

Diabetes Prediction model: NHANES data 2013-2014

Hannah Rosenblum, James Ng, Purnima Sharma

Contents

Introduction	2
EDA	2
summary statistics	3
Density plots	5
Bar plots	6
Missing data	6
Models	6
Linear models	6
Non-linear models	6
Ensemble methods	6
Support vector machines	7
Model comparison	7
Final model	7
Model prediction performance	7
Limitations	7
Conclusion	7

Introduction

This project aimed to study any association between diabetes and several covariates in participants ages 1 and older, using NHANES data, and selecting an optimal prediction model among linear, non-linear, parametric and non-parametric models. The main objective was building a binary classification model with supervised learning. Certain factors of special interest were any association with participant's race, age, cholesterol and lifestyle factors, among others. Data was extracted for the year 2013 - 2014 from the cdc.gov website, <https://wwwn.cdc.gov/nchs/nhanes/continuousnhanes/default.aspx?BeginYear=2013>. Specifically, association was assessed between diabetes and the following covariates:

- Gender: Participant's gender (male or female)
- Age: Age at screening, with possible values of 0 to 79, or 80+ (years)
- Race: 6 categories for race include Mexican American, other Hispanic, White, Black, Asian and other.
- bmi: body mass index (kg/m^2)
- hdl: High-density lipoprotein (mg/dL)
- ldl: Low-density lipoprotein (mg/dL)
- Triglycerides (mg/dL): laboratory test results for serum levels of triglycerides
- Insulin (uU/mL): measured using serum insulin methods
- Glucose (mg/dL): plasma glucose value measured 2 hours after calibrated oral dose
- Blood pressure (mm Hg): Both systolic and diastolic, first-round measurements
- waist: Waist circumference measurement (cm)
- Sedentary activity (lifestyle, minutes): time spent sitting in a given day, not including sleeping.
- Education level: highest degree of adults 20+ years of age, with 7 categories.
- Marital status: Categories include married, widowed, divorced, separated, never married, living with partner, refused, and don't know
- Depression: severity on a scale of 0 to 3 treated as a continuous variable, with 0 as not at all depressed
- Sleep: amount of sleep in hours on a given night on weekdays or workdays

The outcome of "diabetes" dependent-variable was based on classification of the participants into two groups of those with diabetes and those who did not have diabetes. Individuals answered the question "other than during pregnancy, have you ever been told by a doctor or health professional that you have diabetes or sugar diabetes?", and were classified as having diabetes if they answered yes.

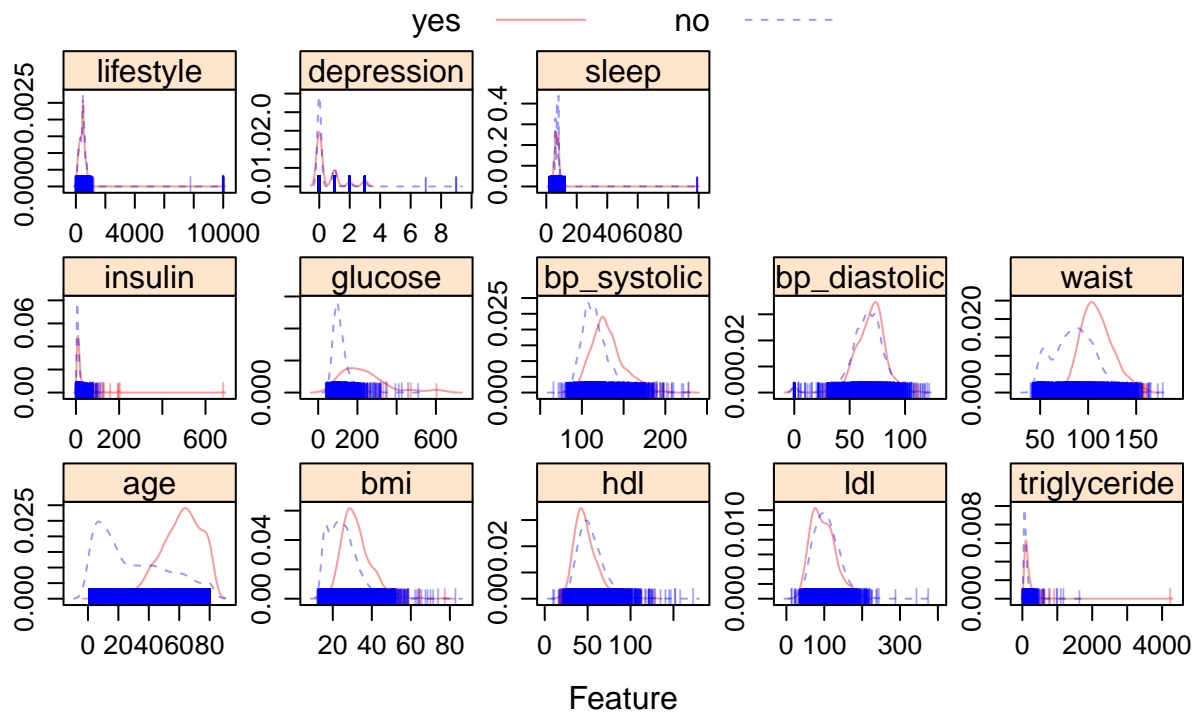
Motivation was provided by the fact that diabetes is one of the major leading causes of death in the United States. As stated by the CDC site's National Diabetes Statistics Report of 2020, 34.2 million Americans are diabetic, while 7.3 million were undiagnosed. Furthermore, increase in type 2 diabetes among children is a growing concern according to the CDC. 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. Gender, race, education level, marital status 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 dataset contained 9,578 observations of 15 variables, including the binary outcome variable diabetes.

