## PREDICTION INTERVALS FOR NEW OBSERVATIONS, $Y$
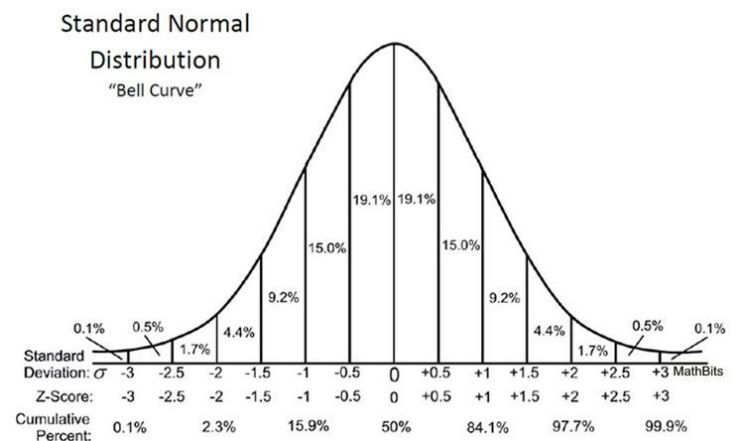
**Goal:**

Therefore, we need a **prediction _interval_**.

**PREDICTION INTERVAL:** PARAMETERS KNOWN

Assuming a Normal error regression model with known parameter values $\beta_0$, $\beta_1$, $E(Y_h)$, and $\sigma^2$,

the **lower and upper limits** of the

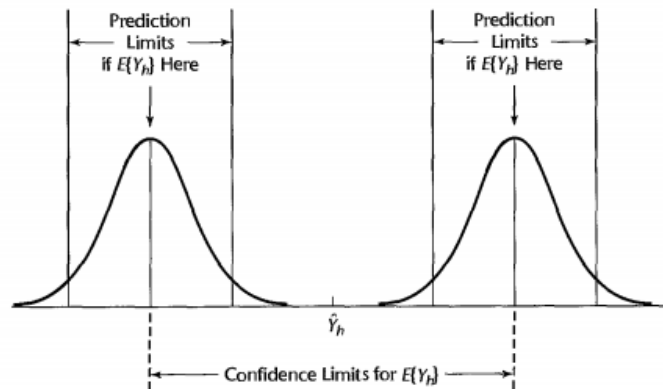$(1 - \alpha) \times 100\%$ **prediction interval** for $Y_{h(new)}$ are:



**PREDICTION INTERVAL:** PARAMETERS UNKNOWN

- The parameter values $\beta_0$, $\beta_1$, $E(Y_h)$, and $\sigma^2$ are unknown, so the mean and variance of the distribution of must be estimated.

    - **MEAN**:

    - **VARIANCE:**

**Two major considerations:**

1. Variation in the possible center of the distribution of $Y$ (since the center, $E(Y_h)$ , is estimated by a confidence interval based on $\widehat{Y_h}$ ).

2. Variation within the probability distribution of $Y$ (determined by prediction limits based on the estimated spread of the distribution of $Y$).

**FIGURE 2.5**
**Prediction of** $Y_{h(new)}$ **when Parameters Unknown.**

So, assuming a Normal error regression model, the **prediction limits for** $Y_{h(new)}$ **are based on** :

-      is a point estimate for the unknown mean $E(Y_h)$

- $pred =$          =prediction error:

    How far might the new observation vary from the estimated mean?

The lower and upper limits of the $(1 - \alpha) \times 100\%$ **prediction interval** for $Y_{h(new)}$ are:

- The **variance o**f the **prediction error**, $pred = (Y_{h(new)} - \widehat{Y_h})$, is:

$$\sigma^2(pred) =$$

- It is estimated, *unbiasedly*, by: $s^2(pred) =$

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

Previously, our regression model predicted Calorie Rate (in calories per day) based on LBM (in kg).  Suppose a new woman is added to the study whose lean body mass is $X_h = 50 \ kg$.  Find a 96% prediction interval for the calorie rate for this new subject.

We need to know:

- From previous example: $t$ **critical values** for **C=0.96** (from R), value of the estimated/predicted mean $(\widehat{Y}_h)$, value of error mean square($MSE$), and value of the estimated standard deviation $s(\widehat{Y}_h)$.

  - $\alpha = 1 - 0.96$=0.04 and $1 - \frac{\alpha}{2} = 0.98$

    $df = 12 - 2 = 10$→qt(**0.98**,10)→[1] 2.359315→   $\pm t(0.98; 10) =$

  - $\widehat{Y}_h$= **201.1616+ 24.0260666(50)=**

  - $MSE =$

  - $S(\widehat{Y}_h) = \sqrt{9040.352[\frac{1}{12} + \frac{48.53445}{518.926667}]} =$

- The value of $s(pred)$.

