

# Impact of Diabetes on Physical Activity, BMI, and Demographic Factors in a Large-scale Population Study

Data to Paper

January 12, 2024

## Abstract

Diabetes is a prevalent chronic health condition with significant implications for public health. Understanding the impact of diabetes on physical activity, body mass index (BMI), and demographic factors is crucial for effective interventions. However, there is a lack of comprehensive analysis in a large-scale population. This study aims to address this research gap by analyzing data from the Behavioral Risk Factor Surveillance System (BRFSS), a national survey conducted by the Centers for Disease Control and Prevention (CDC). Our analysis includes over 250,000 observations from 2015. We found distinct patterns in physical activity levels, BMI, age distribution, and sex proportions between individuals with and without diabetes. Additionally, we developed a multiple linear regression model to examine the impact of these factors on glycemic control among individuals with diabetes. Our results indicate that physical activity, BMI, age, and sex are significant predictors of glycemic control. These findings provide important insights into the influence of diabetes on physical activity, BMI, and demographic factors. They have significant implications for policy development and intervention design aimed at improving diabetes management and population health. Future research should focus on identifying novel strategies to address the underlying determinants of physical activity and BMI among individuals with diabetes.

## Introduction

Understanding the dynamic relationships between diabetes, demographic factors, body mass index (BMI), and physical activity is critical for designing effective population-level interventions and health policies [1]. Despite the

vast body of research focusing on diabetes as a significant global public health issue [2, 3], a comprehensive and large-scale analysis elucidating the interactions between diabetes and these variables remains sparse. This study aims to address this research gap.

Existing research underscores the importance of regular physical activity in managing glycemic control and decreasing diabetes-related comorbidities and cardiovascular risks [4]. Concurrently, elevated BMI has been linked to increased diabetes incidences, emphasizing the role of weight management in diabetes control [5, 6]. These studies, while informative, have often focused on these factors independently. Our understanding of how these variables collectively interact with diabetes, especially when evaluated through a demographic lens, is still incomplete.

In this study, we delve into this challenge by utilizing the rich dataset provided by the CDC’s Behavioral Risk Factor Surveillance System (BRFSS) 2015 survey [7, 8]. This large-scale, population-based study offers a unique perspective, enabling the investigation of broad patterns and interplay between diabetes, BMI, physical activity, and demographic components [7, 9].

Anchoring our methodology in detailed data preprocessing, we implement a multiple linear regression model using the ordinary least squares method [10, 11]. This approach allows us to rigorously examine the impact of physical activity, BMI, and demographic factors on glycemic control among individuals diagnosed with diabetes. Ultimately, our work provides valuable insights into these relationships, paving the way for strategic policy and intervention planning aimed at achieving better diabetes management and overall public health.

## Results

We conducted a comprehensive analysis using data from the Behavioral Risk Factor Surveillance System (BRFSS) 2015 survey to investigate the impact of diabetes on physical activity, body mass index (BMI), and demographic factors and to examine their impact on glycemic control among individuals with diabetes.

First, we examined the mean and standard deviation of physical activity, BMI, age category, and sex stratified by diabetes status (Table 1). Individuals with diabetes had a lower mean physical activity level (**0.631**) compared to those without diabetes (**0.777**). Furthermore, individuals with diabetes had a higher mean BMI (**31.9**) and were older on average (mean age category of **9.38**) compared to those without diabetes. There was also a slightly

higher proportion of males in the diabetes group (**0.479**) compared to the non-diabetes group (**0.434**). These findings highlight distinct patterns in physical activity levels, BMI, age distribution, and sex proportions between individuals with and without diabetes.

Table 1: Mean and Std Dev of P. Act., BMI, Age Cat., and Sex stratified by Diab. Status

Diabetes_Status	P. Act.		BMI		Age Cat.		Sex	
	Mean	Std Dev	Mean	Std Dev	Mean	Std Dev	Mean	Std Dev
<b>Diabetes</b>	0.631	0.483	31.9	7.36	9.38	2.33	0.479	0.5
<b>No Diabetes</b>	0.777	0.416	27.8	6.29	7.81	3.1	0.434	0.496

