

Data Dictionary & Project Plan // SYDE 780 - Benjamin Luo (20890448)

Data: <https://www.kaggle.com/datasets/kamilpytlak/personal-key-indicators-of-heart-disease>

Full Data (2022): https://www.cdc.gov/brfss/annual_data/annual_2022.html

Data dictionary and preliminary preprocessing

The Kaggle dataset is actually an excerpt of data from the CDC BRFSS annual survey with only features relevant to heart disease (n=44/328). The original data is made available through SAS files which can be read through Python's Pandas library for preprocessing. The Kaggle dataset renames the original feature names without documenting them so I instead use the original data with 445132 data points and 328 features. The methodology used to preprocess is as follows:

1	<pre>df22 = pd.read_sas('~/data/LLCP2022.XPT', format='xport') ✓ 0.0s</pre> <table border="1"><thead><tr><th>_STATE</th><th>FMONTH</th><th>IDATE</th><th>IMONTH</th><th>IDAY</th><th>IYEAR</th><th>DISPCODE</th><th>SEQNO</th><th>_PSU</th><th>CTELENMI</th></tr></thead><tbody><tr><td>0</td><td>1.0</td><td>1.0</td><td>b'02032022'</td><td>b'02'</td><td>b'03'</td><td>b'2022'</td><td>1100.0</td><td>b'2022000001'</td><td>2.022000e+09</td><td>1.0</td></tr><tr><td>1</td><td>1.0</td><td>1.0</td><td>b'02042022'</td><td>b'02'</td><td>b'04'</td><td>b'2022'</td><td>1100.0</td><td>b'2022000002'</td><td>2.022000e+09</td><td>1.0</td></tr><tr><td>2</td><td>1.0</td><td>1.0</td><td>b'02022022'</td><td>b'02'</td><td>b'02'</td><td>b'2022'</td><td>1100.0</td><td>b'2022000003'</td><td>2.022000e+09</td><td>1.0</td></tr><tr><td>3</td><td>1.0</td><td>1.0</td><td>b'02032022'</td><td>b'02'</td><td>b'03'</td><td>b'2022'</td><td>1100.0</td><td>b'2022000004'</td><td>2.022000e+09</td><td>1.0</td></tr><tr><td>4</td><td>1.0</td><td>1.0</td><td>b'02022022'</td><td>b'02'</td><td>b'02'</td><td>b'2022'</td><td>1100.0</td><td>b'2022000005'</td><td>2.022000e+09</td><td>1.0</td></tr><tr><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td></tr><tr><td>445127</td><td>78.0</td><td>11.0</td><td>b'12192022'</td><td>b'12'</td><td>b'19'</td><td>b'2022'</td><td>1100.0</td><td>b'2022001527'</td><td>2.022002e+09</td><td>NaN</td></tr><tr><td>445128</td><td>78.0</td><td>11.0</td><td>b'12212022'</td><td>b'12'</td><td>b'21'</td><td>b'2022'</td><td>1100.0</td><td>b'2022001528'</td><td>2.022002e+09</td><td>NaN</td></tr><tr><td>445129</td><td>78.0</td><td>11.0</td><td>b'11292022'</td><td>b'11'</td><td>b'29'</td><td>b'2022'</td><td>1100.0</td><td>b'2022001529'</td><td>2.022002e+09</td><td>NaN</td></tr><tr><td>445130</td><td>78.0</td><td>11.0</td><td>b'12082022'</td><td>b'12'</td><td>b'08'</td><td>b'2022'</td><td>1100.0</td><td>b'2022001530'</td><td>2.022002e+09</td><td>NaN</td></tr><tr><td>445131</td><td>78.0</td><td>11.0</td><td>b'12142022'</td><td>b'12'</td><td>b'14'</td><td>b'2022'</td><td>1100.0</td><td>b'2022001531'</td><td>2.022002e+09</td><td>NaN</td></tr></tbody></table>	_STATE	FMONTH	IDATE	IMONTH	IDAY	IYEAR	DISPCODE	SEQNO	_PSU	CTELENMI	0	1.0	1.0	b'02032022'	b'02'	b'03'	b'2022'	1100.0	b'2022000001'	2.022000e+09	1.0	1	1.0	1.0	b'02042022'	b'02'	b'04'	b'2022'	1100.0	b'2022000002'	2.022000e+09	1.0	2	1.0	1.0	b'02022022'	b'02'	b'02'	b'2022'	1100.0	b'2022000003'	2.022000e+09	1.0	3	1.0	1.0	b'02032022'	b'02'	b'03'	b'2022'	1100.0	b'2022000004'	2.022000e+09	1.0	4	1.0	1.0	b'02022022'	b'02'	b'02'	b'2022'	1100.0	b'2022000005'	2.022000e+09	1.0	445127	78.0	11.0	b'12192022'	b'12'	b'19'	b'2022'	1100.0	b'2022001527'	2.022002e+09	NaN	445128	78.0	11.0	b'12212022'	b'12'	b'21'	b'2022'	1100.0	b'2022001528'	2.022002e+09	NaN	445129	78.0	11.0	b'11292022'	b'11'	b'29'	b'2022'	1100.0	b'2022001529'	2.022002e+09	NaN	445130	78.0	11.0	b'12082022'	b'12'	b'08'	b'2022'	1100.0	b'2022001530'	2.022002e+09	NaN	445131	78.0	11.0	b'12142022'	b'12'	b'14'	b'2022'	1100.0	b'2022001531'	2.022002e+09	NaN	Download the 2022 dataset (.XPT) from the CDC and read it into a Pandas dataframe																																																
