

Diabetes Prediction Model

Using Machine Learning Regression Algorithms and Neural Networks



Index

Presentation by Nastaran Ashoori & Daniella Vuletić
CodeOp DSPT10 - July 2023



1- Introduction

**2- Data Pre-
Processing & Feature
Engineering**

**3- Exploratory
Data Analysis
(EDA)**

**4- Machine
Learning (ML)**

5- Conclusion

Diabetes Prediction Model

Using Machine Learning Regression Algorithms and Neural Networks



Presentation by Nastaran Ashoori & Daniella Vuletić
CodeOp DSPT10 - July 2023

Diabetes Prediction Model

Using Machine Learning Regression Algorithms and Neural Networks



Index

Presentation by Nastaran Ashoori & Daniella Vuletić
CodeOp DSPT10 - July 2023



1- Introduction

**2- Data Pre-
Processing & Feature
Engineering**

**3- Exploratory
Data Analysis
(EDA)**

**4- Machine
Learning (ML)**

5- Conclusion

Index

Index

1 - Introduction

Index

1 - Introduction

2 - Data Pre-processing & Feature Engineering

Index

1 - Introduction

2 - Data Pre-processing & Feature Engineering

3 - Exploratory Data Analysis (EDA)

Index

1 - Introduction

2 - Data Pre-processing & Feature Engineering

3 - Exploratory Data Analysis (EDA)

4 - Machine Learning (ML)

Index

1 - Introduction

2 - Data Pre-processing & Feature Engineering

3 - Exploratory Data Analysis (EDA)

4 - Machine Learning (ML)

5 - Conclusions

Diabetes Prediction Model

Using Machine Learning Regression Algorithms and Neural Networks



Index

Presentation by Nastaran Ashoori & Daniella Vuletić
CodeOp DSPT10 - July 2023



1- Introduction

**2- Data Pre-
Processing & Feature
Engineering**

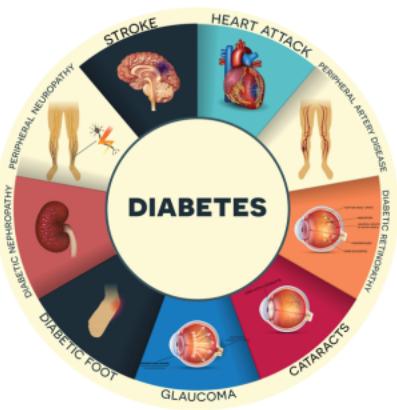
**3- Exploratory
Data Analysis
(EDA)**

**4- Machine
Learning (ML)**

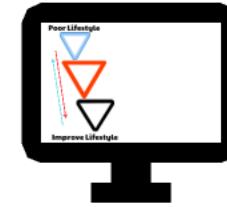
5- Conclusion

1- INTRODUCTION

What is Diabetes?



What can be done?



Why Study Diabetes



Diabetes prevalence, 2000
The share of people aged 20-79 who have diabetes¹.



- 527 million adults (20-79 years) living with diabetes – 1 in 10.
- Number of adults with diabetes is projected to rise to 643 million by 2030 and 783 million by 2045.
- One in 4 adults with diabetes live in low- and middle-income countries.
- Diabetes was responsible for 6.7 million deaths in 2013 – 1 every 6 seconds.

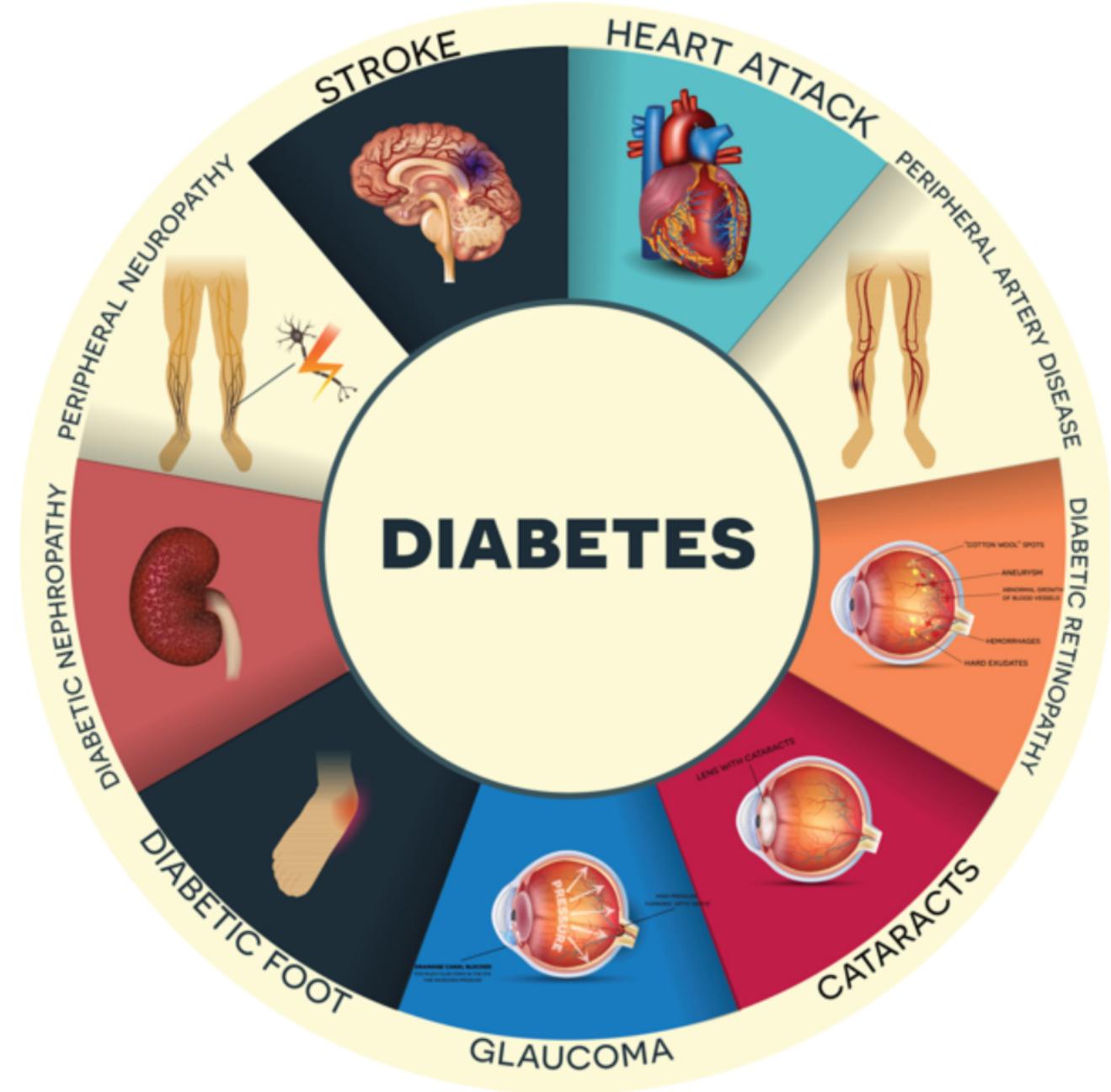
Source: International Diabetes Federation (via World Bank)

OurWorldInData.org/united-of-diabetes | CC BY

Study Objective

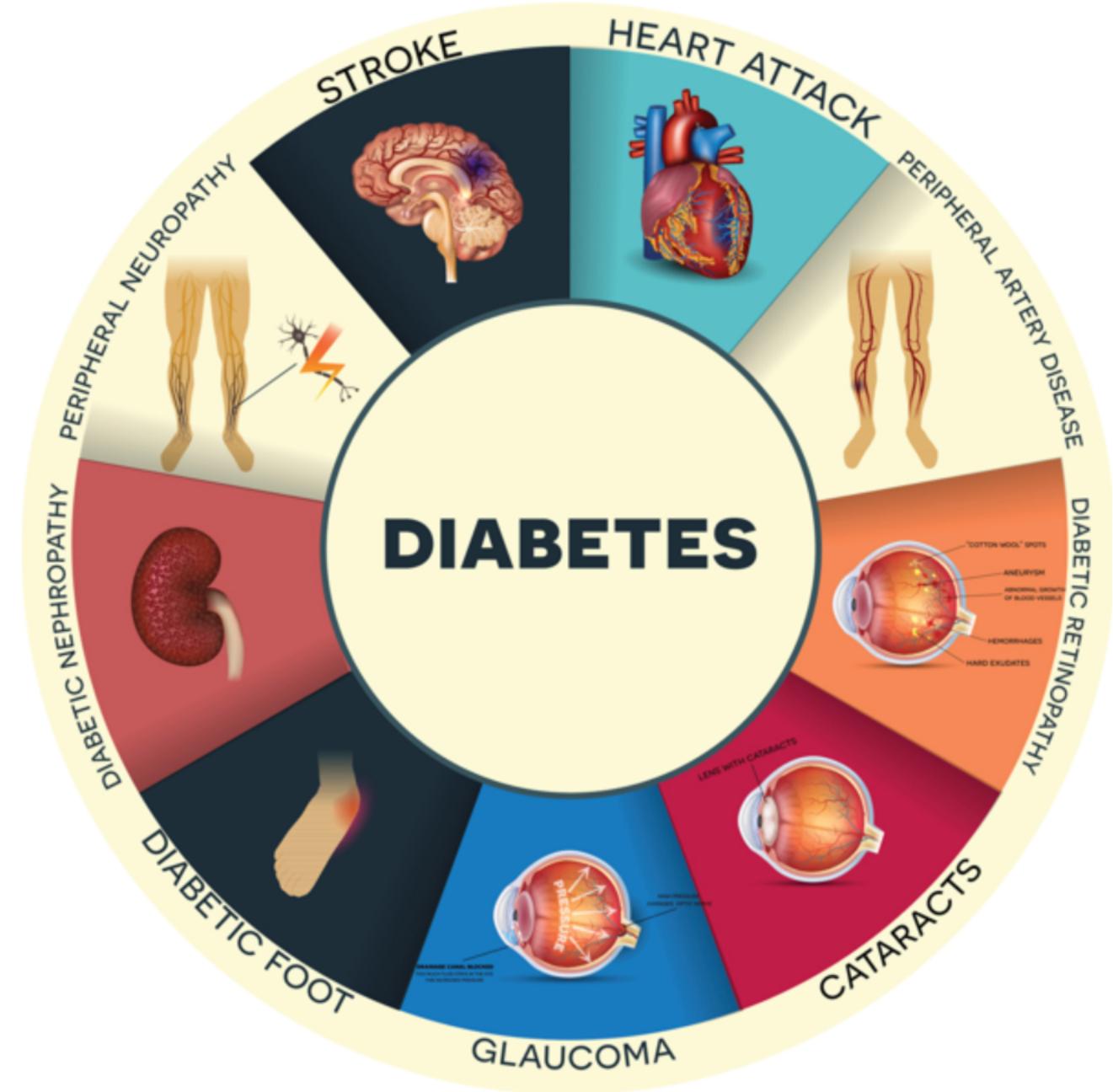


What is Diabetes?



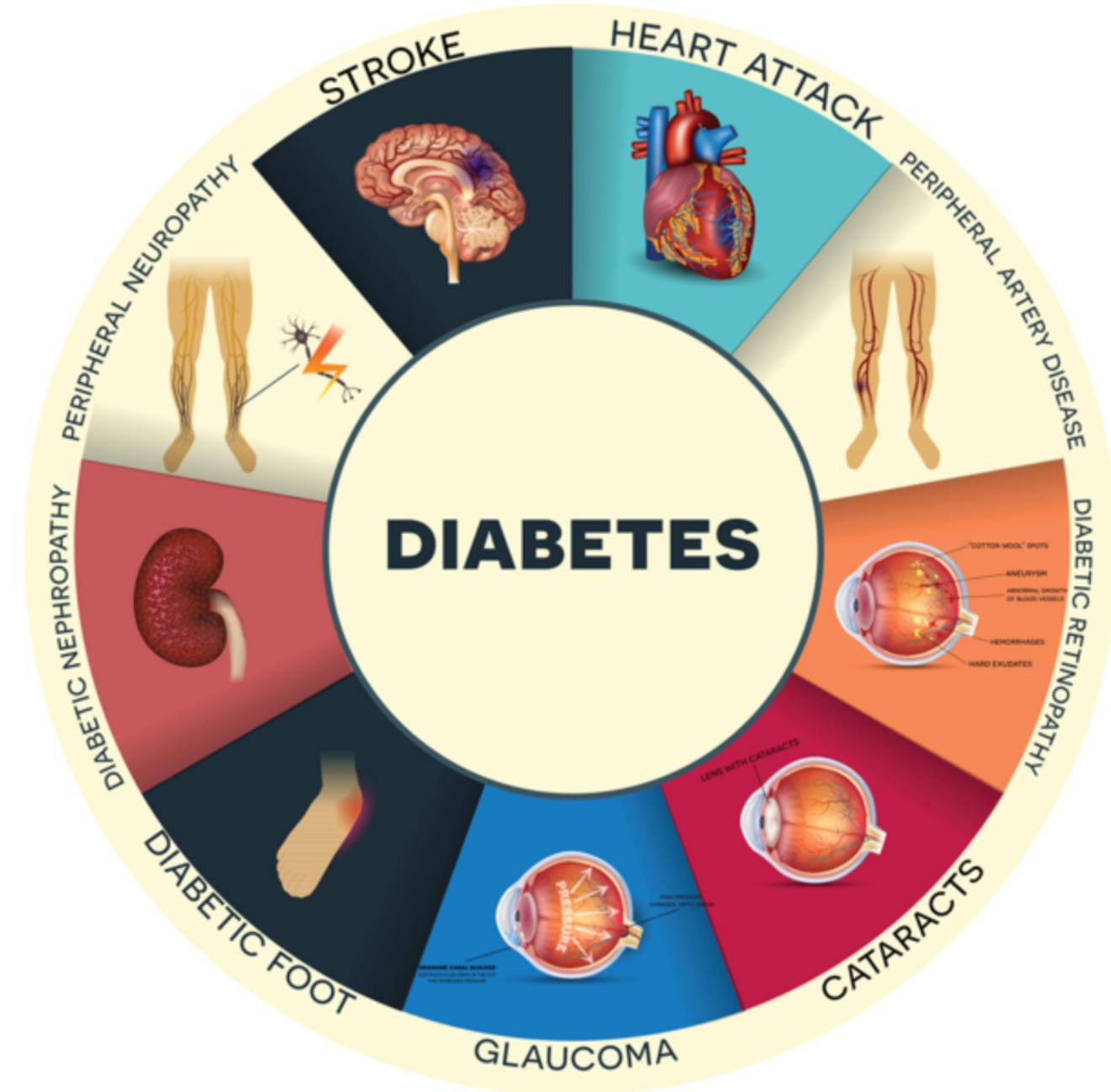
What is Diabetes?

- Endocrine disease



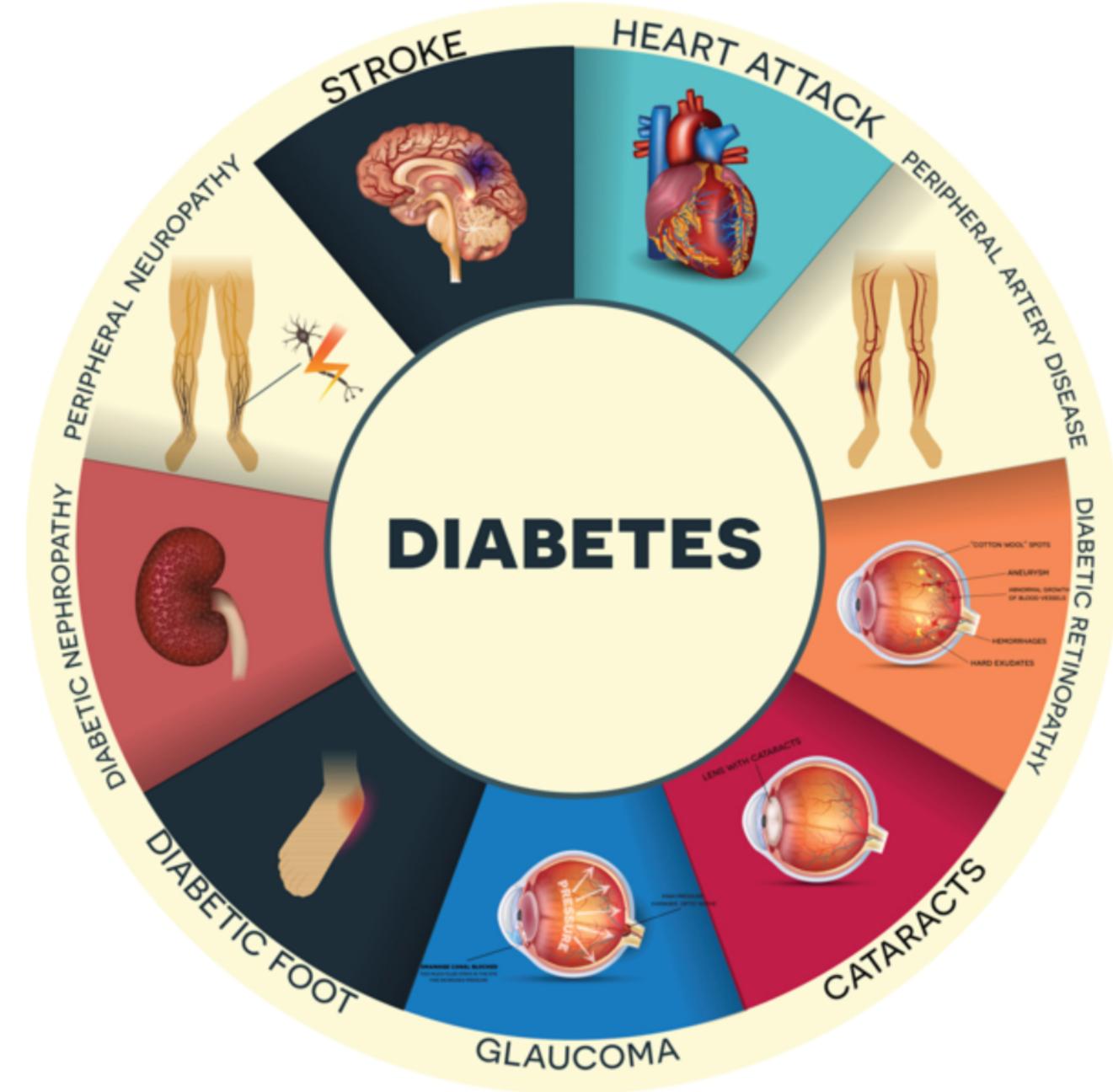
What is Diabetes?

- Endocrine disease
- Sustained **high** blood sugar



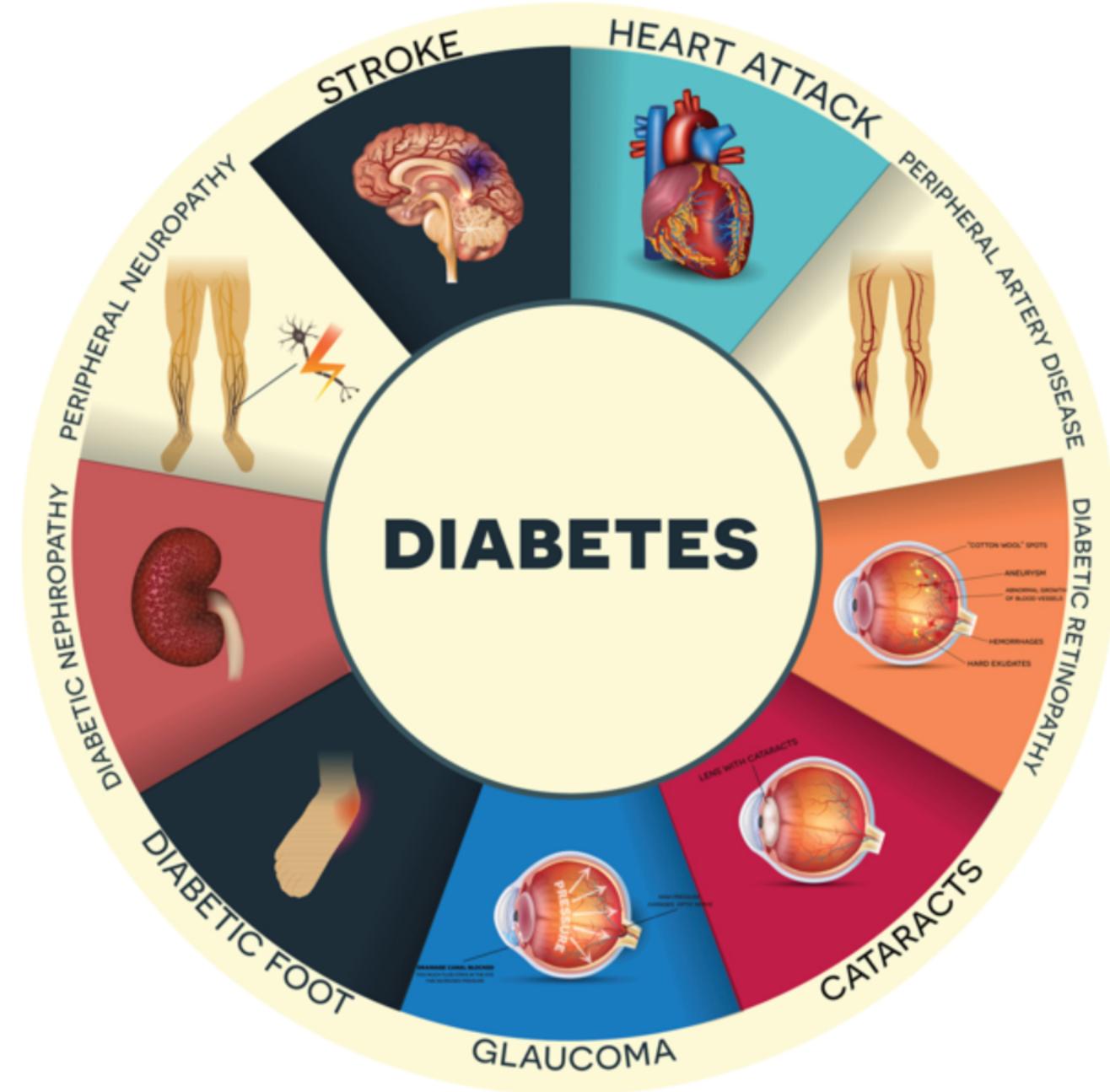
What is Diabetes?