Mean values are likely to be altered due to approximations

**P. Act.:** Physical Activity in past 30 days (0=no, 1=yes)

**BMI:** Body Mass Index

**Age Cat.:** 13-level age category in intervals of 5 years (1=18-24, 2=25-29, ..., 12=75-79, 13=80 or older)

**Sex:** Sex (0=female, 1=male)

Next, we performed a multiple linear regression analysis to predict glycemic control among individuals with diabetes, adjusting for age, sex, and BMI (Table 2). The regression analysis revealed several significant predictors of glycemic control. Decreased physical activity (**coefficient = -0.436**) and higher BMI (**coefficient = 0.0148**) were associated with poorer glycemic control in individuals with diabetes. In addition, older age (**coefficient = -0.0174**) and being male (**coefficient = -0.0979**) were also associated with poorer glycemic control. These coefficients indicate the change in the dependent variable (glycemic control) for a one-unit change in the respective independent variables, while holding other variables constant. The intercept value (**3.3**) represents the expected glycemic control among individuals with diabetes who have a physical activity level of 0, BMI of 0, age of 0, and identify as female. The units of measurement for the variables in the regression model were not explicitly mentioned in the data, but they should be mentioned in the original dataset.

In summary, our analysis reveals important insights into the impact of diabetes on physical activity, BMI, and demographic factors. Lower physical activity levels and higher BMI were observed among individuals with diabetes, indicating potential areas for intervention and diabetes management strategies. The multiple linear regression model underscored the role of physical activity, BMI, age, and sex as significant predictors of glycemic

Table 2: Multiple Linear Regression Model predicting glycemic control among individuals with diabetes, adjusting for age, sex, and BMI.

index	Coeff.	Std Err	t-stat	P-val	CI Lower	CI Upper
<b>Intercept</b>	3.3	0.0376	87.9	$<10^{-6}$	3.23	3.38
<b>P. Act.</b>	-0.436	0.0109	-39.9	$<10^{-6}$	-0.457	-0.414
<b>BMI</b>	0.0148	0.000732	20.2	$<10^{-6}$	0.0134	0.0162
<b>Age Cat.</b>	-0.0174	0.00229	-7.61	$<10^{-6}$	-0.0219	-0.0129
<b>Sex</b>	-0.0979	0.0104	-9.38	$<10^{-6}$	-0.118	-0.0774

P-val denotes P-value for given coefficients. CI Lower & CI Upper are boundaries of 95% Confidence Interval.

**P. Act.:** Physical Activity in past 30 days (0=no, 1=yes)

**BMI:** Body Mass Index

**Age Cat.:** 13-level age category in intervals of 5 years (1=18-24, 2=25-29, ..., 12=75-79, 13=80 or older)

**Sex:** Sex (0=female, 1=male)

**Coeff.:** Coefficient Estimate

control among individuals with diabetes. These findings have significant implications for policy development and intervention design aimed at improving diabetes management and population health.

## Discussion

This analysis investigated the role of diabetes in physical activity, BMI, and demographic factors, recognizing these variables' and diabetes' joint significance as key public health determinants [1, 2, 3]. Our study utilized the formidable BRFSS 2015 dataset, implemented rigorous data preprocessing protocols, and structured a multiple linear regression model using the ordinary least squares method [10, 11].

Our findings confirmed an inverse relationship between diabetes and physical activity levels, and an increased BMI amongst diabetic individuals, mirroring existing studies that have made similar correlations [4, 5, 12]. Furthermore, we found that age and sex have significant impacts on glycemic control, with older individuals and males with diabetes exhibiting poorer control [5]. These findings re-emphasize the need for interventions that specifically address these stratifiers.

However, this study is not immune to limitations: the nature of self-reported health data used here may introduce personal bias or recording

inaccuracies. Also, in the absence of follow-up data, it is challenging to assess how changes in physical activity, BMI, and demographic factors over time can influence glycemic control. The dataset predominantly centers U.S. participants, limiting the broader applicability of the results in a global context. Moreover, the study could not bear out other potential influences like dietary habits or genetic predispositions due to data unavailability.

Nonetheless, our study has laid down clear evidence supporting the correlation between diabetes and factors like physical activity, BMI, and certain demographics. These insights have implications for facilitating more personalized diabetic interventions, taking into consideration factors such as age, sex, BMI, and physical activity level, as shown by our findings. The evidence provided may inform more strategic policy-making and intervention designs, especially where diabetes management and public health are concerned, further substantiating the potential of such large-scale population surveys like BRFSS [6].

For future research, we suggest a holistic approach, involving diet, genetics, and a wider range of demographic variables in the analysis. Furthermore, a longitudinal design capturing data over extended periods may offer a more dynamic understanding of diabetes and the phenomenon's interaction with key health and demographic factors.

## Methods

### Data Source

The data used in this study was derived from the Behavioral Risk Factor Surveillance System (BRFSS), a health-related telephone survey conducted annually by the Centers for Disease Control and Prevention (CDC) in the United States. The dataset utilized for this analysis contained information related to diabetes and various health indicators collected in the year 2015. The original BRFSS dataset collects responses from over 400,000 Americans on health-related risk behaviors, chronic health conditions, and the use of preventative services.

