

Diabetes Prediction

Alexander Wen, Raymond Wong, Michael Eirikson

Table of contents

Summary	1
Introduction	1
Methods	2
EDA	2
Modeling	8
Discussion	8
References	8

Summary

In this project we attempt to build a model to predict diabetes disease. We compared a decision tree model and naive bayes model and found the decision tree is stronger in this context. We used f2-score as our scoring function because detecting diabetes is the priority: a false negative could be much worse than a false positive.

In the test dataset: the decision tree model correctly detected 8283 of 10604 positive cases (recall rate is about 78%). This result does come at a fairly significant cost in terms of false positives (precision rate is about 30%) with 19650 false positives. Depending on the actual cost of false positive this may need significant improvement to be a viable screening model.

Introduction

In Canada and the USA approximately 10% of people are living with diabetes. In Canada in 2023 approximately 3.7 million people were living with diabetes and in the USA in 2021 approximately 38.4 million people were living with diabetes. (“Snapshot of Diabetes in Canada, 2023” (2023)) In the USA it is the 8th leading cause of death. (Rios et al. (2017)) Globally an estimated 44% of people living with diabetes are undiagnosed. (Stafford et al. (2025))

In this project we try to predicted diabetes disease based on common health factors. A reliable model could help to prescreen people and recommend following up with a physician for people who are at risk. Given the large number of people living with undiagnosed diabetes this could potentially have a significant positive impact of world health.

The analysis uses the American CDC Behavioural Risk Factor Surveillance System (BRFSS) 2015 Diabetes Health Indicators dataset (UCI ID 891), containing 253,680 survey responses with 21 health-related features and a binary diabetes outcome (0 = no diabetes/pre-diabetes, 1 = diabetes). (Dane and Teboul (2021))

No missing values were present and all features were already encoded numerically. The target classes is imbalanced (86% non-diabetic, 14% diabetic).

Methods

This analysis was performed in Python 3.11.6 (*Python 3.11.6 Documentation* 2021-2025). Additionally, here is a list of the Python packages used within the analysis with brief explanation:

Package	Version	Use case	Reference
<code>numpy</code>	1.26.4	General analysis use	<i>NumPy Documentation</i> (2008-2022)
<code>pandas</code>	2.1.2	Data management/processing	team (2020), McKinney (2010)
<code>pandera</code>	0.27.0	Data validation	Bantilan (2020)
<code>altair</code>	5.1.2	Generating plots	VanderPlas et al. (2018), Satyanarayanan et al. (2017)
<code>scikit-learn</code>	1.3.2	Model creation and evaluation	Pedregosa et al. (2011)
<code>ucimlrepo</code>	0.0.7	Data extraction	Kelly, Longjohn, and Nottingham (2021)
<code>deepchecks</code>	0.18.1	Data validation	Chorev et al. (2022)
<code>click</code>	8.3.1	Script tool	Pallets (2020)
<code>quarto</code>	1.8.26	Report creation	Allaire et al. (2025)
<code>tabulate</code>	0.9.0	Table formatting	Astanin (2025)

Modeling Approach

Results

EDA

Data Summary

First, here is a sample of the training data showing the first few entries and last few entries in the dataset in Table 2 and Table 3.

Table 2: First few rows of the training data.

	Un-	named:	High-	CB	B	Cho	EM	Scho	Sto	lack	it	Fruit	ges	sum	pcare	Cost	Gen	Mth	Phy	Wt	Ex	Age	com	be
0	0	0	1	1	23	0	0	0	1	1	1	0	1	0	1	0	0	0	0	0	10	6	8	0
1	1	0	0	1	25	0	0	0	1	1	1	0	1	0	3	0	30	0	0	12	6	7	0	
2	2	1	1	1	28	0	0	0	1	0	1	0	1	0	2	15	2	0	1	6	5	6	0	
3	3	0	0	1	25	0	0	0	0	1	1	0	1	0	2	0	0	0	0	8	5	7	0	
4	4	1	0	1	30	1	0	0	1	1	1	0	1	1	4	30	15	0	0	8	4	4	0	

Table 3: Last few rows of the training data.