- Endocrine disease
- Sustained **high** blood sugar
- Type 1 :
 - 10% of cases
 - No insulin production
 - Autoimmune disease, complex pathogenesis, poorly understood



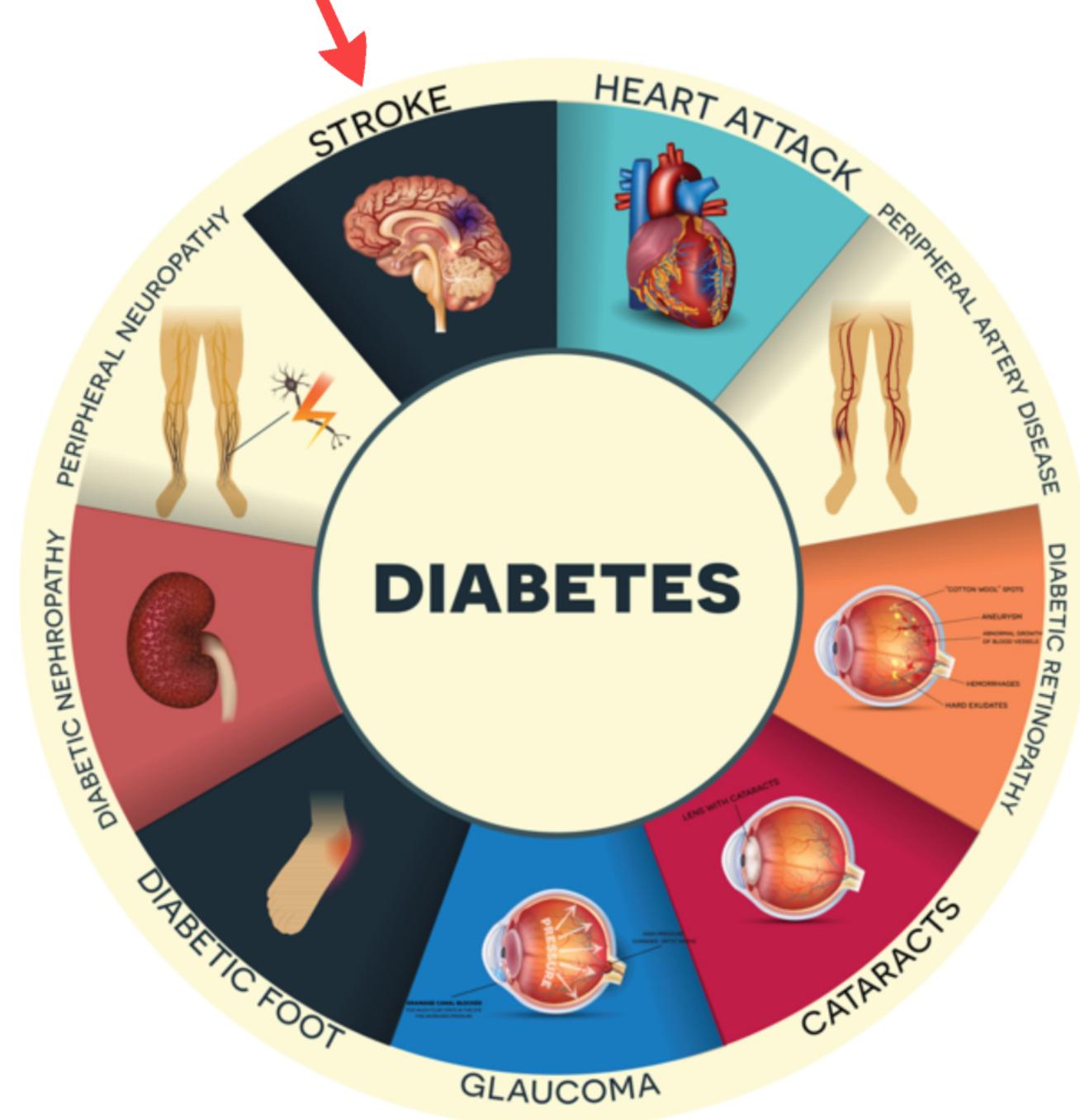
What is Diabetes?

- Endocrine disease
- Sustained **high** blood sugar
- Type 1 :
 - 10% of cases
 - No insulin production
 - Autoimmune disease, complex pathogenesis, poorly understood
- Type 2 :
 - > 90% of cases
 - High insulin production
 - No insulin response by cells
 - Primarily a lifestyle disease, but genetics and family history play a role



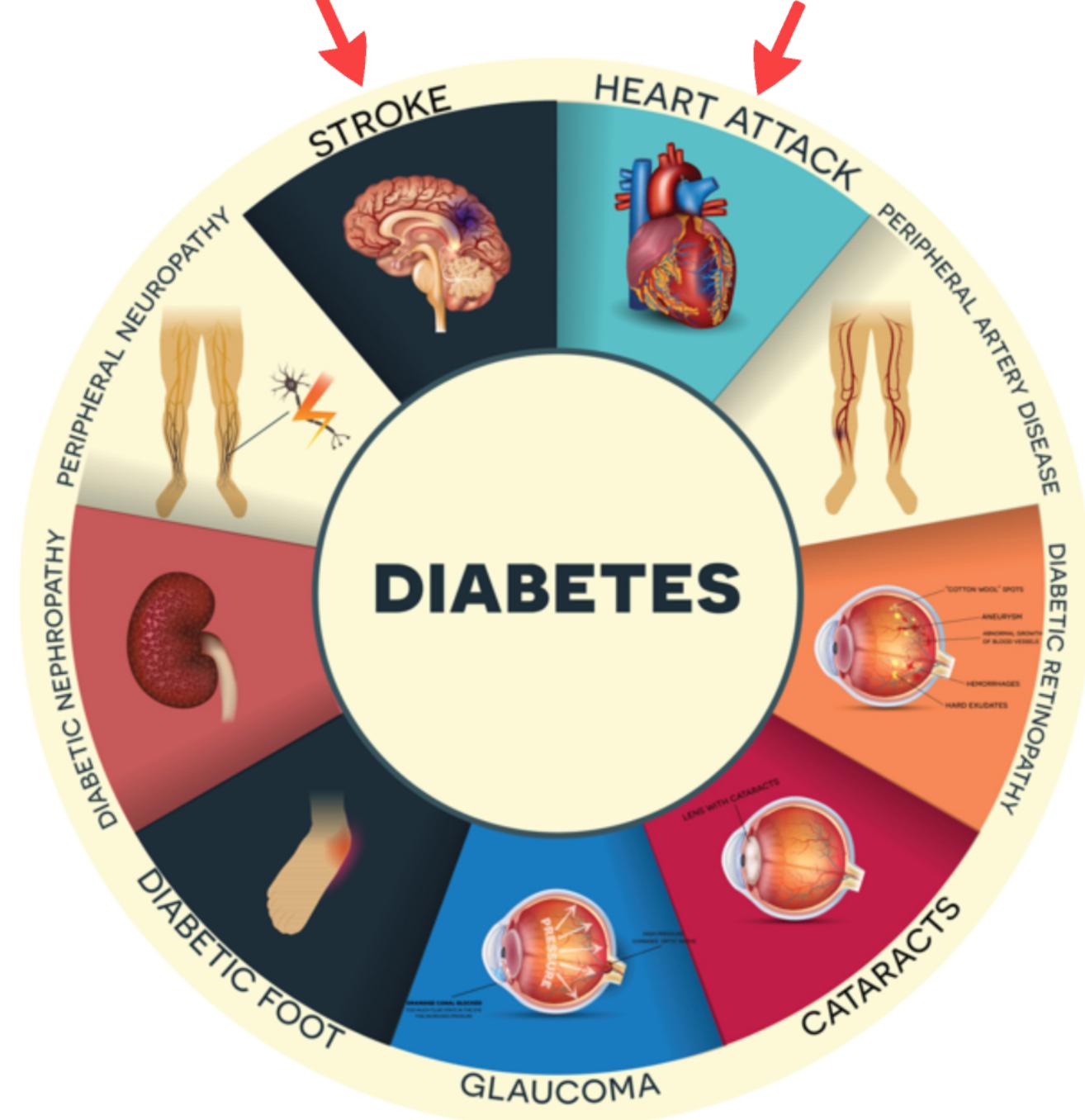
What is Diabetes?

- Endocrine disease
- Sustained **high** blood sugar
- Type 1 :
 - 10% of cases
 - No insulin production
 - Autoimmune disease, complex pathogenesis, poorly understood
- Type 2 :
 - > 90% of cases
 - High insulin production
 - No insulin response by cells
 - Primarily a lifestyle disease, but genetics and family history play a role



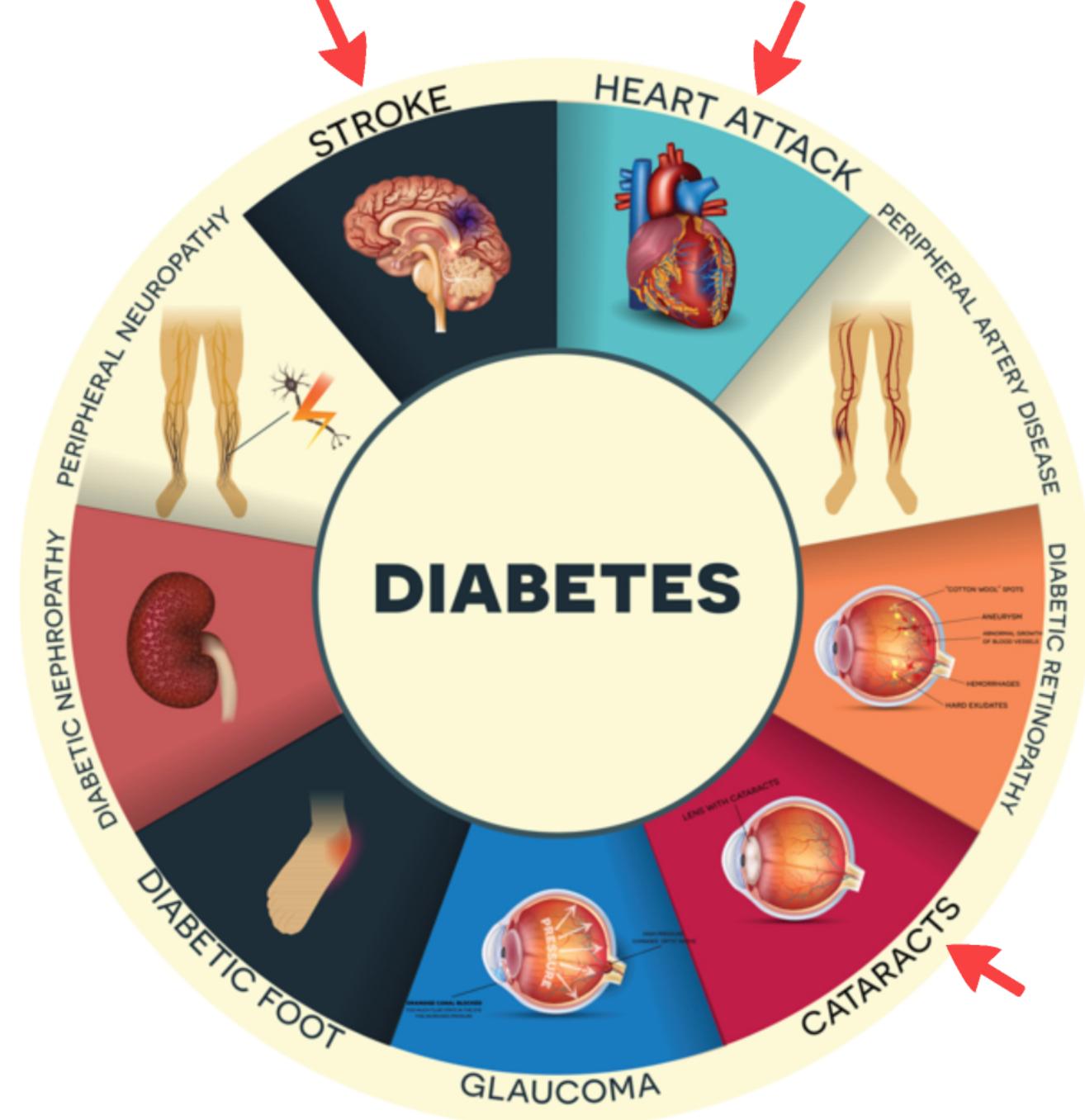
What is Diabetes?

- Endocrine disease
- Sustained **high** blood sugar
- Type 1 :
 - 10% of cases
 - No insulin production
 - Autoimmune disease, complex pathogenesis, poorly understood
- Type 2 :
 - > 90% of cases
 - High insulin production
 - No insulin response by cells
 - Primarily a lifestyle disease, but genetics and family history play a role



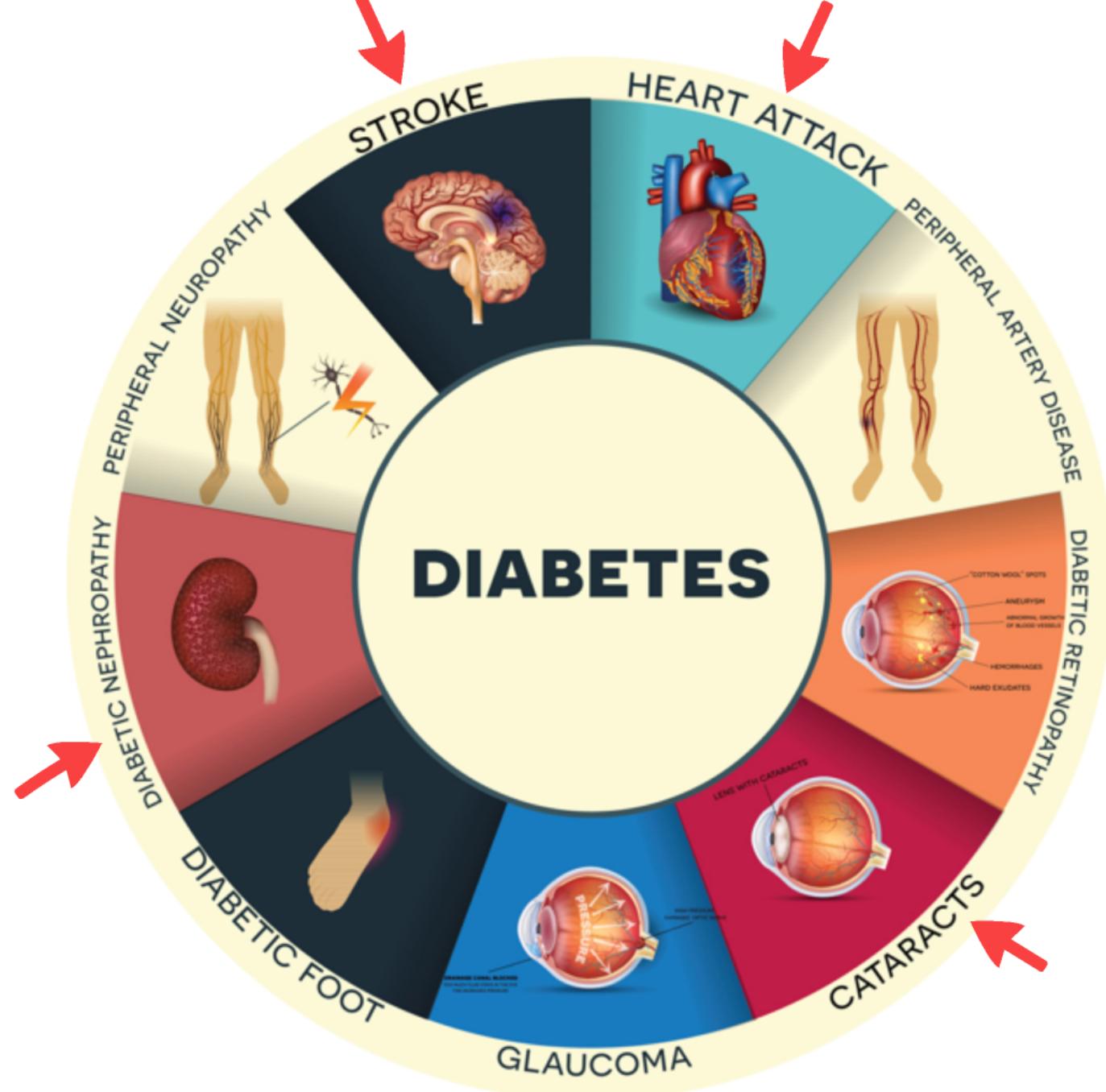
What is Diabetes?

- Endocrine disease
- Sustained **high** blood sugar
- Type 1 :
 - 10% of cases
 - No insulin production
 - Autoimmune disease, complex pathogenesis, poorly understood
- Type 2 :
 - > 90% of cases
 - High insulin production
 - No insulin response by cells
 - Primarily a lifestyle disease, but genetics and family history play a role



What is Diabetes?

- Endocrine disease
- Sustained **high** blood sugar
- Type 1 :
 - 10% of cases
 - No insulin production
 - Autoimmune disease, complex pathogenesis, poorly understood
- Type 2 :
 - > 90% of cases
 - High insulin production
 - No insulin response by cells
 - Primarily a lifestyle disease, but genetics and family history play a role



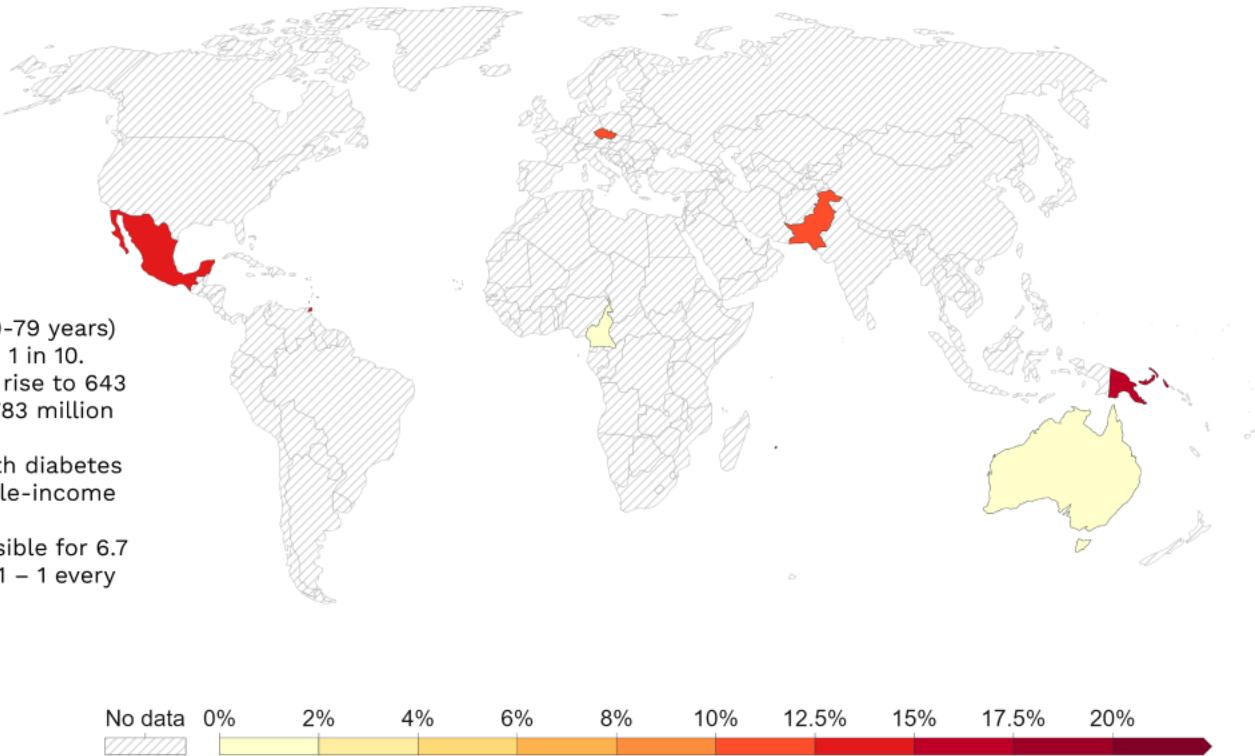
Why Study Diabetes



- 537 million adults (20-79 years) living with diabetes – 1 in 10. Number predicted to rise to 643 million by 2030 and 783 million by 2045.
- Over 3 in 4 adults with diabetes live in low- and middle-income countries.
- Diabetes was responsible for 6.7 million deaths in 2021 – 1 every 5 seconds.

Diabetes prevalence, 2000
The share of people aged 20-79 who have diabetes¹.

Our World
in Data

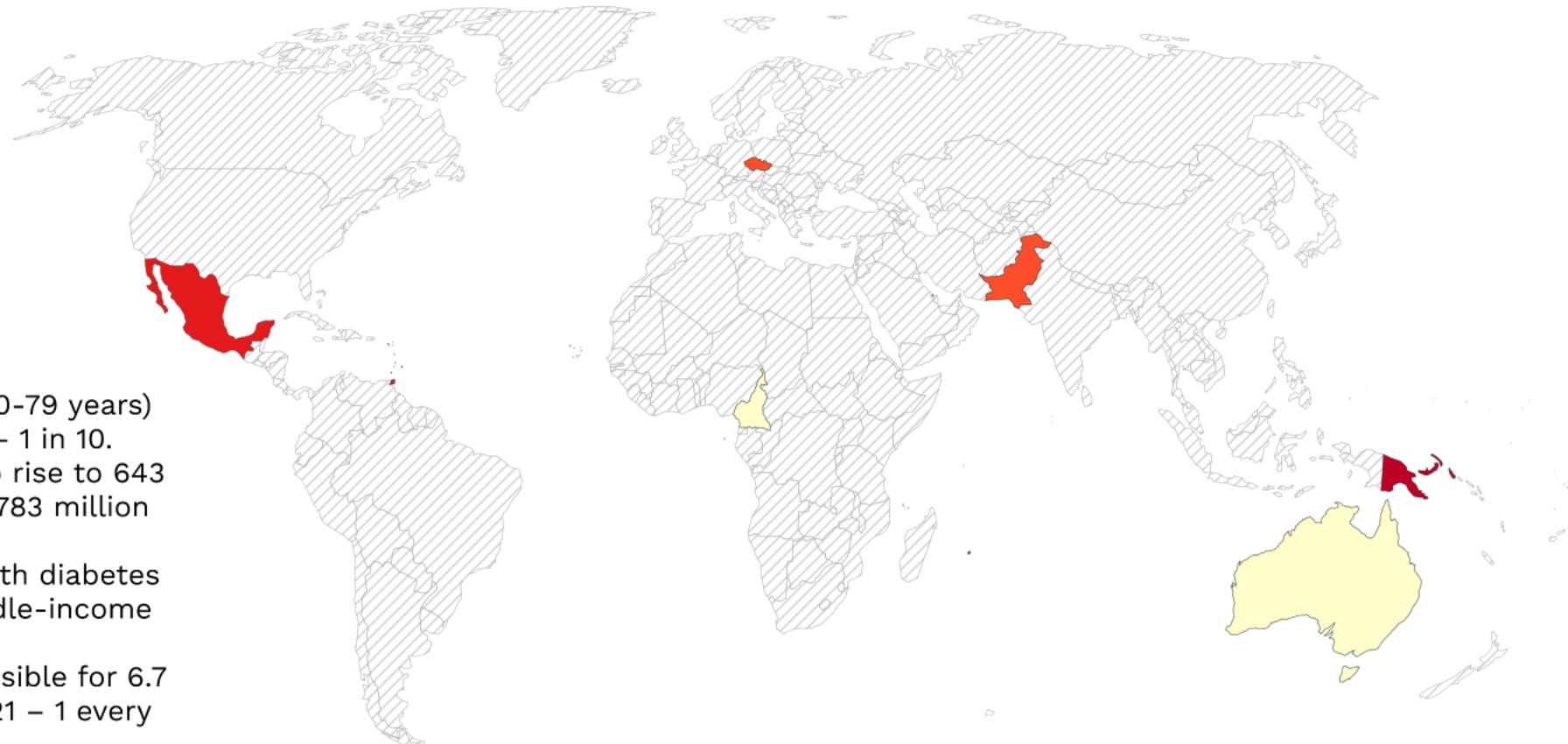


Source: International Diabetes Federation (via World Bank)

OurWorldInData.org/burden-of-disease • CC BY

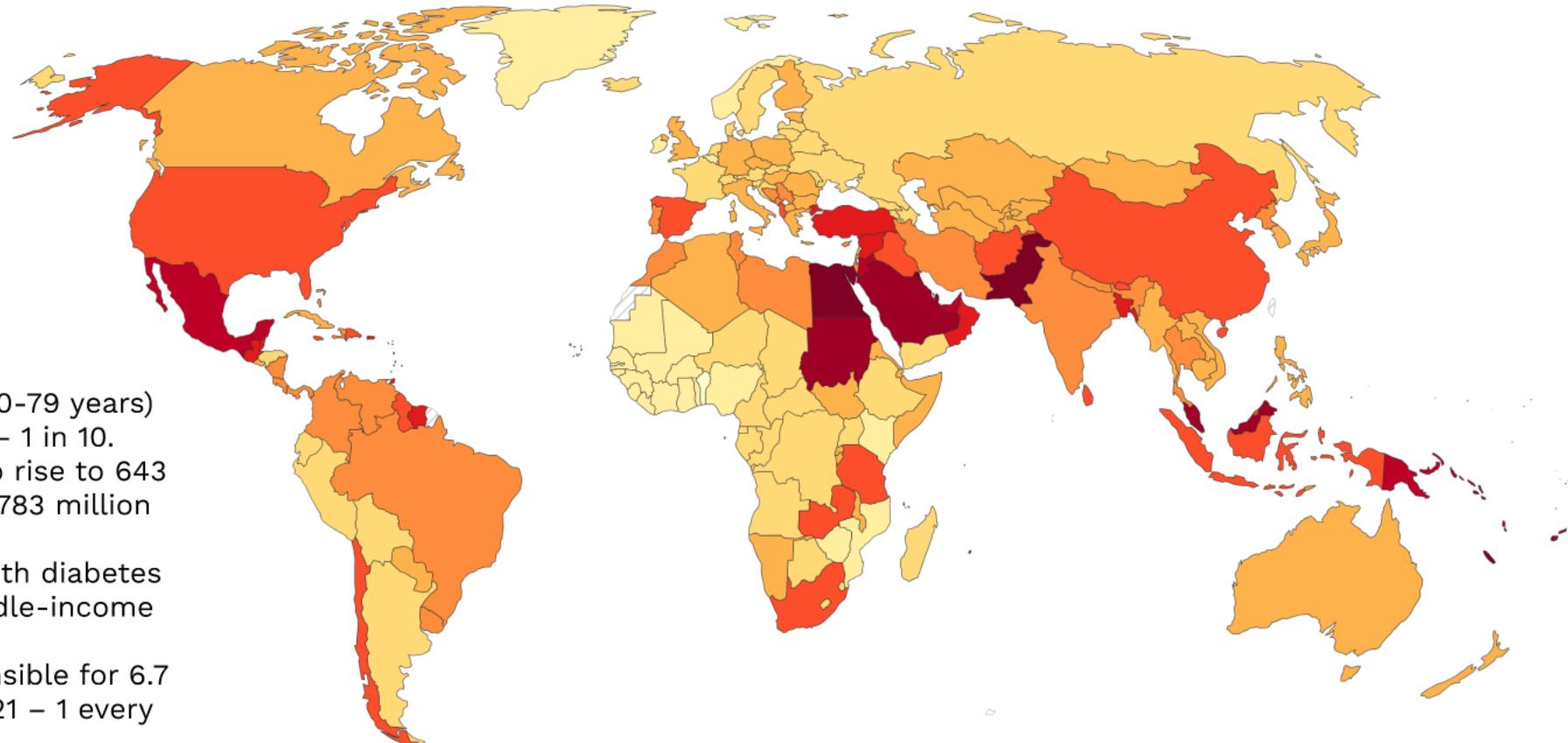
Diabetes prevalence, 2000

The share of people aged 20-79 who have diabetes¹.



