# MODULE THREE: DETERMINING CAUSE AND MAKING RELIABLE FORECASTS

# TOPIC 9: INTRODUCTION TO MULTIPLE REGRESSION

**+**

# Learning Objectives

At the completion of this topic, you should be able to:

- construct a multiple regression model and analyse model output
- differentiate between independent variables and decide which ones to include in the regression model, and determine which impendent variables are more important in predicting a dependent variable
- incorporate categorical and interactive variables in regression model
- detect collinearity

# **+The Multiple Regression Model**

**Idea:** Examine the linear relationship between
1 dependent (Y) and 2 or more independent variables ($X_i$)

Multiple Regression Model with k Independent Variables:

Y-intercept          Population slopes          Random Error

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \ldots + \beta_k X_{ki} + \varepsilon_i$$

# **+Multiple Regression Equation**

**Multiple regression equation with k independent variables:**

Estimated (or predicted) value of Y

Estimated intercept

Estimated slope coefficients

$$\hat{Y}_i = b_0 + b_1 X_{1i} + b_2 X_{2i} + \ldots + b_k X_{ki}$$

In this topic we will use Excel to obtain the regression slope coefficients and other regression summary measures

# +Pie Sales Example:

| Week | Pie Sales | Price ($) | Advertising ($100s) |
|------|-----------|-----------|---------------------|
| 1 | 350 | 5.50 | 3.3 |
| 2 | 460 | 7.50 | 3.3 |
| 3 | 350 | 8.00 | 3.0 |
| 4 | 430 | 8.00 | 4.5 |
| 5 | 350 | 6.80 | 3.0 |
| 6 | 380 | 7.50 | 4.0 |
| 7 | 430 | 4.50 | 3.0 |
| 8 | 470 | 6.40 | 3.7 |
| 9 | 450 | 7.00 | 3.5 |
| 10 | 490 | 5.00 | 4.0 |
| 11 | 340 | 7.20 | 3.5 |
| 12 | 300 | 7.90 | 3.2 |
| 13 | 440 | 5.90 | 4.0 |
| 14 | 450 | 5.00 | 3.5 |
| 15 | 300 | 7.00 | 2.7 |

A distributor of frozen dessert pies wants to evaluate factors thought to influence demand

**Dependent variable:**
Pie sales (units per week)
**Independent variables:**
Advertising ($100s), Price (in $)
Data are collected for 15 weeks

**Multiple regression equation:**

$$\widehat{Sales} = b_0 + b_1 (Price) + b_2 (Advertising)$$

# + Multiple Regression Output

| Regression Statistics | |
|---|---|
| Multiple R | 0.72213 |
| R Square | 0.52148 |
| Adjusted R Square | 0.44172 |
| Standard Error | 47.46341 |
| Observations | 15 |

$$Sales = 306.526 - 24.975(Price) + 74.131(Advertising)$$

| ANOVA | df | SS | MS | F | Significance F |
|---|---|---|---|---|---|
| Regression | 2 | 29460.027 | 14730.01 | 6.53861 | 0.01201 |
| Residual | 12 | 27033.306 | 2252.776 | | |
| Total | 14 | 56493.333 | | | |

| | Coefficients | Standard Error | t Stat | P-value | Lower 95% | Upper 95% |
|---|---|---|---|---|---|---|
| Intercept | 306.52619 | 114.25389 | 2.68285 | 0.01993 | 57.58835 | 555.46404 |
| Price | -24.97509 | 10.83213 | -2.30565 | 0.03979 | -48.57626 | -1.37392 |
| Advertising | 74.13096 | 25.96732 | 2.85478 | 0.01449 | 17.55303 | 130.70888 |

# +The Multiple Regression Equation

$$\widehat{Sales} = 306.526 - 24.975(Price) + 74.131(Advertising)$$

Where:

- Sales is in number of pies per week
- Price is in $
- Advertising is in $100s

**$b_1$ = -24.975:** sales will decrease, on average, by 24.975 pies per week for each $1 increase in selling price, net of the effects of changes due to advertising

**$b_2$ = 74.131:** sales will increase, on average, by 74.131 pies per week for each $100 increase in advertising, net of the effects of changes due to price

