

# Linear Regression

In engineering we often have more than one variable in an application and it is known that there exists some inherent relationship among the variables. We will study the case with 2 variables.

Example : Variable 1: Distance to transmitter :  $X$

Variable 2: Wireless signal strength :  $Y$

Lets assume a linear relationship between  $Y$  and  $x$  is reasonable:

$$Y = \alpha + \beta X$$

↑      ↑      ↑  
 Dependent variable    intercept    slope      Independent variable  
 (Regressor)

Note : Don't confuse  $\alpha, \beta$  with the type I & II error probabilities in hypothesis testing.

If the relationship between  $Y$  and  $X$  is exact, then it is a deterministic relationship.

In real applications, there are many sources of randomness:

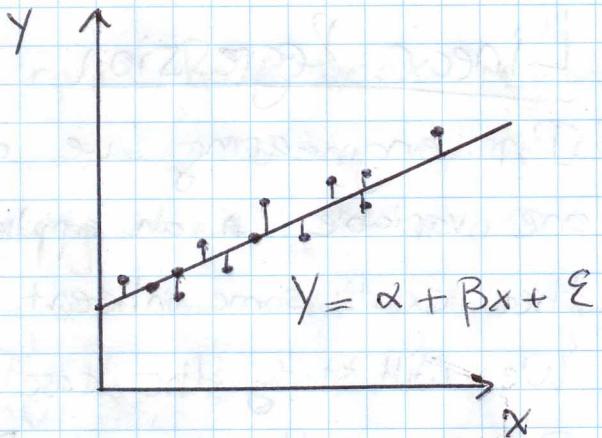
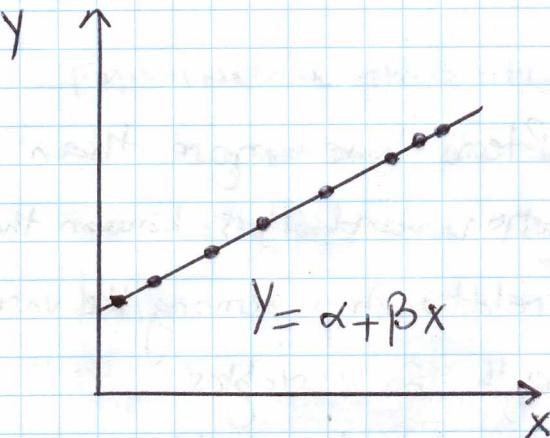
- measurement noise
- the linear regression model might be an approximation to a much more complicated and possibly unknown relationship
- other factors (variables) not considered in model

Randomness means the same value of  $X$  does not always give the same value of  $Y$  (non-deterministic)

## Simple Linear Regression Model :

$$Y = \alpha + \beta X + \varepsilon$$

↑      ↑      ↑  
 Dependent variable    Unknown intercept    Regressor  
 Random variable       $E[\varepsilon] = 0$   
 Unknown slope       $\sigma^2$  Variance



Important:  $\epsilon$  is not a fixed number, it is a random variable which takes on a different value at each data point. The  $\epsilon$  value at each point is shown as bars to the line  $Y = \alpha + \beta x$  in the plot on the right above.

The only conditions on  $\epsilon$  are  $E[\epsilon] = 0$

$$\text{Var}[\epsilon] = \sigma^2$$

- \*  $E[\epsilon]$  implies that at a specific  $x$ , the  $Y$  values are distributed around the true (population) regression line.
- \* In practice  $\alpha$  and  $\beta$  are unknown and must be estimated from data.