Diabetes prevalence, 2021

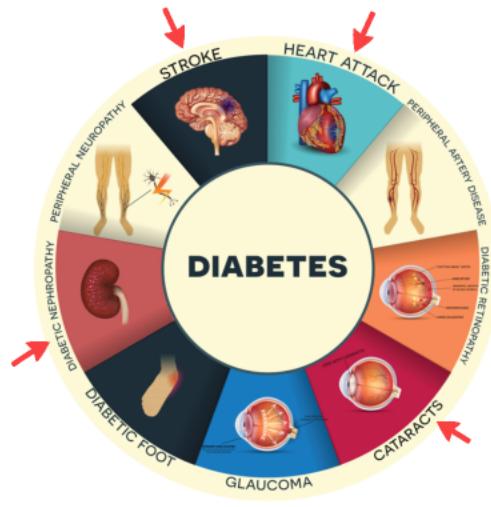
The share of people aged 20-79 who have diabetes¹.



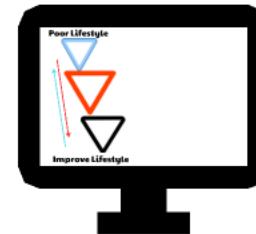
1- INTRODUCTION

What is Diabetes?

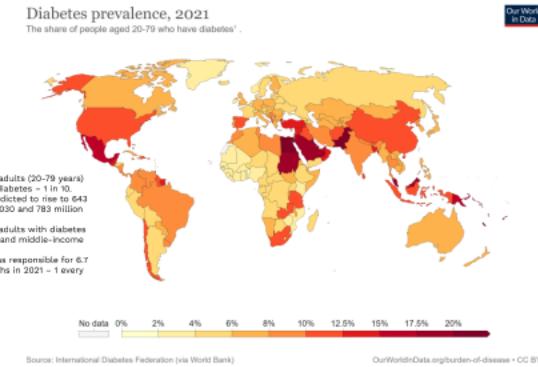
- Endocrine disease
- Sustained **high** blood sugar
- Type 1 :
 - 10% of cases
 - No insulin production
 - Autoimmune disease, complex pathogenesis, poorly understood
- Type 2 :
 - > 90% of cases
 - High insulin production
 - No insulin response by cells
 - Primarily a lifestyle disease, but genetics and family history play a role



What can be done?



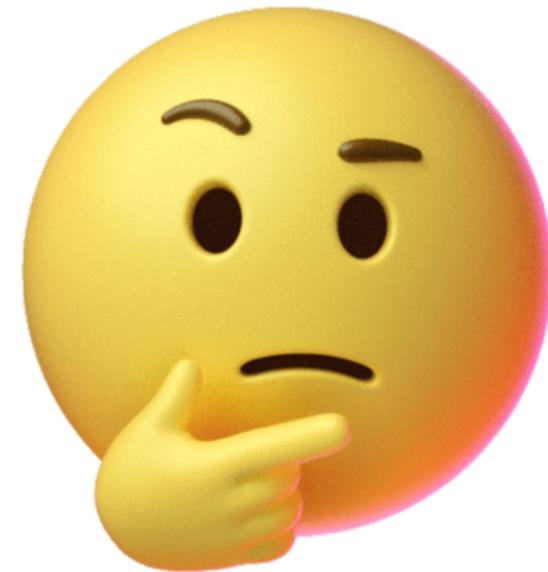
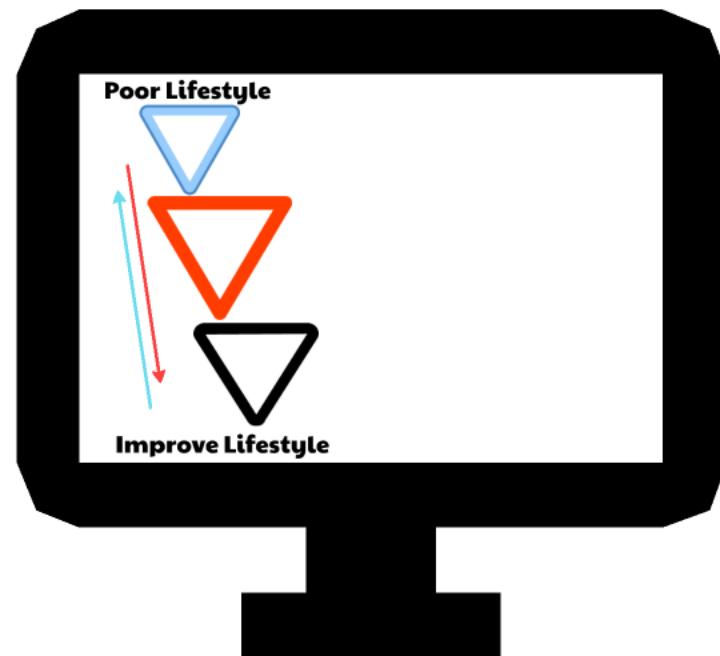
Why Study Diabetes



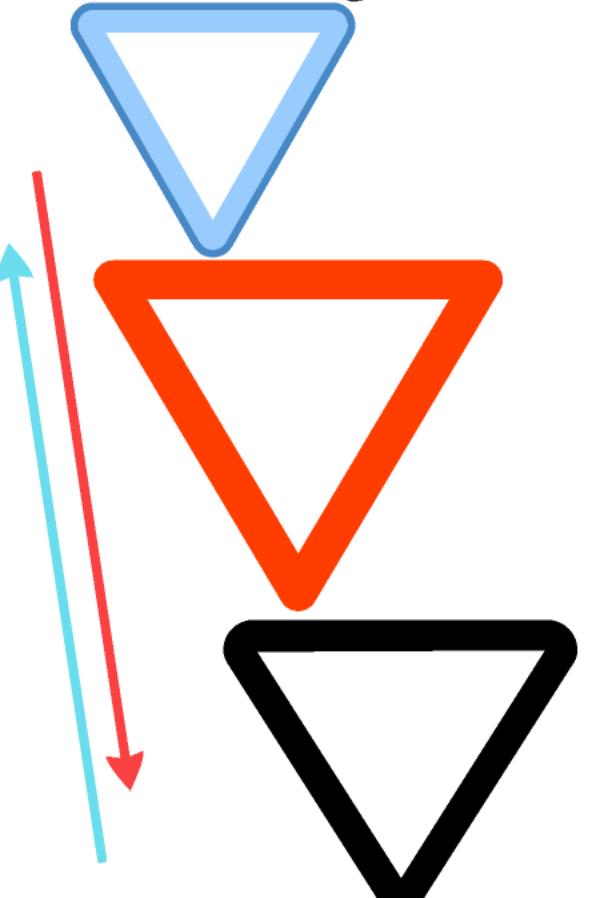
Study Objective



What can be done?

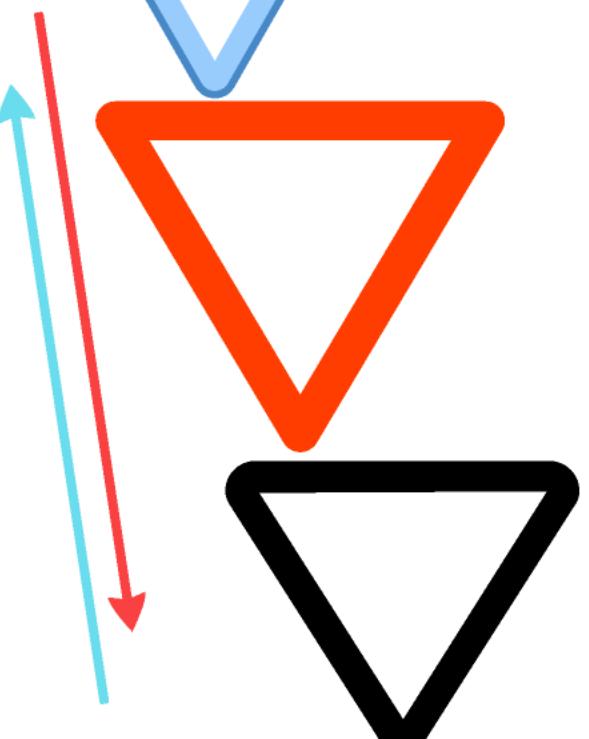


Poor Lifestyle



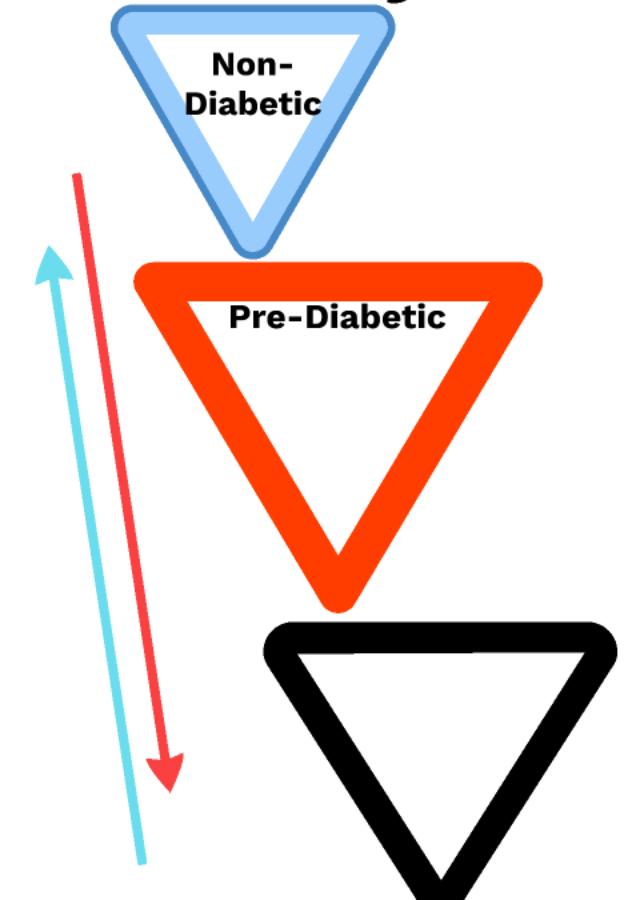
Improve Lifestyle

Poor Lifestyle



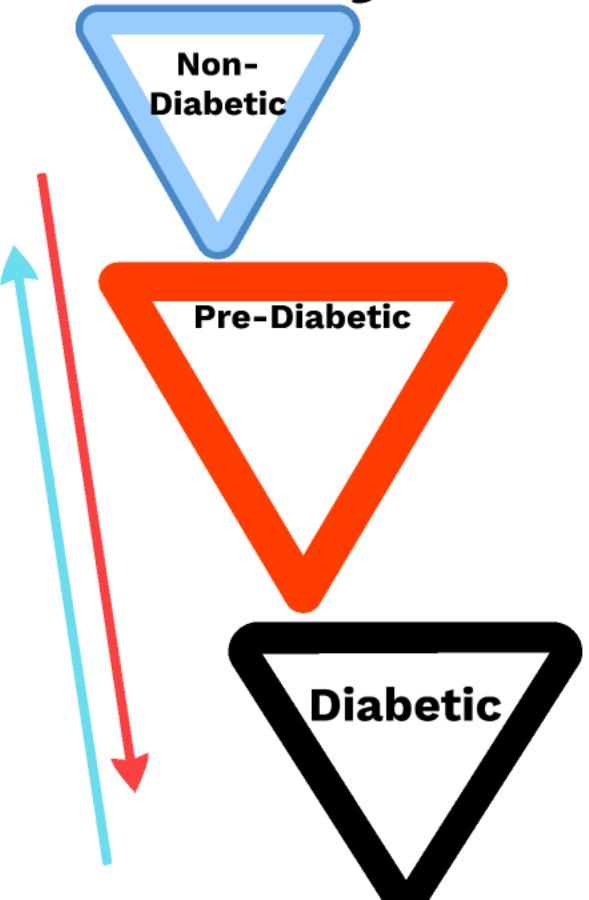
Improve Lifestyle

Poor Lifestyle



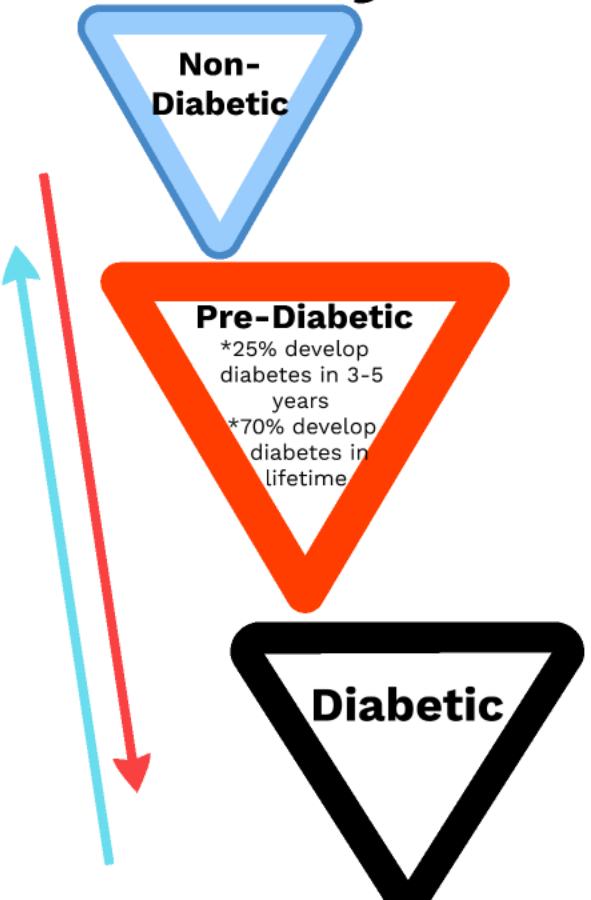
Improve Lifestyle

Poor Lifestyle



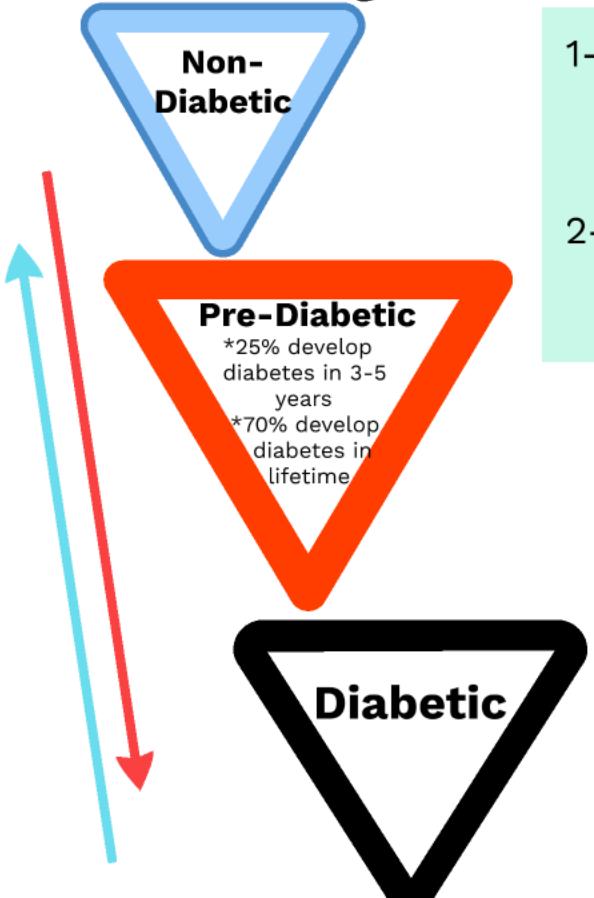
Improve Lifestyle

Poor Lifestyle



Improve Lifestyle

Poor Lifestyle



Improve Lifestyle

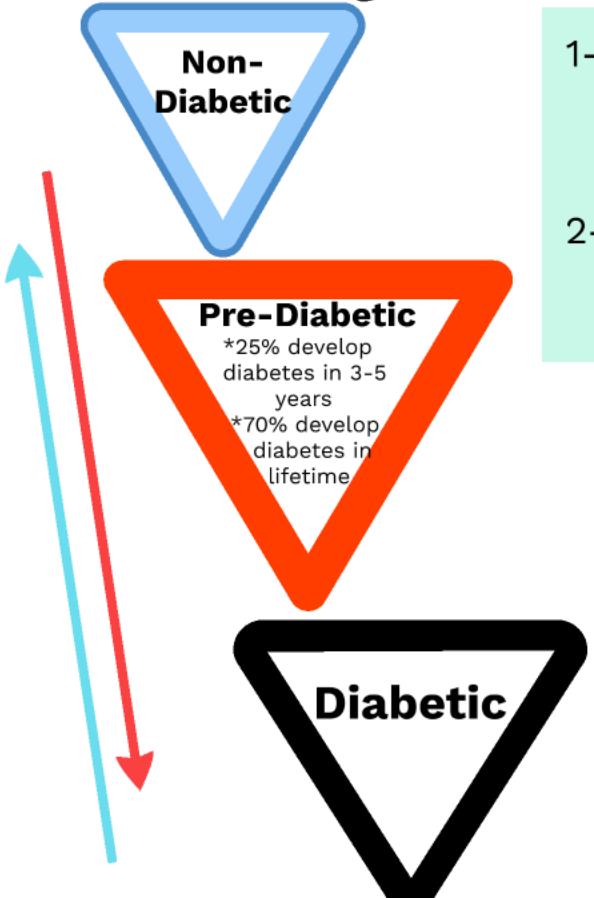
1- Blood work markers:

- Blood glucose (HbA1c)
- Fat profile (TG, Col, etc)

2- Physical changes:

- Body Mass Index (BMI)
- Age

Poor Lifestyle

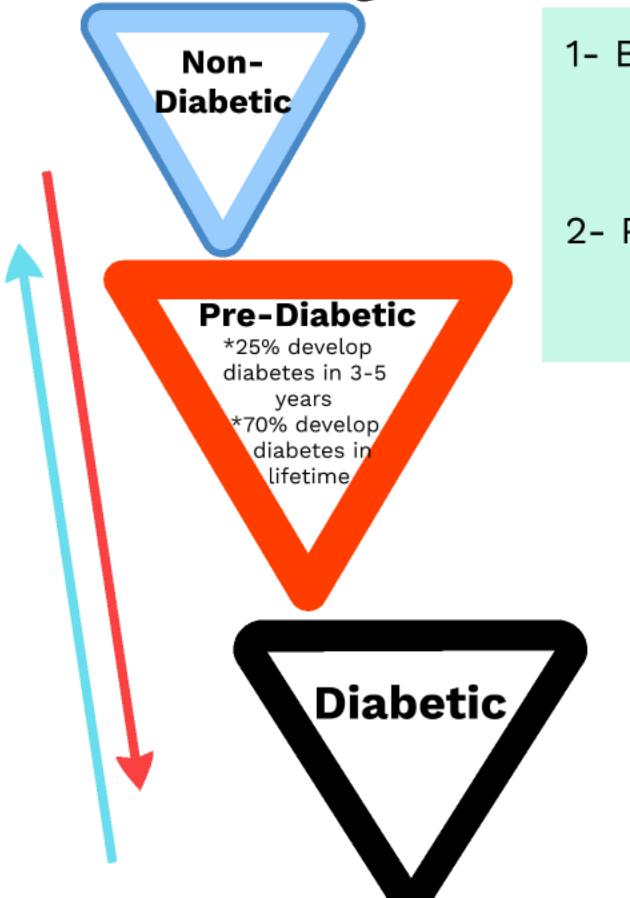


Improve Lifestyle

- 1- Blood work markers:
- Blood glucose (HbA1c)
 - Fat profile (TG, Col, etc)

- 2- Physical changes:
- Body Mass Index (BMI)
 - Age

Poor Lifestyle



- 1- Blood work markers:
 - Blood glucose (HbA1c)
 - Fat profile (TG, Col, etc)

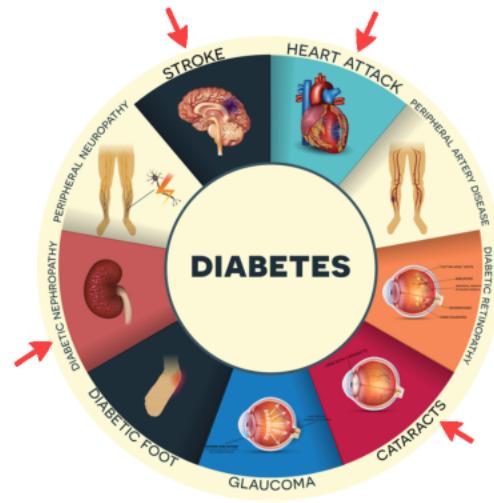
- 2- Physical changes:
 - Body Mass Index (BMI)
 - Age

- Identifying **what factors** put an individual at risk
- **Early accurate diagnosis** is especially important:
 - According to CDC **1 in 3 Americans** are pre-diabetic
 - More than **80% don't even know** they have it!
- Progression of pre-diabetes to diabetes can be **slowed** or **reversed** by lifestyle changes or medication

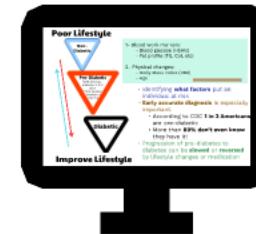
1- INTRODUCTION

What is Diabetes?

