

# Learnig and Estimation of Dynamical Systems

Dante Piotto

Spring semester 2023



# Contents

<b>1</b>	<b>Introduction</b>	<b>5</b>
1.1	What is a system? . . . . .	5
1.1.1	Learning models from data . . . . .	6
1.2	Types of learning . . . . .	6
1.3	fields related to learning from data . . . . .	6
1.3.1	System identification . . . . .	6
1.4	Learning steps . . . . .	7
<b>2</b>	<b>Stochastic Processes</b>	<b>9</b>
2.1	stationary stochastic processes . . . . .	9
2.1.1	cross-correlation and cross-covariance . . . . .	9
2.2	vector stochastic processes . . . . .	9
2.3	gaussian processes . . . . .	9
2.4	white processes . . . . .	9
<b>3</b>	<b>Stochastic models</b>	<b>11</b>
3.1	Matlab stuff . . . . .	11
<b>4</b>	<b>Estimation problem</b>	<b>13</b>
<b>5</b>	<b>Linear regression</b>	<b>15</b>
5.1	The Least Squares Method . . . . .	15
5.1.1	Geometrical interpretation of the LS estimate . . . . .	16
5.1.2	Statistical properties of the LS estimator . . . . .	16
5.1.3	LS estimation for ARX models . . . . .	17
5.1.4	ARX optimal (one step ahead) predictor . . . . .	18
5.1.5	LS estimation of AR models . . . . .	18
5.2	Recursive least squares . . . . .	18
5.2.1	Asymptotic behaviour of the RLS algorithm . . . . .	21
5.2.2	Recursive weighted least squares . . . . .	21
<b>6</b>	<b>Prediction error methods</b>	<b>23</b>
6.0.1	Identification of ARMAX models . . . . .	25
6.0.2	Statistical properties of PEM estimators . . . . .	25
6.0.3	MISO ARX models . . . . .	26
<b>7</b>	<b>Statistical hypothesis testing</b>	<b>27</b>
<b>8</b>	<b>Model complexity selection and regularization</b>	<b>29</b>
8.0.1	The F-test . . . . .	29
8.0.2	The final prediction error (FPE) criterion . . . . .	30
8.0.3	Criteria with complexity terms . . . . .	31
8.0.4	Akaike information criterion (AIC) . . . . .	31

<b>9</b>	<b>model assesment (validation)</b>	<b>33</b>
9.1	Whiteness test . . . . .	33
9.2	Test of cross-correlation . . . . .	34
<b>10</b>	<b>Maximum likelihood estimation</b>	<b>35</b>
10.1	the gaussian case . . . . .	35
10.1.1	the Cramér-Rao lower bound . . . . .	36
<b>11</b>	<b>Classification: probabilistic models</b>	<b>37</b>
11.1	The Bayes classifier . . . . .	37
11.2	Logistic regression . . . . .	38
11.2.1	The gradient descent algorithm . . . . .	39
11.2.2	Multiclass problems . . . . .	40
11.2.3	Dealing with nonlinear boundaries . . . . .	41
11.2.4	Linear discriminant analysis . . . . .	41
11.2.5	Gaussian class densities with common covariance matrix . . . . .	42
11.2.6	Gaussian class denisities with different covariance matrices . . . . .	42
<b>12</b>	<b>Classification: deterministic models</b>	<b>43</b>
12.0.1	Separating hyperplanes . . . . .	44
12.0.2	the maximum margin classifier . . . . .	44
12.1	Support vector machine . . . . .	47
12.1.1	Dealing with nonlinear boundaries: the kernel trick . . . . .	49
12.1.2	Multiclass problems . . . . .	49
<b>13</b>	<b>Regularization</b>	<b>51</b>
13.1	Regularized least squares: the Ridge regression . . . . .	51

# Chapter 1

## Introduction

Goal of the course: build mathematical models of dynamical systems from data

Password of slides: LED\$2023

Office appointments to be decided with professor as needed

### 1.1 What is a system?

A slice of reality whose evolution in time can be described by a certain number of measurable attributes (variables)

Inputs: independent variables (causes) which describe the action the surrounding environment on the system

Outputs: dependent variables (effects) which describe the reaction of the system

Mathematical model: a set of quantitative relationships between the system variables

Solving problems in scientific disciplines by means of mathematical models:

1. determination of a mathematical model of the system
2. solution of the problem by using the model (i.e. in the mathematical world)
3. implementation of the obtained solution on the real process

Competent model: a good model for solving a given problem in a certain problem context

- different mathematical models can be associated with the same system
- classification of models based on the modeling objectives

modeling objectives:

- inference
- control
- prediction
- filtering
- diagnosis
- predictive maintenance
- simulation
- speech and image recognition

### 1.1.1 Learning models from data

Data (set of samples)

$$u(1), u(2), \dots, u(N) \quad y(1), y(2), \dots, y(N) \quad u(t) \in \mathcal{U}, y(t) \in \mathcal{Y}$$

Target model(function)

$$\mathcal{M}_p(\theta)$$

$\mathcal{M}_p(\cdot)$  represents a function linking input and output samples,  $\theta$  is a set of parameters and  $p$  is a set of hyper-parameters

Static models:  $f : \mathcal{U} \rightarrow \mathcal{Y}$

Dynamic models:  $f : (\mathcal{U}, \mathcal{Y}) \rightarrow \mathcal{Y}$

## 1.2 Types of learning

Supervised learning ( $u, y$  known):

- $\mathcal{Y}$  is discrete: classification
- $\mathcal{Y}$  is continuous: regression

Unsupervised learning ( $u, y$  unknown):

- $\mathcal{Y}$  is discrete: clustering
- $\mathcal{Y}$  is continuous: dimensionality reduction

Classification: assign the input to one of a finite number of classes

Regression: find an input-output relation

Reinforcement learning: finding suitable actions to take in a given situation in order to maximize a reward.

Both  $u, y$  are known but we have to find the "optimal" output for a given input

This course deals with Supervised learning

## 1.3 fields related to learning from data

- Machine learning
- Pattern recognition
- Statistical learning
- Data mining
- System identification

### 1.3.1 System identification

System identification is the art and science of building mathematical models of dynamic systems from observed input-output data

Learning from data is a *data-driven(black box) approach*: a model is selected within a specified model class by using a selection criterion on the only basis of experimental (observed) data. No reference to the physical nature of the system is made

- the obtained models have limited validity
- the model parameters may lack any physical meaning
- models relatively easy to construct and use
- ability to extract only some relevant aspects from complex frameworks

In contrast, *physical modeling* is a white box approach. The system is partitioned into subsystems that are described by using known laws of physics. Then, the model of the system is obtained by joining such relations.

**Grey box approach**

It often happens that a model based on physical modeling contains a number of unknown parameters: identification(learning) methods can be applied to estimate the unknown parameters.

**1.4 Learning steps**

The planned use of the model is important in designing the experiment to collect data (when possible).





## Chapter 2

# Stochastic Processes

Let us consider a random experiment, with sample space  $\Omega$ , and let us associate to each event  $\omega_i$  in the sample space a signal  $x(t, \omega_i)$ . With a fixed  $\omega$  we have a function of time  $x(t)$ , and at each fixed  $t_i$ ,  $x(\omega)$  is a random variable

Definition (discrete time stochastic process):

A function  $x(t, \omega)$  where  $t \in \{\dots, -2, -1, 0, 1, 2, 1 \dots\}$  is time and  $\omega \in \Omega$  is an outcome of the sample space.  $x(t, \omega_i)$  is called a *realization* of the stochastic process Given  $t = t_1$  the first order cdf and pdf are:

$$F(x; t_1) = P(x(t_1) \leq x) \quad f(x; t_1) = \frac{\delta F(x; t_1)}{\delta x}$$

Autocovariance: relation proven using linearity of the expectation operator

### 2.1 stationary stochastic processes

A stochastic process is stationary if

$$F(x_1, x_2, \dots, x_k; t_1, t_2, \dots, t_k) = F(x_1, x_2, \dots, x_k; t_1 + \tau, t_2 + \tau, \dots, t_k + \tau) \quad (2.1)$$

$$\forall \tau, \forall k, \forall \{t_1, t_2, \dots, t_k\} \quad (2.2)$$

this property also holds for the pdf Consequences:

$$\mu_x(t) = \mu_x \quad (2.3)$$

$$\sigma_x^2(t) = \sigma_x^2 \quad (2.4)$$

$$r_x(t_1, t_2) = r_x(t_1 - t_2) = r_x(\tau) \quad (2.5)$$

$$c_x(t_1, t_2) = c_x(t_1 - t_2) = c_x(\tau) \quad (2.6)$$

A process is weakly stationary (or wide-sense stationary) if the 4 properties above hold

Toeplitz matrix: symmetric matrix with all elements belonging to a diagonal being equal.

Cross-correlation and cross-covariance can only be defined for stationary stochastic processes

#### 2.1.1 cross-correlation and cross-covariance

### 2.2 vector stochastic processes

### 2.3 gaussian processes

### 2.4 white processes

it is not possible to define a continuous time white process



## Chapter 3

# Stochastic models

PSD= Power Spectral Density

Moving Average noise is also called pink noise Matlab considers the normalized autocorrelation.

### 3.1 Matlab stuff

`rand()`: generates white noise vector

`cov(X)`: gives the variance of vector X

`randn()`: generates gaussian distributed white noise vector

`autocorr(X)`: gives the normalized autocorrelation

to get the non normalized autocorrelation: `autocorr(X)*cov(X)`

`filter()`: useful to generate AR, MA, and ARMA processes



## Chapter 4

# Estimation problem

w.p.1: with probability 1 A biased estimator with small variance and a sufficiently small bias may be preferable to an unbiased estimator with high variance. *bias-variance tradeoff*



## Chapter 5

# Linear regression

Let us consider the static model

$$y(t) = f(u(t)) + e(t) \quad \text{model class } \mathcal{M}_p(\theta)$$

If the function  $f$  is linear in the parameters (elements of  $\theta$ ), the model can be written in the *linear regression* form:

$$y(t) = \varphi^T(t)\theta + e(t)$$

### 5.1 The Least Squares Method

available data set:

$$y(1), y(2), \dots, y(N), u(1), u(2), \dots, u(N)$$

If an estimate  $\hat{\theta}$  were available, we would compute the misfit between  $y(t)$  and its 'prediction'  $\hat{y}(t)$ :

$$\varepsilon(t) = y(t) - \hat{y}(t) = y(t) - \varphi^T(t)\hat{\theta}$$

where the error term  $\varepsilon(t)$  is the residual. The LS method finds the estimate  $\hat{\theta}$  that minimizes the loss function

$$J(\theta) = \sum_{t=1}^N \varepsilon^2(t) = \sum_{t=1}^N (y(t) - \varphi^T(t)\hat{\theta})^2 \quad (5.1)$$

In matrix form:

$$\varepsilon = Y - \Phi\theta$$

where

$$\varepsilon = \begin{bmatrix} \varepsilon(1) \\ \varepsilon(2) \\ \vdots \\ \varepsilon(N) \end{bmatrix} \quad y = \begin{bmatrix} y(1) \\ y(2) \\ \vdots \\ y(N) \end{bmatrix} \quad \varphi = \begin{bmatrix} \varphi(1) \\ \varphi(2) \\ \vdots \\ \varphi(\text{epsilon}N) \end{bmatrix}$$

so that

$$J(\theta) = \sum_{t=1}^N \varepsilon^2(t) = \|\varepsilon\|^2 = \|Y - \Phi\theta\|^2 \quad (5.2)$$

The optimization problem to be solved is

$$\min_{\theta \in \mathcal{M}_p(\theta)} J(\theta)$$

The solution can be found by using the following relations:

$$\frac{\partial A^T x}{\partial x} = A, \quad \frac{\partial x^T A}{\partial x} = A, \quad \frac{\partial x^T A x}{\partial x} = (A + A^T)x$$

where  $A$  is an  $n \times n$  matrix and  $x$  is a  $n \times 1$  vector

From the above relations we obtain the *normal equations*

$$\Phi^T \Phi \theta = \Phi^T Y$$

*proof was covered in class*

$$\begin{aligned} J(\theta) &= Y^T Y - 2Y^T \Phi \theta + \theta^T \Phi^T \Phi \theta \\ \frac{\partial J(\theta)}{\partial \theta} &= 0 - 2\Phi^T Y + (\Phi^T \Phi + \Phi^T \Phi) \theta \\ \frac{\partial J(\theta)}{\partial \theta} &= -2\Phi^T Y + 2\Phi^T \Phi \theta \\ \frac{\partial J(\theta)}{\partial \theta} &= 0 \implies \Phi^T \Phi \theta = \Phi^T Y \end{aligned}$$

To complete the proof we must prove that the second derivative of the matrix is positive definite

$$\frac{\partial^2 J(\theta)}{\partial \theta^2} = 2\Phi^T \Phi$$

If the  $N \times p$  matrix  $\Phi$  is tall ( $N > p$ ) and full rank, the LS estimate is given by

$$\hat{\theta} = (\Phi^T \Phi)^{-1} \Phi^T Y \quad (5.3)$$

### 5.1.1 Geometrical interpretation of the LS estimate

The estimated function is such that the sum of the squares of the distances, evaluated along the  $y$ -axis, between  $y(t)$  and its 'prediction'  $\hat{y}(t)$  is minimized.