EDA

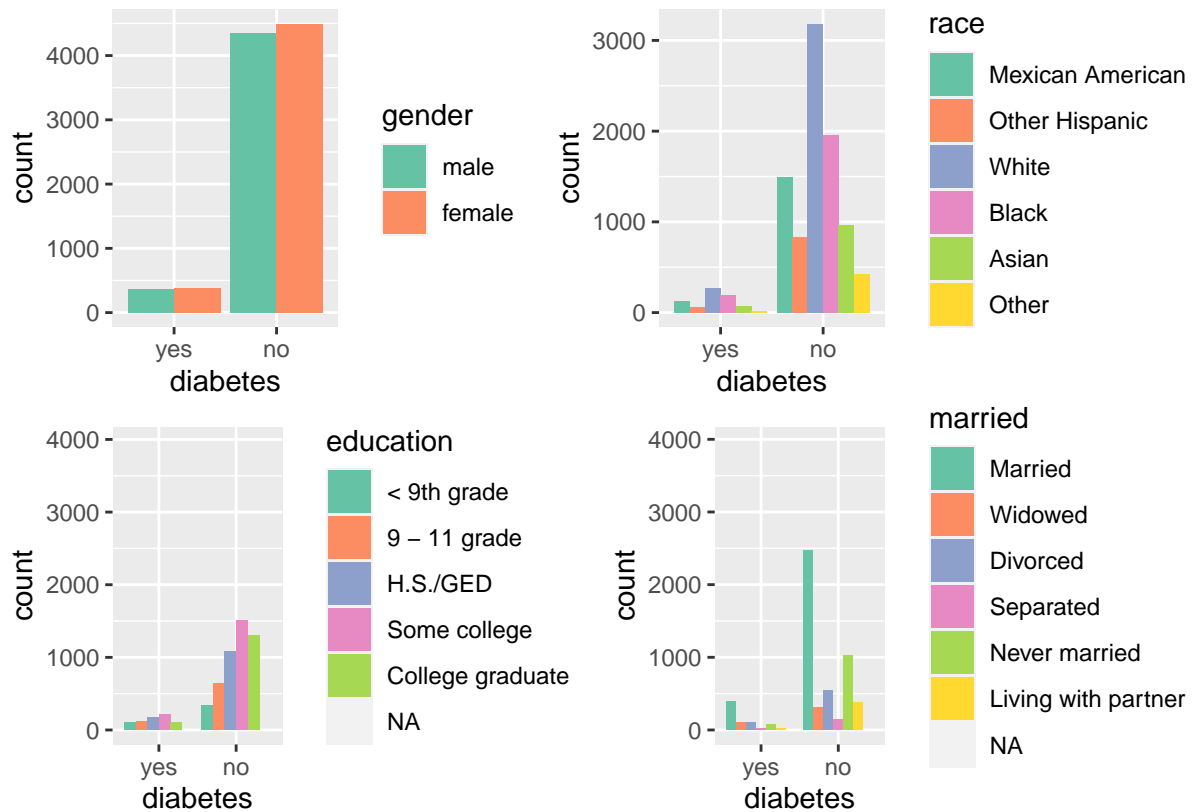
Exploratory data analysis was performed using density plots and bar graphs. Summary statistics were analyzed for all variables, to get an overview of the data and check for extent of missing values. Density plots were used to check for relationships between diabetes and other numeric variables. Categorical variables were visualized separately, using bar graphs instead.

Density plots



Density plots of several numeric covariates showed differences in distributions. Plots of 2 hour glucose tolerance where the participants after the initial fasting were asked to drink a calibrated dose of glucose solution and tested after 2 hours, systolic blood pressure, waist circumference measurement, age, and body mass index showed significant differences between the two groups of participants, those with diabetes and those without. Most significant difference seemed to be among different age groups, with density curve of responses with no diabetes showing right-skewness, and those with diabetes skewed to the left along with a strong shift towards higher age.

Bar plots



Bar plots were analyzed for categorical variables gender, race, education level, and marital status. Presence of diabetes did not seem to be gender-dependent, and with slight differences based on education level. Proportion of positive diabetes cases did seem to vary among different races, and based on marital status. There seemed to be significantly higher proportion of non-diabetics among individuals who were never-married or divorced.

Missing data

Certain variables with high proportion of missing data were retained. The dataset had close to 60% of missing values for the 4 variables which were considered significant for the given response. Assuming that the data was missing at random, and that single imputation might lead to bias and might not preserve relationships between variables; for those reasons imputation was not considered and the missing values were removed. The final sample consisted of 1,564 participants.

Models

Linear models

Non-linear models

Ensemble methods

Support vector machines

Model comparison

Final model

Model prediction performance

Limitations

Conclusion