_STATE	FMONTH	IDATE	IMONTH	IDAY	IYEAR	DISPCODE	SEQNO	_PSU	CTELENMI																																																																																																																																																																											
0	1.0	1.0	b'02032022'	b'02'	b'03'	b'2022'	1100.0	b'2022000001'	2.022000e+09	1.0																																																																																																																																																																										
1	1.0	1.0	b'02042022'	b'02'	b'04'	b'2022'	1100.0	b'2022000002'	2.022000e+09	1.0																																																																																																																																																																										
2	1.0	1.0	b'02022022'	b'02'	b'02'	b'2022'	1100.0	b'2022000003'	2.022000e+09	1.0																																																																																																																																																																										
3	1.0	1.0	b'02032022'	b'02'	b'03'	b'2022'	1100.0	b'2022000004'	2.022000e+09	1.0																																																																																																																																																																										
4	1.0	1.0	b'02022022'	b'02'	b'02'	b'2022'	1100.0	b'2022000005'	2.022000e+09	1.0																																																																																																																																																																										
...																																																																																																																																																																											
445127	78.0	11.0	b'12192022'	b'12'	b'19'	b'2022'	1100.0	b'2022001527'	2.022002e+09	NaN																																																																																																																																																																										
445128	78.0	11.0	b'12212022'	b'12'	b'21'	b'2022'	1100.0	b'2022001528'	2.022002e+09	NaN																																																																																																																																																																										
445129	78.0	11.0	b'11292022'	b'11'	b'29'	b'2022'	1100.0	b'2022001529'	2.022002e+09	NaN																																																																																																																																																																										
445130	78.0	11.0	b'12082022'	b'12'	b'08'	b'2022'	1100.0	b'2022001530'	2.022002e+09	NaN																																																																																																																																																																										
445131	78.0	11.0	b'12142022'	b'12'	b'14'	b'2022'	1100.0	b'2022001531'	2.022002e+09	NaN																																																																																																																																																																										
2	<p>Label: Cellular Telephone Section Name: Land Line Introduction Section Number: 0 Question Number: 5 Column: 67 Type of Variable: Num SAS Variable Name: CELPHON1 Question Prologue: Variable only on the land line survey Question: Is this a cell telephone?</p> <table border="1"><thead><tr><th>Value</th><th>Value Label</th><th>Frequency</th><th>Percentage</th><th>Weighted Percentage</th></tr></thead><tbody><tr><td>1</td><td>Yes, it is a cell phone-Terminate Phone Call</td><td>2</td><td>0.00</td><td>0.00</td></tr><tr><td>2</td><td>Not a cell phone-Go to LL.06, LADULT1</td><td>96,050</td><td>100.00</td><td>100.00</td></tr><tr><td>BLANK</td><td>Missing or Cell Phone Interview Notes: QSTVER >= 20</td><td>349,080</td><td>-</td><td>-</td></tr></tbody></table>	Value	Value Label	Frequency	Percentage	Weighted Percentage	1	Yes, it is a cell phone-Terminate Phone Call	2	0.00	0.00	2	Not a cell phone-Go to LL.06, LADULT1	96,050	100.00	100.00	BLANK	Missing or Cell Phone Interview Notes: QSTVER >= 20	349,080	-	-	Download the 2022 codebook (.HTM) that explains each of the 328 variables. Crawl this HTM file in Python using BeautifulSoup to form a dataframe of variables, data types, descriptions, and possible values. Use string processing to remove unnecessary data (ex. 'Notes')																																																																																																																																																														