### Data Preprocessing

The data preprocessing steps were performed using Python programming language. First, the dataset was loaded into memory using the pandas library. The data file consisted of 253,680 rows and 22 columns representing different variables related to diabetes and health indicators. It is important

to note that any rows with missing values were removed from the dataset prior to analysis.

## **Data Analysis**

The analysis was conducted to investigate the association between physical activity levels and glycemic control among individuals with diabetes, considering potential confounding factors such as age, sex, and BMI. The analysis steps can be summarized as follows:

1. **Descriptive Statistics:** To gain initial insights into the dataset, descriptive statistics were calculated for physical activity levels, BMI, age, and sex, stratified by diabetes incidence. The mean and standard deviation of these variables were computed separately for individuals with and without diabetes.

2. **Multiple Linear Regression Model:** A multiple linear regression model was developed to examine the impact of physical activity, BMI, age, and sex on glycemic control among individuals with diabetes. The model considered glycemic control, as measured by self-reported health, as the dependent variable. Physical activity levels, BMI, age, and sex were included as independent variables. The model was fitted using the ordinary least squares (OLS) method from the statsmodels library. Model coefficients, standard errors, p-values, and R-squared values were extracted from the model summary.

The use of these analysis steps allowed us to explore the relationship between physical activity levels and glycemic control among individuals with diabetes, while accounting for potential confounding factors.

## **Code Availability**

Custom code used to perform the data preprocessing and analysis, as well as the raw code outputs, are provided in Supplementary Methods.

## A Data Description

Here is the data description, as provided by the user:

The dataset includes diabetes related factors extracted from the CDC's Behavioral Risk Factor Surveillance System (BRFSS), year 2015.

The original BRFSS, from which this dataset is derived, is a health-related telephone survey that is collected annually by the CDC.

Each year, the survey collects responses from over 400,000 Americans on health-related risk behaviors, chronic health conditions, and the use of preventative services. These features are either questions directly asked of participants, or calculated variables based on individual participant responses.

1 data file:

"diabetes\_binary\_health\_indicators\_BRFSS2015.csv"

The csv file is a clean dataset of 253,680 responses (rows) and 22 features (columns).

All rows with missing values were removed from the original dataset; the current file contains no missing values.

The columns in the dataset are:

- #1 `Diabetes\_binary`: (int, bool) Diabetes (0=no, 1=yes)
- #2 `HighBP`: (int, bool) High Blood Pressure (0=no, 1=yes)
- #3 `HighChol`: (int, bool) High Cholesterol (0=no, 1=yes)
- #4 `CholCheck`: (int, bool) Cholesterol check in 5 years (0=no, 1=yes)
- #5 `BMI`: (int, numerical) Body Mass Index
- #6 `Smoker`: (int, bool) (0=no, 1=yes)
- #7 `Stroke`: (int, bool) Stroke (0=no, 1=yes)
- #8 `HeartDiseaseorAttack`: (int, bool) coronary heart disease (CHD) or myocardial infarction (MI), (0=no, 1=yes)
- #9 `PhysActivity`: (int, bool) Physical Activity in past 30 days (0=no, 1=yes)
- #10 `Fruits`: (int, bool) Consume one fruit or more each day (0=no, 1=yes)
- #11 `Veggies`: (int, bool) Consume one Vegetable or more each day (0=no, 1=yes)
- #12 `HvyAlcoholConsump`: (int, bool) Heavy drinkers (0=no, 1=yes)
- #13 `AnyHealthcare`: (int, bool) Have any kind of health care coverage (0=no, 1=yes)
- #14 `NoDocbcCost`: (int, bool) Was there a time in the past 12 months when you

needed to see a doctor but could not because of cost? (0=no, 1=yes)

#15 ``GenHlth`` (int, ordinal) self-reported health (1=excellent, 2=very good, 3=good, 4=fair, 5=poor)

#16 ``MentHlth`` (int, ordinal) How many days during the past 30 days was your mental health not good? (1-30 days)

#17 ``PhysHlth`` (int, ordinal) How many days during the past 30 days was your physical health not good? (1-30 days)

#18 ``DiffWalk`` (int, bool) Do you have serious difficulty walking or climbing stairs? (0=no, 1=yes)

#19 ``Sex`` (int, categorical) Sex (0=female, 1=male)

