

The problem is thus  $\min_{\beta_0, \beta_1} \left\{ \sum_{i=1}^n [y_i - (\beta_0 + \beta_1 x_i)]^2 = f(\beta_0, \beta_1) \right\}$

$$\frac{\partial f}{\partial \beta_0} = \sum_{i=1}^n 2[y_i - (\beta_0 + \beta_1 x_i)]$$

$$\frac{\partial f}{\partial \beta_0} = 0 \Leftrightarrow \sum_{i=1}^n y_i - (\beta_0 + \beta_1 x_i) = 0$$

$$\Leftrightarrow \sum_{i=1}^n y_i = n\beta_0 + \beta_1 \sum_{i=1}^n x_i$$

$$\frac{\partial f}{\partial \beta_1} = \sum_{i=1}^n 2[y_i - (\beta_0 + \beta_1 x_i)] x_i$$

$$\frac{\partial f}{\partial \beta_1} = 0 \Leftrightarrow \sum_{i=1}^n [y_i - (\beta_0 + \beta_1 x_i)] x_i = 0$$

$$\Leftrightarrow \sum_{i=1}^n x_i y_i = \beta_0 \sum_{i=1}^n x_i + \beta_1 \sum_{i=1}^n x_i^2$$

After some algebra we get the following results for  $\beta_0$  and  $\beta_1$ :

$$\hat{\beta}_1 = \frac{n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \sum_{i=1}^n y_i}{n \sum_{i=1}^n x_i^2 - \left( \sum_{i=1}^n x_i \right)^2} \quad \text{and}$$

$$\hat{\beta}_0 = \frac{\sum_{i=1}^n y_i \sum_{i=1}^n x_i^2 - \sum_{i=1}^n x_i \sum_{i=1}^n x_i y_i}{n \sum_{i=1}^n x_i^2 - \left( \sum_{i=1}^n x_i \right)^2} = \frac{\sum_{i=1}^n y_i - \beta_1 \sum_{i=1}^n x_i}{n}$$

Example. 4 bulbs are tested for relationship between lifetime and operating voltage.

$i$	1	2	3	4
$V_i x_i$	105	110	115	120
$t_i y_i$	2000	1700	1120	950

To find  $\beta_0$  and  $\beta_1$  we need:

$$\sum_{i=1}^n x_i = 450$$

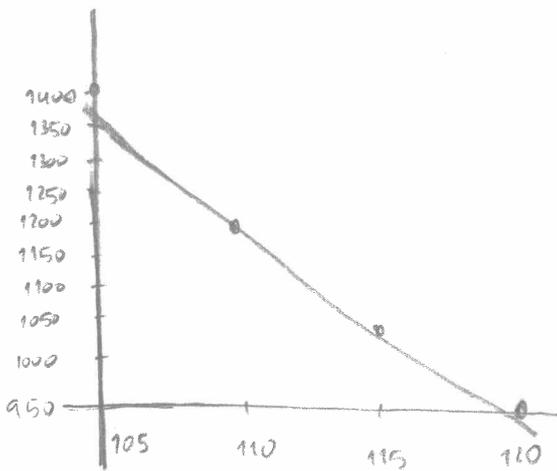
$$\sum_{i=1}^n x_i^2 = 50750$$

$$\sum_{i=1}^n y_i = 4670$$

$$\sum_{i=1}^n x_i y_i = 521800$$

$$\hat{\beta}_1 = \frac{n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \cdot \sum_{i=1}^n y_i}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i\right)^2} = \frac{4 \cdot 521800 - 450 \cdot 4670}{4 \cdot 50750 - (450)^2} = -28.6$$

$$\hat{\beta}_0 = \frac{\sum_{i=1}^n y_i - \beta_1 \sum_{i=1}^n x_i}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i\right)^2} = \frac{4670 + 28.6 \cdot 450}{4} = 4385$$



The regression line equation is

$$y = 4385 - 28.6x$$

As  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$  are the observed values of a sample

