



Linear models

Gabriela Ciolek

Gabriela.ciolek@telecom-paristech.fr

Two variables

- y = explained variable, dependent
- x = explanatory variable, independent
- u = non-observed variable, error, disturbance

$$y = \beta_0 + \beta_1 x + u.$$

- β_0 = constant, intercept
- β_1 = partial effect = $\partial y / \partial x$ (slope parameter in the relationship between y and x holding the other factors in u fixed)

Example 1

-
- y = salary
 - x = education
 - u = labor force experience, productivity, tenure with current employer

$$wage = \beta_0 + \beta_1 educ + u.$$

- Ceteris paribus ? Explain link x - u

Example 2

- y = yield
- x = fertilizer
- u = quality of land, climat etc.

$$yield = \beta_0 + \beta_1 fertilizer + u,$$

- Partial effect, ceteris paribus. $\Delta yield = \beta_1 \Delta fertilizer.$
- Ceteris paribus ? Explain link $x-u$

Fundamental hypothesis

- Error variables have zero mean
- Conditional expectation = unconditional expectation
- X and u are uncorrelated

$$E(u|x) = E(u) = 0,$$

- Education (x) and tenure with current employer (u)

Population regression function

$$E(y|x) = \beta_0 + \beta_1 x$$

- $E(y|x)$ is a linear function of x (the linearity means that a one-unit increase in x changes the expected value of y by β_1)

Predictions

- Fitted value

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i,$$

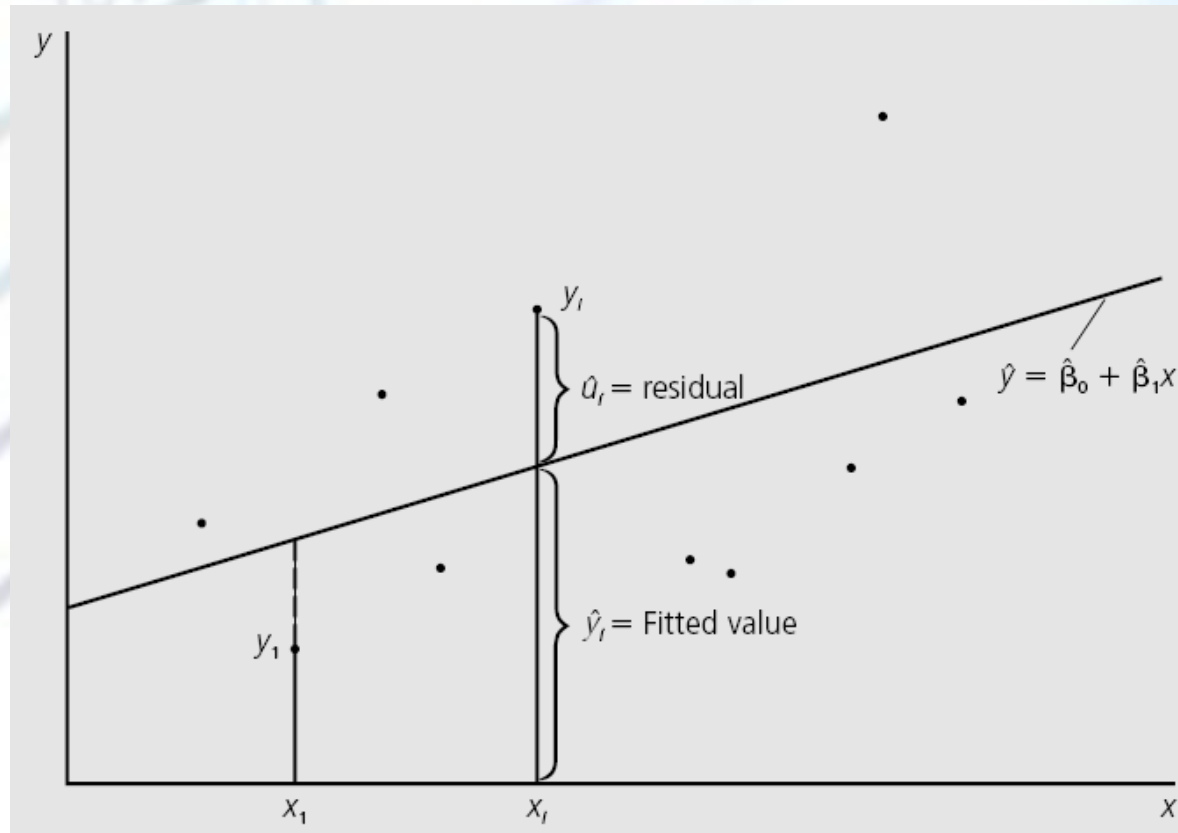
- Residual

$$\hat{u}_i = y_i - \hat{y}_i = y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i.$$

