

System Identification and Data Analysis

Federico A. Ramponi

Università di Brescia, dipartimento di ingegneria dell'informazione

Written test
September 11th, 2013
Duration: 2 hours.



Question 1 [5%] Suppose that the variables $y_i \in \mathbb{R}$ are explained by the regressors $\varphi_i \in \mathbb{R}^p$ according to the following model:

$$y_i = \varphi_i^\top \theta^o + \varepsilon_i, \quad i = 1, \dots, N.$$

Write the normal equations of the least squares method to estimate the parameter $\theta^o \in \mathbb{R}^p$. Discuss (max. 5 lines) the uniqueness of a solution to the normal equations; when the solution $\hat{\theta}_{\text{LS}}$ is unique, provide an explicit expression for it.

Question 2 [15%] Suppose that the variables $y_i \in \mathbb{R}$ are explained by the regressors $\varphi_i \in \mathbb{R}^p$ according to the model $y_i = \varphi_i^\top \theta^o + \varepsilon_i$. Suppose, moreover, that

1. $\{\varphi_i\}_{i=1}^\infty$ are deterministic vectors chosen in such a way that

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N \varphi_i \varphi_i^\top = \Sigma,$$

where Σ is symmetric and *positive definite* ($\Sigma > 0$), hence *invertible*.

2. $\{\varepsilon_i\}_{i=1}^\infty$ are i.i.d. random variables, with mean $E[\varepsilon_i] = 0$.

Show that $\hat{\theta}_{\text{LS}} \rightarrow \theta^o$ almost surely as $N \rightarrow \infty$. How is this property called?

Only a sketch of the proof is requested. Here are some hints:

- If a sequence of matrices $\{\Sigma_i\}$ converges to a matrix Σ , and Σ is invertible, then for large N the matrices $\{\Sigma_i\}$ are invertible as well.
- Strong law of large numbers: let $\{x_i\}_{i=1}^\infty$ be i.i.d. random variables, or vectors, with mean $E[x_i] = \mu$. Then $\frac{1}{N} \sum_{i=1}^N x_i \rightarrow \mu$ almost surely.



Question 3 [20%] Explain briefly (max. $\frac{3}{4}$ of a page) what is the method of instrumental variables, and why we introduced it in class.

Hint. We started from this example: consider the dynamical model

$$y(t) = a^o y(t-1) + b^o u(t-1), \quad y_m(t) = y(t) + e(t),$$

where $a^o, b^o \in \mathbb{R}$ are constants, $\{u(t)\}$ is an (“exogenous”) input process, $\{e(t)\}$ are i.i.d. random variables with zero mean, independent from the past of $y(\cdot)$ and from $u(\cdot)$. The experimenter wants to estimate a^o and b^o ; he/she can impose $u(t)$ and measure $y_m(t)$, but cannot measure $y(t)$. Substituting, we get

$$\begin{aligned} y_m(t) - e(t) &= a^o(y_m(t-1) - e(t-1)) + b^o u(t-1), \\ y_m(t) &= a^o y_m(t-1) + b^o u(t-1) + (e(t) - a^o e(t-1)), \\ y_m(t) &= a^o y_m(t-1) + b^o u(t-1) + r(t), \end{aligned}$$

which has the same structure of an ARX model. Hence one may think to apply the least squares method with $y_t = y(t)$, $\varphi_t = \begin{bmatrix} y_m(t-1) \\ u(t-1) \end{bmatrix}$, $\theta = \begin{bmatrix} a \\ b \end{bmatrix}$, and $\varepsilon_t = r(t)$ in order to get an estimate $\hat{\theta}_{\text{LS}}$ of $\theta^o = \begin{bmatrix} a^o \\ b^o \end{bmatrix}$. But this method does *not* work ($\hat{\theta}_{\text{LS}} \not\rightarrow \theta^o$), because...

OK

Question 4 [20%] We have three measures of the variables u, y as follows:

u_i	1	2	3
y_i	3.2	5.8	9.1

Suppose that we want to explain the measures with a linear model $y_i = \theta u_i + \varepsilon_i$ comprising noise terms ε_i , which are independent, continuous random variables with density symmetric around 0. Use the LSCR method to provide a confidence interval for θ^o .

Hint:

	1	2	3
I_1	•	•	○
I_2	•	○	•
I_3	○	•	•
I_4	○	○	○

OK

Question 5 [15%] Consider a binary classification problem in which 20 classifier functions (threshold, interval, ... whatever) are assigned:

$$\hat{f}_c : \mathbb{R} \rightarrow \{0, 1\} \quad c = 1, \dots, 20.$$

Suppose that $(U_1, Y_1), \dots, (U_N, Y_N)$ are i.i.d. input-output pairs, where $U_i \in \mathbb{R}$ and $Y_i \in \{0, 1\}$. Recall that the *true* error function is $\bar{J}(c) = P[\hat{f}_c(U_i) \neq Y_i]$, and that the *empirical* error function is $\hat{J}_N(c) = \frac{1}{N} \sum_{i=1}^N \mathbb{1}_{\{\hat{f}_c(U_i) \neq Y_i\}}$. Using Hoeffding's inequality,

$$P[|S_N - E[S_N]| > \varepsilon] \leq 2e^{-2N\varepsilon^2},$$

compute how many (N) measures are needed to assert, with confidence $1 - 10^{-4}$, that the empirical error approximates the true one with accuracy $\varepsilon = 5\%$, for all the classifiers ($\forall c$) simultaneously.

OK

Question 6 [5%] Suppose that x_1, \dots, x_{10} are independent and identically distributed random variables, distributed according to an unknown density. Consider the random interval

$$I = \left[\min_{i=1, \dots, 10} x_i, \max_{i=1, \dots, 10} x_i \right].$$

What is the probability that a “future” independent observation x_{11} , distributed according to the same density, belongs to I ?