Consider the linear map described by  $\Phi$ :

$$\Phi : \mathbb{R}^p \rightarrow \mathbb{R}^N$$

Let  $\Phi_1, \Phi_2, \dots, \Phi_p$  be the columns of matrix  $\Phi$ . The LS problem consists in finding the linear combination of  $\Phi_1, \Phi_2, \dots, \Phi_p$  that approximates  $Y$  as closely as possible. Therefore, the solution is given by the orthogonal projection of  $Y$  onto the subspace spanned by  $\Phi_1, \Phi_2, \dots, \Phi_p$ , i.e. the orthogonal projection of  $Y$  onto the image of  $A$

$\hat{\theta}$  is such that  $\varepsilon = Y - \Phi \hat{\theta}$  is orthogonal to  $\Phi_1, \Phi_2, \dots, \Phi_p$

$$\implies \varepsilon^T \Phi = 0$$

$$\begin{aligned} \varepsilon &= Y - \Phi \hat{\theta} = Y - \hat{Y} \\ \varepsilon^T \Phi &= (Y - \Phi \hat{\theta})^T \Phi = Y^T \Phi - \hat{\theta}^T \Phi^T \Phi \\ &= Y^T \Phi - Y^T \Phi (\Phi^T \Phi)^{-1} \Phi^T \Phi = 0 \end{aligned}$$

The LS solution  $\hat{\theta}$  can be obtained by considering the pseudoinverse of  $\Phi$ :

$$\hat{\theta} = \Phi^\dagger Y$$

If  $\Phi$  is full rank, its pseudoinverse is just given by  $\Phi^\dagger = (\Phi^T \Phi)^{-1} \Phi^T$

### 5.1.2 Statistical properties of the LS estimator

Assume the true model

$$y(t) = \varphi^T(t) \theta^* + w(t)$$



where  $w(t)$  is a zero mean white process with variance  $\sigma_w^2$ . Let  $\hat{\theta}_N$  be an estimate obtained by using  $N$  input-output samples

It follows that

$$\begin{aligned} E[\hat{\theta}] &= E[(\Phi^T \Phi)^{-1} \Phi^T Y] \\ Y &= \Phi \theta^* + w \\ E[\hat{\theta}]E &= [(\Phi^T \Phi)^{-1} \Phi^T Y = \Phi \theta^* + w] = E[\theta^* + (\Phi^T \Phi)^{-1} \Phi^T w] \\ &= \theta^* + (\Phi^T \Phi)^{-1} \Phi^T E[w] = \theta^* \end{aligned}$$

therefore the LS estimator is unbiased.

### 5.1.3 LS estimation fo ARX models

The closed form solution of the LS problem is obtained in a similar fashion to the FIR model case.

#### Identifiability

- $u(t)$  has to be persistently exciting of order  $\geq n$
- In order to have a well-condition matrix  $H$  (i.e.) with a low condition number), it is also necessary that  $n$  is not greater than the minimal order of an ARX model compatible with the data

#### Statistical properties

True model:

$$y(t) = \varphi^T(t) \theta^* + w(t)$$

It follows that

$$\begin{aligned} \hat{\theta} &= \left( \frac{1}{N} \sum_{t=1}^N \varphi(t) \varphi^T(t) \right)^{-1} \frac{1}{N} \sum_{t=1}^N \varphi(t) (\varphi^T(t) \theta^* + w(t)) \\ \hat{\theta} &= \theta^* + \left( \frac{1}{N} \sum_{t=1}^N \varphi(t) \varphi^T(t) \right)^{-1} \frac{1}{N} \sum_{t=1}^N \varphi(t) w(t) \end{aligned}$$

The estimate is biased:  $E[\hat{\theta}] \neq \theta^*$  Consistency:

$$\begin{aligned} \lim_{N \rightarrow \infty} \hat{\theta}_N &= \theta^* + \lim_{N \rightarrow \infty} \left( \frac{1}{N} \sum \varphi(t) \varphi^T(t) \right)^{-1} \frac{1}{N} \sum \varphi(t) w(t) \\ \lim_{N \rightarrow \infty} \hat{\theta}_N &= \theta^* + E[\varphi(t) \varphi^T(t)]^{-1} E[\varphi(t) w(t)] = \theta^* - \Sigma_\varphi^{-1} r_{\varphi w} \end{aligned}$$

where  $\Sigma_\varphi^{-1} = E(\varphi(t) \varphi^T(t))$  and  $r_{\varphi w}$  is the cross correlation. Because  $u(t)$  is pe of order  $n$ , and  $w(t)$  is white,  $\Sigma_\varphi$  is invertible, and  $r_{\varphi w} = 0$  because  $y(t-n)$  does not depend on  $w(t)$  for  $n > 0$  and  $u(t)$  in general does not depend on  $w$ , therefore

$$\lim_{N \rightarrow \infty} \hat{\theta} = \theta^*$$

It is also possible to prove that

$$\sqrt{N}(\hat{\theta} - \theta^*) \rightarrow \mathcal{N}(0, P), \text{ for } N \rightarrow \infty$$

with

$$P = \sigma_w^2 \Sigma_\varphi^{-1}$$

A consistent estimate of  $\sigma_w^2$  is given by

$$\sigma_w^2 = J(\hat{\theta})$$

Then,  $cov(\hat{\theta})$  can be estimated as

$$\frac{\hat{\sigma}_w^2 \hat{\Sigma}_\varphi^{-1}}{N} = \hat{\sigma}_w^2 (H^T H)^{-1}$$

The equivalence can be derived considering

$$\hat{\Sigma}_\varphi = \frac{1}{N} \sum \varphi(t) \varphi^T(t) = \frac{H^T H}{N}$$

### 5.1.4 ARX optimal (one step ahead) predictor

True model:

$$y(t) = \varphi^T(t)\theta^* + w(t)$$

Problem: find the optimal (minimal variance) prediction of  $y(t)$  given the past data  $y(t-1), u(t-1), y(t-2), u(t-2), \dots$ . It is easy to show that the optimal predictor  $\hat{y}(t|t-1)$  is given by

$$y(t|t-1) = \varphi^T(t)\theta^*$$

we know that

$$y(t) = -a_1^*y(t-1) - \dots - a_n^*y(t-n) + b_1^*u(t-1) + \dots + b_n^*u(t-n) + w(t)$$

because  $w(t)$  is a white process, the best estimate we can make for it is its mean, 0.

For any other predictor  $\hat{y}(t)$

$$E[(y(t) - \hat{y}(t))^2] \geq E[(y(t) - \hat{y}(t|t-1))^2] = \sigma_w^2$$

As a consequence, the LS estimation  $\hat{\theta}$  leads to a predictive model and the residual  $\varepsilon(t)$  can be seen as a prediction error. In fact, from

$$\lim_{N \rightarrow \infty} \hat{\theta} = \theta^*$$

it follows that

$$\varepsilon(t) = y(t) - \varphi^T(t)\hat{\theta} \xrightarrow{N \rightarrow \infty} \theta^*$$

### 5.1.5 LS estimation of AR models

$$y(t) + a_1y(t-1) + \dots + a_ny(t-n) = e(t)$$

Linear regression form:

$$y(t) = \varphi^T(t)\theta + e(t)$$

where:

$$\begin{aligned} \varphi(t) &= [-y(t-1) \quad -y(t-2) \quad \dots \quad -y(t-n)]^T \\ \theta &= [a_1 \quad \dots \quad a_n]^T \\ Y &= -H_y(n)\theta + \varepsilon \end{aligned}$$

LS estimate:

$$\hat{\theta} = -(H_y^T(n)H_y(n))^{-1}H_y^T(n)Y$$

Statistical properties and optimale predictor: similar considerations to those of the ARX case.

## 5.2 Recursive least squares

Let  $\hat{\theta}(t-1)$  be a LS estimate obtained from data collected up to time  $t-1$ :

$$\hat{\theta}(t-1) = \left( \sum_{k=1}^{t-1} \varphi(k)\varphi^T(k) \right)^{-1} \sum_{k=1}^{t-1} \varphi(k)y(k) \quad (1)$$

Recursive identification methods consist of updating  $\hat{\theta}(t-1)$  by some "simple modification" once data at time  $t$  becomes available to compute  $\hat{\theta}(t)$

Starting from (1) and

$$\begin{aligned}
 S(t) &= HH^T = \sum_{k=1}^t \varphi(k)\varphi^T(k) = S(t-1) + \varphi(t)\varphi^T(t) \\
 \varphi(k)y(k) &= \sum_{k=1}^t \varphi(k)y(k) + \varphi(t)y(t) \\
 \hat{\theta}(t) &= S(t)^{-1} \sum_{k=1}^t \varphi(k)y(k)
 \end{aligned}$$

It is possible to obtain

$$\hat{\theta}(t) = \hat{\theta}(t-1) + K(t)\varepsilon(t)$$

where

$$K(t) = S(t)^{-1}\varphi(t)$$

and we can recall that

$$\varepsilon(t) = y(t) - \varphi^T(t)\hat{\theta}(t-1)$$

is the prediction error. computation was developed in class This leads to the following recursive least squares algorithm:

#### RLS I

1.  $S(t) = S(t-1) + \varphi(t)\varphi^T(t)$
2.  $K(t) = S(t)^{-1}\varphi(t)$
3.  $\varepsilon(t) = y(t) - \varphi^T(t)\hat{\theta}(t-1)$
4.  $\hat{\theta}(t) = \hat{\theta}(t-1) + K(t)\varepsilon(t)$

By defining  $R(t) = \frac{S(t)}{t}$  it is possible to derive a RLS algorithm for

$$\hat{\theta}(t) = \left( \frac{1}{t} \sum_{k=1}^t \varphi(k)\varphi^T(k) \right)^{-1} \frac{1}{t} \sum_{k=1}^t \varphi(k)y(k)$$

in fact,  $R(t)$  can be easily updated

$$R(t) = \frac{S(t)}{t} = \frac{S(t-1) + \varphi(t)\varphi^T(t)}{t} = \frac{t-1}{t}R(t-1) + \frac{\varphi(t)\varphi^T(t)}{t}$$

Following the same steps used to derive RLS I we get (computation was developed in class):

#### RLS II

1.  $R(t) = \frac{t-1}{t}R(t-1) + \frac{1}{t}\varphi(t)\varphi^T(t)$
2.  $K(t) = \frac{1}{t}R(t)^{-1}\varphi(t)$
3.  $\varepsilon(t) = y(t) - \varphi^T(t)\hat{\theta}(t-1)$
4.  $\hat{\theta}(t) = \hat{\theta}(t-1) + K(t)\varepsilon(t)$

In order to not compute a matrix inversion at each step and instead updating the inverse itself, we can rely on the following result

**Matrix inversion lemma (Woodbury identity)**

Let  $A, C$  be square and invertible. Then

$$(A + BCD)^{-1} = A^{-1} - A^{-1}B(C^{-1} + DA^{-1}B)^{-1}DA^{-1}$$

$$S(t)^{-1} = (S(t-1) + \varphi(t)\varphi^T(t))^{-1}$$

$$A = S(t-1)$$

$$B = \varphi(t)$$

$$C = 1$$

$$D = \varphi^T(t)$$

$$S(t)^{-1} = S(t-1)^{-1} - S(t-1)^{-1}\varphi(t)(1 + \varphi^T(t)S(t-1)^{-1}\varphi(t))^{-1}\varphi^T(t)S(t-1)^{-1}$$

note that  $(1 + \varphi^T(t)S(t-1)^{-1}\varphi(t))$  is a scalar, so we can write:

$$S(t)^{-1} = S(t-1)^{-1} - \frac{S(t-1)^{-1}\varphi(t)\varphi^T(t)S(t-1)^{-1}}{1 + \varphi^T(t)S(t-1)^{-1}\varphi(t)}$$

We can now modify RLS I and RLS II in order to avoid matrix inversion

**RLS III**

1.  $S(t) = S(t-1)^{-1} - \frac{S(t-1)^{-1}\varphi(t)\varphi^T(t)S(t-1)^{-1}}{1 + \varphi^T(t)S(t-1)^{-1}\varphi(t)}$
2.  $K(t) = S(t)^{-1}\varphi(t)$
3.  $\varepsilon(t) = y(t) - \varphi^T(t)\hat{\theta}(t-1)$
4.  $\hat{\theta}(t) = \hat{\theta}(t-1) + K(t)\varepsilon(t)$

**RLS IV**

1.  $P(t) = \frac{t}{t-1}P(t-1) - \frac{t}{t-1} \frac{P(t-1)^{-1}\varphi(t)\varphi^T(t)P(t-1)^{-1}}{t-1 + \varphi^T(t)P(t-1)^{-1}\varphi(t)}$
2.  $K(t) = P(t)\frac{1}{t}\varphi(t)$
3.  $\varepsilon(t) = y(t) - \varphi^T(t)\hat{\theta}(t-1)$
4.  $\hat{\theta}(t) = \hat{\theta}(t-1) + K(t)\varepsilon(t)$

with  $P(t) = R(t)^{-1}$

**Initialization**

The RLS algorithm needs to be initialized. One possibility is to start with an initial batch estimate, which will be better the more data is available. It is also possible to start with some other type of guess on the model parameters, potentially also very bad, e.g. a vector of zeros. However, a guess for the gain matrix (either  $S(0), R(0), P(0)$  depending on the specific algorithm) is necessary. With the RLS IV, a good initialization would be

$$P(0) = \alpha I_p$$

If we are confident about the initial guess of the parameters, a small value of  $\alpha$  is appropriate. On the flipside, if we expect the initial estimator to be bad, the value of  $\alpha$  shall be quite large.

