

**Chapter 8**  
**Hypothesis Testing**  
**And**  
**Linear Regression**

## Hypothesis Testing

A statistical hypothesis is an assertion or conjecture concerning one or more populations.

To prove that a hypothesis is true with absolute certainty (or to prove that it is false with absolute certainty), we would have to examine the entire population.

Instead, take a random sample from the population and use it as ~~evidence~~ evidence that either supports or does not support the hypothesis.

$H_0$  : Null hypothesis. The hypothesis we wish to test.

The rejection of  $H_0$  leads to the acceptance of an alternate hypothesis  $H_1$ . 2 possible outcomes:

- ① reject  $H_0$  : in favor of  $H_1$  because of sufficient evidence in the sample.
- ② fail to reject  $H_0$  : because of insufficient evidence

Note : There is no formal outcome that says accept  $H_0$ .  $H_0$  often represents "status quo" in opposition to a new idea  $H_1$ .

Example : Jury Trial     $H_0$  : defendant is innocent

$H_1$  : " " is guilty

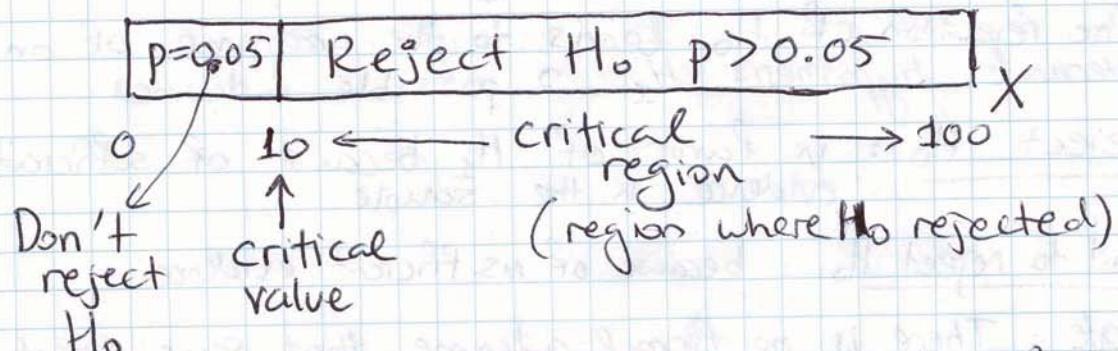
$H_0$  is rejected if  $H_1$  is supported by evidence beyond a reasonable doubt. Failure to reject  $H_0$  doesn't imply innocence, only that the evidence was not sufficient to reject it. Therefore, failure to reject  $H_0$  does not necessarily mean accept  $H_0$ .

Example: A company manufacturing RAM chips claims the defective rate of the population is 5%. Let  $p$  denote the true defective probability.

$$\begin{array}{l} H_0 : p = 0.05 \\ H_1 : p > 0.05 \end{array} \quad \left\{ \begin{array}{l} \text{Take a sample of 100} \\ \text{RAM chips off the production} \\ \text{line and test. } \cancel{H_0} \end{array} \right.$$

Let  $X$  denote the number of defectives in the sample of 100. Reject  $H_0$  if  $X \geq 10$   
(Note  $H_0$  is chosen arbitrarily in this case)

$X$  is called the test statistic



Q: Why did we choose a critical value of 10 in this example?

A: Since we know this is a Bernoulli Process, the expected value of defectives is  $np$ . So if we believe that  $p = 0.05$ , we should expect  $100 \times 0.05 = 5$  defectives in a sample of 100. Therefore, 10 defectives would be strong evidence that  $p > 0.05$ . Later we will learn how to choose the critical value based on the desired level of significance for the hypothesis test.

## Possible situations

	$H_0$ is true	$H_0$ is false
Do not reject $H_0$	Correct decision	Type II error
Reject $H_0$	Type I error	Correct decision

Defn: Rejection of  $H_0$  when it is true is called a Type I error.

Example: Convicting the defendant when he is innocent.

Defn: The probability of committing a type I error, denoted by  $\alpha$ , is called the level of significance.

Example (RAM chips continued)

$$\alpha = P(\text{Type I error}) = P(\underbrace{X \geq 10}_{\text{critical region}} \text{ when } p = 0.05) \quad H_0 \text{ is true.}$$

$$= \sum_{x=10}^{100} b(x; n=100, p=0.05) \quad \begin{matrix} \text{Binomial} \\ \text{Distribution} \end{matrix}$$

$$= \sum_{x=10}^{100} \binom{100}{x} 0.05^x 0.95^{100-x} = 0.0282$$

Level of significance is  $\alpha = 0.0282$ .

The lower the  $\alpha$ , the less likely we are to commit a type I error. Would like small  $\alpha$  values. (0.05 or smaller generally used)

Defn: Nonrejection of  $H_0$  when it is false is called a type II error. The probability of committing a type II error is denoted  $\beta$ .

Note =  $\beta$  is impossible to compute unless we have a specific alternate hypothesis.

Example continued : Can't compute  $\beta$  for

$H_1: p > 0.05$  but can compute it for testing

$H_0: p = 0.05$  against the alternate hypothesis that  $p = 0.1$  for instance.

$$\beta = P(\text{Type II error}) = P(X \leq 10 \text{ when } p=0.1)$$

$$= \sum_{x=0}^g b(x; n=100, p=0.1) = 0.4513$$

$\uparrow$   
 $H_1 (H_0 \text{ False})$

This high probability suggests we are likely to fail to reject  $H_0$  if the true  $p$  is 0.1.