OK

Question 7 [20%] Consider a convex problem of this form:

$$\begin{aligned} & \text{minimize } \gamma \\ & \text{subject to } |\theta - x_1| \leq \gamma \\ & \quad \vdots \\ & \quad |\theta - x_{10}| \leq \gamma \\ & \quad \theta \in \mathbb{R}, \gamma \in \mathbb{R}, \end{aligned}$$

where x_i are random observations as in Question 6. Suppose that its solution is (θ^*, γ^*) . Write down the definition of a *support constraint*. Discuss briefly (max. $\frac{1}{2}$ of a page) how this problem is linked to the interval I of Question 6, and to the probability requested therein.

Hint: make sure that you understand *clearly* what is the (epi-)graph of the function $\gamma = g_i(\theta) = |\theta - x_i|$. Now draw an instance of the optimization problem on the (θ, γ) -plane.

System Identification and Data Analysis

Federico A. Ramponi

Università di Brescia, dipartimento di ingegneria dell'informazione

Written test

September 3rd, 2013

Duration: 2 hours.

Question 1 [5%] Suppose that the variables $y_i \in \mathbb{R}$ are explained by the regressors $\varphi_i \in \mathbb{R}^p$ according to the following model:

$$y_i = \varphi_i^\top \theta^o + \varepsilon_i, \quad i = 1, \dots, N.$$

1. Write the normal equations of the least squares method to estimate the parameter $\theta^o \in \mathbb{R}^p$.
2. Discuss (max. 5 lines) the uniqueness of a solution to the normal equations; when the solution $\hat{\theta}_{LS}$ is unique, provide an explicit expression for it.

Question 2 [15%] Suppose that the variables $y_i \in \mathbb{R}$ are explained by the regressors $\varphi_i \in \mathbb{R}^p$ according to the model $y_i = \varphi_i^\top \theta^o + \varepsilon_i$. Suppose, moreover, that

1. $\{\varphi_i\}_{i=1}^\infty$ are i.i.d. random vectors, with a covariance matrix $\Sigma = \mathbb{E}[\varphi_i \varphi_i^\top]$ which is *positive definite*, hence *invertible*;
2. $\{\varepsilon_i\}_{i=1}^\infty$ are i.i.d. random variables, with mean $\mathbb{E}[\varepsilon_i] = 0$;
3. ε_i is independent from φ_i for all i (it follows that $\mathbb{E}[\varphi_i \varepsilon_i] = \mathbb{E}[\varphi_i] \mathbb{E}[\varepsilon_i]$).

Show that $\hat{\theta}_{LS} \rightarrow \theta^o$ almost surely as $N \rightarrow \infty$. How is this property called?

Only a sketch of the proof is requested. Here are some hints:

- if a sequence of matrices $\{\Sigma_i\}$ converges to a matrix Σ , and Σ is invertible, then for large N the matrices $\{\Sigma_i\}$ are invertible as well;
- Strong law of large numbers: Let $\{x_i\}_{i=1}^\infty$ be i.i.d. random variables (or vectors) with mean $\mathbb{E}[x_i] = \mu$. Then $\frac{1}{N} \sum_{i=1}^N x_i \rightarrow \mu$ almost surely. 10/10

Question 3 [15%] Consider the following dynamical model:

$$y(t) = ay(t-1) + bu(t-1) + e(t),$$

where $a, b \in \mathbb{R}$ are constants, $\{u(t)\}$ is an ("exogenous") input process, $\{e(t)\}$ are i.i.d. random variables with zero mean, and $e(t)$ is independent from $y(t-1), y(t-2), \dots$

1. Show how to apply the method of least squares in order to estimate the constants a and b , as seen in class.
2. Suppose now that $u(t) = ky(t)$ for all t (output feedback). Show that the solution to the corresponding least squares problem cannot be unique. Draw some conclusion about the identifiability of a and b .

Question 4 [15%] We have three measures of the variables u, y as follows:

u_i	1	2	3
y_i	1.2	3.9	6.1

15%

Suppose that we want to explain the measures with a linear model comprising a noise term:

$$y_i = \theta u_i + \varepsilon_i,$$

where ε_i are independent, continuous random variables with density symmetric around 0. Use the LSCR method to provide a confidence interval for θ^* .

Hint:

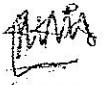
	1	2	3
I_1	•	•	○
I_2	•	○	•
I_3	○	•	•
I_4	○	○	○

Question 5 [5%] Suppose that X_1, \dots, X_N, \dots is a sequence of i.i.d. random variables with common distribution $F(x) = P[X_i \leq x]$. Write down the definition of their *empirical distribution* $\hat{F}_N(x)$ and show, using the strong law of large numbers, that for any fixed $x \in \mathbb{R}$, $\hat{F}_N(x) \rightarrow F(x)$ almost surely.

15%

Question 6 [15%] Suppose that N students compile a questionnaire and reply to a question that allows for 6 possible answers. Assume that the respective answers are i.i.d. random variables (taking 6 possible values). We estimate the probability p_i of a certain answer A_i as

$$\hat{p}_i = \frac{\text{number of answers equal to } A_i}{N}.$$



Using Hoeffding's inequality,

$$P[|S_N - E[S_N]| > \varepsilon] \leq 2e^{-2N\varepsilon^2},$$

compute how many questionnaires are needed to claim, with confidence 99%, that the true "mass distribution" $\{p_1, \dots, p_6\}$ has been estimated with an uniform accuracy of 10%.



Question 7 [5%] Consider a convex problem of this form:

$$\begin{aligned} & \text{minimize } \gamma \\ & \text{subject to } g_1(\theta) - \gamma \leq 0 \\ & \quad \vdots \\ & \quad g_n(\theta) - \gamma \leq 0, \\ & \quad \theta \in \mathbb{R}^d, \gamma \in \mathbb{R}. \end{aligned}$$

15%

Suppose that its solution is (θ^*, γ^*) . Write down the definition of a *support constraint*.

