### Table of Contents I

- Least squares estimators
- mathematical formulation
- basic properties
- linear least squares
  - Most important python code for this sub-module
  - Self-assessment material

 this is the table of contents of this document; each section corresponds to a specific part of the course

- 1



### Contents map

developed content units	taxonomy levels
least squares	u1, e1

prerequisite content units	taxonomy levels
dataset	u1, e1



- Least squares estimators 2

### Main ILO of sub-module <u>"Least squares estimators"</u>

Describe the concept of least squares in geometrical perspectives

Derive and use the normal equations for solving separable least squares problems



- Least squares estimators 3

data generation model: $y_t = f(u_t; \theta) + v_t$ dataset: $\mathcal{D} = \{(u_t, y_t)\}_{t=1,...,N}$ hypothesis space: $\theta \in \Theta$ 

Problem: find a  $\widehat{\theta} \in \Theta$  that "best explains"  $\mathcal{D}$ 

notes \_

- assume that we have selected a model structure, that means that we hypothesize that the data is generated according to a certain data generation model
- importantly, we do not know  $\theta$  and we want to estimate it
- remember: the choice of the model structure is an assumption, i.e., we think that the data is
  generated following the model structure that we chose. We can practically never know if this
  is actually true though. However sometimes the model structure that we choose is a good
  guess, so that even if there is some mismatch we practically are happy with it
- beyond having chosen the model structure, we also have collected a dataset (this is given, typically there is no "user choice" here – exceptions being if we want to pre-process the dataset because we want to remove outliers or filter some noise)
- we finally know that the unknown parameter that we want to estimate lives in a certain set
- what we want to do is to estimate the unknown parameter in the sense (to be clarified soon) that we want to "explain" the dataset

- Least squares estimators 4

### Geometrical interpretation





- Least squares estimators 5

Consider

$$\begin{bmatrix} f(u_1;\theta)\\ \vdots\\ f(u_N;\theta) \end{bmatrix};$$

varying  $u_1, \ldots, u_N$  but keeping  $\theta$  fixed corresponds in general to find:

Potential answers:		
I: (wrong) II: (wrong) III: (correct) IV: (wrong)	a scalar a vector a manifold I do not know	

- Least squares estimators 6





### Intuitions, towards a mathematical formulation



 let's re-see the geometrical intuition we developed before, where the dot is the vector of measurements, the manifold is all the potential predicted vectors given all the potential parameters

- mathematical formulation 2

Mathematical formulation





### Example: regression line



- mathematical formulation 4

- let's see a geometrically inspectable example, i.e., the LS estimate of a regression line
- the dataset is this one
- the data generation model is this one
- and the hypothesis space is this one
- very important! The hypothesis space must always be specified!
- we will indeed see how the topological properties of the hypothesis space have huge influences on the mathematical properties of the various estimators
- then by using the definition of LS and interpretation of residuals, the optimization problem becomes this one
- note how this specific problem can be solved by putting the first derivative of the cost to zero
  and then checking if the Hessian is positive definite
- we will see soon how actually for this type of linear systems there is the possibility of computing the solutions through closed formulas

### ${\small Question} \ 2$

### Consider

$$f(u;\theta) = \sum_{k=0}^{2} \theta_k u^k \qquad \qquad \mathcal{D} = \{(0,0), (1,1)\} \qquad \qquad \Theta = \mathbb{R}^2.$$

How many solutions will the LS problem have?

### **Potential answers:**

l: <b>(wr</b>	ong)	0
ll: <b>(wr</b>	ong)	1
III: ( <u>co</u>	rrect)	$+\infty$
IV: ( <u>wr</u>	ong)	l do not knov





The concepts behind LS are simple, so it is simple to compute analytically  $\widehat{ heta}_{\mathrm{LS}}$ 

Potential answ	vers:
I: (wrong)	true
II: (correct)	false
III: (wrong)	I do not know



### Example: computing the LS may be numerically infeasible

 $u_t \in \mathbb{R}^{10^6}$ 

 $f(u_t; \theta)$  extremely nonlinear  $\mathcal{D} = \{(u_t, y_t)_{t=1,...,N}\}, N = 10^{12}$  $\theta \in \text{ very non-convex set}$ 



- basic properties 3

### Question 4

The LS estimate  $\widehat{ heta}_{\mathrm{LS}}$  always exists

Potential answers:		
true false I do not know		
	<b>vers:</b> true false I do not know	



$$\widehat{\theta}_{\text{LS}} = \arg \min_{\theta \in (0,1)} \sum_{t=1}^{100} r_t^2(\theta) \qquad r_t(\theta) =$$

$$r_t(\theta) = \frac{1}{\theta}$$

- it may also be that there is no estimate, in the same way that some function may do not admit a minimum
- note that this is something that has a lot to do with the openness of the hypothesis space
- as you may remember, a continuous function over a compact always admits minimum and maximum (Weierstrass)
- in this specific case the residuals form a continuous function, but the domain is not compact, and this makes the overall problem have no solution
- this is the first time we see that the topological properties of the hypothesis space play an important role. Things like this one will happen again in the next modules

- basic properties 5

### Question 5

When it exists, the LS estimate  $\widehat{\theta}_{\mathrm{LS}}$  is unique

Potential answers:	
: (wrong)   : ( <u>correct</u> )    : (wrong)	true false I do not know



### Example: the LS estimate is not unique

How many quadratics fit perfectly this dataset?





