned}\hat{t}(\mathbf{x}, \mathbf{w}) &= \text{sign}(\mu(\mathbf{x}, \mathbf{w})) \\ &= \text{sign}(\mathbf{w}^\top \mathbf{x} + w_0)\end{aligned}$$

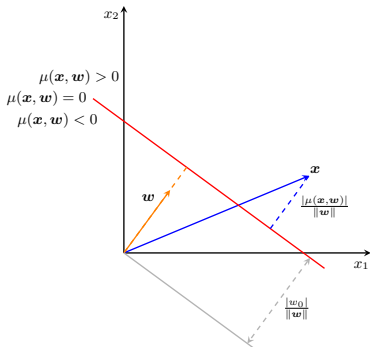
$$\mathbf{x} \in \mathcal{C}_1 \text{ if } \mu(\mathbf{x}, \mathbf{w}) > 0$$

$$\mathbf{x} \in \mathcal{C}_2 \text{ if } \mu(\mathbf{x}, \mathbf{w}) < 0$$



# Geometric interpretation

**Decision rule:** Defines a **hyperplane** that separates the domain points classified as belonging to either of the two classes.



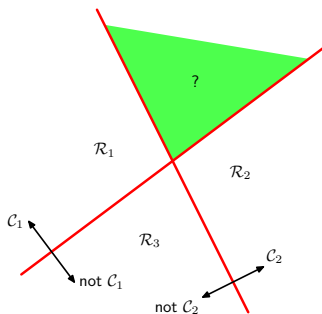
- $\mathbf{w}$  determines **orientation** of decision hyperplane
- $w_0$  determines **location** of decision surface
- **Classification margin:**  
 $|\mu(\mathbf{x}, \mathbf{w})| / \|\mathbf{w}\|$

# Multi-class classification

Combining 2-class discriminant functions:

**One-versus the rest classifier:** Separating points in  $C_k$  from points not in  $C_k$ .

- Leads to regions of input space that are **ambiguously classified**

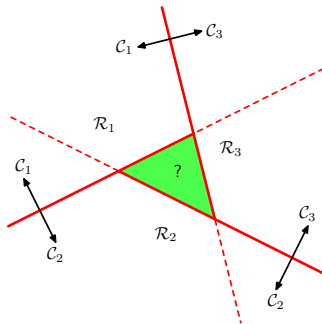


# Multi-class classification

Combining 2-class discriminant functions:

**One-versus-one classifier:**  $\binom{K}{2}$  binary discriminant functions. Each point classified according to a majority vote.

- Leads to regions of input space that are **ambiguously classified**



# Multi-class classification

A single  $K$ -class discriminant with  $K$  linear functions

$$\mu_k(\mathbf{x}) = w_{k,0} + \mathbf{w}_k^\top \mathbf{x}$$

and

$$\hat{t}(\mathbf{x}) = \arg \max_k \mu_k(\mathbf{x})$$

$\mu_k(\mathbf{x})$  can be interpreted as  $p(t = k|\mathbf{x})$ .

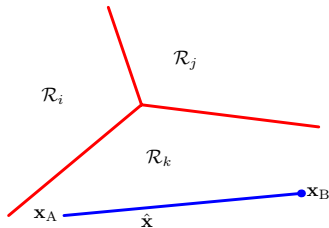
Decision boundary between  $C_k$  and  $C_j$ :

$$\{\mathbf{x} : (\mathbf{w}_k - \mathbf{w}_j)^\top \mathbf{x} + (w_{k,0} - w_{j,0}) = 0\}$$

# Multi-class classification

A single  $K$ -class discriminant with  $K$  linear functions

$$\mu_k(\mathbf{x}) = w_{k,0} + \mathbf{w}_k^\top \mathbf{x}$$



- Decision regions singly connected and convex
- Parameters can be optimized via, e.g., least squares  
( $\mathbf{W}^* = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{T}$ )

# Discriminative probabilistic models

Discriminative probabilistic models: Model  $p(C_k|\mathbf{x})$  via a parametrized model.

