

# Predicting Diabetes Risk using the 2015 BRFSS Survey



DSC 510 - Health Data Science  
Alex Teboul

# Diabetes Health Indicators Dataset

253,680 survey responses from cleaned BRFSS 2015 + balanced dataset



Alex Teboul • updated 4 months ago

[Data](#) [Code \(17\)](#) [Discussion \(1\)](#) [Activity](#) [Metadata](#) [Settings](#)

Download (52 MB)

New Notebook

Usability 10.0

License CC0: Public Domain

Tags health, beginner, classification, diabetes, public health

## Links

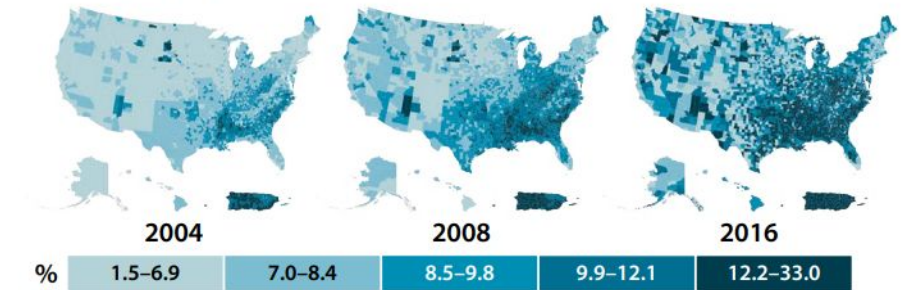
1. [Diabetes Health Indicators Dataset](#) - Kaggle Open Source
2. [Data Cleaning Notebook](#) - Kaggle Open Source
3. [Full Code Notebook with Model Building](#)
4. [Paper](#)

# Introduction - Diabetes

- **34.2 million** Americans have diabetes<sup>2</sup>
  - **88 million** American adults have prediabetes<sup>2</sup>
  - **7th leading cause of death** in the United States<sup>2</sup>
  - **\$400+ Billion** annually<sup>5</sup>
- **1 in 5 are unaware** they are diabetic<sup>2</sup>



Figure 3. Age-adjusted, county-level prevalence of diagnosed diabetes among adults aged 20 years or older, United States, 2004, 2008, and 2016



Note: Data were unavailable for some US territories.

Data sources: US Diabetes Surveillance System; Behavioral Risk Factor Surveillance System.

# Literature Review - Machine Learning in Diabetes

---

## Diabetes Machine Learning Systematic Review (Kavakiotis et al., 2017)

- Applications in diagnosis
- Clinical datasets
- Naive bayes, logistic regression, support vector machines, random forests, and neural networks.

## Diabetes Healthcare Systems Hadoop ML (Yuvaraj & SriPreethaa, 2019)






- Robust Clinical Dataset
- 94% Accuracy (RF)
- 94% PPV (RF)
- 88% Sensitivity (RF)
- 91% F-measure (RF)
- Not applicable as screening tool.

## BRFSS 2014 Type II Diabetes Prediction, CDC (Xie et al., 2019)

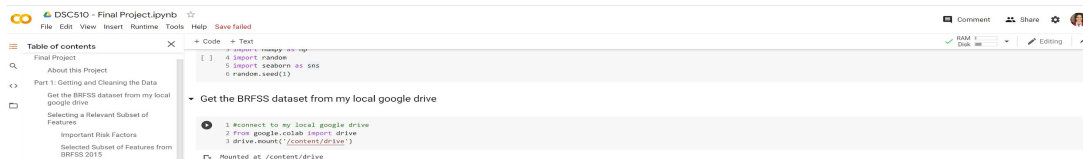
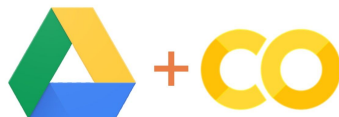
- 82.4% Accuracy (NN)
- 90.2% Specificity (NN)
- 37.8% Sensitivity (NN)
- 51.6% Sensitivity (DT)
- Concluded decision tree model could provide initial population screening.

# Methods - Research Questions

1. **Predicting Diabetes:** To what extent can a subset of survey questions from the BRFSS be used to effectively predict type II diabetes risk?