- basic properties 7



If you don't remember how to do computations with matrices and vectors...

the matrix cookbook

### Separable problems



- linear least squares 2

- linear least squares 3

LS for **unconstrained** separable problems  $\implies$  normal equations

 $\boldsymbol{y} = \Phi(\boldsymbol{u})\boldsymbol{\theta} + \boldsymbol{e}, \quad \boldsymbol{\theta} \in \mathbb{R}^n \qquad \qquad \widehat{\theta}_{\mathrm{LS}} = \arg\min_{\boldsymbol{\theta} \in \mathbb{R}^n} \|\boldsymbol{y} - \Phi(\boldsymbol{u})\boldsymbol{\theta}\|^2$ 

Ideally, 
$$\hat{\theta}_{LS}$$
 is s.t.  $\Phi(\boldsymbol{u})\hat{\theta}_{LS} = \boldsymbol{y}!$   
normal equations:  $\Phi(\boldsymbol{u})^T \Phi(\boldsymbol{u})\hat{\theta}_{LS} = \Phi(\boldsymbol{u})^T \boldsymbol{y}$ 





- we now will see how solving <u>unconstrained</u> separable problems leads to what are called the normal equations
- it is important to notice that we are having an <u>unconstrained</u> problem here, i.e., the hypothesis space is the whole R<sub>θ</sub>
- so let's assume that we have this problem
- ideally, what we would like to happen is that the LS solution perfectly explains the dataset
- this happens only if  $\boldsymbol{y}$  is in the range of  $\boldsymbol{\Phi}$
- in general this does not happen and if it happens it may be that there is more than one solution, as we will see soon
- by solving the optimization problem in a matricial form we end up with these equations here, that are called normal equations
- Discuss the pros and cons of using direct methods in comparison to iterative methods

### Exercise

Compute the solution of

$$\arg\min_{\theta\in\mathbb{R}^n} \left( \boldsymbol{y} - \Phi(\boldsymbol{u})\theta \right)^T W \left( \boldsymbol{y} - \Phi(\boldsymbol{u})\theta \right)^T$$

- linear least squares 5

### Question 6

### Starting from

 $\Phi(\boldsymbol{u})^T \Phi(\boldsymbol{u}) \widehat{\theta}_{\mathrm{LS}} = \Phi(\boldsymbol{u})^T \boldsymbol{y}$ 

we can always set

$$\widehat{\theta}_{\text{LS}} = \left( \Phi(\boldsymbol{u})^T \Phi(\boldsymbol{u}) \right)^{-1} \Phi(\boldsymbol{u})^T \boldsymbol{y}$$

### Potential answers:

I: (wrong) true

- II: (correct) false III: (wrong) I do not know
  - g) I do not know



### Using the pseudoinverse when necessary

what if  $\Phi(\boldsymbol{u})^T \Phi(\boldsymbol{u})$  does not have an inverse?

### Definition (Moore-Penrose pseudoinverse of a matrix)

Given  $A \in \mathbb{R}^{m \times n}$ ,  $A^{\dagger}$  is its pseudoinverse if

 $AA^{\dagger}A = A$  $A^{\dagger}AA^{\dagger} = A^{\dagger}$  $(AA^{\dagger})^{H} = AA^{\dagger}$  $(A^{\dagger}A)^{H} = A^{\dagger}A$ 

more in http://www.math.ucla.edu/~laub/33a.2.12s/mppseudoinverse.pdf - linear least squares 7

### as hinted in the previous exercise, that inverse may not exist – then we can use the Moore-Penrose pseudoinverse

- this is defined in this way
- and there are a lot of things that may be said about it
- unfortunately we cannot discuss about it too much, so: for the interested student, please check here

Using the pseudoinverse when necessary

strong connections with singular values decompositions!



We can always solve the normal equations for every unconstrained separable LS problem

Potential answ	ers:
I: ( <u>correct</u> )	true
II: (wrong)	false
III: (wrong)	I do not know



- linear least squares 9

### Question 8

We can always solve the normal equations for every separable LS problem, even for constrained ones  $% \left( {{{\rm{C}}_{\rm{B}}}} \right)$ 

Potential answ	/ers:
I: (wrong)	true
II: ( <u>correct</u> )	false
III: (wrong)	I do not know



- linear least squares 10

### LS for constrained separable problems $\implies$ normal equations

 $\boldsymbol{y} = \Phi(\boldsymbol{u})\boldsymbol{\theta} + \boldsymbol{e}, \quad \boldsymbol{\theta} \in \Theta \qquad \qquad \widehat{\theta}_{\mathrm{LS}} = \arg\min_{\boldsymbol{\theta} \in \Theta} \|\boldsymbol{y} - \Phi(\boldsymbol{u})\boldsymbol{\theta}\|^2$ 

Ideally, looking for  $\theta^*$  s.t.  $\Phi(u)\theta^* - y = 0$ , but it may happen that  $\theta^* \notin \Theta$ !