#20 ``Age`` (int, ordinal) Age, 13-level age category in intervals of 5 years (1=18-24, 2=25-29, ..., 12=75-79, 13=80 or older)

#21 ``Education`` (int, ordinal) Education level on a scale of 1-6 (1=Never attended school, 2=Elementary, 3=Some high school, 4=High school, 5=Some college, 6=College)

#22 ``Income`` (int, ordinal) Income scale on a scale of 1-8 (1=<=10K, 2=<=15K, 3=<=20K, 4=<=25K, 5=<=35K, 6=<=50K, 7=<=75K, 8=>75K)

## B Data Exploration

### B.1 Code

The Data Exploration was carried out using the following custom code:

```

1
2 import pandas as pd
3
4 # Load dataset
5 data =
  ↪ pd.read_csv('diabetes_binary_health_indicators_BRFSS2015.csv')
6
7 # Open output file
8 with open('data_exploration.txt', 'w') as f:
9
10     ## Data size
11     f.write('# Data Size\n')
12     f.write('Number of Rows: {}\n'.format(data.shape[0]))
13     f.write('Number of Columns: {}\n\n'.format(data.shape[1]))

```



```

14
15     ## Summary Statistics for key numerical variables
16     keys = ['BMI', 'GenHlth', 'MentHlth', 'PhysHlth', 'Age',
17             ↪ 'Education', 'Income']
18     summary = data[keys].describe()
19
20     f.write('# Summary Statistics\n')
21     f.write('{}\n\n'.format(summary))
22
23     ## Categorical Variables
24     f.write('# Categorical Variables\n')
25     categorical_columns =
26     ↪ data.select_dtypes(include=['object', 'bool']).columns
27
28     for column in categorical_columns:
29         most_common = data[column].value_counts().idxmax()
30         f.write('{}: Most common value is {}\n'.format(column,
31             ↪ most_common))
32
33     f.write('\n')
34
35     ## Missing Values
36     missing_values = data.isnull().sum()
37     f.write('# Missing Values\n')
38     f.write('{}\n\n'.format(missing_values))
39
40     # Unknown/undefined values
41     f.write('# Unknown/undefined values\n')
42     for column in data.columns:
43         count = data[data[column] == 99].shape[0] # assuming
44         ↪ 99 is the unknown/undefined value based on the
45         ↪ dataset description
46         f.write('{}: Number of unknown/undefined values:
47             ↪ {}\n'.format(column, count))
48
49     f.close()
50

```

## B.2 Code Description

The provided code performs a data exploration analysis on the diabetes related factors dataset. The analysis involves several steps to gain insights into the dataset's characteristics and identify any patterns or issues within the data.

First, the code loads the dataset into a Pandas DataFrame. Then, it opens an output file to store the results of the data exploration analysis.

The code starts by calculating the number of rows and columns in the dataset, which reflects the size of the data.

Next, summary statistics are computed for key numerical variables, including BMI (Body Mass Index), self-reported health (GenHlth), number of days with poor mental health (MentHlth), number of days with poor physical health (PhysHlth), age, education level, and income. These statistics provide an overview of the distribution of these variables in terms of measures such as mean, standard deviation, minimum, maximum, and quartiles.

The code then focuses on categorical variables in the dataset. It identifies the most common value for each categorical column, providing insights into the prevalence of different values within each category.

The analysis also checks for missing values in the dataset. The code calculates the count of missing values for each column and reports it in the output file. This information helps identify the completeness of the dataset and any potential issues related to missing data.

Finally, the code identifies unknown or undefined values in the dataset. It iterates over each column and counts the occurrences of a specific value (assumed to be 99 in this case) that indicates an unknown or undefined value. This information can be crucial in understanding the quality and reliability of the data.

The code writes all the analysis results into the "data\_exploration.txt" file. The file includes information such as the data size, summary statistics for numerical variables, most common values for categorical variables, counts of missing values, and counts of unknown/undefined values. This file serves as a comprehensive report summarizing the findings of the data exploration analysis.

Overall, the code provides useful insights into the dataset's characteristics, which can inform further analysis and decision-making based on the data.