- Potentially more powerful than deterministic models since allow to model **sources of uncertainty** in the label assignment to input variables

We will consider:

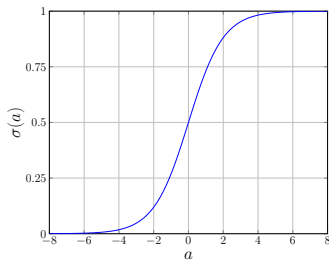
- Models **linear** in the **parameters**
- Applying a nonlinear transformation  $\phi = (\phi_0, \dots, \phi_M)$  to  $\mathbf{x}$ , then work with  $\phi(\mathbf{x})$  as our features

**Decision boundaries** **linear** in the feature space  $\phi$ , **nonlinear** in original space  $\mathbf{x}$ .

# Logistic regression

Binary classification:

$$\begin{aligned} p(\mathbf{t} = C_1 | \phi, \mathbf{x}, \mathbf{w}) &= p(\mathbf{t} = C_1 | \phi(\mathbf{x}), \mathbf{w}) \\ &= \frac{\exp(\mathbf{w}^\top \phi(\mathbf{x}))}{1 + \exp(\mathbf{w}^\top \phi(\mathbf{x}))} \\ &= \frac{1}{1 + e^{-(\mathbf{w}^\top \phi(\mathbf{x}))}} \triangleq \sigma(\mathbf{w}^\top \phi(\mathbf{x})) \end{aligned}$$



For an  $M$ -dimensional feature space  $\phi$ ,  $M$  adjustable parameters.

# Logistic regression

$$p(\mathbf{t} = \mathcal{C}_1 | \phi, \mathbf{w}) = \frac{1}{1 + e^{-(\mathbf{w}^\top \phi(\mathbf{x}))}} \triangleq \sigma(\mathbf{w}^\top \phi(\mathbf{x}))$$

Inference:

Misclassification error minimized by

$$\mathbf{x} \in \mathcal{C}_1 \text{ if } p(\mathcal{C}_1 | \phi, \mathbf{w}) > 1/2$$

$$\mathbf{x} \in \mathcal{C}_2 \text{ if } p(\mathcal{C}_1 | \phi, \mathbf{w}) < 1/2$$

Equivalently,

$$\mathbf{x} \in \mathcal{C}_1 \text{ if } \mathbf{w}^\top \phi(\mathbf{x}) > 0$$

$$\mathbf{x} \in \mathcal{C}_2 \text{ if } \mathbf{w}^\top \phi(\mathbf{x}) < 0$$



# Logistic regression: Learning

ML learning:

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} p(t_{\mathcal{D}} | x_{\mathcal{D}}, \mathbf{w})$$

Likelihood function:

$$\begin{aligned} p(t_{\mathcal{D}} | x_{\mathcal{D}}, \mathbf{w}) &= \prod_{i=1}^N p(t_i | \mathbf{x}_i, \mathbf{w}) = \prod_{i=1}^N p(t_i | \phi(\mathbf{x}_i), \mathbf{w}) \\ &= \prod_{i=1}^N y_i^{t_i} (1 - y_i)^{1-t_i} \end{aligned}$$

with  $y_i = p(t_i = \mathcal{C}_1 | \phi(\mathbf{x}_i), \mathbf{w})$

Negative log-likelihood:

$$-\ln p(t_{\mathcal{D}} | x_{\mathcal{D}}, \mathbf{w}) = - \sum_{i=1}^N \left( t_i \ln(y_i) + (1 - t_i) \ln(1 - y_i) \right)$$

# Logistic regression: Learning

ML learning:

$$\begin{aligned}\mathbf{w}^* &= \arg \min_{\mathbf{w}} -\ln p(t_{\mathcal{D}} | x_{\mathcal{D}}, \mathbf{w}) \\ &= \arg \min_{\mathbf{w}} - \sum_{i=1}^N \left( t_i \ln(y_i) + (1 - t_i) \ln(1 - y_i) \right)\end{aligned}$$

(cross-entropy loss criterion)

No simple formula to find optimal  $\mathbf{w}^*$ .

However:

1. NLL has unique minimum unless classes perfectly separated by hyperplane
2. The negative log-likelihood surface is concave
3. The derivatives and the second derivatives can be easily computed

Efficient methods (SGD, Newton-Raphson) exist.