<b>High Blood Pressure</b> Do you have high blood pressure? 	<b>High Cholesterol</b> Do you have high cholesterol? 	<b>General Health</b> Rate your health: <input type="radio"/> Excellent <input type="radio"/> Very Good <input type="radio"/> Good <input type="radio"/> Fair <input type="radio"/> Poor 
<b>Age</b> What is your age? 	<b>Body Mass Index (BMI)</b> What is your BMI? 	

2. **Research Tool:** Could this serve as a screening tool and can we produce an open source Google Colab notebook to allow researchers or students to clean BRFSS datasets and run machine learning models on them?



# Methods - Dataset



A health-related telephone survey that is collected annually by the CDC.



- **Diabetes** as Dependent Variable
- **21 Other Variables** selected from survey
- **Data Cleaning** in Colab Notebook

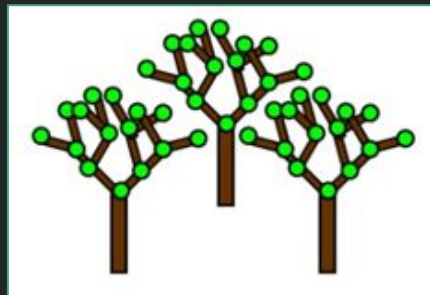


Table 2: Cleaned Datasets

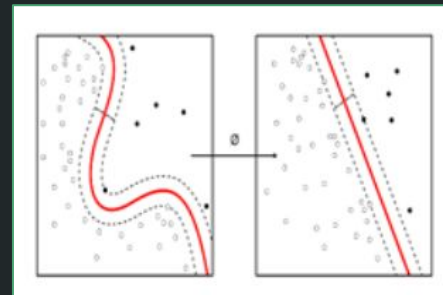
Dataset	Participants with Diabetes	Participants without Diabetes	Total
Binary Unbalanced	35,346	218,334	253,680
Binary Balanced	35,346	35,346	70,692

# Models Tested

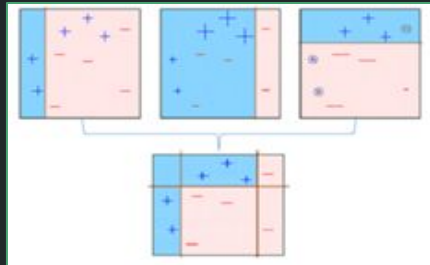
## Random Forest



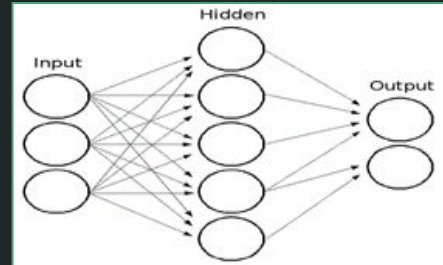
## Gradient Boosting



## AdaBoost



## Neural Network



# Model Testing

- 5-fold cross validation
- Feature selection tested with Random Forest Entropy and Gradient Boosting
- Unbalance & Balanced Dataset
- Performance Metrics
  - Accuracy
  - Sensitivity
  - Specificity
  - Positive Predictive Value (PPV)
  - Negative Predictive Value (NPV)

Table 3

MODEL FEATURES		
#	Renamed Features	Categories
*1	Diabetes	Response Variable
2	<u>HighBP</u>	High Blood Pressure
3	<u>HighChol</u>	High Cholesterol
4	<u>CholCheck</u>	High Cholesterol
5	BMI	BMI
6	Smoker	Smoking History
7	Stroke	Chronic Health Conditions
8	<u>HeartDiseaseorAttack</u>	Chronic Health Conditions
9	<u>PhysActivity</u>	Physical Activity
10	Fruits	Diet
11	Veggies	Diet
12	<u>HvyAlcoholConsump</u>	Alcohol Consumption
13	<u>AnyHealthcare</u>	Health Care Access
14	<u>NoDocbcCost</u>	Health Care Access
15	<u>GenHlth</u>	General Health & Wellbeing
16	<u>MentHlth</u>	General Health & Wellbeing
17	<u>PhysHlth</u>	General Health & Wellbeing
18	<u>DiffWalk</u>	General Health & Wellbeing
19	Sex	Demographics
20	Age	Demographics
21	Education	Demographics
22	Income	Demographics