More advanced topics:

a) If there are more than one regressor variables

$$Y = \alpha + \beta_1 x_1 + \beta_2 x_2 + \epsilon$$

Example: Electromagnetic absorption in head due to cell phone use :  $Y$

Signal frequency =  $x_1$

Signal strength =  $x_2$

Head size =  $x_3$

b) If the relationship is nonlinear, There are many non-linear models, but a simple one would look like:

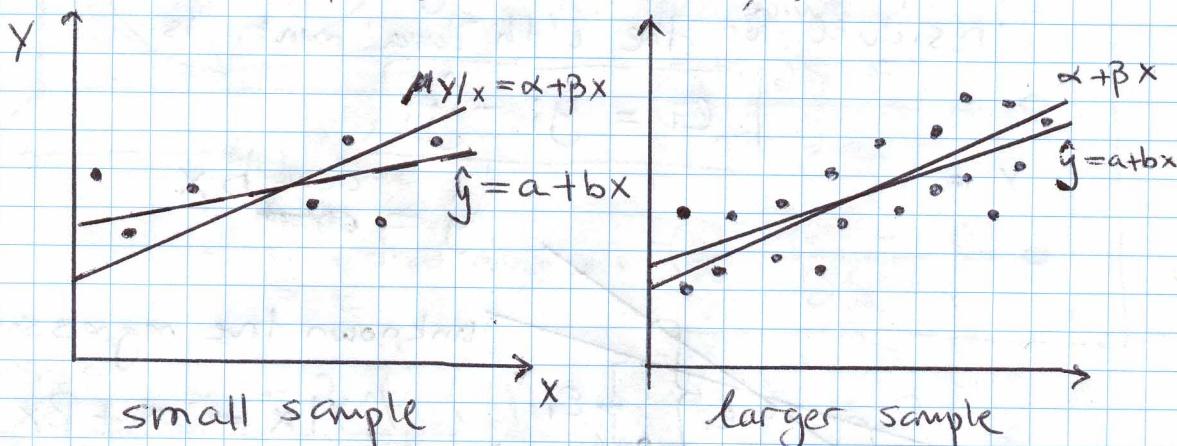
$$Y = \alpha + \beta x + \gamma x^2 + \epsilon$$

This is a quadratic relationship.

Finding the parameters  $\alpha$  and  $\beta$  from data

Example: In the example with  $Y$  signal strength and  $X$  distance to transmitter, we are given the pairs of observations (sample)  $(x_i, y_i)$   $i=1$  to  $n$  and we want to estimate  $\alpha, \beta$ .

Let  $a$  be our estimate of the true population parameter  $\alpha$ . Let  $b$  be our estimate of the true population parameter  $\beta$ . Just like the sample mean  $\bar{x}$ ,  $a$  and  $b$  depend on the particular sample (they are random!) The larger the sample size the closer  $a$  should be to  $\alpha$  and  $b$  to  $\beta$  (just like  $\bar{x}$  to  $\mu$ )



$E(Y|X) = \alpha + \beta X$  : expected value of  $Y$  given  $X$  (true regression)  
 $y = \alpha + \beta X$

