

Lecture 14 (Ch. 3)

Last time we learned about regression (or fitting).

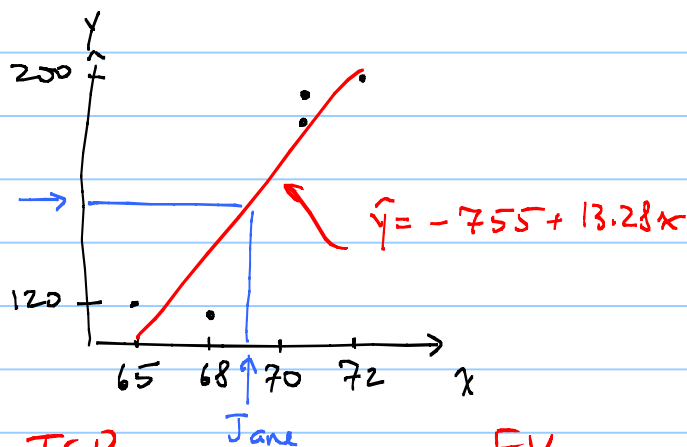
We learned that given data on x and y , the "best" fit is

$$f(x) = \hat{\alpha} + \hat{\beta}x \quad \text{where} \quad \hat{\beta} = \frac{\overline{xy} - \bar{x}\bar{y}}{\overline{x^2} - \bar{x}^2}, \quad \hat{\alpha} = \bar{y} - \hat{\beta}\bar{x}.$$

Here was the example:

Example

	height (x)	weight (y)
	72	200
Joe:	70	180
	65	120
	68	118
	70	190



⇒ We can now predict everyone's ~~weight~~ ^{ICP}, from their ~~height~~ ^{FV}:

Height (x)	weight (y)		\hat{y}	$(y - \hat{y})$
72	200	...	201.5	-1.5
Joe = 70	180		174.9	5.1
65	120		108.5	11.5
68	118		148.3	-30.3
70	190		174.9	15.1

$$\hat{y} = \hat{\alpha} + \hat{\beta}x$$

predicted y

⇒ For the people in the data set, we can also find their error/residual

⇒ For people outside the data set (eg. Jane) we can predict their y from their x , but we cannot compute error, because we don't know their true y . In Ch. 11, we'll address this issue.

Now, let's derive the eqns for $\hat{\alpha}$, $\hat{\beta}$: ↴

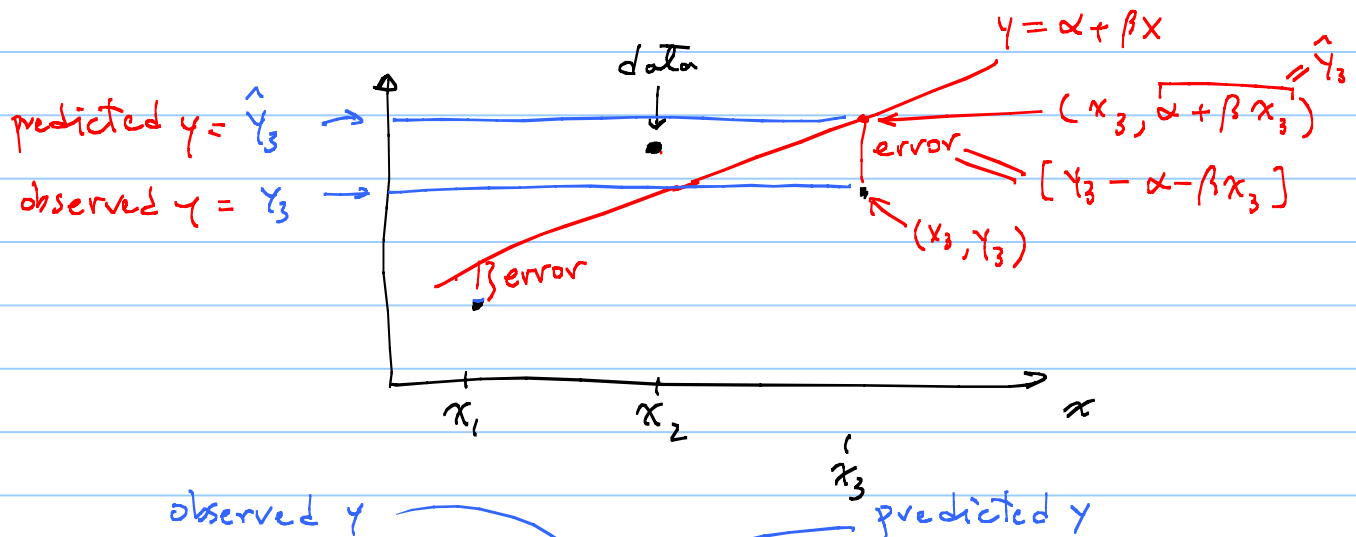
Derivation :