# +Using The Equation to Make Predictions

Predict sales for a week in which the selling price is $5.50 and advertising is $350:

$$\widehat{\text{Sales}} = 306.526 - 24.975(\text{Price}) + 74.131(\text{Advertising})$$
$$= 306.526 - 24.975\ (5.50) + 74.131\ (3.5)$$
$$= 428.62$$

**Predicted sales is 428.62 pies**

*Note:* Advertising is in $100s, so $350 means that $X_2 = 3.5$

# **+Coefficient of Multiple Determination**

Reports the proportion of total variation in Y explained by all X variables taken together

$$r^2 = \frac{SSR}{SST} = \frac{\text{regression sum of squares}}{\text{total sum of squares}}$$

# + Coefficient of Multiple Determination (Cont)

| Regression Statistics | |
|---|---|
| Multiple R | 0.72213 |
| R Square | 0.52148 |
| Adjusted R Square | 0.44172 |
| Standard Error | 47.46341 |
| Observations | 15 |

$$r^2 = \frac{SSR}{SST} = \frac{29460.0}{56493.3} = .52148$$

52.1% of the variation in pie sales is explained by the variation in price and advertising

| ANOVA | df | SS | MS | F | Significance F |
|---|---|---|---|---|---|
| Regression | 2 | 29460.027 | 14730.01 | 6.53861 | 0.01201 |
| Residual | 12 | 27033.306 | 2252.776 | | |
| Total | 14 | 56493.333 | | | |

| | Coefficients | Standard Error | t Stat | P-value | Lower 95% | Upper 95% |
|---|---|---|---|---|---|---|
| Intercept | 306.52619 | 114.25389 | 2.68285 | 0.01993 | 57.58835 | 555.46404 |
| Price | -24.97509 | 10.83213 | -2.30565 | 0.03979 | -48.57626 | -1.37392 |
| Advertising | 74.13096 | 25.96732 | 2.85478 | 0.01449 | 17.55303 | 130.70888 |

# +Adjusted r$^2$

$r^2$ never decreases when a new X variable is added to the model - this can be a disadvantage when comparing models

What is the net effect of adding a new variable?

- we lose a degree of freedom when a new X variable is added

- did the new X variable add enough explanatory power to offset the loss of one degree of freedom?

# **+Adjusted r² (Cont)**

Shows the proportion of variation in Y explained by all X variables adjusted for the number of X variables used

$$r^2_{adj} = 1 - \left[ (1 - r^2)\left( \frac{n-1}{n-k-1} \right) \right]$$

(where: n = sample size, k = number of independent variables)

- Penalises excessive use of unimportant independent variables
- Smaller than r²
- Useful in comparing among models

# +Adjusted r² (Cont)

| Regression Statistics | |
|---|---|
| Multiple R | 0.72213 |
| R Square | 0.52148 |
| Adjusted R Square | 0.44172 |
| Standard Error | 47.46341 |
| Observations | 15 |

$$r_{adj}^2 = .44172$$

44.2% of the variation in pie sales is explained by the variation in price and advertising, taking into account the sample size and number of independent variables

| ANOVA | df | SS | MS | F | Significance F |
|---|---|---|---|---|---|
| Regression | 2 | 29460.027 | 14730.01 | 6.53861 | 0.01201 |
| Residual | 12 | 27033.306 | 2252.776 | | |
| Total | 14 | 56493.333 | | | |

| | Coefficients | Standard Error | t Stat | P-value | Lower 95% | Upper 95% |
|---|---|---|---|---|---|---|
| Intercept | 306.52619 | 114.25389 | 2.68285 | 0.01993 | 57.58835 | 555.46404 |
| Price | -24.97509 | 10.83213 | -2.30565 | 0.03979 | -48.57626 | -1.37392 |
| Advertising | 74.13096 | 25.96732 | 2.85478 | 0.01449 | 17.55303 | 130.70888 |

# +Is the Model Significant?

F Test for Overall Significance of the Model

Shows if there is a linear relationship between all of the X variables considered together and Y

Hypotheses:

$H_0$: $\beta_1 = \beta_2 = \ldots = \beta_k = 0$  (no linear relationship)

$H_1$: at least one $\beta_i \neq 0$ (at least one independent variable affects Y)

# **+F Test for Overall Significance**

Test statistic

$$F = \frac{MSR}{MSE} = \frac{\dfrac{SSR}{k}}{\dfrac{SSE}{n-k-1}}$$

where F has: (numerator) = k, and

(denominator) = (n – k - 1) degrees of freedom

# +F Test for Overall Significance (Cont)

| Regression Statistics | |
|---|---|
| Multiple R | 0.72213 |
| R Square | 0.52148 |
| Adjusted R Square | 0.44172 |
| Standard Error | 47.46341 |
| Observations | 15 |

$$F = \frac{MSR}{MSE} = \frac{14730.0}{2252.8} = 6.5386$$

**With 2 and 12 degrees of freedom**

**P-value for the F Test**

| ANOVA | df | SS | MS | F | Significance F |
|---|---|---|---|---|---|
| Regression | 2 | 29460.027 | 14730.01 | 6.53861 | 0.01201 |
| Residual | 12 | 27033.306 | 2252.776 | | |
| Total | 14 | 56493.333 | | | |

| | Coefficients | Standard Error | t Stat | P-value | Lower 95% | Upper 95% |
|---|---|---|---|---|---|---|
| Intercept | 306.52619 | 114.25389 | 2.68285 | 0.01993 | 57.58835 | 555.46404 |
| Price | -24.97509 | 10.83213 | -2.30565 | 0.03979 | -48.57626 | -1.37392 |
| Advertising | 74.13096 | 25.96732 | 2.85478 | 0.01449 | 17.55303 | 130.70888 |

# +F Test for Overall Significance (Cont)

$H_0: \beta_1 = \beta_2 = 0$

$H_1: \beta_1$ and $\beta_2$ not both zero

$\alpha = .05$

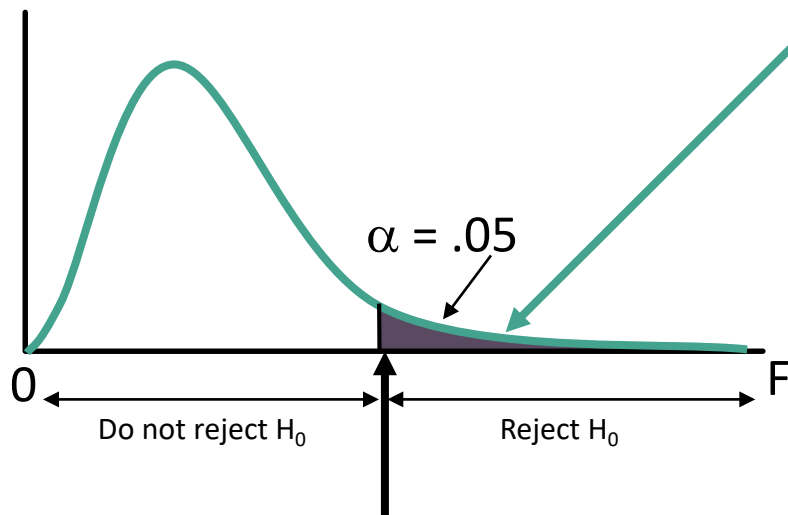$df_1 = 2 \qquad df_2 = 12$

**Test Statistic:**

$$F = \frac{MSR}{MSE} = 6.5386$$

**Decision:**
Since F test statistic is in the rejection region (p-value < .05), reject $H_0$

**Conclusion:**
There is evidence that at least one independent variable affects Y



$\alpha = .05$

0

Do not reject $H_0$

Reject $H_0$

F

**Critical Value: $F_\alpha = 3.885$**

# +Are Individual Variables Significant?

Shows if there is a linear relationship between the variable $X_j$ and Y

Hypotheses:

$H_0$: $\beta_j = 0$ (no linear relationship)

$H_1$: $\beta_j \neq 0$ (linear relationship does exist)

Use t tests of individual variable slopes (between $X_j$ and Y)