### 5.2.1 Asymptotic behaviour of the RLS algorithm

$$y(t) = \varphi^T(t)\theta^* + w(t)$$

$$\lim_{t \rightarrow \infty} \hat{\theta}(t) = \theta^* \quad \text{w.p. 1}$$

$R(t)$  in the RLS II algo is an estimate of the covariance matrix. Looking at the RLS II we can observe that from step 2,  $K(t)$  approaches 0. Therefore, asymptotically, there is no correction.

$$\lim_{t \rightarrow \infty} K(t) = \frac{\Sigma_{\varphi}^{-1}}{\infty} \varphi(t) = 0$$

### 5.2.2 Recursive wighted least squares

A modification of the RLS algorithm aimed at tracking parameter variations by giving less importance to past data and more importance to recenet data.

$$J(\theta) = \sum_{t=1}^N \lambda^{N-t} \varepsilon^2(t) = \varepsilon^T W \varepsilon$$

where

$$W = \text{diag} [\lambda^{N-1} \quad \lambda^{N-2} \quad \dots \quad \lambda \quad 1]$$

and  $\lambda, 0 < \lambda < 1$  is the *forgetting factor*. We have

$$\hat{\theta} = \left( \sum_{t=1}^N \lambda^{N-t} \varphi(t) \varphi^T(t) \right)^{-1} \sum_{t=1}^N \lambda^{N-t} \varphi(t) y(t)$$

The scalar  $\lambda$  should be chosen by a trade-off between the ability to track parameter changes on one hand, and good estimation accuracy on the other hand.

Recursive weighted LS: determine a recursive form of

$$\hat{\theta}(t) = \left( \sum_{k=1}^t \lambda^{t-k} \varphi(k) \varphi^T(k) \right)^{-1} \sum_{k=1}^t \lambda^{t-k} \varphi(k) y(k)$$

Define

$$S(t) = \sum_{k=1}^t \lambda^{t-k} \varphi(k) \varphi^T(k) = \lambda S(t-1) + \varphi(t) \varphi^T(t)$$

Following the same reasoning used to derive the previous recursive algorithms we can derive the RWLS algorithms

#### RWLS I

1.  $S(t) = \lambda S(t-1) + \varphi(t) \varphi^T(t)$
2.  $K(t) = S(t)^{-1} \varphi(t)$
3.  $\varepsilon(t) = y(t) - \varphi^T(t) \hat{\theta}(t-1)$
4.  $\hat{\theta}(t) = \hat{\theta}(t-1) + K(t) \varepsilon(t)$



## Chapter 6

# Prediction error methods

Let us consider the ARMAX model

$$A(z^{-1})y(t) = B(z^{-1})u(t) + C(z^{-1})w(t)$$

or

$$y(t) + a_1y(t-1) + \dots + a_ny(t-n) = b_1u(t-1) + \dots + b_nu(t-n) + w(t) + c_1w(t-1) + \dots + c_nw(t-n)$$

Then

$$y(t) = \varphi^T(t)\theta + w(t)$$

where

$$\varphi(t) = [-y(t-1) \quad \dots \quad -y(t-n) \quad u(t-1) \quad \dots \quad u(t-n) \quad w(t-1) \quad \dots \quad w(t-n)]^T$$

In a real data setting we may assume the model

$$A(z^{-1})y(t) = B(z^{-1})u(t) + C(z^{-1})\varepsilon(t)$$

where  $\varepsilon(t)$  is the residual. Problem: the residual can be computed once the estimate  $\hat{\theta}$  is available, so that it is a function of  $\theta, \varepsilon(t, \theta)$

$$y(t) = \varphi^T(t, \theta)\theta + \varepsilon(t)$$

This is no longer a linear regression and the least squares method cannot be applied. If we consider

$$\frac{C(z^{-1})}{A(z^{-1})}$$

as a disturbance  $d(t)$  and only try to estimate the plant, we obtain

$$y(t) = \bar{\varphi}^T(t)\bar{\theta} + e(t)$$

and we can use the least squares method. However,  $e(t)$  is now a coloured process and this estimator is not

unbiased. Let us assume a true model exists  $\theta^*$ .

$$\begin{aligned}
y(t) &= \bar{\varphi}^T(t)\theta^+ + e(t) \\
\hat{\theta}_{LS} &= \left( \frac{1}{N} \sum_{t=1}^N \bar{\varphi}(t)\bar{\varphi}^T(t) \right) \frac{1}{N} \sum_{t=1}^N \bar{\varphi}(t)y(t) \\
&= \theta^+ + \left( \frac{1}{N} \sum_{t=1}^N \bar{\varphi}(t)\bar{\varphi}^T(t) \right)^{-1} \frac{1}{N} \sum_{t=1}^N \bar{\varphi}(t)e(t) \\
\lim_{N \rightarrow \infty} \hat{\theta}_{LS} &= \theta^* + \Sigma_{\bar{\varphi}}^{-1} r_{\bar{\varphi}e} \\
r_{\bar{\varphi}e} &= E[\bar{\varphi}(t)e(t)] = E \left[ \begin{bmatrix} -y(t-1) \\ \vdots \\ -y(t-n) \\ u(t-1) \\ \vdots \\ u(t-n) \end{bmatrix} e(t) \right] \\
e(t) &= w(t)c_1w(t-1) + \dots + c_nw(t-n) \\
r_{\bar{\varphi}e} &= \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \text{ with } n \text{ zeros}
\end{aligned}$$

$$y(t-1) = f(w(t-1), w(t-2), \dots)$$

$e(t) = f(w(t), w(t-1), \dots, w(t-n))$  so there is a correlation between the terms of  $y$  and those of  $e$ , so the first elements of the above vector are not zero, therefore

$$\lim_{N \rightarrow \infty} \hat{\theta}_{LS} = \theta^* + \Sigma_{\bar{\varphi}}^{-1} r_{\bar{\varphi}e} \neq 0$$

so the estimator is biased.

General model structure:

$$y(t) = G(z^{-1})u(t) + H(z^{-1})w(t)$$

where  $w(t)$  is a zero mean white process with variance  $\sigma_w^2$  and uncorrelated with  $u(t)$

Available measurements:

$$y(1-n), y(2-n), \dots, y(N), u(1-n), u(2-n), \dots, u(N)$$

Prediction error method: find the estimate  $\hat{\theta}$  that minimizes the loss function

$$J(\theta) = \frac{1}{N} \sum_{t=1}^N \varepsilon^2(t, \theta) = \frac{1}{N} \sum_{t=1}^N (y(t) - \hat{y}(t|t-1, \theta))^2$$

where  $\hat{y}(t|t-1, \theta)$  is the *optimal one step ahead prediction* of  $y(t)$ . Optimal (minimal variance) prediction: a prediction of  $y(t)$  given  $\theta$  and the past I/O data up to time  $t-1$  s.t. the variance of the prediction error  $\varepsilon(t)$  is minimal. The model can be rewritten as

$$y(t) = G(z^{-1})u(t) + (H(z^{-1}) - 1)w(t) + w(t)$$

by replacing  $w(t)$  with  $\frac{1}{H(z^{-1})}y(t) - \frac{G(z^{-1})}{H(z^{-1})}u(t)$ , after some steps we get

$$y(t) = \left( 1 - \frac{1}{H(z^{-1})} \right) y(t) + \frac{G(z^{-1})}{H(z^{-1})} u(t) + w(t)$$



from which it is easy to prove

$$\hat{y}(t|t-1, \theta) = \left(1 - \frac{1}{H(z^{-1})}\right) y(t) + \frac{G(z^{-1})}{H(z^{-1})} u(t)$$

For any other predictor  $y^p(t)$  we have

$$\begin{aligned} E[(y(t) - y^p(t))^2] &\geq E[(y(t) - \hat{y}(t|t-1, \theta))^2] \\ E[(y(t) - y^p(t))^2] &= E[(Gu(t) + Hw(t) - y^p(t))^2] \\ &= E\left[\left(\left(1 - \frac{1}{H}\right)y + \frac{G}{H}u + w(t) - y^p(t)\right)^2\right] \end{aligned}$$

We can notice that  $w(t)$  is uncorrelated with all other members of the sum as  $y^p(t)$  depends only on the first  $t-1$  samples, therefore

$$= E\left[\left(\left(1 - \frac{1}{H}\right)y + \frac{G}{H}u - y^p(t)\right)^2\right] + E[w(t)^2] = E\left[\left(\left(1 - \frac{1}{H}\right)y + \frac{G}{H}u - y^p(t)\right)^2\right] + \sigma_w^2$$

$$\lim_{N \rightarrow \infty} J''(\theta) = 2E[\psi(t, \theta)[\psi^T(t, \theta)]] - 2E[\varepsilon(t, \theta) \frac{\partial^2 \varepsilon(t, \theta)}{\partial \theta^2}]$$

assume  $\theta^*$  exists, then

### 6.0.1 Identification of ARMAX models

$$\varepsilon(t, \theta) = -c_1 \varepsilon(t-1, \theta) - c_2 \varepsilon(t-2, \theta), \dots - c_n \varepsilon(t-n, \theta) + y(t) + a_1 y(t-1) + \dots + b_n u(t-n)$$

Identification of ARARX models: the optimal predictor is

$$\hat{y}(t|t-1\theta) = (1 - a(z^{-1})D(z^{-1}))y(t) + B(z^{-1})D(z^{-1})u(t)$$

which can be seen as

$$\hat{y}(t|t-1\theta) = (1 - \bar{A}(z^{-1})) - \bar{B}(z^{-1})u(t)$$

which is equivalent to an ARX model. One could apply LQ estimation to the model and try to identify the model by finding common roots between  $\bar{A}$  and  $\bar{B}$  or if just a predictive model is required LQ estimation itself is sufficient.

### 6.0.2 Statistical properties of PEM estimators

Assume that a true model exists and

- the input is persistently exciting of sufficiently high order
- The Hessian  $J''(\theta)$  is nonsingular at least locally around the minimum points of  $J(\theta)$

In this case, the PEM estimate is consistent (proof skipped)

$$\lim_{N \rightarrow \infty} \hat{\theta} = \theta^*$$

Moreover

$$\sqrt{N}(\hat{\theta} - \theta^*) \rightarrow \mathcal{N}(0, P), \quad \text{for } N \rightarrow \infty$$

with

$$P = \sigma_w^2 (E[\psi(t, \theta^*) \psi(t, \theta^*)^T])^{-1} = \sigma_w^2 \Sigma_\psi^{-1}$$

If  $w(t)$  is gaussian distributed

$$w(t) \sim \mathcal{N}(0, \sigma_w^2)$$

The PEM estimate is also asymptotically efficient, therefore  $P$  is the lowest covariance matrix of the estimate that can be obtained.

### 6.0.3 MISO ARX models

$$u(t) \in \mathbb{R}^r \quad u(t) = \begin{bmatrix} u_1(t) \\ u_2(t) \\ \vdots \\ u_r(t) \end{bmatrix}$$

$$A(z^{-1})y(t) = B_1(z^{-1})u_1(t) + \cdots + B_ru_r(t) + e(t)$$

$$B_1(z^{-1}) = b_{11}z^{-1} + b_{12}z^{-2} + \cdots + b_{1n}z^{-n}$$

$$B_2(z^{-1}) \text{ similar } P = n + rn = (r+1)n$$

$$H = [-H_y(n) \quad H_{u_1}(n) \quad \cdots \quad H_{u_n}(n)]$$

$$\hat{\theta}_{LS} \left( \frac{H^T H}{N} \right)^{-1} \frac{H^T Y}{N}$$

$$H^T H \theta = H^T Y$$

## Chapter 7

# Statistical hypothesis testing



## Chapter 8

# Model complexity selection and regularization

Model complexity selection: given a model class  $\mathcal{M}(\theta)$ , estimate the model complexity  $p$ , i.e. choose the best model class  $\mathcal{M}_p(\theta)$

*Training set*: the set of data used for learning the model

*Underfitting*: the model is not rich enough to fit the data well

*Overfitting*: the model is too rich and adapts too closely to the training data

*Validation set*: the set of data used for evaluating the predictive capabilities of the models obtained with the training set in order to estimate the model complexity.

- General rule: 60-70% of the data used for training set, rest for validation
- when using the training set, the higher the model complexity, the better the data fitting. the prediction error is thus underestimated
- when dealing with real data, the loss function exhibits a monotone decrease in the training set and a *U-shape* in the validation set.