This might be OK if we only want to make sure we don't fail to reject if the true  $p$  is 0.15

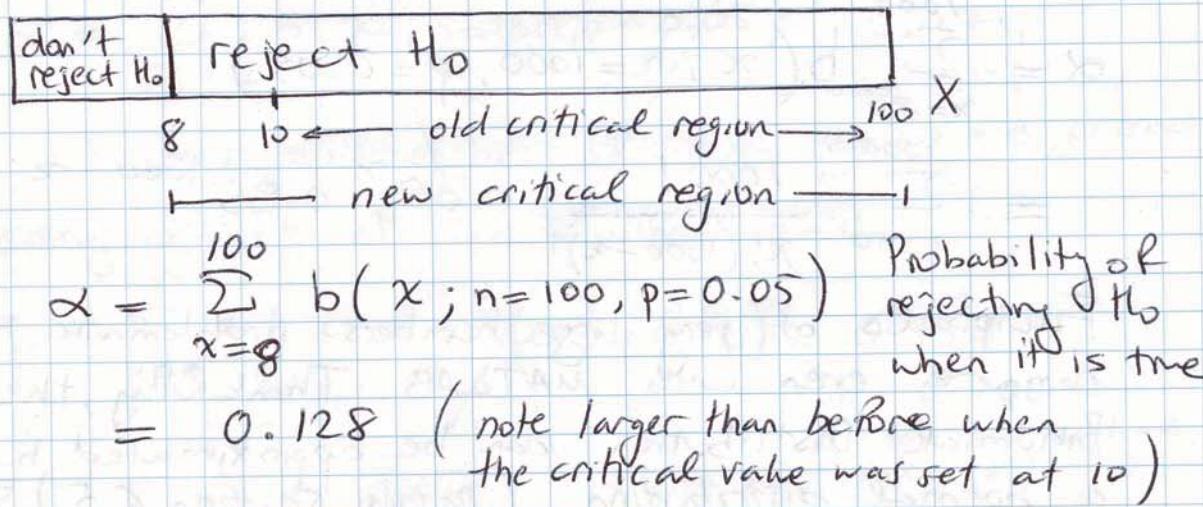
$$\beta = P(\text{Type II error}) = P(X \leq 10 \text{ when } p=0.15)$$

$$= \sum_{x=0}^g b(x; n=100, p=0.15) = 0.0551$$

## Effect of critical value

Moving the critical value provides a trade-off between  $\alpha$  and  $\beta$ . A reduction in  $\beta$  is always possible by ~~reducing~~ increasing the size of the critical region, but this increases  $\alpha$ .

Example continued : Lets change the critical value from 10 to 8. So reject  $H_0$  if  $X \geq 8$ .



testing against alternate hypothesis  $p=0.1$

$$\beta = \sum_{x=0}^7 b(x; n=100, p=0.1) \quad \text{Probability of not rejecting } H_0 \text{ when } p=0.1$$

$$= 0.206 \quad (\text{lower than before})$$

testing against alternate hypothesis  $p=0.15$

$$\beta = \sum_{x=0}^7 b(x; n=100, p=0.15) \quad \text{Probability of not rejecting } H_0 \text{ when } p=0.15$$

$$= 0.012 \quad (\text{again lower than before})$$

## Effect of sample size

Both  $\alpha$  and  $\beta$  can be reduced simultaneously by increasing the sample size.

Example continued : Sample size  $n=150$ , critical value 12. Reject  $H_0$  if  $X \geq 12$  ( $X$  is the # of defectives in sample of 150)

$$\alpha = \sum_{x=12}^{150} b(x; n=150, p=0.05) = 0.074$$

(was 0.128 for  $n=100$  and critical value 8)

testing against alternate hypothesis  $p=0.1$

$$\beta = \sum_{x=0}^{\text{---}} b(x; n=150, p=0.1) = 0.171$$

(was 0.206 for  $n=100$  and critical value 8)

Factorials of very large numbers problematic to compute accurately even with MATLAB. Thankfully, the Binomial Distribution can be approximated by the normal distribution - Details Section 6.5

Theorem : If  $X$  is a binomial random variable with  $n$  trials and probability of success of each trial  $p$ , then the limiting form of the distribution of  $Z = \frac{X - np}{\sqrt{np(1-p)}}$  as  $n \rightarrow \infty$  is the standard normal distribution.

Approximation is good when  $n$  is large and  $p$  is not extremely close to 0 or 1.

Lets recompute our  $\alpha$  with the normal approximation.

$$\begin{aligned}\alpha &= P(\text{Type I error}) = \sum_{x=12}^{150} b(x; n=150, p=0.05) \\ &\approx P\left(Z \geq \frac{12 - 150 \times 0.05}{\sqrt{150 \times 0.05 \times 0.95}}\right) = P(Z \geq 1.69) \\ &= 1 - P(Z \leq 1.69) = 1 - 0.9545 = 0.0455 \quad \text{not too bad.}\end{aligned}$$

What if we increase the sample size to  
 $n = 500$  and the critical value to 40?