Value	Value Label	Frequency	Percentage	Weighted Percentage																																																																																																																																																																																
1	Yes, it is a cell phone-Terminate Phone Call	2	0.00	0.00																																																																																																																																																																																
2	Not a cell phone-Go to LL.06, LADULT1	96,050	100.00	100.00																																																																																																																																																																																
BLANK	Missing or Cell Phone Interview Notes: QSTVER >= 20	349,080	-	-																																																																																																																																																																																
3	<p>Independent variables CVDINFR4: (Ever told) you had a heart attack, also called a myocardial infarction? CVDCRHD4: (Ever told) (you had) angina or coronary heart disease? _MICHD: Respondents that have ever reported having coronary heart disease (CHD) or myocardial infarction (MI).</p> <p>Demographic and socioeconomic _STATE: State of residence (geographic location can influence access to healthcare and environmental factors). SEXVAR: Sex of the respondent (cardiovascular risk varies by sex). _AGEGR4: Age group (cardiovascular risk increases with age). _RACEGR4: Race/ethnicity (certain racial/ethnic groups have higher cardiovascular risk). _EDUCAG: Education level (lower education levels are associated with higher cardiovascular risk). _INCOMG1: Income level (lower income is associated with higher cardiovascular risk). MARITAL: Marital status (social support can influence cardiovascular health). EMPLOY1: Employment status (stress and job type can impact cardiovascular health).</p>	Use an LLM to identify and cluster relevant variables (n=53). Manually validate results																																																																																																																																																																																		
4	<pre>df22c = df22c.dropna(axis=1, thresh=int(0.75 * len(df22)) + 1) removed_cols = df22c.columns.difference(df22c.columns) df22c</pre> <table border="1"><thead><tr><th>CVDINFR4</th><th>CVDCRHD4</th><th>_MICHD</th><th>_STATE</th><th>SEXVAR</th><th>_AGEGR4</th><th>_RACEGR4</th><th>_EDUCAG</th><th>_INCOMG1</th><th>MARITAL</th><th>...</th><th>HEIGHT3</th><th>_BMIS6</th><th>_BMISCAT</th></tr></thead><tbody><tr><td>0</td><td>2.0</td><td>2.0</td><td>2.0</td><td>10</td><td>2.0</td><td>13.0</td><td>1.0</td><td>4.0</td><td>9.0</td><td>10</td><td>...</td><td>9999.0</td><td>NaN</td><td>NaN</td></tr><tr><td>1</td><td>2.0</td><td>2.0</td><td>2.0</td><td>10</td><td>2.0</td><td>13.0</td><td>1.0</td><td>2.0</td><td>3.0</td><td>3.0</td><td>...</td><td>503.0</td><td>2653.0</td><td>3.0</td></tr><tr><td>2</td><td>2.0</td><td>2.0</td><td>2.0</td><td>10</td><td>2.0</td><td>8.0</td><td>1.0</td><td>4.0</td><td>6.0</td><td>10</td><td>...</td><td>502.0</td><td>2561.0</td><td>3.0</td></tr><tr><td>3</td><td>2.0</td><td>2.0</td><td>2.0</td><td>10</td><td>2.0</td><td>14.0</td><td>1.0</td><td>2.0</td><td>9.0</td><td>10</td><td>...</td><td>505.0</td><td>2330.0</td><td>2.0</td></tr><tr><td>4</td><td>2.0</td><td>2.0</td><td>2.0</td><td>10</td><td>2.0</td><td>5.0</td><td>1.0</td><td>3.0</td><td>10.0</td><td>10</td><td>...