- Endocrine disease
- Sustained **high** blood sugar
- Type 1 :
 - 10% of cases
 - No insulin production
 - Autoimmune disease, complex pathogenesis, poorly understood
- Type 2 :
 - > 90% of cases
 - High insulin production
 - No insulin response by cells
 - Primarily a lifestyle disease, but genetics and family history play a role



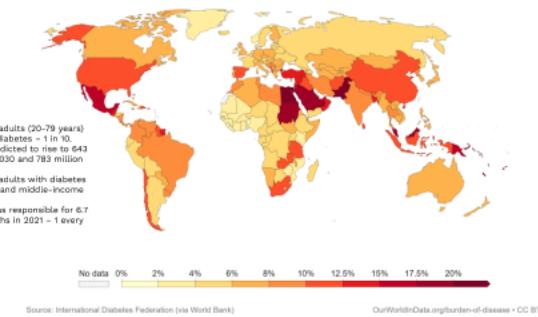
What can be done?



Why Study Diabetes



Diabetes prevalence, 2021
The share of people aged 20-79 who have diabetes¹.



Study Objective



Study Objective





Study Objective

"**Design a Predictive Diabetes Model** that can estimate the **likelihood of diabetes** development based on the provided factors".





Study Objective

"**Design a Predictive Diabetes Model** that can estimate the **likelihood of diabetes** development based on the provided factors".

Such a model can provide **valuable insights** for **healthcare professionals** and **researchers** in developing targeted interventions and preventative strategies.



The Dataset

- The data were collected from the **Iraqi society** – The laboratory of Medical City Hospital (Specialist Center for Endocrinology and Diabetes-Al-Kindy Teaching Hospital).
- Patients' files were taken and data extracted from them and entered in to the database to construct the diabetes dataset. The data consists of patients' medical information and laboratory analysis.
- The data is recorded as a table of **1000** patients represented in **rows** and **14** measured variables in **columns**.

ID	No_Pation	Gender	AGE	Urea	Cr	HbA1c	Chol	TG	HDL	LDL	VLDL	BMI	CLASS	
0	502	17975	F	50	4.7	46	4.9	4.2	0.9	2.4	1.4	0.5	24.0	N
1	735	34221	M	26	4.5	62	4.9	3.7	1.4	1.1	2.1	0.6	23.0	N
2	420	47975	F	50	4.7	46	4.9	4.2	0.9	2.4	1.4	0.5	24.0	N
3	680	87656	F	50	4.7	46	4.9	4.2	0.9	2.4	1.4	0.5	24.0	N
4	504	34223	M	33	7.1	46	4.9	4.9	1.0	0.8	2.0	0.4	21.0	N
...	
995	200	454317	M	71	11.0	97	7.0	7.5	1.7	1.2	1.8	0.6	30.0	Y
996	671	876534	M	31	3.0	60	12.3	4.1	2.2	0.7	2.4	15.4	37.2	Y
997	669	87654	M	30	7.1	81	6.7	4.1	1.1	1.2	2.4	8.1	27.4	Y
998	99	24004	M	38	5.8	59	6.7	5.3	2.0	1.6	2.9	14.0	40.5	Y
999	248	24054	M	54	5.0	67	6.9	3.8	1.7	1.1	3.0	0.7	33.0	Y

1000 rows × 14 columns

	ID	No_Pation	Gender	AGE	Urea	Cr	HbA1c	Chol	TG	HDL	LDL	VLDL	BMI	CLASS
0	502	17975	F	50	4.7	46	4.9	4.2	0.9	2.4	1.4	0.5	24.0	N
1	735	34221	M	26	4.5	62	4.9	3.7	1.4	1.1	2.1	0.6	23.0	N
2	420	47975	F	50	4.7	46	4.9	4.2	0.9	2.4	1.4	0.5	24.0	N
3	680	87656	F	50	4.7	46	4.9	4.2	0.9	2.4	1.4	0.5	24.0	N
4	504	34223	M	33	7.1	46	4.9	4.9	1.0	0.8	2.0	0.4	21.0	N
...
995	200	454317	M	71	11.0	97	7.0	7.5	1.7	1.2	1.8	0.6	30.0	Y
996	671	876534	M	31	3.0	60	12.3	4.1	2.2	0.7	2.4	15.4	37.2	Y
997	669	87654	M	30	7.1	81	6.7	4.1	1.1	1.2	2.4	8.1	27.4	Y
998	99	24004	M	38	5.8	59	6.7	5.3	2.0	1.6	2.9	14.0	40.5	Y
999	248	24054	M	54	5.0	67	6.9	3.8	1.7	1.1	3.0	0.7	33.0	Y

1000 rows × 14 columns



Variable	Description
ID	Patient ID
No_Pation	Patient number
Gender	Male or Female
Age	Age of patient in years
Urea	Nitrogenous end-product derived from dietary protein and tissue protein turnover
Cr	Creatinine is the nitrogenous end-product from muscle and protein metabolism and is an important indicator of kidney health
HbA1c	A form of haemoglobin chemically linked to sugar and is a measure to test the average blood sugar over the past 90 days
Chol	Cholesterol is a type of lipid (fat) that helps the body perform many important functions and comes from our liver and our diet - it travels through the blood on lipoproteins HDL and LDL
TG	Triglycerides come from our diet and are a type of lipid that circulate in our blood that the body uses for energy

	Range
	-
	-
	-
	-
tissue protein turnover	Normal range is 1.8 to 7.1 mmol urea/l
protein metabolism and is an measure to test the average blood	Normal range is 61.9 to 114.9 µmol/L for men is 53 to 97.2 µmol/L for women Measured as a percentage, a normal reading is between 4.0-5.6%, pre-diabetic 5.7-6.4%, and diabetic over 6.4%
form many important functions and blood on lipoproteins HDL and LDL	A normal healthy reading for total serum cholesterol should be below 5 mmol/L
it circulate in our blood that the body	Less than 1.7 mmol/L is considered normal, 1.8-2.2 mmol/L borderline high, 2.3-5.6 mmol/L high, and very high 5.7 mmol/L and above

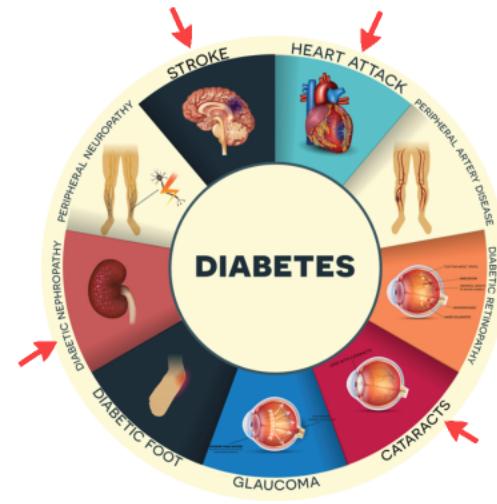
HbA1c	A form of haemoglobin chemically linked to sugar and is a measure to test the average blood sugar over the past 90 days
Chol	Cholesterol is a type of lipid (fat) that helps the body perform many important functions and comes from our liver and our diet - it travels through the blood on lipoproteins HDL and LDL
TG	Triglycerides come from our diet and are a type of lipid that circulate in our blood that the body uses for energy
HDL	High-density lipoprotein is known as the “good” cholesterol - it absorbs cholesterol in the blood and takes it to the liver to be flushed out the body
LDL	Low-density lipoprotein is known as the “bad” cholesterol - too much can build up in the walls of our blood vessels and cause health problems
VLDL	Very-low-density-lipoprotein is a type of LDL and carries triglycerides to the body’s cells and tissues - too much VLDL can build up in the walls of our blood vessels
BMI	Body mass index is a measure of body fat based on height and weight measured in kg/m ²
Class	Refers to diagnosis as N as non-diabetic, P as predicted-diabetic (pre-diabetic), and Y as diabetic

measure to test the average blood	Measured as a percentage, a normal reading is between 4.0-5.6%, pre-diabetic 5.7-6.4%, and diabetic over 6.4%
from many important functions and blood on lipoproteins HDL and LDL	A normal healthy reading for total serum cholesterol should be below 5 mmol/L
it circulate in our blood that the body	Less than 1.7 mmol/L is considered normal, 1.8-2.2 mmol/L borderline high, 2.3-5.6 mmol/L high, and very high 5.7 mmol/L and above
HDL - it absorbs cholesterol in the	HDL should ideally be above 1 mmol/L for men and 1.2 mmol/L for women
- too much can build up in the walls	LDL should ideally be below 3 mmol/L
triglycerides to the body's cells and blood vessels	A normal VLDL level is anything up to 0.77 mmol/L
and weight measured in kg/m ²	A BMI of less than 18.5 is underweight, 18.5-24.9 is healthy, 25.0-29.9 is overweight, and greater than 30.0 is obese
hyperglycemic (pre-diabetic), and Y as diabetic	

1- INTRODUCTION

What is Diabetes?

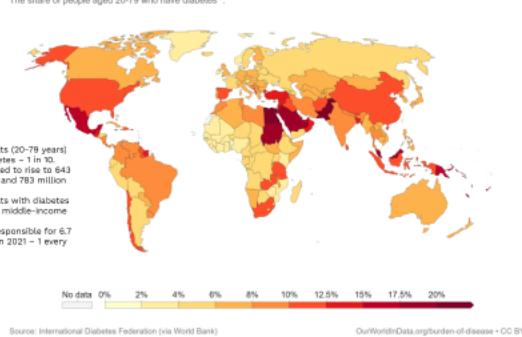
- Endocrine disease
- Sustained **high** blood sugar
- Type 1 :
 - 10% of cases
 - No insulin production
 - Autoimmune disease, complex pathogenesis, poorly understood
- Type 2 :
 - > 90% of cases
 - High insulin production
 - No insulin response by cells
 - Primarily a lifestyle disease, but genetics and family history play a role



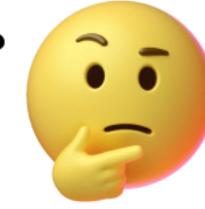
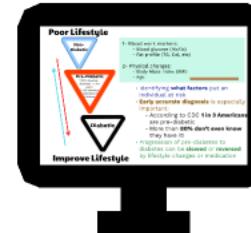
Why Study Diabetes



Diabetes prevalence, 2021
The share of people aged 20-79 who have diabetes¹.



What can be done?



Study Objective



"Design a Predictive Diabetes Model that can estimate the likelihood of diabetes development based on the provided factors".

Such a model can provide valuable insights for healthcare professionals and researchers in developing targeted interventions and preventative strategies.



Diabetes Prediction Model

Using Machine Learning Regression Algorithms and Neural Networks



Index

Presentation by Nastaran Ashoori & Daniella Vuletić
CodeOp DSPT10 - July 2023



1- Introduction

**2- Data Pre-
Processing & Feature
Engineering**

**3- Exploratory
Data Analysis
(EDA)**

**4- Machine
Learning (ML)**

5- Conclusion

2- DATA PRE-PROCESSING & FEATURE ENGINEERING

1. Overview of dataset

data.info()

• 2000 entries, 24 columns
• 2 object data types

2. We were particularly interested to see **how many categories** exist in "CLASS" & "Gender" columns as **object** data types, and **number of patients** in each category.

"CLASS" column = **"target variable"** classifies patients in 3 categories:
Non-diabetic (N), Pre-diabetic (P) and Diabetic(Y)

"Gender" column classifies patients in 2 categories:
Female (F) and Male (M)

value_counts()

CLASS	N	P	Y
count	984	102	12
gender	M	F	
dtype	int64	int64	int64

Redundant categories were found

3. Redundant categories were found, due to typing error!

So we replaced:

Any spaces (" ") with an empty string ("") **str.replace()**

CLASS	N	P	Y
count	984	102	12
gender	M	F	
dtype	int64	int64	int64

Lowercase letters with uppercase ones **str.upper()**

CLASS	N	P	Y
count	988	102	12
gender	M	F	
dtype	int64	int64	int64

4. For easier readability of the code we replaced the values in "CLASS" and "Gender" as following:

"N" to "Non-Diabetic"
"P" to "Pre-Diabetic"
"Y" to "Diabetic"
"F" to "Female"
"M" to "Male"

map()

CLASS	Non-Diabetic	Pre-Diabetic	Diabetic
count	988	102	12
gender	M	F	
dtype	int64	int64	int64

5. We then converted the "CLASS" and "Gender" columns from 'object' to 'category' data type:

df[["CLASS", "Gender"]].apply(lambda x: x.astype("category"))

CLASS	Non-Diabetic	Pre-Diabetic	Diabetic
count	988	102	12
gender	Male	Female	
dtype	category	category	category

6. Two numerical columns as "CLASS_Category" from the "CLASS" column and "Gender_Category" from "Gender" column were added for feasibility of further analysis and machine learning algorithms.

The code we used allowed us to explicitly define the mapping between categories and numerical values:

`{'Non-Diabetic': 0, 'Pre-Diabetic': 1, 'Diabetic': 2}`

With this code numerical values were automatically assigned to each category:

`{0: 'Female', 1: 'Male'}`

7. Remove unnecessary columns - the "ID" and "No-Patient" columns.

df.drop(["ID", "No-Patient"], axis=1, inplace=True)

Age	BP	Skin	Insulin	BFS	CLASS	CLASS_Category
0	39	130	39	350	Non-Diabetic	0
1	40	120	43	430	Non-Diabetic	0
2	38	135	43	430	Non-Diabetic	0
3	37	133	43	430	Non-Diabetic	0
4	35	120	43	300	Non-Diabetic	0
...
988	31	120	39	350	Non-Diabetic	0
989	31	120	39	350	Non-Diabetic	0
990	33	130	43	350	Non-Diabetic	0
991	33	130	43	350	Non-Diabetic	0
992	33	130	43	350	Non-Diabetic	0
993	33	130	43	350	Non-Diabetic	0
994	33	130	43	350	Non-Diabetic	0
995	33	130	43	350	Non-Diabetic	0
996	33	130	43	350	Non-Diabetic	0
997	33	130	43	350	Non-Diabetic	0
998	33	130	43	350	Non-Diabetic	0
999	33	130	43	350	Non-Diabetic	0

1000 rows * 14 columns

1. Overview of dataset

data.info()



```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000 entries, 0 to 999
Data columns (total 14 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   ID          1000 non-null    int64  
 1   No_Pation   1000 non-null    int64  
 2   Gender       1000 non-null    object  
 3   AGE          1000 non-null    int64  
 4   Urea         1000 non-null    float64 
 5   Cr           1000 non-null    int64  
 6   HbAlc        1000 non-null    float64 
 7   Chol         1000 non-null    float64 
 8   TG           1000 non-null    float64 
 9   HDL          1000 non-null    float64 
 10  LDL          1000 non-null    float64 
 11  VLDL         1000 non-null    float64 
 12  BMI          1000 non-null    float64 
 13  CLASS         1000 non-null    object  
dtypes: float64(8), int64(4), object(2)
```



- 1000 entries, 14 columns
- No null values
- 2 objects, 4 int, and 8 float data types

2. We were particularly interested to see **how many categories** exist in “**CLASS**” & “**Gender**” columns as ‘**object**’ data types, and **number of patients** in each category.

“**CLASS**” column = “**target variable**” classifies patients in 3 categories:
Non-diabetic (N), **Pre-diabetic (P)** and **Diabetic(Y)**

“**Gender**” column classifies patients in 2 categories:
Female (F) and **Male (M)**

`value_counts()`



CLASS	
Y	840
N	102
P	53
Y	4
N	1

Name: count, dtype: int64



Gender	
M	565
F	434
f	1

Name: count, dtype: int64

Redundant categories
were found

3. Redundant categories were found, due to typing error!

Note: We see here that we have an unbalanced data set N=103, P=53, and Y=844

So we replaced:

Any spaces (" ") with an empty string ("")

str.strip()



CLASS	
Y	844
N	103
P	53
Name: count, dtype: int64	

Lowercase letters with uppercase ones

str.upper()



Gender	
M	565
F	435
Name: count, dtype: int64	

4. For easier readability of the code we replaced the values in “CLASS” and “Gender as following:
“N” to “Non-Diabetic”
“P” to “Pre-Diabetic”
“Y” to “Diabetic”
“F” to “Female”
“M” to “Male”

.map()

The diagram illustrates the use of the `.map()` function in Python's pandas library. It shows two red-bordered boxes representing data frames. The top box contains the 'CLASS' column with three categories: 'Diabetic' (count 844), 'Non-Diabetic' (count 103), and 'Pre-Diabetic' (count 53). The bottom box contains the 'Gender' column with two categories: 'Male' (count 565) and 'Female' (count 435). Two purple arrows point from the text "For easier readability of the code we replaced the values in ‘CLASS’ and ‘Gender as following:” to the respective data frames. A green box labeled ".map()" is positioned between the two boxes, indicating the transformation process.

CLASS	
Diabetic	844
Non-Diabetic	103
Pre-Diabetic	53

Name: count, dtype: int64

Gender	
Male	565
Female	435

Name: count, dtype: int64

5. We then converted the “CLASS” and “Gender” columns from ‘object’ to ‘category’ data type.

.astype()



ID	int64
No_Pation	int64
Gender	category
AGE	int64
Urea	float64
Cr	int64
HbAlc	float64
Chol	float64
TG	float64
HDL	float64
LDL	float64
VLDL	float64
BMI	float64
CLASS	category
dtype: object	

6. Two numerical columns as “CLASS_Category” from the “CLASS” column and “Gender_Category” from “Gender” column were added for feasibility of further analysis and machine learning algorithms.

```
CLASS      category  
dtype: object
```

6. Two numerical columns as “CLASS_Category” from the “CLASS” column and “Gender_Category” from “Gender” column were added for feasibility of further analysis and machine learning algorithms.

The code we used allowed us to explicitly define the mapping between categories and numerical values:

```
{'Non-Diabetic': 0, 'Pre-Diabetic': 1, 'Diabetic': 2}
```

.map().astype(int)



CLASS	CLASS_Category
Non-Diabetic	0

With this code numerical values were automatically assigned to each category:

```
{0: 'Female', 1: 'Male'}
```

.cat.codes



Gender_Category
0
1
0
0
1

7. Remove unnecessary columns - the “ID” and “No-Pation” columns.

.drop()

	Gender	AGE	Urea	Cr	HbA1c	Chol	TG	HDL	LDL	VLDL	BMI	CLASS	CLASS_Category	Gender_Category
0	Female	50	4.7	46	4.9	4.2	0.9	2.4	1.4	0.5	24.0	Non-Diabetic	0	0
1	Male	26	4.5	62	4.9	3.7	1.4	1.1	2.1	0.6	23.0	Non-Diabetic	0	1
2	Female	50	4.7	46	4.9	4.2	0.9	2.4	1.4	0.5	24.0	Non-Diabetic	0	0
3	Female	50	4.7	46	4.9	4.2	0.9	2.4	1.4	0.5	24.0	Non-Diabetic	0	0
4	Male	33	7.1	46	4.9	4.9	1.0	0.8	2.0	0.4	21.0	Non-Diabetic	0	1
...
995	Male	71	11.0	97	7.0	7.5	1.7	1.2	1.8	0.6	30.0	Diabetic	2	1
996	Male	31	3.0	60	12.3	4.1	2.2	0.7	2.4	15.4	37.2	Diabetic	2	1
997	Male	30	7.1	81	6.7	4.1	1.1	1.2	2.4	8.1	27.4	Diabetic	2	1
998	Male	38	5.8	59	6.7	5.3	2.0	1.6	2.9	14.0	40.5	Diabetic	2	1
999	Male	54	5.0	67	6.9	3.8	1.7	1.1	3.0	0.7	33.0	Diabetic	2	1

1000 rows × 14 columns

	Gender	AGE	Urea	Cr	HbA1c	Chol	TG	HDL	LDL	VLDL	BMI	CLASS	CLASS_Category	Gender_Category
0	Female	50	4.7	46	4.9	4.2	0.9	2.4	1.4	0.5	24.0	Non-Diabetic	0	0
1	Male	26	4.5	62	4.9	3.7	1.4	1.1	2.1	0.6	23.0	Non-Diabetic	0	1
2	Female	50	4.7	46	4.9	4.2	0.9	2.4	1.4	0.5	24.0	Non-Diabetic	0	0
3	Female	50	4.7	46	4.9	4.2	0.9	2.4	1.4	0.5	24.0	Non-Diabetic	0	0
4	Male	33	7.1	46	4.9	4.9	1.0	0.8	2.0	0.4	21.0	Non-Diabetic	0	1
...
995	Male	71	11.0	97	7.0	7.5	1.7	1.2	1.8	0.6	30.0	Diabetic	2	1
996	Male	31	3.0	60	12.3	4.1	2.2	0.7	2.4	15.4	37.2	Diabetic	2	1
997	Male	30	7.1	81	6.7	4.1	1.1	1.2	2.4	8.1	27.4	Diabetic	2	1
998	Male	38	5.8	59	6.7	5.3	2.0	1.6	2.9	14.0	40.5	Diabetic	2	1
999	Male	54	5.0	67	6.9	3.8	1.7	1.1	3.0	0.7	33.0	Diabetic	2	1

1000 rows × 14 columns

Diabetes Prediction Model

Using Machine Learning Regression Algorithms and Neural Networks



Index

Presentation by Nastaran Ashoori & Daniella Vuletić
CodeOp DSPT10 - July 2023



1- Introduction

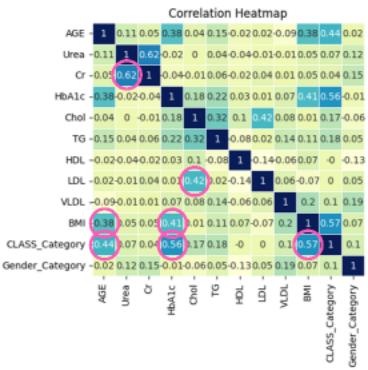
**2- Data Pre-
Processing & Feature
Engineering**

**3- Exploratory
Data Analysis
(EDA)**

**4- Machine
Learning (ML)**

5- Conclusion

3- EXPLORATORY DATA ANALYSIS (EDA)

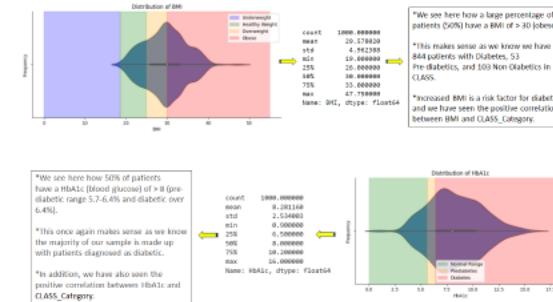


We used a Pearson's correlation with a heatmap to look at correlations and strength of correlations between our variables.