The normal approximation should be even better  
since  $n$  larger

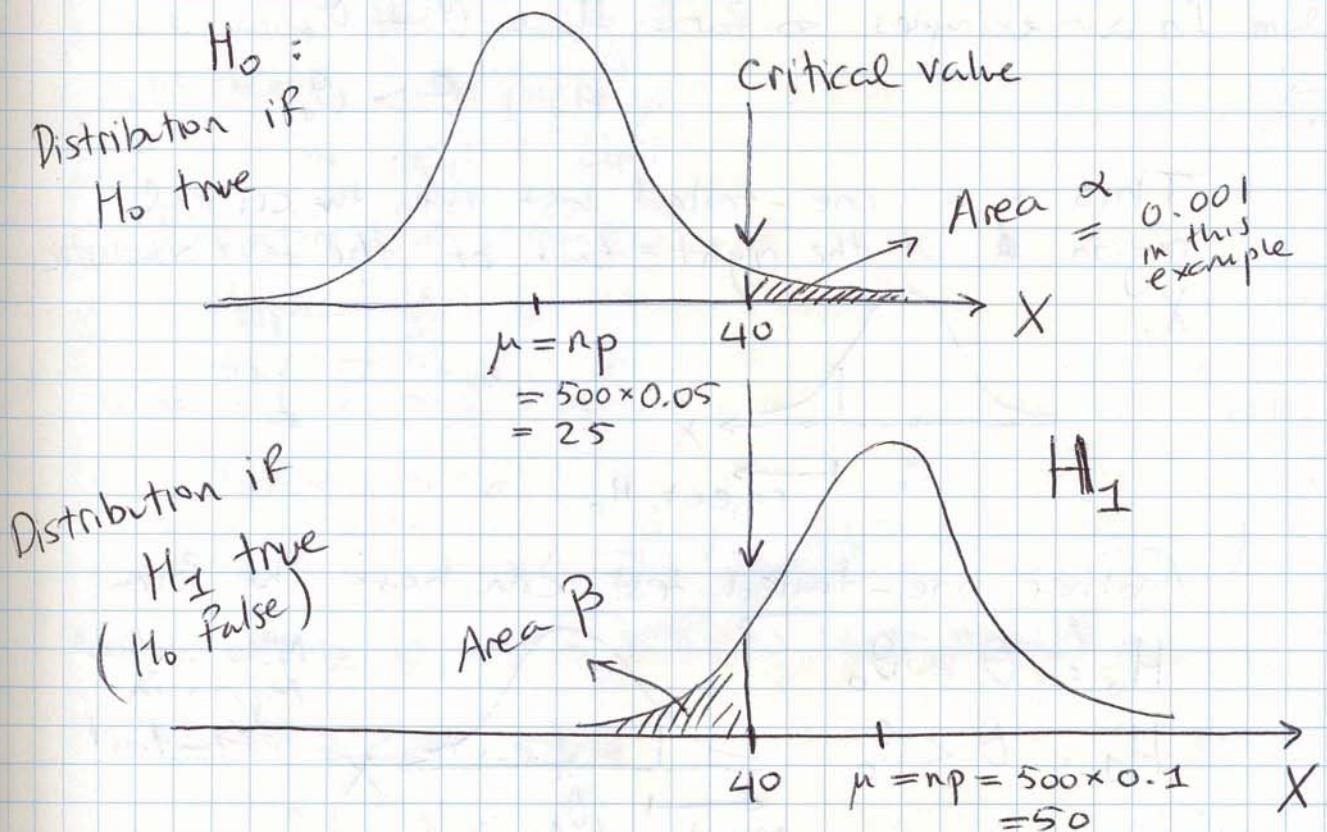
$$\begin{aligned}\alpha &\approx P\left(Z \geq \frac{40 - 500 \times 0.05}{\sqrt{500 \times 0.05 \times 0.95}}\right) = P(Z \geq 3.08) \\ &= 1 - P(Z \leq 3.08) = 1 - 0.999 = 0.001\end{aligned}$$

Very unlikely to commit type I error.

testing against alternate hypothesis  $p = 0.1$

$$\begin{aligned}\beta &= \sum_{x=0}^{39} b(x; n=500, p=0.1) \\ &\approx P\left(Z \leq \frac{39 - 500 \times 0.1}{\sqrt{500 \times 0.1 \times 0.9}}\right) = P(Z \leq -1.49) \\ &= 0.0681\end{aligned}$$

## Visual interpretation with normal approximation



Defn : The power of a test is the probability of rejecting  $H_0$  given that a specific alternate is true. Power =  $1 - \beta$ .

### ~~Properties of hypothesis testing~~

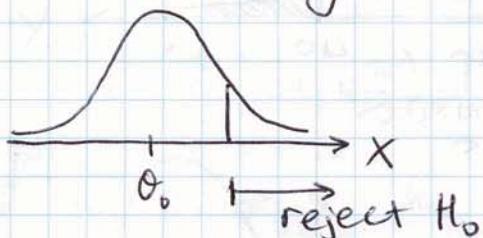
- ①  $\alpha$  and  $\beta$  are related. Decreasing one generally increases the other.
- ②  $\alpha$  can be set to a desired value by adjusting the critical value.  $\alpha$  typically set at 0.05 or 0.01.
- ③ Increasing  $n$  decreases  $\alpha$  and  $\beta$  both.
- ④  $\beta$  decreases as the distance between the true value and hypothesized ( $H_1$ ) value increases.

## One-tailed vs. two-tailed tests

In our examples so far  $H_0 : \theta = \theta_0$

$$H_1 : \theta > \theta_0$$

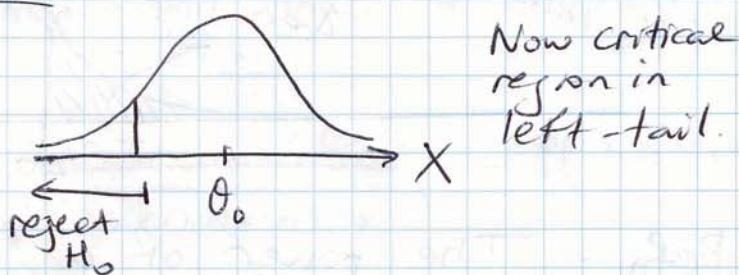
This is a one-tailed test with the critical region in the right-tail of the test statistic  $X$ .



