



# MULTIPLE REGRESSION ANALYSIS

# MULTIPLE REGRESSION

Many practical situations involve analyzing the relationships among three or more variables.

When multiple independent variables are to be included in an analysis simultaneously, multiple regression is very useful.



## Multiple Regression Model Population

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_k x_k + \varepsilon$$

where:

$\beta_0$  = Population's regression constant

$\beta_j$  = Population's regression coefficient for each variable  $x_j = 1, 2, \dots, k$

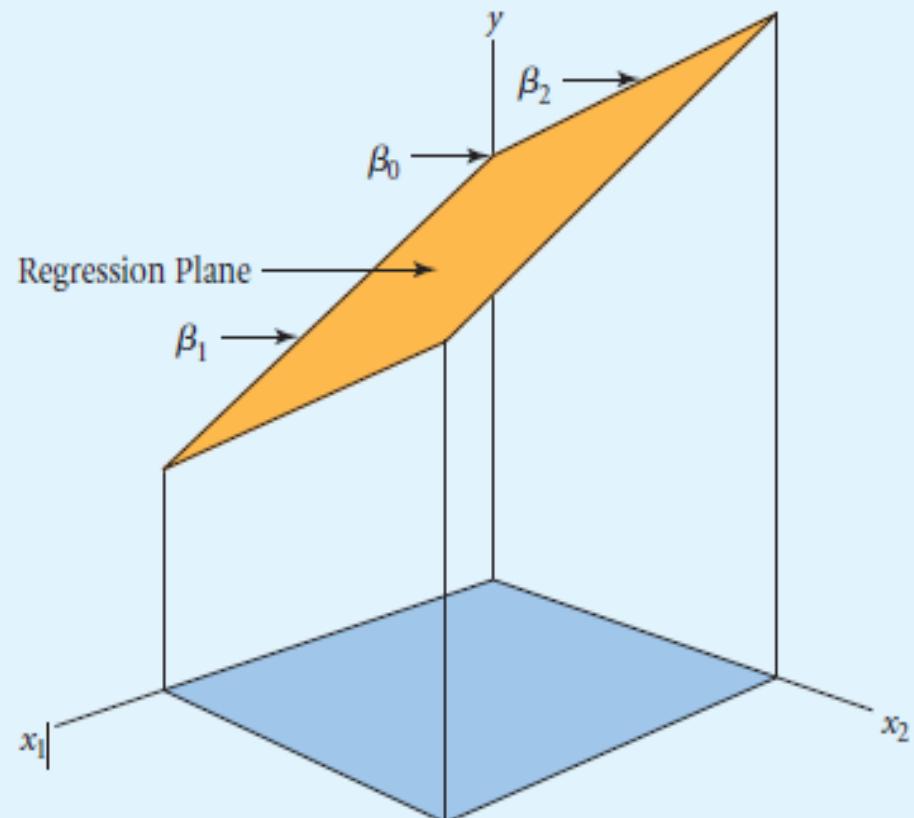
$k$  = Number of independent variables

$\varepsilon$  = Model error

## Estimated Multiple Regression Model

$$\hat{y} = b_0 + b_1 x_1 + b_2 x_2 + \cdots + b_k x_k$$

Figure: Multiple Regression Hyperplane  
for Population



# CONDUCTING MULTIPLE REGRESSION

Hospital readmissions within 30 days of discharge have become a critical quality metric in healthcare systems globally. High readmission rates not only reflect negatively on patient care outcomes but also lead to increased healthcare costs. Hospital administrators seek to understand the key factors that contribute to readmission risks to improve service quality and reduce unnecessary costs.

Given a large dataset of discharged patients—including demographic information, comorbidities, treatment duration, discharge instructions compliance, medication types, lab results, and hospital resource utilization—this study aims to build a multiple linear regression model to predict the likelihood of a 30-day readmission.

However, due to potential interdependencies among predictors (e.g., certain comorbidities often occurring together or overlapping with lab values), multicollinearity is anticipated and must be addressed to ensure model validity and interpretability.

# **CONDUCTING MULTIPLE REGRESSION**

## **Primary Question:**

What are the key predictors of 30-day hospital readmission among discharged patients, and how effectively can a multiple regression model estimate this likelihood based on patient-level and treatment-related variables?

## **Supporting Questions:**

1. Which independent variables have the strongest association with 30-day readmission?
2. Is there evidence of multicollinearity among predictors, and how does it affect the stability and interpretability of the regression model?
3. What strategies can be used to mitigate multicollinearity in the dataset?



# Variable Descriptions (Healthcare Readmission Dataset)

## Age

- Patient's age in years (whole number).

## Length\_of\_Stay

- Number of days the patient stayed in the hospital during their latest admission.

## Comorbidities

- Number of co-existing medical conditions (e.g., diabetes, hypertension).

## Lab\_Abnormalities

- A derived score indicating abnormalities found in lab tests (often correlated with comorbidities).

## Medication\_Complexity

- A score reflecting the complexity and number of medications prescribed.

## Discharge\_Instruction\_Score

- Score (1 to 10) evaluating patient understanding of post-discharge instructions.

## Resource\_Utilization

- A proxy score combining factors such as staff time, equipment, and space used during treatment.

## Readmission\_Risk

- Predicted risk score (continuous) estimating likelihood of readmission within 30 days.

# CONDUCTING MULTIPLE REGRESSION



## Business Requirements:

Building a **multiple linear regression model** to predict the likelihood of a 30-day readmission to a hospital.



## Data Set:

Readmission\_Data\_Final.xlsx

Variables Settings: Xs = ?, Y=?

Correlation Analysis: Measuring the strength of the relationship

Multicollinearity Issue: Check which variables



## Regression Model Development:

Excel, SPSS, R-Programming,  
Python

## Regression Model Diagnostic:

Statistical Test, R-Squared, RMSE

# MODEL SPECIFICATION

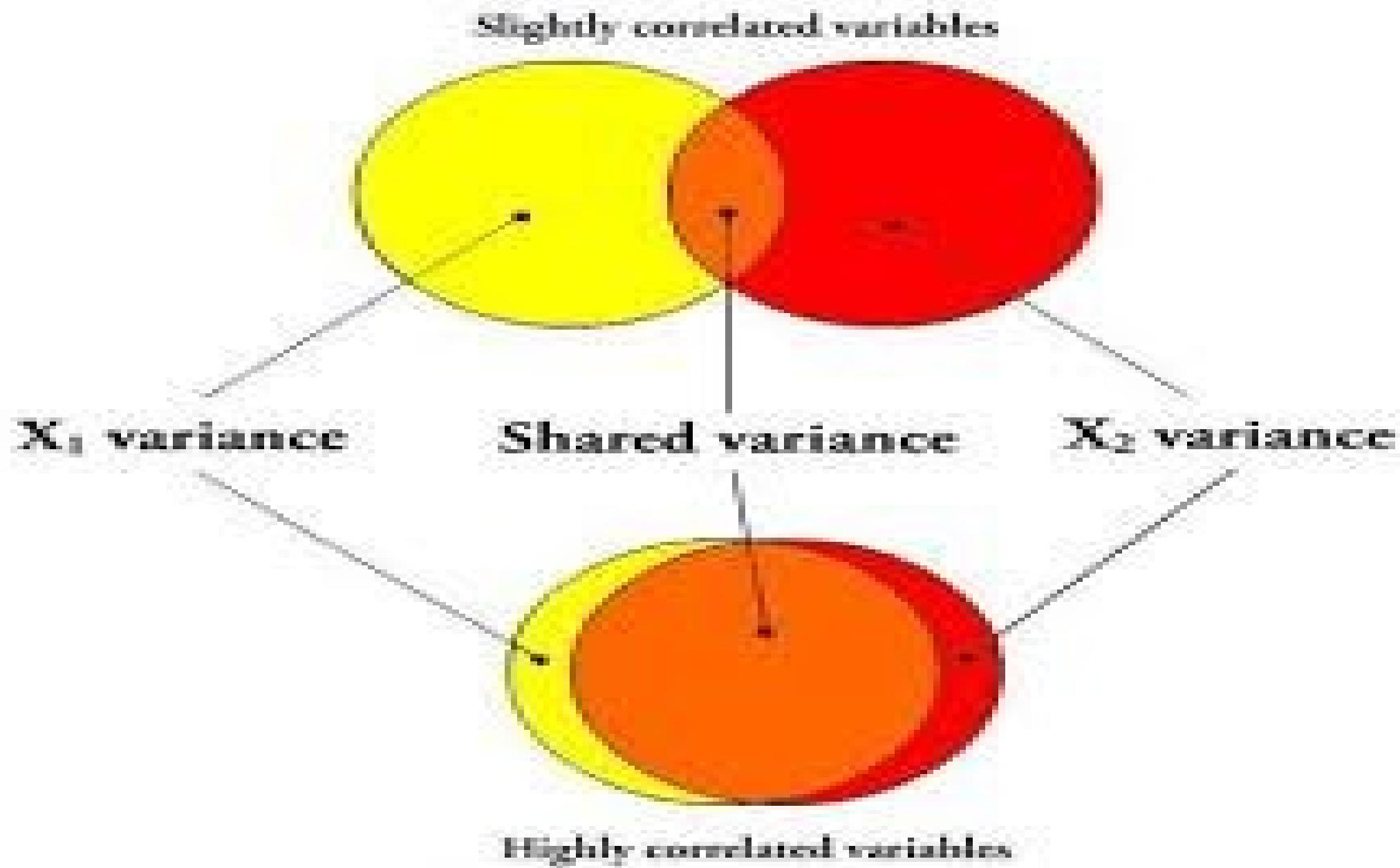
- Construct correlation matrix
- Checking the relationship between Y and X(s)
- Check the **multicollinearity** between X(s)



# MULTICOLLINEARITY

two explanatory variables that have high correlation.

correlation between two explanatory variables less than -0.7 or greater than 0.7 may be cause for concern





## Example of multicollinearity situation:

Problem 1: Prediction analysis on sales of lemonade for XYZ Café. Variables that might help to explain lemonade sales are, outside temperature and air-conditioning bills.

If the researcher includes both explanatory variables in the model, he may get results that are a little strange, because the two explanatory variables are themselves highly correlated.

As temperatures increase, so do air-conditioning bills. It would be meaningless to include both variables in the model because they are both doing the same job when it comes to explaining lemonade sales.



If exists explanatory variables that are highly correlated among themselves watch out for strange results in the regression output.

The strange results:

- Getting estimates of slope coefficients that are the opposite sign of what we would expect.

$x_1$	$x_2$	$y$
11.4	-9.7	14.7
12.5	-11.5	38.8
16.4	-15.9	42.9
14.4	-13.9	45.7
15.3	-14.2	52.3
18	-18.5	55.9
19.5	-21.2	60.1
25.2	-27.2	72.6

Example: Refer to the dataset given and answer the questions:

- (a) Find the correlation matrix among all three variables.
- (b) Find the least-squares regression model using both  $x_1$  and  $x_2$  as explanatory variables.
- (c) Comment on the effect that including both  $x_1$  and  $x_2$  has on the t-test statistics.

Solution:

(a) The correlation matrix is as below.

	$x_1$	$x_2$
$x_2$	-0.99607612	
$y$	0.89091722	-0.89445738

An extremely high correlation exists between  $x_1$  and  $x_2$ , so multicollinearity exists between the two variables.

Solution:

(b) Find the least-squares regression model using both  $x_1$  and  $x_2$  as explanatory variables.

Parameter estimates:						
Parameter	Estimate	Std. Err.	Alternative	DF	T-Stat.	P-value
Intercept	3.1966894	35.467992	$\neq 0$	5	0.090128851	0.9317
$x_1$	-0.015176864	8.8287255	$\neq 0$	5	-0.0017190323	0.9987
$x_2$	-2.7209724	6.8440747	$\neq 0$	5	-0.39756615	0.7074

Analysis of variance table for multiple regression model:					
Source	DF	SS	MS	F-stat	P-value
Model	2	1645.8514	822.92568	10.003384	0.0179
Error	5	411.32364	82.264729		
Total	7	2057.175			

Summary of fit:	
Root MSE:	9.0699906
R-squared:	0.8001
R-squared (adjusted):	0.7201

The  $P$ -value for the  $F$ -test statistic is 0.0179, indicating that the regression model is significant. However, if we look at each individual  $t$ -test statistic, we see that all regression parameter has a very high  $P$ -value indicating that neither coefficient is different from zero.

(c) Comment on the effect that including both  $x_1$  and  $x_2$  has on the  $t$ -test statistics.

The contradictory results of the regression output occur because both  $x_1$  and  $x_2$  are related to the response variable  $y$ , as indicated by the correlation matrix.

However,  $x_1$  and  $x_2$  are also related to each other. So, with  $x_1$  in the model,  $x_2$  adds little explanation. Likewise, with  $x_2$  in the model,  $x_1$  adds little explanation.

**The solution** is to use only one explanatory variable. Which explanatory variable we choose is up to you. We can choose either the explanatory variable with the lower  $P$  value or the explanatory variable that has the higher correlation with the response.

A black and white photograph of a woman with long dark hair, wearing a polka-dot dress, looking intently at a tablet device she is holding in her hands. The background is blurred, suggesting an indoor setting like a library or office.

# VIF – VARIANCE INFLATION FACTORS

VIF is a statistical measure that helps detect multicollinearity in a multiple regression model.

## 🔍 How VIF Works:

For each predictor variable, VIF is calculated using this formula:

$$VIF_I = \frac{1}{1-R_i^2}$$

Where:

$R_i^2$  is the R-squared value from a regression of variable  $i$  against all the other predictors.

A black and white photograph of a woman with long dark hair, wearing a polka-dot dress, looking intently at a tablet device she is holding in her hands. The background is blurred, suggesting an indoor setting like a library or office.

# VIF – VARIANCE INFLATION FACTORS

VIF Value	Interpretation
< 1	No multicollinearity
1–5	Moderate multicollinearity (acceptable)
>5	High multicollinearity (caution advised)
>10	Serious multicollinearity (problematic)

## 🛠 How to Address High VIFs:

- Drop one of the highly correlated variables
- Combine them (e.g., using PCA)
- Use regularized models like Ridge or Lasso regression

## How to Run VIF in SPSS

### Step 1: Open Your Data

Open your dataset in SPSS.

### Step 2: Run Linear Regression

Go to Analyze > Regression > Linear...

Set your dependent variable (e.g., Readmission\_Risk)

Move all your independent variables into the Independent(s) box

### Step 3: Request Collinearity Diagnostics

Click on the "Statistics..." button

Check Collinearity diagnostics

Click Continue, then OK

### Step 4: Interpret Output

In the output table, under Coefficients. you will see two columns:

Tolerance  
VIF



# MODEL SPECIFICATION

LET CHECK: Multicollinearity

- Construct correlation matrix
- Checking the relationship between Y and X(s)
- Check the multicollinearity between X(s)

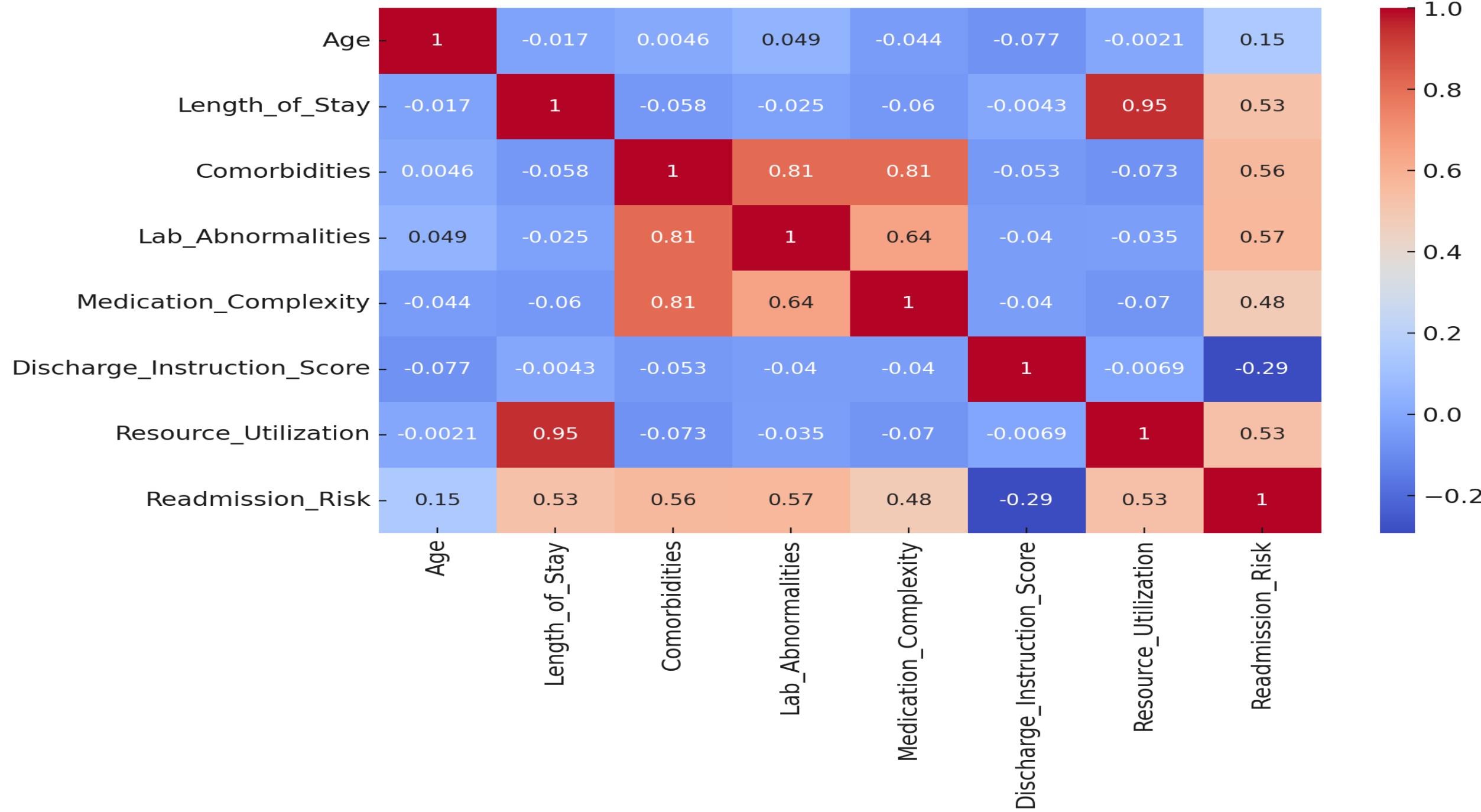


Do you find any multicollinearity among the explanatory variable?

Correlations									
	Age	Length_of_Stay	Comorbidities	Lab_Abnormalities	Medication_Complexity	Discharge_Instruction_Score	Resource_Utilization	Readmission_Risk	
Age	Pearson Correlation	1	-.014	.005	.049	-.042	-.079	-.002	.152**
	Sig. (2-tailed)		.751	.911	.278	.352	.079	.972	<.001
	N	500	500	500	500	500	500	500	500
Length_of_Stay	Pearson Correlation	-.014	1	-.052	-.025	-.052	.000	.947**	.531**
	Sig. (2-tailed)	.751		.241	.570	.249	.992	<.001	<.001
	N	500	500	500	500	500	500	500	500
Comorbidities	Pearson Correlation	.005	-.052	1	.798**	.792**	-.049	-.069	.564**
	Sig. (2-tailed)	.911	.241		<.001	<.001	.277	.123	<.001
	N	500	500	500	500	500	500	500	500
Lab_Abnormalities	Pearson Correlation	.049	-.025	.798**	1	.621**	-.032	-.041	.554**
	Sig. (2-tailed)	.278	.570	<.001		<.001	.472	.364	<.001
	N	500	500	500	500	500	500	500	500
Medication_Complexity	Pearson Correlation	-.042	-.052	.792**	.621**	1	-.048	-.063	.482**
	Sig. (2-tailed)	.352	.249	<.001	<.001		.282	.162	<.001
	N	500	500	500	500	500	500	500	500
Discharge_Instruction_Score	Pearson Correlation	-.079	.000	-.049	-.032	-.048	1	-.001	-.292**
	Sig. (2-tailed)	.079	.992	.277	.472	.282		.987	<.001
	N	500	500	500	500	500	500	500	500
Resource_Utilization	Pearson Correlation	-.002	.947**	-.069	-.041	-.063	-.001	1	.534**
	Sig. (2-tailed)	.972	<.001	.123	.364	.162	.987		<.001
	N	500	500	500	500	500	500	500	500
Readmission_Risk	Pearson Correlation	.152**	.531**	.564**	.554**	.482**	-.292**	.534**	1
	Sig. (2-tailed)	<.001	<.001	<.001	<.001	<.001	<.001	<.001	
	N	500	500	500	500	500	500	500	500

\*\* Correlation is significant at the 0.01 level (2-tailed).

Correlation Matrix



Do you find any multicollinearity among the explanatory variable?

Model		Coefficients <sup>a</sup>						Collinearity Statistics	
		B	Std. Error	Standardized Coefficients Beta	t	Sig.	Tolerance	VIF	
1	(Constant)	2.113	3.351		.631	.528			
	Age	.265	.046	.126	5.739	<.001	.981	1.019	
	Length_of_Stay	1.273	.472	.184	2.696	.007	.102	9.775	
	Comorbidities	4.210	.689	.284	6.106	<.001	.220	4.543	
	Lab_Abnormalities	3.106	.439	.257	7.076	<.001	.361	2.768	
	Medication_Complexity	1.478	.426	.125	3.474	<.001	.371	2.699	
	Discharge_Instruction_Score	-1.994	.173	-.253	-11.555	<.001	.991	1.009	
	Resource_Utilization	2.592	.445	.397	5.818	<.001	.102	9.789	

a. Dependent Variable: Readmission\_Risk

# **MICROSOFT EXCEL**



Ms Excel -> Data -> Data Analysis

The image shows two overlapping dialog boxes from Microsoft Excel's Data Analysis add-in. The left dialog is titled 'Data Analysis' and lists various statistical tools, with 'Regression' selected. The right dialog is titled 'Regression' and contains settings for the analysis.

**Data Analysis Dialog (Left):**

- Analysis Tools:
  - F-Test Two-Sample for Variances
  - Fourier Analysis
  - Histogram
  - Moving Average
  - Random Number Generation
  - Rank and Percentile
  - Regression** (selected)
  - Sampling
  - t-Test: Paired Two Sample for Means
  - t-Test: Two-Sample Assuming Equal Variances
- Buttons: OK, Cancel, Help.

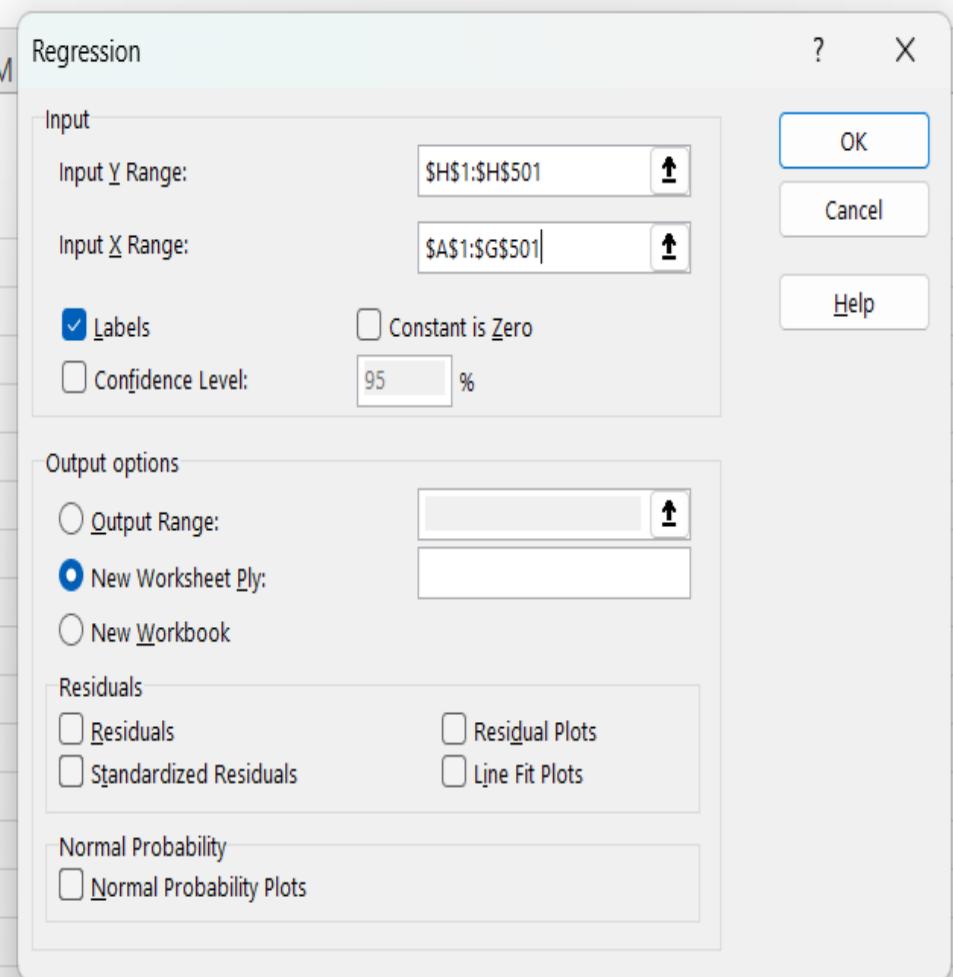
**Regression Dialog (Right):**

- Input:**
  - Input Y Range: (empty)
  - Input X Range: (empty)
  - Labels
  - Confidence Level: (set to 95%)
- Output options:**
  - Output Range:
  - New Worksheet Ply: (selected)
  - New Workbook
- Residuals:**
  - Residuals
  - Standardized Residuals
  - Residual Plots
  - Line Fit Plots
- Normal Probability:**
  - Normal Probability Plots
- Buttons: OK, Cancel, Help.

Annotations with arrows point to specific features:

- An arrow points from the 'Labels' checkbox in the Regression dialog to the 'Make 1<sup>st</sup> row as label' text.
- An arrow points from the 'Confidence Level' input field in the Regression dialog to the 'Changing the confidence level' text.
- An arrow points from the 'Residual Plots' checkbox in the Regression dialog to the 'Output for error (residual) diagnosis' text.
- An arrow points from the 'Normal Probability Plots' checkbox in the Regression dialog to the 'To see either normal or not normal distribution' text.
- An arrow points from the 'New Worksheet Ply' radio button in the Regression dialog to the 'Select col. represent DV' text.
- An arrow points from the 'Residual Plots' checkbox in the Regression dialog to the 'Select col. represent IV' text.

A	B	C	D	E	F	G	H	I	J	K	L	M
Age	Length_of_Stay	Comorbidities	Lab_Abnormalities	Medication_Complexity	Discharge_Instruction_Score	Resource_Utilization	Readmission_Risk					
1	70	3	1	1	-1	8	5	32.53149				
2	64	3	1	0	2	3	4	28.718				
3	71	1	2	1	1	10	0	35.53446				
4	80	13	3	4	2	1	11	91.25385				
5	63	3	3	2	2	6	3	24.00817				
6	63	1	2	4	2	7	1	29.18736				
7	81	0	1	0	0	7	0	-5.45488				
8	73	0	4	5	2	2	-1	59.91099				
9	60	1	1	1	2	9	2	-4.15083				
10	70	1	2	3	3	3	3	46.28458				
11	60	1	2	1	3	7	2	5.221882				
12	60	1	2	2	1	3	1	29.29702				
13	67	1	0	0	0	4	1	32.53378				
14	46	7	4	3	4	6	7	65.90648				
15	48	0	1	2	2	5	1	11.76506				
16	59	2	0	-1	1	2	1	17.38797				
17	55	2	2	2	1	3	2	35.18017				
18												



## Regression Model Building and Diagnostic

SUMMARY OUTPUT						
Regression Statistics						
Multiple R	0.874948681					
R Square	0.765535194					
Adjusted R Square	0.762199313					
Standard Error	10.03168192					
Observations	500					
ANOVA						
	df	SS	MS	F	Significance F	
Regression	7	161659.0822	23094.15	229.4851	1.6361E-150	
Residual	492	49512.24395	100.6346			
Total	499	211171.3261				
	Coefficients	Standard Error	t Stat	P-value	Lower 95%	Upper 95%
Intercept	2.113373435	3.350539192	0.630756	0.528493	-4.46975711	8.696504
Age	0.264513525	0.046088206	5.739289	1.66E-08	0.17395954	0.355068
Length_of_Stay	1.273095669	0.472189717	2.696153	0.007255	0.345338562	2.200853
Comorbidities	4.209721227	0.689416504	6.106209	2.07E-09	2.855157503	5.564285
Lab_Abnormalities	3.10565441	0.438892191	7.076121	5.13E-12	2.243320195	3.967989
Medication_Complexity	1.478398588	0.425511594	3.474403	0.000557	0.642354533	2.314443
Discharge_Instruction_Score	-1.993580038	0.172526552	-11.5552	1.72E-27	-2.332559753	-1.6546
Resource_Utilization	2.592036335	0.44548748	5.818427	1.07E-08	1.716743712	3.467329

Compute the regression equation:

The estimate of the multiple regression model are:

Readmission\_Risk=2.11+0.26(Age)+1.27(Length\_of\_Stay)+4.21(Comorbidities)+3.11(Lab\_Abnormalities)+1.48(Medication\_Complexity)-1.99(Discharge\_Instruction\_Score)+2.99(Resource\_Utilization)

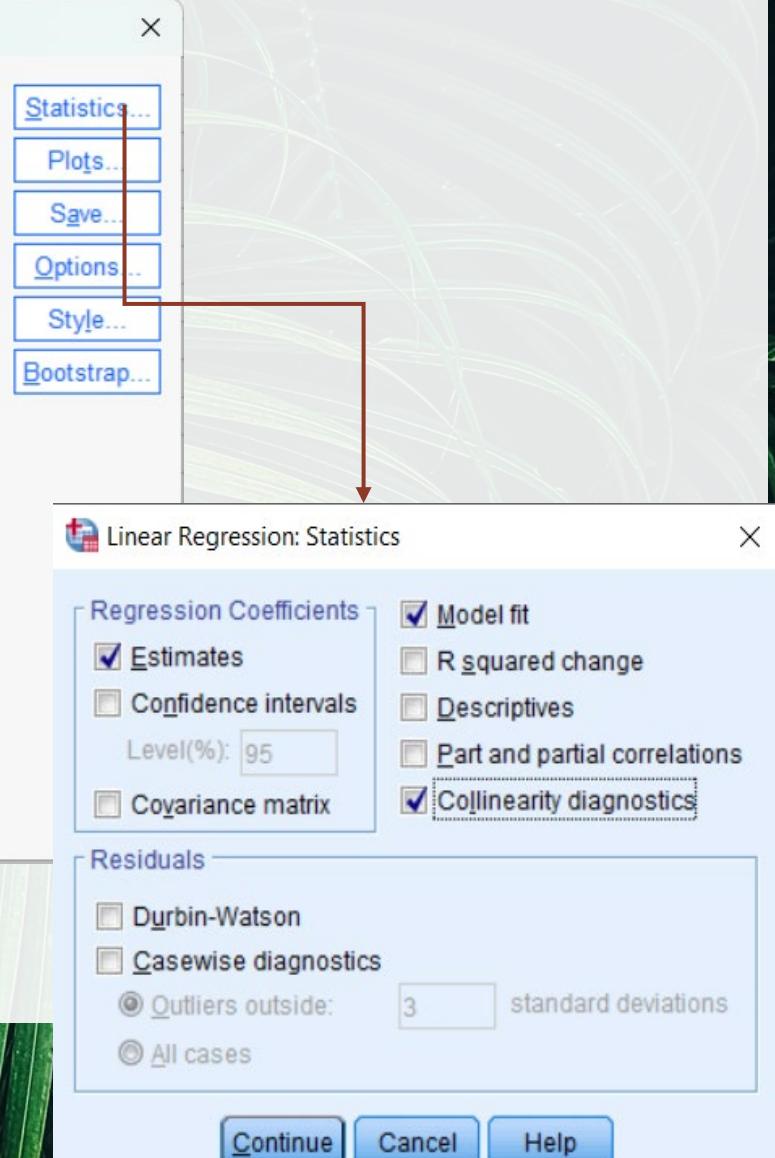
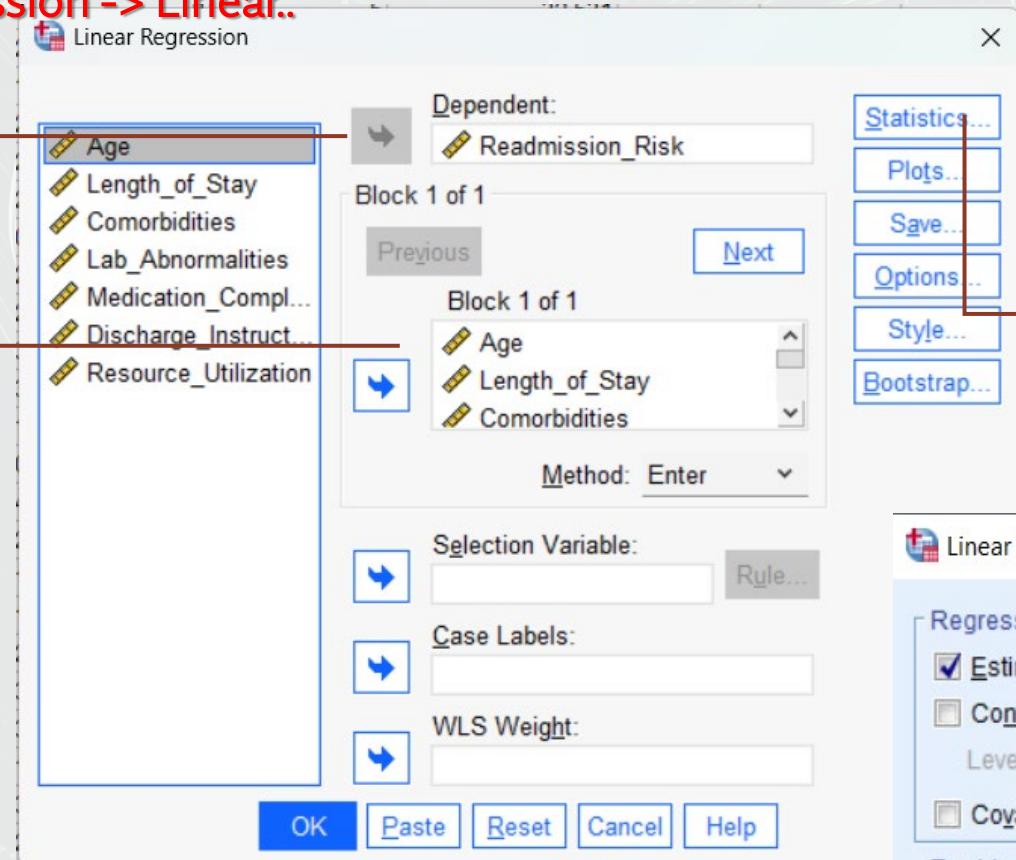
# SPSS



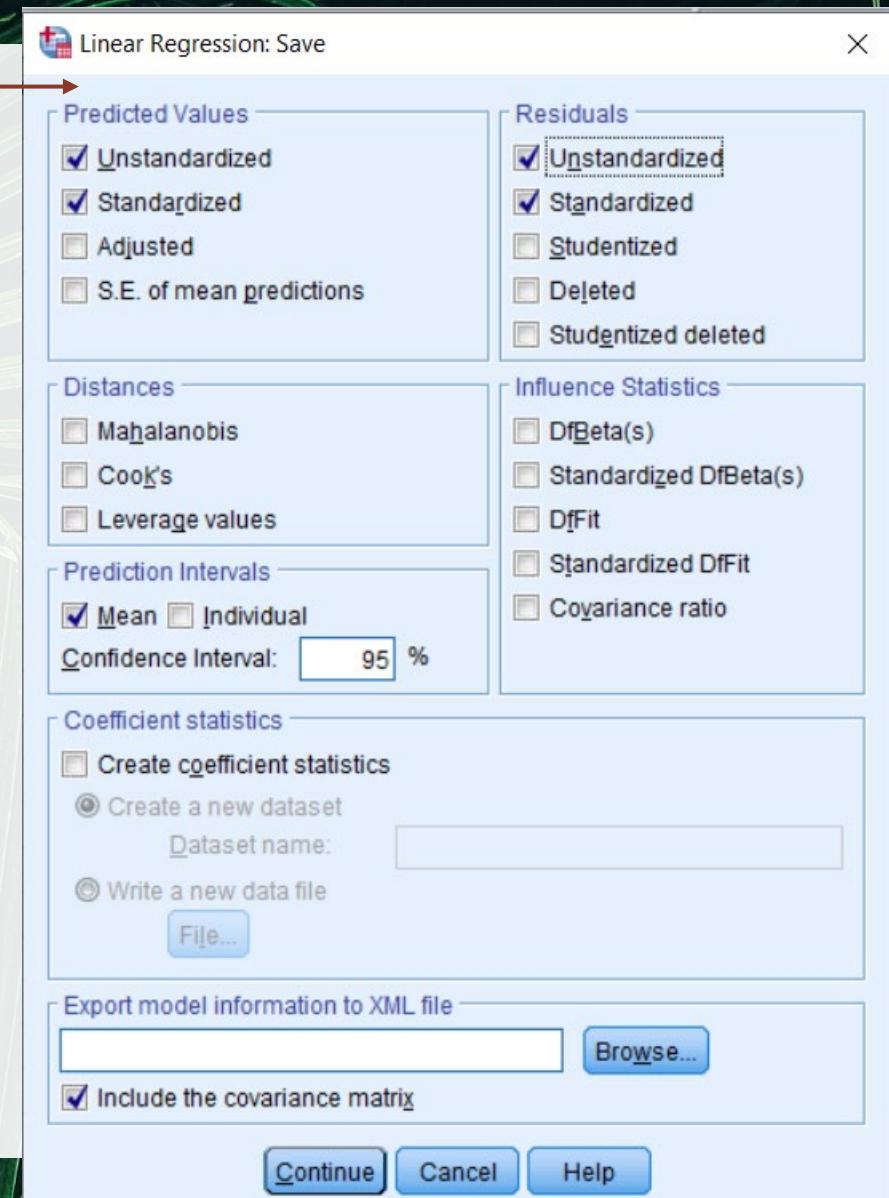
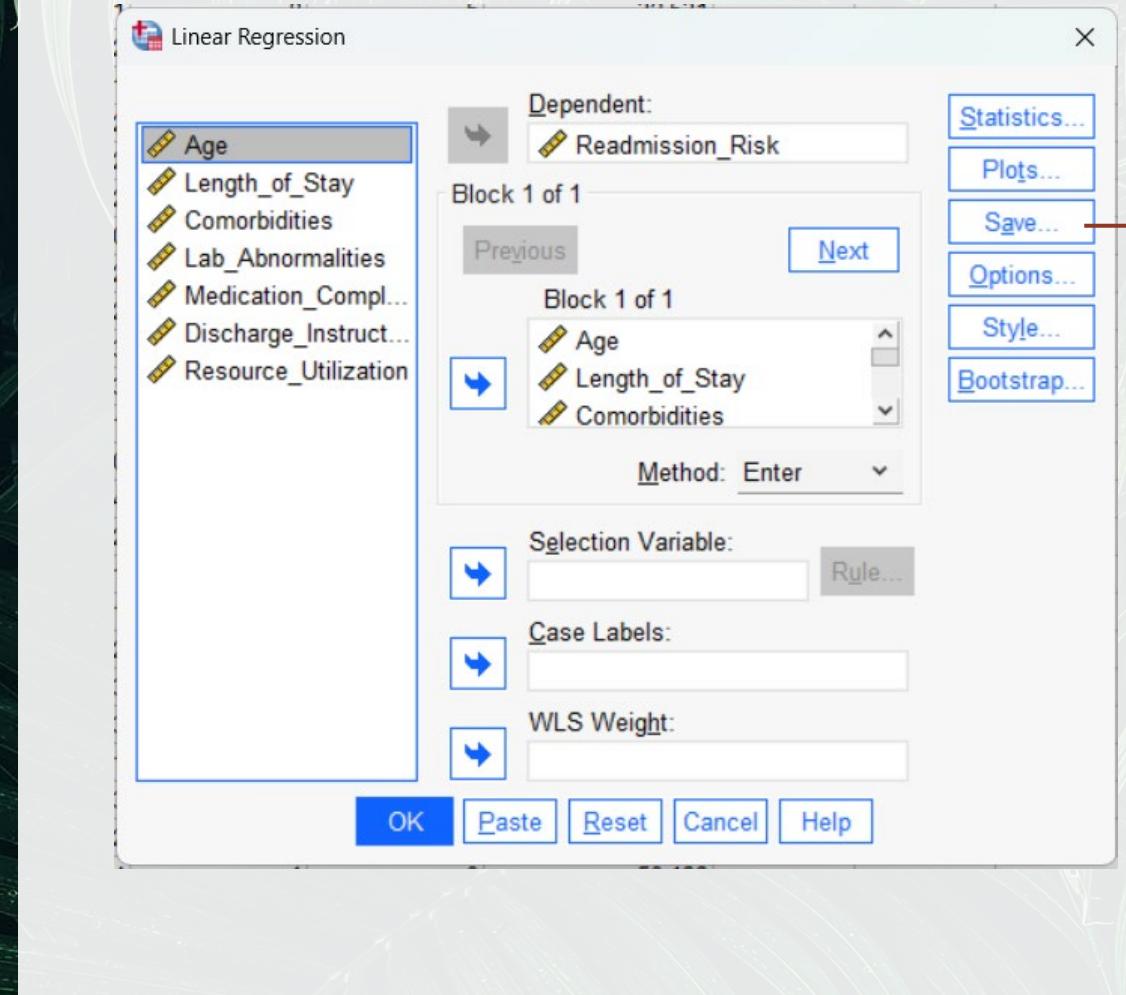
SPSS-> Analyze -> Regression -> Linear.

Dependent (Y) variable

Independent (X) variable



SPSS-> Analyze -> Regression -> Linear..



### Model Summary

Model	R	R Square	Adjusted R Square	Std. Error of the Estimate
1	.875 <sup>a</sup>	.766	.762	10.031682

a. Predictors: (Constant), Resource\_Utilization, Discharge\_Instruction\_Score, Lab\_Abnormalities, Age, Medication\_Complexity, Comorbidities, Length\_of\_Stay

### ANOVA<sup>a</sup>

Model		Sum of Squares	df	Mean Square	F	Sig.
1	Regression	161659.082	7	23094.155	229.485	<.001 <sup>b</sup>
	Residual	49512.244	492	100.635		
	Total	211171.326	499			

a. Dependent Variable: Readmission\_Risk

b. Predictors: (Constant), Resource\_Utilization, Discharge\_Instruction\_Score, Lab\_Abnormalities, Age, Medication\_Complexity, Comorbidities, Length\_of\_Stay

### Coefficients<sup>a</sup>

Model		Unstandardized Coefficients		t	Sig.	Collinearity Statistics	
		B	Std. Error			Beta	Tolerance
1	(Constant)	2.113	3.351	.631	.528		
	Age	.265	.046	.126	5.739	<.001	.981
	Length_of_Stay	1.273	.472	.184	2.696	.007	.102
	Comorbidities	4.210	.689	.284	6.106	<.001	.220
	Lab_Abnormalities	3.106	.439	.257	7.076	<.001	.361
	Medication_Complexity	1.478	.426	.125	3.474	<.001	.371
	Discharge_Instruction_Score	-1.994	.173	-.253	-11.555	<.001	.991
	Resource_Utilization	2.592	.445	.397	5.818	<.001	.102

a. Dependent Variable: Readmission\_RISK

$$\text{Readmission_Risk} = 2.11 + 0.26(\text{Age}) + 1.27(\text{Length\_of\_Stay}) + 4.21(\text{Comorbidities}) + 3.11(\text{Lab\_Abnormalities}) \\ + 1.48(\text{Medication\_Complexity}) - 1.99(\text{Discharge\_Instruction\_Score}) + 2.99(\text{Resource\_Utilization})$$

# REGRESSION MODEL DIAGNOSTIC

## Examine Regression Prediction Model Significance

p-value = 0.000 which is  $<\alpha= 0.05$ , therefore H<sub>0</sub> is rejected.

Conclusion: The regression model does explain a significant proportion of the variation in sales price. Thus, the overall model is statistically significant.



## Measure Model Fitness

Adjusted  $R^2$  value, about 76.2% of the variation in sales price



## Examine Regression Slope Significance

For all reg. coeff, the p-value = 0.000 which is  $<\alpha= 0.05$ , therefore H<sub>0</sub> is rejected. All variables is significance to be included in the model



## Examine the Error Rate

RMSE (residual standard error) = 10.03



## MULTICOLLINEARITY

How does the new model perform when the multicollinearity issue is tackled?



<u>Variable</u>	<u>Coefficient</u>	<u>Meaning</u>
Age	+0.26	Older patients are more likely to be readmitted.
Length_of_Stay	+1.27	Longer hospital stays suggest more serious conditions = higher risk.
Comorbidities	+4.21	More chronic diseases = higher chance of returning to hospital.
Lab_Abnormalities	+3.11	Abnormal lab results suggest complications = increased risk.
Medication_Complexity	+1.48	Complex medication regimens may lead to non-compliance = higher risk.
Resource_Utilization	+2.99	Higher use of resources implies complex care = increased risk.
Discharge_Instruction_Score	-1.99	Better patient understanding of discharge instructions reduces risk.



# ENHANCED REGRESSION MODEL

## Model Summary

Model	R	R Square	Adjusted R Square	Std. Error of the Estimate
1	.866 <sup>a</sup>	.749	.746	10.360553

a. Predictors: (Constant), Discharge\_Instruction\_Score, Length\_of\_Stay, Lab\_Abnormalities, Age, Medication\_Complexity, Comorbidities

## ANOVA<sup>a</sup>

Model		Sum of Squares	df	Mean Square	F	Sig.
1	Regression	158252.188	6	26375.365	245.716	<.001 <sup>b</sup>
	Residual	52919.138	493	107.341		
	Total	211171.326	499			

a. Dependent Variable: Readmission\_Risk

b. Predictors: (Constant), Discharge\_Instruction\_Score, Length\_of\_Stay, Lab\_Abnormalities, Age, Medication\_Complexity, Comorbidities



## Coefficients<sup>a</sup>

Model	1	Unstandardized Coefficients		Standardized Coefficients Beta	t	Sig.	Collinearity Statistics	
		B	Std. Error				Tolerance	VIF
	(Constant)	1.621	3.459		.469	.640		
	Age	.275	.048	.131	5.778	<.001	.982	1.018
	Length_of_Stay	3.876	.156	.560	24.796	<.001	.996	1.004
	Comorbidities	4.088	.712	.276	5.744	<.001	.220	4.539
	Lab_Abnormalities	3.085	.453	.255	6.806	<.001	.361	2.768
	Medication_Complexity	1.506	.439	.127	3.426	<.001	.371	2.698
	Discharge_Instruction_Score	-1.997	.178	-.254	-11.207	<.001	.991	1.009

a. Dependent Variable: Readmission\_Risk



## ADDING QUALITATIVE INDEPENDENT VARIABLE

# **QUALITATIVE INDEPENDENT VARIABLE**

There is situation you may wish to use a **qualitative** (lower level) variable as an explanatory variable in a regression model

Example: use of variable such as; marital status, gender, education level or job performance.

How these variable can be incorporated into a multiple regression analysis?

Dummy (or indicator) variable – a variable that is assigned a value equal to either 0 or 1, depending on whether the observation possesses a given characteristic.

$$x_1 = 1 \text{ if female}$$

$$x_1 = 0 \text{ if male}$$

if more than two mutually exclusive (for example, never married, married, divorced)

$$x_1 = 1 \text{ if never married, 0 if not}$$

$$x_2 = 1 \text{ if married, 0 if not}$$

$$x_3 = 1 \text{ if divorced, 0 if not}$$

## Example:

The population from which the sample was selected consists of executives between the ages of 24 and 60 who are working in U.S. manufacturing businesses. Data for annual salary ( $y$ ) and age ( $x_1$ ) is describe in the table. The objective of the problem is to determine whether a model can be generated to explain the variation in annual salary for business executives given the explanatory variable of age and the qualitative variable ( $x_2$ ) of had a master of business administration (MBA) degree. The dummy variable is hold with this indication:

$$\begin{aligned}x_2 &= 1 \text{ if holds MBA degree} \\x_2 &= 0 \text{ if did not hold MBA degree}\end{aligned}$$

**TABLE 15.2 | Executive Salary Data Including MBA Variable**

Salary(\$)	Age	MBA
65,000	26	0
85,000	28	1
74,000	36	0
83,000	35	0
110,000	35	1
160,000	40	1
100,000	41	0
122,000	42	1
85,000	45	0
120,000	46	1
105,000	50	0
135,000	51	1
125,000	55	0
175,000	50	1
156,000	61	1
140,000	63	0

From the statistical packages, the estimated regression equation are:

$$\hat{y} = 6,974 + 2,055x_1 + 35,236x_2$$

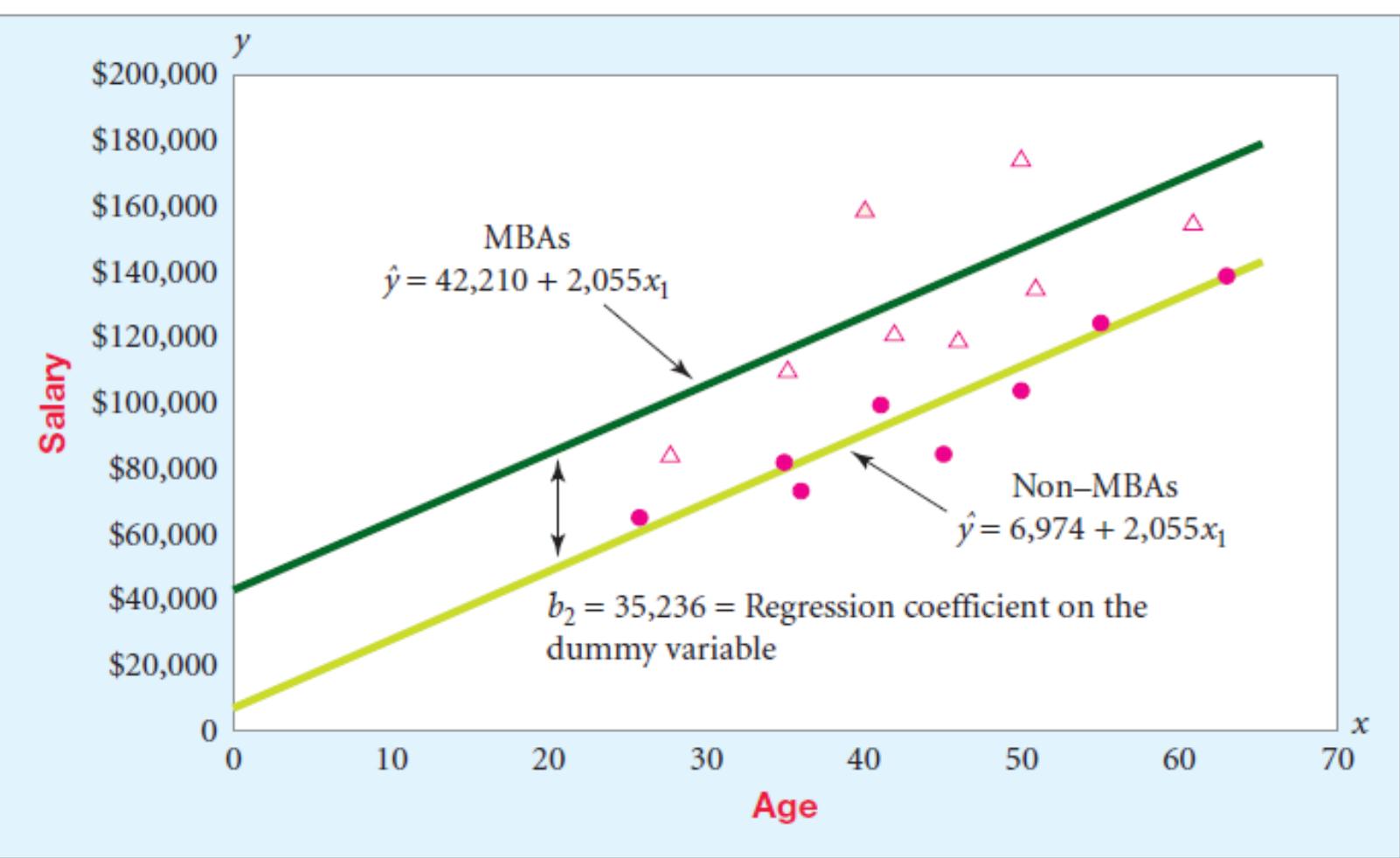
Because the dummy variable,  $x_2$ , has been coded 0 or 1 depending on MBA status, incorporating it into the regression model is like having two simple linear regression lines with the same slope, but different intercept.

For instance when  $x_2= 1$  (respondent who holds MBA degree),  $\hat{y} = 6,974 + 2,055x_1 + 35,236(1)$

$$\hat{y} = 42,210 + 2,055x_1$$

and when  $x_2= 0$  (respondent who did not holds MBA degree),  $\hat{y} = 6,974 + 2,055x_1 + 35,236(0)$

$$\hat{y} = 6,974 + 2,055x_1$$



## **NEW DATA SET:**

Readmission\_data\_final\_category.xlsx

# SPSS



### Model Summary

Model	R	R Square	Adjusted R Square	Std. Error of the Estimate
1	.885 <sup>a</sup>	.782	.779	9.892

a. Predictors: (Constant), Type of Discharge, Medication\_Complexity, Age, Length\_of\_Stay, Discharge\_Instruction\_Score, Lab\_Abnormalities, Comorbidities

### ANOVA<sup>a</sup>

Model		Sum of Squares	df	Mean Square	F	Sig.
1	Regression	173064.357	7	24723.480	252.676	<.001 <sup>b</sup>
	Residual	48140.561	492	97.847		
	Total	221204.918	499			

a. Dependent Variable: Readmission\_Risk

b. Predictors: (Constant), Type of Discharge, Medication\_Complexity, Age, Length\_of\_Stay, Discharge\_Instruction\_Score, Lab\_Abnormalities, Comorbidities

### Coefficients<sup>a</sup>

Model		Unstandardized Coefficients		Standardized Coefficients Beta	t	Sig.	Collinearity Statistics	
		B	Std. Error				Tolerance	VIF
1	(Constant)	4.408	3.321		1.327	.185		
	Age	.232	.045	.108	5.100	<.001	.982	1.018
	Length_of_Stay	4.220	.149	.596	28.252	<.001	.994	1.006
	Comorbidities	2.522	.680	.166	3.712	<.001	.220	4.539
	Lab_Abnormalities	3.922	.433	.317	9.049	<.001	.360	2.777
	Medication_Complexity	2.125	.420	.175	5.063	<.001	.370	2.700
	Discharge_Instruction_Score	-1.997	.170	-.248	-11.738	<.001	.991	1.009
	Type of Discharge	7.584	.991	.162	7.654	<.001	.993	1.007

a. Dependent Variable: Readmission\_Risk

$$\text{Readmission_Risk} = 4.41 + 0.23(\text{Age}) + 4.22(\text{Length\_of\_Stay}) + 2.52(\text{Comorbidities}) + 3.92(\text{Lab\_Abnormalities}) \\ + 2.13(\text{Medication\_Complexity}) - 1.99(\text{Discharge\_Instruction\_Score}) + 7.58(\text{Type o Discharge})$$

# REGRESSION MODEL DIAGNOSTIC

## Examine Regression Prediction Model Significance

p-value = 0.000 which is  $<\alpha= 0.05$ , therefore H<sub>0</sub> is rejected.

Conclusion: The regression model does explain a significant proportion of the variation in sales price. Thus, the overall model is statistically significant.



## Measure Model Fitness

Adjusted  $R^2$  value, about 77.9% of the variation in sales price



## Examine Regression Slope Significance

For all reg. coeff, the p-value = 0.000 which is  $<\alpha= 0.05$ , therefore H<sub>0</sub> is rejected. All variables is significance to be included in the model



## Examine the Error Rate

RMSE (residual standard error) = 9.89



## MULTICOLLINEARITY

How does the new model perform when the multicollinearity issue is tackled?



# ENHANCED REGRESSION MODEL



## Model Summary

Model	R	R Square	Adjusted R Square	Std. Error of the Estimate
1	.881 <sup>a</sup>	.776	.774	10.019

a. Predictors: (Constant), Type of Discharge, Medication\_Complexity, Age, Length\_of\_Stay, Discharge\_Instruction\_Score, Lab\_Abnormalities

## ANOVA<sup>a</sup>

Model		Sum of Squares	df	Mean Square	F	Sig.
1	Regression	171716.099	6	28619.350	285.102	<.001 <sup>b</sup>
	Residual	49488.819	493	100.383		
	Total	221204.918	499			

a. Dependent Variable: Readmission\_Risk

b. Predictors: (Constant), Type of Discharge, Medication\_Complexity, Age, Length\_of\_Stay, Discharge\_Instruction\_Score, Lab\_Abnormal



## Coefficients<sup>a</sup>

Model		Unstandardized Coefficients		Standardized Coefficients Beta	t	Sig.	Collinearity Statistics	
		B	Std. Error				Tolerance	VIF
1	(Constant)	5.561	3.349		1.660	.098		
	Age	.232	.046	.108	5.036	<.001	.982	1.018
	Length_of_Stay	4.203	.151	.594	27.792	<.001	.995	1.005
	Lab_Abnormalities	4.948	.338	.400	14.629	<.001	.607	1.648
	Medication_Complexity	3.096	.332	.255	9.320	<.001	.606	1.649
	Discharge_Instruction_Score	-2.010	.172	-.250	-11.665	<.001	.991	1.009
	Type of Discharge	7.626	1.004	.162	7.600	<.001	.993	1.007

a. Dependent Variable: Readmission\_Risk

# **THANK YOU**

Assoc. Prof. Dr. Nurulhuda Firdaus  
Bt Mohd Azmi (Huda)

012-5719624

[huda@utm.my](mailto:huda@utm.my)