ARX MODEL:

$$\begin{aligned}\hat{\theta}_{LS} &\rightarrow \theta^* \text{ for } N \rightarrow \infty \\ \varepsilon(t) &\rightarrow w(t) \text{ for } N \rightarrow \infty \\ J(\hat{\theta}_{LS}) &= \frac{1}{N} \sum_{t=1}^N \varepsilon(t, \hat{\theta}_{LS})^2 = \hat{\sigma}_\varepsilon^2 \\ J(\hat{\theta}_{LS}) &\rightarrow \sigma_w^2 \text{ for } N \rightarrow \infty \\ n = 4 \quad \theta^* &= [a_1^* \quad a_3^* \quad a_3^* \quad a_4^* \quad b_1^* \quad b_2^* \quad b_3^* \quad b_4^*]\end{aligned}$$

If the validation set is used repeatedly to estimate the model complexity, the prediction error may be underestimated as well. This happens for very complex models like neural networks. For this reason, when dealing with neural networks the dataset is split into three parts. The third set is called *test set* and is used for model assessment, i.e. for testing the predictive capabilities of the final chosen model. This requires a very large number of available samples. It is possible to design criteria that allow to estimate the model complexity by using the training test

### 8.0.1 The F-test

Let  $\mathcal{M}_{p_1}(\theta), \mathcal{M}_{p_2}(\theta)$  such that  $p_1 < p_2$  ( $p = 2n$  for ARX models,  $p = n$  for AR models etc.) Consider the test quantity

$$x = N \frac{J(\hat{\theta}_N^1) J(\hat{\theta}_N^2)}{J(\hat{\theta}_N^2)}$$

where  $\hat{\theta}_N^1, \hat{\theta}_N^2$  are PEM (or LS) estimates: intuitively:

- $x$  large: the decrease in the loss function is significant, hence  $\mathcal{M}_{p_2}(\theta)$  is better
- $x$  small:  $\mathcal{M}_{p_1}(\theta)$  and  $\mathcal{M}_{p_2}(\theta)$  are almost equivalent so that  $\mathcal{M}_{p_1}(\theta)$  should be chosen according to the parsimony principle

How to quantify "large" and "small"?

- If  $\mathcal{M}_{p_1}(\theta)$  is not large enough to include the true system:

$$J(\hat{\theta}_N^1) - J(\hat{\theta}_N^2) \text{ is } O(1) \quad x \text{ is of magnitude } N$$

- If  $\mathcal{M}_{p_1}(\theta)$  is large enough:

$$x \rightarrow \chi^2(p_2 - p_1), \quad \text{for } N \rightarrow \infty$$

The following statistical test can be performed:

$$\left\{ \begin{array}{l} H_0 : \mathcal{M}_{p_1}(\theta) \text{ is suitable to describe the system} \\ H_1 : \mathcal{M}_{p_1}(\theta) \text{ is not suitable} \end{array} \right.$$

that is:

$$\left\{ \begin{array}{l} H_0 : x \leq \chi^2(p_2 - p_1) \\ H_1 : \text{not } H_0 \end{array} \right.$$

after the choice of the significance level:

$$\left\{ \begin{array}{l} x \leq \chi_\alpha^2(p_2 - p_1) \implies \text{accept } H_0 \\ x > \chi_\alpha^2(p_2 - p_1) \implies \text{accept } H_1 \end{array} \right.$$

### 8.0.2 The final prediction error (FPE) criterion

Let  $\hat{\theta}_N$  be a PEM (or LS) estimate of a model of complexity  $p$  and assume that it is then used to predict future data. Assume also that a true model exists and consider the prediction error variance (the expectation is with respect to future data):

$$V(V\hat{\theta}_N) = E \left[ (y(t) - \hat{y}(t|t-1, \hat{\theta}_N))^2 \right]$$

by replacing  $y(t) = \varphi^T(t)\theta^* + w(t)$  in  $V(\hat{\theta}_N)$  and computing the expectation we get

$$V(\hat{\theta}_N) = \sigma_\omega^2 + (\hat{\theta}_N - \theta^*)^T \Sigma_\varphi (\hat{\theta}_N - \theta^*)$$

consider now the criterion function

$$FPE = E[V(\hat{\theta}_N)]$$

where the expectation is with respect to past data. By taking into account that:

•

$$E[(\hat{\theta}_N - \theta^*)^T \Sigma_\varphi (\hat{\theta}_N - \theta^*)] = E[\text{trace}(\Sigma_\varphi (\hat{\theta}_N - \theta^*) (\hat{\theta}_N - \theta^*)^T)]$$

1

- Asymptotically:  $\sqrt{N}(\hat{\theta}_N - \theta^*) \sim \mathcal{N}(0, \sigma_\Omega^2 \Sigma_\varphi^{-1})$

it is easy to obtain

$$FPE \approx \sigma_\omega^2 \left( 1 + \frac{p}{N} \right)$$

in practice we need an asymptotically unbiased estimate of  $\sigma_\omega^2$

$$\hat{\sigma}_\omega^2 = \frac{1}{N-p} \sum_{t=1}^N \varepsilon(t, \hat{\theta}_N)^2$$

so that

$$FPE = \hat{\sigma}_e^2 \frac{N+p}{N} = \frac{N+p}{N-p} J(\hat{\theta}_N) = \frac{N \left( 1 + \frac{p}{N} \right)}{N \left( 1 - \frac{p}{N} \right)} J(\hat{\theta}_N) = J(\hat{\theta}_N) + \frac{2 \frac{p}{N}}{1 - \frac{p}{N}} J(\hat{\theta}_N)$$

---

<sup>1</sup>  $E[\text{trace}(v^T A v)] = E[\text{trace}(A v v^T)]$

Note that for large  $N$ :

$$FPE \approx J(\hat{\theta}_N) + \frac{2p}{N}J(\hat{\theta}_N)$$

for large  $N$ , this criterion belongs to the family of *criteria with complexity terms* (terms that penalise complex models).

### 8.0.3 Criteria with complexity terms

These criteria are obtained by penalizing in some way the decrease of  $J(\hat{\theta}_N)$  with increasing orders. The order giving the smallest value for the criterion is selected. General form:

$$V(\hat{\theta}_N) = N \log J(\hat{\theta}_N) + f(N, p)$$

where  $f(N, p)$  penalizes high order models.

### 8.0.4 Akaike information criterion (AIC)

$$AIC = N \log J(\hat{\theta}_N) + 2p$$

AIC and FPE are asymptotically equivalent. They do not give consistent estimates of  $n$  (the probability of overestimating the order is non-null). To get consistent estimates, the penalizing function must be such that

$$\begin{cases} f(N, p) = kp g(N) \\ \lim_{N \rightarrow \infty} g(N) = \infty \\ \lim_{N \rightarrow \infty} \frac{g(N)}{N} = 0 \end{cases}$$

### Minimum description length (MDL) criterion

$$MDL = N \log J(\hat{\theta}_N) + 2p \log(N)$$

MDL leads, in general, to models of lower complexity wrt AIC and FPE. Even though the derivation is different, the MDL approach is formally equivalent to the *bayesian information criterion (BIC)*





## Chapter 9

# model assesment (validation)

Consists in evaluating the capability of the identified model to describe the process that has generated the data in a way compatible with its planned use.

Linear regression models:

$$y(t) = \varphi^T(t)\theta^* + w(t), \quad w(t) \text{ zero mean and white} \quad E[u(t)w(t-\tau)] = 0, \forall \tau$$

The PEM estimate  $\hat{\theta}_N$  is consistent, then

$$\hat{\theta}_N \rightarrow \theta^* \text{ for } N \rightarrow \infty, \quad \varepsilon(t, \hat{\theta}_N) = y(t) - \varphi^T(t)\hat{\theta}_N \rightarrow w(t) \text{ for } N \rightarrow \infty$$

If we assume that our real data are well described by a linear regression model, we can make the following assumptions about the residual  $\varepsilon(t, \hat{\theta}_N)$ :

1.  $\varepsilon(t, \hat{\theta}_N)$  is a zero mean white process
2.  $\varepsilon(t, \hat{\theta}_N)$  is uncorrelated with the input signal  $u(t)$

It is thus possible to perform (on the training set) the following *test on residuals*:

- test of whiteness of  $\varepsilon(t, \hat{\theta}_N)$
- test of cross-correlation between  $\varepsilon(t, \hat{\theta}_N)$  and  $u(t)$

### 9.1 Whiteness test

Sequence of residuals:  $\varepsilon(1, \hat{\theta}_N), \varepsilon(2, \hat{\theta}_N), \dots, \varepsilon(N, \hat{\theta}_N)$

$$\begin{cases} H_0 : \varepsilon(t, \hat{\theta}_N) \text{ is a zero mean white process} \\ H_1 : \text{not } H_0 \end{cases}$$

Consider the sample variance  $\hat{r}_\varepsilon(0)$  and the first  $m$  sample autocorrelations of  $\varepsilon(t)$  and define the vector

$$\hat{r}_\varepsilon = \begin{bmatrix} \hat{r}_\varepsilon(1) \\ \hat{r}_\varepsilon(2) \\ \vdots \\ \hat{r}_\varepsilon(m) \end{bmatrix}$$

Under  $H_0$ :

$$\hat{r}_\varepsilon(0) \rightarrow \sigma_w^2 \text{ for } N \rightarrow \infty, \quad \hat{r}_\varepsilon(\tau) \rightarrow 0 \text{ for } N \rightarrow \infty, \forall \tau \neq 0$$

and it is possible to prove that

$$\sqrt{N} \xrightarrow[N \rightarrow \infty]{\sim} \mathcal{N}(0, P), \quad P = \lim_{N \rightarrow \infty} E[N\hat{r}_\varepsilon\hat{r}_\varepsilon^T] = \sigma_w^4 I \quad (*)$$

as a consequence

$$x = N \frac{\hat{r}_\varepsilon^T \hat{r}_\varepsilon}{\hat{r}_\varepsilon^2(0)} \xrightarrow[N \rightarrow \infty]{\sim} \chi^2(m)$$

This leads to the statistical test

$$\left\{ x \leq \chi_\alpha^2(m) \implies \text{accept } H_0 \right. \quad \left. x > \chi_\alpha^2(m) \implies \text{accept } H_1 \right.$$

where  $\alpha$  is the chosen significance level. Define the normalized text quantities

$$\hat{\gamma}(\tau) = \frac{\hat{r}_\varepsilon(\tau)}{\hat{r}_\varepsilon(0)}, \quad \tau = 1, 2, \dots, m$$

From (\*) it follows that:

$$\sqrt{N} \hat{\gamma}(\tau) \xrightarrow[N \rightarrow \infty]{\sim} \mathcal{N}(0, 1), \quad \tau = 1, 2, \dots, m$$

so that a set of  $m$  gaussian tests can also be performed. The whiteness model test can be of help in model order estimation.

## 9.2 Test of cross-correlation

$$\begin{cases} H_0 : \varepsilon(t\hat{\theta}_N) \text{ and } u(t) \text{ are uncorrelated} \\ H_1 : \text{not } H_0 \end{cases}$$

Consider the following vector of sample cross correlations:

$$\hat{r}_{\varepsilon u} = \begin{bmatrix} \hat{r}_{\bar{\tau}+1} \\ \hat{r}_{\bar{\tau}+2} \\ \vdots \\ \hat{r}_{\bar{\tau}+m} \end{bmatrix}$$

Consider also the vector

$$\varphi_u(t, m) = [u(t-1) \quad u(t-2) \quad \dots \quad u(t-m)]^T$$

and the sample autocorrelation matrix  $\hat{\Sigma}_u(m)$ , which is an estimate of  $E[\varphi_u(t, m)\varphi_u^T(t, m)]$

It is possible to prove that

$$\sqrt{N} \hat{r}_{\varepsilon u} \xrightarrow[N \rightarrow \infty]{\sim} \mathcal{N}(0, P), \quad P = \lim_{N \rightarrow \infty} E[N \hat{r}_{\varepsilon u} \hat{r}_{\varepsilon u}^T] = \sigma_w^2 \Sigma_u(m)$$

As a consequence

$$x = N \frac{\hat{r}_{\varepsilon u}^T \hat{\Sigma}_u^{-1} \hat{r}_{\varepsilon u}}{\hat{r}_{\varepsilon u}(0)} \xrightarrow[N \rightarrow \infty]{\sim} \chi^2(m)$$

This leads to the statistical test

$$\left\{ H_0 : x \leq \chi_\alpha^2(m) \implies \text{accept } H_0 \right. \quad \left. x > \chi_\alpha^2(m) \implies \text{accept } H_1 \right.$$

where  $\alpha$  is the chosen significance level.  $\bar{\tau}$  must be chosen with care. In some methods  $\hat{r}_{\varepsilon u}(\tau)$  is constrained to be zero for some values of  $\tau$  by construction. Let us consider the quantity it follows that

We define the normalized test quantities

$$\hat{\gamma}(\tau) = \frac{\hat{r}_{\varepsilon u}(\tau)}{\sqrt{\hat{r}_\varepsilon(0)\hat{r}_u(0)}}, \quad \tau = 1, 2, \dots, m$$

It follows that

$$\sqrt{N} \hat{\gamma} \xrightarrow[N \rightarrow \infty]{\sim} \mathcal{N}(0, 1), \quad \tau = 1, 2, \dots, m$$

so that a set of  $m$  gaussian tests can also be performed. Once the significance level  $\alpha$  has been chosen,  $m$  gaussian tests are performed and the final decision is taken by considering the Anderson test.