# Discriminative probabilistic models

Other discriminative probabilistic models:

- Support vector machine
- Decision tree
- Random forest

# Bayesian logistic regression

We want to obtain

$$p(t|\mathcal{D}, \mathbf{x}) = \int \underbrace{p(\mathbf{w}|\mathcal{D})}_{\text{posterior dist. of } \mathbf{w}} p(t|\mathbf{x}, \mathbf{w}) d\mathbf{w}$$

**Bayesian approach:**  $p(t|\mathbf{x}, \mathbf{w})$  associated with each value of  $\mathbf{w}$  weighted by the posterior belief

$$p(\mathbf{w}|\mathcal{D}) = \frac{p(\mathbf{w})p(t_{\mathcal{D}}|x_{\mathcal{D}}, \mathbf{w})}{p(t_{\mathcal{D}}|x_{\mathcal{D}})}$$

# Bayesian logistic regression

We want to compute the full posterior  $p(\mathbf{w}|\mathcal{D})$ .

**Problem:** Exact Bayesian inference for logistic regression **not possible** as exact evaluation of  $p(\mathbf{w}|\mathcal{D})$  **intractable**!

$$\begin{aligned} p(\mathbf{w}|\mathcal{D}) &\propto p(\mathbf{w})p(t_{\mathcal{D}}|x_{\mathcal{D}}, \mathbf{w}) \\ &= p(\mathbf{w}) \prod_{i=1}^N p(t_i|\phi(\mathbf{x}_i), \mathbf{w}) \end{aligned}$$

**Bayesian logistic regression:** More complex than for linear regression models (cannot integrate exactly over  $\mathbf{w}$  since posterior not Gaussian)  $\rightarrow$  need to introduce some **approximations**.

**Idea:** Use **Laplace approximation**

# The Laplace approximation

**Laplace approximation:** Aims to find a Gaussian approximation to a posterior distribution.

Case of a single continuous random variable  $z$ :

$$p(z) = \frac{1}{Z} f(z)$$

$Z$ : (unknown) normalization constant ( $Z = \int f(z) dz$ )

**Goal:** Find a Gaussian approximation  $q(z)$  centered on a mode of  $p(z)$ .

# The Laplace approximation

**Goal:** Find a Gaussian approximation  $q(z)$  centered on a mode of  $p(z)$ .

1. Find a **mode** of  $p(z)$  ( $p'(z_0) = 0$ )

**Observation:** The logarithm of a Gaussian is a quadratic function of the variables.

2. **Taylor series expansion** of  $\ln f(z)$  centered on mode  $z_0$ :

$$\begin{aligned}\ln f(z) &\approx \ln f(z_0) + (z - z_0) \left. \frac{d}{dz} \ln f(z) \right|_{z=z_0} + (z - z_0)^2 \frac{1}{2} \cdot \left. \frac{d^2}{dz^2} \ln f(z) \right|_{z=z_0} \\ &= \ln f(z_0) + (z - z_0)^2 \frac{1}{2} \cdot \left. \frac{d^2}{dz^2} \ln f(z) \right|_{z=z_0} \\ &= \ln f(z_0) - \frac{1}{2} A (z - z_0)^2\end{aligned}$$

with

$$A = - \left. \frac{d^2}{dz^2} \ln f(z) \right|_{z=z_0}$$

# The Laplace approximation

$$\ln f(z) \approx \ln f(z_0) - \frac{1}{2}A(z - z_0)^2$$

with

$$A = - \left. \frac{d^2}{dz^2} \ln f(z) \right|_{z=z_0}$$

We obtain:

$$f(z) \approx f(z_0)e^{-\frac{A}{2}(z-z_0)^2}$$

and normalized distribution

$$\begin{aligned} q(z) &= \left( \frac{A}{2\pi} \right)^{1/2} e^{-\frac{A}{2}(z-z_0)^2} \\ &= \mathcal{N}(z|z_0, A^{-1}) \end{aligned}$$