</td><td>502.0</td><td>2177.0</td><td>2.0</td></tr><tr><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td></tr><tr><td>445127</td><td>2.0</td><td>2.0</td><td>2.0</td><td>78.0</td><td>2.0</td><td>1.0</td><td>2.0</td><td>2.0</td><td>5.0</td><td>...</td><td>...</td><td>505.0</td><td>2563.0</td><td>3.0</td></tr><tr><td>445128</td><td>2.0</td><td>2.0</td><td>2.0</td><td>78.0</td><td>2.0</td><td>7.0</td><td>2.0</td><td>4.0</td><td>6.0</td><td>10</td><td>...</td><td>507.0</td><td>2868.0</td><td>3.0</td></tr><tr><td>445129</td><td>2.0</td><td>2.0</td><td>2.0</td><td>78.0</td><td>2.0</td><td>10.0</td><td>9.0</td><td>2.0</td><td>9.0</td><td>10</td><td>...</td><td>507.0</td><td>1723.0</td><td>1.0</td></tr><tr><td>445130</td><td>1.0</td><td>2.0</td><td>1.0</td><td>78.0</td><td>1.0</td><td>11.0</td><td>2.0</td><td>3.0</td><td>5.0</td><td>10</td><td>...</td><td>600.0</td><td>3255.0</td><td>4.0</td></tr><tr><td>445131</td><td>2.0</td><td>2.0</td><td>2.0</td><td>78.0</td><td>1.0</td><td>5.0</td><td>2.0</td><td>1.0</td><td>2.0</td><td>3.0</td><td>...</td><td>506.0</td><td>2260.0</td><td>2.0</td></tr></tbody></table>	CVDINFR4	CVDCRHD4	_MICHD	_STATE	SEXVAR	_AGEGR4	_RACEGR4	_EDUCAG	_INCOMG1	MARITAL	...	HEIGHT3	_BMIS6	_BMISCAT	0	2.0	2.0	2.0	10	2.0	13.0	1.0	4.0	9.0	10	...	9999.0	NaN	NaN	1	2.0	2.0	2.0	10	2.0	13.0	1.0	2.0	3.0	3.0	...	503.0	2653.0	3.0	2	2.0	2.0	2.0	10	2.0	8.0	1.0	4.0	6.0	10	...	502.0	2561.0	3.0	3	2.0	2.0	2.0	10	2.0	14.0	1.0	2.0	9.0	10	...	505.0	2330.0	2.0	4	2.0	2.0	2.0	10	2.0	5.0	1.0	3.0	10.0	10	...	502.0	2177.0	2.0	445127	2.0	2.0	2.0	78.0	2.0	1.0	2.0	2.0	5.0	505.0	2563.0	3.0	445128	2.0	2.0	2.0	78.0	2.0	7.0	2.0	4.0	6.0	10	...	507.0	2868.0	3.0	445129	2.0	2.0	2.0	78.0	2.0	10.0	9.0	2.0	9.0	10	...	507.0	1723.0	1.0	445130	1.0	2.0	1.0	78.0	1.0	11.0	2.0	3.0	5.0	10	...	600.0	3255.0	4.0	445131	2.0	2.0	2.0	78.0	1.0	5.0	2.0	1.0	2.0	3.0	...	506.0	2260.0	2.0	Remove features where more than 75% of the responses are NaN. In retrospect, I could have first removed users where the majority of responses were NaN before removing features.
CVDINFR4	CVDCRHD4	_MICHD	_STATE	SEXVAR	_AGEGR4	_RACEGR4	_EDUCAG	_INCOMG1	MARITAL	...	HEIGHT3	_BMIS6	_BMISCAT																																																																																																																																																																							
0	2.0	2.0	2.0	10	2.0	13.0	1.0	4.0	9.0	10	...	9999.0	NaN	NaN																																																																																																																																																																						
1	2.0	2.0	2.0	10	2.0	13.0	1.0	2.0	3.0	3.0	...	503.0	2653.0	3.0																																																																																																																																																																						
2	2.0	2.0	2.0	10	2.0	8.0	1.0	4.0	6.0	10	...	502.0	2561.0	3.0																																																																																																																																																																						
3	2.0	2.0	2.0	10	2.0	14.0	1.0	2.0	9.0	10	...	505.0	2330.0	2.0																																																																																																																																																																						
4	2.0	2.0	2.0	10	2.0	5.0	1.0	3.0	10.0	10	...	502.0	2177.0	2.0																																																																																																																																																																						
...																																																																																																																																																																							
445127	2.0	2.0	2.0	78.0	2.0	1.0	2.0	2.0	5.0	505.0	2563.0	3.0																																																																																																																																																																						
445128	2.0	2.0	2.0	78.0	2.0	7.0	2.0	4.0	6.0	10	...	507.0	2868.0	3.0																																																																																																																																																																						