Called Ordinary Least Squares (OLS)

The very common selection criterion is to take the fit (line) that has the smallest Sum of Squared Errors (SSE)

or equivalently Mean " " " (MSE = $\frac{1}{n}$ SSE)

Suppose we have n cases of data: (x_i, y_i) $i=1, 2, 3, \dots, n$



$$\text{MSE} = \frac{1}{n} \text{SSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \alpha - \beta x_i)^2$$

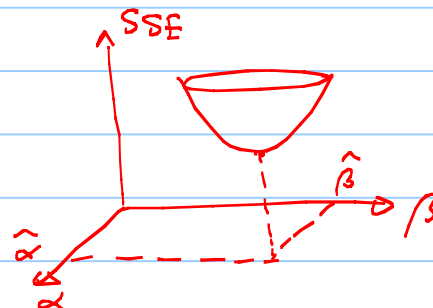
of cases

Minimize MSE \Rightarrow differentiate w.r.t. α, β ; set to zero; solve for the critical values of $\alpha, \beta \Rightarrow \boxed{\hat{\alpha}, \hat{\beta}}$

The specific values of α, β that minimize SSE are called OLS estimates of α, β , and denoted $\hat{\alpha}, \hat{\beta}$:

$$\frac{\partial}{\partial \alpha} \text{MSE}(\alpha, \beta) \Big|_{\alpha=\hat{\alpha}, \beta=\hat{\beta}} = 0$$

$$\frac{\partial}{\partial \beta} \text{MSE}(\alpha, \beta) \Big|_{\alpha=\hat{\alpha}, \beta=\hat{\beta}} = 0$$



If you are not familiar with partial derivatives, $\frac{\partial}{\partial \alpha}$, then just think of them as total derivatives. Let's do one:

Walk
thru
these

$$\begin{aligned}\frac{\partial}{\partial \beta} \text{MSE} &= \frac{1}{n} \sum_i \frac{\partial}{\partial \beta} [y_i - \alpha - \beta x_i]^2 \\&= \frac{1}{n} 2 \sum_i [y_i - \alpha - \beta x_i] [-x_i] \\&= -\frac{2}{n} \sum_i [x_i y_i - \alpha x_i - \beta x_i^2] \\&= -2 \left[\frac{1}{n} \sum_i x_i y_i - \alpha \frac{1}{n} \sum_i x_i - \beta \frac{1}{n} \sum_i x_i^2 \right] \\&= -2 \left[\bar{xy} - \alpha \bar{x} - \beta \bar{x^2} \right] \\&\therefore \boxed{\bar{xy} - \hat{\alpha} \bar{x} - \hat{\beta} \bar{x^2} = 0}\end{aligned}$$

That's 1 eqn for 2 unknowns $(\hat{\alpha}, \hat{\beta})$. But there is $\frac{\partial}{\partial \alpha}$:

$$\frac{\partial}{\partial \alpha} \text{MSE} \Big|_{\hat{\alpha}, \hat{\beta}} = 0 \Rightarrow \boxed{\bar{y} - \hat{\alpha} - \hat{\beta} \bar{x} = 0} \quad \text{See how, below.}$$

Now we have 2 eqns for 2 unknowns. Solve!

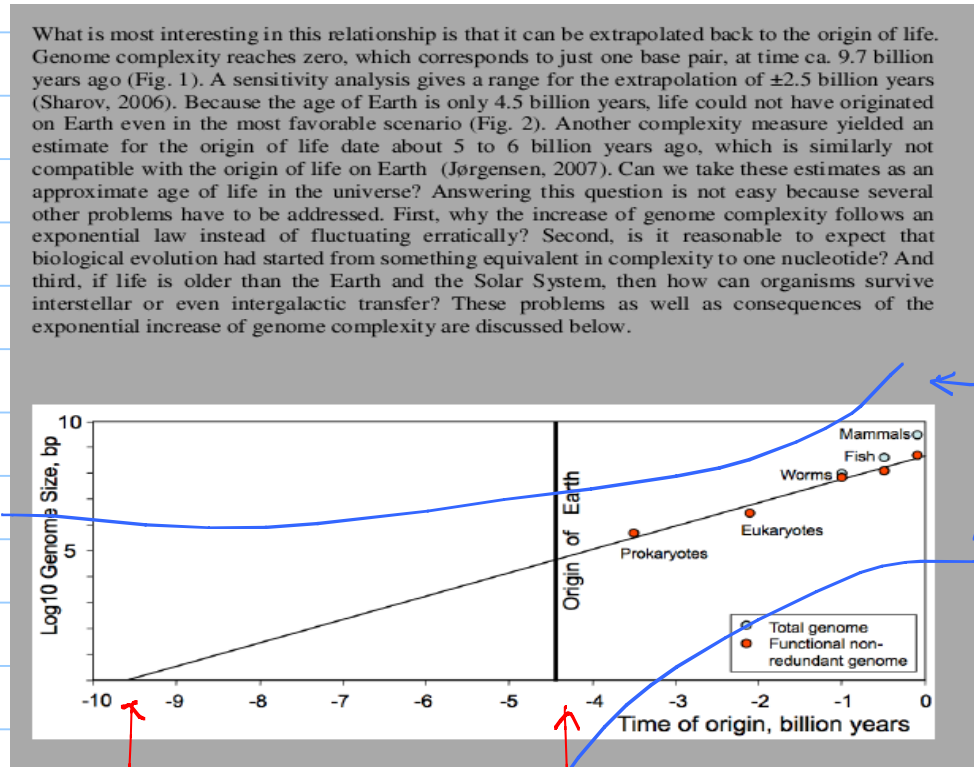
$$\boxed{\hat{\beta} = \frac{\bar{xy} - \bar{x} \bar{y}}{\bar{x^2} - \bar{x}^2}, \quad \hat{\alpha} = \bar{y} - \hat{\beta} \bar{x}}$$