# +Are Individual Variables Significant? (Cont)

| Regression Statistics | |
|---|---|
| Multiple R | 0.72213 |
| R Square | 0.52148 |
| Adjusted R Square | 0.44172 |
| Standard Error | 47.46341 |
| Observations | 15 |

t-stat for Price is: t = -2.306, with p-value .0398

t-stat for Advertising is: t = 2.855, with p-value .0145

| ANOVA | df | SS | MS | F | Significance F |
|---|---|---|---|---|---|
| Regression | 2 | 29460.027 | 14730.01 | 6.53861 | 0.01201 |
| Residual | 12 | 27033.306 | 2252.776 | | |
| Total | 14 | 56493.333 | | | |

| | Coefficients | Standard Error | t Stat | P-value | Lower 95% | Upper 95% |
|---|---|---|---|---|---|---|
| Intercept | 306.52619 | 114.25389 | 2.68285 | 0.01993 | 57.58835 | 555.46404 |
| Price | -24.97509 | 10.83213 | -2.30565 | 0.03979 | -48.57626 | -1.37392 |
| Advertising | 74.13096 | 25.96732 | 2.85478 | 0.01449 | 17.55303 | 130.70888 |

# +Are Individual Variables Significant? (Cont)

**From Excel output:**

$H_0$: $\beta_i = 0$

$H_1$: $\beta_i \neq 0$

d.f. = 15-2-1 = 12
$\alpha$ = .05 $t_{\alpha/2}$ = 2.1788

|  | Coefficients | Standard Error | t Stat | P-value |
|---|---|---|---|---|
| Price | -24.97509 | 10.83213 | -2.30565 | 0.03979 |
| Advertising | 74.13096 | 25.96732 | 2.85478 | 0.01449 |

**Decision:**

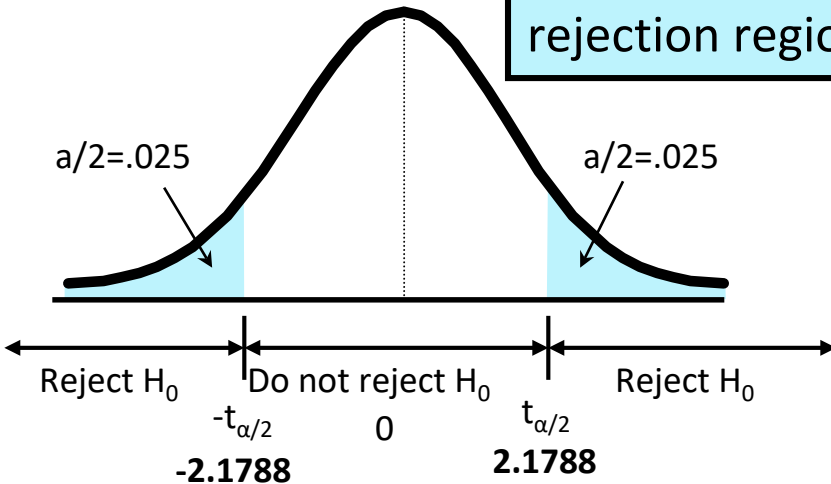The test statistic for each variable falls in the rejection region (p-values < .05)

a/2=.025

a/2=.025

Reject $H_0$    Do not reject $H_0$    Reject $H_0$

$-t_{\alpha/2}$   0   $t_{\alpha/2}$

**-2.1788**     **2.1788**

**Conclusion:**
Reject $H_0$ for each variable.
There is evidence that both Price and Advertising affect pie sales at $\alpha$ = .05

# +Confidence Interval Estimate for the Slope

Confidence interval for the population slope $\beta_j$

$$b_j \pm t_{n-k-1}S_{b_j}$$ Where t has: $(n - k - 1)$ d.f.

|  | Coefficients | Standard Error |
|---|---|---|
| Intercept | 306.52619 | 114.25389 |
| Price | -24.97509 | 10.83213 |
| Advertising | 74.13096 | 25.96732 |

Here, t has: $(15 - 2 - 1) = 12$ d.f.