## Chapter 10

# Maximum likelihood estimation

Let  $y(t)$  be an observed stochastic process whose distribution depends on an unknown parameter vector  $\theta$ . The *likelihood function*  $L(\theta)$  is the pdf of  $y(t)$  given  $\theta$ :

$$L(\theta) = f(y(t)|\theta)$$

The maximum likelihood estimate of  $\theta$  is

$$\hat{\theta}_{ML} = \arg \max_{\theta} L(\theta)$$

- $\hat{\theta}_{ML}$  is the vector  $\theta$  that makes the a posteriori probability of the observations as large as possible
- ML estimates require the knowledge of the conditional pdf and are often difficult to calculate analytically
- In the gaussian case the problem often becomes more tractable. For dynamic systems, the problem can be simplified if the effect of initial conditions is neglected.

### 10.1 the gaussian case

$$y(t) = G(z^{-1})u(t) + H(z^{-1})w(t) \quad w(t) \sim \mathcal{N}(0, \sigma_w^2)$$

Neglecting the effect of the initial conditions, the pdf  $f(y(t)|\theta)$  can be replaced with  $f(w(t)|\theta)$ . By introducing also the approximation  $\varepsilon(t, \theta) \approx w(t)$  we can write

$$L(\theta) = \frac{1}{\sqrt{(2\pi\sigma_\varepsilon^2)^N}} \exp \left( -\frac{1}{2\sigma_\varepsilon^2} \sum_{t=1}^N \varepsilon^2(t, \theta) \right)$$

By considering the logarithmic likelihood function (log-likelihood function)

$$\log L(\theta) = -\frac{N}{2} \log 2\pi - \frac{N}{2} \log \sigma_\varepsilon^2 - \frac{1}{2\sigma_\varepsilon^2} \sum_{t=1}^N \varepsilon^2(t, \theta) \quad (*)$$

Assume now that  $\sigma_\varepsilon^2$  is an additional parameter to be estimated  $\implies \log L(\theta, \sigma_\varepsilon^2)$ . First maximize wrt  $\sigma_\varepsilon^2$  by solving

$$\frac{\partial \log L(\theta, \sigma_\varepsilon^2)}{\partial \sigma_\varepsilon^2} = 0$$

We get

$$\sigma_\varepsilon^2 = \frac{1}{N} \sum_{t=1}^N \varepsilon^2(t, \theta) = J(\theta)$$

The inserting the above expression in (\*) we find

$$\log L(\theta) = -\frac{N}{2} \log 2\pi - \frac{N}{2} - \frac{N}{2} \log J(\theta)$$

$\hat{\theta}_{ML}$  is thus obtained by minimizing  $J(\theta) \implies$  the PEM estimate can be interpreted as the ML estimate provided that  $w(t)$  is gaussian distributed

- As a consequence, the least squares estimate of static models, FIR models, ARX models and AR models, can be interpreted as the ML estimate provided that  $w(t)$  is gaussian distributed
- Strictly speaking, for dynamic models  $L(\theta)$  is not the exact likelihood function, because of the effect of the initial conditions.
- When the number of samples is large, the difference between the exact and the approximated likelihood functions becomes small.

### 10.1.1 the Cramér-Rao lower bound

Let  $Y$  be a stochastic vector and  $L(\theta^*) = g(Y|\theta)$ . For an arbitrary unbiased estimate  $\hat{\theta}(Y)$  of  $\theta^*$  determined from  $Y$  we have

$$\text{cov}(\hat{\theta}(Y)) \geq \left[ E \left( \frac{\partial \log L(\theta^*)}{\partial \theta^*} \right)^T \frac{\partial \log L(\theta^*)}{\partial \theta^*} \right]^{-1} = \left[ E \frac{\partial^2 \log L(\theta^*)}{\partial \theta^{*2}} \right]^{-1} = F^{-1}$$

Where  $F$  is called the *Fisher information matrix*. The matrix  $F^{-1}$  is thus a lower bound on the covariance matrix of the estimate and is called *Cramér-Rao lower bound (CRLB)*.

- if an estimate attains the CRLB, it is efficient
- the maximum likelihood estimator is (in general) consistent and efficient
- The least squares estimate of static models, FIR models, ARX models and AR models is (asymptotically) efficient provided that  $w(t)$  is gaussian distributed
- the PEM estimate of ARMAX, ARMA and ARARX models is asymptotically efficient provided that  $w(t)$  is gaussian distributed.

# Chapter 11

## Classification: probabilistic models

Classification: assigning the input  $u(t) \in \mathcal{U}$  to one of  $M$  classes  $C_1, C_2, \dots, C_M$ . The input is numerical and the output is categorical. Available data:

$$(u(1), y(1)), (u(2), y(2)), \dots, (u(N), y(N))$$

A possible way to represent class labels numerically is

$$y(t) = y^i \quad \text{if } u(t) \in C_i$$

where  $y^i$  is a  $M \times 1$  vector whose elements are zero except for the  $i$ -th element, which is equal to 1. For a two-class problem we will use

$$y(t) = \begin{cases} 1 & \text{if } u(t) \in C_1 \\ 0 & \text{if } u(t) \in C_2 \end{cases}$$

The input space  $\mathcal{U}$  is divided into  $M$  regions labeled according to the classification rule. The boundaries of these regions are called *decision boundaries*.

Deterministic models: the decision boundaries are explicitly modeled. This involves the construction of a *discriminant function* that directly assigns each input  $u(t)$  to a specific class.

Probabilistic models: first, the posterior probabilities  $P(C_i|u(t)) = P(y(t) = y^i|u(t))$ ,  $i = 1, 2, \dots, M$  are estimated. Then, these probabilities are used to decide to which class a given input  $u(t)$  belongs.

The performance for a classifier can be evaluated by computing the error rate:

$$Er = \frac{1}{N} \sum_{t=1}^N I(y(t), \hat{y}(t))$$

where  $\hat{y}(t)$  is the prediction of  $y(t)$  provided by the classifier and  $I(y(t), \hat{y}(t))$  is an *indicator variable*

$$I(y(t), \hat{y}(t)) = \begin{cases} 1 & \text{if } y(t) \neq \hat{y}(t) \\ 0 & \text{if } y(t) = \hat{y}(t) \end{cases}$$

The error rate is thus the fraction of misclassified observations.

- The error rate can be evaluated on the training set (training error rate), on the validation set (validation error rate) and on the test set (test error rate)
- A good classifier is one for which the validation (test) error rate is small
- As for regression problems, choosing the model complexity by using the training error rate leads to overfitting.

### 11.1 The Bayes classifier

The Bayes classifier assigns each input to the most likely class. Given an input  $u(t) \in \mathcal{U}$ , it will be assigned to the class  $C_i$  with largest posterior probability.

$$C_i : \arg \max_{i=1,2,\dots,M} \{P(C_i|u(t)) = P(y(t) = y^i|u(t))\}$$

It maximizes the conditional probability that  $y(t) = y^i$  given the observed input value  $u(t)$ .

- In a two-class problem, the Bayes classifier consists in assigning  $u(t)$  to class  $C_1$  if  $P(C_1|u(t)) > 0.5$  and to class  $C_2$  otherwise
- It is possible to show that this classifier leads to the minimum classification (prediction) error on a future unseen dataset (validation or test set)
- When dealing with real data, the conditional distribution of  $\mathcal{Y}$  given  $\mathcal{U}$  is usually unknown so that the Bayes classifier cannot be used in practice
- Probabilistic models use estimated probabilities  $\hat{P}(C_i|u(t)), i = 1, 2, \dots, M$  rather than the true ones

## 11.2 Logistic regression

Consider a two-class classification problem with  $r$  inputs ( $u(t) = [u_1(t) \dots u_r(t)]^T$ ) and assume a linear decision boundary. This boundary is a hyperplane in  $\mathbb{R}^r$  described by the equation

$$\beta_0 + \beta_1 u_1(t) + \beta_2 u_2(t) + \dots + \beta_r u_r(t) = 0$$

that is

$$\varphi^T(t)\theta = 0$$

with  $\varphi(t) = [1 \ u_1(t) \ u_2(t) \ \dots \ u_r(t)]^T$  and  $\theta = [\beta_0 \ \beta_1 \ \dots \ \beta_r]^T$  How to model  $P(C_1|u(t))$ ? We need a function  $f(z(t))$ , where  $z(t) = \varphi^T(t)\theta$ , such that

- $f(z)$  takes values in the range  $(0, 1)$
- $f(0) = 0.5$
- $f(z) > 0.5 (f(z) < 0.5)$  when  $z > 0 (z < 0)$
- $f(z) \rightarrow 1$  when  $z \rightarrow \infty$
- $f(z) \rightarrow 0$  when  $z \rightarrow -\infty$

Logistic sigmoid function:

$$f(z) = \frac{e^z}{1 + e^z} = \frac{1}{1 + e^{-z}}$$

We have

$$P(C_1|u(t)) = f(z(t)), \quad P(C_2|u(t)) = 1 - f(z(t)) = \frac{1}{1 + e^z}$$

and

$$z(t) = \log \frac{P(C_1|u(t))}{P(C_2|u(t))} = \log \frac{f(z)}{1 - f(z)}$$

How to estimate  $\theta$ ?

Find the estimate  $\hat{\theta}$  that maximizes the posterior probability of the observation  $Y = [y(1) \ y(2) \ \dots \ y(N)]^T \implies$  maximum likelihood estimation

$$P(Y|\theta) = \prod_{t=1}^N P(C_1|u(t))^{y(t)} P(C_2|u(t))^{1-y(t)} = \prod_{t=1}^N f(z(t))^{y(t)} (1 - f(z(t)))^{1-y(t)}$$

Log-likelihood function:

$$\log P(Y|\theta) = \sum_{t=1}^N [y(t) \log f(z(t)) + (1 - y(t)) \log(1 - f(z(t)))]$$

Maximizing the above function is equivalent to minimizing its negative, so we have to find the estimate  $\hat{\theta}$  minimizing

$$J(\theta) = -\log P(Y|\theta)$$

To solve this optimization problem we can use the Newton-Raphson algorithm:

$$\hat{\theta}^{k+1} = \hat{\theta}^k - \eta^k J''(\hat{\theta}^k)^{-1} J'(\hat{\theta}^k)^T$$

By using

$$\frac{\partial f(z)}{\partial z} = f(z)(1 - f(z))$$

It is possible to show that

$$J'(\theta) = \frac{J(\theta)}{\partial \theta} = \sum_{t=1}^N [f(z(t)) - y(t)] \varphi(t) = \Phi^T (F(\theta) - Y)$$

where

$$\Phi = \begin{bmatrix} 1 & u_1(1) & \cdots & u_r(1) \\ 1 & u_1(2) & \cdots & u_r(2) \\ \vdots & \vdots & \ddots & \vdots \\ 1 & u_1(N) & \cdots & u_r(N) \end{bmatrix} \quad Y = \begin{bmatrix} y(1) \\ y(2) \\ \vdots \\ y(N) \end{bmatrix} \quad F(\theta) = \begin{bmatrix} f(z(1)) \\ f(z(2)) \\ \vdots \\ f(z(N)) \end{bmatrix}$$

and

$$J''(\theta) = \frac{\partial^2 J(\theta)}{\partial \theta^2} = \sum_{t=1}^N f(z(t))(1 - f(z(t))) \varphi(t) \varphi^T(t) = \Phi^T W(\theta) \Phi$$

where  $W(\theta)$  is the diagonal matrix

$$W(\theta) = \text{diag} [f(z(1))(1 - f(z(1))) \quad f(z(2))(1 - f(z(2))) \quad \cdots \quad f(z(N))(1 - f(z(N)))]$$

The iterative algorithm is

$$\hat{\theta}^{k+1} = \hat{\theta}^k - \eta^k (\Phi^T W(\theta) \Phi)^{-1} \Phi^T (F(\theta) - Y)$$

- the elements of  $W(\theta)$  are in the range  $(0, 1)$ , therefore the hessian  $\Phi^T W(\theta) \Phi$  is positive definite and the loss function  $J(\theta)$  has a unique minimum.
- $\hat{\theta}^0$  is often chosen as a vector of zeros. It can also be chosen randomly.