	Un-	named:	High-	CB	B	Cho	EM	Scho	Sto	lack	it	Fruit	ges	sum	pcare	Cost	Gen	Mth	Phy	Wt	Ex	Age	com	be
0	177571	1	1	29	0	0	0	1	1	1	0	1	0	2	0	2	0	0	4	6	5	0		
1	177572	0	1	22	0	0	0	0	1	1	0	1	0	3	0	0	1	0	13	4	1	0		
2	177573	1	1	25	1	0	0	1	1	1	0	1	0	2	0	0	0	1	6	5	7	0		
3	177574	1	1	24	1	0	0	1	1	1	0	1	0	3	2	1	0	0	4	4	8	0		
4	177575	1	1	31	1	0	0	0	1	1	0	1	0	2	0	2	0	0	8	4	4	0		

All features of the dataset are numeric, and further EDA shows there are no null values in the dataset.

Table 4: Description of the training data.

	Heart-		HvyAl-			
	Dis-	HvyAl-	co-		Ed-	
	ease-Phys-	co-	hol-	Any-	u-	di-
Un-	o-	Ac-	hol-	Any-	u-	di-
named:	High-	rAt-	tiv-	Veg-Con-Heal	NoDocbc-	Dif-
0	High	CH	Cho	KM	Stro	Fruit
1	count	775	776	51767	51767	51767
2	mean	0.4202472079626	7385440346070940	356313505120205402954028325141798332501186018767830074699504139332		
3	std	0.49097900189559064907897629187520308043905530184527782058408632037003934558803479316294				
4	min	0	0	0	12	0
5	25%	0	1	24	0	0
6	50%	0	1	27	0	0
7	75%	1	1	31	1	0
8	max	1	1	98	1	1
9						
10						
11						
12						
13						
14						
15						
16						
17						
18						
19						
20						
21						
22						
23						
24						
25						
26						
27						
28						
29						
30						
31						
32						
33						
34						
35						
36						
37						
38						
39						
40						
41						
42						
43						
44						
45						
46						
47						
48						
49						
50						
51						
52						
53						
54						
55						
56						
57						
58						
59						
60						
61						
62						
63						
64						
65						
66						
67						
68						
69						
70						
71						
72						
73						
74						
75						
76						
77						
78						
79						
80						
81						
82						
83						
84						
85						
86						
87						
88						
89						
90						
91						
92						
93						
94						
95						
96						
97						
98						
99						
100						
101						
102						
103						
104						
105						
106						
107						
108						
109						
110						
111						
112						
113						
114						
115						
116						
117						
118						
119						
120						
121						
122						
123						
124						
125						
126						
127						
128						
129						
130						
131						
132						
133						
134						
135						
136						
137						
138						
139						
140						
141						
142						
143						
144						
145						
146						
147						
148						
149						
150						
151						
152						
153						
154						
155						
156						
157						
158						
159						
160						
161						
162						
163						
164						
165						
166						
167						
168						
169						
170						
171						
172						
173						
174						
175						
176						
177						
178						
179						
180						
181						
182						
183						
184						
185						
186						
187						
188						
189						
190						
191						
192						
193						
194						
195						
196						
197						
198						
199						
200						
201						
202						
203						
204						
205						
206						
207						
208						
209						
210						
211						
212						
213						
214						
215						
216						
217						
218						
219						
220						
221						
222						
223						
224						
225						
226						
227						
228						
229						
230						
231						
232						
233						
234						
235						
236						
237						
238						
239						
240						
241						
242						
243						
244						
245						
246						
247						
248						
249						
250						
251						
252						
253						
254						
255						
256						
257						
258						
259						
260						
261						
262						
263						
264						
265						
266						
267						
268						
269						
270						
271						
272						
273						
274						
275						
276						
277						
278						
279						
280						
281						
282						
283						
284						
285						
286						
287						
288						
289						
290						
291						
292						
293						
294						
295						
296						
297						
298						
299						
300						
301						
302						
303						
304						
305						
306						
307						
308						
309						
310						
311						
312						
313						
314						
315						
316						
317						
318						
319						
320						
321					</	