- Sum of squared residuals

$$\sum_{i=1}^n \hat{u}_i^2 = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2,$$

Graphical illustration



K explanatory variables

- Multiple linear regression model

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \dots + \beta_k x_k + u,$$

- Fundamental hypothesis

$$E(u|x_1, x_2, \dots, x_k) = 0.$$

(all factors in the unobserved error term be uncorrelated with the explanatory variables)

Obtaining the OLS Estimates

- We minimize the sum of squared estimators

- Matrix notation:
$$\min_{b_0, b_1, \dots, b_k} \sum_{i=1}^n (y_i - b_0 - b_1 x_{i1} - \dots - b_k x_{ik})^2.$$

- $\text{Min } (y - x\beta)'(y - x\beta)$

- First order conditions (concavity criteria)

- $(y - xb)' x = 0$

- Remark: we assume $x_0 = 1$

Matrix notation

-
- $X'(y - Xb) = 0$
 - $X'y - X'Xb = 0$
 - $b = (X'X)^{-1}X'y$
 - Problem if $X'X \approx 0$: multicollinearity

(we call this problem collinearity: it looks like we have p different predictor variables, but really some of them are linear combinations of the others, so they don't add any information)

- $X'X$ = "variance"
- $X'y$ = "covariance"

Unbiased estimators

- The OLS estimators are unbiased estimators of the population parameters
- $E(b) = E(X'X)^{-1}X'y = E(X'X)^{-1}X'(X\beta + u) = \beta$
- Fundamental hypothesis
- $E(u | X) = 0$



The background image is a blurred photograph of a document. It features a table with numerical data and a bar chart. The table has three rows with values 100,000, 10,000, and 10,000. Below the table, there is a bar chart with a single bar labeled 5,205. The text 'Data analysis with Matlab' is overlaid on the center of the image.