What we see here, and relevance to analysis:

- BMI and CLASS_Category showed a moderate positive correlation (**0.57**) – this is not unexpected due to the risk of increased BMI and diabetes.
- HbA1c and CLASS_Category showed a moderate positive correlation (**0.56**) – this too is not unexpected because of high blood sugar characteristic in diabetes patients.
- Age and CLASS_Category showed a moderate positive correlation (**0.44**).

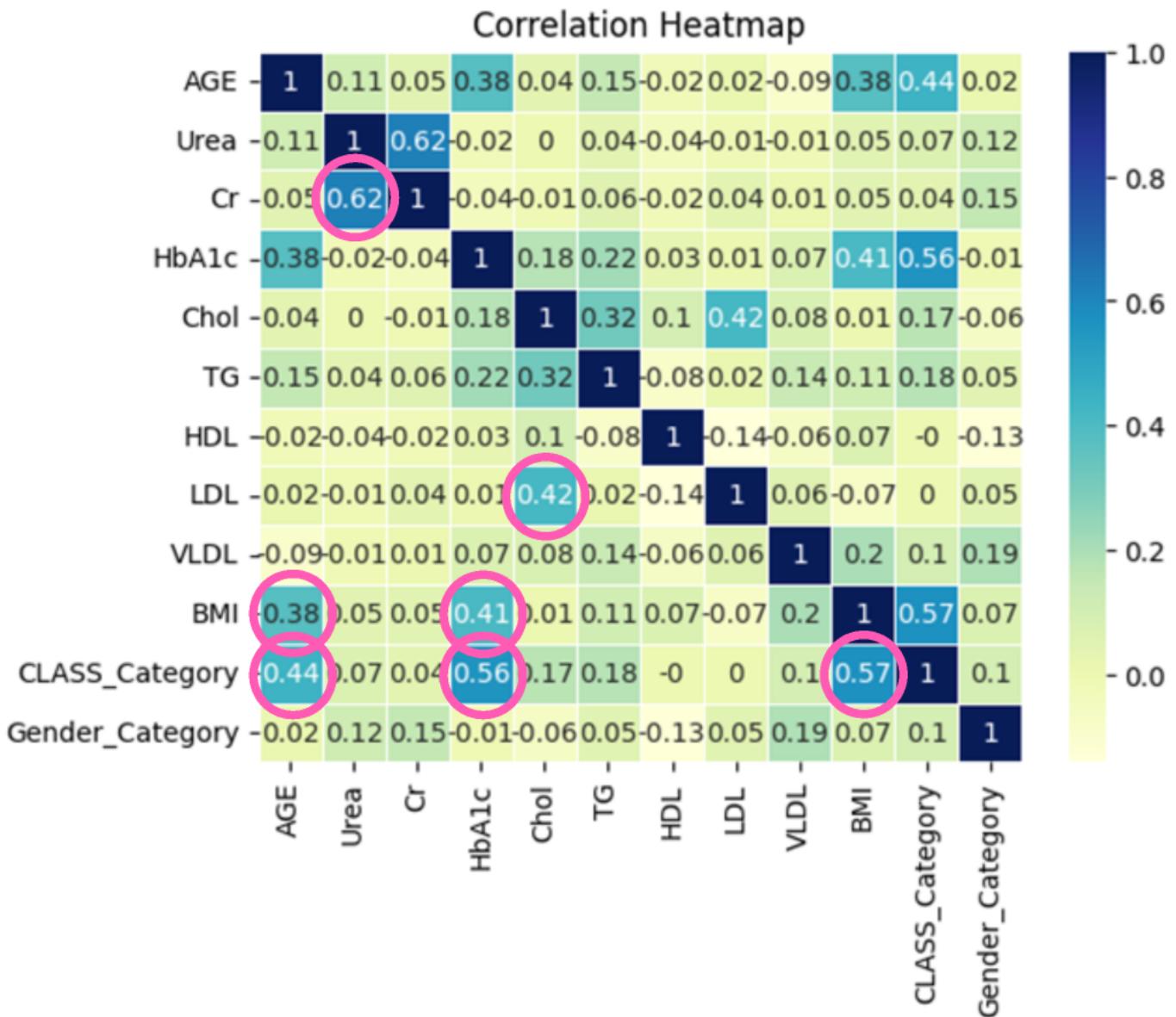
Let's take a look at the distribution and normal ranges of BMI and HbA1c in a little more detail:



Conclusions from EDA...

Interactive Plots With Diabetes Dataset





We used a Pearson's correlation with a heatmap to look at correlations and strength of correlations between our variables.

What we see here, and relevance to analysis:

- **BMI** and **CLASS_Category** showed a moderate positive correlation (**0.57**) – this is not unexpected due to the risk of increased BMI and diabetes.
- **HbA1c** and **CLASS_Category** showed a moderate positive correlation (**0.56**) – this too is not unexpected because of high blood sugar characteristic in diabetes patients.
- **Age** and **CLASS_Category** showed a moderate positive correlation (**0.44**).

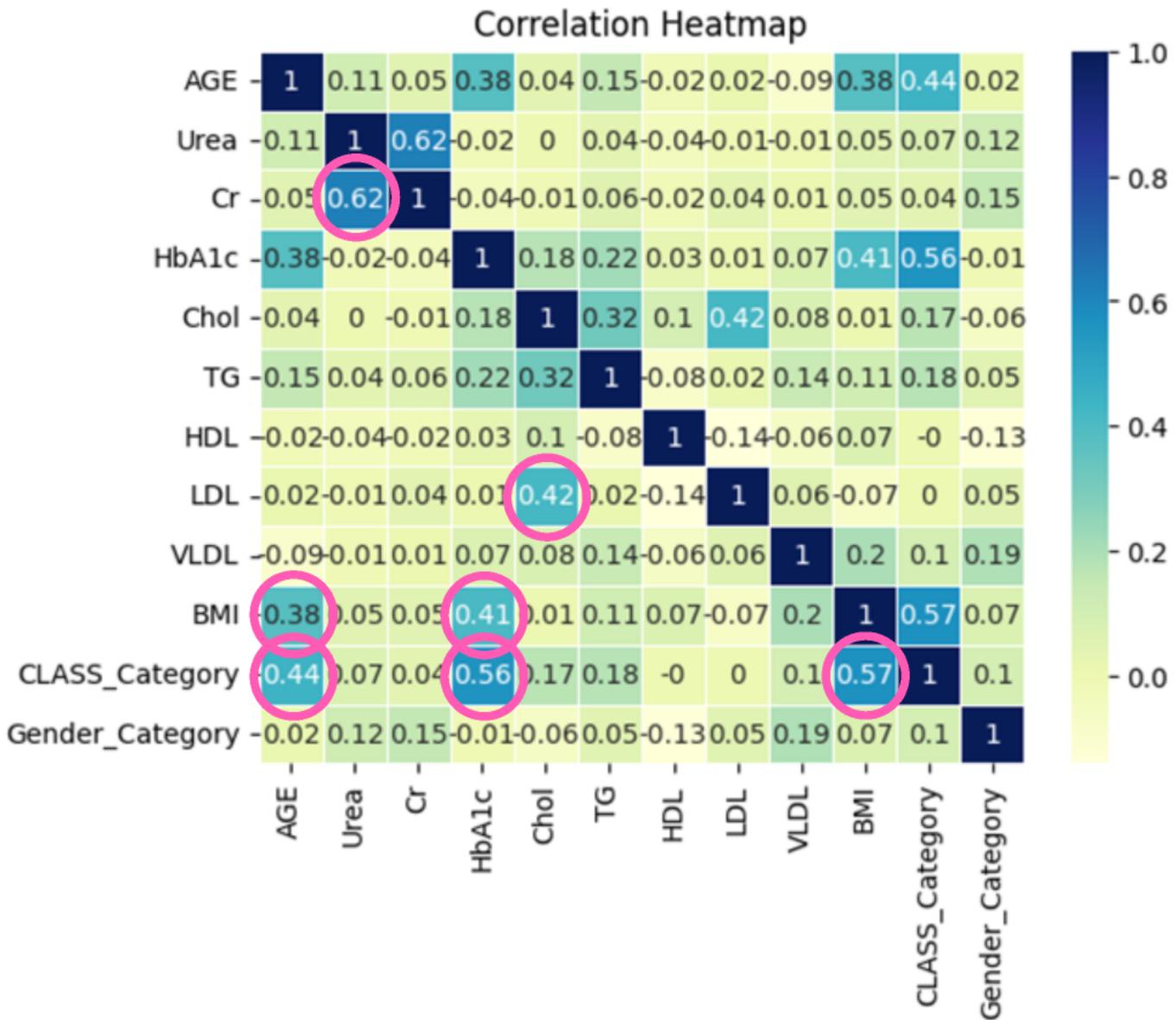
- 0.4

- 0.2

- 0.0

- **BMI** and **CLASS_Category** showed a moderate positive correlation (**0.57**) – this is not unexpected due to the risk of increased BMI and diabetes.
- **HbA1c** and **CLASS_Category** showed a moderate positive correlation (**0.56**) – this too is not unexpected because of high blood sugar characteristic in diabetes patients.
- **Age** and **CLASS_Category** showed a moderate positive correlation (**0.44**).

	AGE	Urea	Cr	HbA1c	Chol	TG	HDL	LDL	VLDL	BMI	CLASS_Category	Gender_Category
AGE	-	-	-	-	-	-	-	-	-	-	-	-
Urea	-0.02	0.12	0.15	-0.01	-0.06	0.05	-0.13	0.05	0.19	0.07	0.1	-0.02
Cr	-	-	-	-	-	-	-	-	-	-	-	-
HbA1c	-	-	-	-	-	-	-	-	-	-	-	-0.01
Chol	-	-	-	-	-	-	-	-	-	-	-	-
TG	-	-	-	-	-	-	-	-	-	-	-	-
HDL	-	-	-	-	-	-	-	-	-	-	-	-
LDL	-	-	-	-	-	-	-	-	-	-	-	-
VLDL	-	-	-	-	-	-	-	-	-	-	-	-
BMI	-	-	-	-	-	-	-	-	-	-	-	-
CLASS_Category	-0.38	0.05	0.05	0.41	0.01	0.11	0.07	-0.07	0.2	1	0.57	0.07
Gender_Category	-0.44	0.07	0.04	0.56	0.17	0.18	-0	0	0.1	0.57	1	0.1

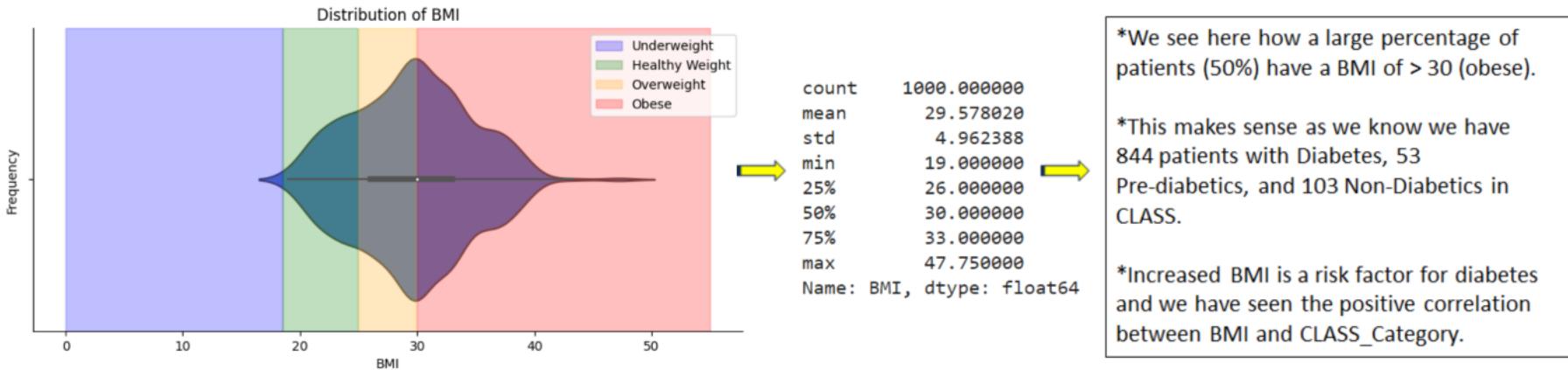


We used a Pearson's correlation with a heatmap to look at correlations and strength of correlations between our variables.

What we see here, and relevance to analysis:

- **BMI** and **CLASS_Category** showed a moderate positive correlation (**0.57**) – this is not unexpected due to the risk of increased BMI and diabetes.
- **HbA1c** and **CLASS_Category** showed a moderate positive correlation (**0.56**) – this too is not unexpected because of high blood sugar characteristic in diabetes patients.
- **Age** and **CLASS_Category** showed a moderate positive correlation (**0.44**).

Let's take a look at the distribution and normal ranges of BMI and HbA1c in a little more detail:

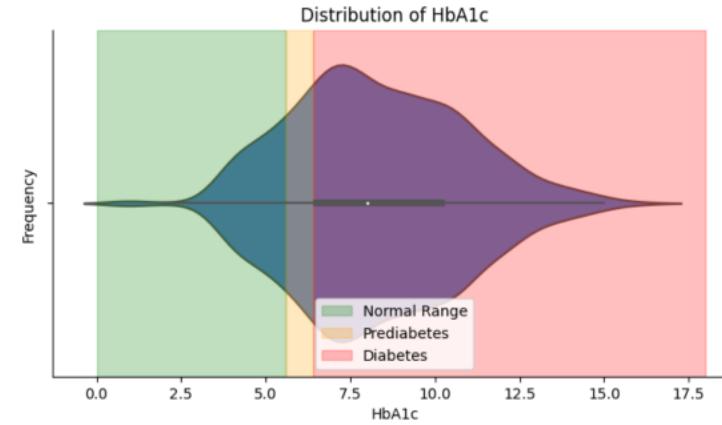


*We see here how 50% of patients have a HbA1c (blood glucose) of > 8 (pre-diabetic range 5.7-6.4% and diabetic over 6.4%).

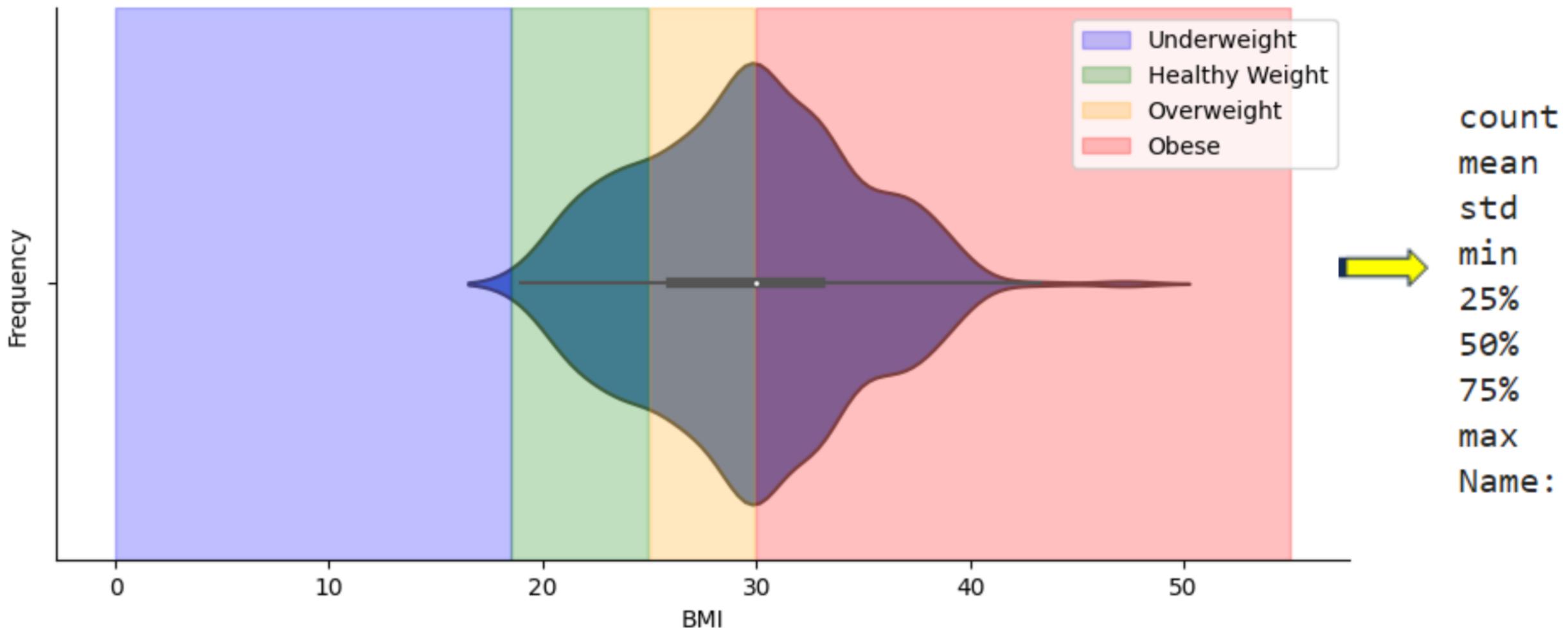
*This once again makes sense as we know the majority of our sample is made up with patients diagnosed as diabetic.

*In addition, we have also seen the positive correlation between HbA1c and CLASS_Category.

count 1000.000000
mean 8.281160
std 2.534003
min 0.900000
25% 6.500000
50% 8.000000
75% 10.200000
max 16.000000
Name: HbA1c, dtype: float64



Distribution of BMI



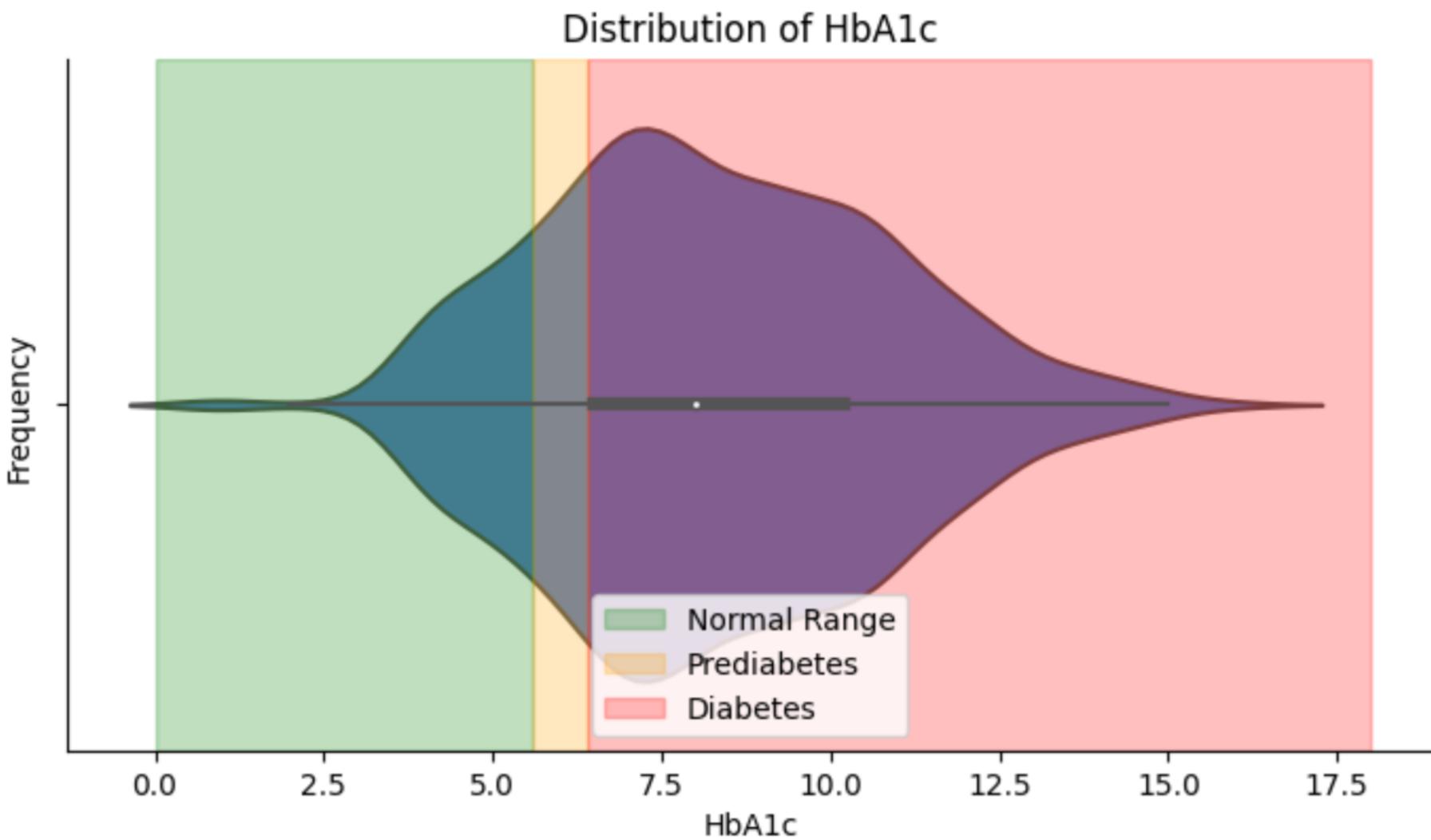
```
count      1000.000000
mean       29.578020
std        4.962388
min        19.000000
25%       26.000000
50%       30.000000
75%       33.000000
max        47.750000
Name: BMI, dtype: float64
```

*We see here how a large percentage of patients (50%) have a BMI of > 30 (obese).

*This makes sense as we know we have 844 patients with Diabetes, 53 Pre-diabetics, and 103 Non-Diabetics in CLASS.

*Increased BMI is a risk factor for diabetes and we have seen the positive correlation between BMI and CLASS_Category.

```
000000  
281160  
534003  
900000  
500000  
000000  
200000  
000000  
type: float64
```



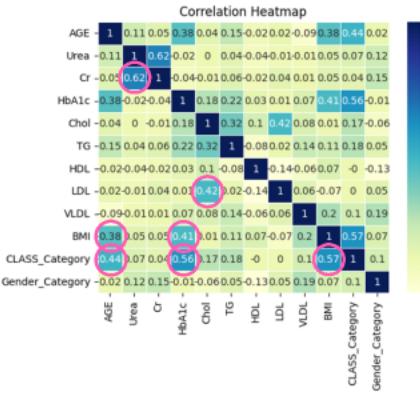
*We see here how 50% of patients have a HbA1c (blood glucose) of > 8 (prediabetic range 5.7-6.4% and diabetic over 6.4%).

*This once again makes sense as we know the majority of our sample is made up with patients diagnosed as diabetic.

*In addition, we have also seen the positive correlation between HbA1c and CLASS_Category.

```
count      1000.000000
mean       8.281160
std        2.534003
min        0.900000
25%        6.500000
50%        8.000000
75%        10.200000
max       16.000000
Name: HbA1c, dtype: float64
```

3- EXPLORATORY DATA ANALYSIS (EDA)

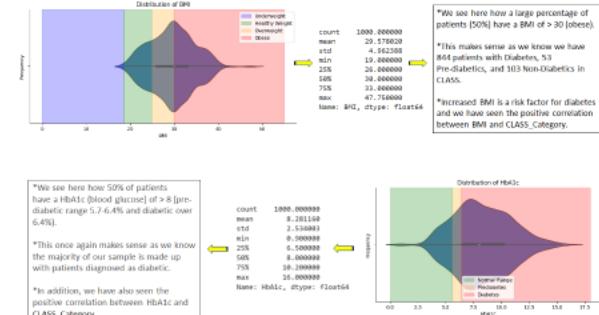


We used a Pearson's correlation with a heatmap to look at correlations and strength of correlations between our variables.