Table 4 displays a numerical distribution of the dataset. Most features are binary in nature, with the exceptions of To better visualize we then plotted frequency counts of the target labels:

Visualizations

Figure 1

Figure 2

Figure 3

Figure 4

Figure 5

Count of Diabetes vs Non-Diabetes Records in Dataset

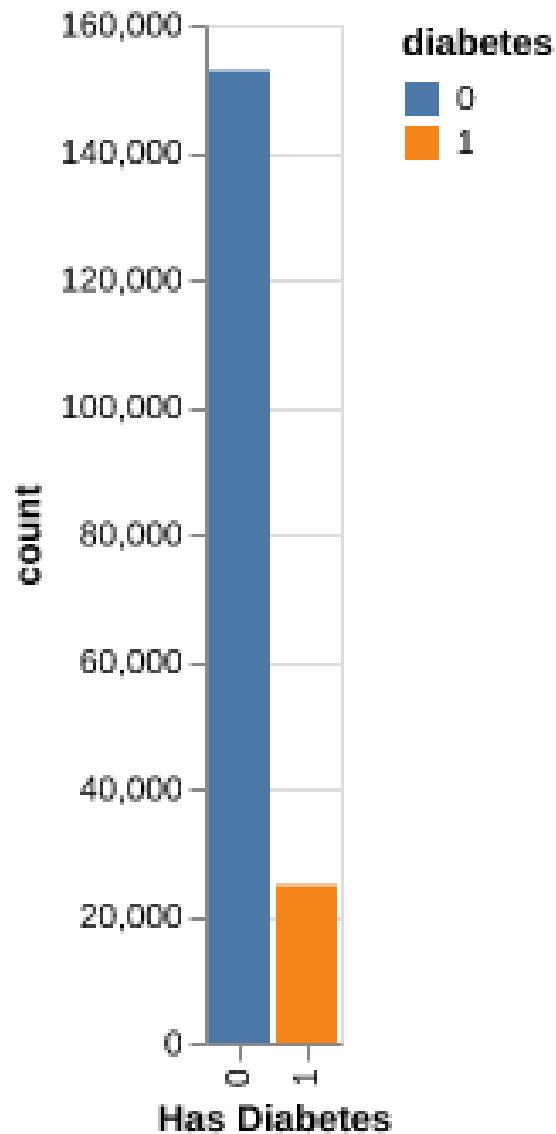


Figure 1: Frequency bar graph of the labels.