**Example:** Form a 95% confidence interval for the effect of changes in price ($X_1$) on pie sales: $-24.975 \pm (2.1788)(10.832)$

So the interval is $(-48.576, -1.374)$
(This interval does not contain zero, so price has a significant effect on sales)

# + Confidence Interval Estimate for the Slope (Cont)

Confidence interval for the population slope $\beta_i$

| | Coefficients | Standard Error | ... | Lower 95% | Upper 95% |
|---|---|---|---|---|---|
| Intercept | 306.52619 | 114.25389 | ... | 57.58835 | 555.46404 |
| Price | -24.97509 | 10.83213 | ... | -48.57626 | -1.37392 |
| Advertising | 74.13096 | 25.96732 | ... | 17.55303 | 130.70888 |

**Example:** Excel output also reports these interval endpoints:

Weekly sales are estimated to be reduced by between 1.37 to 48.58 pies for each increase of $1 in the selling price

# **+Using Dummy Variables**

A dummy variable is a categorical explanatory variable with two levels:

- yes or no, on or off, male or female

- coded as 0 or 1

Regression intercepts are different if the variable is significant

Assumes equal slopes for other variables

If more than two levels, the number of dummy variables needed is number of levels minus 1

# +Dummy Variable Example (with 2 Levels):

$$\hat{Y} = b_0 + b_1 X_1 + b_2 X_2$$

Let: Y = pie sales

$X_1$ = price

$X_2$ = holiday (dummy variable)

($X_2$ = 1 if a holiday occurred during the week)
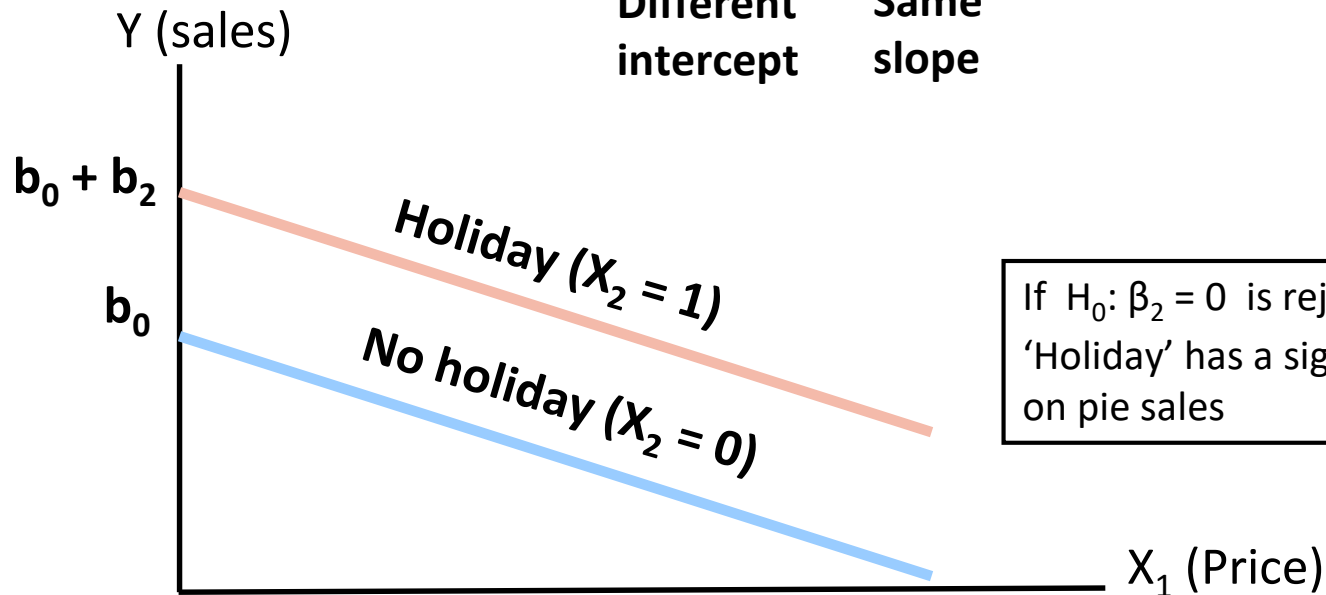
($X_2$ = 0 if there was no holiday that week)