15%

Question 8 [15%] Describe briefly (max. $\frac{3}{4}$ of a page) what is the purpose of the IPM (Interval Predictor Models) method, and how it works. You may provide an example of your choice. If you decide to state some results, no proofs are requested.

System Identification and Data Analysis, Spring 2013

Federico A. Ramponi

Università di Brescia, dipartimento di ingegneria dell'informazione

Written test
 July 8th, 2013
 Duration: 2 hours.

Question 1 [5%] Suppose that the variables $y_i \in \mathbb{R}$ are explained by the regressors $\varphi_i = \varphi(x_i) \in \mathbb{R}^p$ according to the following model:

$$y_i = \varphi_i^\top \theta^o + \varepsilon_i, \quad i = 1, \dots, N. \quad (1)$$

Write down the normal equations of the least squares method to estimate the parameter $\theta^o \in \mathbb{R}^p$. Assuming some more hypotheses if necessary, provide an explicit expression for the estimate $\hat{\theta}_{\text{LS}}$.

Question 2 [15%] Assume that (1) is the “true” model according to which the $\{y_i\}$ are generated. Suppose that the regressors $\{\varphi_i\}$ are deterministic, and that the errors $\{\varepsilon_i\}$ are i.i.d. random variables with mean 0 and variance σ^2 . Under the further hypotheses that you have assumed in Question 1, show that $E[\hat{\theta}_{\text{LS}}] = \theta^o$ (how is this property of $\hat{\theta}_{\text{LS}}$ called?) and that $\text{Var}[\hat{\theta}_{\text{LS}}] = \sigma^2 (\sum_{i=1}^n \varphi_i \varphi_i^\top)^{-1}$.

Question 3 [15%] We have four measures of the variables x, y as follows:

x_i	1	2	3	4
y_i	2.1	0.9	0.7	0.5

Say that we have good reasons to suspect an inverse proportionality between them, as follows:

$$y_i = \frac{\theta}{x_i};$$

on the other hand, we also know that the measures y_i are corrupted by noise. Compute the least squares estimate of the parameter $\theta \in \mathbb{R}$.

Question 4 [15%] Suppose that the variables y_1, y_2, y_3 , are generated by the model

$$y_i = \theta^o u_i + \varepsilon_i,$$

where $u_i \in \mathbb{R}$, $u_i > 0$, and $\theta^o \in \mathbb{R}$ is a parameter. Explain briefly what is a *confidence interval* for θ^o . Make some assumptions on $\varepsilon_1, \varepsilon_2, \varepsilon_3$, and show how to obtain a confidence interval for θ^o (having *guaranteed* confidence) with the LSCR method.

Hint:

	1	2	3
I_1	•	•	○
I_2	•	○	•
I_3	○	•	•
I_4	○	○	○

Question 5 [20%] Describe *briefly* (do not exceed $\frac{3}{4}$ of a page) what are, in a binary classification problem in which the classifier is indexed by the parameter c , the “true” error $\bar{J}(c)$, the empirical error $\hat{J}_N(c)$, and what is the “hope”, in the machine learning theory seen in class, when we minimize the latter with respect to c (or vice-versa, what can go wrong if the family of classifiers indexed by c is too complex).

Question 6 [5%] Write down the statement (*not* the proof) of the Glivenko/Cantelli theorem. (Provide a definition of the functions involved in the statement.)

→ **Question 7 [5%]** Consider a convex problem of this form:

$$\begin{aligned} & \text{minimize } \gamma \\ & \text{subject to } g_1(\theta) - \gamma \leq 0 \\ & \quad \vdots \\ & \quad g_i(\theta) - \gamma \leq 0, \\ & \quad \vdots \\ & \quad g_n(\theta) - \gamma \leq 0, \\ & \theta \in \mathbb{R}^d, \gamma \in \mathbb{R}. \end{aligned}$$

Suppose that its solution is (θ^*, γ^*) . Write down the definition of a *support constraint*.

→ **Question 8 [15%]** Consider now a convex problem of this form:

$$\begin{aligned} & \text{minimize } \gamma \\ & \text{subject to } |\theta - y_1| - \gamma \leq 0 \\ & \quad \vdots \\ & \quad |\theta - y_n| - \gamma \leq 0, \\ & \theta \in \mathbb{R}, \gamma \in \mathbb{R}, \end{aligned}$$

where $y_1, \dots, y_n \in \mathbb{R}$. Sketch a picture of such a problem in the (θ, γ) -plane and mark its solution on the picture. Say if, *in general*, it is possible that: A) the problem has no support constraints at all; B) the problem has 1 support constraint; C) the problem has 2 support constraints; D) the problem has 3 support constraints. (If you believe that a situation is actually possible, sketch an example; otherwise justify your belief.)

→ **Question 9 [5%]** Consider the same problem as in Question 8, and suppose that y_1, \dots, y_n are i.i.d. continuous random variables with uniform density in the interval $[0, 1]$. (Thus, the constraints are now random themselves.) What can you deduce, from the probabilistic standpoint, about the number of support constraints?

System Identification and Data Analysis

Federico A. Ramponi
 Università di Brescia, dipartimento di ingegneria dell'informazione

Written test
 July 29th, 2013
 Duration: 2 hours.

✓**Question 1 [15%]** Let $y_i \in \mathbb{R}$ and $\varphi_i \in \mathbb{R}^p$ for $i = 1, \dots, N$, and let $\theta \in \mathbb{R}^p$.