### Example = fitting a convex quadratic here:



- linear least squares 11

### note that the solution through the normal equations is something that is guaranteed to be the actual solution only whenever we have unconstrained problems! the intuition is this one: through the normal equations, one gets a point. But that point may not be in the hypothesis space if that happens, one has to solve a constrained optimization problem thus for constrained separable problems one may do as follows: compute a potential solution through the normal equations check whether this potential solution is within Θ if so, then ok, the problem is solved if not, one has to solve the constrained optimization problem

Summarizing

Describe the concept of least squares in geometrical perspectives

Derive and use the normal equations for solving separable least squares problems

- visualize the dataset in an opportune multidimensional plot
- if we have separability, we can use linear algebra to arrive at  $X^T X \hat{\theta} = Xy$



Most important python code for this sub-module

notes

- linear least squares 1

### Illustrative example

### https:

//scikit-learn.org/stable/auto\_examples/linear\_model/plot\_ols.html



- linear least squares 2

Self-assessment material

- linear least squares 1

### Question 9

In the geometric interpretation of least squares, what does the vector  $\boldsymbol{y}$  represent?

Potential answers:	
I: (wrong)	The model parameters to be estimated
II: (correct)	The fixed vector of measured output values
III: (wrong)	The manifold of all possible model predictions
IV: (wrong)	The noise affecting the measurements
V: (wrong)	I do not know

### Solution 1:

In the geometric interpretation, y is the fixed vector of measured output values from the dataset. The least squares problem aims to find the point on the model manifold (determined by  $\Phi(u)\theta$ ) that is closest to y. notes
 see the associated solution(s), if compiled with that ones :)

What is the fundamental assumption required to derive the normal equations for least squares?

Potential answers:	
I: (wrong)	The noise must be Gaussian distributed
II: (wrong)	The model must be nonlinear in parameters
III: (correct)	The problem must be linear in parameters (separable)
IV: (wrong)	The hypothesis space must be constrained
V: (wrong)	l do not know

### Solution 1:

The normal equations  $\Phi^T \Phi \theta = \Phi^T \mathbf{y}$  can only be derived for problems that are linear in their parameters (separable problems). This allows the analytical solution reast squares 3 through matrix operations.



### Question 11

When is the Moore-Penrose pseudoinverse required in least squares problems?

### Solution 1:

The pseudoinverse is needed when  $\Phi^T \Phi$  is singular (not invertible), which occurs when the columns of  $\Phi$  are linearly dependent. It provides a generalized inverse that gives the minimum-norm solution.



What guarantees the existence of a unique least squares solution?

### **Potential answers:**

I:	(wrong)	Having more parameters than measurements
II:	( <u>correct</u> )	$\Phi$ having full column rank and unconstrained parameters
III:	(wrong)	The hypothesis space being compact
IV:	(wrong)	The noise being normally distributed
V:	(wrong)	l do not know

### Solution 1:

A unique solution exists when  $\Phi$  has full column rank (making  $\Phi^T \Phi$  invertible) and the parameters are unconstrained. This ensures the normal equations have exactly one solution.

## see the associated solution(s), if compiled with that ones :)

### Question 13

What is a key difference between constrained and unconstrained least squares problems?

Potential answers:			
l: (wrong) ll: (correct)	Constrained problems always have unique solutions The normal equations may give solutions outside the constraint		
set			
III: (wrong)	Only unconstrained problems can use the pseudoinverse		
IV: (wrong)	Constrained problems require nonlinear optimization		
V: (wrong)	l do not know		

### Solution 1:

For constrained problems, the solution from normal equations may violate the reast squares 6 constraints, requiring additional optimization techniques. Unconstrained problems can be solved directly via normal equations when  $\Phi^T \Phi$  is invertible.



### Recap of sub-module "linear least squares"

- Least squares aims to minimize the squared residuals between model predictions and observed data
- The geometric interpretation views system identification as finding the closest point on a model manifold to measurement vectors
- Normal equations provide an analytical solution for unconstrained linear least squares problems through  $\Phi^T \Phi \theta = \Phi^T y$
- The pseudoinverse generalizes solutions for rank-deficient systems and connects with singular value decomposition
- Existence and uniqueness of LS solutions depend on hypothesis space topology and model structure identifiability
- Constrained LS problems require different approaches than normal equations when parameters must satisfy domain restrictions

- linear least squares 7

# • the most important remarks from this sub-module are these ones