# Results - Predicting Diabetes

Table 2: Binary Unbalanced Dataset - Model Results					
Model	Accuracy	Sensitivity	Specificity	PPV	NPV
Random Forest	85.0%	18.5%	95.8%	41.8%	87.9%
Gradient Boosting	87.0%	14.0%	98.3%	57.2%	87.6%
AdaBoost	87.0%	15.0%	98.2%	56.1%	87.6%
Neural Network	87.0%	12.3%	98.5%	57.6%	87.4%

- High accuracy and specificity but low sensitivity → Bad
- Need to optimize Sensitivity and PPV if applicable as screening tool.



# Results - Predicting Diabetes






Table 3: Binary Balanced Dataset - Model Results					
Model	Accuracy	Sensitivity	Specificity	PPV	NPV
Random Forest	71.0%	74.5%	68.3%	70.2%	72.8%
Gradient Boosting	74.0%	78.9%	70.0%	72.4%	76.9%
AdaBoost	74.0%	76.7%	71.9%	73.2%	75.5%
Neural Network	74.0%	79.6%	69.3%	72.2%	77.2%

- Improved sensitivity, but lower accuracy and specificity → Better
- CDC BRFSS Diabetes model 51.6% sensitivity & 82.4% accuracy (Xie et al., 2019)

# Discussion - Predicting Diabetes

---

- To what extent can a subset of survey questions from the BRFSS be used to effectively predict type II diabetes risk?
- Performance not high enough to be used in place of medical diagnosis.
- Performance comparable to models with BRFSS survey data (Xie et al. 2019)
- High Blood Pressure, High Cholesterol, BMI, Age, and General Health.

<b>High Blood Pressure</b> Do you have high blood pressure? 	<b>High Cholesterol</b> Do you have high cholesterol? 	<b>General Health</b> Rate your health: <input type="radio"/> Excellent <input type="radio"/> Very Good <input type="radio"/> Good <input type="radio"/> Fair <input type="radio"/> Poor 
<b>Age</b> What is your age? 	<b>Body Mass Index (BMI)</b> What is your BMI? 	

# Discussion - Research Tool

- Could this serve as a screening tool and can we produce an open source Google Colab notebook to allow researchers or students to clean BRFSS datasets and run machine learning models on them?

The screenshot shows a Google Colab notebook interface. The title bar reads 'DSC510 - Final Project.ipynb'. The left sidebar contains a 'Table of contents' with sections: 'Final Project', 'About this Project', 'Part 1: Getting and Cleaning the Data' (with sub-items like 'Get the BRFSS dataset from my local google drive', 'Selecting a Relevant Subset of Features', 'Important Risk Factors', 'Selected Subset of Features from BRFSS 2015', 'Get Subset of Features'), 'Cleaning the Data' (with sub-items like 'Missing Values', 'Modifying Values', 'Make Feature Names More Readable', 'Save Finalized Dataset to CSV', 'Create Binary Dataset for diabetes vs. no diabetes'), 'Part 2: Data Exploration', and 'Part 3: Model Building'. The main content area has a title 'DSC 510: Health Data Science' and a subtitle 'Building Predictive Models for Diabetes Using the 2015 BRFSS'. It lists the author as Alex Teboul and the professor as Stephanie Besser. The data source is cited as <https://www.kaggle.com/cdc/behavioral-risk-factor-surveillance-system#2015.csv>. Under 'About this Project', the objective is stated: 'The goal of this project is to build predictive models for diabetes using the 2015 BRFSS dataset. This project was completed for DSC 510: Health Data Science at DePaul University. In this Google Colab notebook, we go through the process of getting the 2015 BRFSS dataset, selecting features for exploration, and building predictive models based on the risk factors identified in past diabetes research.' Research questions are listed: 1. Can survey questions from the BRFSS provide accurate predictions of whether an individual has diabetes? 2. What risk factors are most predictive of diabetes risk? 3. Can we use a subset of the risk factors to accurately predict whether an individual has diabetes? 4. Can we create a short form of questions from the BRFSS using feature selection to accurately predict if someone might have diabetes or is at high risk of diabetes? The bottom of the notebook shows a status bar with a message: 'The Data was updated remotely or in another tab. Show diff'.

