

# Statistical Learning – week 3.3

Joris Bierkens

Delft University of Technology, The Netherlands

29 February 2024

# Outline

## 1 Maximum Likelihood Estimation

## 2 Classification

- Linear regression

- Decision boundaries

- Quadratic and Linear Discriminant Analysis

- Logistic regression

## Recap learning objectives week 3.2

- Parametric models (G)
- The linear regression model and the least squares estimator
- Linear regression in the over-parametrized case; Moore-Penrose inverse
- The Gauss-Markov theorem
- **Refresher:** Lagrange Multipliers (G)
- Regularized linear regression: ridge regression and the LASSO

# Maximum Likelihood Estimation

- Consider a family of densities  $\{p_{\theta}(\mathbf{x}) : \theta \in \Theta\}$ .
- Each  $p_{\theta}(\mathbf{x})$  is a model for the observations  $\mathbf{x}$ .
- Here  $\mathbf{x}$  may (or may not) consist of multiple independent observations:  $\mathbf{x} = (x_1, \dots, x_n)$ .
- In **frequentist statistics** or **classical statistics** we typically assume there is a 'true' value  $\theta_0$ .
- We may use the **maximum likelihood method** to determine an estimator  $\hat{\theta}$  for  $\theta_0$ .

$$\hat{\theta} = \arg \max p_{\theta}(\mathbf{x}).$$

- Usually easier to work with the **log likelihood**  $\ell(\theta) = \log p_{\theta}(\mathbf{x})$ .
- **uncertainty quantification** using **confidence intervals**.

## Example: linear regression

Consider the linear regression model

$$y_i = \mathbf{x}_i^T \boldsymbol{\beta} + \varepsilon_i, \quad \varepsilon_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2) \quad i = 1, \dots, n.$$

Suppose we wish to estimate  $(\boldsymbol{\beta}, \sigma^2)$  from the data using maximum likelihood.

The likelihood is given by

$$L(\boldsymbol{\beta}) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-(y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 / (2\sigma^2)\right).$$

The log likelihood is

$$\ell(\boldsymbol{\beta}) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2.$$

### Exercise

The maximum likelihood estimator is given by

$$\hat{\boldsymbol{\beta}}_{\text{MLE}} = \hat{\boldsymbol{\beta}}_{\text{OLS}} \quad \text{and} \quad \hat{\sigma}_{\text{MLE}}^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{x}_i^T \hat{\boldsymbol{\beta}}_{\text{OLS}})^2 = \frac{1}{n} \text{RSS}(\hat{\boldsymbol{\beta}}_{\text{OLS}}).$$

## Maximum Likelihood Estimation: consistency

For observations  $\mathbf{X}_i$ ,  $i = 1, \dots, n$  and a model family

$$\{p_{\boldsymbol{\theta}}(\mathbf{x}) : \boldsymbol{\theta} \in \Theta\},$$

write the log likelihood as

$$\ell_n(\boldsymbol{\theta}) = \sum_{i=1}^n \log p_{\boldsymbol{\theta}}(\mathbf{X}_i).$$

### Theorem

Under mild conditions, the MLE is consistent.

This means that, if  $\mathbf{X}_i \sim p_{\boldsymbol{\theta}_0}(\mathbf{x})$ , for  $i = 1, \dots, n$ , and

$$\hat{\boldsymbol{\theta}}_n = \arg \max_{\boldsymbol{\theta} \in \Theta} \ell_n(\boldsymbol{\theta}),$$

then  $\hat{\boldsymbol{\theta}}_n \rightarrow \boldsymbol{\theta}_0$  in probability.

# Maximum Likelihood Estimation: asymptotic normality

For observations  $X_1, \dots, X_n$  and a family of distributions  $\{p_\theta : \theta \in \Theta\}$ , write  $\ell_n(\theta) = \sum_{i=1}^n \log p_\theta(X_i)$  for the log likelihood and  $\hat{\theta}_n = \arg \max_{\theta \in \Theta} \ell_n(\theta)$ .

## Theorem

Under mild conditions,

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \mathcal{I}_{\theta_0}^{-1}),$$

where the **Fisher information** is defined as

$$\mathcal{I}_\theta = \mathbb{E}_\theta \left[ (\nabla_\theta \log p_\theta(X)) (\nabla_\theta \log p_\theta(X))^T \right].$$

The Fisher information is often not available exactly:

- We do not know the value of  $\theta_0$ ;
- The required expectation may be difficult to compute.