### B.3 Code Output

#### data\_exploration.txt

##### # Data Size

Number of Rows: 253680

Number of Columns: 22

##### # Summary Statistics

	BMI	GenHlth	MentHlth	PhysHlth	Age	Education	Income
count	253680	253680	253680	253680	253680	253680	253680
mean	28.38	2.511	3.185	4.242	8.032	5.05	6.054
std	6.609	1.068	7.413	8.718	3.054	0.9858	2.071
min	12	1	0	0	1	1	1
25%	24	2	0	0	6	4	5
50%	27	2	0	0	8	5	7
75%	31	3	2	3	10	6	8
max	98	5	30	30	13	6	8

##### # Categorical Variables

##### # Missing Values

Diabetes_binary	0
HighBP	0
HighChol	0
CholCheck	0
BMI	0
Smoker	0
Stroke	0
HeartDiseaseorAttack	0
PhysActivity	0
Fruits	0
Veggies	0
HvyAlcoholConsump	0
AnyHealthcare	0
NoDocbcCost	0
GenHlth	0
MentHlth	0
PhysHlth	0
DiffWalk	0

```
Sex          0
Age          0
Education    0
Income       0
dtype: int64
```

```
# Unknown/undefined values
```

```
Diabetes_binary: Number of unknown/undefined values: 0
HighBP: Number of unknown/undefined values: 0
HighChol: Number of unknown/undefined values: 0
CholCheck: Number of unknown/undefined values: 0
BMI: Number of unknown/undefined values: 0
Smoker: Number of unknown/undefined values: 0
Stroke: Number of unknown/undefined values: 0
HeartDiseaseorAttack: Number of unknown/undefined values: 0
PhysActivity: Number of unknown/undefined values: 0
Fruits: Number of unknown/undefined values: 0
Veggies: Number of unknown/undefined values: 0
HvyAlcoholConsump: Number of unknown/undefined values: 0
AnyHealthcare: Number of unknown/undefined values: 0
NoDocbcCost: Number of unknown/undefined values: 0
GenHlth: Number of unknown/undefined values: 0
MentHlth: Number of unknown/undefined values: 0
PhysHlth: Number of unknown/undefined values: 0
DiffWalk: Number of unknown/undefined values: 0
Sex: Number of unknown/undefined values: 0
Age: Number of unknown/undefined values: 0
Education: Number of unknown/undefined values: 0
Income: Number of unknown/undefined values: 0
```

## C Data Analysis

### C.1 Code

The Data Analysis was carried out using the following custom code:

```
1
2 # IMPORT
3 import pandas as pd
```

```

4 import numpy as np
5 import pickle
6 import statsmodels.api as sm
7
8 # LOAD DATA
9 df =
    ↪ pd.read_csv("diabetes_binary_health_indicators_BRFSS2015.csv")
10
11 # DATASET PREPARATIONS
12 # No dataset preparations are needed.
13
14 # DESCRIPTIVE STATISTICS
15 ## Table 0: "Mean and STD of physical activity, BMI, age and
    ↪ sex stratified by diabetes incidence"
16 df["Diabetes_Status"] = df["Diabetes_binary"].map({0: 'No
    ↪ Diabetes', 1: 'Diabetes'})
17
18 grouped = df.groupby(by="Diabetes_Status")
19 df0 = grouped[["PhysActivity", "BMI", "Age",
    ↪ "Sex"]].agg([np.mean, np.std])
20
21 df0.to_pickle('table_0.pkl')
22
23 # PREPROCESSING
24 # No preprocessing is needed, because of the simplicity of the
    ↪ model and the operations required.
25
26 # ANALYSIS
27 ## Table 1: "Multiple Linear Regression Model predicting
    ↪ glycemic control among individuals with diabetes,
    ↪ adjusting for age, sex, and BMI."
28 diabetes_df = df[df["Diabetes_binary"]==1]
29
30 X = diabetes_df[["PhysActivity", "BMI", "Age", "Sex"]] #
    ↪ independent variables
31 X = sm.add_constant(X) # adding a constant
32 Y = diabetes_df["GenHlth"] # dependent variable
33
34 model = sm.OLS(Y, X).fit()

```

```

35 df1 = pd.DataFrame(model.summary2().tables[1]) # Extracting
    ↳ the model estimates table
36 df1.reset_index(inplace=True)
37 df1['index'] = df1['index'].astype(str) # Ensure index is not
    ↳ purely numeric
38 df1.set_index('index', inplace=True)
39 df1.to_pickle('table_1.pkl')
40
41 # SAVE ADDITIONAL RESULTS
42 additional_results = {
43     'Total number of observations': df.shape[0],
44     'R-squared of regression model': model.rsquared
45 }
46
47 with open('additional_results.pkl', 'wb') as f:
48     pickle.dump(additional_results, f)
49

```

## C.2 Code Description

The code performs an analysis on a dataset that includes diabetes related factors extracted from the CDC's Behavioral Risk Factor Surveillance System (BRFSS) for the year 2015. The dataset contains information on various health indicators such as diabetes status, high blood pressure, high cholesterol, physical activity, body mass index (BMI), and other demographic characteristics.