445129	2.0	2.0	2.0	78.0	2.0	10.0	9.0	2.0	9.0	10	...	507.0	1723.0	1.0																																																																																																																																																																						
445130	1.0	2.0	1.0	78.0	1.0	11.0	2.0	3.0	5.0	10	...	600.0	3255.0	4.0																																																																																																																																																																						
445131	2.0	2.0	2.0	78.0	1.0	5.0	2.0	1.0	2.0	3.0	...	506.0	2260.0	2.0																																																																																																																																																																						
5	<pre>dd22c = dd22c.loc[~dd22c.index.isin(removed_cols)] dd22c.head(3)</pre> <table border="1"><thead><tr><th>Format</th><th>Description</th><th>Survey Question</th><th>Possible values</th><th>Example Num</th></tr></thead><tbody><tr><td>Data Element Name</td><td></td><td></td><td></td><td></td></tr><tr><td>CVDINFR4</td><td>Num</td><td>Ever Diagnosed with Heart Attack</td><td>(Ever told) you had a heart attack, also calle...</td><td>{1: 'Yes', 2: 'No', 7: 'Don't know/Not s...</td></tr><tr><td>CVDCRHD4</td><td>Num</td><td>Ever Diagnosed with Angina or Coronary Heart D...</td><td>(Ever told) (you had) angina or coronary heart...</td><td>{1: 'Yes', 2: 'No', 7: 'Don't know/Not s...</td></tr><tr><td>_MICHD</td><td>Num</td><td>Ever had CHD or MI</td><td>Respondents that have ever reported having cor...</td><td>{1: 'Reported having MI or CHD', 2: 'Did n...</td></tr></tbody></table>	Format	Description	Survey Question	Possible values	Example Num	Data Element Name					CVDINFR4	Num	Ever Diagnosed with Heart Attack	(Ever told) you had a heart attack, also calle...	{1: 'Yes', 2: 'No', 7: 'Don't know/Not s...	CVDCRHD4	Num	Ever Diagnosed with Angina or Coronary Heart D...	(Ever told) (you had) angina or coronary heart...	{1: 'Yes', 2: 'No', 7: 'Don't know/Not s...	_MICHD	Num	Ever had CHD or MI	Respondents that have ever reported having cor...	{1: 'Reported having MI or CHD', 2: 'Did n...	Prune the data dictionary dataframe to match the dataset's remaining features. 36 features remain																																																																																																																																																									
Format	Description	Survey Question	Possible values	Example Num																																																																																																																																																																																
Data Element Name																																																																																																																																																																																				
CVDINFR4	Num	Ever Diagnosed with Heart Attack	(Ever told) you had a heart attack, also calle...	{1: 'Yes', 2: 'No', 7: 'Don't know/Not s...																																																																																																																																																																																
CVDCRHD4	Num	Ever Diagnosed with Angina or Coronary Heart D...	(Ever told) (you had) angina or coronary heart...	{1: 'Yes', 2: 'No', 7: 'Don't know/Not s...																																																																																																																																																																																
_MICHD	Num	Ever had CHD or MI	Respondents that have ever reported having cor...	{1: 'Reported having MI or CHD', 2: 'Did n...																																																																																																																																																																																

The final dataset contains all numeric values, with most of the values being identifiers that map to a natural language description using the 'Possible values' header in the data dictionary dataframe.