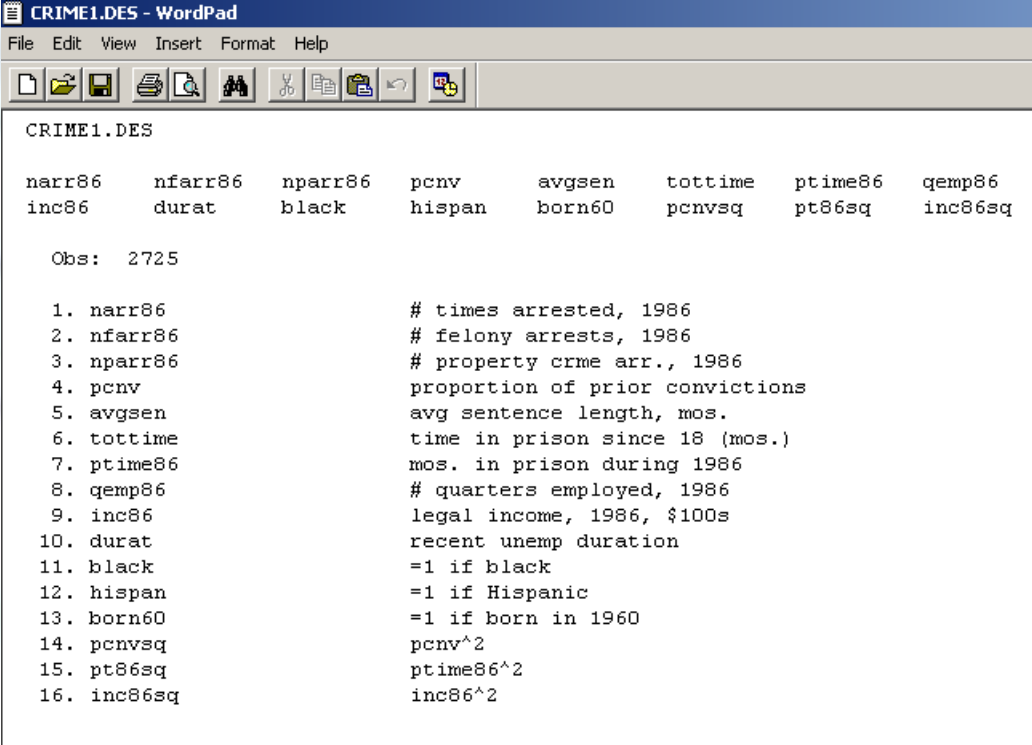
Data analysis with Matlab

Loading the data

load CRIME1.raw

Description of the data

■CRIME1.des



```
CRIME1.DES

narr86      nfarr86      nparr86      pcnv      avgsen      tottime      ptime86      qemp86
inc86      durat      black      hispan      born60      pcnvsq      pt86sq      inc86sq

Obs:  2725

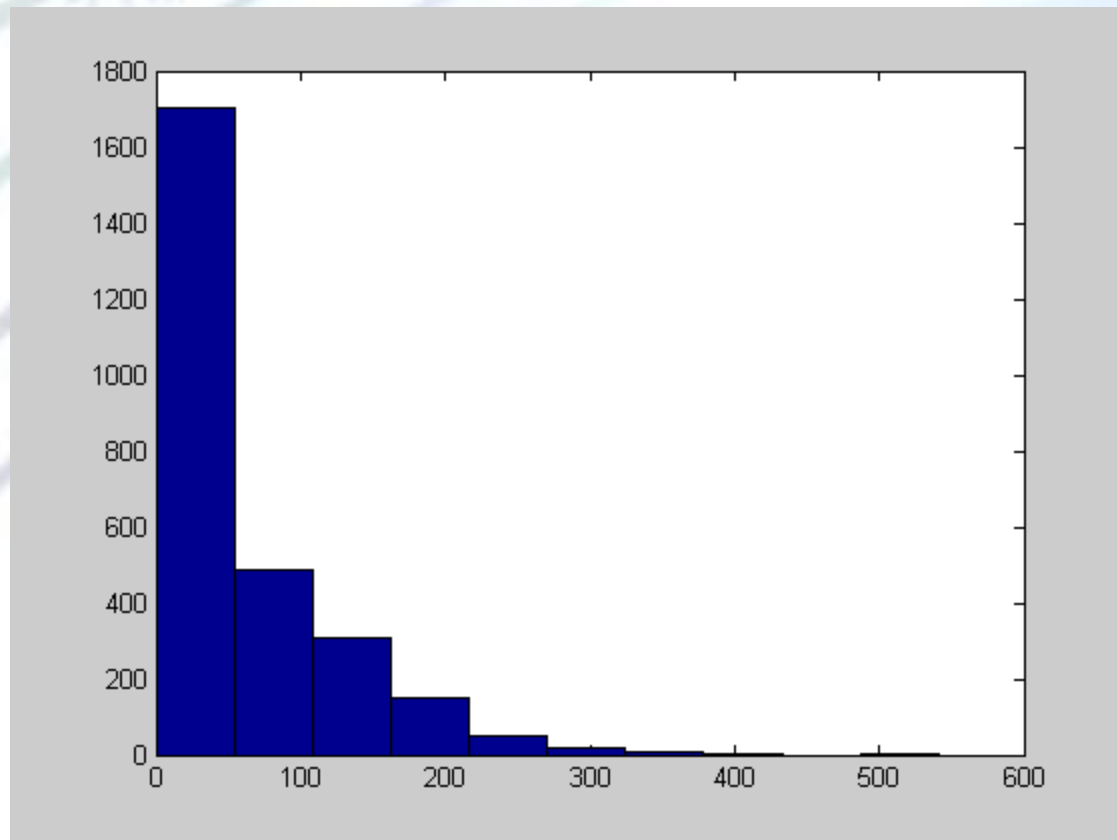
1. narr86      # times arrested, 1986
2. nfarr86      # felony arrests, 1986
3. nparr86      # property crime arr., 1986
4. pcnv      proportion of prior convictions
5. avgsen      avg sentence length, mos.
6. tottime      time in prison since 18 (mos.)
7. ptime86      mos. in prison during 1986
8. qemp86      # quarters employed, 1986
9. inc86      legal income, 1986, $100s
10. durat      recent unemp duration
11. black      =1 if black
12. hispan      =1 if Hispanic
13. born60      =1 if born in 1960
14. pcnvsq      pcnv^2
15. pt86sq      ptime86^2
16. inc86sq      inc86^2
```

Exercise

- Draw a histogram of incomes in dollars in 1986.