What we see here, and relevance to analysis:

- **BMI** and **CLASS_Category** showed a moderate positive correlation (**0.57**) – this is not unexpected due to the risk of increased BMI and diabetes.
- **HbA1c** and **CLASS_Category** showed a moderate positive correlation (**0.56**) – this too is not unexpected because of high blood sugar characteristic in diabetes patients.
- **Age** and **CLASS_Category** showed a moderate positive correlation (**0.44**).

Let's take a look at the distribution and normal ranges of BMI and HbA1c in a little more detail:



*We see here how a large percentage of patients (20%) have a BMI of > 30 (obese).

*This makes sense as we know we have 864 patients with Diabetes, 53 Pre-diabetics, and 103 Non-Diabetics In CLASS.

*Increased BMI is a risk factor for diabetes and we have seen the positive correlation between BMI and CLASS_Category.

Conclusions from EDA...

Interactive Plots With Diabetes Dataset

Scatterplot



Histogram

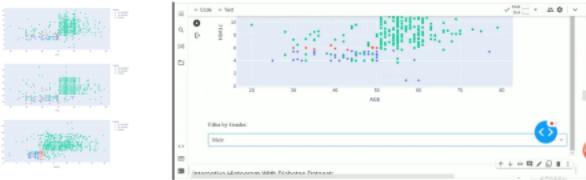


Catplot

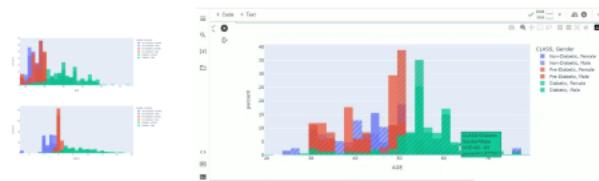


Interactive Plots With Diabetes Dataset

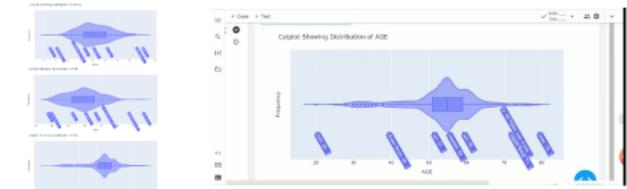
Scatterplot



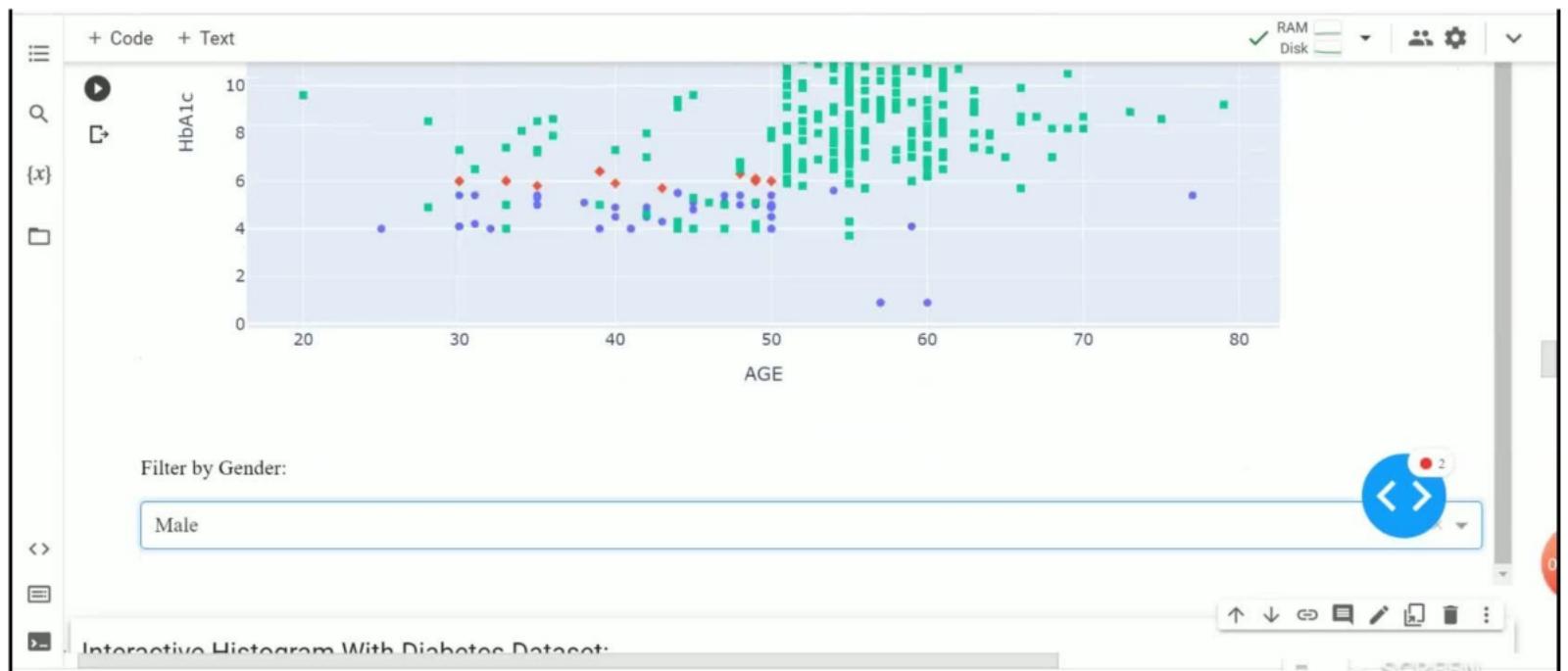
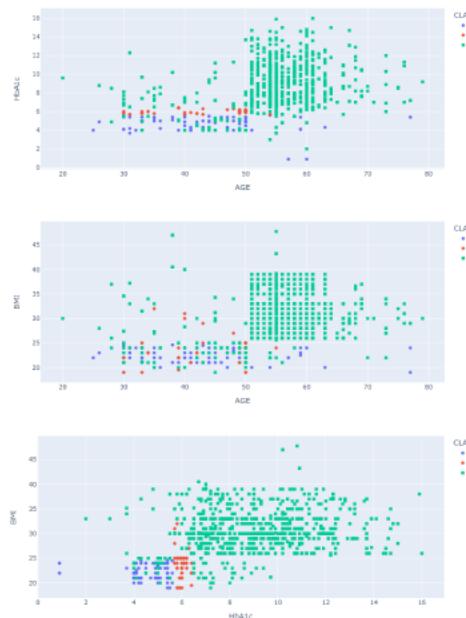
Histogram



Catplot

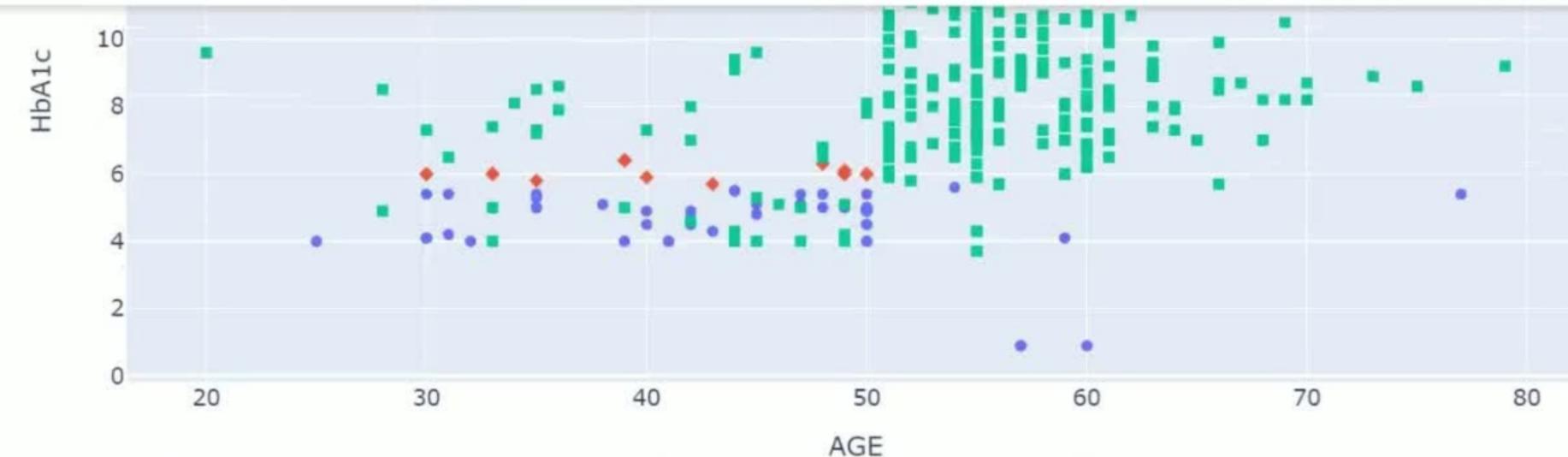


Scatterplot



+ Code + Text

RAM
Disk



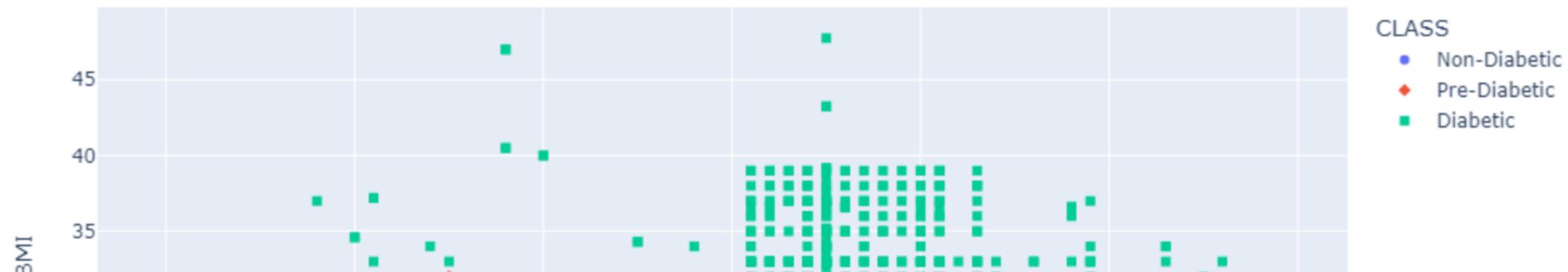
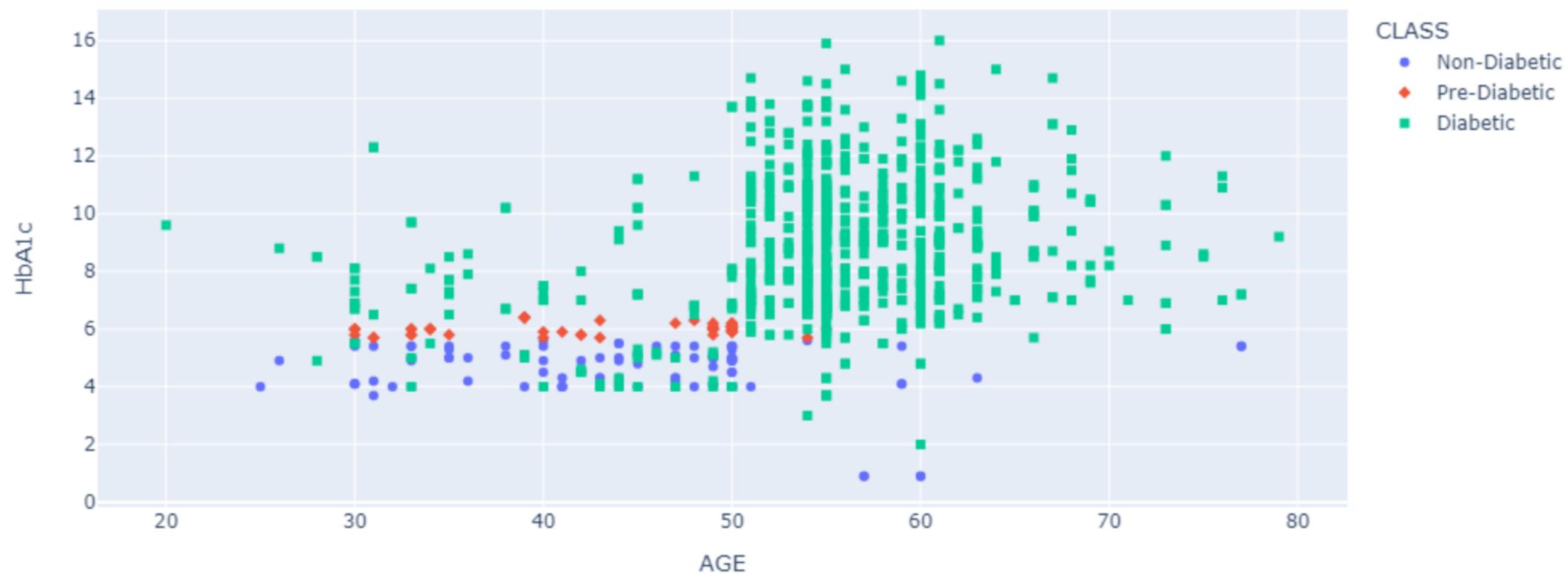
Filter by Gender:

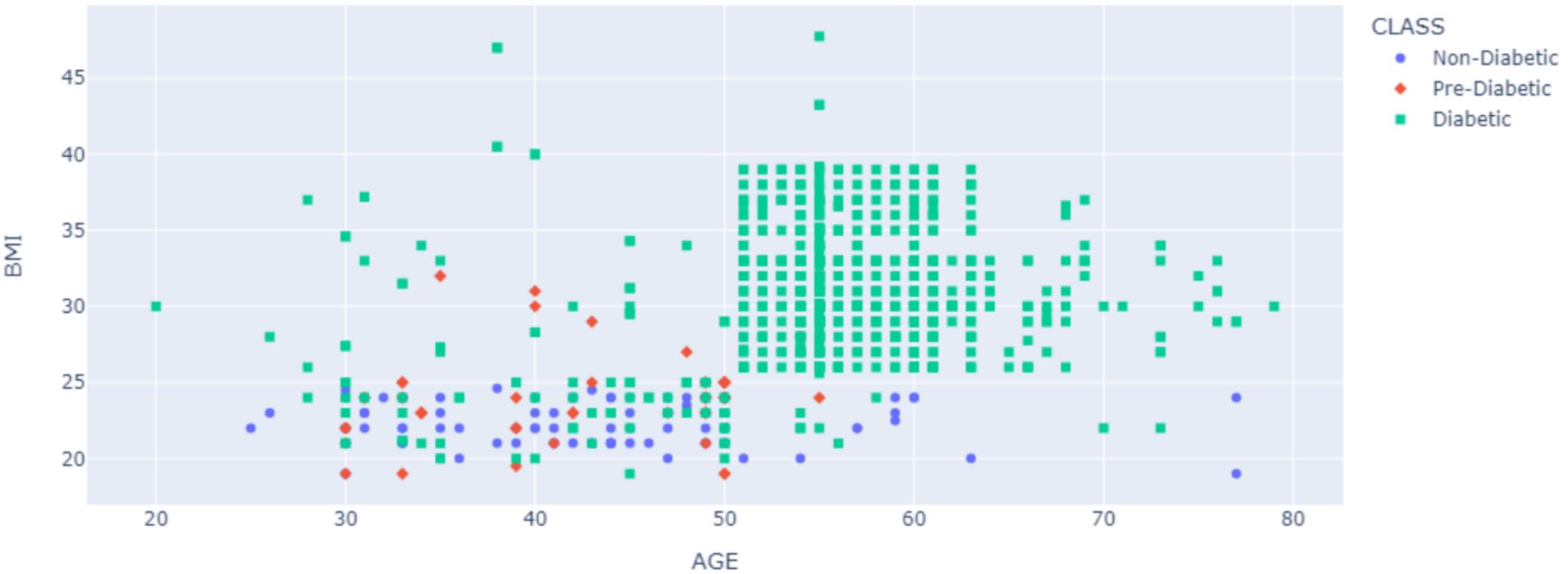
Male

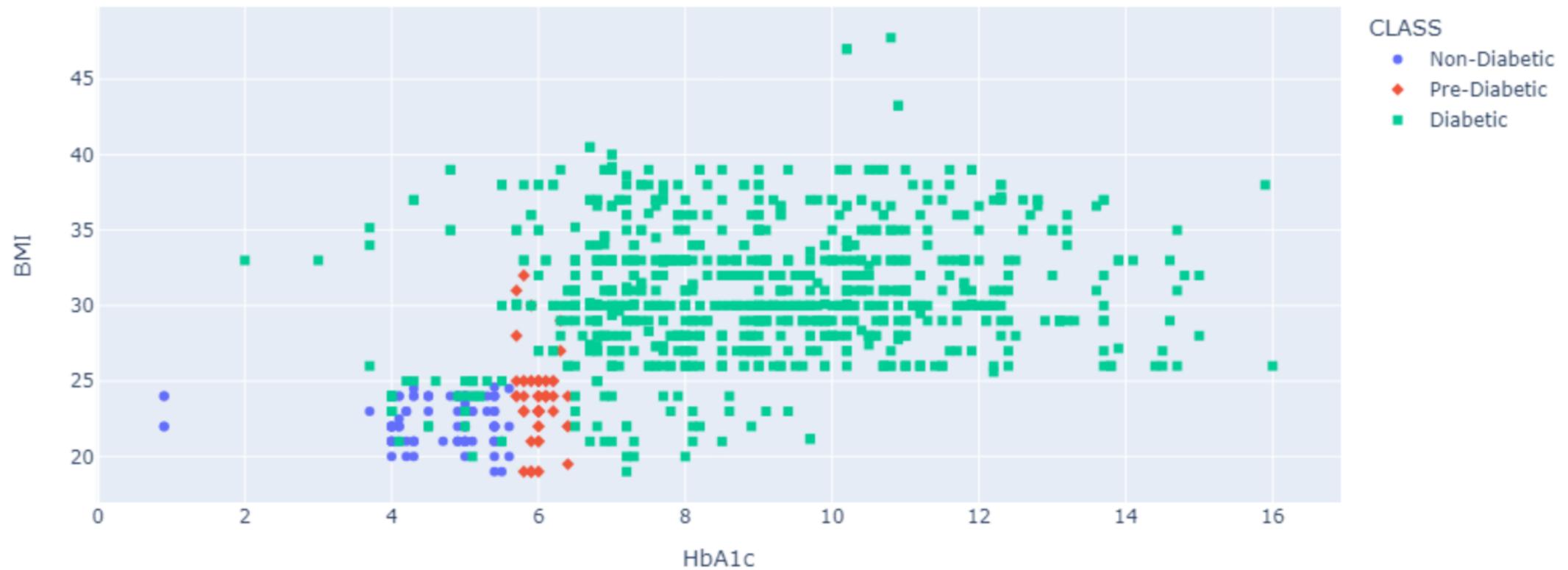
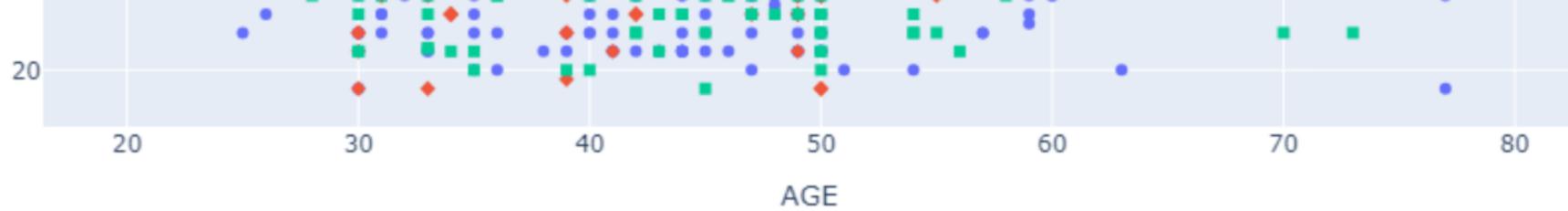


Interactive Histogram With Diabetes Dataset:



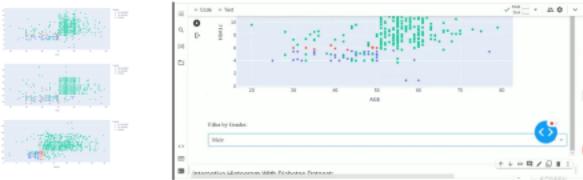




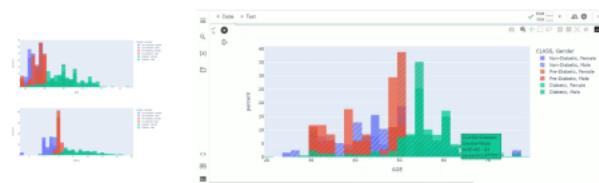


Interactive Plots With Diabetes Dataset

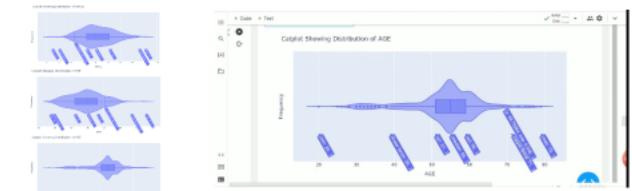
Scatterplot



Histogram

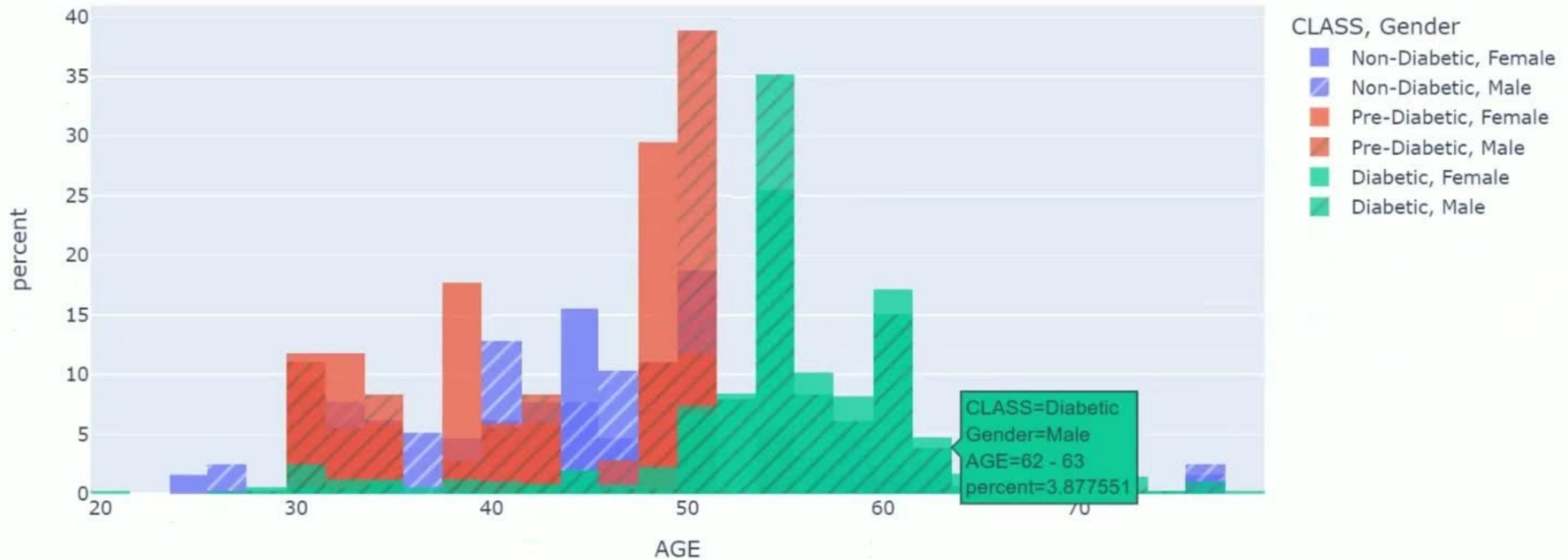


Catplot

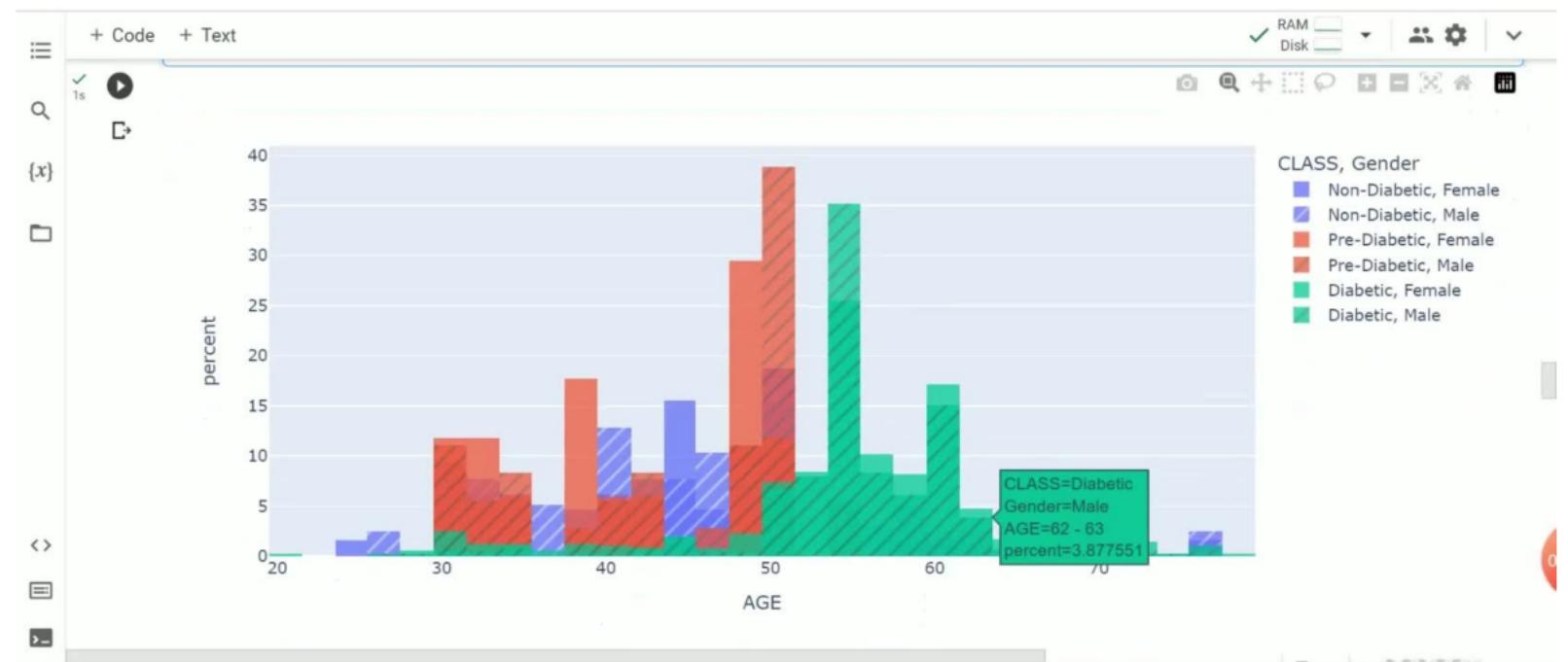
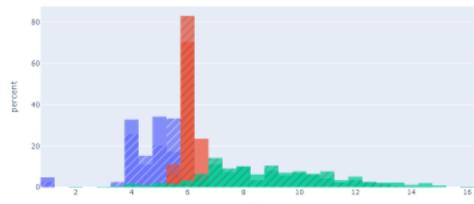
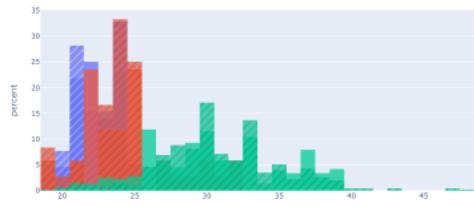