Bar Plots of Binary Features

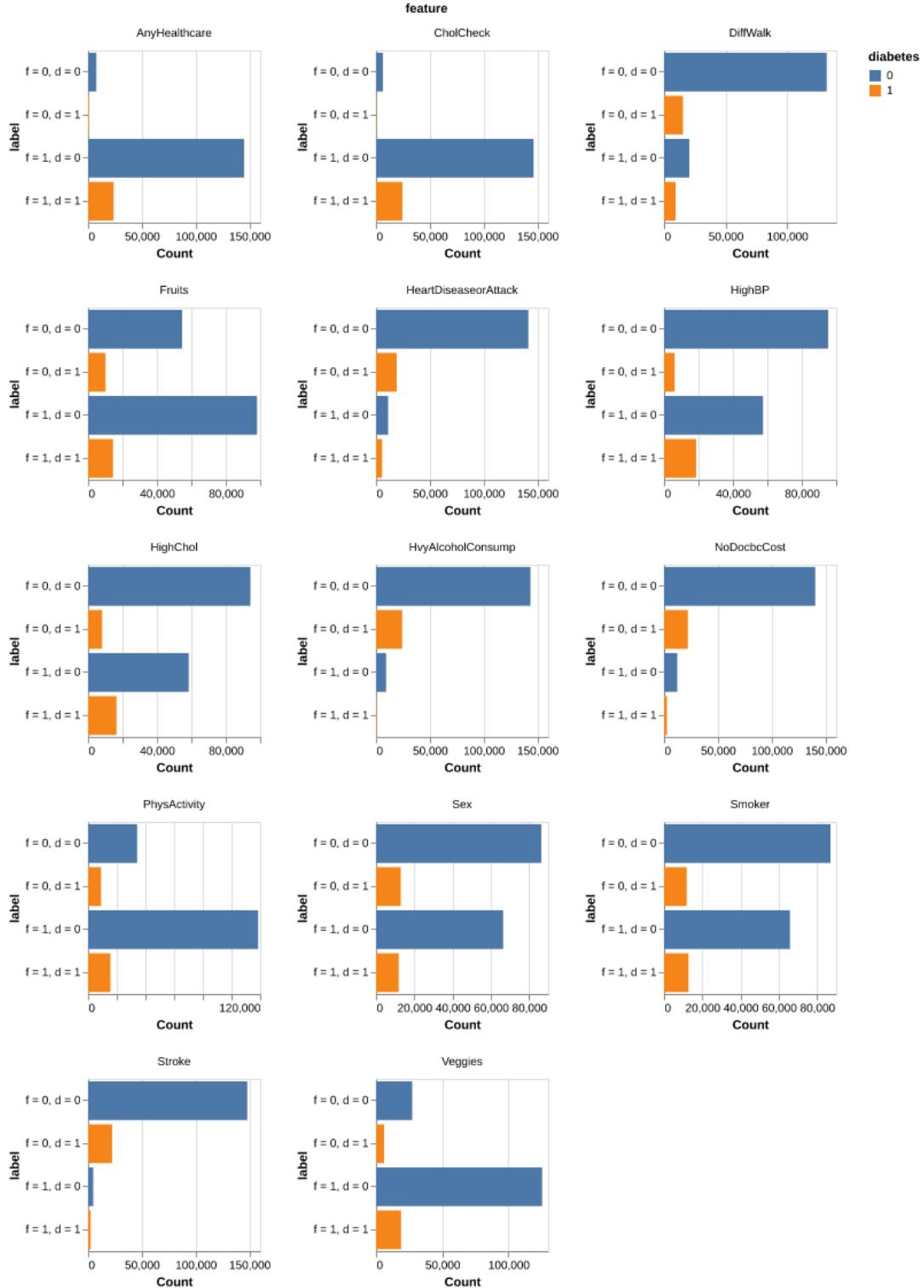


Figure 2: Frequency bar graphs of the binary features.