Histogram

```
hist(crime1(:,9))
```

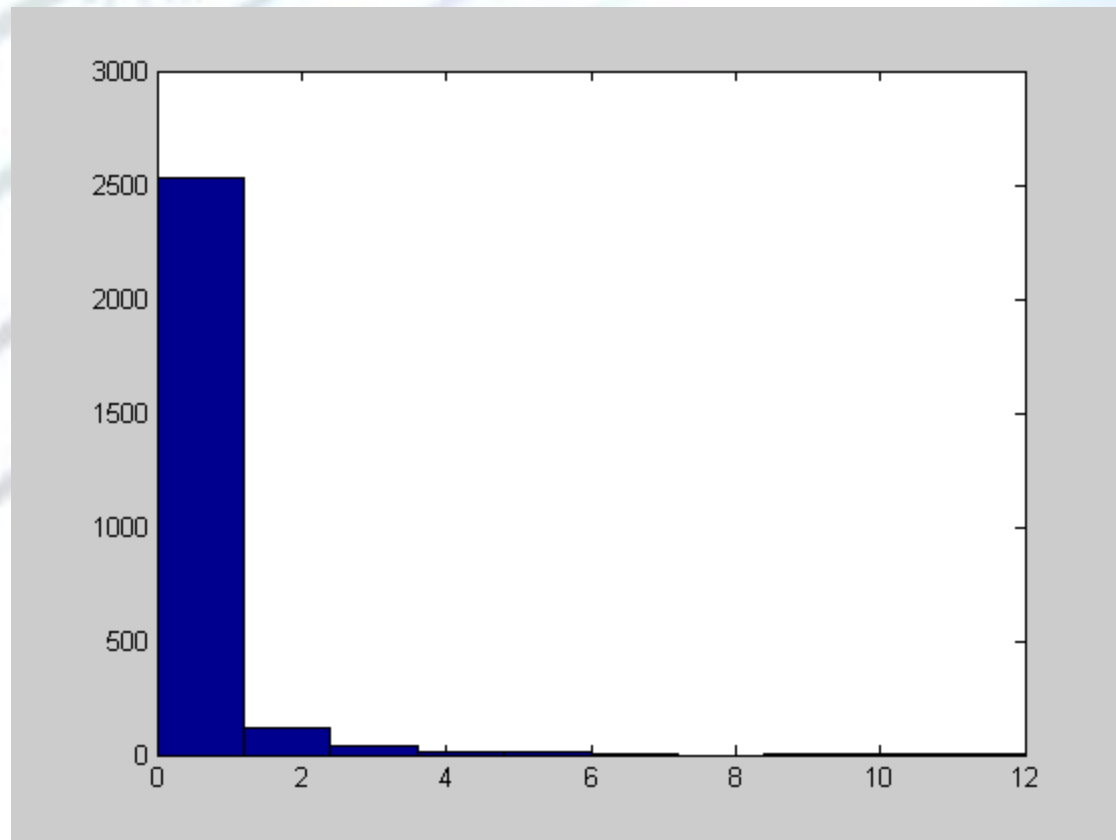


Exercise

- Draw a histogram of number of times when an individual was arrested in 1986.

Histogramme

```
hist(crime1(:,1))
```



Exercise

Calculate:

- expectation
- variance
- correlation

between the number of times when the individual was arrested and other variables.

Descriptive statistics

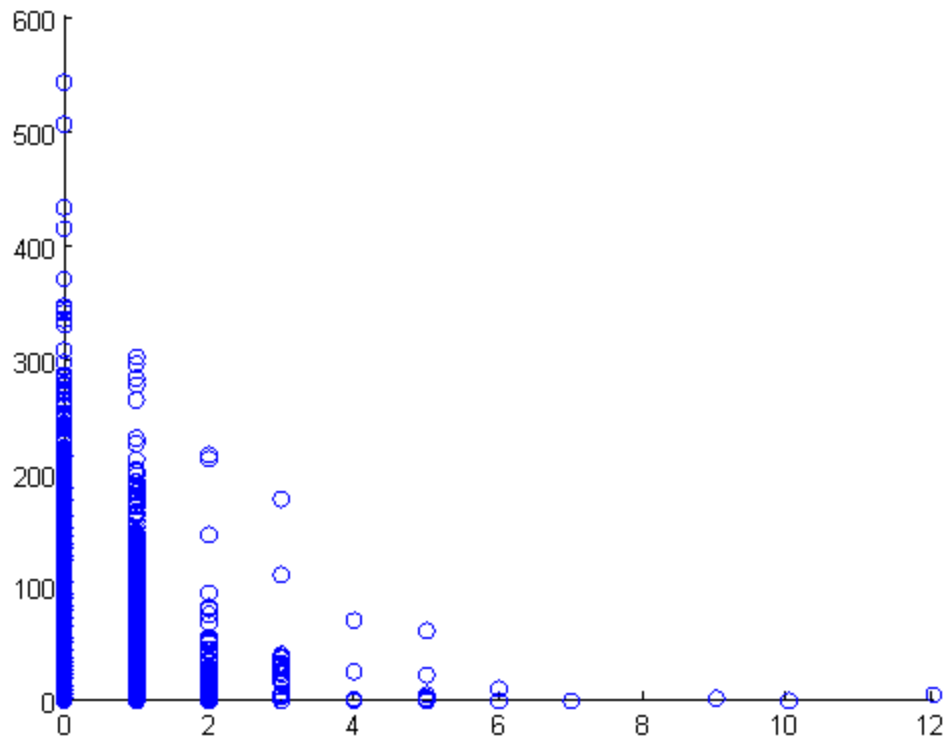
```
mean(crime1)'  
mean(crime1(:,1:13))'  
std(crime1(:,1:13))'  
min(crime1)'  
max(crime1)'  
cov(crime1)'  
corrcoef(crime1(:,1:13))'
```