- Show that any solution to the LS problem, that is to find the $\arg \min_{\theta \in \mathbb{R}^p} \sum_{i=1}^N (y_i - \varphi_i^\top \theta)^2$, must satisfy the normal equations:

$$\left(\sum_{i=1}^N \varphi_i \varphi_i^\top \right) \theta = \sum_{i=1}^N \varphi_i y_i. \quad (1)$$

In other words, show how the normal equations were obtained in class. Hint: *pay attention to the compatibility of multiplications!* φ_i and θ are supposed to be *column vectors*, hence expressions like " $\theta \varphi_i$ " do not make sense.

- Discuss (max. 5 lines) the uniqueness of a solution of (1).

✓**Question 2 [15%]** We have five measures of the variables x, y as follows:

x_i	1	2	3	4	5
y_i	1.2	3.9	9.1	15.6	25.3

We have good reasons to suspect a quadratic relation between them, as follows:

$$y_i = \theta x_i^2;$$

on the other hand, we also know that the measures y_i are corrupted by noise. Compute the least squares estimate of the parameter $\theta \in \mathbb{R}$.

✓ **Question 3 [15%]** Consider the following dynamical model:

$$y(t) = ay(t-1) + bu(t-1) + e(t),$$

where $a, b \in \mathbb{R}$ are constants, $\{u(t)\}$ is an ("exogenous") input process, $\{e(t)\}$ are i.i.d. random variables with zero mean, and $e(t)$ is independent from $y(t-1), y(t-2), \dots$. Show how to apply the method of least squares in order to estimate the constants a and b . Describe a pathological situation in which it is *not possible* to estimate a and b (because the least squares solution is intrinsically not unique).

✓**Question 4 [15%]** Describe *briefly* (max. $\frac{3}{4}$ of a page) what the purpose of the LSCR method is, the hypotheses under which it is usually applied, and how it works. You may provide an example of your choice. If you decide to state some results, no proofs are requested.

Hint:

	1	2	3
I_1	•	•	○
I_2	•	○	•
I_3	○	•	•
I_4	○	○	○

✓ **Question 5 [5%]** Suppose that X_1, \dots, X_N, \dots is a sequence of i.i.d. random variables with common distribution $F(x) = P[X_i \leq x]$. Write down the definition of their *empirical distribution* $\hat{F}_N(x)$ and show, using the strong law of large numbers, that *for any fixed* $x \in \mathbb{R}$, $\hat{F}_N(x) \rightarrow F(x)$ *almost surely.*

✓ **Question 6 [15%]** Consider a binary classification problem in which 20 classifier functions (threshold, interval, ... whatever) are assigned:

$$\hat{f}_c : \mathbb{R} \rightarrow \{0, 1\} \quad c = 1, \dots, 20.$$

Suppose that $(U_1, Y_1), \dots, (U_N, Y_N)$ are i.i.d. input-output pairs, where $U_i \in \mathbb{R}$ and $Y_i \in \{0, 1\}$. Recall that the *true* error function is $\bar{J}(c) = P[\hat{f}_c(U_i) \neq Y_i]$, and that the *empirical* error function is $\hat{J}_N(c) = \frac{1}{N} \sum_{i=1}^N \mathbb{1}_{\{\hat{f}_c(U_i) \neq Y_i\}}$. Using Hoeffding's inequality,

$$P[|S_N - E[S_N]| > \varepsilon] \leq 2e^{-2N\varepsilon^2},$$

compute how many (N) measures are needed to assert, with confidence $1 - 10^{-4}$, that the empirical error approximates the true one with accuracy $\varepsilon = 5\%$, for all the classifiers ($\forall c$) simultaneously.

✓ **Question 7 [5%]** Consider a convex problem of this form:

$$\begin{aligned} &\text{minimize } \gamma \\ &\text{subject to } g_1(\theta) - \gamma \leq 0 \\ &\quad \vdots \\ &\quad g_n(\theta) - \gamma \leq 0, \\ &\quad \theta \in \mathbb{R}^d, \gamma \in \mathbb{R}. \end{aligned}$$

Suppose that its solution is (θ^*, γ^*) . Write down the definition of a *support constraint*.

✓ **Question 8 [15%]** Recall that an optimization problem is called *feasible* if there exists at least one point (θ, γ) that satisfies all the constraints. Consider now a 2-dimensional convex problem with four linear constraints:

$$\begin{aligned} &\text{minimize } \gamma \\ &\text{subject to } a_1\theta + b_1 \leq \gamma \\ &\quad \vdots \\ &\quad a_4\theta + b_4 \leq \gamma \\ &\quad \theta \in \mathbb{R}, \gamma \in \mathbb{R}, \end{aligned}$$

where $a_1, b_1, \dots, a_4, b_4 \in \mathbb{R}$. Say if, *in general*, it is possible that:

- A) the problem is not feasible;
- B) the problem is feasible, but does not admit a solution (θ^*, γ^*) ;
- C) the problem admits a solution, but such solution is not unique;
- D) the problem has a unique solution and no support constraints at all;
- E) the problem has a unique solution and 1 support constraint;
- F) the problem has a unique solution and 2 support constraints;
- G) the problem has a unique solution and 3 support constraints.

If you believe that a situation is actually possible, sketch an example; otherwise justify your belief.

Exercise 1

1. Consider the process $\{y(t)\}_{t=0}^{\infty}$ generated by the following dynamic model:

$$y(t+1) = a_0 y(t) + a_1 y(t-1) + w(t)$$

where $\{w(t)\}_{t=0}^{\infty}$ are independent and identically distributed random variables, with zero mean and with a certain unknown variance $\sigma^2 > 0$, and $w(t)$ is independent from $y(s)$ for each $s \leq t$. Given $N+2$ measures $y(0), \dots, y(N+1)$ of the process, write the normal equations to estimate the parameters a_0 and a_1 with the least squares method.

2. Explain briefly what does it mean that, under appropriate hypothesis, the least squares method is consistent; say if the consistency is verified in the previous point, justifying briefly the answer.