Snapshot of the data dictionary

Data Element Name	Format	Description	Possible values	Example Char	Example	Renamed Variable
CVDINF4	Char	Ever Diagnosed with Heart Attack	{'1': 'Yes', '2': 'No', '7': 'Don't know/Not sure', '9': 'Refused', 'BLANK': 'Not asked'}	Yes	2	heart_attack
CVDCRHD4	Char	Ever Diagnosed with Angina or Coronary Heart Disease	{'1': 'Yes', '2': 'No', '7': 'Don't know/Not sure', '9': 'Refused', 'BLANK': 'Not asked'}	Yes	2	angina_or_chd
_MICHD	Char	Ever had CHD or MI	{'1': 'Reported having MI or CHD', '2': 'Did not report having MI or CHD', 'BLANK': 'Not asked'}	Reported having MI or CHD	2	chd_or_mi
_STATE	Char	State FIPS Code	{'1': 'Alabama', '2': 'Alaska', '4': 'Arizona', '5': 'Arkansas', '6': 'California', '8': 'Colorado', '9': 'Hawaii', '0': 'Illinois', '1': 'Indiana', '2': 'Iowa', '3': 'Michigan', '4': 'Minnesota', '5': 'Mississippi', '6': 'Missouri', '7': 'North Dakota', '8': 'South Dakota', '9': 'Tennessee', 'A': 'Texas', 'B': 'Utah', 'C': 'Vermont', 'D': 'Washington', 'E': 'West Virginia', 'F': 'Wyoming'}	Alabama	1	state
SEXVAR	Char	Sex of Respondent	{'1': 'Male', '2': 'Female'}	Male	2	sex
_AGEG5YR	Char	Reported age in five-year age categories calculated variable	{'1': 'Age 18 to 24', '2': 'Age 25 to 29', '3': 'Age 30 to 34', '4': 'Age 35 to 39', '5': 'Age 40 to 44', '6': 'Age 45 to 49', '7': 'Age 50 to 54', '8': 'Age 55 to 59', '9': 'Age 60 to 64', '0': 'Age 65 to 69', '1': 'Age 70 to 74', '2': 'Age 75 to 79', '3': 'Age 80 to 84', '4': 'Age 85 to 89', '5': 'Age 90 to 94'}	Age 18 to 24	13	age
_RACEGR4	Char	Computed Five level race/ethnicity category.	{'1': 'White only, Non-Hispanic', '2': 'Black only, Non-Hispanic', '3': 'Other race or ethnicity, Non-Hispanic', '4': 'White only, Hispanic', '5': 'Black only, Hispanic', '6': 'Other race or ethnicity, Hispanic'}	White only, Non-Hispanic	1	ethnicity
_EDUCAG	Char	Computed level of education completed categories	{'1': 'Did not graduate High School', '2': 'Graduated High School', '3': 'Attended College', '4': 'Completed College', '5': 'Postsecondary'}	Did not graduate High School	2	education
_INCOMG1	Char	Computed income categories	{'1': 'Less than \$15,000', '2': '\$15,000 to < \$25,000', '3': '\$25,000 to < \$35,000', '4': '\$35,000 to < \$50,000', '5': '\$50,000 to < \$75,000', '6': '\$75,000 to < \$100,000', '7': '\$100,000 or more'}	Less than \$15,000	3	income
MARITAL	Char	Marital Status	{'1': 'Married', '2': 'Divorced', '3': 'Widowed', '4': 'Separated', '5': 'Never married'}	Married	3	marital
EMPLOY1	Char	Employment Status	{'1': 'Employed for wages', '2': 'Self-employed', '3': 'Out of work for 1 year or more', '4': 'Out of work less than 1 year', '5': 'Student', '6': 'Retired', '7': 'Disabled', '8': 'Homemaker', '9': 'Other'}	Employed for wages	2	employment
SMOKE100	Char	Smoked at Least 100 Cigarettes	{'1': 'Yes', '2': 'No', '7': 'Don't know/Not Sure', '9': 'Refused', 'BLANK': 'Not asked'}	Yes	2	smoker_past
USENOW3	Char	Use of Smokeless Tobacco Products	{'1': 'Every day', '2': 'Some days', '3': 'Not at all', '7': 'Don't know/Not Sure', '9': 'Refused', 'BLANK': 'Not asked'}	Every day	3	smoker_now
ECIGNOW2	Char	Do you now use e-cigarettes, or vaping products everyday	{'1': 'Never used e-cigarettes in your entire life', '2': 'Use them every day', '3': 'Use them some days', '7': 'Don't know/Not Sure', '9': 'Refused', 'BLANK': 'Not asked'}	Never used e-cigarettes in your entire life	1	vaper
ALCDAY4	Char	Days in past 30 had alcoholic beverage	{'101 - 199': 'Days per week', '201 - 299': 'Days in past 30 days', '777': 'Don't know/Not Sure', '999': 'Refused', 'BLANK': 'Not asked'}	Days per week	888	alcohol
EXERANY2	Char	Exercise in Past 30 Days	{'1': 'Yes', '2': 'No', '7': 'Don't know/Not Sure', '9': 'Refused', 'BLANK': 'Not asked'}	Yes	2	exercise