First, the code loads the dataset into a pandas DataFrame. The dataset is a clean version of the original BRFSS dataset, with missing values removed.

Next, the code conducts descriptive statistics on physical activity, BMI, age, and sex stratified by diabetes incidence. The code groups the data by diabetes status and calculates the mean and standard deviation for each variable of interest. The results are stored in a DataFrame called "df0" and saved as a pickle file named "table\_0.pkl". This table provides an overview of the average physical activity, BMI, age, and sex for individuals with and without diabetes.

After preparing the dataset, the code performs an analysis using a multiple linear regression model to predict glycemic control among individuals with diabetes. The independent variables used in the model are physical activity, BMI, age, and sex. A constant term is added to the independent

variables using the "sm.add\_constant()" function from the statsmodels library. The dependent variable is self-reported general health status. The code fits the regression model using the ordinary least squares (OLS) method provided by the statsmodels library.

The code generates a summary of the model using the "summary2()" function from the statsmodels library and extracts the table containing the model estimates. This table, which includes coefficients, standard errors, t-values, and p-values, is stored in a DataFrame called "df1". The DataFrame is then saved as a pickle file named "table\_1.pkl". This table provides insights into the relationship between physical activity, BMI, age, sex, and glycemic control among individuals with diabetes.

Finally, the code saves additional results in a pickle file named "additional\_results.pkl". These results include the total number of observations in the dataset and the R-squared value of the regression model. The R-squared value represents the proportion of the variance in the dependent variable that is explained by the independent variables in the regression model.

In conclusion, the code performs data analysis to explore the relationship between various health indicators, demographic factors, and glycemic control among individuals with diabetes. The descriptive statistics and regression results provide valuable insights into the associations between these factors, which can inform future research, interventions, and public health initiatives related to diabetes management and prevention.

### C.3 Code Output

**table\_0.pkl**

	PhysActivity		BMI		Age		Sex	
	mean	std	mean	std	mean	std	mean	std
Diabetes_Status								
Diabetes	0.6305	0.4827	31.94	7.363	9.379	2.33	0.4791	0.4996
No Diabetes	0.7769	0.4163	27.81	6.291	7.814	3.101	0.4341	0.4956

**table\_1.pkl**

	Coef.	Std.Err.	t	P> t	[0.025	0.975]
index						
const	3.303	0.03758	87.9	0	3.229	3.377
PhysActivity	-0.4355	0.01091	-39.93	0	-0.4569	-0.4141
BMI	0.0148	0.0007316	20.23	1.851e-90	0.01336	0.01623

Age	-0.01741	0.002287	-7.612	2.763e-14	-0.0219	-0.01293
Sex	-0.09791	0.01044	-9.376	7.226e-21	-0.1184	-0.07744

additional\_results.pkl

```
{
    'Total number of observations': 253680,
    'R-squared of regression model': 0.06896
}
```

## D LaTeX Table Design

### D.1 Code

The LaTeX Table Design was carried out using the following custom code:

```
1
2 # IMPORT
3 import pandas as pd
4 from typing import Dict, Tuple, Optional
5 from my_utils import to_latex_with_note, format_p_value
6
7 Mapping = Dict[str, Tuple[Optional[str], Optional[str]]]
8
9 def split_mapping(d: Mapping):
10     abbrs_to_names = {abbr: name for abbr, (name, definition)
11                       ↪ in d.items() if name is not None}
12     names_to_definitions = {name or abbr: definition for abbr,
13                             ↪ (name, definition) in d.items() if definition is not
14                             ↪ None}
15     return abbrs_to_names, names_to_definitions
16
17 # PREPARATION FOR ALL TABLES
18 shared_mapping: Mapping = {
19     'PhysActivity': ('P. Act.', 'Physical Activity in past 30
20     ↪ days (0=no, 1=yes)'),
21     'BMI': ('BMI', 'Body Mass Index'),
22     'Age': ('Age Cat.', '13-level age category in intervals of
23     ↪ 5 years (1=18-24, 2=25-29, ..., 12=75-79, 13=80 or
24     ↪ older)'),
```



```

19     'Sex': ('Sex', 'Sex (0=female, 1=male)'),
20     'Diabetes_Status': ('Diab. Status', 'Diabetes (0=no,
    ↪ 1=yes)'),
21 }
22
23 # TABLE 0:
24 df = pd.read_pickle('table_0.pkl')
25
26 # Renaming Abbreviated Columns and Rows
27 mapping = {k: v for k, v in shared_mapping.items() if k in
    ↪ df.columns or k in df.index}
28 mapping |= {
29     'mean': ('Mean', None),
30     'std': ('Std Dev', None),
31 }
32 abbrs_to_names, legend = split_mapping(mapping)
33 df = df.rename(columns=abbrs_to_names, index=abbrs_to_names)
34
35 # Save as latex:
36 to_latex_with_note(
37     df,
38     'table_0.tex',
39     caption="Mean and Std Dev of P. Act., BMI, Age Cat., and Sex
    ↪ stratified by Diab. Status",
40     label='table:table_0',
41     note="Mean values are likely to be altered due to
    ↪ approximations",
42     legend=legend)
43
44 # TABLE 1:
45 df = pd.read_pickle('table_1.pkl')
46
47 # Formatting P-Values
48 df['P>|t|'] = df['P>|t|'].apply(format_p_value)
49
50 # Renaming Abbreviated Columns and Rows
51 mapping = {k: v for k, v in shared_mapping.items() if k in
    ↪ df.columns or k in df.index}
52 mapping |= {
53     'const': ('Intercept', None),