Another one-tailed test can have the form

$$H_0 : \theta = \theta_0$$

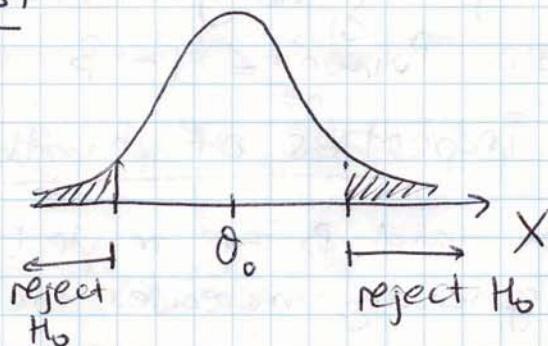
$$H_1 : \theta < \theta_0$$



## Two-sided test

$$H_0 : \theta = \theta_0$$

$$H_1 : \theta \neq \theta_0$$



Example: Production line of resistors that are supposed to be 100 ohms. Assume  $\sigma^2 = 8$

$$H_0 : \mu = 100$$

Let  $\bar{X}$  be the sample mean for a sample of size  $n = 100$

$$H_1 : \mu \neq 100$$

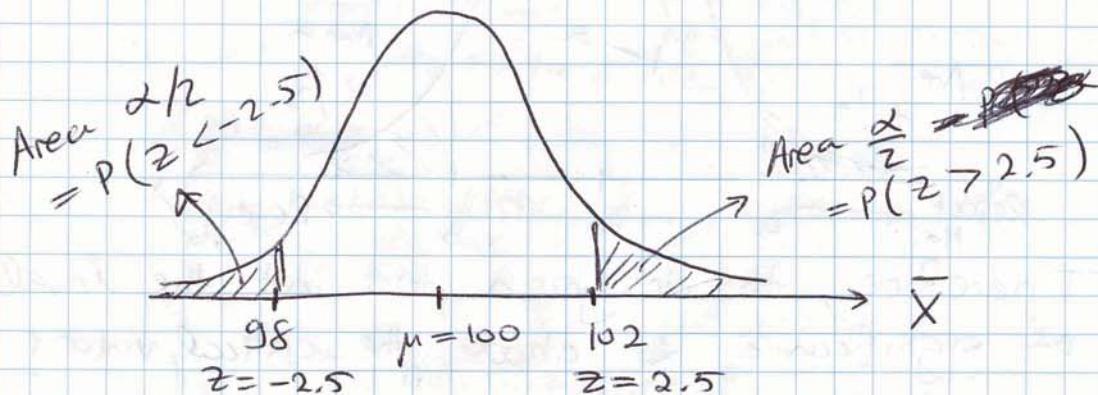
Reject $H_0$	Do not reject $H_0$	Reject $H_0$
98		102

The test statistic is the sample mean in this case.

We know the sampling distribution of  $\bar{X}$  is a normal distribution with mean  $\mu$  and standard deviation  $\sigma/\sqrt{n}$  due to the central limit theorem. ( $n=100 \geq 30$ )

Therefore we can compute the probability of a type I error as

$$\begin{aligned}\alpha &= P(\bar{X} < 98 \text{ when } \mu=100) + P(\bar{X} > 102 \text{ when } \mu=100) \\ &= P\left(Z < \frac{98-100}{8/\sqrt{100}}\right) + P\left(Z > \frac{102-100}{8/\sqrt{100}}\right) \\ &= P(Z < -2.5) + P(Z > 2.5) \\ &= P(Z < -2.5) + (1 - P(Z < 2.5)) \\ &= 2P(Z < -2.5) = 2 \times 0.0062 = 0.0124\end{aligned}$$



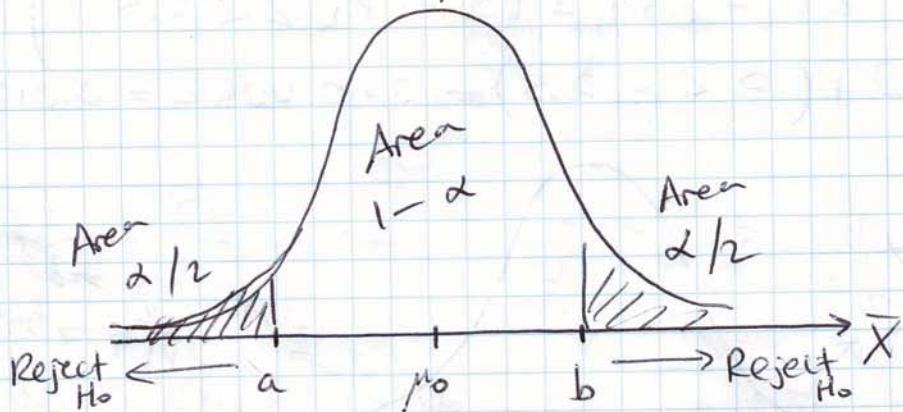
\* Testing  $H_0: \mu = \mu_0$  against  $H_1: \mu \neq \mu_0$  at a significance level  $\alpha$  is equivalent to computing a  $100(1-\alpha)\%$  confidence interval for  $\mu$  and rejecting  $H_0$  if  $\mu_0$  is outside the confidence interval.