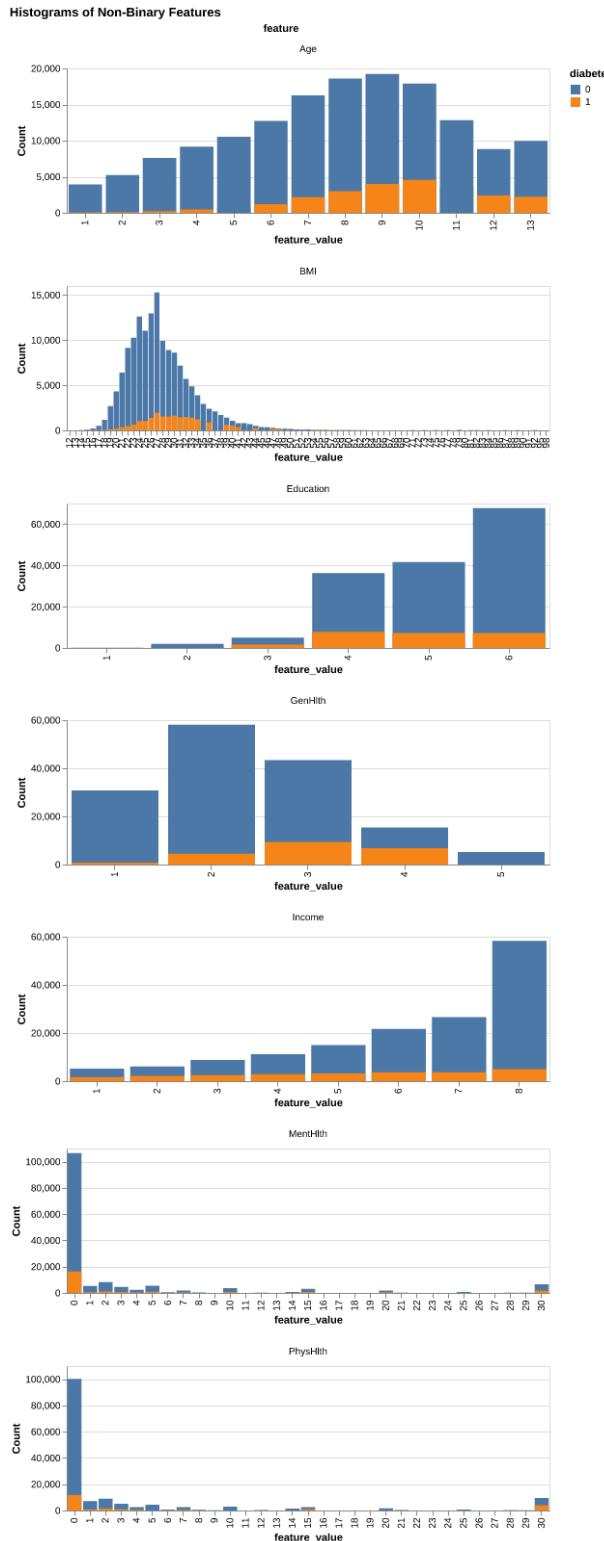


Figure 3: Histograms of the non-binary numeric features.

Boxplots for Non-Binary Features (Sample size n = 1000)

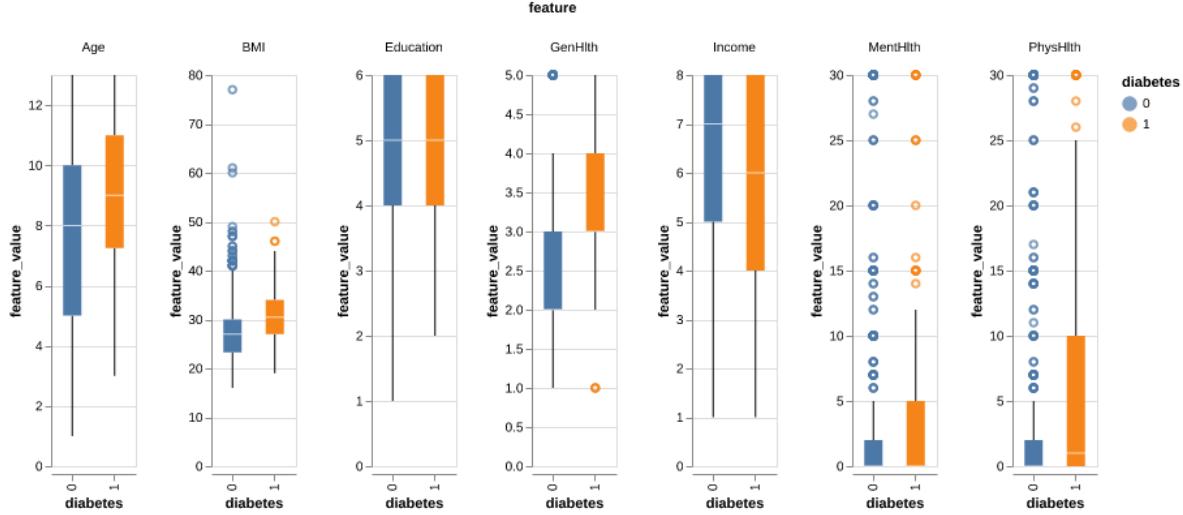


Figure 4: Boxplots of the non-binary numeric features.