### How to classify a new input $u(t)$ ?

Once an estimate  $\hat{\theta}$  has been computed and a new input  $u(t)$  becomes available:

1. Compute  $\hat{z}(t) = \varphi^T(t) \hat{\theta}$
2. Compute  $f(\hat{z}(t)) = \hat{P}(C_1 | u(t)) = \frac{e^{\hat{z}}}{1 + e^{\hat{z}}}$ , which is an estimate of  $P(C_1 | u(t))$
3. If  $f(\hat{z}(t)) > 0.5$  assign  $u(t)$  to class  $C_1$  ( $\hat{y}(t) = 1$ ), otherwise assign it to class  $C_2$  ( $\hat{y}(t) = 0$ )

Depending on the specific classification problem, the decision criterion could be different from that described in the above step 3

#### 11.2.1 The gradient descent algorithm

$$\hat{\theta}^{k+1} = \hat{\theta}^k - \eta^k J'(\hat{\theta}^k)$$

- It only requires computation of the gradient of the loss function
- At each step  $\hat{\theta}^k$  is moved in the direction of the greatest rate of decrease of the loss function  $J(\theta)$
- The step size  $\eta^k$  (also called *learning rate*) plays a very important role. Both very large and very small values can lead to inefficiency. A variable step size is often adopted
- The algorithm can converge to a local minimum
- The procedure can be repeated by starting from different initial guesses
- *Stopping criterion*:  $\|\hat{\theta}^{k+1} - \hat{\theta}^k\| < \delta$  or  $|J(\hat{\theta}^{k+1}) - J(\hat{\theta}^k)| < \delta$  or  $\|J'(\hat{\theta}^k)\| < \delta$  where  $\delta$  is a small positive number

For the logistic regression problem, the iterative algorithm is

$$\hat{\theta}^{k+1} = \hat{\theta}^k - \eta^k \Phi^T (F(\theta) - Y)$$

### Stochastic gradient descent

Instead of using all samples, a smaller number of samples is randomly picked for gradient computation at each iteration.

### Confusion matrix

Suppose a training set is used to obtain an estimate on a two-class classification problem, and a validation set is used to verify the predictive capabilities of the model, with 5000 samples corresponding to class  $C_1$  (H) and 1000 samples belonging to class  $C_2$  (F) the confusion matrix is:

	H	F	Total
H	4950	250	5200
F	50	750	800
Total	5000	1000	6000

The error rate is:

$$Er = \frac{300}{6000} = 5\%$$

The error rate concerning the first class is

$$Er_H = \frac{50}{5000} = 1\%$$

The error rate concerning the second class is

$$Er_H = \frac{250}{1000} = 25\%$$

If these errors are tolerable, we may keep the chosen decision rule, otherwise we may change it. New decision rule:  $\hat{P}(C|u(t)) > 0.75 \implies u(t) \in C_1$  leading to a new confusion matrix:

	H	F	Total
H	4750	110	4860
F	250	890	1140
Total	5000	1000	6000

The error rate is:

$$Er = \frac{360}{6000} = 6\%$$

The error rate concerning the first class is

$$Er_H = 5\%$$

The error rate concerning the second class is

$$Er_H = 11\%$$

### 11.2.2 Multiclass problems

One-vs-all: the multiclass problem is split into  $M$  two-class problems. The generic  $i$ -th binary problem involves the classes  $C_i$  and  $\bar{C}_i = \{C_1, \dots, C_{i-1}, C_{i+1}, \dots, C_M\}$ . Therefore,  $M$  conditional probabilities  $\hat{P}(C_i|u(t))$ ,  $i = 1, 2, \dots, M$  are estimated. A new input  $u(t)$  is assigned to the class  $C_k$  such that

$$k = \arg \max_{i \in \{1, 2, \dots, M\}} \hat{P}(C_i|u(t))$$

Note that the sum of the posterior probabilities  $\hat{P}(C_i|u(t))$  is not 1 as multiple two-class classification problems are being considered. Multiclass logistic regression: the posterior probabilities are modelled using the *normalized exponential* (or *softmax function*)

$$P(C_i|u(t)) = \frac{e^{z_i(t)}}{e^{z_1(t)} + e^{z_2(t)} + \dots + e^{z_M(t)}}, \quad i = 1, 2, \dots, M$$



where

$$z_i(t) = \varphi^T(t)\theta_i, \quad i = 1, 2, \dots, M$$

Note that  $\sum_{i=1}^M P(C_i|u(t)) = 1$

The input space  $\mathcal{U} = \mathbb{R}^r$  can thus be divided into  $M$  regions defined by a set of  $M$  hyperplanes. Each hyperplane represents the linear decision boundary between two classes. More precisely, the hyperplane separating classes  $k$  and  $j$  is described by the equation

$$(z_k(t) - z_j(t)) = 0 \quad \text{that is} \quad \varphi^T(t)(\theta_k - \theta_j) = 0$$

An alternative consists in selecting a single class as the baseline. If we select the  $M$ -th class the posterior probabilities are given by:

$$P(C_i|u(t)) = \frac{e^{z_i(t)}}{1 + e^{z_1(t)} + e^{z_2(t)} + \dots + e^{z_{M-1}(t)}}, \quad i = 1, 2, \dots, M-1$$

and

$$P(C_M|u(t)) = \frac{1}{1 + e^{z_1(t)} + e^{z_2(t)} + \dots + e^{z_{M-1}(t)}}$$

Both alternatives lead to the same results. For both methods, the parameter vectors  $\theta_1, \theta_2, \dots$  are estimated by using maximum likelihood and Newton's method (or gradient descent).

### 11.2.3 Dealing with nonlinear boundaries

By means of a nonlinear transformation in the input space

$$u(t) \in \mathcal{U} \rightarrow \bar{u}(t) = g(u(t)) \in \bar{\mathcal{U}}$$

a nonlinear boundary in  $\mathcal{U}$  can be mapped into a linear one in  $\bar{\mathcal{U}}$ . The logistic regression is then applied to the transformed inputs.

- In general, the dimension of  $\bar{\mathcal{U}}$  is greater than that of the original input space  $\mathcal{U}$  and the complexity of the model increases
- Nonlinear transformations cannot remove class overlap. They can increase the level of overlap or create overlap where none existed in  $\mathcal{U}$

### 11.2.4 Linear discriminant analysis

The posterior probabilities are estimated through the Bayes theorem:

$$P(C_i|u(t)) = \frac{P(u(t)|C_i)P(C_i)}{P(u(t))} = \frac{P(u(t)|C_i)P(C_i)}{\sum_{k=1}^M P(u(t)|C_k)P(C_k)}$$

if  $u(t)$  is a continuous random variable,  $P(u(t)|C_i)$  is replaced with the conditional pdf  $f_u(u|C_i)$ . In fact,  $f_u(u|C_i)du$  is the probability that  $u(t) \in [u, u + du]$ .

How does it work?

- Assume a specific probability density function for  $P(u(t)|C_i), i = 1, 2, \dots, M$  and estimate the parameters characterizing such distribution
- Estimate the prior probabilities  $P(C_1), P(C_2), \dots, P(C_M)$ .
- Compute estimates of the posterior probabilities  $P(C_i|u(t)), i = 1, 2, \dots, M$
- Assign a new input  $u(t)$  to the class  $C_k$  for which  $\hat{P}(C_i|u(t)), i = 1, 2, \dots, M$  is largest

### 11.2.5 Gaussian class densities with common covariance matrix

Two class classification problem

$$f_u(u|C_i) = \frac{1}{\sqrt{(2\pi)^r \det(\Sigma)}} \exp\left(-\frac{(u(t) - \mu_i)^T \Sigma^{-1} (u(t) - \mu_i)}{2}\right), \quad i = 1, 2$$

It follows that

$$P(C_1|u(t)) = \frac{f_u(u|C_1)P(C_1)}{f_u(u|C_1)P(C_1) + f_u(u|C_2)P(C_2)} = \frac{1}{1 + e^{-z(t)}}$$

where

$$z(t) = u^T(t)\Sigma^{-1}(\mu_1 - \mu_2) + \frac{1}{2}\mu_2^T \Sigma^{-1} \mu_2 - \frac{1}{2}\mu_1^T \Sigma^{-1} \mu_1 + \log \frac{P(C_1)}{P(C_2)} = \varphi^T(t)\theta$$

with

$$\varphi(t) = [1 \quad u_1(t) \quad u_2(t) \quad \cdots \quad u_r(t)]^T, \quad \theta = [\beta_0 \quad \beta_1 \quad \cdots \quad \beta_r]^T = [\beta_0 \beta^T]^T$$

and

$$\beta_0 = \frac{1}{2}\mu_2^T \Sigma^{-1} \mu_2 - \frac{1}{2}\mu_1^T \Sigma^{-1} \mu_1 + \log \frac{P(C_1)}{P(C_2)}, \quad \beta = \Sigma^{-1}(\mu_1 - \mu_2)$$

As for the logistic regression, the decision boundary is given by the hyperplane

$$z(t) = 0 \quad \text{that is} \quad \varphi^T(t)\theta = 0$$

The estimates required in steps 1 and 2 can be computed as

$$\begin{aligned} \hat{\mu}_1 &= \frac{1}{N_1} \sum_{t=1}^N y(t)u(t) & \hat{\mu}_2 &= \frac{1}{N_2} \sum_{t=1}^N y(t)u(t) \\ \hat{\Sigma} &= \frac{1}{N} \sum_{t=1}^N y(t)(u(t) - \hat{\mu}_1)(u(t) - \hat{\mu}_1)^T + \frac{1}{N} \sum_{t=1}^N y(t)(u(t) - \hat{\mu}_2)(u(t) - \hat{\mu}_2)^T \\ \hat{P}(C_1) &= \frac{N_1}{N_1 + N_2}, & \hat{P}(C_2) &= 1 - \hat{P}(C_1) \end{aligned}$$

where  $N_1$  is the number of inputs belonging to class  $C_1$  and  $N_2 = N - N_1$  is the number of inputs belonging to class  $C_2$ . A new input  $u(t)$  is assigned to class  $C_1$  if  $\hat{P}(C_1|u(t)) > 0.5$  and to class  $C_2$  otherwise. As for logistic regression, a different decision criterion can be used.

### Multiclass problem

The input space  $\mathcal{U} = \mathbb{R}^r$  can thus be divided into  $M$  regions defined by a set of  $M$  hyperplanes. Each hyperplane represents the linear decision boundary between two classes. More precisely, the hyperplane separating classes  $k$  and  $j$  is described by the equation

$$(z_k(t) - z_j(t)) = 0 \quad \text{that is} \quad \varphi^T(t)(\theta_k - \theta_j) = 0$$

The required estimates are computed as follows A new input  $u(t)$  is assigned to the class  $C_k$  such that

$$k = \arg \max_{i \in \{1, 2, \dots, M\}} \hat{P}(C_i|u(t))$$

### 11.2.6 Gaussian class densities with different covariance matrices

In this case, the input space  $\mathcal{U}$  can be divided into  $M$  regions defined by  $M$  quadratic decision boundaries. The quadratic boundary separating classes  $k$  and  $j$  is described by the equation

$$(z_k(t) - z_j(t)) = 0 \quad \text{that is} \quad \varphi^T(t)(\theta_k - \theta_j) = 0$$

For this reason in this case the approach is called *quadratic discriminant analysis*

## Chapter 12

# Classification: deterministic models

The decision boundaries are explicitly modeled. This involves the construction of a discriminant function that directly assigns each input  $u(t)$  to a specific class

Assumptions:

1. two-class classification problems
2. Linear decision boundaries

Because of 2, the discriminant function  $z(t)$  is such that

$$z(t) = \varphi^T(t)\theta = u^T(t)\beta + \beta_0$$

where  $u(t) = [u_1(t) \ u_2(t) \ \cdots \ u_r(t)]^T$  and  $\beta = [\beta_1 \ \beta_2 \ \cdots \ \beta_r]^T$ . For the output  $y(t)$  we will use

$$y(t) = \begin{cases} 1 & \text{if } u(t) \in C_1 \quad (z(t) > 0) \\ -1 & \text{if } u(t) \in C_2 \quad (z(t) < 0) \end{cases}$$

The equation  $z(t) = 0$ , i.e.

$$u^T(t)\beta + \beta_0 = 0$$

defines a hyperplane  $\mathcal{H} \in \mathcal{U} = \mathbb{R}^r$   $u_A, u_B \in \mathcal{H} \implies \beta^T(u_A - u_B) = 0$  implies that  $\beta$  is orthogonal to  $\mathcal{H}$

$$u_i \in \mathcal{H} \implies u_i^T \beta + \beta_0 = 0$$

so that

$$\frac{\beta^T u}{\|\beta\|} = -\frac{\beta_0}{\|\beta\|}$$

is the perpendicular distance of  $\mathcal{H}$  from the origin

$$u \in \mathcal{U} \implies z(t) = \beta^T u + \beta_0$$

so that

$$\frac{z(t)}{\|\beta\|}$$

is the signed distance of  $u(t)$  to  $\mathcal{H}$  finally, note that, due to the classification rule

$$y(t) \frac{z(t)}{\|\beta\|} \geq 0$$

is the distance (unsigned) of  $u(t)$  to the hyperplane  $\mathcal{H}$  ( $y(t)z(t) \geq 0 \forall t$ )

### 12.0.1 Separating hyperplanes

A dataset whose classes can be separated exactly by linear decision boundaries are said to be linearly separable. These linear decision boundaries are called separating hyperplanes. In general, if the dataset is linearly separable, there exists an infinite number of separating hyperplanes. Once a separating hyperplane  $z(t) = 0$  has been found, a new input  $u(t)$  is classified as follows:

$$u(t) \in \begin{cases} C_1 & \text{if } z(t) > 0 \\ C_2 & \text{if } z(t) < 0 \end{cases} \implies y(t) = \begin{cases} 1 & \text{if } z(t) > 0 \\ 2 & \text{if } z(t) < 0 \end{cases}$$

### 12.0.2 the maximum margin classifier

Margin: the smallest distance from the observations to the separating hyperplane

Optimal (maximum margin) hyperplane: it is the separating hyperplane for which the margin is largest. It is found by solving the optimization problem

$$\arg \max_{\beta_0, \beta} \left\{ \min_t \left[ y(t) \frac{z(t)}{\|\beta\|} \right] \right\} = \arg \max_{\beta_0, \beta} \left\{ \min_t \left[ y(t) \frac{\beta_0 + u^T(t)\beta}{\|\beta\|} \right] \right\}$$

that, in turn, is equivalent to

$$\max_{\beta_0, \beta} D \text{ subject to } y(t) \frac{z(t)}{\|\beta\|} \geq D, \quad t = 1, 2, \dots, N$$

The optimal hyperplane can be seen as the mi-line of the widest "slab" that can be inserted between two classes. Although the identification of the maximum margin requires all the available data, the optimal hyperplane depends directly on only a small subset of the data points, known as *support vectors*.