How to deal with this?

- We may use a  $\hat{\theta}_n$  as a **plugin-estimator** for  $\theta_0$ .
- The Fisher information may be approximated using the **observed Fisher information**,

$$\mathcal{I}_{\theta_0} \approx \mathcal{I}_n = - \frac{1}{n} \sum_{i=1}^n \left( \nabla_\theta^2 \log p_{\hat{\theta}}(X_i) \right) \Big|_{\hat{\theta}_n}.$$

## Approximate confidence intervals

We have seen that the MLE  $\hat{\boldsymbol{\theta}}_n$  (under some conditions) admits the asymptotic distribution

$$\sqrt{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \mathcal{I}_{\boldsymbol{\theta}_0}^{-1}).$$

So if  $U_\alpha \subset \mathbb{R}^p$  is such that

$$\mathbb{P}(\mathcal{N}(\mathbf{0}, \mathbf{I}_p) \in U_\alpha) = 1 - \alpha,$$

then for

$$\mathcal{C}_\alpha := \{\boldsymbol{\theta} \in \mathbb{R}^p : \sqrt{n}\mathcal{I}_{\boldsymbol{\theta}_0}^{-1/2}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}) \in U_\alpha\},$$

we find that

$$\mathbb{P}(\boldsymbol{\theta}_0 \in \mathcal{C}_\alpha) \approx 1 - \alpha.$$



# Linear regression for classification

## Indicator response

- Classification: outcomes in a finite set  $\mathcal{Y}$ .
- Observations  $(\mathbf{x}_i, y_i)$  for  $\mathbf{x}_i \in \mathbb{R}^d$  and  $y_i \in \mathcal{Y}$ .
- Equivalently: observations  $(\mathbf{x}_i, \mathbf{z}_i)$  where  $\mathbf{z}_i \in \{0, 1\}^{\mathcal{Y}}$ .

$$\mathbf{z}_i(k) = \begin{cases} 1 & (y_i = k), \\ 0 & (y_i \neq k). \end{cases}$$

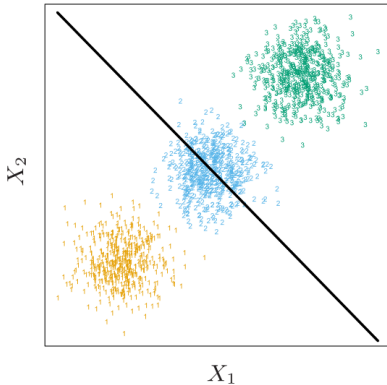
- $|\mathcal{Y}|$  regression problems: for  $k \in \mathcal{Y}$ ,

$$\mathbf{z}_i(k) = \mathbf{x}_i^T \boldsymbol{\beta}_k + \varepsilon_{i,k}.$$

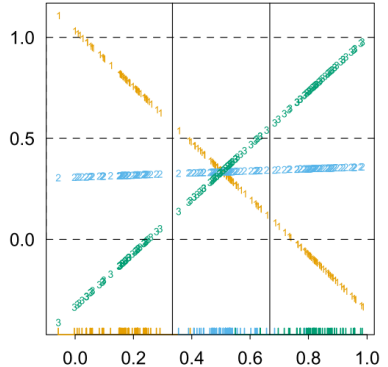
- Classify new input  $\mathbf{x}$  by finding

$$\hat{y}(\mathbf{x}) = \arg \max_{k \in \mathcal{Y}} \mathbf{x}^T \boldsymbol{\beta}_k.$$

# Masking



(a) A three-cluster model and the linear regression decision boundary



(b) Linear regression slopes along the diagonal

**Take away:** linear regression is not suitable for classification

# Decision boundaries in classification

Recall the **generative model for classification**: For  $\mathbf{x} \in \mathbb{R}^d$  and  $y \in \mathcal{Y}$ ,

- class-conditional densities  $p(\mathbf{x} \mid y)$ , and
- prior class probabilities  $\pi_k := p(y = k)$ .