EXERCISE 2

Show briefly (using maximum half of a page) what are, in a binary classification problem where the classifier is indexed by the parameter c , the "true" error $\bar{J}(c)$, the empirical error $\hat{J}_N(c)$ and which problems can arise minimizing the second with the idea of approximating the minimum of the first.

EXERCISE 3

A big company receives N complaint phone calls t_1, \dots, t_N and for each call t_i , the company records the origin region r_i . Suppose that r_i are independent and identically distributed random variables; each of them can assume values in the set of the 20 regions (Piemonte, Lombardia, ..., Sicilia) with the respective probabilities $P = (p(\text{Piemonte}), p(\text{Lombardia}), \dots, p(\text{Sicilia}))$ (that depend on the population of the regions, quality of service, etc.). Using Hoeffding inequality:

$$P[|S_N - E[S_N]| \geq \epsilon] \leq 2e^{-2N\epsilon^2}$$

compute how many phone calls should be collected in order to estimate the probability "distribution" P in order to be sure, with confidence at least $1 - 10^{-10}$, that for all the regions simultaneously the probability estimation error is at most $\epsilon = 1\%$.

EXERCISE 4

1. Consider a convex optimization problem like the following:

$$\min_{x,y} y \text{ such that } g_1(x) \leq y, \dots, g_n(x) \leq y, x \in \mathbb{R}, y \in \mathbb{R}$$

and (\bar{x}, \bar{y}) is its solution. Explain briefly what does it mean that a certain constraint $g_1(x) \leq y$ is of support for the solution.

2. Consider an optimization problem of the same kind, with three constraints delimited by parables facing upwards:

$$\begin{aligned} \min_{x,y} y \text{ such that } & a_1 x^2 + b_1 x + c_1 \leq y \\ & a_2 x^2 + b_2 x + c_2 \leq y \\ & a_3 x^2 + b_3 x + c_3 \leq y \\ & x \in \mathbb{R}, y \in \mathbb{R} \end{aligned}$$

where $a_i > 0$ for $i=1,2,3$. Consider the following situations:

- no constraints are of support
- only one support constraint is present
- two support constraints are present
- all the constraints are of support

in each of the upper cases say if the situation is possible and if it is so, draw the constraints graph of an instance of the problem in which the situation exists (is not required the millimetric precision, a sketch is enough).

EXERCISE OPTIONAL

Note: this exercise will be corrected in order to improve the mark if and only if all the previous exercises have been solved; it's suggested so not to spend time for this before having solved the others.

Below are reported the lengths of some trades of the Milano–Venezia railway and the ticket price of a "fast regional" train that covers them, approximated at the nearest multiple of 1€.

Trade	Length	Ticket Price
Milano–Brescia	83 Km	7 €
Milano–Verona	148 Km	11 €
Milano–Padova	229 Km	15 €

Give a reasonable model that explains the price according to the length of the trade and of a fixed price of the ticket; estimate the price of a ticket from Milano to Vicenza, knowing that the relative trade length is 199 Km.

SIDA

Exam of 28-06-2011

EXERCISE 1

We have some observations $(u_i, y_i), i=1, \dots, N$, where $u_i \in \mathbb{R}, y_i \in \mathbb{R}$.

In order to describe the phenomena that relates u to y , it's used a model with the structure $\hat{y} = \theta u^2$ and the parameter θ is estimated at the least squares.

1. Write the normal equations that allow to determine θ .
2. Given the three observations $(1,1), (1.5,1), (1.8,1.5)$ determine the parameter θ_{LS} that solves the normal equations. Draw the model overlapped to the data in the (u, y) plane.
3. Suppose that the data are generated with the mechanism $y = 0.6u^2 + n$ where u is uniformly distributed in $[0,2]$ and n is Gaussian with zero average and unitary variance. The different data are generated with inputs u independents one from each other and the noise n is independent from u . Say what the least square estimation tends to when the number of observations N tends to infinite.

EXERCISE 2

1. Describe briefly (using maximum 1/3 of a page) what the LSCR method is.
2. Consider again the 3 observations at the point 2 of the previous exercise. Determine a region Θ for the parameter θ of the model $\hat{y} = \theta u^2$ using the LSCR method.

EXERCISE 3

1. Hoeffding inequality is written: $\hat{P}(A) - P(A) \leq \epsilon$, with probability $1 - \beta$ where $\beta = e^{-2N\epsilon^2}$, where $\hat{P}(A)$ is the sample estimate of the probability of an event A . Explain the the meaning of ϵ and β (using maximum 1/4 of page).
2. Over a 1 meter long bar are extracted N points according to an unknown distribution where each extraction is independent from the others. Starting from the extracted points is selected the smallest interval quantized to a millimeter (that is, the infimum of the interval is x mm and the supremum is y mm, where x and y are integers) that contains all the extractions. Using Hoeffding inequality say how many extractions is sufficient to have in order to be sure with confidence $\beta = 10^{-10}$ that the selected interval contains at least the 95% of the probability distribution.

EXERCISE 4 (optional)

Two investments funds A and B give a yearly financial return r^A and r^B ($\text{financial return} = \text{final value}/\text{initial value}$) a priori unknown and variable according to the investment year. So, if an investor invests 1\$, x \$ on A and $(1-x)$ \$ on B, the money that the investor will have at the end of the year will be $xr^A + (1-x)r^B$.

1. You have N independent observations $r_i^A, r_i^B, i=1, \dots, N$, of the financial returns of A and B in the past years (historical data). Consider the following optimization program:

$$\max \text{ of } r \text{ such that: } r \leq xr_i^A + (1-x)r_i^B, i=1, \dots, N$$

and denote with x^*, r^* its solution.

Give an interpretation of the meaning of x^*, r^* .

2. Say what is maximum number of the support constraints for the problem at point 1.
3. Suppose to invest the next year 1\$, where x^* \$ on A and $(1-x^*)$ \$ on B. Supposing that the financial returns are independent random variables in different years and that their distribution doesn't vary over the years, determine the probability that the money that the investor will have at the end of the next year is less than r^* .