Normal equations of regression.
R: `lm(y ~ x)`

Q1: For the above example, with $\hat{y} = -755 + 13.28x$, the SSE is given by $(-1.5)^2 + (5.1)^2 + \dots + (15.1)^2$. The SSE of any other line ($\hat{y} = \dots$) will be

A) greater B) smaller C) equal to D) Will depend!
 ↑ We minimized SSE to get $\hat{y} = -755 + 13.28x$. I.e. none of the above.

Here is an example of regression where the x-intercept is of interest. Sharov & Gordon (2013) "Life Before Earth":



origin of life

Earth's age

From This \rightarrow They conclude That Life predates Earth, and that life must have been formed on some other planet, then transported to Earth.

In a follow-up paper (Marzban et al. (2014): "Earth Before Life", Biology Direct 9:1) we showed that there are (at least) 2 problems with that analysis

- 1) Extrapolation is bad!
- 2) Uncertainty ("Confidence Bands", in Ch. 11) must be considered.

There is a different (more useful) way of looking at regression, via Variance. This way, we will arrive at quantities called R^2 and s_e , which together assess how good the fit is.

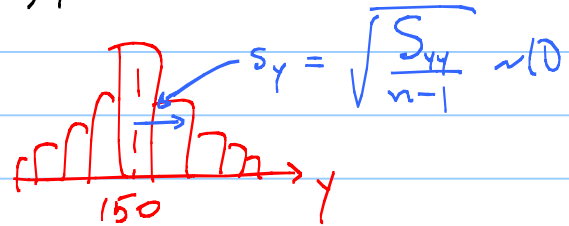
Let me motivate it:

→ Suppose we measure Table Length, y .

→ Repeat, and histogram:

→ One may report:

True length = 150 ± 10 cm



→ Now, suppose you are unhappy with the large s_y . low precision.

→ You may wonder, could some of that variability be due to something else that is varying everytime you make a measurement of y . $x = \text{temperature? humidity?}$

If so, then by measuring y and x , we may be able to reduce the \pm of our report, by specifying y at a given x .

The (scary) math will be in the next lecture.

hw-lect14-1: Show that $\frac{\partial}{\partial \alpha} \text{MSE} |_{\alpha, \hat{\beta}} = 0$ implies $\bar{y} - \hat{\alpha} - \hat{\beta} \bar{x} = 0$

hw-lect14-2: Show that $\hat{\beta}$ as defined by $\frac{\overline{xy} - \bar{x}\bar{y}}{\overline{x^2} - \bar{x}^2}$ or $\frac{S_{xy}}{S_{xx}}$ can be written as $\hat{\beta} = r \frac{S_y}{S_x}$ where $S_x = \text{sample std. dev. of } x$.
 $S_y = \text{sample std. dev. of } y$.

hw-lect14-3:

Suppose data on x and y fall on a straight line $y_i = b + mx_i$. If we perform a linear fit $y = \alpha + \beta x$ to this data, what is the value of the OLS estimate of β ?

hw-lect14-4: According to the OLS fit, what is the predicted value of y , when $x = \bar{x}$ (i.e. when $x = \text{sample mean of } x$).
Hint: All you need is $\hat{\alpha}$.

hw-lect14-5: Come up with another example of x and y (like FV and ICP), where regression can help in predicting y from x in a situation where without regression the "cost" of measuring y directly is extremely high (like ICP).

This document was created with Win2PDF available at <http://www.win2pdf.com>.
The unregistered version of Win2PDF is for evaluation or non-commercial use only.
This page will not be added after purchasing Win2PDF.