## Tests concerning sample mean (Variance known)

$$\begin{array}{l} H_0: \mu = \mu_0 \\ H_1: \mu \neq \mu_0 \end{array} \quad \left\{ \begin{array}{l} \text{Sample } X_1, \dots, X_n \\ \text{Sample mean } \bar{X} \\ \text{Known population variance } \sigma^2 \end{array} \right.$$

Under  $H_0 \mu = \mu_0$  so  $P(\text{type I error})$  is computed using the sampling distribution of  $\bar{X}$  which is normal due to the central limit theorem with mean  $\mu$  and standard deviation  $\sigma/\sqrt{n}$ . From confidence intervals we know that

$$P\left(-z_{\alpha/2} < \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}} < z_{\alpha/2}\right) = 1 - \alpha$$



Therefore, to design a test at the level of significance  $\alpha$  choose the critical values  $a$  and  $b$  as

$$a = \mu_0 - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

$$b = \mu_0 + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

Then collect sample, compute the sample mean  $\bar{X}$  reject  $H_0$  if  $\bar{X} < a$  OR  $\bar{X} > b$

## Steps in hypothesis testing

- ① State the null and alternative hypothesis
- ② Choose a significance level  $\alpha$
- ③ Choose test statistic and establish critical region
- ④ Collect sample, compute test statistic on sample  
Reject  $H_0$  if test statistic in critical region.  
Otherwise do not reject  $H_0$

Example : A batch of 100 resistors have an average of 102 Ohms. Assuming a population standard deviation of 8 Ohms. Test whether the population mean is 100 Ohms at a significance level  $\alpha = 0.05$ .

- Soln = ①  $H_0 : \mu = 100$      $H_1 : \mu \neq 100$   
This is  $H_0$     (Note unless otherwise stated we use a two-tailed test)
- ②  $\alpha = 0.05$
  - ③ Test statistic  $\bar{X}$ , Reject  $H_0$  if  
 $\bar{X} < a$  or  $\bar{X} > b$

$$\begin{cases} z_{\alpha/2} = \\ z_{0.025} = 1.96 \end{cases}$$

$$a = \mu_0 - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} = 100 - 1.96 \frac{8}{10} = 98.432$$

$$b = \mu_0 + z_{\alpha/2} \frac{\sigma}{\sqrt{n}} = 100 + 1.96 \frac{8}{10} = 101.568$$

- ④  $\bar{X} = 102 > b$  therefore  
Reject  $H_0$ .

## One-sided tests of the sample mean

(A)  $H_0: \mu = \mu_0$  Reject  $H_0$  at significance level  $\alpha$

$$H_1: \mu > \mu_0 \quad \text{if } \bar{X} > \mu_0 + z_\alpha \frac{\sigma}{\sqrt{n}}$$

$\uparrow$

Note that  $z_\alpha$  appears instead of  $z_{\alpha/2}$  just like in one-tailed confidence intervals

(B)  $H_0: \mu = \mu_0$  Reject  $H_0$  at  $\alpha$  significance level

$$H_1: \mu < \mu_0 \quad \alpha \text{ if } \bar{X} < \mu_0 - z_\alpha \frac{\sigma}{\sqrt{n}}$$

Example: A quality control engineer finds that a sample of 100 light bulbs had an average life-time of 470 hours. Assuming a population standard deviation  $\sigma = 25$  hours, test whether the population mean is 480 hours vs. the alternative hypothesis  $\mu < 480$  at a significance level  $\alpha = 0.05$

Soln: ①  $H_0: \mu = 480$   $H_1: \mu < 480$  (one-tailed test)

②  $\alpha = 0.05$

③ Test statistic  $\bar{X}$ . Reject  $H_0$  if

$$\bar{X} < \mu_0 - z_\alpha \frac{\sigma}{\sqrt{n}} = 480 - 1.645 \frac{25}{\sqrt{10}}$$

$$\bar{X} < 475.9$$

④ Since  $\bar{X} = 470$  Reject  $H_0$

Note : Table A.3 has no entry with  $P(Z < z) = 0.05$  exactly. The closest ones are  $P(Z < -1.64) = 0.505$  and  $P(Z < -1.65) = 0.495$

The book (and in this example we) use the average of  $-1.64$  and  $-1.65$  as the  $z$ -value for which  $P(Z > z) = 0.05$ . Alternatively, we could just use whichever is closest to the  $\alpha$  value we are looking for.

### Tests concerning sample mean (Variance unknown)

$H_0: \mu = \mu_0$  } Sample  $X_1, \dots, X_n$  from a normal population. Unknown  $\sigma^2$ .

$H_1: \mu \neq \mu_0$  Sample mean  $\bar{X}$ , sample variance  $S^2$

We know that in this case the sampling distribution for  $\bar{X}$  is the t-distribution.

Critical region at significance level  $\alpha$  is

$$\bar{X} < a \text{ OR } \bar{X} > b \text{ reject } H_0$$

where  $a = \mu_0 - t_{\alpha/2} \frac{s}{\sqrt{n}}$  where  $t_{\alpha/2}$  is from

$$b = \mu_0 + t_{\alpha/2} \frac{s}{\sqrt{n}}$$
 row with  $v=n-1$

OR equivalently let  $T = \frac{\bar{X} - \mu_0}{s/\sqrt{n}}$

reject  $H_0$  if  $T < -t_{\alpha/2}$  or  $T > t_{\alpha/2}$   
for  $v=n-1$  degrees of freedom.

In one-sided test  $t_{\alpha/2}$  is replaced by  $t_\alpha$  as usual.

