

NORMAL ERROR REGRESSION MODEL

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$$

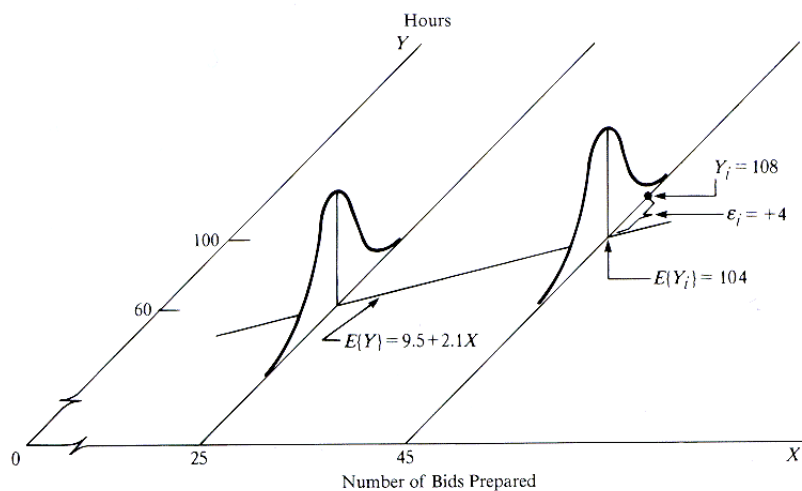
The random error term  $\varepsilon_i$  has:

- Mean:  $E(\varepsilon_i) =$
- Variance:  $\sigma^2(\varepsilon_i) =$

AND  $\varepsilon_i$

- - Are independent  $\forall i, j$  such that  $i \neq j$
- (The size of the error term for each trial has no effect on the size for any other.)
- - The model implies that  $Y_i$  are also independent normal random variables (as shown below).

FIGURE 1.6 Illustration of Simple Linear Regression Model (1.1).



<http://www.nielsen.sites.oasis.unc.edu/soci708/m15/m1005.gif>

Note: Both distributions shown are **Normally** distributed with:

- Mean:  $E(Y_i) =$
- $E(Y|X = 25) =$
- $E(Y|X = 45) =$
- Variance:  $\sigma^2(Y_i) =$
- Spread is equal in both

\*\*These are two distributions we showed in last week's notes!

Note: Unless specified, this model is assumed in the remainder of Chapters 2-5!

**INFERENCE & SAMPLING DISTRIBUTIONS**

As reviewed in the first two weeks, methods of inference are based on sampling distributions of estimators.

For each sampling distribution discussed, we will look at:

- **Shape:**
- **Center:**
- **Spread:**

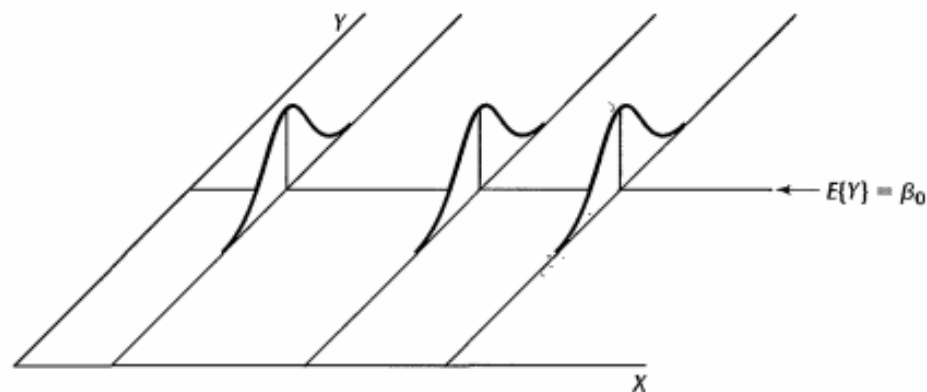
**INFERENCE ABOUT THE SLOPE,  $\beta_1$** 

Component	Interpretation
(slope) $\beta_1$	Change in the expected (average) value of Y per unit change in X

**Implications of  $\beta_1 = 0$** 

- 
- With the Normal Error Regression Model (shown below)
  - 
  -

**FIGURE 2.1**  
Regression  
Model (2.1)  
when  $\beta_1 = 0$ .



Sampling Distribution of  $b_1$ 

- $b_1 =$  is the *point estimator* of  $\beta_1$
- The sampling distribution of  $b_1$  means that we are looking at all the different values of  $b_1$  in repeated samples, *while holding the level of the predictor variable constant from sample to sample.*
- Assuming a Normal Error Regression Model, we have:
  - **SHAPE:**
  - **CENTER (MEAN):**  $E(b_1) =$
  - **SPREAD (VARIANCE):**  $\sigma^2(b_1) =$

**SHAPE: WHY NORMAL?** ( $b_1$  is a linear estimator of  $\beta_1$ )

- A linear combination of the  $Y_i$  means that  $Y_i$  is not multiplied by itself or another variable, but may be multiplied by constants, and combined by addition or subtraction.