  - $s(pred) = \sqrt{MSE + s^2(\widehat{Y}_h)} =$

    **\*\*Note: The sample size is still 12 since our prediction is based on the original regression model.**

The **lower and upper limits** are:

**\*\*Note: This **interval is wider** than the interval found previously for the mean:

$$1308.125 \leq E(Y_h) \leq 1496.805$$

since the prediction interval incorporates the variability of estimated mean calorie rate from sample to sample as well as the variability of calorie rates from one woman to the next (for women with LBM=50 kg).

**Final Comments about Prediction Intervals for One Observation**

- As we saw in the last example, the **prediction interval** for the same level of X and same confidence level will always be **wider than** the **confidence interval** since prediction intervals must consider the variability (spread) of the distribution of $Y$ in addition to the sampling variability of $\widehat{Y}_h$.

- Unlike confidence limits, **prediction limits are sensitive to distinct non-normality** of the distribution of error terms.

- **Prediction limits** apply to **one prediction** based on the sample data.

- A **prediction interval** is **statement about** the **value taken on by a random variable**.

(This is not the same as a confidence interval which is intended to capture the true value of a parameter.)

## PREDICTION OF M NEW OBSERVATIONS FOR GIVEN X

**Goal:**


The **lower** and **upper limits** of the $(1 - \alpha) \times 100\%$ prediction interval are:


Where:

$$s^2(predmean) =$$

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

Suppose a 4 new women are added to the study with lean body mass is $X_h = 50\ kg$. Find a 96% prediction interval for the mean calorie rate for these new subjects.

- The value of **s(predmean)**.

$$s(predmean) = \sqrt{\frac{MSE}{m} + s^2(\widehat{Y}_h)} =$$

**The lower and upper limits** are:

**Note:** This **interval is wider** than the interval found previously for the mean:
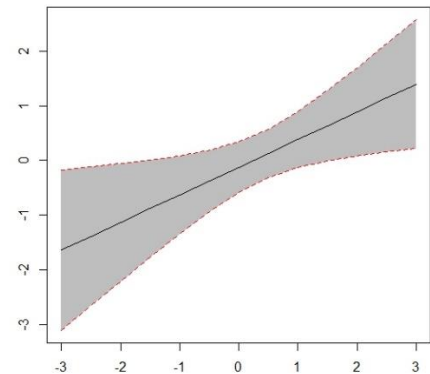
$$1308.125 \leq E(Y_h) \leq 1496.805$$

but **narrower** than the interval found previously for one new observation:

$$1159.109 \leq Y_{h(new)} \leq 1645.821$$
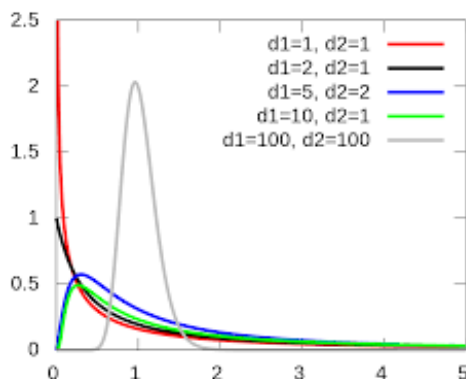
## CONFIDENCE BANDS FOR THE REGRESSION LINE

**Goal:**

This region is called a **confidence band**.

Assuming the Normal error regression model, the Working-Hotelling $(1-\alpha) \times 100\%$ **confidence band** for $E(Y_h) = \beta_0 + \beta_1 X_h$ at any level of $X_h$ has **boundary points**:

Where:

$$W^2 = \qquad\qquad \text{is used as the critical value rather than } t\left(1 - \frac{\alpha}{2}; n - 2\right).$$

**Example:** Lean Body Mass (LBM) and Calorie Rate-How accurately are we able to estimate the regression function? Find a 96% confidence band for the regression line.
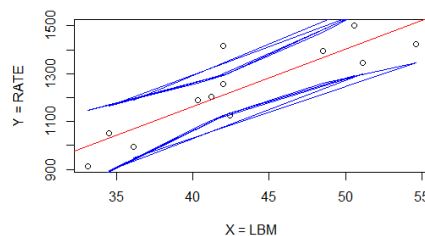
To show the calculation for $X_h = 50 \, kg$ we need to know:

- **W critical values** for **C=0.96** (based on $F$ and from R).