# Limitations & Future Work

---

- Prediabetics included in non-diabetic group because small sample size and didn't match CDC prevalence.
- Survey data is self-reported and subject to recall bias.
- Model is predicting if patient has been told by a doctor that they have diabetes - 1 in 5 are undiagnosed (Centers for Disease Control and Prevention, 2020).
- Future work could involve adding in race variables or improving modeling notebook.

# Conclusion

---



- Need for effective screening tools
- BRFSS survey predictive models comparable to clinical data models for Type II Diabetes risk.
- Kaggle & Open-source

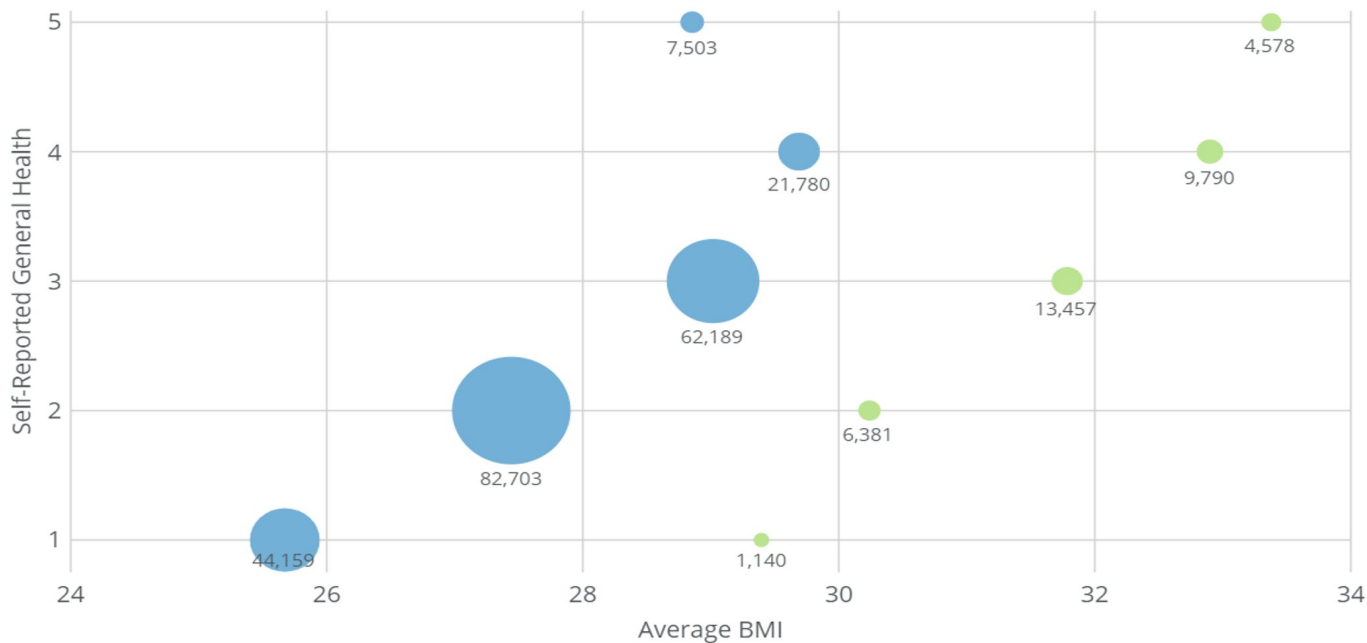
# References

---

1. American Diabetes Association. (2019). 2. Classification and diagnosis of diabetes: standards of medical care in diabetes—2019. *Diabetes care*, 42(Supplement 1), S13-S28.
2. Centers for Disease Control and Prevention. National Diabetes Statistics Report, 2020. Atlanta, GA: Centers for Disease Control and Prevention, U.S. Dept of Health and Human Services; 2020.  
<https://www.cdc.gov/diabetes/pdfs/data/statistics/national-diabetes-statistics-report.pdf>
3. Hippisley-Cox, J., & Coupland, C. (2017). Development and validation of QDiabetes-2018 risk prediction algorithm to estimate future risk of type 2 diabetes: cohort study. *bmj*, 359, j5019.
4. Kavakiotis, I., Tsave, O., Salifoglou, A., Maglaveras, N., Vlahavas, I., & Chouvarda, I. (2017). Machine learning and data mining methods in diabetes research. *Computational and structural biotechnology journal*, 15, 104-116.
5. O'Connell, J. M., & Manson, S. M. (2019). Understanding the economic costs of diabetes and prediabetes and what we may learn about reducing the health and economic burden of these conditions. *Diabetes care*, 42(9), 1609-1611..
6. Xie, Z., Nikolayeva, O., Luo, J., & Li, D. (2019). Peer Reviewed: Building Risk Prediction Models for Type 2 Diabetes Using Machine Learning Techniques. *Preventing chronic disease*, 16.
7. Yuvaraj, N., & SriPreethaa, K. R. (2019). Diabetes prediction in healthcare systems using machine learning algorithms on Hadoop cluster. *Cluster Computing*, 22(1), 1-9.