If  $Y_1, Y_2, \dots, Y_n$  are independent normally distributed random variables, the linear combination  $a_1Y_1 + a_2Y_2 + \dots + a_nY_n$  is normally distributed.

•

- Let's define  $k_i = \frac{X_i - \bar{X}}{\sum_{i=1}^n (X_i - \bar{X})^2} = \frac{X_i - \bar{X}}{s_{xx}}$ .
- Then we can write the estimated slope as  $b_1 =$  .
- Some properties of  $k_i$ :
  - $\sum_{i=1}^n k_i =$
  - $\sum_{i=1}^n k_i X_i =$
  - $\sum_{i=1}^n k_i^2 =$

**CENTER (MEAN):** WHY IS  $E(b_1) = \beta_1$  ? ( $b_1$  is an unbiased estimator of  $\beta_1$ )

**SPREAD (VARIANCE):** WHY IS  $\sigma^2(b_1) = \frac{\sigma^2}{s_{xx}}$ ?

Recall:  $Y_i$  are independent random variables with constant variance  $\sigma^2(Y_i) = \sigma^2$ .

**ESTIMATED VARIANCE** of  $b_1$

Replacing  $\sigma^2$  with MSE we get:

$$s^2(b_1) =$$

- $s^2(b_1)$  is an unbiased estimator of  $\sigma^2(b_1)$ .
- The estimated standard deviation is  $s(b_1) =$  and is an unbiased estimator of  $\sigma(b_1)$ .
- Recall:

$$MSE = \frac{SSE}{n-2} = \frac{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}{n-2} = \frac{\sum_{i=1}^n e_i^2}{n-2}$$

Sampling Distribution of  $\frac{(b_1 - \beta_1)}{s(b_1)}$

- We know that the distribution of  $b_1$  is Normal. So, then we have that the standardized statistic
- As usual, we need to estimate the **unknown** value of  $\sigma(b_1)$  using  $s(b_1)$ . This introduces more uncertainty and so,
  - $\frac{(b_1 - \beta_1)}{s(b_1)}$  follows a ***t-distribution*** with  $df = n - 2$  (we are estimating **two** parameters,  $\beta_0$  and  $\beta_1$ , in our regression model)

CONFIDENCE INTERVAL FOR  $\beta_1$

lower and upper limits of the  $(1 - \alpha)\%$  confidence interval for  $\beta_1$  are:

How is this based on the sampling distribution of  $\frac{(b_1 - \beta_1)}{s(b_1)}$ ?

**Example:** Lean Body Mass (LBM) and Calorie Rate

In our regression model predicting a Calorie Rate (in calories per day) based on LBM (in kg), find a 98% confidence interval for the slope,  $\beta_1$ .

We need to know:

- t critical values for **C=0.98** (from R)
  - $\alpha =$                       and  $1 - \frac{\alpha}{2} =$                       ;  $df =$
- Values of the estimated slope ( $b_1$ ), error mean square (**MSE**), sum of squares for X ( $S_{xx}$ ), and estimated standard deviation of slope  $s(b_1)$ .