## The Laplace approximation for higher dimensions

$$p(\mathbf{z}) = \frac{1}{Z} f(\mathbf{z})$$

1. Find a **mode** of  $p(\mathbf{z})$  ( $\nabla f(\mathbf{z}) = 0$ )
2. **Taylor series expansion** of  $\ln f(\mathbf{z})$  centered on mode  $\mathbf{z}_0$ :

$$\ln f(\mathbf{z}) \approx \ln f(\mathbf{z}_0) - \frac{1}{2}(\mathbf{z} - \mathbf{z}_0)^T \mathbf{A}(\mathbf{z} - \mathbf{z}_0)$$

with  $\mathbf{A}$  the  $d \times d$  Hessian matrix (evaluated at  $\mathbf{z} = \mathbf{z}_0$ )

$$\mathbf{A} = -\nabla \nabla \ln f(\mathbf{z}) \Big|_{\mathbf{z}=\mathbf{z}_0} = \begin{pmatrix} \frac{\partial^2 \ln f(\mathbf{z})}{\partial z_1^2} & \frac{\partial^2 \ln f(\mathbf{z})}{\partial z_1 \partial z_2} & \cdots & \frac{\partial^2 \ln f(\mathbf{z})}{\partial z_1 \partial z_d} \\ \frac{\partial^2 \ln f(\mathbf{z})}{\partial z_2 \partial z_1} & \frac{\partial^2 \ln f(\mathbf{z})}{\partial z_2^2} & \cdots & \frac{\partial^2 \ln f(\mathbf{z})}{\partial z_2 \partial z_d} \\ \vdots & \vdots & \vdots & \vdots \\ \frac{\partial^2 \ln f(\mathbf{z})}{\partial z_d \partial z_1} & \frac{\partial^2 \ln f(\mathbf{z})}{\partial z_d \partial z_2} & \cdots & \frac{\partial^2 \ln f(\mathbf{z})}{\partial z_d^2} \end{pmatrix}$$

# The Laplace approximation for higher dimensions

We obtain:

$$f(\mathbf{z}) \approx f(\mathbf{z}_0) e^{-\frac{1}{2}(\mathbf{z}-\mathbf{z}_0)^{\top} \mathbf{A}(\mathbf{z}-\mathbf{z}_0)}$$

and normalized distribution

$$\begin{aligned} q(\mathbf{z}) &= \frac{|\mathbf{A}|^{1/2}}{(2\pi)^{d/2}} e^{-\frac{1}{2}(\mathbf{z}-\mathbf{z}_0)^{\top} \mathbf{A}(\mathbf{z}-\mathbf{z}_0)} \\ &= \mathcal{N}(\mathbf{z}|\mathbf{z}_0, \mathbf{A}^{-1}) \end{aligned}$$

# Bayesian logistic regression

We want to compute the full posterior  $p(\mathbf{w}|\mathcal{D})$ .

**Problem:** Exact Bayesian inference for logistic regression **not possible** as exact evaluation of  $p(\mathbf{w}|\mathcal{D})$  **intractable!**

$$\begin{aligned} p(\mathbf{w}|\mathcal{D}) &\propto p(\mathbf{w})p(t_{\mathcal{D}}|x_{\mathcal{D}}, \mathbf{w}) \\ &= p(\mathbf{w}) \prod_{i=1}^N p(t_i|x_i, \mathbf{w}) \end{aligned}$$

**Idea:** Use Laplace approximation

Also need to select a prior for  $\mathbf{w}$   $\rightarrow$  a Gaussian prior,

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_0, \mathbf{S}_0)$$

# Bayesian logistic regression

Log-likelihood function:

$$\ln p(\mathbf{w}|\mathcal{D}) = -\frac{1}{2}(\mathbf{w} - \mathbf{m}_0)^\top \mathbf{S}_0^{-1}(\mathbf{w} - \mathbf{m}_0) + \sum_{i=1}^N \left( t_i \ln(y_i) + (1 - t_i) \ln(1 - y_i) \right) \\ + \text{const.}$$

We proceed as follows:

1. We maximize  $p(\mathbf{w}|\mathcal{D})$  (MAP solution  $\mathbf{w}_{\text{MAP}}$ )  $\rightarrow$  defines the mean of the Gaussian
2. Covariance given by

$$\mathbf{S}_N^{-1} = -\nabla \nabla \ln p(\mathbf{w}|\mathcal{D}) = \mathbf{S}_0^{-1} + \sum_{i=1}^N y_i(1 - y_i) \phi_i \phi_i^\top$$

Laplace approximation of  $p(\mathbf{w}|\mathcal{D})$ :

$$q(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{w}_{\text{MAP}}, \mathbf{S}_N)$$

# Bayesian logistic regression

Why Laplace approximation?