8. Zou, Q., Qu, K., Luo, Y., Yin, D., Ju, Y., & Tang, H. (2018). Predicting diabetes mellitus with machine learning techniques. *Frontiers in genetics*, 9, 515.

# Diabetes Risk: Self-Reported General Health vs. BMI

**253,680** Total Survey Participants



Diabetes (0 = No; 1 = Yes)

**Self-Reported General Health:** 1=Excellent; 2=Very Good; 3=Good; 4=Fair; 5=Poor

● 0

● 1







# Appendix

Table 3: Final Model Parameters

Model	Parameters
Random Forest	<b>Parameters:</b> n_estimators = 200, max_depth = None, min_samples_split = 3, criterion = entropy, Train-Test Split: 70%-30%, cross validation = 5-fold <b>Feature Selection Parameters:</b> n_estimators = 200, max_depth = None, min_samples_split = 3, criterion = entropy <b>Selected Features:</b> HighBP, BMI, GenHlth, MentHlth, PhysHlth, Age, Education, Income
Gradient Boosting	<b>Parameters:</b> n_estimators = 200, loss = deviance, learning_rate = 0.1, max_depth = 3, min_samples_split = 3 Train-Test Split = 70%-30%, cross validation = 5-fold <b>Feature Selection Parameters:</b> Wrapper selection method using Gradient Boosting Classifier, n_estimators=200, loss=deviance, learning_rate=0.1, max_depth=3, min_samples_split=3 <b>Selected Features:</b> HighBP, HighChol, BMI, GenHlth, Age
AdaBoost	<b>Parameters:</b> n_estimators = 200, base_estimator = None, learning_rate = 0.1 Train-Test Split = 70%-30%, cross validation = 5-fold <b>Feature Selection Parameters:</b> Wrapper selection method using Gradient Boosting Classifier, n_estimators=200, loss=deviance, learning_rate=0.1, max_depth=3, min_samples_split=3 <b>Selected Features:</b> HighBP, HighChol, BMI, GenHlth, Age
Neural Network	<b>Parameters:</b> activation = logistic, solver = adam, alpha = 0.0001, max_iter = 1000, hidden_layer_sizes = (10,) Train-Test Split = 70%-30%, cross validation = 5-fold <b>Feature Selection Parameters:</b> Wrapper selection method using Gradient Boosting Classifier, n_estimators=200, loss=deviance, learning_rate=0.1, max_depth=3, min_samples_split=3 <b>Selected Features:</b> HighBP, HighChol, BMI, GenHlth, Age