ES. 1

Si hanno a disposizione alcune osservazioni $(u_i, y_i) \quad i=1 \dots N$ reali ($u_i, y_i \in \mathbb{R}$). Al fine di descrivere le relazioni che lega u e y si usa il modello

$\hat{y} = \theta u^2$, e il parametro θ viene stimato ai minimi quadrati.

1) Si scrivono le equazioni parametriche per la determinazione di θ .

$$\sum \varphi_i \varphi_i^T \theta = \sum \varphi_i y_i \Rightarrow \sum_{i=1}^N u_i^4 \theta = \sum_{i=1}^N u_i^2 y_i$$

$$\text{Qui } \varphi_i = u^2, \quad \varphi_i^T = u^2$$

2) Sono date le 3 osservazioni

$$\begin{aligned} &(1, 1) \\ &(1.5, 1) \\ &(1.8, 1.5) \end{aligned}$$

Determinare le stime ai minimi quadrati di θ e disegnare il modello corrispondente (cioè se curva è una parabola).

u_i	u_i^4	u_i^2	$u_i^2 y_i$	$\sum_{i=1}^N u_i^2 y_i = 8.11$
1	1	1	1	
1.5	5.06	2.25	2.25	
1.8	10.5	3.24	4.86	

$$\theta = \frac{8.11}{16.56} \approx 0.49$$

$$16.56$$

1.5

1.25

0.5

1

1.5

1.8

2

modello

$i=1$

$$\sum_{i=1}^N u_i^4 = 16.56$$

$$\text{modello: } \hat{y} = 0.49 u^2$$

$$\text{con } u=1 \quad \hat{y} \approx 0.5$$

$$\text{con } u=1.5 \quad \hat{y} \approx 1.125$$

3) Si supponga che i dati siano generati dal meccanismo $y = 0.6 u^2 + n$, dove u è distribuito uniformemente in $[0, 2]$, n è gaussiana $N(0, 1)$. I diversi u sono generati con ingressi u indipendenti gli uni dagli altri e il numero n è indipendente da u . Si dice a cosa tende la stima ai minimi quadrati quando $\#$ osservazioni $\rightarrow \infty$.

1) Se il modello è $y_i = 0.6 u_i^2 + n_i$

Converge a 0.6 sicuramente. n_i ha media nulla e indipendente dalle u_i

Per es. peggi forte dei grandi numeri $\rightarrow 0.6$

2) Se n_i non è una successione, ma è una variabile

~~$\hat{\theta}_k = \frac{1}{k} \sum u_i^4$~~

$$y_i := \varphi_i^T \theta_0 + r_i$$

$$\hat{y}_i := \varphi_i^T \theta$$

$$y_i - \hat{y}_i = \varphi_i^T (\theta - \theta_0) + r_i$$

è un resto

L'errore di polarizzazione è:

Non c'è più la media

$\sum \varphi_i \rightarrow 1$ sono uniformi

$\sum \varphi_i \varphi_i^T \rightarrow$ matrice di positivo

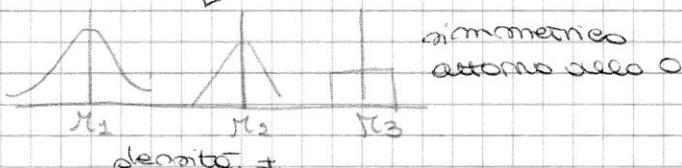
riportare le stime

$$(\sum \varphi_i \varphi_i^T)^{-1} \sum \varphi_i r_i$$

1) Si descrive brevemente (≤ 1 di pagina) a parole a cosa serve il metodo LSCR (Leave out 3 Sign-dominant Correlation Region).

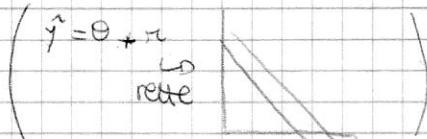
Scopo: Fornire un intervallo di confidenza con probabilità certificata per il parametro θ in un modello lineare in θ nelle ipoteze che il rumore sia additivo e costituito da variabili indipendenti dotate di densità simmetrica rispetto all'origine (non necessariamente identicamente distribuite).

$$y_i = \mu_i \theta + r_i$$



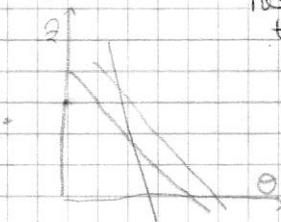
2) Si considerano ancora le 3 osservazioni del punto 1) delle es. precedente. Si determini una regione Ω per il parametro θ del modello $y_i = \mu_i \theta + r_i$ con il metodo LSCR.

- Continuavano un grappo



Qui: $z = y_1 - \mu_1^2 \theta$ rette con una sola y e una sola μ

non hanno +
tutte = pendenze



Scegli 2 rette:

Oss₁ Oss₂ Oss₃

Medie di
2 osservazioni

Nelle es. in area

1-2-4-5 somma

1-3-4-6 2-3-5-6

2-3-5-6

1-2-6-7

1-3-5-7

2-3-4-7

6-5-6-7

elemento neutro

1-2-4-5

2-3-5-6

1-2-6-7

1-3-5-7

2-3-4-7

6-5-6-7

elemento neutro

Ogni riga ha = n di elementi.