   - $C = 1 - \alpha = 0.96; \, df_2 = 12 - 2 = 10 \rightarrow$

The **boundary points** at level $X_h = 50 \, kg$ are:

**\*\*Note:** This **interval is wider** than the interval found previously for the mean:

$$1308.125 \leq E(Y_h) \leq 1496.805$$



The whole confidence band from R.

**Final Comments about Confidence Bands**

- **W is larger than t** and thus the **boundary points are farther apart** than confidence interval limits.

- Boundary points will still be farther apart the further $X_h$ is from the mean $\bar{X}$.

- $(1 - \alpha) \times 100\%$ represents

## ANOVA APPROACH TO REGRESSION ANALYSIS

This is a different approach to testing the same hypothesis, $H_0$:

- Based on partitioning the <u>sums of squares</u> and <u>degrees of freedom</u> associated with $Y$.

**Partitioning Total Variation**

$$SST \quad = \quad SSR \quad + \quad SSE$$

$$SST = \sum(Y_i - \overline{Y})^2 \qquad\qquad SSR = \sum(\widehat{Y}_i - \overline{Y})^2 \qquad\qquad SSE = \sum(Y_i - \widehat{Y}_i)^2$$

$$= b_1^2 s_{XX}$$

**Example: Lean body mass (LBM) and calorie rate**

| $x_i$ (LBM) | $y_i$ (Rate) | $\widehat{y}_i$ | $(y_i - \overline{y})^2$ | $(\widehat{y}_i - \overline{y})^2$ | $(y_i - \widehat{y}_i)^2$ |
|---|---|---|---|---|---|
| 36.1 | 995 | 1068.502604 | 57640.0069 | 27749.1393 | 5402.63283 |
| 54.6 | 1425 | 1512.984836 | 36068.3403 | 77229.24538 | 7741.33143 |
| 48.5 | 1396 | 1366.42583 | 25894.1736 | 17250.85146 | 874.631525 |
| 42 | 1418 | 1210.256397 | 33458.5069 | 616.3767578 | 43157.4045 |
| 50.6 | 1502 | 1416.88057 | 71244.5069 | 33050.23525 | 7245.31737 |
| 42 | 1256 | 1210.256397 | 437.506944 | 616.3767578 | 2092.4772 |
| 40.3 | 1189 | 1169.412084 | 2123.67361 | 4312.712992 | 383.686454 |
| 33.1 | 913 | 996.4244045 | 103737.674 | 56958.08433 | 6959.63126 |
| 42.4 | 1124 | 1219.866824 | 12339.5069 | 231.5421612 | 9190.44791 |
| 34.5 | 1052 | 1030.060898 | 33519.5069 | 42034.19911 | 481.32421 |
| 51.1 | 1347 | 1428.893603 | 12525.3403 | 37562.42073 | 6706.56225 |
| | | | | | |
| 41.2 | 1204 | -1.833333333 | 966.173611 | 1940.207752 | 168.077121 |
| | SUM | 14820.99999 | 389954.917 | 299551.392 | 90403.5241 |

**The Regression Sum of Squares:**

$$SSR = \sum_{i=1}^{n} (\widehat{Y}_i - \overline{Y})^2 =$$

**The Error Sum of Squares:**

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

**Total Sum of Squares:**

$$SST = \sum_{i=1}^{n}(Y_i - \overline{Y})^2 =$$

Why does $SST = SSR + SSE$ follow from $Y_i - \overline{Y} = (\widehat{Y}_i - \overline{Y}) + (Y_i - \widehat{Y}_i)$ ?

**Partitioning of Degrees of Freedom**

$$SST \quad = \quad SSR \quad + \quad SSE$$

**Mean Squares**

Previously we have worked with $MSE$, error mean square.

In general, a **mean square** is                    .

So, $MSR =$            $= SSR$ and $MSE =$        .

<div align="center">

**One-Way ANOVA Table for Regression**

</div>

| Source of Variation | Sum of Squares (SS) | $df$ | Mean Squares (MS) | $E(MS)$ |
|---|---|---|---|---|
| REGRESSION | $SSR$ $= \sum (\hat{Y}_i - \bar{Y})^2$ | 1 | $MSR = \dfrac{MSR}{1}$ | $\sigma^2 + \beta_1^2 (X_i - \bar{X})^2$ |
| ERROR | $SSE$ $= \sum\limits_{i=1}^{n} (Y_i - \hat{Y}_i)^2$ | $n-2$ | $MSE = \dfrac{SSE}{n-2}$ | $\sigma^2$ |
| TOTAL | $SST$ $= \sum (Y_i - \bar{Y})^2$ | $n-1$ | | |

**Expected Value of Mean Squares**

- $E(MSE) =$

  - Center(mean) of the sampling distribution of MSE is $\sigma^2$ even if $\beta_1 \neq 0$.