- **The original BRFSS 2015 .csv can be downloaded here:**  
<https://www.kaggle.com/cdc/behavioral-risk-factor-surveillance-system#2015.csv>
- **The BRFSS 2015 Codebook is available here:**  
[https://www.cdc.gov/brfss/annual\\_data/2015/pdf/codebook15\\_1lcp.pdf](https://www.cdc.gov/brfss/annual_data/2015/pdf/codebook15_1lcp.pdf)
- **The open-source Google Colab notebook with Python code is available here:**  
[https://colab.research.google.com/drive/1HUYgcxhmgzv5zELcsnM0da\\_gwo1Zuiuo?usp=sharing](https://colab.research.google.com/drive/1HUYgcxhmgzv5zELcsnM0da_gwo1Zuiuo?usp=sharing)
- **The datasets used for model building, created in the Google Colab notebook are available in this Google Drive folder:**  
[https://drive.google.com/drive/folders/1yoEQqCn75TxKknWGkWOvDrVn2O\\_qYitl?usp=sharing](https://drive.google.com/drive/folders/1yoEQqCn75TxKknWGkWOvDrVn2O_qYitl?usp=sharing)
- **Video Presentation Link**

# Appendix

Variable Specifics		
Index	Variable	BRFSS Question
 0	<b>Diabetes_Binary</b> <i>BRFSS: <u>DIABETE3</u></i>	(Ever told) you have diabetes (If "Yes" and respondent is female, ask "Was this only when you were pregnant?").
 1	<b>HighBP</b> <i>BRFSS: <u>RFHYPE5</u></i>	Adults who have been told they have high blood pressure by a doctor, nurse, or other health professional.
 2	<b><u>HighChol</u></b> <i>BRFSS: <u>TOLDHI2</u></i>	Have you EVER been told by a doctor, nurse or other health professional that your blood cholesterol is high?
3	<b>CholCheck</b> <i>BRFSS: <u>CHOLCHK</u></i>	Cholesterol check within past five years
 4	<b>BMI</b> <i>BRFSS: <u>BMIS</u></i>	Body Mass Index (BMI)



= Dependent Variable



= Selected in best Model

# Appendix

---

5	<b>Smoker</b> <i>BRFSS: SMOKE100</i>	Have you smoked at least 100 cigarettes in your entire life? [Note: 5 packs = 100 cigarettes]
6	<b>Stroke</b> <i>BRFSS: CVDSTRK3</i>	(Ever told) you had a stroke.
7	<b>HeartDiseaseorAttack</b> <i>BRFSS: _MICH</i>	Respondents that have ever reported having coronary heart disease (CHD) or myocardial infarction (MI)
8	<b>PhysActivity</b> <i>BRFSS: _TOTINDA</i>	Adults who reported doing physical activity or exercise during the past 30 days other than their regular job
9	<b>Fruits</b>	Consume Fruit 1 or more times per day

# Appendix

10	<b>Veggies</b> <i>BRFSS: _VEGLT1</i>	Consume Vegetables 1 or more times per day
11	<b>HvyAlcoholConsump</b> <i>BRFSS: _RFDRHV5</i>	Heavy drinkers (adult men having more than 14 drinks per week and adult women having more than 7 drinks per week)
12	<b>AnyHealthcare</b> <i>BRFSS: HLTHPLN1</i>	Do you have any kind of health care coverage, including health insurance, prepaid plans such as HMOs, or government plans such as Medicare, or Indian Health Service?
13	<b>NoDocbcCost</b> <i>BRFSS: MEDCOST</i>	Was there a time in the past 12 months when you needed to see a doctor but could not because of cost?
14	<b>GenHlth</b> <i>BRFSS: GENHLTH</i>	Would you say that in general your health is:
15	<b>MentHlth</b> <i>BRFSS: <u>MENTHLTH</u></i>	Now thinking about your mental health, which includes stress, depression, and problems with emotions, for how many days during the past 30 days was your mental health not good?

# Appendix

16	<b>PhysHlth</b> <i>BRFSS: <u>PHYSHLTH</u></i>	Now thinking about your physical health, which includes physical illness and injury, for how many days during the past 30 days was your physical health not good?
17	<b>DiffWalk</b> <i>BRFSS: DIFFWALK</i>	Do you have serious difficulty walking or climbing stairs?
18	<b>Sex</b> <i>BRFSS: SEX</i>	Indicate sex of respondent.
19	<b>Age</b> <i>BRFSS: _AGEG5YR</i>	Fourteen-level age category
20	<b>Education</b> <i>BRFSS: EDUCA</i>	What is the highest grade or year of school you completed?
21	<b>Income</b> <i>BRFSS: INCOME2</i>	Is your annual household income from all sources: (If respondent refuses at any income level, code "Refused.")