SLEPTIM1	Num	How Much Time Do You Sleep	{'1 - 24': 'Number of hours [1-24]', '77': 'Don't know/Not Sure', '99': 'Refused', 'BLANK': 'Not asked'}	1 - 24	6	sleep
GENHLTH	Char	General Health	{'1': 'Excellent', '2': 'Very good', '3': 'Good', '4': 'Fair', '5': 'Poor', '7': 'Don't know/Not Sure', '9': 'Refused', 'BLANK': 'Not asked'}	Excellent	1	general_health
PHYSHLTH	Num	Number of Days Physical Health Not Good	{'1 - 30': 'Number of days', '88': 'None', '77': 'Don't know/Not sure', '99': 'Refused', 'BLANK': 'Not asked'}	1 - 30	88	physical_health_days
MENTHLTH	Num	Number of Days Mental Health Not Good	{'1 - 30': 'Number of days', '88': 'None', '77': 'Don't know/Not sure', '99': 'Refused', 'BLANK': 'Not asked'}	1 - 30	88	mental_health_days
DIABETE4	Char	(Ever told) you had diabetes	{'1': 'Yes', '2': 'Yes, but female told only during pregnancy', '3': 'No', '4': 'No, pre-diabetes', '7': 'Don't know/Not Sure', '9': 'Refused', 'BLANK': 'Not asked'}	Yes	3	diabetes
HAVARTH4	Char	Told Had Arthritis	{'1': 'Yes', '2': 'No', '7': 'Don't know/Not Sure', '9': 'Refused', 'BLANK': 'Not asked'}	Yes	2	arthritis
CHCCOPD3	Char	Ever told you had C.O.P.D. emphysema or chronic bronchitis	{'1': 'Yes', '2': 'No', '7': 'Don't know / Not sure', '9': 'Refused', 'BLANK': 'Not asked'}	Yes	2	copd
ADDEPEV3	Char	(Ever told) you had a depressive disorder	{'1': 'Yes', '2': 'No', '7': 'Don't know/Not sure', '9': 'Refused', 'BLANK': 'Not asked'}	Yes	2	depression
ASTHMA3	Char	Ever Told Had Asthma	{'1': 'Yes', '2': 'No', '7': 'Don't know/Not Sure', '9': 'Refused', 'BLANK': 'Not asked'}	Yes	2	asthma
WEIGHT2	Num	Reported Weight in Pounds	{'50 - 0776': 'Weight (pounds)', '7777': 'Don't know/Not sure', '9023 - 9352': 'Refused', 'BLANK': 'Not asked'}	50 - 0776	150	weight
HEIGHT3	Num	Reported Height in Feet and Inches	{'200 - 711': 'Height (ft/inches)', '7777': 'Don't know/Not sure', '9061 - 9998': 'Refused', 'BLANK': 'Not asked'}	200 - 711	503	height
_BMIS5	Num	Computed body mass index	{'1 - 9999': '1 or greater', 'BLANK': 'Don't know/Refused/Missing'}	1 or greater	2657	bmi
_BMIS5CAT	Char	Computed body mass index categories	{'1': 'Underweight', '2': 'Normal Weight', '3': 'Overweight', '4': 'Obese', 'BLANK': 'Not asked'}	Underweight	3	bmi_category
PRIMINSR	Char	What is Primary Source of Health Insurance?	{'1': 'A plan purchased through an employer or union (including plans purchased through a union)', '2': 'A plan purchased through an employer or union (including plans purchased through a union)', '3': 'Other', '4': 'Refused', '5': 'Don't know/Not Sure', '6': 'Not asked'}	A plan purchased through an employer or union (including plans purchased through a union)	3	health_insurance
MEDCOST1	Char	Could Not Afford To See Doctor	{'1': 'Yes', '2': 'No', '7': 'Don't know/Not sure', '9': 'Refused', 'BLANK': 'Not asked'}	Yes	2	low_ses
CHECKUP1	Num	Length of time since last routine checkup	{'1': 'Within past year (anytime less than 12 months ago)', '2': 'Within past 2 years', '3': 'Within past 5 years', '4': 'Within past 10 years', '5': 'Within past 15 years', '6': 'Within past 20 years', '7': 'Within past 25 years', '8': 'Within past 30 years', '9': 'Within past 35 years', '0': 'Within past 40 years', '1': 'Within past 45 years', '2': 'Within past 50 years', '3': 'Within past 55 years', '4': 'Within past 60 years', '5': 'Within past 65 years', '6': 'Within past 70 years', '7': 'Within past 75 years', '8': 'Within past 80 years', '9': 'Within past 85 years', '0': 'Within past 90 years', '1': 'Within past 95 years', '2': 'Within past 100 years', '3': 'Don't know/Not sure', '4': 'Refused', '5': 'Not asked'}	Within past year (anytime less than 12 months ago)	8	checkup_days