# +Dummy Variable Example (with 2 Levels):

| | | | |
|---|---|---|---|
| $\hat{Y} = b_0 + b_1 X_1 + b_2(1) =$ | $(b_0 + b_2)$ | $+ b_1 X_1$ | Holiday |
| $\hat{Y} = b_0 + b_1 X_1 + b_2(0) =$ | $b_0$ | $+ b_1 X_1$ | No holiday |

**Different intercept**  **Same slope**

Y (sales)

$b_0 + b_2$

$b_0$

*Holiday (X$_2$ = 1)*

*No holiday (X$_2$ = 0)*

If $H_0$: $\beta_2 = 0$ is rejected, then 'Holiday' has a significant effect on pie sales

$X_1$ (Price)

# +Interpreting the Dummy Variable Coefficient - with 2 Levels

$$\widehat{Sales} = 300 - 30(Price) + 15\ (Holiday)$$

Sales: number of pies sold per week

Price:  pie price in $

Holiday: $\begin{cases} 1 & \text{If a holiday occurred during the week} \\ 0 & \text{If no holiday occurred} \end{cases}$

$b_2 = 15$: on average, sales were 15 pies greater in weeks with a holiday than in weeks without a holiday, given the same price

# +Dummy Variable Models - more than 2 Levels

The number of dummy variables is **one less than the number of levels**

Example:

Y = apartment price

$X_1$ = size of apartment in hundreds of square metres

If number of bedrooms is incorporated:

Bedrooms = one, two, three

Three levels, so two dummy variables are needed

# +Dummy Variable Models - more than 2 Levels (Cont)

**Example:**

Let '1-bedroom' be the default category, and let X2 and X3 be used for the other two categories

$Y$  = apartment price
$X_1$ = size in hundreds of square metres
$X_2$ = 2 bedroom, 0 otherwise
$X_3$ = 3 bedroom, 0 otherwise

The multiple regression equation is:

$$\hat{Y} = b_0 + b_1 X_1 + b_2 X_2 + b_3 X_3$$

# +Dummy Variable Models - more than 2 Levels (Cont)

Consider the regression equation:

$$\hat{Y} = 20.43 + 0.045X_1 + 18.84X_2 + 33.53X_3$$

For 1-bedroom: $X_2 = X_3 = 0$

$$\hat{Y} = 20.43 + 0.045X_1$$

For 2-bedroom: $X_2 = 1$; $X_3 = 0$

$$\hat{Y} = 20.43 + 0.045X_1 + 18.84$$

With the same size in hundreds of square meters, a 2-bedroom will have an estimated average price of 18.84 thousand dollars more than a 1-bedroom apartment

For 3-bedroom: $X_2 = 0$; $X_3 = 1$

$$\hat{Y} = 20.43 + 0.045X_1 + 33.53$$

With the same size in hundreds of square meters, a 3-bedroom will have an estimated average price of 33.53 thousand dollars more than a 1-bedroom

# +Collinearity

High correlation exists among two or more independent variables

This means the correlated variables contribute redundant information to the multiple regression model

Including two highly correlated independent variables can adversely affect the regression results

No new information provided:

- Can lead to unstable coefficients (large standard error and low t-values)

- Coefficient signs may not match prior expectations

# +Some Indications of Strong Collinearity

- Incorrect signs on the coefficients

- Large change in the value of a previous coefficient when a new variable is added to the model

- A previously significant variable becomes non-significant when a new independent variable is added

- The estimate of the standard deviation of the model increases when a variable is added to the model

# + Measuring Collinearity Variance Inflationary Factor

The variance inflationary factor $VIF_j$ can be used to measure collinearity:

$$VIF_j = \frac{1}{1-R_j^2}$$

Where: $R_j^2$ is the coefficient of multiple determination of independent variable $X_j$ with all other X variables

**If:** $VIF_j = 1$, $X_j$ **is uncorrelated with the other Xs**

**If:** $VIF_j > 10$, $X_j$ **is highly correlated with the other Xs (conservative estimate reduces this to $VIF_j > 5$)**