elemento neutro
del grappo
(riga 0)

$$z = y_1 - \mu_1^2 \theta$$

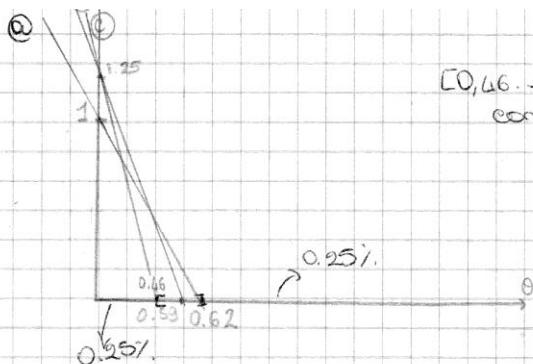
$$z = y_2 - \mu_2^2 \theta$$

$$z = y_3 - \mu_3^2 \theta$$

@ 1^a media: $z = \frac{y_1 + y_2}{2} - \frac{\mu_1^2 + \mu_2^2}{2} \theta$ $z = 1 - 1,62 \theta$

⑥ 2^a " $z = \frac{y_1 + y_3}{2} - \frac{\mu_1^2 + \mu_3^2}{2} \theta$ $z = 1,25 - 2,12 \theta$

⑦ 3^a " $z = \frac{y_2 + y_3}{2} - \frac{\mu_2^2 + \mu_3^2}{2} \theta$ $z = 1,25 - 2,74 \theta$



[0.46, 0.62] intervallo di confidenza

con probabilità = 50%.

(cioè 4 regioni in cui suddivisione è omogenea, tutte equiprobabili)

prima (es. 1) $\theta = 0.48$ è all'interno

ES. 3

1. La diseguaglianza di Hoeffding si scrive $\hat{P}(A) - P(A) \leq \epsilon$

con probabilità (1 - β) $\geq 1 - \epsilon^2 / 2N^2$ (armonico cioè $\epsilon^2 / 2N^2$ bound!!)

dove $\beta = e^{-\epsilon^2 / 2N^2}$

(Hoeffding: $P[... \leq \epsilon] \geq \alpha$ alfa di piccolezza).

$$P[S_N - E[S_N] \geq \epsilon] \leq \exp\left(\frac{-2N\epsilon^2}{\sum(b_i - a_i)^2}\right)$$

$$P[S_N - E[S_N] \leq \epsilon] \leq \exp(...)$$

dove $\hat{P}(A)$ è la stima campionaria delle probabilità di un evento A. Si spieghi le parole: è significativo di ϵ e β . ($\leq 1/2$ pagina)

Qui $S_N = \frac{1}{N} \sum_{i=1}^N I_{\{x_i \in A\}}$ funzione indicatrice
 numero di casi in cui A si verifica
 $E[S_N] = P[x \in A]$ numero totale di casi (= N)
 probabilità dell'evento

$\hat{P}(A)$ ex. Stima che un italiano sia biondo \rightarrow quando 5 persone e ne vedo 2.

ϵ = numero deciso a priori x stabilire quanto la probabilità campionaria si avvicina a quella vera \rightarrow accuratezza con cui voglio approssimare la distribuzione empirica. ϵ è piccolo, la differenza è piccola.

1 - β : prima di fare l'esperimento è una probabilità. Dopo l'esperimento ho numeri, quindi è una confidenza.

È piccolo con probabilità. 1 - β = confidenza

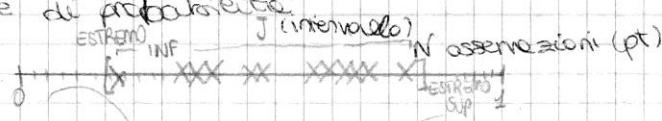
$$P[S_N - E[S_N] \leq \epsilon] \geq 1 - \exp\left(-\frac{2N\epsilon^2}{\sum(b_i - a_i)^2}\right)$$

Graebiger:

$$P[|x - \mu| \geq \epsilon] \leq \frac{\sigma^2}{\epsilon^2}$$

2. Vengono estratti N punti su un'asse di lunghezza 200 cm, in accordo ad una distribuzione non nota, dove ogni estrazione è indipendente dalle altre. A partire dalle osservazioni, la selezione è più piccolo intervalli quantizzati al millimetro, che si compongono tutti.

L'intervalllo x millimetri - y millimetri, dove x e y sono interi]. Un esempio: la diseguaglianza di Hoeffding si dice quante estrazioni sono sufficienti per essere sicuri con confidenza 1 - β , dove $\beta = 10^{-10}$, che è intervallo contenuto contenuto nel 99,9% della distribuzione di probabilità.



to stare in accordo ad una distribuzione non nota

È un bound sulla confidenza, non è una vera confidenza, su cui le dirò subito.

Se N non è abbastanza grande il bound esce dall'intervallo (0, 1) \rightarrow occorrono sufficienti osservazioni, delimitando più veramente ben > 1.

J denotato su $J \geq 0.95$

Quante estrazioni = quanto è grande N

Avrò, ma sbagliato:

$$e^{-2N \cdot 0.05^2} = 10^{-10} \quad \epsilon = 5\%.$$

$P(A) = 1$ se mi controviamo appartenente all'intervalllo

$$1 - P(A) \leq \epsilon \quad P(A) \geq 1 - \epsilon \\ \Rightarrow \epsilon = 0.05\% \\ P(A) \geq 95\%.$$

Risolvere per N :

$$e^{-2N \cdot 0.05^2} = e^{-10} \quad \text{risolvere } N$$

E' sbagliato perché necessaria disegualanza di Hoeffding
è evento deve essere fisso, così come le sue intervalli.
Qui l'intervalllo min-Max è casuale, non fisso.

$$\text{Sv} = \frac{1}{N} \sum_{i=1}^N I_A(x)$$

l'intervalllo fisso, evento fisso!

Qui l'intervalllo fisso a priori sono tutti gli intervallli possibili.

Sono $\frac{1000 \cdot 1000}{2} \approx 1.000.000 = 10^6$ intervallli possibili

(estremo inf fra tutti

i < 1000 pt possibili,

così come l'estremo sup.).

$$i \quad 10^6 \rightarrow$$

Nello spazio campionario, f ciascun intervalllo ha una certa probabilità di sbagliare. (vedo sotto una certa soglia $\bar{\beta}$)

$P(A)$ vale qualcosa

Per l'intervalllo i_2 , la disegualanza $P(A) - P(A)$ non ^{confidenza} vale per $\bar{\beta}$.

