

Last Lecture

① Final Exam

(i) Friday → 9:00am - ...

(ii) Two sections

a) Zybooks programs
(3-4)

b) Data Science
Problems
(2)

→ Send to me via
email!!

② Basic Plotting - Full Example.

③ Linear Regression - Complex Data Example

(a) with numpy & scikit-learn

(a) with statsmodels
(b) with pandas & statsmodels

Linear Regression

Complex Data :

<u>x1</u>	<u>x2</u>	<u>x3</u>	<u>x4</u>	<u>y</u>

Example :
y → housing prices
x1 → sq. footage
x2 → zip code
x3 → # of bedrooms
x4 → # of bath rooms

$x_4 \rightarrow$

We would like, in the end, to determine an equation like:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4$$

i.e. an equation which is linear in the possible factors.

N.B. To visualize this, we would need to make a FIVE dimensional plot !!!
 \rightarrow obviously, impossible !!

This is why this is a hard problem to solve and understand.

Key I deal: We need to

determine

$$\beta_0 \pm \delta \beta_0$$

$$\beta_1 \pm \delta \beta_1$$

\vdots

$$\beta_n \pm \delta \beta_n$$

If some $\beta_k \pm \delta \beta_k$ is constant
with zero, then y does
not depend on this factor.

Example : maybe housing prices
do not depend on
the number of bedrooms.

Key Idea 2 :

The "Explore" stage of this
analysis is important!! We
should first plot :

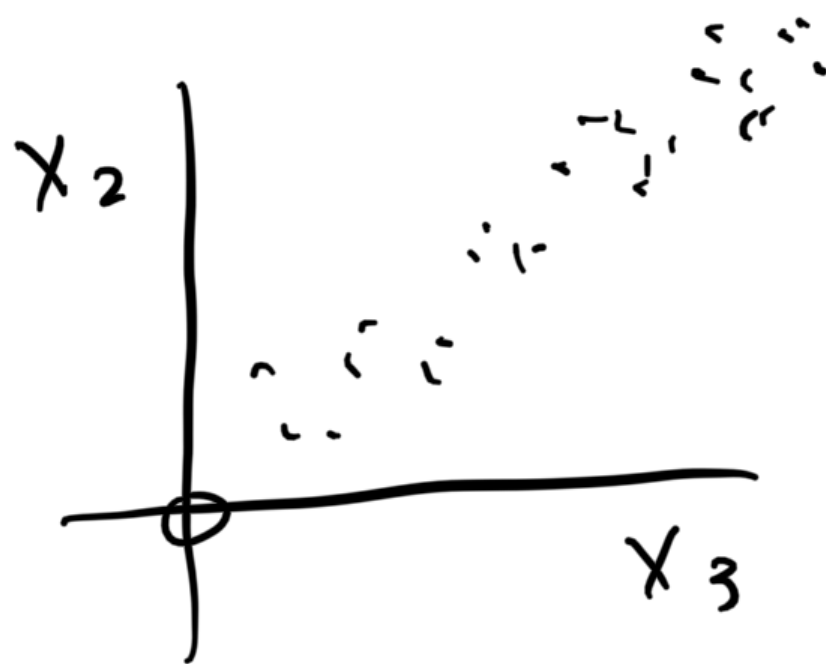
y vs. x_1 , y vs. x_2 , ...

... and some "basic"

to understand the
behavior of the data.

Key Idea # 3 :

We are ignoring correlations.



This can cause misleading results.

N.B.

We need to be
cognizant of

Correlations vs. Causation

Let's start with a simpler
example, so that we can
understand the engine of
in Python.

Linear regression in Python

Step 1 : Get the data into appropriate data structures.

TWO CHOICES :

- ① numpy arrays
- ② Pandas data frames.

(i) numpy arrays.

Aside on Linear Algebra.

Suppose we have :

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

x_1 x_2 y

1

[illegible]

—
—
—

——
——
——

——
——

N data points.

We can write:

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik}$$

$$y_2 = \beta_0 + \beta_1 x_{12} + \beta_2 x_{22}$$

$$y_3 = \beta_0 + \beta_1 x_{13} + \beta_2 x_{23}$$

•

$$y_N = \beta_0 + \beta_1 x_{1N} + \beta_2 x_{2N}$$

Transform this to a matrix equation:

$$\begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix} = \begin{pmatrix} \beta_0 \\ \beta_0 \\ \vdots \\ \beta_0 \end{pmatrix} +$$

easy to
determine
from
average of y data

$$\begin{pmatrix} x_{11} & x_{21} \\ x_{12} & x_{22} \\ x_{13} & x_{23} \\ \vdots & \vdots \\ x_{1n} & x_{2n} \end{pmatrix} \quad \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix}$$

$$\begin{pmatrix} y_1 - \beta_0 \\ y_2 - \beta_0 \\ \vdots \\ y_N - \beta_0 \end{pmatrix} = \begin{pmatrix} x_{11} & x_{21} \\ \vdots & \vdots \\ x_{1N} & x_{2N} \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix}$$

$$\therefore \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix} = \begin{pmatrix} x_{11} & x_{21} \\ \vdots & \vdots \\ x_{1N} & x_{2N} \end{pmatrix}^{-1} \begin{pmatrix} y_1 - \beta_0 \\ \vdots \\ y_N - \beta_0 \end{pmatrix}$$

inverse!!

This is, in a broad sense, how linear regression algorithms work!!

For our purposes, the main point is that the:

X data should be a 2-D Matrix

Y data should be
a 1-D Matrix

Example:

$x_i = [1, 2, 3, 4, 5]$

$y_i = [1.1, 2.1, 2.9, 3.9, 5.1]$



$X = \text{numpy.array}(x_i).$
 $\text{reshape}((-1, 1))$

$y = \text{numpy.array}(y_i)$

print(X)

$\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{bmatrix}$

print(y)

[1.1 2.1 2.9 3.9 5.1]

N.B.

numpy arrays → printed
with no commas !!
(as opposed to Python lists)

Step 2:

Clean the data

- plot y vs. x
- look for bad / missing data.
- think about appropriate relationships.