```

```

54   'Coef.': ('Coeff.', 'Coefficient Estimate'),
55   'Std.Err.': ('Std Err', None),
56   't': ('t-stat', None),
57   'P>|t|': ('P-val', None),
58   '[0.025]': ('CI Lower', None),
59   '[0.975]': ('CI Upper', None),
60 }
61 abbrs_to_names, legend = split_mapping(mapping)
62 df = df.rename(columns=abbrs_to_names, index=abbrs_to_names)
63
64 # Save as latex:
65 to_latex_with_note(
66     df,
67     'table_1.tex',
68     caption="Multiple Linear Regression Model predicting glycemic
        ↪ control among individuals with diabetes, adjusting for
        ↪ age, sex, and BMI.",
69     label='table:table_1',
70     note="P-val denotes P-value for given coefficients. CI Lower
        ↪ & CI Upper are boundaries of 95% Confidence Interval.",
71     legend=legend)
72

```

## D.2 Code Output

table\_0.tex

```

\begin{table}[h]
\caption{Mean and Std Dev of P. Act., BMI, Age Cat., and Sex stratified by Diab.
        Status}
\label{table:table_0}
\begin{threeparttable}
\renewcommand{\TPTminimum}{\linewidth}
\makebox[\linewidth]{%
\begin{tabular}{lrrrrrrrr}
\toprule
& \multicolumn{2}{r}{P. Act.} & \multicolumn{2}{r}{BMI} & &
\multicolumn{2}{r}{Age Cat.} & \multicolumn{2}{r}{Sex} \\
& Mean & Std Dev & Mean & Std Dev & Mean & Std Dev & Mean & Std Dev & \\
Diabetes\_Status & & & & & & & & & \\

```

```

\midrule
\textbf{Diabetes} & 0.631 & 0.483 & 31.9 & 7.36 & 9.38 & 2.33 & 0.479 & 0.5 \\
\textbf{No Diabetes} & 0.777 & 0.416 & 27.8 & 6.29 & 7.81 & 3.1 & 0.434 & 0.496 \\
\\
\bottomrule
\end{tabular}}
\begin{tablenotes}
\footnotesize
\item Mean values are likely to be altered due to approximations
\item \textbf{P. Act.}: Physical Activity in past 30 days (0=no, 1=yes)
\item \textbf{BMI}: Body Mass Index
\item \textbf{Age Cat.}: 13-level age category in intervals of 5 years (1=18-24,
2=25-29, ..., 12=75-79, 13=80 or older)
\item \textbf{Sex}: Sex (0=female, 1=male)
\end{tablenotes}
\end{threeparttable}
\end{table}

```

table\_1.tex

```

\begin{table}[h]
\caption{Multiple Linear Regression Model predicting glycemic control among
individuals with diabetes, adjusting for age, sex, and BMI.}
\label{table:table_1}
\begin{threeparttable}
\renewcommand{\TPTminimum}{\linewidth}
\makebox[\linewidth]{%
\begin{tabular}{lrrrrlrr}
\toprule
& Coeff. & Std Err & t-stat & P-val & CI Lower & CI Upper \\
index & & & & & & \\
\midrule
\textbf{Intercept} & 3.3 & 0.0376 & 87.9 &  $<1e-06$  & 3.23 & 3.38 \\
\textbf{P. Act.} & -0.436 & 0.0109 & -39.9 &  $<1e-06$  & -0.457 & -0.414 \\
\textbf{BMI} & 0.0148 & 0.000732 & 20.2 &  $<1e-06$  & 0.0134 & 0.0162 \\
\textbf{Age Cat.} & -0.0174 & 0.00229 & -7.61 &  $<1e-06$  & -0.0219 & -0.0129 \\
\textbf{Sex} & -0.0979 & 0.0104 & -9.38 &  $<1e-06$  & -0.118 & -0.0774 \\
\bottomrule
\end{tabular}}
\end{threeparttable}

