

Diabetes Prediction model: NHANES data 2013-2014

Hannah Rosenblum, James Ng, Purnima Sharma

Contents

Introduction	2
EDA	2

Introduction

This project aimed to study any association between diabetes and several covariates in participants, ages 12 and older, using NHANES data. Effort was made to build an optimal prediction model to identify key variables within the dataset for a categorical outcome of diabetes, objective being a binary classification model with supervised learning. Certain factors of special interest were any association with participant's race, age, cholesterol and fasting glucose levels, among others. Specifically, association was assessed between diabetes and plasma fasting glucose, body mass index, blood pressure, hours of fasting, weekly workload, cholesterol, triglycerides, insulin levels, 2 hour glucose tolerance, glucose challenge and other demographic factors such as age, gender, race, using the NHANES data for the year 2013 - 2014 from the cdc.gov website, <https://wwwn.cdc.gov/nchs/nhanes/continuousnhanes/default.aspx?BeginYear=2013> .

Motivation was provided by the fact that diabetes is one of the major leading causes of death in the United States. As stated by the CDC site's National Diabetes Statistics Report of 2020, 34.2 million Americans are diabetic, while 7.3 million were undiagnosed. With prevalence of diabetes and prediabetes on the rise, it was of interest to find factors that might affect the diabetes status. Later years post-2013 were tried for the data, however were unavailable for the variables of interest possibly due to continuing updates.

After extracting and merging the necessary files by participant's Id number, variables of interest were retained in a dataframe, and saved as a dataset file within the project for easy access. Gender, Race and the response variable Diabetes were converted to factors from numeric data type. Missing entries for the response of diabetes status were removed. 185 "borderline" reported cases, 5 with "don't know" responses and 1 with "refused" response were also removed given the small scale of these categories, which accounted for less than 2% of the data, and in order to focus on the majority of binary responses of presence or absence of diabetes. The cleaned data contained 9,578 observations of 15 variables, including the binary outcome variable diabetes.

EDA

...

Data Frame Summary

raw_data
Dimensions: 9578 x 15
Duplicates: 506

No	Variable	Stats / Values	Freqs (% of Valid)	Graph	Missing
1	gender [factor]	1. 1 2. 2	4706 (49.1%) 4872 (50.9%)	IIIIIIII IIIIIIII	0 (0.0%)
2	age [numeric]	Mean (sd) : 32.4 (23.9) min < med < max: 1 < 28 < 80 IQR (CV) : 41 (0.7)	80 distinct values	: : : : : . : : : : : : : : : : : : : : : :	0 (0.0%)
3	race [factor]	1. 1 2. 2 3. 3 4. 4 5. 6 6. 7	1616 (16.9%) 893 (9.3%) 3449 (36.0%) 2148 (22.4%) 1033 (10.8%) 439 (4.6%)	III I IIIIII IIII II II	0 (0.0%)

No	Variable	Stats / Values	Freqs (% of Valid)	Graph	Missing
14	glucose [numeric]	Mean (sd) : 104.6 (31.2) min < med < max: 51 < 98 < 421 IQR (CV) : 15 (0.3)	184 distinct values	: : : . : : : : .	6490 (67.8%)
15	diabetes [factor]	1. yes 2. no	737 (7.7%) 8841 (92.3%)	I IIIIIIIIIIIIIIIIIIII	0 (0.0%)