$\hat{y} = a + b x$  : fitted regression

↳ predicted (fitted) value

## Conceptual model errors $\epsilon_i$

$$Y = \alpha + \beta X + \epsilon$$

$$\therefore E[Y] = E[\alpha + \beta X + \epsilon]$$

$$= \alpha + \beta X + \underbrace{E[\epsilon]}_{\text{these are not random}}$$

○ assumption about  $\epsilon$

$$\mu Y|x = E[Y] = \alpha + \beta x$$

Let  $x_i$  be the data points for  $X$

$$y_i = \alpha + \beta x_i + \epsilon_i$$

$$\epsilon_i = y_i - (\alpha + \beta x_i)$$

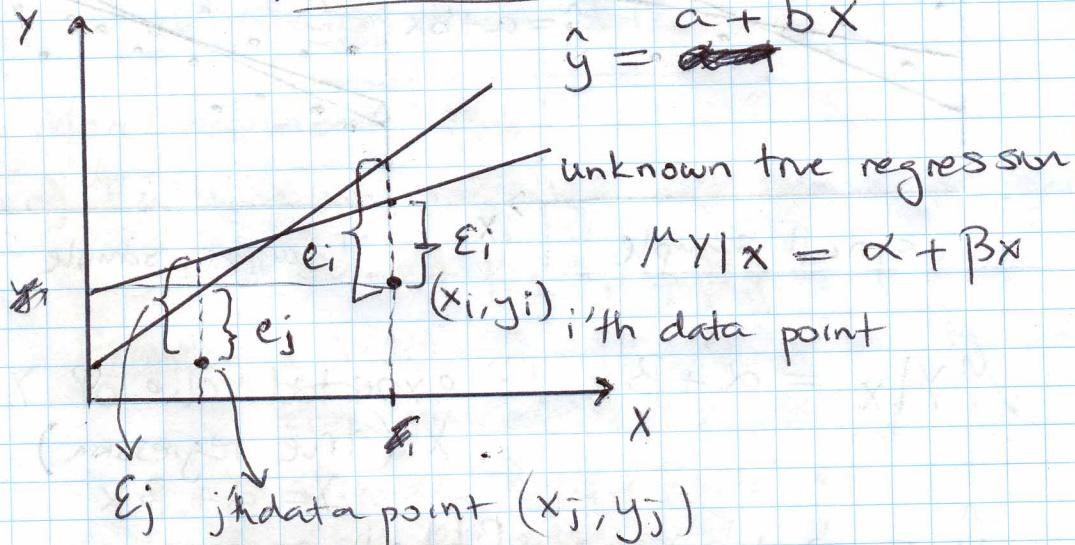
are the conceptual errors (realizations of the random var  $\epsilon$ )

## Residuals : error in fit

Given a set of data (sample)  $(x_i, y_i)$   $i=1$  to  $n$  and a fitted model  $\hat{y}_i = a + b x_i$  then the residual for the  $i$ 'th data point is

$$e_i = y_i - \hat{y}_i$$

$$\hat{y} = a + b x$$



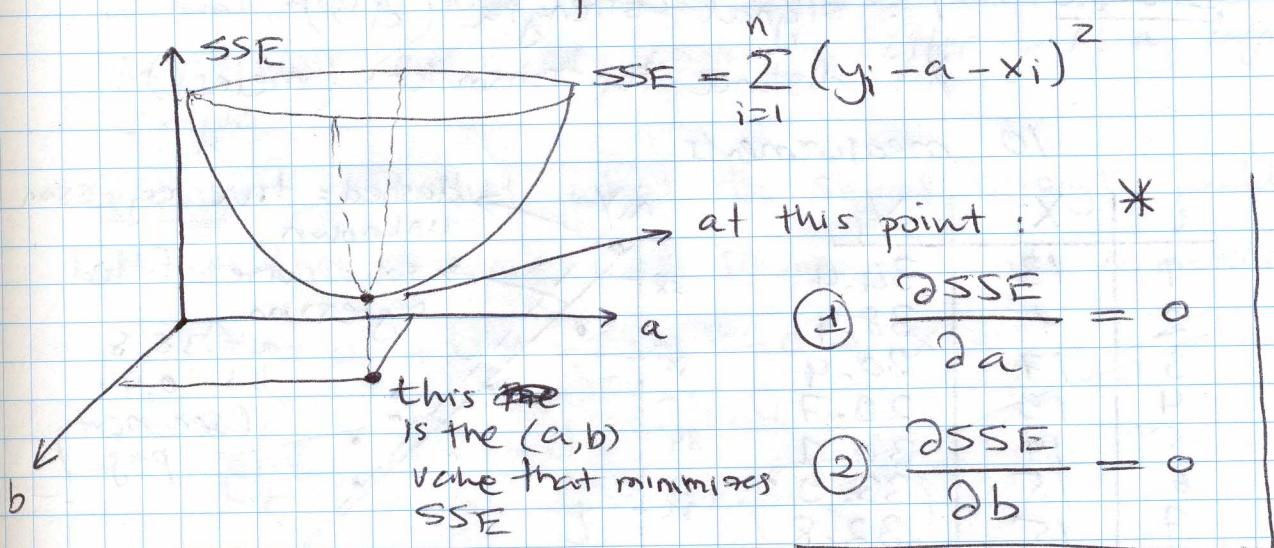
## Least squares fitting method

One way to choose reasonable values for  $a$  and  $b$  is to choose them to minimize the sum of squared residuals.

$$SSE(\text{sum of squared errors}) = \sum_{i=1}^n e_i^2 \quad \leftarrow \text{Note this is } e_i \text{ not } E_i$$

$$= \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - a - b x_i)^2$$

To minimize SSE with respect to  $a$  and  $b$ , from calculus we know that the partial derivatives of SSE with respect to  $a$  and  $b$  must be 0.