,  $\beta_0$  and  $\beta_1$  are indeed random variables and

estimators of the real population parameters  $\beta_0, \beta_1$  of  $\beta_0, \beta_1$  in  $y = \beta_0 + \beta_1 x$  (the real population regression line). If we can find the distribution

of  $\hat{\beta}_0$  and  $\hat{\beta}_1$ , then we can compute confidence intervals for  $\beta_0$  and  $\beta_1$  and test hypothesis on different values. To do that we need a probabilistic approach to the problem. Assume the sample  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$  and the model

$y_i = \beta_0 + \beta_1 x_i + e_i$  where  $i = 1, 2, \dots, n$  and the  $e_i$  are random noise such that they are independent,  $E[e_i] = 0$  and  $\text{Var}(e_i) = \sigma^2$ . we consider the  $x_i$  as fixed but the  $y_i$  are random because of the noise  $e_i$ .

Observation:  $y_i - (\beta_0 + \beta_1 x_i) = e_i$  so the least-squares method seeks to minimize the noise.

Theorem. Under the above assumptions  $E[\hat{\beta}_0] = \beta_0$  and  $E[\hat{\beta}_1] = \beta_1$

Proof. Under the statistical model

$$\hat{\beta}_0 = \frac{\sum_{i=1}^n x_i^2 \sum_{i=1}^n y_i - \sum_{i=1}^n x_i \sum_{i=1}^n x_i y_i}{n \sum_{i=1}^n x_i^2 - \left( \sum_{i=1}^n x_i \right)^2} \quad \text{then}$$

$$E[\hat{\beta}_0] = \frac{\sum_{i=1}^n x_i^2 \sum_{i=1}^n E[y_i] - \sum_{i=1}^n x_i \sum_{i=1}^n x_i E[y_i]}{n \sum_{i=1}^n x_i^2 - \left( \sum_{i=1}^n x_i \right)^2}$$

$$= \frac{\sum_{i=1}^n x_i^2 \sum_{i=1}^n (\beta_0 + \beta_1 x_i) - \sum_{i=1}^n x_i \sum_{i=1}^n x_i (\beta_0 + \beta_1 x_i)}{n \sum_{i=1}^n x_i^2 - \left( \sum_{i=1}^n x_i \right)^2}$$

$$= \frac{\sum_{i=1}^n x_i^2 \left( n\beta_0 + \beta_1 \sum_{i=1}^n x_i \right) - \sum_{i=1}^n x_i \left( \beta_0 \sum_{i=1}^n x_i + \beta_1 \sum_{i=1}^n x_i^2 \right)}{n \sum_{i=1}^n x_i^2 - \left( \sum_{i=1}^n x_i \right)^2}$$

$$= \frac{n\beta_0 \sum_{i=1}^n x_i^2 - \beta_0 \left( \sum_{i=1}^n x_i \right)^2}{n \sum_{i=1}^n x_i^2 - \left( \sum_{i=1}^n x_i \right)^2}$$

$$= \frac{\left[ n \sum_{i=1}^n x_i^2 - \left( \sum_{i=1}^n x_i \right)^2 \right] \beta_0}{n \sum_{i=1}^n x_i^2 - \left( \sum_{i=1}^n x_i \right)^2}$$

$$= \beta_0 \quad \text{And similarly for } \beta_1.$$

Theorem. Under the above assumptions

$$\text{Var}(\hat{\beta}_0) = \frac{\sigma^2 \sum_{i=1}^n x_i^2}{n \sum_{i=1}^n x_i^2 - \left( \sum_{i=1}^n x_i \right)^2} = \sigma_{\hat{\beta}_0}^2, \quad \text{Var}(\hat{\beta}_1) = \frac{n\sigma^2}{n \sum_{i=1}^n x_i^2 - \left( \sum_{i=1}^n x_i \right)^2} = \sigma_{\hat{\beta}_1}^2$$

Theorem. If  $\epsilon_i \sim N(0, \sigma^2)$   $i = 1, 2, \dots, n$  independent, then

$$\hat{\beta}_0 \sim N(\beta_0, \hat{\sigma}_{\hat{\beta}_0}^2) \quad \text{and} \quad \hat{\beta}_1 \sim N(\beta_1, \hat{\sigma}_{\hat{\beta}_1}^2)$$

If we know  $\sigma^2$  this would be enough to compute confidence intervals and test hypothesis, however that won't be always the case. We need to estimate  $\sigma^2$  from the data.