Modeling

Classification Analysis

Result Visualizations

Discussion

References

- Allaire, J. J., Charles Teague, Carlos Scheidegger, Yihui Xie, Christophe Dervieux, and Gordon Woodhull. 2025. “Quarto.” <https://doi.org/10.5281/zenodo.5960048>.
- Astanin, Sergio. 2025. “Tabulate: Pretty-Print Tabular Data in Python.” <https://github.com/astanin/python-tabulate>.
- Bantilan, Niels. 2020. “Pandera: Statistical Data Validation of Pandas Dataframes.” In *Proceedings of the 19th Python in Science Conference*, edited by Meghann Agarwal, Chris Calloway, Dillon Niederhut, and David Shupe, 116–24. <https://doi.org/10.25080/Majora-342d178e-010>.
- Chorev, Shir, Philip Tannor, Dan Ben Israel, Noam Bressler, Itay Gabbay, Nir Hutnik, Jonatan Liberman, Matan Perlmutter, Yurii Romanyshyn, and Lior Rokach. 2022. “Deepchecks: A Library for Testing and Validating Machine Learning Models and Data.” <https://arxiv.org/abs/2203.08491>.

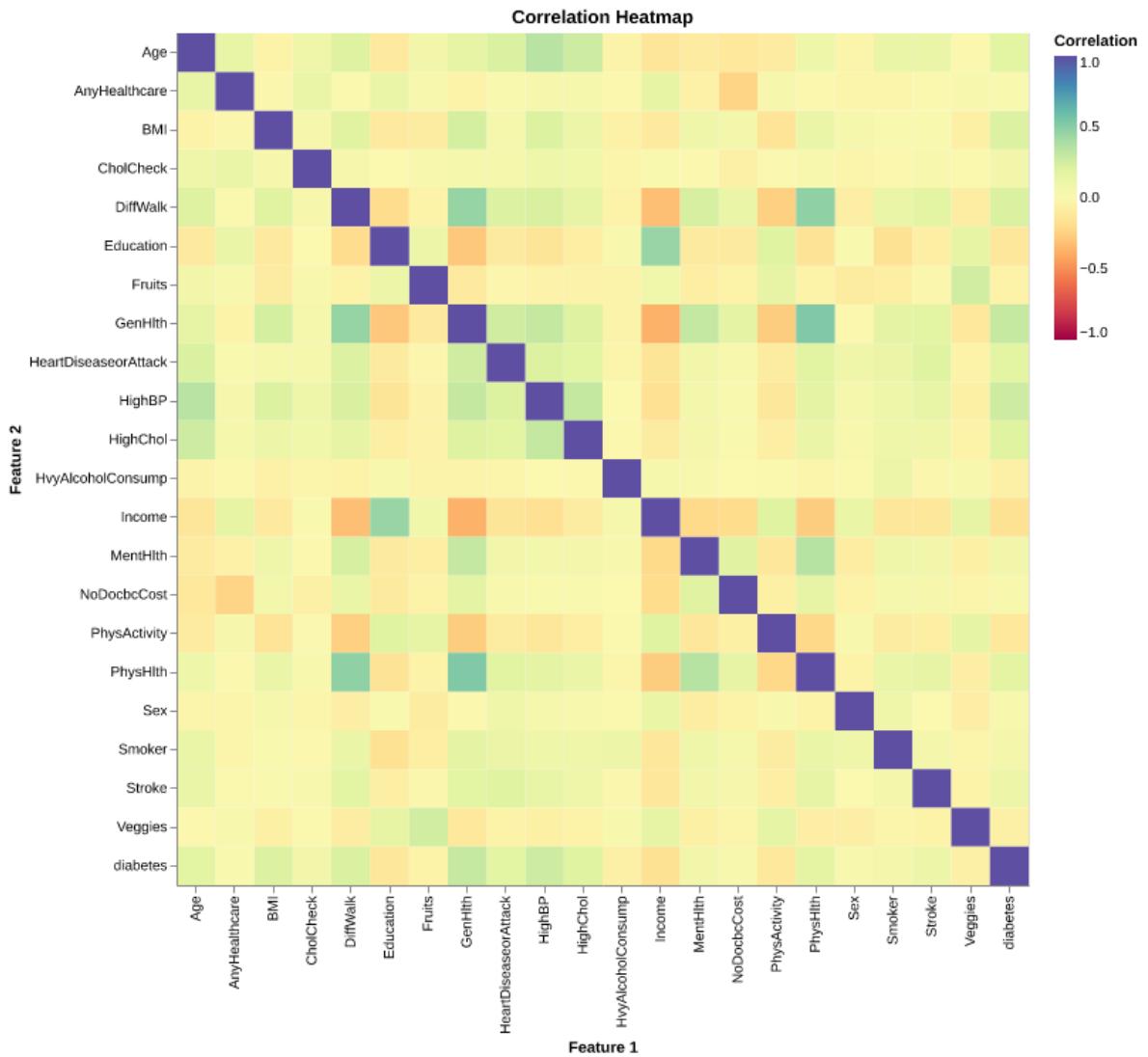


Figure 5: Feature-feature correlation plot of all features.

- Dane, Sohier, and Alex Teboul. 2021. “Diabetes Health Indicators Dataset.” <https://www.kaggle.com/datasets/alextreboul/diabetes-health-indicators-dataset/data>; Center of Disease Control.
- Kelly, Markelle, Rachel Longjohn, and Kolby Nottingham. 2021. “The UCI Machine Learning Repository.” University of California, Irvine, School of Information; Computer Sciences. <http://archive.ics.uci.edu/ml>.
- McKinney, Wes. 2010. “Data Structures for Statistical Computing in Python.” In *Proceedings of the 9th Python in Science Conference*, edited by Stéfan van der Walt and Jarrod Millman, 56–61. <https://doi.org/10.25080/Majora-92bf1922-00a>.
- NumPy Documentation*. 2008-2022. <https://numpy.org/doc/1.26/>.
- Pallets. 2020. *Click*. <https://click.palletsprojects.com/>.