- $E(MSR) =$

  - Show that $MSR$ is an unbiased estimator of $\sigma^2 + \beta_1^2 s_{XX}$ .



  - When $\beta_1 = 0$ $E(MSR) = \sigma^2$ and when $\beta_1 \neq 0$ $E(MSR) > \sigma^2$ and hence $MSR$ will tend to be larger than $MSE$.

**F TEST for TESTING $H_0$: $\beta_1 = 0$ vs $H_A$: $\beta_1 \neq 0$.**

- The **test statistic** is:

$$F_0 =$$

   - Large values (far from 1) support $H_A$, so the test is really an                    .

   - Sampling distribution: Assuming $H_0$: $\beta_1 = 0$,                    .

- **Decision Rule:**

   - Reject $H_0$ if                    .

   - Otherwise, fail to reject $H_0$.

**Example: Lean body mass (LBM) and calorie rate**

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

| Source of Variation | Sum of Squares (SS) | $df$ | Mean Squares (MS) | $F-statistic$ |
|---|---|---|---|---|
| REGRESSION | SSR = 299551.392 | 1 | MSR = | $F_0 =$ |
| ERROR | SSE = 90403.5241 | $n-2$ | MSE = | |
| TOTAL | SST = 389954.917 | $n-1$ | | |

 (Note: The two-sided p-value for the previous part with the t test was $2 * P(t(10)) = 2 * 9.178927e\text{-}05 = 0.0001835785$.)
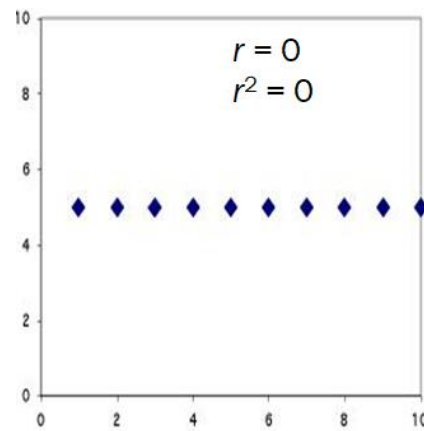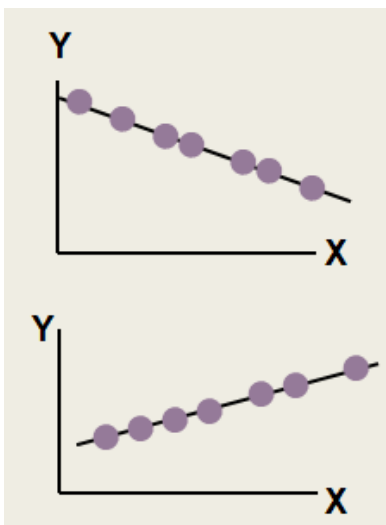
**Descriptive Measures of Linear Association**

- No single descriptive statistic can fully explain whether a particular regression model is useful.

    - Two common statistics:

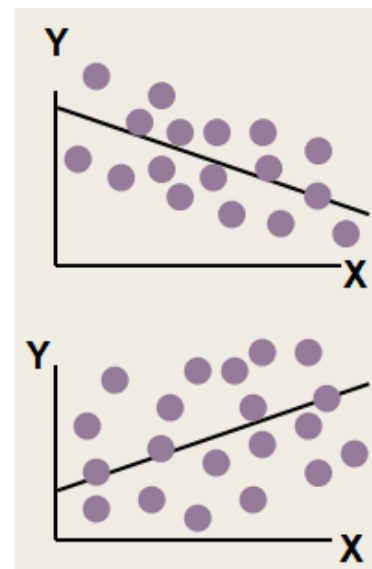        - **Coefficient of Determination**, $R^2$ (or $r^2$ )

        $$R^2 = (r)^2 =$$

The **coefficient of determination is** the proportion of the total variation in the response variable ($Y$) that is explained by variation in the explanatory variable ($X$), or the proportion of variation in Y that is explained by the regression model.

$R^2$=1                                         $0 < R^2 < 1$



- **Correlation Coefficient**, $r$

$$r =$$

    –   *The closer to –1, the stronger the negative linear relationship*

    –   *The closer to 1, the stronger the positive linear relationship*

    –   *The closer to 0, the weaker the linear relationship*