- Posteriors more “Gaussian like” as sample-size increases  $\rightarrow$  Good approximation for sufficiently large  $N$

# Predictive distribution

Want to **predict**  $t$ !  $\longrightarrow$  Want to learn **posterior predictive distribution**  
 $p(t|\mathcal{D}, x) = p(t|\mathcal{D}, \phi)$

Predictive distribution for class  $\mathcal{C}_1$ ,  $p(t = \mathcal{C}_1|\mathcal{D}, \phi)$ :

$$p(\mathcal{C}_1|\mathcal{D}, \phi) = \int p(\mathcal{C}_1|\phi, w)p(w|\mathcal{D})dw$$

Given  $\phi$  (or  $x$ ), we can forecast  $t$  via the **predictive distribution**.

# Predictive distribution

We use

$$p(\mathbf{t} = \mathcal{C}_1 | \phi, \mathbf{w}) = \frac{1}{1 + e^{-(\mathbf{w}^\top \phi(\mathbf{x}))}} \triangleq \sigma(\mathbf{w}^\top \phi(\mathbf{x}))$$

and

$$q(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{w}_{\text{MAP}}, \mathbf{S}_N)$$

We obtain:

$$\begin{aligned} p(\mathcal{C}_1 | \mathcal{D}, \phi) &\approx \int \sigma(\mathbf{w}^\top \phi(\mathbf{x})) q(\mathbf{w}) d\mathbf{w} \\ &= \int \sigma(a) \mathcal{N}(a | \mu_a, \sigma_a^2) da \end{aligned}$$

with

$$\begin{aligned} a &= \mathbf{w}^\top \phi \\ \mu_a &= \mathbb{E}[a] = \mathbf{w}_{\text{MAP}}^\top \phi \\ \sigma_a^2 &= \text{Var}[a] = \phi^\top \mathbf{S}_N \phi \end{aligned}$$

## Predictive distribution

$$p(\mathcal{C}_1|\mathcal{D}, \phi) \approx \int \sigma(a) \mathcal{N}(a|\mu_a, \sigma_a^2) da$$

Cannot be evaluated **analytically**

**Idea:** exploit that  $\sigma(a)$  is similar to the **probit function**  $\Phi(a)$  to obtain an approximation

$$\Phi(a) \triangleq \int_{-\infty}^a \mathcal{N}(x|0, 1) dx$$

$\sigma(a)$  **well approximated** by  $\Phi(\lambda a)$  with  $\lambda = \sqrt{\pi/8}$ .

$$\int \Phi(\lambda a) \mathcal{N}(a|\mu, \sigma^2) da = \Phi\left(\frac{\mu}{(\lambda^{-2} + \sigma^2)^{1/2}}\right)$$



## Predictive distribution

$$\int \Phi(\lambda a) \mathcal{N}(a|\mu, \sigma^2) da = \Phi\left(\frac{\mu}{(\lambda^{-2} + \sigma^2)^{1/2}}\right)$$

Then:

$$\begin{aligned} p(\mathcal{C}_1|\mathcal{D}, \phi) &\approx \int \sigma(a) \mathcal{N}(a|\mu_a, \sigma_a^2) da \\ &\approx \Phi\left(\frac{\mu_a}{(\lambda^{-2} + \sigma_a^2)^{1/2}}\right) \\ &\approx \sigma\left(\frac{\mu_a}{(1 + \pi\sigma_a^2/8)^{1/2}}\right) \end{aligned}$$

# Generative probabilistic models

Generative probabilistic models: Model  $p(\mathbf{x}, t)$  specifying  $p(t)$  ( $p(C_k)$ ) and  $p(\mathbf{x}|t)$  ( $p(\mathbf{x}|C_k)$ ). Then,

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

- Make more assumptions about the data by considering **distribution of covariates**  $\mathbf{x} \rightarrow$  may suffer from **bias** when model incorrectly selected
- Capability to capture properties of  $\mathbf{x}$  can improve learning if  $p(\mathbf{x}|t)$  has significant **structure**

# Generative probabilistic models

Binary classification:

$$\begin{aligned} p(\mathcal{C}_1|\mathbf{x}) &= \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1) + p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)} \\ &= \frac{1}{1 + e^{-a}} \triangleq \sigma(a) \end{aligned}$$

with

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

Multi-class classification ( $K > 2$ ):

$$p(\mathcal{C}_k|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)}{\sum_j p(\mathbf{x}|\mathcal{C}_j)p(\mathcal{C}_j)} = \frac{e^{a_k}}{\sum_j e^{a_j}} \quad (\text{softmax})$$

with

$$a_k = \ln p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)$$

# Generative probabilistic models

**Model:** Different models for  $p(\mathbf{x}|C_i)$  lead to different generative models.  
Generative models typically defined by assuming

$$\mathbf{x}|\mathbf{t} = t \sim \text{exponential}(\eta_t)$$

**Exponential family:**  $p(\mathbf{x}|\boldsymbol{\eta}) = \frac{1}{Z(\boldsymbol{\eta})} h(\mathbf{x}) \exp(\boldsymbol{\eta}^\top \mathbf{u}(\mathbf{x}))$

**Gaussian discriminant analysis:** Class-conditional density modeled by a multivariate Gaussian

- Linear discriminant analysis
- Quadratic discriminant analysis

# Linear discriminant analysis

**Principle:**  $\mathbf{x} \in \mathbb{R}^d$  is Gaussian distributed for each class  $C_k$ , with covariance matrix  $\Sigma_k = \Sigma$ ,

$$\mathbf{x}|\mathbf{t} = C_k \sim \mathcal{N}(\boldsymbol{\mu}_k, \Sigma)$$

Equivalently,

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

Fitting an LDA model to data  $\equiv$  estimating  $\Sigma$  and  $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K$ .

LDA gives rise to linear decision boundaries.

## Linear discriminant analysis

$$\begin{aligned}\log \frac{p(\mathbf{t} = C_i | \mathbf{x})}{p(\mathbf{t} = C_j | \mathbf{x})} &= \log \frac{p(\mathbf{x} | C_i) p(\mathbf{t} = C_i)}{p(\mathbf{x} | C_j) p(\mathbf{t} = C_j)} \\&= \log \frac{p(\mathbf{x} | C_i)}{p(\mathbf{x} | C_j)} + \log \frac{p(\mathbf{t} = C_i)}{p(\mathbf{t} = C_j)} \\&= -\frac{1}{2} \boldsymbol{\mu}_i^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_i + \frac{1}{2} \boldsymbol{\mu}_j^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_j + \mathbf{x}^\top \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_i - \boldsymbol{\mu}_j) + \log \frac{p(\mathbf{t} = C_i)}{p(\mathbf{t} = C_j)} \\&= w_0 + \mathbf{w}^\top \mathbf{x}\end{aligned}$$

with

$$\begin{aligned}\mathbf{w} &= \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_i - \boldsymbol{\mu}_j) \\w_0 &= -\frac{1}{2} \boldsymbol{\mu}_i^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_i + \frac{1}{2} \boldsymbol{\mu}_j^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_j + \log \frac{p(C_i)}{p(C_j)}\end{aligned}$$

**Decision boundaries:** points for which  $p(\mathbf{t} = C_i | \mathbf{x}) = p(\mathbf{t} = C_j | \mathbf{x}) \longrightarrow \{\mathbf{x} : w_0 + \mathbf{w}^\top \mathbf{x} = 0\}$

# Quadratic discriminant analysis

**Principle:**  $\mathbf{x} \in \mathbb{R}^D$  is Gaussian distributed for each class  $C_k$ , with covariance matrix  $\Sigma_k$ ,

$$\mathbf{x}|\mathbf{t} = C_k \sim \mathcal{N}(\boldsymbol{\mu}_k, \Sigma_k)$$

Equivalently,

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

**Discriminant functions:**

$$\log \frac{p(\mathbf{t} = C_i|\mathbf{x})}{p(\mathbf{t} = C_j|\mathbf{x})} = -\frac{1}{2} \log |\Sigma_i| - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^\top \Sigma_i^{-1}(\mathbf{x} - \boldsymbol{\mu}_i) + \log \frac{p(C_i)}{p(C_j)}$$

The **decision boundaries** are given by a **quadratic** equation.