$x_i$ (LBM)	$y_i$ (Rate)	$x_i - \bar{x}$	$y_i - \bar{y}$	$(x_i - \bar{x})^2$	$(y_i - \bar{y})^2$	$(x_i - \bar{x})(y_i - \bar{y})$
36.1	995	-6.933333333	-240.0833333	48.07111111	57640.0069	1664.577778
54.6	1425	11.56666667	189.9166667	133.7877778	36068.3403	2196.702778
48.5	1396	5.466666667	160.9166667	29.88444444	25894.1736	879.6777778
42	1418	-1.033333333	182.9166667	1.067777778	33458.5069	-189.0138889
50.6	1502	7.566666667	266.9166667	57.25444444	71244.5069	2019.669444
42	1256	-1.033333333	20.91666667	1.067777778	437.506944	-21.61388889
40.3	1189	-2.733333333	-46.08333333	7.471111111	2123.67361	125.9611111
33.1	913	-9.933333333	-322.0833333	98.67111111	103737.674	3199.361111
42.4	1124	-0.633333333	-111.0833333	0.401111111	12339.5069	70.35277778
34.5	1052	-8.533333333	-183.0833333	72.81777778	33519.5069	1562.311111
51.1	1347	8.066666667	111.9166667	65.07111111	12525.3403	902.7944444
41.2	1204	-1.833333333	-31.08333333	3.361111111	966.173611	56.98611111
SUM		2.84217E-14	9.09495E-13	518.926667	389954.917	12467.76667

$$b_1 = \frac{S_{xy}}{S_{xx}} =$$

$$S_{xx} =$$

$$SSE = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \sum_{i=1}^n e_i^2 = \mathbf{90403.5241}$$

$$s^2 = \mathbf{MSE} =$$

$$s(b_1) = \sqrt{\frac{MSE}{S_{xx}}} =$$

The **lower** and **upper limits** are:

**T TESTS FOR  $\beta_1$** 

We test the null hypothesis  $H_0$ : .

- Represents ***no linear association*** between  $X$  and  $Y$ .

The alternative hypothesis, as usual, can be two-sided or one-sided.

- **Two-sided**

$$H_A:$$

- Represents ***a linear association*** between  $X$  and  $Y$ .

- **One-sided (lower-tail)**

$$H_A:$$

- Represents a ***negative*** slope (and thus negative association).

- **One-sided (upper-tail)**

$$H_A:$$

- Represents a ***positive*** slope (and thus positive association).

- The test statistic is:

$$t_0 =$$

- Note: This follows the usual formula of  $\frac{\text{Estimator} - \text{Hypothesized Parameter Value}}{\text{SD of Estimator}}$ .

**Example:** Lean body mass (LBM) and Calorie Rate

Is there a **linear association** between LBM and Calorie Rate? Use  $\alpha = 0.02$ .

$$H_0:$$

$$H_A:$$

From our earlier example, we had a 98% confidence interval:

Is there a **positive association** between LBM and Calorie Rate? Use  $\alpha = 0.02$ .

$H_0$ :                       $H_A$ :

The test statistic is:

$t_0 =$

The p-value is:

$p - value =$

### INFERENCE ABOUT THE INTERCEPT, $\beta_0$

Component	Interpretation
(Y-intercept) $\beta_0$	Mean value of Y when $X=0$ ( <b>meaningful when <math>X=0</math> is within the range of the model</b> )

### SAMPLING DISTRIBUTION OF $b_0$

- $b_0 =$                       is the **point estimator** of  $\beta_0$
- The sampling distribution of  $b_0$  means that we are looking at all the different values of  $b_0$  in repeated samples, **while holding the level of the predictor variable constant from sample to sample**.
- Assuming a Normal Error Regression Model, we have:
  - **SHAPE:**
  - **CENTER (MEAN):**  $E(b_0) =$
  - **SPREAD (VARIANCE):**  $\sigma^2(b_0) =$
  - 
  - **ESTIMATED VARIANCE:**  $S^2(b_0) =$
  - **ESTIMATED STANDARD DEVIATION:**  $S(b_0) =$



### Sampling Distribution of $\frac{(b_0 - \beta_0)}{s(b_0)}$

- We know that the distribution of  $b_0$  is Normal.
- As usual, we need to estimate the **unknown** value of  $\sigma(b_0)$  using  $s(b_0)$ .
- $\frac{(b_0 - \beta_0)}{s(b_0)}$  also follows a ***t-distribution*** with  $df = n - 2$

### Confidence Interval for $\beta_0$

The **lower** and **upper limits** of the  $(1 - \alpha)\%$  confidence interval for  $\beta_0$  are:

### Considerations for Inference About $\beta_0$ and $\beta_1$

- We are working with  $t$  distributions, so must remember that:
  - If the distributions of  $Y_i$  are not Normal:
    - Sampling distributions of  $b_0$  and  $b_1$  will be when there is not clear non-Normality in distributions of  $Y_i$ .
    - Sampling distributions of  $b_0$  and  $b_1$  *approach* Normal as sample size increases.
- In summary:

INTERVAL ESTIMATION OF PREDICTED VALUES,  $E(Y_h)$ Inference About  $E(Y_h)$ 

- **GOAL:**

- Let  $X_h =$

(may be an observed value of  $X$  occurring in the data or any value within the domain/scope of the model)

- $E(Y_h) =$

Sampling Distribution Of  $\hat{Y}_h$ 

- $\hat{Y}_h =$  is the **point estimator** of  $E(Y_h) =$
- The sampling distribution of  $\hat{Y}_h$  means that we are looking at all the different values of  $\hat{Y}_h$  in repeated samples, **while holding the level of the predictor variable constant from sample to sample.**
- Assuming a Normal Error Regression Model, we have:

- **SHAPE:**

(follows from fact that  $\hat{Y}_h$  is a linear combination of  $Y_i$ )

- **CENTER (MEAN):**  $E(\hat{Y}_h) =$

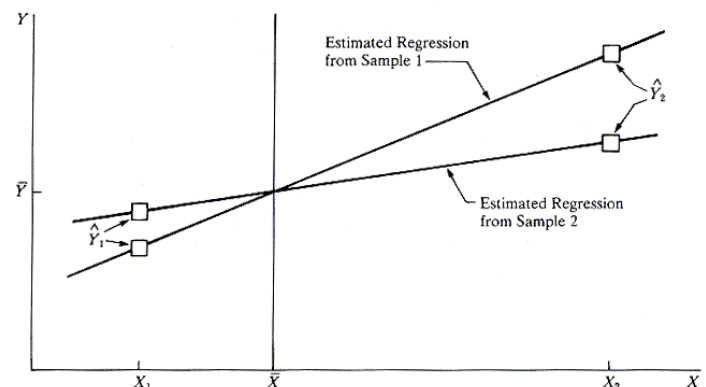
( $\hat{Y}_h$  is an unbiased estimator of  $E(Y_h)$ )

- **SPREAD (VARIANCE):**  $\sigma^2(\hat{Y}_h) =$

**Variation** in  $\hat{Y}_h$  values will be **greater** from sample to sample when  $X_h$  is **far from the mean** and **smaller**

when  $X_h$  is **closer to the mean**.

FIGURE 2.3 Effect on  $\hat{Y}_h$  of Variation in  $b_1$  from Sample to Sample in Two Samples with Same Means  $\bar{Y}$  and  $\bar{X}$ .



- Assuming a Normal Error Regression Model, we have:
  - ESTIMATED VARIANCE:  $S^2(\hat{Y}_h) =$
  - 
  - ESTIMATED STANDARD DEVIATION:  $S(\hat{Y}_h) =$

**Sampling Distribution of**  $\frac{(\hat{Y}_h - E(Y_h))}{s(\hat{Y}_h)}$

$\frac{(\hat{Y}_h - E(Y_h))}{s(\hat{Y}_h)}$  also follows a ***t-distribution*** with  $df = n - 2$

**Confidence Interval For**  $\hat{Y}_h$

The **lower** and **upper limits** of the  $(1 - \alpha)\%$  confidence interval for  $E(Y_h)$  are:

**Example:** Lean Body Mass (LBM) and Calorie Rate

In our regression model predicting a Calorie Rate (in calories per day) based on LBM (in kg), find a 96% confidence interval for the mean number of calories burned per day,  $E(Y_h)$ , for a woman with LBM of

$$X_h = 50 \text{ kg}.$$

We need to know:

- **t critical values** for **C=0.96** (from R)
  - $\alpha =$                       and  $1 - \frac{\alpha}{2} =$     ;  $df =$
- Values of the estimated/predicted mean ( $\hat{Y}_h$ ), error mean square (**MSE**), sum of squares for X ( $S_{xx}$ ), deviation from the mean of X ( $(X_h - \bar{X})^2$ ), and estimated standard deviation  $S(\hat{Y}_h)$ .
- The equation of the LSRL is:

$$\hat{Y} = 201.1616 + 24.0260666X_i$$

$$\bullet \quad \hat{Y}_h =$$

$$MSE = 9040.352 \quad S_{xx} = 518.926667$$

$$\bar{X} =$$

$$(X_h - \bar{X})^2 =$$

$$S(\hat{Y}_h) = \sqrt{MSE \left[ \frac{1}{n} + \frac{(X_h - \bar{X})^2}{S_{xx}} \right]} =$$

The lower and upper limits are: