

---

## *Uncovering the Diabetes Puzzle:*

### *Exploring the Interplay of Inactivity and Obesity*

---

The issues:

CDC (Center for Disease Control and Prevention) gathered a dataset in 2018 of diabetes, inactivity and obesity. They have instigated the relationship between them. We come across the following issues.

1. What are the main reasons behind the increase in diabetes cases?
2. Can we see a connection between the number of people who are overweight and the number of people with diabetes?
3. How does the amount of physical activity people do relate to the number of diabetes cases?
4. Which is more responsible for diabetes, being inactive or being overweight?
5. Does having both a lack of physical activity and being overweight make diabetes more likely?
6. Can we use data about obesity and inactivity to predict diabetes rates in the future?
7. How do we choose the best way to analyze the relationships between these health factors, like whether it's better to use a simple analysis or a more complex one?

8. Which health factors, like inactivity and obesity, should we include in our analysis?
9. How do we check if our analysis is a good fit for the data, and what measures can we use to see how well our analysis explains the patterns we observe?
10. When presenting our findings from this data, what are the most effective ways to show the information to others?

## Findings:

The findings from the analysis "Uncovering the Diabetes Puzzle: Exploring the Interplay of Inactivity and Obesity" relevant to the issues raised are as follows:

1. The study revealed a moderate positive correlation (0.57) between diabetes and inactivity, suggesting that increased physical inactivity is a significant factor in higher diabetes rates. Obesity also showed a positive correlation with diabetes, though less pronounced (0.39).
2. A positive correlation exists between obesity and diabetes, although it's weaker compared to the correlation between inactivity and diabetes.
3. Physical inactivity showed a strong and significant effect on diabetes rates, with a notable correlation and higher predictive power in regression models.
4. Inactivity has a more significant impact on diabetes prevalence than obesity, as indicated by higher t-values and stronger correlations in the analysis.
5. The combined analysis of inactivity and obesity in predicting diabetes rates showed that the multiple linear regression model, incorporating both factors, had a higher explanatory power (R-squared of 0.52) compared to models considering these factors individually.

6. The regression models used in the study, especially the multiple linear regression model, can be effective in predicting future diabetes rates based on data about obesity and inactivity.
7. The analysis employed both simple and multiple linear regressions. The multiple linear regression model, considering both inactivity and obesity, provided a stronger fit (higher R-squared values) and better predictive power for diabetes rates.
8. Both inactivity and obesity were included as independent variables in the regression models, with the combined model showing better predictive performance.
9. The goodness of fit was assessed using R-squared values, which indicated how well the independent variables explained the variation in the dependent variable (diabetes rates). Model performance was also evaluated using metrics like MAE, MSE, and RMSE, with lower errors in the test set compared to the training set, indicating good generalization.
10. The findings were visually presented through scatter plots, emphasizing the multifactorial nature of diabetes rates and the relationships between the variables.

These findings address the issues by showing the significant role of inactivity and obesity in diabetes prevalence, with inactivity having a more pronounced effect. The combined analysis of these factors provides a comprehensive understanding of their impact on diabetes rates.

## Discussions:

The findings from the "Uncovering the Diabetes Puzzle: Exploring the Interplay of Inactivity and Obesity" study have several important implications. Firstly, the strong correlation between physical inactivity and diabetes rates highlights the urgent need for public health initiatives to prioritize increasing physical activity among populations. This could significantly reduce diabetes prevalence. Although obesity shows a weaker correlation with diabetes compared to inactivity, it remains a considerable factor. Hence, public health policies must continue addressing obesity as part of comprehensive diabetes prevention and management strategies.

Simultaneously addressing both inactivity and obesity could lead to more effective diabetes prevention programs. This combined approach suggests that interventions promoting physical activity alongside healthy eating could be more impactful. The study's regression models, which effectively predict diabetes rates based on inactivity and obesity data, can be instrumental in planning future health interventions. These models enable health authorities to identify at-risk populations and implement targeted prevention programs.

Additionally, these results call for a reassessment of health education initiatives, emphasizing the dual effects of obesity and inactivity on diabetes. Awareness campaigns should educate the public about the importance of both maintaining regular physical activity and a healthy weight. For clinical practice, these findings provide valuable insights, particularly in advising patients about the

importance of physical activity and weight management in diabetes prevention and control.

The significant correlations and predictive power of inactivity and obesity for diabetes rates also call for more research into the effectiveness of specific interventions targeting these factors. This can help develop more effective diabetes prevention and management programs. Additionally, these findings can guide healthcare policymakers in resource allocation, suggesting that prioritizing interventions targeting physical inactivity might yield greater benefits in reducing diabetes prevalence. Interventions could be tailored to different populations based on their specific risk profiles, allowing for more personalized healthcare approaches.

## Appendix A-Methodology:

### **Data Source:**

We obtained the dataset from the CDC (Center for Disease Control and Prevention) in the year 2018. The indicators for obesity and inactivity represent the percentage of individuals within a county in an American state who are classified as obese or have engaged in no physical activity, respectively. This dataset was organized into three separate worksheets in an Excel format, dedicated to Diabetes, Inactivity, and Obesity.

### **Data Description:**

Additional indicators encompassed a range of factors, including food availability and gym accessibility. However, we made the decision to exclude these supplementary indicators from our analysis due to the limited number of data points available, which totaled approximately 5460. Although these factors are theoretically important contributors, incorporating them using common data points would have substantially reduced the dataset to only around 50 data points, resulting in a significant impact on the accuracy of our model. Therefore, we chose to work with a more extensive dataset to ensure the robustness of our analysis.

### **Data Preparation:**

Following the filtration process, our dataset was reduced to 3142 data points for diabetes, 1370 for inactivity, and 363 for obesity. Notably, only 352 data points contained values for all three variables: diabetes, inactivity, and obesity. So, our initial step involved data cleaning to remove missing data points. We also utilized the 'FIPS' column as a reference point. To integrate common data points across these datasets, we employed the VLOOKUP function, ensuring that we had a comprehensive and accurate dataset.

### **Variable Creation:**

Within this dataset, we considered three primary variables:

- % DIABETIC rates represent the prevalence of diabetes in a given area.
- % Inactivity indicates the level of physical inactivity in that same area.
- % Obesity represents the obesity rates in the same geographic regions.

These variables were central to our analysis.

## **Analytical Methods:**

Our analytical approach consisted of several key methods:

1. **Correlation Analysis**: We calculated a correlation matrix to understand how these variables are related to each other. This analysis provided insights into the interplay between diabetes rates, obesity rates, and physical inactivity.
2. **Simple Linear Regression**: Simple linear regression is a statistical method that examines the relationship between two variables, where one variable is used to predict or explain the variation in another variable. We conducted separate simple linear regression analyses. One focused on % Inactivity as the independent variable to predict % DIABETIC rates, isolating the impact of physical inactivity on diabetes prevalence. The other analysis used % Obesity as the independent variable, focusing on obesity as a predictor of diabetes.
3. **Multiple Linear Regression**: Multiple linear regression extends simple linear regression by considering the influence of two or more independent variables on a dependent variable, allowing for a more comprehensive analysis of relationships and predictions. We employed multiple linear regression, incorporating both % Inactivity and % Obesity as independent variables. This comprehensive approach allowed us to assess their combined impact on % DIABETIC rates, considering the joint influence of these two factors.
4. **Train-Test Split**: In all our analyses, we utilized the train-test split method. We divided the data into two parts, with 75% (265 data points) allocated to the train set and 25% (89 data points) to the test set. This approach ensured robust model evaluation.

5. **Model Performance Metrics**: We assessed model performance metrics, including Mean Absolute Error (MAE), Mean Squared Error (MSE), and Root Mean Squared Error (RMSE). These metrics provided insights into the accuracy and generalization capability of our multiple linear regression model.
6. **Heteroscedasticity Detection**: Our analysis detected heteroscedasticity in the model residuals, as indicated by the Breusch-Pagan test. This test identified variations in error variability across different levels of independent variables.
7. **T-Test**: We conducted t-tests on the coefficients of % Inactivity and % Obesity in the multiple linear regression model to assess their statistical significance in predicting % DIABETIC rates.
8. **R-Squared Values**: R-squared values were calculated to assess how well the independent variables in our regression models explain the variation in the dependent variable (% DIABETIC rates).

In summary, this methodology encompasses data cleaning, variable creation, and a range of statistical procedures such as correlation analysis, linear regression, train-test split, model performance evaluation, and tests for heteroscedasticity and statistical significance. These methods were crucial for our analysis of the relationships between diabetes, obesity, and physical inactivity.

## Appendix B- Results:

### **Correlation matrix:**



	% DIABETIC	%Obesity	%Inactivity
% DIABETIC	1.000000	0.389941	0.567104
%Obesity	0.389941	1.000000	0.472656
%Inactivity	0.567104	0.472656	1.000000

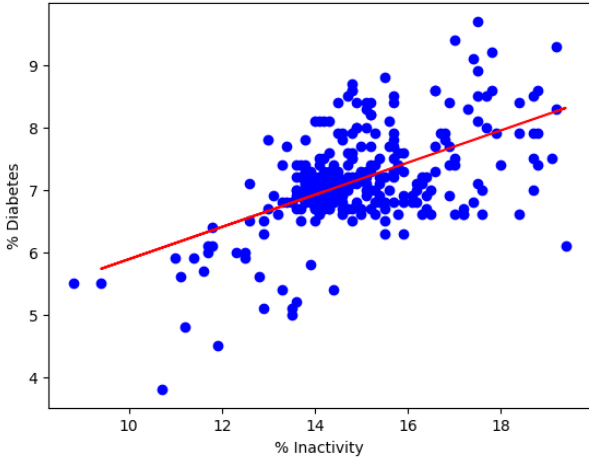
The correlation matrix reveals the following correlation coefficients:

% DIABETIC and % Obesity: 0.389941

% DIABETIC and % Inactivity: 0.567104

% Obesity and % Inactivity: 0.472656

**Simple Linear Regression (% Inactivity as Independent Variable):**

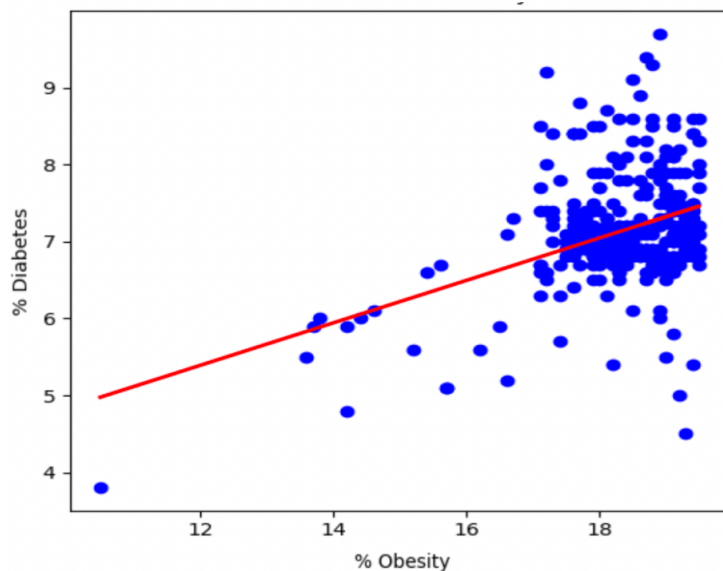


In this analysis, we used % Inactivity as the sole independent variable to predict % DIABETIC rates. The train set R-squared of 0.27 suggests that approximately 27% of the variation in % DIABETIC rates in the training data can be explained by % Inactivity alone. The train set R-squared of 0.27 suggests that

approximately 27% of the variation in % DIABETIC rates in the training data can be explained by % Inactivity alone.

The higher R-squared score in the test set compared to the train set is generally a positive sign. It suggests that the model based on % Inactivity generalizes well to new, unseen data, as it explains more variability in the test set. The R-squared score of 0.48 for the test set indicates that % Inactivity is a strong and meaningful predictor of % DIABETIC rates.

### **Simple Linear Regression (% Obesity as Independent Variable):**



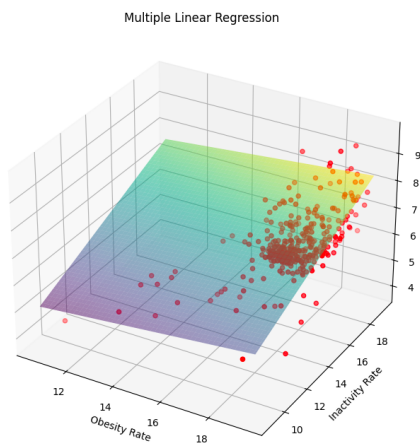
In this analysis, we used % Obesity as the sole independent variable to predict % DIABETIC rates. The test set R-squared of 0.30 indicates that approximately 30%

of the variation in % DIABETIC rates in the test data can be explained by % Obesity alone. The train set R-squared of 0.10 suggests that approximately 10% of the variation in % DIABETIC rates in the training data can be explained by % Obesity alone.

Similar to the % Inactivity model, the higher R-squared score in the test set compared to the train set is a positive sign. It suggests that the model based on % Obesity generalizes reasonably well to new data. The R-squared score of 0.30 for the test set indicates that % Obesity is a predictor of % DIABETIC rates, but its predictive power is weaker compared to % Inactivity.

This finding highlights the role of obesity in diabetes prevalence but suggests that other factors, such as physical inactivity, may have a stronger influence.

### **Multiple Linear Regression (Inactivity and Obesity as Independent Variables):**



In this analysis, you utilized both % Inactivity and % Obesity as independent variables in a multiple linear regression model to predict % DIABETIC rates. The test set R-squared of 0.52 indicates that approximately 52% of the variation in % DIABETIC rates in the test data can be explained by the combination of % Inactivity

and % Obesity. The train set R-squared of 0.28 suggests that approximately 28% of the variation in % DIABETIC rates in the training data can be explained by % Inactivity and % Obesity.

The multiple linear regression model outperforms the single linear regression models in terms of R-squared. It explains a higher proportion of the variability in % DIABETIC rates, both in the test and train sets. This suggests that considering both % Inactivity and % Obesity together improves the model's predictive power.

We have performed model performance metrics, including Mean Absolute Error (MAE), Mean Squared Error (MSE), and Root Mean Squared Error (RMSE), indicated lower errors in the test set compared to the train set, suggesting good generalization.

#### Breusch-Pagan Test:

The Breusch-Pagan test returned a small p-value of 0.0018, which is below the usual significance threshold of 0.05. When the p-value is less than this threshold, it indicates strong evidence to reject the null hypothesis. In this context, it signifies that the data exhibits heteroscedasticity.

#### T-test:

In the multiple linear regression model, we conducted t-tests on the coefficients of % Inactivity and % Obesity. The results indicated that both % Inactivity and % Obesity have statistically significant relationships with % DIABETIC rates. % Inactivity had a higher t-value (7.9) and a nearly zero p-value, underscoring its statistically significant effect on diabetes rates. % Obesity, while still statistically

significant, had a lower t-value (1.7). This confirms that both factors are meaningful predictors of % DIABETIC rates, with % Inactivity being a stronger predictor.