The posterior class probabilities are given by Bayes rule

$$p(y \mid \mathbf{x}) = \frac{p(\mathbf{x}, y)}{p(\mathbf{x})} = \frac{p(\mathbf{x} \mid y)p(y)}{\sum_{y' \in \mathcal{Y}} p(\mathbf{x} \mid y')p(y')}.$$

## Discriminant function

**Discriminant functions**  $\delta_k : \mathbb{R}^d \rightarrow \mathbb{R}$ , for  $k \in \mathcal{Y}$ , are functions satisfying

$$p(y = k \mid \mathbf{x}) = \frac{\exp(\delta_k(\mathbf{x}))}{\sum_{k'} \exp(\delta_{k'}(\mathbf{x}))}.$$

Equivalently  $\delta_k(\mathbf{x}) = \log p(y = k \mid \mathbf{x}) + h(\mathbf{x})$  for some function  $h(\mathbf{x})$ .

The **decision boundaries** are hypersurfaces in  $\mathbb{R}^d$  given by

$$\{\mathbf{x} \in \mathbb{R}^d : \delta_k(\mathbf{x}) = \delta_\ell(\mathbf{x})\}, \quad k, \ell \in \mathcal{Y}.$$

# Softmax and logistic function

## The softmax function

The **softmax function**  $\sigma : \mathbb{R}^{\mathcal{Y}} \rightarrow [0, 1]^{\mathcal{Y}}$  is given by

$$\sigma(\mathbf{a})_k = \frac{\exp(a_k)}{\sum_{k'} \exp(a_{k'})}, \quad k \in \mathcal{Y}.$$

It is invariant under transformation  $a'_k = a_k + c$ .

Using the softmax function, our posterior class probabilities may be written as

$$p(y = k \mid \mathbf{x}) = \frac{\exp(\delta_k(\mathbf{x}))}{\sum_{k'} \exp(\delta_{k'}(\mathbf{x}))} = \sigma(\delta(\mathbf{x}))_k.$$

If  $|\mathcal{Y}| = 2$ , say  $\mathcal{Y} = \{0, 1\}$ , this simplifies to

$$p(y = 1 \mid \mathbf{x}) = \sigma(\delta(\mathbf{x})),$$

where

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

is the **logistic function**.

# Quadratic Discriminant Analysis

Suppose for  $k \in \mathcal{Y}$  the class-conditional probabilities are given by a  $\mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$  distribution, i.e.

$$p(\mathbf{x} \mid y = k) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}_k|^{1/2}} \exp \left( -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right).$$

This gives **quadratic discriminant functions**

$$\delta_k(\mathbf{x}) = -\frac{1}{2} \log |\boldsymbol{\Sigma}_k| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) + \log \pi_k.$$

# Quadratic discriminant analysis

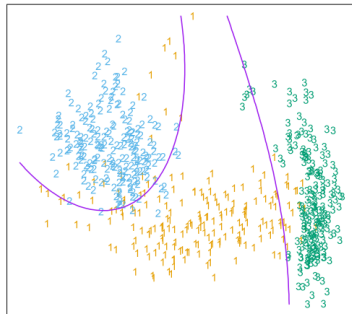
$$\delta_k(x) = -\frac{1}{2} \log |\Sigma_k| - \frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) + \log \pi_k$$

- $(\Sigma_k)_{k \in \mathcal{Y}}$ ,  $(\mu_k)_{k \in \mathcal{Y}}$  and  $(\pi_k)_{k \in \mathcal{Y}}$  are **parameters** of the model.
- Unbiased estimators

$$\hat{\pi}_k = n_k / n, \quad (n_k = \sum_{i=1}^n \mathbb{1}_{y_i=k}),$$

$$\hat{\mu}_k = \frac{1}{n_k} \sum_{i:y_i=k} \mathbf{x}_i,$$

$$\hat{\Sigma}_k = \frac{1}{n_k - 1} \sum_{i:y_i=k} (\mathbf{x}_i - \hat{\mu}_k)(\mathbf{x}_i - \hat{\mu}_k)^T.$$

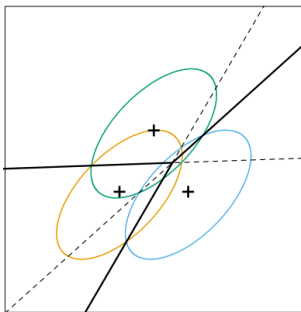


# Linear discriminant analysis

Linear discriminant analysis (LDA) arises by assuming that  $\Sigma_k = \Sigma$  does not depend on  $k$ .

In this case the decision boundaries are **hyperplanes** determined by the **linear discriminant functions**

$$\delta_k(x) = \mathbf{x}^T \Sigma^{-1} \boldsymbol{\mu}_k - \frac{1}{2} \boldsymbol{\mu}_k^T \Sigma^{-1} \boldsymbol{\mu}_k + \log \pi_k.$$

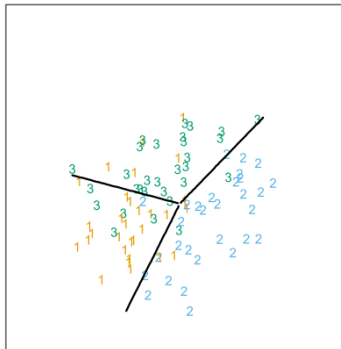


# Linear discriminant analysis (LDA)

$$\delta_k(\mathbf{x}) = \mathbf{x}^T \Sigma^{-1} \boldsymbol{\mu}_k - \frac{1}{2} \boldsymbol{\mu}_k^T \Sigma^{-1} \boldsymbol{\mu}_k + \log \pi_k.$$

- Again,  $(\pi_k)_{k \in \mathcal{Y}}$ ,  $(\boldsymbol{\mu}_k)_{k \in \mathcal{Y}}$  and  $\Sigma$  are parameters.
- May use the same estimators as before for  $\pi_k$  and  $\boldsymbol{\mu}_k$ .
- An unbiased estimator for  $\Sigma$  is given by

$$\hat{\Sigma} = \frac{1}{n - |\mathcal{Y}|} \sum_{k \in \mathcal{Y}} \sum_{i: y_i = k} (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_k)(\mathbf{x}_i - \hat{\boldsymbol{\mu}}_k)^T$$





## From LDA to logistic regression

In LDA the posterior class probabilities are fully determined by the discriminant functions of the form  $\delta_k(\mathbf{x}) = \beta_{k,0} + \beta_k^T \mathbf{x}$ , i.e.,

$$p(y = k \mid \mathbf{x}) = \sigma(\delta(\mathbf{x}))_k = \frac{\exp(\delta_k(\mathbf{x}))}{\sum_{k'} \exp(\delta_{k'}(\mathbf{x}))}.$$

We therefore parametrize using  $(\beta_{k,0}, \beta_k)_{k \in \mathcal{Y}}$  instead of  $(\mu_k, \pi_k)_{k \in \mathcal{Y}}$  and  $\Sigma$ .

This is still (slightly) overparametrized:  $\sigma(\delta(\mathbf{x}))$  is invariant under addition of  $\alpha_0 + \alpha^T \mathbf{x}$  to each function  $\delta_k(\mathbf{x})$ . Therefore we set  $\delta_K(\mathbf{x}) = 0$  for a single  $K \in \mathcal{Y}$ .

This gives the **multinomial logistic regression model**

$$p(y = k \mid \mathbf{x}) = \frac{\exp(\beta_{k,0} + \beta_k^T \mathbf{x})}{1 + \sum_{k \in \mathcal{Y}; k \neq K} \exp(\beta_{k,0} + \beta_k^T \mathbf{x})},$$
$$p(y = K \mid \mathbf{x}) = \frac{1}{1 + \sum_{k \in \mathcal{Y}; k \neq K} \exp(\beta_{k,0} + \beta_k^T \mathbf{x})}.$$


How to estimate  $(\beta_{k,0}, \beta_k)_{k \in \mathcal{Y}, k \neq K}$ ? Maximum likelihood!

## Binary logistic regression

In binary classification we often choose  $\mathcal{Y} = \{0, 1\}$ .

This gives the **binary logistic regression model**

$$p_{\beta}(y = 1 \mid \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{x}^T \boldsymbol{\beta})}.$$

The log likelihood is given by 

$$\begin{aligned}\ell(\boldsymbol{\beta}) &= \sum_{i=1}^n y_i \log p_{\beta}(y_i = 1 \mid \mathbf{x}_i) + (1 - y_i) \log p_{\beta}(y_i = 0 \mid \mathbf{x}_i) \\ &= \sum_{i=1}^n y_i \mathbf{x}_i^T \boldsymbol{\beta} - \log \left( 1 + e^{\mathbf{x}_i^T \boldsymbol{\beta}} \right)\end{aligned}$$

## Learning objectives week 3.3

- Discriminant functions in classification
- Quadratic Discriminant Analysis and Linear Discriminant Analysis
- **Refresher:** Maximum likelihood method (G)
- Logistic regression

Lecture note exercises: 2.1, 2.2, 2.3, 4.1, 4.2, 4.3 and 4.4.

First assignment is now available.