If we make the rescaling  $\beta \rightarrow k\beta$  and  $\beta_0 \rightarrow k\beta_0$ , the distance from any point  $u(t)$  to the decision boundary is unchanged, therefore, we can arbitrarily set  $\|\beta\| = 1/D$ . As a consequence, the optimization problem is equivalent to

$$\max_{\beta_0, \beta} D \text{ subject to } y(t)z(t) \geq 1, \quad t = 1, 2, \dots, N$$

that is

$$\max_{\beta_0, \beta} \|\beta\|^2 \text{ subject to } y(t)z(t) \geq 1, \quad t = 1, 2, \dots, N$$

This is a convex optimization problem with  $N$  linear inequality constraints. The constraints define an empty slab around the linear decision boundary of thickness  $2D = 2\|\beta\|$

### Constrained optimization problems: the Lagrange multipliers

Consider the following optimization problem with equality constraint:

$$\min_x f(x) \text{ subject to } g(x) = 0$$

where  $x \in \mathbb{R}^p$ . The constraint  $g(x) = 0$  represents a  $(p-1)$ -dimensional hypersurface  $\mathcal{S} \in \mathbb{R}^p$ . We denote the gradients of  $f(x)$  and  $g(x)$  as  $\nabla f(x), \nabla g(x)$

Since

1. at any point  $x \in \mathcal{S}$  the gradient  $\nabla g(x)$  is orthogonal to  $\mathcal{S}$
2. for a point  $x^* \in \mathcal{S}$  that minimizes  $f(x)$ ,  $\nabla f(x^*)$  is orthogonal to  $\mathcal{S}$ <sup>1</sup>

the gradient  $\nabla f(x^*)$  and  $\nabla g(x^*)$  are parallel or anti-parallel. Then, there exists a scalar  $\lambda$  called *Lagrange multiplier* such that

$$\nabla f(x^*) = \lambda \nabla g(x^*)$$

We introduce the Lagrangian function

$$L(x, \lambda) = f(x) - \lambda g(x)$$

---

<sup>1</sup>this is due to the fact that if the gradient were not orthogonal to  $\mathcal{S}$  by moving in the direction given by the opposite of the gradient projected onto  $\mathcal{S}$  we would reach a lower point for  $f(x)$  belonging to  $\mathcal{S}$

The original problem can be reformulated as

$$\min_x L(x, \lambda)$$

in fact:

$$\begin{aligned} \frac{\partial L(x, \lambda)}{\partial x} = 0 &\implies \nabla f(x) - \lambda \nabla g(x) = 0 \\ \frac{\partial L(x, \lambda)}{\partial \lambda} = 0 &\implies g(x) = 0 \end{aligned}$$

The minimization problem with  $N_e$  equality constraints

$$\min_x f(x) \quad \text{subject to} \quad g_j(x) = 0, \quad j = 1, 2, \dots, N_e$$

can be reformulated as

$$\min_x L(x, \lambda)$$

where  $L(x, \lambda)$  is the Lagrangian function

$$L(x, \lambda) = f(x) - \sum_{j=1}^{N_e} \lambda_j g_j(x)$$

and  $\lambda_1, \lambda_2, \dots, \lambda_{N_e}$  are Lagrange multipliers.

Consider the following optimization problem with inequality constraints:

$$\min_x f(x) \quad \text{subject to} \quad g(x) \geq 0$$

where  $x \in \mathbb{R}^p$ . Again the solution will be based on a Lagrange multiplier  $\lambda$  and the associated Lagrangian function. Let  $x^*$  be a solution of the problem. Two possibilities exist:

1.  $g(x^*) > 0$ : In this case the constraint is *inactive* and the solution is found by setting  $\nabla f(x) = 0$  so that  $g(x)$  plays no role  $\implies \lambda = 0$
2.  $g(x^*) = 0$ : In this case the constraint is *active* so that  $\lambda = 0$ . Moreover, the descending direction  $-\nabla f(x^*)$  must be oriented away from the region  $g(x) > 0$  otherwise  $f(x^*)$  does not correspond to a minimum. Consequently, we have  $\nabla f(x^*) = \lambda \nabla g(x^*)$  for some  $\lambda > 0$

Note that both the above cases lead to  $\lambda g(x) = 0$ . We also have that  $\lambda \geq 0$ . The original problem can be reformulated as follows:

$$\min_x L(x, \lambda) = f(x) - \lambda g(x) \quad \text{subject to} \quad \begin{cases} g(x) \geq 0 \\ \lambda \geq 0 \\ \lambda g(x) = 0 \end{cases}$$

The above constraints are known as the Karush-Kuhn-Tucker (KKT) conditions.

- Note that  $L(x, \lambda)$  has to be minimized wrt  $x$  and maximized wrt  $\lambda$  to prevent  $L(x, \lambda) \rightarrow -\infty$

Consider now the more general optimization problem with  $N_i$  inequality constraints

$$\min_x f(x) \quad \text{subject to} \quad g_j(x) \geq 0 \quad j = 1, 2, \dots, N_i$$

Define the Lagrangian function

$$L(x, \lambda) = f(x) - \sum_{j=1}^{N_i} \lambda_j g_j(x)$$

where  $\lambda = [\lambda_1 \quad \lambda_2 \quad \dots \quad \lambda_{N_i}]^T$ . The problem can be reformulated as

$$\min_x L(x, \lambda) = f(x) - \lambda g(x) \quad \text{subject to} \quad \begin{cases} g_j(x) \geq 0 & j = 1, \dots, N_i \\ \lambda_j \geq 0 & j = 1, \dots, N_i \\ \lambda_j g_j(x) = 0 & j = 1, \dots, N_i \end{cases}$$

The number of KKT conditions is thus  $3N_i$ . The Lagrangian function  $L(x, \lambda)$  has to be minimized wrt  $x$  and maximized wrt  $\lambda$

**Finding the maximum margin classifier by using the Lagrange multipliers**

The problem

$$\min_{\beta_0, \beta} \frac{1}{2} \|\beta\|^2 \quad \text{subject to} \quad y(t)z(t) \geq 1, \quad t = 1, 2, \dots, N$$

becomes

$$\min_{\beta_0, \beta} J(\beta_0, \beta, \lambda) = \frac{1}{2} \|\beta\|^2 - \sum_{t=1}^N \lambda_t (y(t)z(t) - 1) = \frac{1}{2} \|\beta\|^2 - \sum_{t=1}^N \lambda_t [y(t)(\beta_0 + u^T(t)\beta) - 1]$$

subject to the constraints

$$\begin{aligned} y(t)z(t) - 1 &\geq 0 \\ \lambda_t &\geq 0 \quad t = 1, 2, \dots, N \\ \lambda_t (y(t)z(t) - 1) &= 0 \end{aligned}$$

- the points for which  $\lambda_t = 0$  do not contribute to the estimated model  $\hat{\theta} = [\hat{\beta}_0 \quad \hat{\beta}^T]^T$ . The remaining points, for which  $y(t)z(t) = 1$ , are the *support vectors* and lie on the edges of the margin. This means that the estimated function  $\hat{z}(t)$  will depend only on a small part of the training data.

Setting to zero the derivatives of  $J$  wrt  $\beta, \beta_0$  we get

$$\frac{\partial J}{\partial \beta} = 0 \implies \beta = \sum_{t=1}^N \lambda_t y(t) u(t) \quad (12.1)$$

$$\frac{\partial J}{\partial \beta_0} = 0 \implies 0 = \sum_{t=1}^N \lambda_t y(t) \quad (12.2)$$

Substituting the above relations in  $J(\beta_0, \beta, \lambda)$  we obtain the *dual objective function*

$$J(\lambda) = \sum_{t=1}^N \lambda_t - \sum_{t=1}^N \sum_{k=1}^N \lambda_t \lambda_k y(t) y(k) u^T(t) u(k)$$

to be maximized wrt  $\lambda$  only subject to

$$\begin{aligned} \lambda_t &\geq 0, \quad t = 1, 2, \dots, N \\ \sum_{t=1}^N \lambda_t y(t) &= 0 \end{aligned}$$

This is a simpler convex optimization problem, involving a quadratic function of  $\lambda$ .

By replacing 12.1 in  $z(t) = \beta_0 + u^T(t)\beta$  we get

$$z(t) = \beta_0 + u^T(t) \sum_{k=1}^N \lambda_k y(k) u(k) = \beta_0 + \sum_{k=1}^N \lambda_k y(k) u^T(t) u(k)$$

Since only the support vectors have  $\lambda_k \neq 0$  we have

$$z(t) = \beta_0 + \sum_{k \in \mathcal{S}_v} \lambda_k y(k) u^T(t) u(k)$$

where  $\mathcal{S}_v$  is the set of indices of the support vectors.

Once the estimate  $\hat{\theta}, \hat{\lambda}$  has been determined, a new input  $u(t)$  is classified as in slide 3 by using

$$\hat{z}(t) = \varphi^T(t) \hat{\theta} = \hat{\beta}_0 + \sum_{k \in \mathcal{S}} \hat{\lambda}_k y(k) u^T(t) u(k)$$

The above relation shows how the obtained solution only depends on a limited number of scalar products  $u^T(t) u(k)$

- Although none of the training inputs are misclassified by construction, this will not necessarily be the case for future unseen data points. Hopefully, a large margin on the training data will lead to a good separation on the validation (test) data
- The distance of a point from the hyperplane can be seen as a measure of the confidence that the observation was correctly classified.
- The maximum margin hyperplane can be extremely sensitive to changes in the available data. Even a change in a single observation could lead to a dramatic change in the optimal hyperplane.
- If the training set is not linearly separable, the optimization problem has no feasible solution

## 12.1 Support vector machine

Support vector (soft margin) classifier: the previous approach is modified so that some of the training points are allowed to be on the "wrong" side of the margin or even on the wrong side of the hyperplane but with a penalty that increases with the distance from the hyperplane. In this way, the non separable case can also be tackled.

To solve this problem we introduce a set of *slack variables*  $\xi_1, \xi_2, \dots, \xi_N$  with  $\xi_t \geq 0, \forall t$ . More precisely, each input  $u(t)$  is associated with a slack variable as follows:

$$\xi_t = \begin{cases} 0 & \text{for data points that are on the correct side of the margin or on the correct edge of the margin} \\ |y(t) - z(t)| & \text{for other points} \end{cases}$$

It follows that

- $\xi_t = 0 \implies u(t)$  is correctly classified with distance  $\geq D$
- $0 < \xi_t < 1 \implies u(t)$  is correctly classified with distance  $< D$
- $\xi_t = 1 \implies u(t)$  is on the decision boundary
- $\xi_t > 1 \implies u(t)$  is misclassified

The optimization problem is

$$\max_{\beta_0, \beta, \xi} D \quad \text{subject to} \quad \begin{cases} y(t) \frac{z(t)}{\|\beta\|} \geq D(1 - \xi) \\ \xi_t \geq 0 \\ \sum_{t=1}^N \xi_t \leq K \end{cases} \quad t = 1, 2, \dots, N$$

where  $\xi = [\xi_1 \ \xi_2 \ \dots \ \xi_N]^T$  and  $K > 0$  is a tuning parameter that controls the trade-off between the slack variable penalty and the margin. It is an upper bound on the number of misclassified points. As  $K$  increases, the margin increases.