### Example (Example 10.5 textbook)

Claim: A certain electrical appliance (vacuum cleaner) expends 46 kWh per year. A random sample of 12 homes indicate that vacuum cleaners expend an average of 42 kWh per year with standard deviation 11.9 kWh. Does this suggest at the 0.05 level of significance that, on the average, vacuum cleaners expend less than 46 kWh per year? Assume population of kWh to be normally distributed.

$$\text{So In: } \textcircled{1} \quad H_0 = \mu = 46 \text{ kWh}$$

$$H_1 = \mu < 46 \text{ kWh}$$

$$\textcircled{2} \quad \alpha = 0.05$$

$$\textcircled{3} \quad \begin{aligned} \text{Test statistic } T &= \frac{\bar{X} - \mu_0}{s/\sqrt{n}} \end{aligned} \quad \left. \begin{array}{l} \text{reject } H_0 \text{ if} \\ \text{---} \\ T < -t_{0.05} \end{array} \right\}$$

for  $v = n - 1 = 11$   
degrees of freedom.  
 $t_{0.05} = 1.796$   
so reject if  
 $T < -1.796$

$$\textcircled{4} \quad \bar{X} = 42 \quad s = 11.9 \quad n = 12$$

$$T = \frac{42 - 46}{11.9/\sqrt{12}} = -1.16 \quad \text{Do not reject } H_0$$

## P-values

Preselection of significance level  $\alpha$  has its roots in the philosophy that making a type-I error should be controlled.

The P-value approach aims to give more information to the user, especially when the test-statistic is close to the critical region.

Defn : A P-value is the lowest level of significance at which the observed value of a test statistic is significant (rejects  $H_0$ ).

Example : A batch of 100 resistors have an average of 101.5 Ohms. Assuming a population standard deviation of 5 Ohms

- a) Test whether the population mean is 100 Ohms at level of significance  $\alpha = 0.05$
- b) Compute p-value.

Soln a)  $H_0: \mu = 100$   $H_1: \mu \neq 100$

Test statistic  $\bar{X}$ . Reject  $H_0$  if

$$\bar{X} < 100 - z_{0.025} \frac{\sigma}{\sqrt{n}} = 100 - 1.96 \times \frac{5}{10} \\ = 99.02$$

or

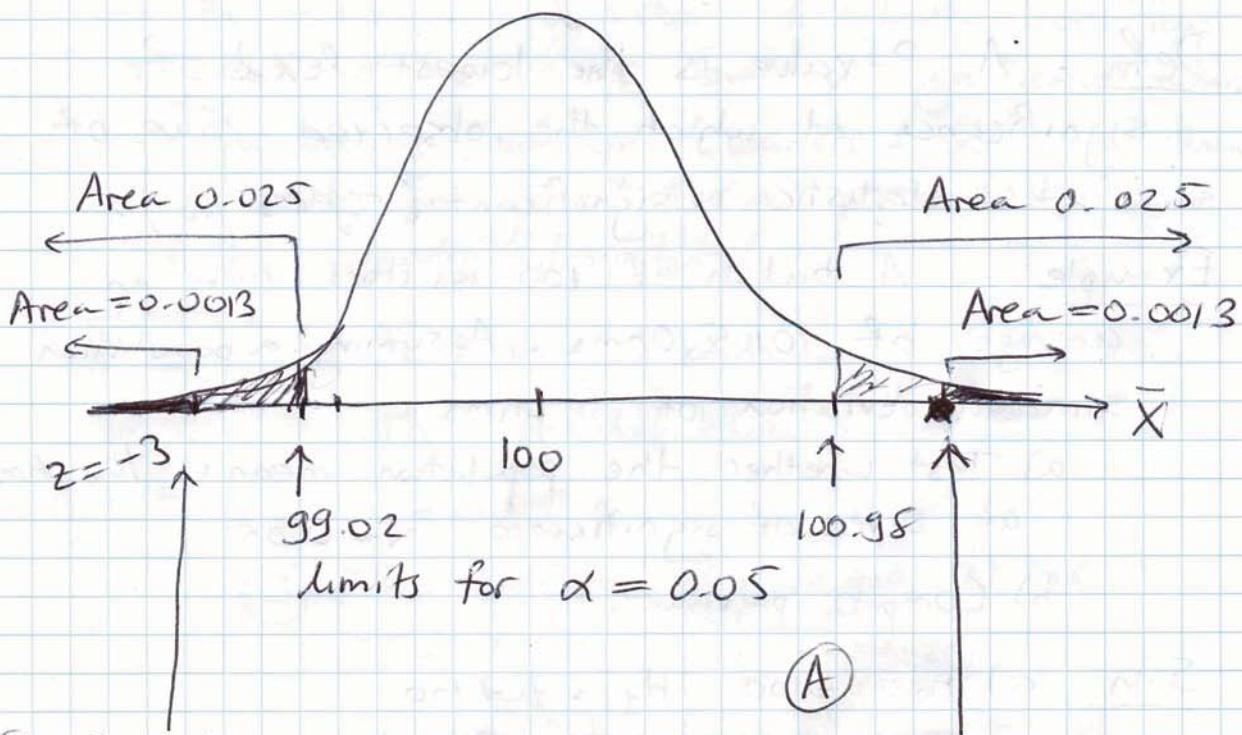
$$\bar{X} > 100 + z_{0.025} \frac{\sigma}{\sqrt{n}} = 100 + 1.96 \times \frac{5}{10} \\ = 100.98$$

$\bar{X} = 101.5$  therefore reject  $H_0$

b)  $Z = \frac{\bar{X} - 100}{\sigma/\sqrt{n}} = \frac{101.5 - 100}{5/10} = 3$  observed z-value

$$P = 2P(Z > 3 \text{ when } \mu = 100) = 2 \times 0.0013 = 0.0026$$

This means that  $H_0$  could have been rejected at significance level  $\alpha = 0.0026$  which is much stronger than rejecting it at  $\alpha = 0.05$ . Hence, the P-value gives more information than rejecting or not rejecting  $H_0$  at a fixed  $\alpha$ .