COVIDPOS	Char	Have you ever been told you tested positive for COVID 19?	{'1': 'Yes', '2': 'No', '3': 'Tested positive using home test without health professional', '4': 'Tested positive using professional test', '5': 'Don't know/Not sure', '6': 'Refused', '7': 'Not asked'}	Yes	2	covid
FLUSHOT7	Char	Adult flu shot/spray past 12 mos	{'1': 'Yes', '2': 'No', '7': 'Don't know/Not Sure', '9': 'Refused', 'BLANK': ''}	Yes	2	flu_shot
PNEUVAC4	Char	Pneumonia shot ever	{'1': 'Yes', '2': 'No', '7': 'Don't know/Not Sure', '9': 'Refused', 'BLANK': 'Not asked'}	Yes	2	pneumonia_shot
TETANUS1	Char	Received Tetanus Shot Since 2005?	{'1': 'Yes, received Tdap', '2': 'Yes, received tetanus shot, but not Tdap', '3': 'Yes, i Yes, received Tdap', '4': 'Yes, received tetanus shot, but not Tdap', '5': 'Don't know/Not sure', '6': 'Refused', '7': 'Not asked'}	Yes, received Tdap	4	tetanus_shot

Project Overview

Objective	<p>To examine the prevalence and distribution of cardiovascular diseases (CVD) in the United States across demographic groups and evaluate the influence of socioeconomic factors and lifestyle choices on CVD outcomes.</p> <p>RQ1 (Demographics): What is the prevalence of CVD across different age, sex, and ethnic groups?</p> <p>RQ2 (Socioeconomics): Are there significant differences in CVD prevalence based on socioeconomic factors such as income, education, or insurance status?</p> <p>RQ3 (Lifestyle): How do behavioral and lifestyle choices such as smoking, alcohol use, and physical activity influence the prevalence of CVD?</p>
Hypothesis	<p>RQ1 (Demographics):</p> <ul style="list-style-type: none"> - CVD prevalence <u>increases with age</u> as the body naturally weakens over its lifetime. - CVD prevalence is <u>higher for females</u> with children as the heart undergoes physiological changes during pregnancy that may be detrimental to heart health. - CVD prevalence is <u>higher for the Black and Hispanic populations</u> compared to the White population due to increased socioeconomic disparities that deter access to healthcare as well as potential implicit bias in medical care. <p>RQ2 (Socioeconomics): CVD prevalence is inversely correlated with socioeconomic status as it limits access to healthcare.</p> <p>RQ3 (Lifestyle): Increased exercise and decreased substance use leads to decreased CVD prevalence. Exercise is known to improve heart physiology while substance use has adverse effects on the brain which negatively affects the heart through hormone and neurotransmitter release.</p>
Motivation	<p>CVD comprises some of the world's leading causes of death. By studying cardiovascular diseases, we can develop a deeper understanding of how demographic, socioeconomic, and lifestyle factors impact heart health. This knowledge can help:</p> <p>(RQ1) Identify high-risk groups in healthcare settings</p> <p>(RQ2) Inform policies that promote health equity and access to care</p> <p>(RQ3) Inspire individuals to make healthier lifestyle choices</p>
Approach	<p>Exploratory analysis will be done to determine variable inclusion, variable distributions, and causes of missing data. All approaches will involve controlling for confounding variables through preprocessing. Final results will report the p-values, effect size, and confidence intervals.</p> <p>For the modeling, I want to cover the broad topics taught in this course from statistics (descriptive, inferential) to machine learning (decision trees)</p> <p>RQ1 (Demographics): Use descriptive statistics to describe CVD prevalence rates across each demographic group (age, sex, and ethnicity)</p> <p>RQ2 (Socioeconomics): Train tree-based model(s) to rank socioeconomic factors by their impact on CVD prevalence. Emphasize the explainability of decision trees for describing model outputs through examples and visualization</p> <p>RQ3 (Lifestyle): Use multiple regression to model the relationship between lifestyle factors and CVD prevalence</p>