```

```

\begin{tablenotes}
\footnotesize
\item P-val denotes P-value for given coefficients. CI Lower \& CI Upper are
    boundaries of 95\% Confidence Interval.
\item \textbf{P. Act.}: Physical Activity in past 30 days (0=no, 1=yes)
\item \textbf{BMI}: Body Mass Index
\item \textbf{Age Cat.}: 13-level age category in intervals of 5 years (1=18-24,
    2=25-29, ..., 12=75-79, 13=80 or older)
\item \textbf{Sex}: Sex (0=female, 1=male)
\item \textbf{Coeff.}: Coefficient Estimate
\end{tablenotes}
\end{threeparttable}
\end{table}

```

## References

- [1] F. Hill-Briggs, N. Adler, Seth A. Berkowitz, M. Chin, T. Gary-Webb, A. Navas-Acien, Pamela L. Thornton, and D. Haire-Joshu. Social determinants of health and diabetes: A scientific review. *Diabetes Care*, 44:258 – 279, 2020.
- [2] Moien A. B. Khan, M. J. Hashim, J. King, R. Govender, H. Mustafa, and J. Al Kaabi. Epidemiology of type 2 diabetes global burden of disease and forecasted trends. *Journal of Epidemiology and Global Health*, 10:107 – 111, 2019.
- [3] D. Mozaffarian. Dietary and policy priorities for cardiovascular disease, diabetes, and obesity: A comprehensive review. *Circulation*, 133:187225, 2016.
- [4] B. Bohn, A. Herbst, M. Pfeifer, D. Krakow, S. Zimny, F. Kopp, A. Melmer, J. Steinacker, and R. Holl. Impact of physical activity on glycemic control and prevalence of cardiovascular risk factors in adults with type 1 diabetes: A cross-sectional multicenter study of 18,028 patients. *Diabetes Care*, 38:1536 – 1543, 2015.
- [5] Ying Chen, Xiao-Ping Zhang, Jie Yuan, Bo zhi Cai, Xiao-Li Wang, Xiao li Wu, Yueqi Zhang, Xiao-Yi Zhang, T. Yin, Xiaohui Zhu, Y. Gu, S. Cui, Zhidong Lu, and Xiao ying Li. Association of body mass index

and age with incident diabetes in chinese adults: a population-based cohort study. *BMJ Open*, 8, 2018.

- [6] T. Jafar, N. Chaturvedi, and G. Pappas. Prevalence of overweight and obesity and their association with hypertension and diabetes mellitus in an indo-asian population. *Canadian Medical Association Journal*, 175:1071 – 1077, 2006.
- [7] Masayuki Kato, Mitsuhiro Noda, T. Mizoue, A. Goto, Yoshihiko Takahashi, Y. Matsushita, A. Nanri, H. Iso, M. Inoue, N. Sawada, and S. Tsugane. Diagnosed diabetes and premature death among middle-aged japanese: results from a large-scale population-based cohort study in japan (jphc study). *BMJ Open*, 5, 2015.
- [8] Juhua Luo, M. Iwasaki, M. Inoue, S. Sasazuki, T. Otani, W. Ye, S. Tsugane, and for the Japan Public Health Center-based Prospective S Group. Body mass index, physical activity and the risk of pancreatic cancer in relation to smoking status and history of diabetes: a large-scale population-based cohort study in japan the jphc study. *Cancer Causes & Control*, 18:603–612, 2007.
- [9] Marc A. Pitasi, Emeka Oraka, H. Clark, Machell Town, and E. Dinunno. Hiv testing among transgender women and men 27 states and guam, 2014–2015. *Morbidity and Mortality Weekly Report*, 66:883 – 887, 2017.
- [10] L. D. Robinson and N. Jewell. Some surprising results about covariate adjustment in logistic regression models. *International Statistical Review*, 59:227–240, 1991.
- [11] J. Faraway. *Linear models with r*. 2014.
- [12] J. Batsis, A. Zbehlik, L. Barre, J. Bynum, D. Pidgeon, and S. Bartels. Impact of obesity on disability, function, and physical activity: data from the osteoarthritis initiative. *Scandinavian Journal of Rheumatology*, 44:495 – 502, 2015.