(B) Since this is a two-tailed test, the critical value on this side mirrors the right tail

Observed  $X = 101.5 \leftarrow z = 3$   
very rare event

Could have moved critical value all the way here and still reject  $H_0$

(C) Therefore

$$P = 2 \times 0.0013 = 0.0026$$

## 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 :

Simple Linear Regression

$$Y = \alpha + \beta X + \varepsilon \leftarrow \begin{array}{l} \text{Random variable} \\ E[\varepsilon] = 0 \end{array}$$

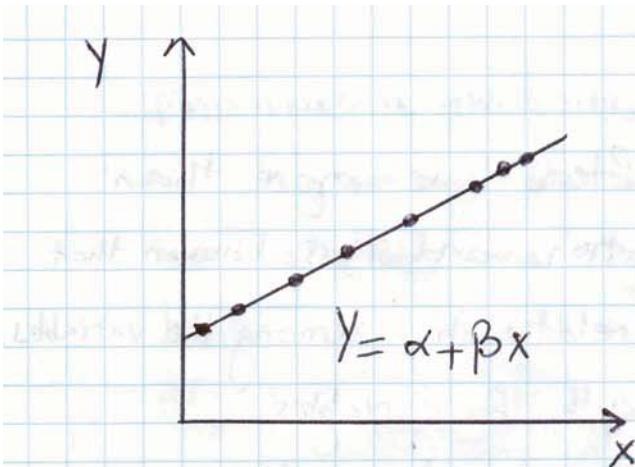
Dependent variable

Unknown intercept

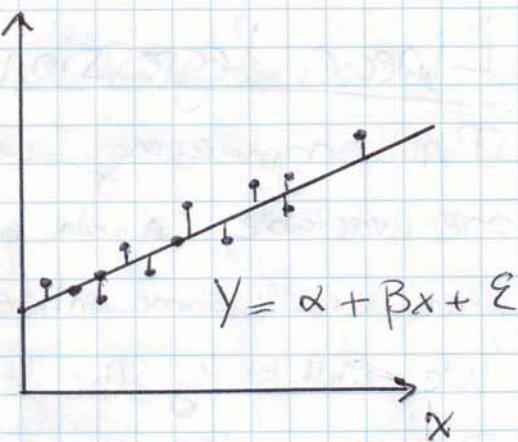
Unknown slope

Regressors

$\sigma^2$  Variance



A deterministic relationship



A non-deterministic relationship

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:

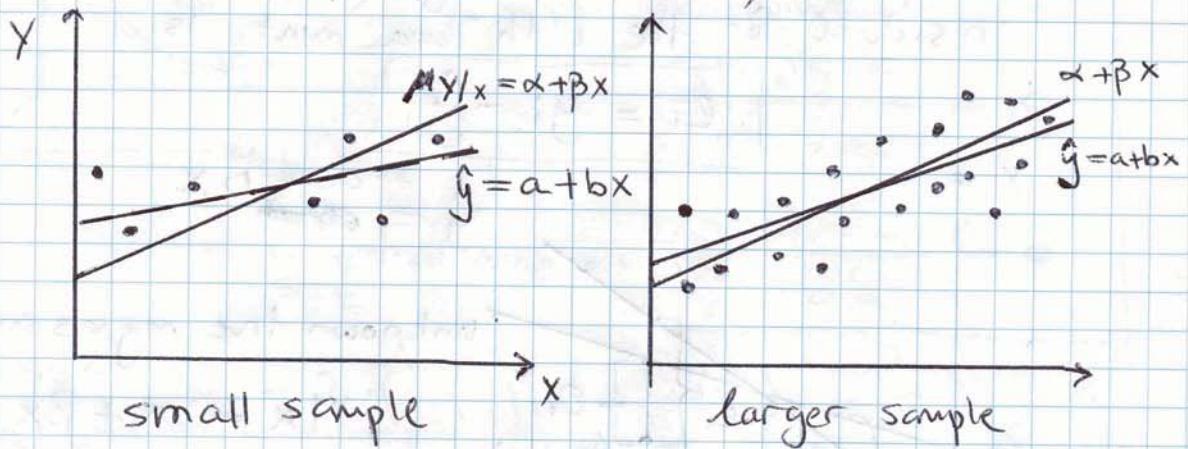
$$Y = \alpha + \beta X + \gamma X^2 + \epsilon$$