Define the residual sum of squares RSS as:

$$RSS = \sum_{i=1}^n [y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i)]^2$$

Proposition.  $s^2 = \frac{RSS}{n-2}$  is an unbiased estimate for  $\sigma^2$

(that is  $E[s^2] = \sigma^2$ ) and

$$s_{\hat{\beta}_0}^2 = \frac{s^2 \sum_{i=1}^n x_i^2}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i\right)^2}$$

$$s_{\hat{\beta}_1}^2 = \frac{n s^2}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i\right)^2}$$

are estimates of  $\hat{\sigma}_{\hat{\beta}_0}^2$  and  $\hat{\sigma}_{\hat{\beta}_1}^2$  respectively

Proposition. Under the normality assumptions

$$\frac{\hat{\beta}_0 - \beta_0}{s_{\hat{\beta}_0}} \sim t_{(n-2)} \quad \text{and} \quad \frac{\hat{\beta}_1 - \beta_1}{s_{\hat{\beta}_1}} \sim t_{(n-2)}$$

An important test is that  $H_0: \beta_1 = 0$  because if  $\beta_1 = 0$  then  
 $H_1: \beta_1 \neq 0$

$y_i = \beta_0 + \epsilon_i$ , that is, the dependent variable isn't really dependent on the prediction variable, but random noise about the horizontal line  $y = \beta_0$ .

In general, with help of the test statistic  $T = \frac{\hat{\beta}_1 - \beta_1^0}{s_{\hat{\beta}_1}}$  we can

test the null hypothesis  $H_0: \beta_1 = \beta_1^0$  at a significance level as follows

$H_1$	Reject $H_0$ if
$\beta_1 \neq \beta_1^0$	$ T_{obs}  > t_{(n-2)}(1 - \frac{\alpha}{2})$
$\beta_1 > \beta_1^0$	$T_{obs} > t_{(n-2)}(1 - \alpha)$
$\beta_1 < \beta_1^0$	$T_{obs} < t_{(n-2)}(\alpha)$

A two-sided confidence interval at  $(100(1-\alpha)\%)$  confidence level for  $\hat{\beta}_1$  is given by  $[\hat{\beta}_1 - t_{(n-2)}(1-\frac{\alpha}{2}) S_{\hat{\beta}_1}, \hat{\beta}_1 + t_{(n-2)}(1-\frac{\alpha}{2}) S_{\hat{\beta}_1}]$

Example. In our previous example of 4 bulbs tested for lifetimes against operating voltage we got  $\sum x_i = 450$   $\sum y_i = 4670$   
 $\sum x_i^2 = 50750$   $\sum x_i y_i = 521800$

Test  $H_0: \beta_1 = 0$   
 $H_1: \beta_1 \neq 0$

First we compute  $S^2 = \frac{RSS}{n-2} = \frac{1}{2} \sum_{i=1}^n (y_i - 4385 + 28.6x_i)^2 = 1215$

$$S_{\hat{\beta}_1}^2 = \frac{nS^2}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2} = \frac{4 \cdot 1215}{4 \cdot 50750 - (450)^2} = 9.72$$

$$T_{obs} = \frac{\hat{\beta}_1 - \beta_1^0}{S_{\hat{\beta}_1}} = \frac{-28.6}{\sqrt{9.72}} = -9.1735 \Rightarrow |T_{obs}| = 9.1735$$

$$\text{At } \alpha = 0.05 \quad t_{(n-2)}(1 - \frac{\alpha}{2}) = t_{(2)}(0.975) = 4.303$$

as  $|T_{obs}| > 4.303$  we reject  $H_0$ . A 95% confidence interval is given by

$$[\hat{\beta}_1 - t_{(n-2)}(1-\frac{\alpha}{2}) S_{\hat{\beta}_1}, \hat{\beta}_1 + t_{(n-2)}(1-\frac{\alpha}{2}) S_{\hat{\beta}_1}] = [-28.6 - 4.303\sqrt{9.72}, -28.6 + 4.303\sqrt{9.72}]$$

$$= [-42.02, -15.19]$$

which confirms that  $H_0: \beta_1 = 0$  must be rejected as 0 is not contained in the confidence interval.