# Gaussian discriminant analysis: Learning

**Need to fit model:** Determine values of  $\mathbf{w}$  and  $p(\mathcal{C}_k)$  via **maximum likelihood** criterion.

Binary case:

- $N$  training examples  $\mathcal{D} = \{(\mathbf{x}_i, t_i)\}$
- $t = 1$  for class  $C_1$  and  $t = 0$  for class  $C_2$
- $p(C_1) = \pi$ ,  $p(C_2) = 1 - \pi$

We have:

$$\begin{aligned}p(\mathbf{x}_i, C_1) &= p(C_1)p(\mathbf{x}_i|C_1) = \pi\mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}_1, \Sigma) \\p(\mathbf{x}_i, C_2) &= p(C_2)p(\mathbf{x}_i|C_2) = (1 - \pi)\mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}_2, \Sigma)\end{aligned}$$

**Likelihood function:**

$$p(\mathcal{D}|\pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \Sigma) = \prod_{i=1}^N (\pi\mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}_1, \Sigma))^{t_i} ((1 - \pi)\mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}_2, \Sigma))^{1-t_i}$$



# Gaussian discriminant analysis: Learning

## Likelihood function

$$p(\mathcal{D}|\pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}) = \prod_{i=1}^N (\pi \mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}_1, \boldsymbol{\Sigma}))^{t_i} ((1 - \pi) \mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}_2, \boldsymbol{\Sigma}))^{1-t_i}$$

## Log-likelihood function:

$$\begin{aligned} \ln p(\mathcal{D}|\pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}) &= \sum_{i=1}^N \left( t_i \ln (\pi \mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}_1, \boldsymbol{\Sigma})) + (1 - t_i) \ln ((1 - \pi) \mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}_2, \boldsymbol{\Sigma})) \right) \\ &= \sum_{i=1}^N \left( t_i \ln (\pi) + (1 - t_i) \ln(1 - \pi) \right) + \sum_{i=1}^N \left( t_i \ln (\mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}_1, \boldsymbol{\Sigma})) \right) \\ &\quad + \sum_{i=1}^N \left( (1 - t_i) \ln (\mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}_2, \boldsymbol{\Sigma})) \right) \end{aligned}$$

Can estimate  $\pi$ ,  $\boldsymbol{\mu}_i$ , and  $\boldsymbol{\Sigma}$  separately!

# Gaussian discriminant analysis: Learning

Maximization with respect to  $\pi$ :

$$\sum_{i=1}^N \left( t_i \ln(\pi) + (1 - t_i) \ln(1 - \pi) \right)$$

Differentiating and equating to zero:

$$\frac{d}{d\pi} \sum_{i=1}^N \left( t_i \ln(\pi) + (1 - t_i) \ln(1 - \pi) \right) = \sum_{i=1}^N t_i \frac{1}{\pi} - (1 - t_i) \frac{1}{1 - \pi}$$

$$\implies \pi = \frac{1}{N} \sum_{i=1}^N t_i = \frac{N_1}{N}$$

$N_1$ : number of data points in  $C_1$

# Gaussian discriminant analysis: Learning

Maximization with respect to  $\mu_1$ :

$$\sum_{i=1}^N \left( t_i \ln (\mathcal{N}(\mathbf{x}_i | \mu_1, \Sigma)) \right) = -\frac{1}{2} \sum_{i=1}^N t_i (\mathbf{x}_i - \mu_1)^\top \Sigma^{-1} (\mathbf{x}_i - \mu_1) + \text{const.}$$

Differentiating and equating to zero:

$$\mu_1 = \frac{1}{N_1} \sum_{i=1}^N t_i \mathbf{x}_i$$

Similarly,

$$\mu_2 = \frac{1}{N_2} \sum_{i=1}^N (1 - t_i) \mathbf{x}_i$$

# Gaussian discriminant analysis: Learning

Maximization with respect to  $\Sigma$ :

$$\Sigma = \frac{N_1}{N} \Sigma_1 + \frac{N_2}{N} \Sigma_2$$

with

$$\Sigma_1 = \frac{1}{N_1} \sum_{\substack{i=1; \\ t_i=C_1}}^N (\mathbf{x}_i - \boldsymbol{\mu}_1)(\mathbf{x}_i - \boldsymbol{\mu}_1)^\top$$

$$\Sigma_2 = \frac{1}{N_2} \sum_{\substack{i=1; \\ t_i=C_2}}^N (\mathbf{x}_i - \boldsymbol{\mu}_2)(\mathbf{x}_i - \boldsymbol{\mu}_2)^\top$$