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 $\epsilon_i$

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

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

$$= \alpha + \beta X + \overline{E[\epsilon]}$$

○ assumption about  $\epsilon$

$$\mu_{Y|x} = E[Y] = \alpha + \beta x$$

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

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

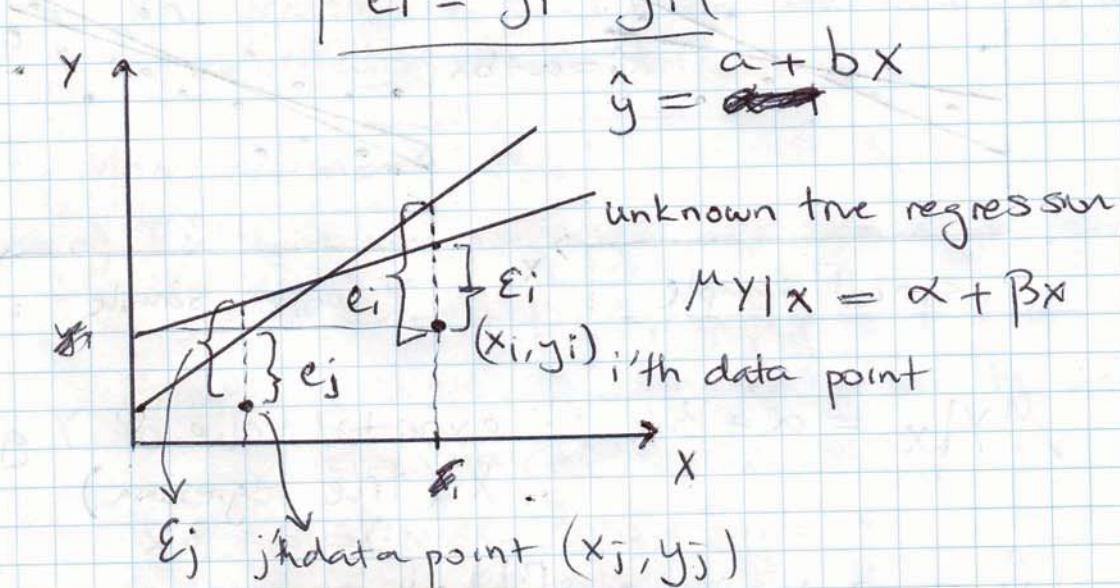
$\epsilon_i = y_i - (\alpha + \beta x_i)$  are the conceptual errors (realizations of the random var  $\epsilon$ )

### 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

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

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



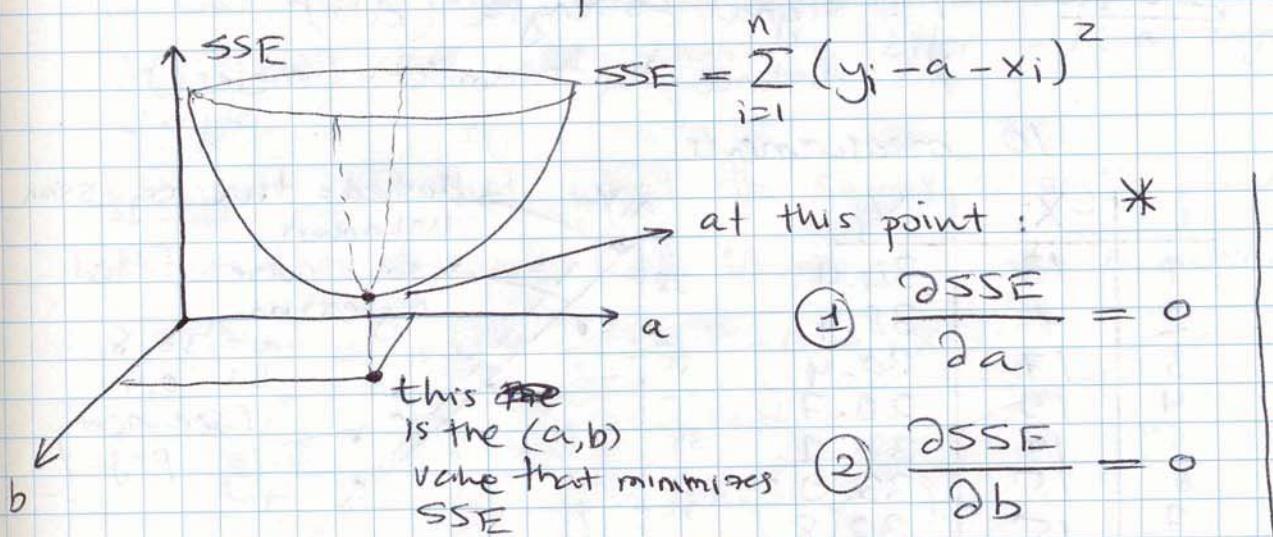
## Least squares fitting method

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.



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

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

Rearranging terms gives:

$$\begin{aligned} ① n a + b \sum_{i=1}^n x_i &= \sum_{i=1}^n y_i \\ ② a \sum_{i=1}^n x_i + b \sum_{i=1}^n x_i^2 &= \sum_{i=1}^n x_i y_i \end{aligned}$$

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

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

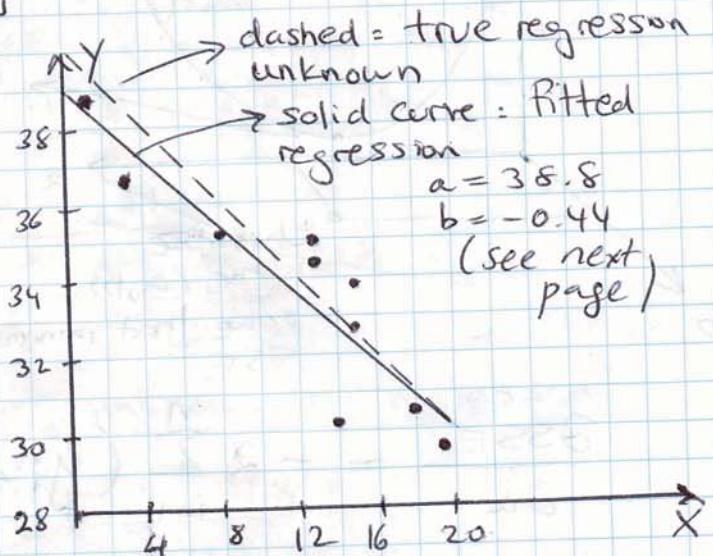
$$\text{and } a = \bar{y} - b \bar{x}$$

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

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

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

$$b = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{-137.58}{315.6} \approx -0.4359$$

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

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.