### Appendix C- Coding:

```
#Multiple linear regression
```

```
import pandas as pa
```

```
import numpy as np
```

```
from sklearn.model_selection import train_test_split
```

```
from sklearn import metrics
```

```
from sklearn.linear_model import LinearRegression
```

```
from sklearn.metrics import r2_score
```

```
import matplotlib.pyplot as plt
```

```
from scipy import stats
```

```
from statsmodels.stats.diagnostic import het_breuschpagan
```

```
from statsmodels.tools.tools import add_constant
```

```
import statsmodels.api as sm
```

```
from mpl_toolkits.mplot3d import Axes3D
```

```
# Load the data from the uploaded Excel file (assuming it's in Excel format)
```

```
# Adjust the filename and sheet name as needed
```

```
df=pa.read_csv("/Users/tysonmukesh/Desktop/MTH-522/Project-1.csv")
```

```
# Select the independent variables (Diabetes Rate and Inactivity Rate) and the dependent variable (Obesity Rate)
```

```
X = df[['%Obesity','%Inactivity']]
```

```
y = df['% DIABETIC']
```

```
# Split the data into training and testing sets
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25, random_state=25)
```

```
# Initialize and fit the linear regression model
```

```
model = LinearRegression()
```

```
model.fit(X_train, y_train)
```

```
# Coefficients
```

```
print('Coefficients:')
```

```
print('Obesity Rate:', model.coef_[0])
```

```
print('Inactivity Rate:', model.coef_[1])
```

```
print('Intercept:', model.intercept_)
```

```
# For Train set
```

```
y_pred = model.predict(X_train)
```

```
residuals = y_train - y_pred
```

```
r2=r2_score(y_train,y_pred)

# Perform the Breusch-Pagan test for heteroscedasticity

# Add a constant (intercept) to the predictors
X_train_with_constant = add_constant(X_train)
_, p_val, _, _ = het_breuschpagan(residuals, X_train_with_constant)
print('p-value from Breusch-Pagan test:', p_val)

X_train = sm.add_constant(X_train)

# Fit the model
model1 = sm.OLS(y_train, X_train).fit()

# Summary of the model
print(model1.summary())

# Mean Absolute Error (MAE)
mae = np.mean(np.abs(y_pred - y_train))
print("MAE for train set:", round(mae,2))

# Mean Squared Error (MSE)
mse = np.mean((y_pred - y_train) ** 2)
```

```
print("MSE for train set:", round(mse,2))

# Root Mean Squared Error (RMSE)
rmse = np.sqrt(mse)
print("RMSE for train set:", round(rmse,2))
print("r2 for train set is: "+str(round(r2,2)))

# For Test set

y_pred = model.predict(X_test)
r2=r2_score(y_test,y_pred)

# Add a constant (intercept) to the predictors
X_test = sm.add_constant(X_test)

# Fit the model
model2 = sm.OLS(y_test, X_test).fit()

# Summary of the model
print(model2.summary())

# Mean Absolute Error (MAE)
mae = np.mean(np.abs(y_pred - y_test))
print("MAE for test set:", round(mae,2))
```



```
# Mean Squared Error (MSE)
mse = np.mean((y_pred - y_test) ** 2)
print("MSE for test set:", round(mse,2))

# Root Mean Squared Error (RMSE)
rmse = np.sqrt(mse)
print("RMSE for test set:", round(rmse,2))
print("r2 for test set is: "+str(round(r2,2)))
```

### Contributions:

- **Mukesh Kumar Karanam Rameshbabu** – Worked on Findings, coding, Discussions, Methods and Results.
- **Sai Sudhamsh Kamisetty** – Worked on issues, coding, t-tests, graphs and results
- **Rohith Rasi Reddy** – Worked on initial cleaning of data using excel, coding and helped with the regression models.
- **Anish Krishna Kalisetti** – Worked on results and report preparation. And helped in Breusch-Pagan test.