$$\textcircled{1} \quad \frac{\partial SSE}{\partial a} = 0$$

$$\textcircled{2} \quad \frac{\partial SSE}{\partial b} = 0$$

$$\frac{\partial SSE}{\partial a} = -2 \sum_{i=1}^n (y_i - a - b x_i) = 0$$

$$\frac{\partial SSE}{\partial b} = -2 \sum_{i=1}^n (y_i - a - b x_i) x_i = 0$$

Rearranging terms gives:

$$\textcircled{1} \quad a + b \sum_{i=1}^n x_i = \sum_{i=1}^n y_i$$

$$\textcircled{2} \quad a \sum_{i=1}^n x_i + b \sum_{i=1}^n x_i^2 = \sum_{i=1}^n x_i y_i$$

Solving ① and ② simultaneously gives the following formulas for  $a$  and  $b$ :

$$b = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

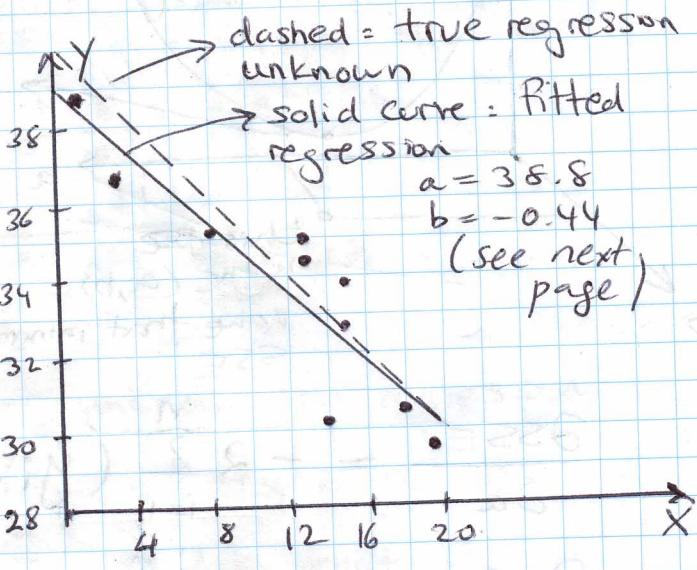
$$\text{and } a = \bar{y} - b \bar{x}$$

$$\text{where } \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

Example  $y$ : signal strength (dB)  
 $x$ : distance to transmitter (meters)

10 measurements

$i$	$x_i$	$y_i$
1	13	34.4
2	1	38.4
3	17	30.4
4	19	29.7
5	14	30.1
6	15	33.9
7	15	32.8
8	8	35.2
9	13	34.9
10	3	36.8



Notice that

for the same  $x$  value we have two different  $y$  values  
(non-deterministic)  
Same for  $x = 13$

## Least squares fitting

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = 11.8 \quad \bar{y} = 33.66$$

$$b = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{-137.58}{315.6} \approx -0.44$$

$$a = \bar{y} - b\bar{x} = 38.8 \Rightarrow \hat{y} = 38.8 - 0.44x$$

Normally, we wouldn't know the true values  $\alpha, \beta$ , but in this case, I generated the data myself according to  $y = 40 - 0.5x + \epsilon$  where  $\epsilon$  was normally distributed with mean 0 and variance 1. So  $\alpha$  was -0.5 and  $\beta$  was 40. Our estimates are quite close, but they could be better with a larger sample.

Question : Predict what the signal strength would be if the distance was 10 meters and 26 meters?

$$\hat{y} = a + bx = 38.8 - 0.44x$$

$$\text{so for } x=10 \quad \hat{y} = 34.4 \text{ dB}$$

$$\text{for } x=26 \quad \hat{y} = 27.36 \text{ dB}$$

Note : Sometimes we have control over for which  $x_i$  we make measurements. For instance, we could have measured signal strength  $y_i$  for  $x_i$  at regularly spaced intervals

- 2, 4, 6, 8, 10, 12, 14, 16, 18, 20 meters.