The problem can be rewritten as

$$\max_{\beta_0, \beta, \xi} D \quad \text{subject to} \quad \begin{cases} y(t)z(t) \geq (1 - \xi_t) \\ \xi_t \geq 0 \\ \sum_{t=1}^N \xi_t \leq K \end{cases} \quad t = 1, 2, \dots, N$$

by rescaling  $\beta_0$  by a factor  $\eta$  such that  $\|\beta\| = 1/D$ . Computationally, it is convenient to exploit the equivalent form

$$\min_{\beta_0, \beta, \xi} \left\{ \frac{1}{2} \|\beta\|^2 + C \sum_{t=1}^N \xi_t \right\} \quad \text{subject to} \quad \begin{cases} y(t)z(t) \geq (1 - \xi_t) \\ \xi_t \geq 0 \end{cases} \quad t = 1, 2, \dots, N$$

Where the tuning parameter  $C > 0$  plays the role of  $1/K$  in the previous form (as  $C$  increases the margin narrows). The scalar  $1/2$  has been included for later convenience. Once again, we have a convex optimization problem that can be solved by using the method of Lagrange multipliers. By introducing the vectors  $\lambda = [\lambda_1 \ \lambda_2 \ \cdots \ \lambda_N]^T$  and  $\mu = [\mu_1 \ \mu_2 \ \cdots \ \mu_N]^T$ , the above problem can be seen as the optimization of the loss function

$$J(\beta_0, \beta, \xi, \lambda, \mu) = \frac{1}{2} \|\beta\|^2 + C \sum_{t=1}^N \xi_t - \sum_{t=1}^N \lambda_t (y(t)z(t) - 1 + \xi_t) - \sum_{t=1}^N \mu_t \xi_t$$

subject to the constraints (KKT conditinos):

$$\begin{aligned} y(t)z(t) - 1 + \xi_t &\geq 0 \\ \lambda_t &\geq 0 \\ \lambda_t (y(t)z(t) - 1 + \xi_t) &= 0 \\ \xi_t &\geq 0 \\ \mu_t &\geq 0 \\ \mu_t \xi_t &= 0 \end{aligned}$$

As before, the points for which  $\lambda_t = 0$  do not contribute to the estimated model  $\hat{\theta}$ . The remaining points, for which  $y(t)z(t) = 1 - \xi_t$  are the support vectors. They may have  $\xi_t = 0$  (lie on the edges of the margin),  $0 < \xi_t \leq 1$  (lie inside the margin) or  $\xi_t > 1$  (misclassified)

Setting to zero the derivatives of  $J$  wrt  $\beta, \beta_0$  and  $\xi_t (t = 1, \dots, N)$  we get

$$\frac{\partial J}{\partial \beta} = 0 \implies \beta = \sum_{t=1}^N \lambda_t y(t) u(t) \quad (12.3)$$

$$\frac{\partial J}{\partial \beta_0} = 0 \implies 0 = \sum_{t=1}^N \lambda_t y(t) \quad (12.4)$$

$$\frac{\partial J}{\partial \xi_t} = 0 \implies \lambda_t = C - \mu_t, \quad t = 1, \dots, N \quad (12.5)$$

Substituting the above relations in  $J(\beta_0, \beta, \xi, \lambda, \mu)$  we obtain the *dual objective function*

$$J(\lambda) = \sum_{t=1}^N \lambda_t - \frac{1}{2} \sum_{t=1}^N \sum_{k=1}^N \lambda_t \lambda_k y(t) y(k) u^T(t) u(k)$$

to be maximized wrt  $\lambda$  subject to

$$\begin{aligned} 0 &\leq \lambda_t \leq C, \quad t = 1, 2, \dots, N \\ \sum_{t=1}^N \lambda_t y(t) &= 0 \end{aligned}$$

By looking at 12.3, it is clear that the solution of the optimization problem depends on input samples only through the scalar products

$$u^T(t) u(k), \quad t, k = 1, 2, \dots, N$$

By replacing 12.3 in  $z(t) = \beta_0 + u^T(t) \beta$  we get

$$z(t) = \beta_0 + u^T(t) \sum_{k=1}^N \lambda_k y(k) u(k) = \beta_0 + \sum_{k=1}^N \lambda_k y(k) u^T(t) u(k)$$

Since only the support vectors have  $\lambda_k \neq 0$  we have

$$z(t) = \beta_0 + \sum_{k \in \mathcal{S}_v} \lambda_k y(k) u^T(t) u(k)$$

Where  $\mathcal{S}_v$  is the set of indices of the support vectors. Once that  $\hat{\lambda}$  and  $\hat{\beta}_0$  have been computed, a new input  $u(t)$  is classified by using

$$\hat{z}(t) \varphi^T(t) \hat{\theta} = \hat{\beta}_0 + \sum_{k \in \mathcal{S}_v} \hat{\lambda}_k y(k) u^T(t) u(k)$$

Again, all we need is scalar products.



### 12.1.1 Dealing with nonlinear boundaries: the kernel trick

Nonlinear decision boundaries can be tackled by introducing a generalization of the scalar product

$$\Psi(u(t), u(k))$$

where  $\Psi(x, y) = \Psi(y, x)$  is a symmetric function called kernel, that involves nonlinear functions of the input samples. The scalar product  $\Psi(x, y) = x^T y$  is a linear kernel.

By using kernels, we can fit a support vector classifier in a higher-dimensional space involving nonlinear functions rather than in the original input space. The resulting classifier is known as *support vector machine*.

#### How does it work?

1. Choose a kernel function  $\Psi(u(t), u(k))$ .
2. Solve the optimization problem

$$\max_k J(\lambda) = \sum_{t=1}^N \lambda_t - \frac{1}{2} \sum_{t=1}^N \sum_{k=1}^N \lambda_t \lambda_k y(t) y(k) \Psi(u(t), u(k))$$

subject to

$$0 \leq \lambda_t \leq C, \quad t = 1, 2, \dots, N$$

$$\sum_{t=1}^N \lambda_t y(t) = 0$$

and find the estimates  $\hat{\lambda}$ . Then, find the estimate  $\hat{\beta}_0$

3. Classify the new inputs by using

$$\hat{z}(t) \varphi^T(t) \hat{\theta} = \hat{\beta}_0 + \sum_{k \in \mathcal{S}} \hat{\lambda}_k y(k) \Psi(u(t), u(k))$$

Two popular choices for the kernel function  $\Psi$  are:

Polynomial kernels

$$\Psi(u(t), u(k)) = (1 + u^T(t) u(k))^q$$

Radial basis (gaussian) kernels

$$\Psi(u(t), u(k)) = \exp(-\gamma \|u(t) - u(k)\|^2)$$

where  $\gamma > 0$  is a hyperparameter.

By using kernels, we can determine complex nonlinear decision boundaries without explicitly working in the enlarged feature space. The enlarged feature space can even be implicit and infinite-dimensional. This is the case for the radial basis kernels.

### 12.1.2 Multiclass problems

#### One-vs-all approach

Consists in solving  $M$  two-class SVM problems. The generic  $i$ -th problem involves the classes  $C_i$  and  $\bar{C}_i = \{C_1, \dots, C_{i-1}, C_{i+1}, \dots, C_M\}$ . Therefore, we get  $M$  different estimates for  $\lambda$  and  $\beta_0$  and, consequently,  $M$  discriminant functions  $\hat{z}^1(t), \hat{z}^2(t), \dots, \hat{z}^M(t)$ . For any new input  $u(t)$ , the values  $\hat{z}^i(u(t)), i = 1, 2, \dots, M$  are computed in order to take the final decision.

This method does not lead to an optimal solution and inconsistent results could be obtained, that is, an input could be assigned to multiple classes simultaneously. For this reason, the most common approach is to assign a new input  $u(t)$  to the class  $k$  such that

$$k = \arg \max_{i \in \{1, 2, \dots, M\}} \hat{z}^i(t)$$

However, this is a heuristic approach that suffers from the problem that the  $M$  obtained classifiers were trained on different tasks.

**One-vs-one approach**

Consists in solving  $M(M - 1)/2$  two-class SVM problems. The generic binary problem involves the classes  $C_i$  and  $C_k$ , so that only part of the training data are used. For any new input  $u(t)$ , the values  $\hat{z}^i(u(t))$ ,  $i = 1, 2, \dots, M(M - 1)/2$  are computed. The new input is then classified according to the class that got the highest number of "votes".

- Again, the obtained solution is not optimal and the approach can lead to ambiguities in the resulting classification
- For large  $M$ , the required computational effort significantly increases.

# Chapter 13

## Regularization

Consists in adding a penalty term to the loss function

$$V(\theta) = J(\theta) + \lambda g(\theta)$$

where  $\lambda > 0$  is a tuning parameter. The solution of the learning from data problem becomes

$$\hat{\theta} = \arg \min_{\theta \in \mathcal{M}_p(\theta)} V(\theta)$$

This technique is used both in regression and classification problems. One of the most common form is *quadratic regularization*

$$V(\theta) = J(\theta) + \lambda \|\theta\|_2^2 = J(\theta) + \lambda \theta^T \theta$$

Why regularization?:

- To counteract the effects of overfitting
- To make the optimization problem better conditioned

Note: the tuning parameter  $\lambda$  must be determined separately as for the model complexity.

### 13.1 Regularized least squares: the Ridge regression

Consider the static model

$$y(t) = \varphi^T(t)\theta + \varepsilon(t), \quad t = 1, 2, \dots, N$$

and the associated quadratic regularized LS problem

$$\min V(\theta) : \quad v(\theta) = \sum_{t=1}^N \varepsilon^2(t) + \lambda \theta^T \theta = \|Y - \Phi\theta\|^2 + \lambda \|\theta\|_2^2$$

This estimation method is called Ridge regression. By setting the gradient of  $v(\theta)$  to zero it is easy to find the normal equations

$$(\Phi^T \Phi + \lambda I_p) \theta = \Phi^T Y$$

where  $p$  is the complexity of the model. The regularized LS estimator is thus given by

$$\hat{\theta} = (\Phi^T \Phi + \lambda I_p)^{-1} \Phi^T Y$$

the penalty term  $\lambda \|\theta\|_2^2$  is also called a *shrinkage penalty* in statistics because it has the effect of shrinking the estimated parameters towards zero.

Assume now the existence of the true model

$$y(t) = \varphi^T(t)\theta^* + w(t) \quad t = 1, 2, \dots, N$$

By analyzing the statistical properties of the Ridge estimator we get

$$\begin{aligned} E[\hat{\theta}] &= E[(\Phi^T \Phi + \lambda I_p)^{-1} \Phi^T Y] \\ Y &= \Phi \theta^* + w \\ E[\hat{\theta}] &= E[(\Phi^T \Phi + \lambda I_p)^{-1} \Phi^T (\Phi \theta^* + w)] = \\ &= E[(\Phi^T \Phi + \lambda I_p)^{-1} \Phi^T \Phi \theta^*] + E[(\Phi^T \Phi + \lambda I_p)^{-1} \Phi^T w] = \\ &= (\Phi^T \Phi + \lambda I_p)^{-1} \Phi^T \Phi \theta^* \neq \theta^* \end{aligned}$$

and

$$\text{cov}(\hat{\theta}) = \sigma_w^2 (\Phi^T \Phi + \lambda I_p)^{-1} \Phi^T \Phi (\Phi^T \Phi + \lambda I_p)^{-1} < \sigma_w^2 (\Phi^T \Phi)^{-1} = \text{cov}(\hat{\theta}_{LS})$$

- The Ridge estimator  $\hat{\theta}$  is biased and  $\text{cov}(\hat{\theta}) < \text{cov}(\hat{\theta}_{LS})$
- as  $\lambda$  increases, the shrinkage of the estimated coefficients leads to a reduction in the variance of the estimates at the expense of an increase in bias  $\implies \lambda$  is a hyperparameter that can manage the bias-variance tradeoff
- If  $\Phi^T \Phi$  is ill-conditioned (or even rank deficient), the use of  $\lambda$  leads to the better conditioned matrix  $\Phi^T \Phi + \lambda I_p$

Regularization can be exploited also for the identification of dynamic models and for classification problems.

In general:

- If the model complexity  $p$  is high (many parameters) it may not be possible to estimate several of the accurately  $\implies$  it is advantageous to pull them towards zero as the ones having the smallest influence on  $J(\theta)$  will be affected most by the shrinkage property  $\implies$  regularization allows complex models to be trained on small data sets without severe overfitting.
- The problem of minimizing  $J(\theta)$  may be ill conditioned, especially when the complexity  $p$  is high, in the sense that the Hessian  $J''(\theta)$  may be ill-conditioned  $\implies$  adding the norm penalty will add  $\lambda I_p$  to this matrix so that it becomes better conditioned.
- The choice of  $\lambda$  is a crucial issue as we may think of  $\lambda$  as a knob to control the bias-variance tradeoff (the larger  $\lambda$  the larger the number of parameters that will be close to zero). The best way consists in choosing the values of  $\lambda$  leading to the smallest value of the loss function evaluated on the validation set.

### Alternative formulation for Ridge regression

It is possible to show that the problem

$$\arg \min_{\theta \in \mathcal{M}_p(\theta)} V(\theta), \quad v(\theta) = J(\theta) + \lambda \theta^T \theta$$

is equivalent to the problem

$$\arg \min_{\theta \in \mathcal{M}_p(\theta)} J(\theta), \quad \text{subject to } \theta^T \theta \leq K$$

- For every  $\lambda$  there is some  $K$  such that the solutions  $\hat{\theta}$  of the two optimization problems are the same. The two approaches can be related using Lagrange multipliers.
- The parameter  $K$  can be seen as a budget for how large the norm of  $\theta$  can be
- Note that  $K$  plays the role of  $1/\lambda$ . In fact, decreasing  $K$  has the effect of shrinking the estimated parameters towards zero.

Look at prof. notes for geometrical interpretation (The solution will lie on the level curve of  $J(\theta)$  that intercepts the constraint curve)