# Gaussian discriminant analysis: Learning

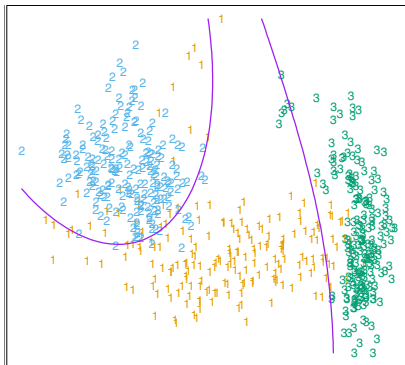
Multi-class classification ( $K > 2$ ):

$$\pi_k = \frac{N_k}{N}$$

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{\substack{i=1: \\ t_i=C_k}}^N t_i \mathbf{x}_i$$

$$\boldsymbol{\Sigma}_k = \frac{1}{N_k} \sum_{\substack{i=1: \\ t_i=C_k}}^N (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^\top$$

# Quadratic discriminant analysis



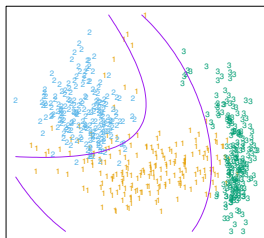
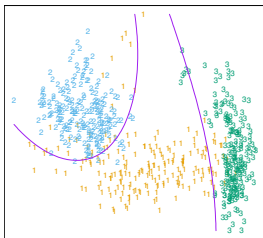
- Data  $\{x_i = (x_{i,1}, x_{i,2})\}$  generated from three different Gaussian mixtures

# Gaussian discriminant analysis

Two **alternative ways** to realize **quadratic** decision boundaries:

1. Apply QDA in the original  $d$ -dimensional space  $\mathbf{x} = (x_1, \dots, x_d)$
2. Apply LDA in an augmented space of dimension  $d^2 + 1$ :  
 $\mathbf{x}' = (x_1, \dots, x_d, x_1^2, \dots, x_d^2, x_1x_2, \dots, x_{d-1}x_d)$

Linear functions in the **augmented space** correspond to **quadratic** functions in the **original  $d$ -dimensional space**.



- Left: **QDA** applied to original 2-dimensional space  $\mathbf{x} = (x_1, x_2)$
- Right: **LDA** applied to 5-dimensional space  $\mathbf{x} = (x_1, x_2, x_1x_2, x_1^2, x_2^2)$

# Gaussian discriminant analysis

**Assumption** in LDA and QDA:  $p(\mathbf{x}|C_i)$  are **Gaussian** distributed  
... and yet perform well on a large variety of classification tasks!



# Gaussian discriminant analysis: Learning

## Observations:

- From  $p(\mathbf{x}|\mathcal{C}_k)$  and  $p(\mathcal{C}_k)$ , we use Bayes' theorem to find  $p(\mathcal{C}_k|\mathbf{x})$
- **Generative probabilistic model**: We can generate synthetic data by drawing values of  $\mathbf{x}$  from  $p(\mathbf{x})$
- **ML learning**: fast and simple...but it may **overfit** in high dimensions

# Discriminative or generative models?

## Easy to fit?

- **Generative**: Easy (counting and averaging)
- **Discriminative (logistic regression)**: Solving a complex optimization problem

## Adding more classes?

- **Generative**: No need to retrain (parameters of each class conditional density estimated independently)
- **Discriminative**: Need to retrain (all parameters interact)

## Feature processing?

- **Generative**: Hard
- **Discriminative**: Easy, we can simply replace  $x$  by  $\phi(x)$

# Reading

“Pattern recognition and machine learning,”

Chapter 4 (Intro, 4.1.1–4.1.3, 4.1.7, 4.2 (until 4.2.3), 4.3 (until 4.3.4), 4.4. (not 4.4.1), 4.5)