Exercise

Draw the cloud of points (X,Y) with X the income and Y the numer of arrests

Scatterplots

```
scatter(crime1(:,1),crime1(:,9))
```



The background is a blurred image of a document. On the left, there is a table with three rows. The first row contains the number '100,000'. The second and third rows contain the number '10,000'. Below the table, there is a bar chart with a single bar labeled '5,205,000'. The text 'Analysis of the data' is overlaid in the center of the image.

Analysis of the data



Import of the data

load WAGE1.raw

Description of the data

■ WAGE1.des

```
WAGE1.DES - WordPad
File Edit View Insert Format Help

WAGE1.DES

wage      educ      exper      tenure      nonwhite      female      married      numdep
smsa      northcen  south      west      construc  ndurman      trcommpu      trade
services  profserv  profocc    clerocc     servocc      lwage        expersq      tenursq

Obs:      526

1. wage      average hourly earnings
2. educ      years of education
3. exper      years potential experience
4. tenure     years with current employer
5. nonwhite   =1 if nonwhite
6. female     =1 if female
7. married    =1 if married
8. numdep     number of dependents
9. smsa       =1 if live in SMSA
10. northcen  =1 if live in north central U.S
11. south     =1 if live in southern region
12. west      =1 if live in western region
13. construc  =1 if work in construc. indus.
14. ndurman   =1 if in nondur. manuf. indus.
15. trcommpu  =1 if in trans, commun, pub ut
16. trade     =1 if in wholesale or retail
17. services  =1 if in services indus.
18. profserv  =1 if in profess. serv. indus.
19. profocc   =1 if in profess. occupation
20. clerocc   =1 if in clerical occupation
21. servocc   =1 if in service occupation
22. lwage     log(wage)
23. expersq   exper^2
24. tenursq   tenure^2
```

Comparison of salaries between men and women

Calculate:

- The average hourly wage for two sexes
- Average hourly wage for each sex separately
- Histograms of the salaries for each sex
- Descriptive statistics for other variables for each gender (expected value, variance, etc.)

Data selection

- Women
- Number

```
s=(wage1(:,6)==1);  
sum(s)  
wage1f=wage1(s,:)'  
mean(wage1(s,1:21))'
```

Question

Is there a wage discrimination against women?

Response

We can not say based on the histograms of salaries of women and men because the other variables are not fixed.

It is possible that women can be less skilled so can be less paid, etc.

Explain wage | educ, exper, tenure

```
load WAGE1.raw
y=wage1(:,1);
[n,k]=size(wage1)
X=[ones(n,1),wage1(:,[2,3,4])];
[n,k]=size(X)
```

Estimation of parameters

Formula

$$\text{beta} = \text{inv}(X' * X) * X' * y$$

Calculate the residuals, variance of residuals and standard deviation of the estimator

$u = y - X * \text{beta}$

$\text{sig2} = u' * u / (n - 4)$ because 3 variables + intercept

$\text{std} = \text{sqrt}(\text{diag}(\text{sig2} * \text{inv}(X' * X)))$

Exercise

What will happen if:

- One will multiply the salary by 1000?
- One will multiply one explanatory variable by 1000?
- The coefficients will change? Are they robust to change unit?

Exercise

- Draw a histogram of residuals u .
- What are the properties of the distribution?

Exercise

- Do a regression when taking the log of salary, calculate the parameter.
- Draw a new histogram for residuals.
- Properties?
- Detect the outliers, remove them and recalculate beta.

Exercise

Find observations such as $u > 2.5$
Remove those observations

```
I = find(u >= 2.5)
s = (u <= 2.5)
sum(s)
X = X(s,:);
y = y(s,:);
```