- Pedregosa, F., G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, et al. 2011. “Scikit-Learn: Machine Learning in Python.” *Journal of Machine Learning Research* 12: 2825–30.
- Python 3.11.6 Documentation*. 2021-2025. <https://docs.python.org/release/3.11.6/reference/index.html>.
- Rios, Nilka Burrows, Isreal Hora, Linda S. Geiss, Edward W. Gregg, and Ann Albright. 2017. “Incidence of End-Stage Renal Disease Attributed to Diabetes Among Persons with Diagnosed Diabetes — United States and Puerto Rico.” *Incidence of End-Stage Renal Disease Attributed to Diabetes Among Persons with Diagnosed Diabetes — United States and Puerto Rico*, no. 66(43): 1165–70.
- Satyanarayan, Arvind, Dominik Moritz, Kanit Wongsuphasawat, and Jeffrey Heer. 2017. “Vega-Lite: A Grammar of Interactive Graphics.” *IEEE Transactions on Visualization and Computer Graphics* 23 (1): 341–50.
- “Snapshot of Diabetes in Canada, 2023.” 2023. <https://www.canada.ca/en/public-health/services/publications/diseases-conditions/snapshot-diabetes-canada-2023.html>; Public Health Agency of Canada.
- Stafford, Lauryn K, Anna Gage, Yvonne Yiru Xu, Madeleine Conrad, Ismael Barreras Beltran, and Edward J Boyko. 2025. “Global, Regional, and National Cascades of Diabetes Care, 2000–23: A Systematic Review and Modelling Analysis Using Findings from the Global Burden of Disease Study.” *The Lancet Diabetes & Endocrinology*, no. 13(11): 924–34.
- team, The pandas development. 2020. “Pandas-Dev/Pandas: Pandas.” Zenodo. <https://doi.org/10.5281/zenodo.3509134>.
- VanderPlas, Jacob, Brian Granger, Jeffrey Heer, Dominik Moritz, Kanit Wongsuphasawat, Arvind Satyanarayan, Eitan Lees, Ilia Timofeev, Ben Welsh, and Scott Sievert. 2018. “Altair: Interactive Statistical Visualizations for Python.” *Journal of Open Source Software* 3 (32): 1057. <https://doi.org/10.21105/joss.01057>.