... AVERE 10^6 intervallli, in cui ognuno esclude una parte di probabilità sempre con = confidenza.

P complessiva che mi sbagli \leq somma delle rispettive (unione di eventi) probabilità ($= \bar{\beta}$).

$$\text{Area complessiva} \leq 10^{-10}$$

$\bar{\beta}$ deve essere preso in modo che la somma sia al massimo 10^{-10}

A deve essere fissato a priori.

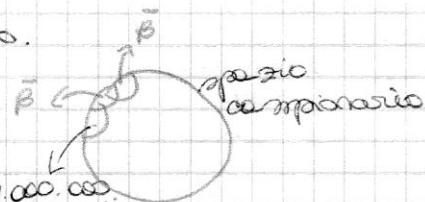
$$\text{L} \bar{\beta} = \frac{10^{-10}}{10^6} = 10^{-16} \quad \text{corretto.}$$

$$\text{poi } e^{-2N \cdot 0.05^2} = 10^{-16} \quad \text{e risolvere per } N$$

vedi
dim.
lunga
simile

Un intervallo, esiste una parte dello spazio campionario di questi intervallli con probab. 10^{-10} che fa di sbagliarsi di Hoeffding della parte 1. non venga.

\Rightarrow Anche se somma delle aree $\leq 10^{-10}$



ES. 4

2 fondi di investimento A e B danno un ritorno finanziario annuale R^A , R^B .
 R^A e R^B sono casuali.

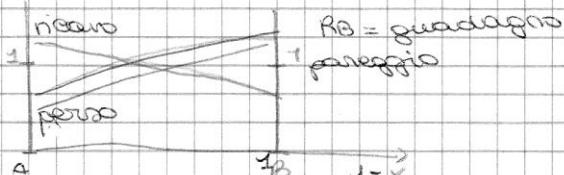
per i ritorni finanziari di A e B negli anni passati
si fanno N osservazioni indipendenti $R_i^A, R_i^B, i=1 \dots N$ (ad es. su anni).

I dobbiamo calcolare la probabilità che i 2 fondi x^A in A
 $1-x^A$ in B dove x sono
centrati.

Se un investitore investe x dollaro, di cui x
su A e $1-x$ su B il dovrà ricevere al ritorno

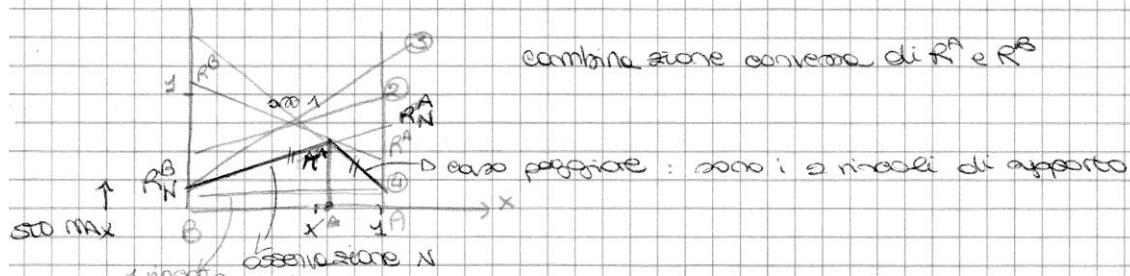
proprio: $R = x \cdot R^A + (1-x) R^B$

investimento x (A) + $(1-x)$ (B)



In N periodi avremo N osservazioni +
coppie.
→ sono tutte rette

Se $x=1$ $1-x=0$ investo tutto in A



Si consideri il problema di ottimizzazione

r_{\max} o r

$$\text{t.e. } r \leq x r_i^A + (1-x) r_i^B \quad i=1 \dots N$$

$$\text{sol. } (x^*, r^*)$$

2 variabili x, r

x questo che investo nel fondo A
 r questo che ne ottengo

x^*, r^* ? Spiegazione a parte:

x^* quanto per cui in passato avei guadagnato di + investendo nel caso peggiore. Anzi massimizzato il minimo ritorno possibile in tutti gli anni (tenendo le stesse).

Lo caso peggiore rispetto a tutte le annate possibili (curva nera)

e lo sto massimizzando (per minimizzare le perdite)

al variare di r e x se investo in x^* guadagno almeno r^* in ogni anno

(dim=2, vincoli di supporto).

affatto al massimo) x^* - prozione di investimento che garantisce il maggiore ritorno r^* negli anni passati.

Si consideri il problema (=ma ribaltato)

$$\min_{\theta} \gamma$$

$$\text{t.e. } \theta g^A + (1-\theta) g^B - \gamma \leq 0$$

$$\theta \in [0,1]$$

Il problema ha dimensione 3.
 \max no vincoli di supporto = 2

$$\dim_{\theta} = 1$$

3) Si supponga di investire le prossime anni x dollaro, di cui x^* su A e $1-x^*$ su B.
Supponendo che R^A e R^B siano indipendenti su anni + e la loro distribuzione
non varia di anno in anno, probabilmente che le ritorno < r^* .

$P[N+1 - \text{esime diversi di supporto}]$

nell'oo. $\min_{\theta} \gamma \# 2$

$P[\max_{\theta} \gamma \# 0]$

potrei avere 1 solo

vincolo di supporto (l'ultimo vincolo osservato) ≤ 2 al max - che i vincoli di supporto sono 2.

con probabilità $\# 0$ nulla (a causa di tali rintralcie)

Sono in x^* → l'anno prossimo ha una nuova retta
di guadagno minore di prima ($< r^*$) e nuovo
vincolo di supporto.

Non posso + dire
che i vincoli di supporto sono 2.