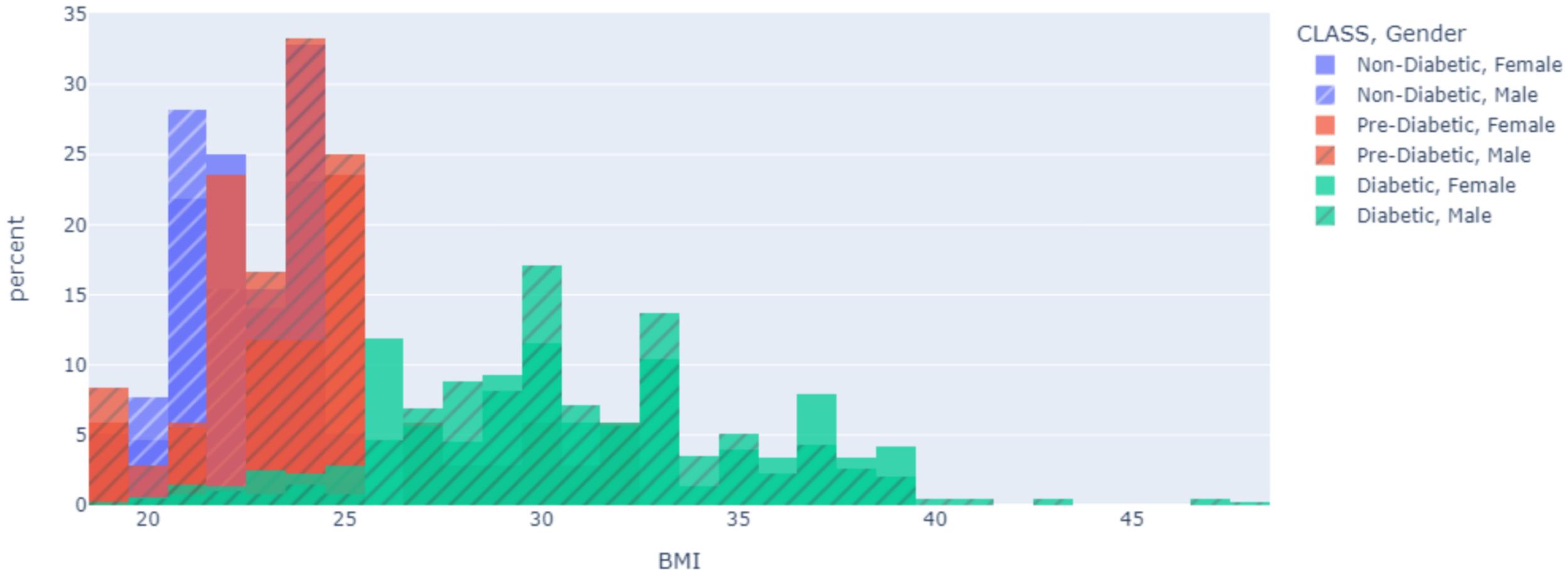
+ Code + Text

✓ RAM Disk    

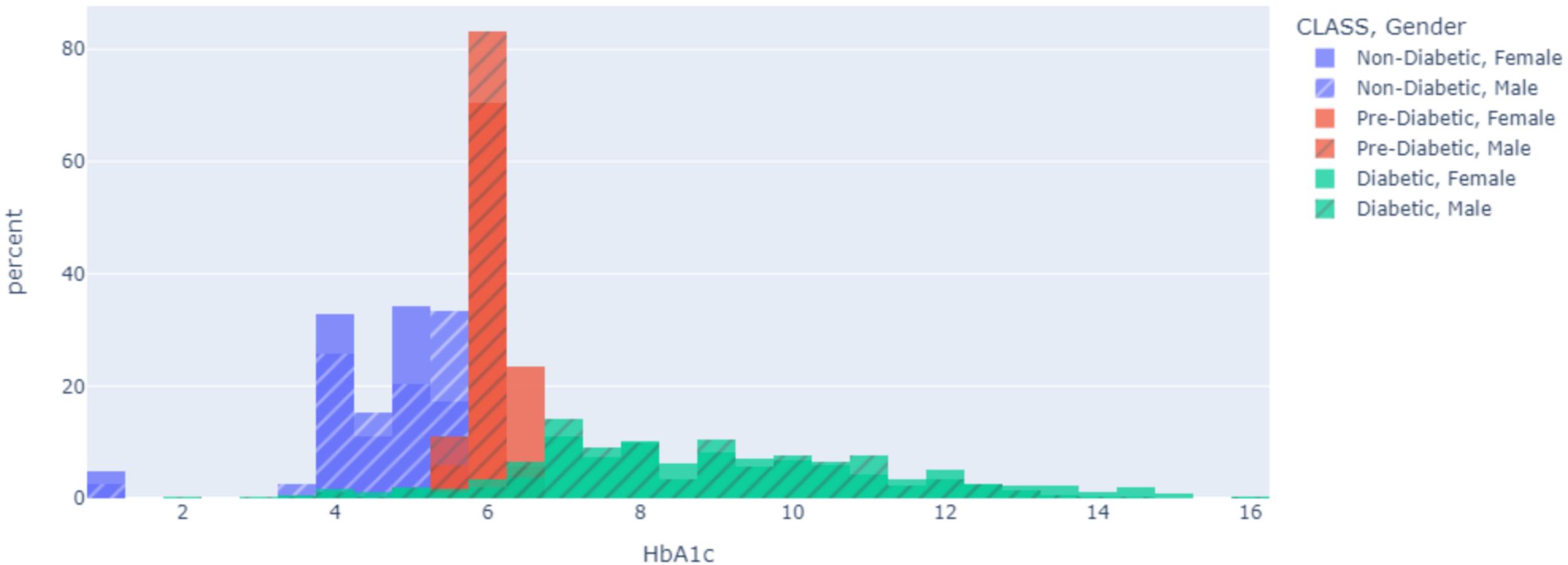


Histogram





BMI

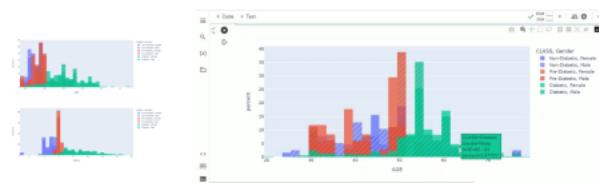


Interactive Plots With Diabetes Dataset

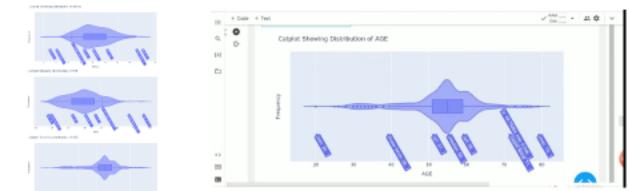
Scatterplot



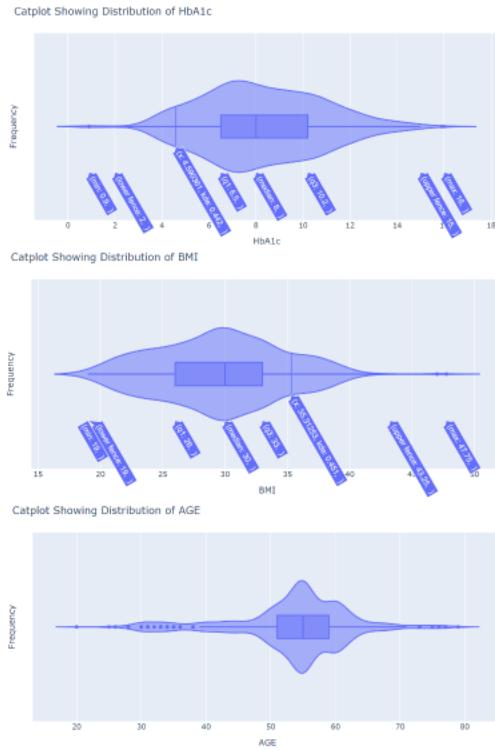
Histogram



Catplot



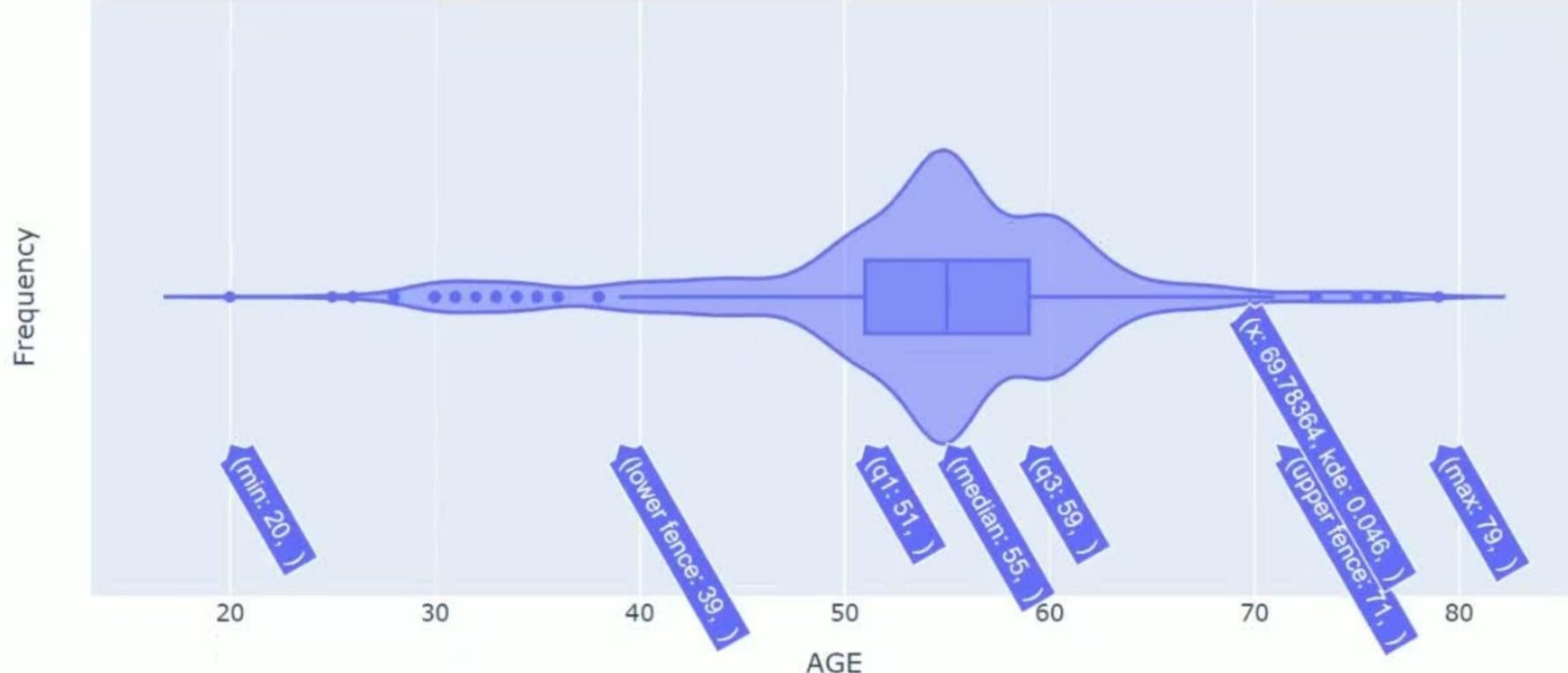
Catplot



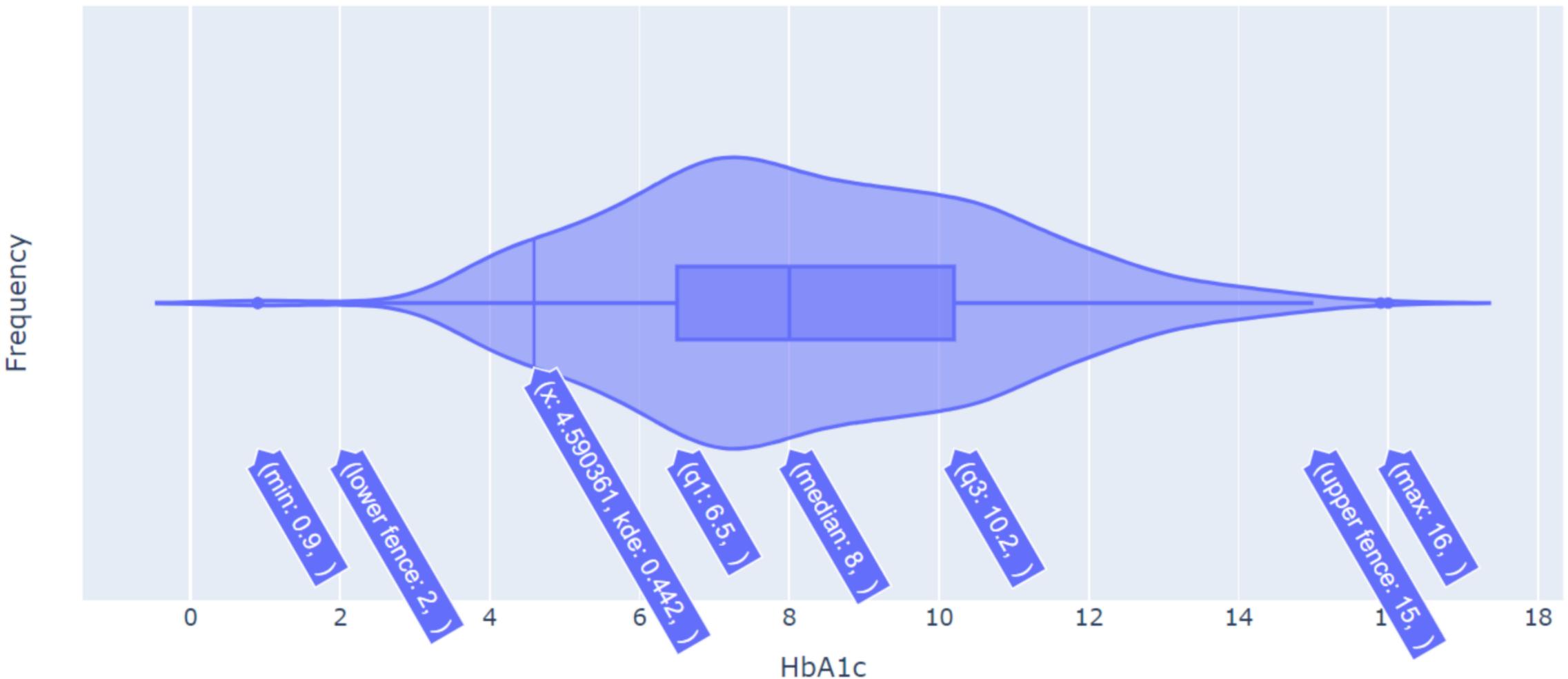
+ Code + Text

RAM
Disk

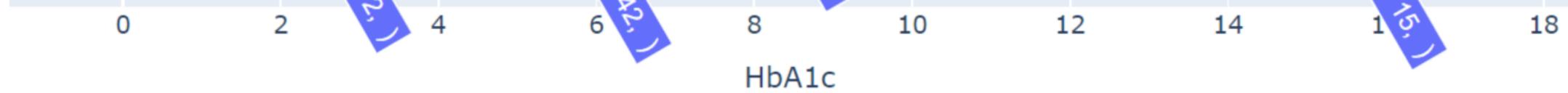
Catplot Showing Distribution of AGE



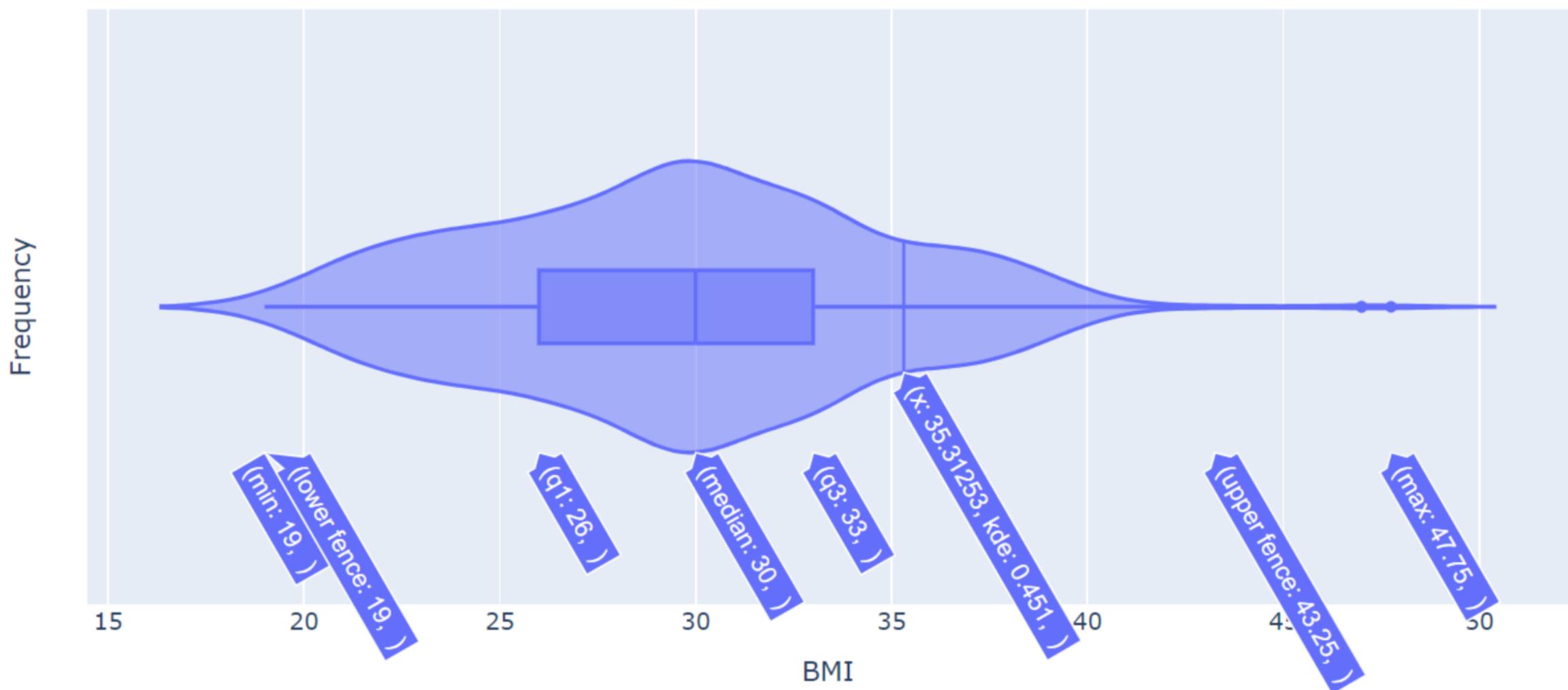
Catplot Showing Distribution of HbA1c



Catplot Showing Distribution of BMI



Catplot Showing Distribution of BMI



Catplot Showing Distribution of AGE

15

20

25

30

35

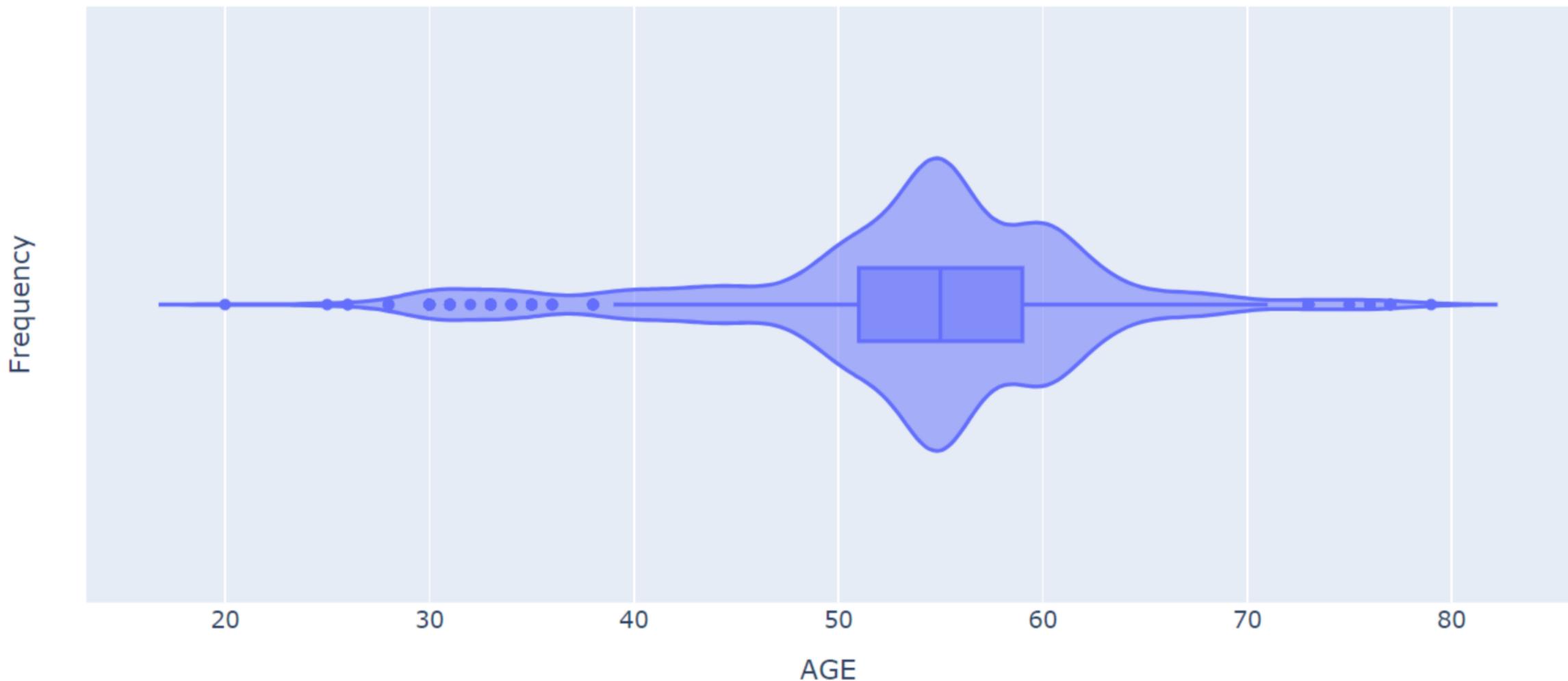
40

45

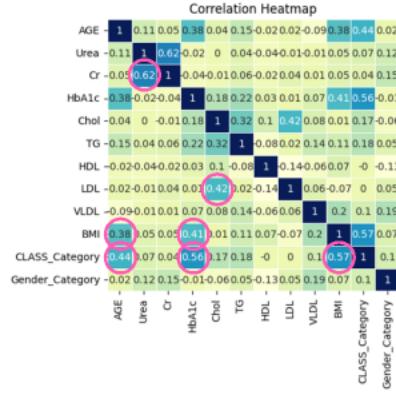
50

BMI

Catplot Showing Distribution of AGE



3- EXPLORATORY DATA ANALYSIS (EDA)

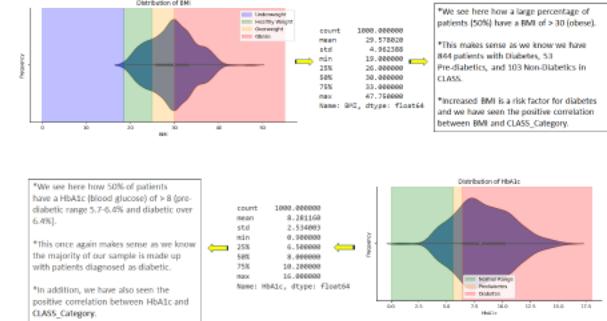


We used a Pearson's correlation with a heatmap to look at correlations and strength of correlations between our variables.

What we see here, and relevance to analysis:

- BMI and CLASS_Category** showed a moderate positive correlation (**0.57**) – this is not unexpected due to the risk of increased BMI and diabetes.
- HbA1c and CLASS_Category** showed a moderate positive correlation (**0.56**) – this too is not unexpected because of high blood sugar characteristic in diabetes patients.
- Age and CLASS_Category** showed a moderate positive correlation (**0.44**).

Let's take a look at the distribution and normal ranges of BMI and HbA1c in a little more detail:



*We see here how 50% of patients have a HbA1c (Blood glucose) of > 8 (pre-diabetic range 5.7-6.4% and diabetic over 6.4%).

*This once again makes sense as we know the majority of our sample is made up with patients diagnosed as diabetic.

*In addition, we have also seen the positive correlation between HbA1c and CLASS_Category.

*We see here how a large percentage of patients (50%) have a BMI of > 30 (obese).

*This makes sense as we know we have 844 patients with Diabetes, 53 Pre-diabetics, and 103 Non-Diabetics in CLASS.

*Increased BMI is a risk factor for diabetes and we have seen the positive correlation between BMI and CLASS_Category.

Conclusions from EDA...

Interactive Plots With Diabetes Dataset

Scatterplot



Histogram



Catplot



Conclusions from EDA...

Conclusions from EDA...

- * Results using Pearson's correlation suggest that **BMI**, **HbA1c**, and **age** are moderately positively correlated with CLASS_Category i.e. **an increase in these variables correlates with an increased risk of diabetes**. This was not unexpected as these variables are all known risk factors for diabetes.

Conclusions from EDA...

* Results using Pearson's correlation suggest that **BMI**, **HbA1c**, and **age** are moderately positively correlated with CLASS_Category i.e. **an increase in these variables correlates with an increased risk of diabetes**. This was not unexpected as these variables are all known risk factors for diabetes.

* Interestingly no correlation was seen between gender and CLASS_Category – suggesting that **males and females are equally as likely to develop diabetes** (literature generally shows type 2 diabetes to be more prevalent in men).

Conclusions from EDA...

- * Results using Pearson's correlation suggest that **BMI**, **HbA1c**, and **age** are moderately positively correlated with CLASS_Category i.e. **an increase in these variables correlates with an increased risk of diabetes**. This was not unexpected as these variables are all known risk factors for diabetes.
- * Interestingly no correlation was seen between gender and CLASS_Category – suggesting that **males and females are equally as likely to develop diabetes** (literature generally shows type 2 diabetes to be more prevalent in men).
- * The dataset has a skewed distribution within the **BMI**, **HbA1c** and **age** variables, with **a large percentage of patients falling into the 'obese', 'diabetic glucose', and 'older' (median 55 years) range**, respectively. This, however, was not unexpected as most of the sample were diagnosed diabetes patients.

Diabetes Prediction Model

Using Machine Learning Regression Algorithms and Neural Networks



Index

Presentation by Nastaran Ashoori & Daniella Vuletić
CodeOp DSPT10 - July 2023



1- Introduction

**2- Data Pre-
Processing & Feature
Engineering**

**3- Exploratory
Data Analysis
(EDA)**

**4- Machine
Learning (ML)**

5- Conclusion

4- MACHINE LEARNING & DEEP LEARNING

4.1 Logistic Regression - Classification

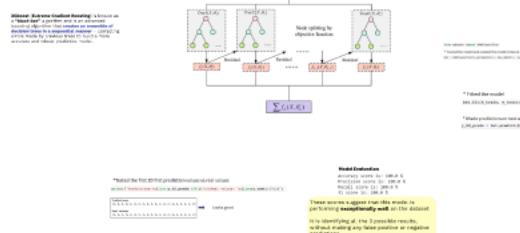


4.1 Logistic Regression - Classification

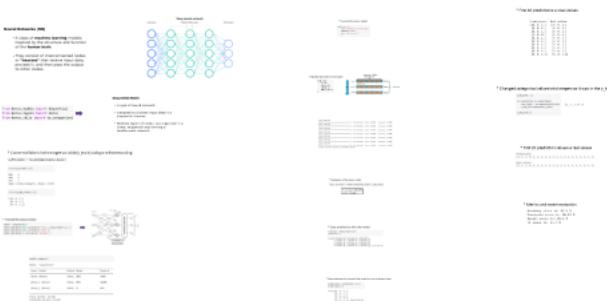
4.2 - Normalization of the data - MinMaxScaler



4.3 XGboost



4.4 NEURAL NETWORKS

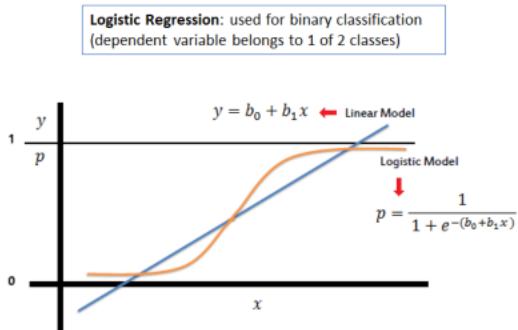


Machine learning (ML)

- is a subfield of **artificial intelligence** that focuses on developing and training algorithms to **learn patterns** and **make decisions** from data without being explicitly trained to do so.
- The aim is to ultimately develop models that can make **accurate predictions or decisions on new unseen data.**

4.1 Logistic Regression - Classification

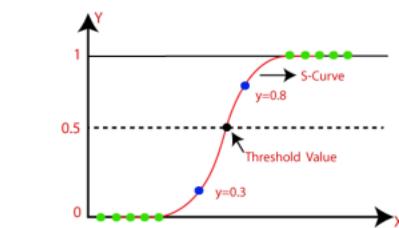
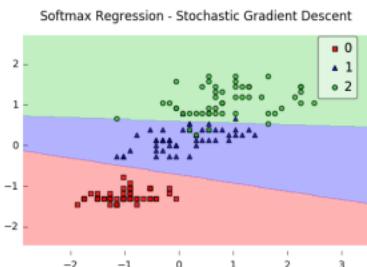
Logistic Regression
• A supervised ML classification technique



Our dataset had 3 possible classes

Multinomial Logistic Regression:
- An extension of logistic regression used to handle **multiclass** problems.

Multinomial Logistic Regression: used for multiclass classification (more than 2 or K classes)



* First, we defined 'X' as the training dataset and 'y' as the target variable

```
x = data.drop(['CLASS_Category', 'CLASS', 'Gender'], axis=1)
y = data['CLASS_Category']
```

* Split the dataset into 80/20 training/testing split with a random state of 50

```
from sklearn.model_selection import train_test_split

print(data.shape)
print(X_train.shape)
print(X_test.shape)
print(y_train.shape)
print(y_test.shape)

(1000, 14)
(800, 13)
(200, 13)
(800, )
(200, )
```

* Trained the model

```
from sklearn.linear_model import LogisticRegression

lr = LogisticRegression(multi_class='multinomial', solver='lbfgs', max_iter=2600)
lr.fit(X_train, y_train)

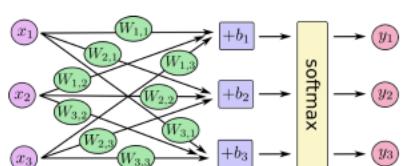
* Predictions on the test set
y_pred = lr.predict(X_test)
```

* Tested the first 20 first prediction values vs real values

```
print("Predictions:\n", list(y_pred[:20]), "\nReal values: \n", list(y_test[:20]))
```

Predictions:	{2, 2, 2, 0, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2}
Real values:	{2, 2, 2, 0, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2}

Looks promising



* Ran metrics and model evaluation

```
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score

Accuracy score is: 93.5 %
Precision score is: 67.51 %
Recall score is: 64.34 %
F1 score is: 65.63 %

Accuracy = TP + FN
TP + TN + FP + FN
Precision = TP
TP + FP
Recall = TP
TP + FN
F1-score = 2 * Precision * Recall
Precision + Recall
```

Logistic Regression

- A **supervised ML classification** technique

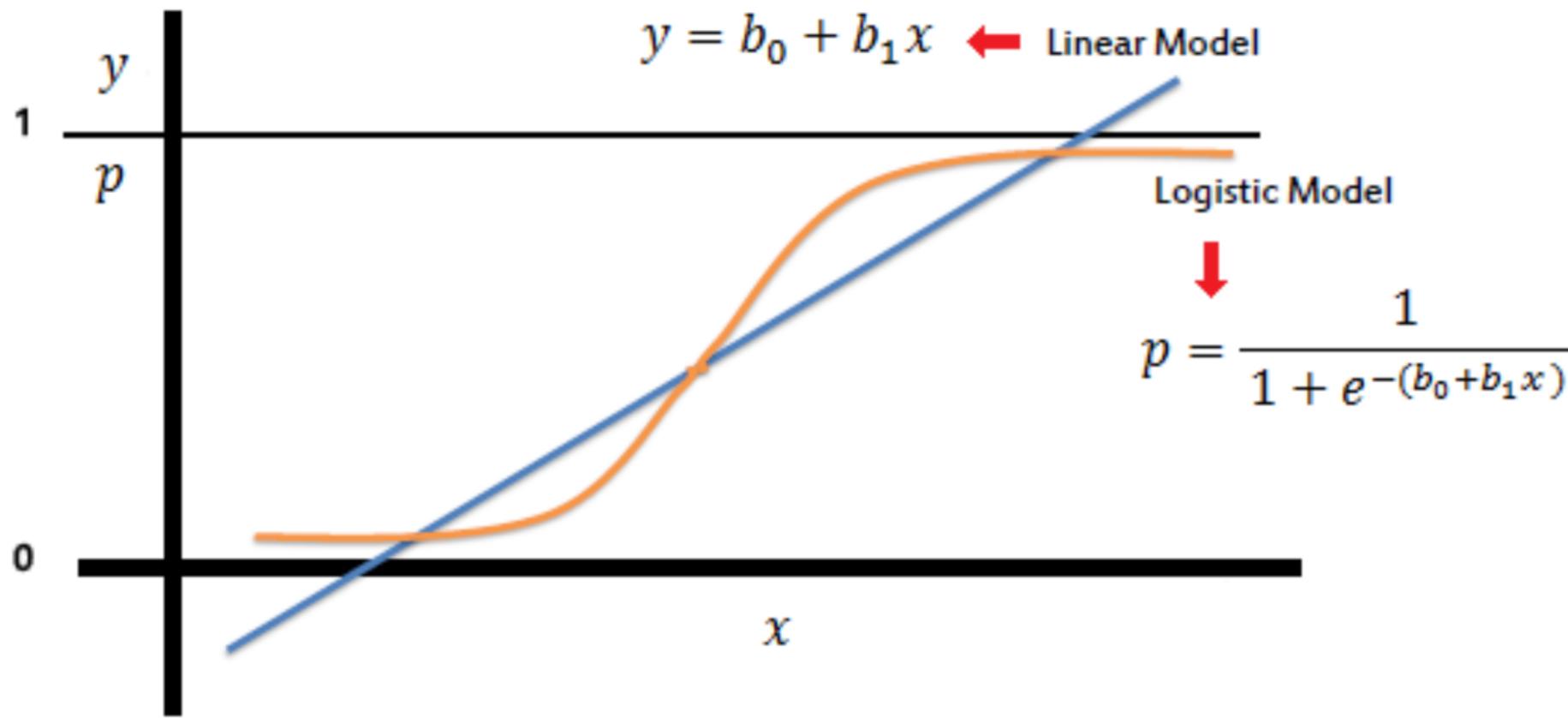
Logistic Regression

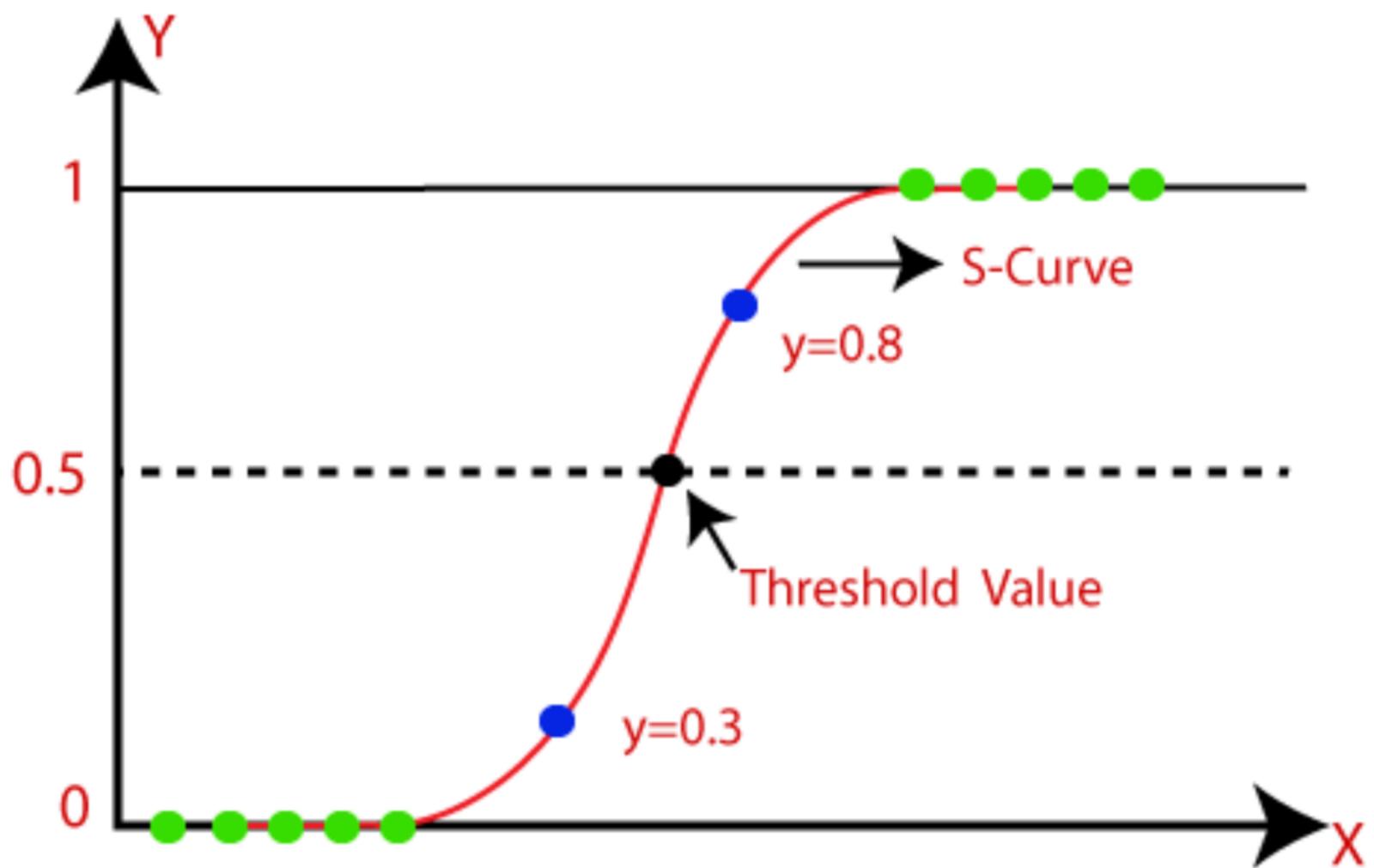
- A **supervised ML classification** technique
- Takes a set of predictor variables
(independent variables)

Logistic Regression

- A **supervised ML classification** technique
- Takes a set of predictor variables (**independent variables**)
- Predicts the **probability** of an observation belonging to a certain class or category (**dependent variables**)

Logistic Regression: used for binary classification
(dependent variable belongs to 1 of 2 classes)





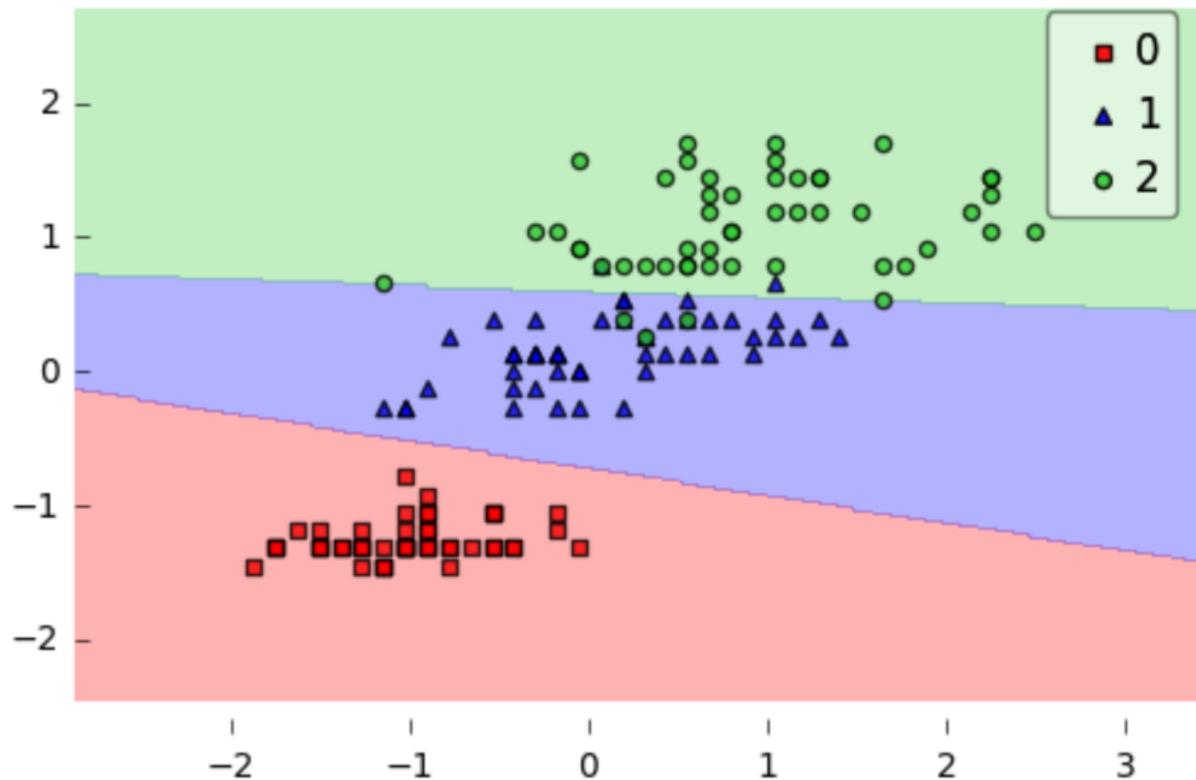
Our dataset had **3** possible **classes**

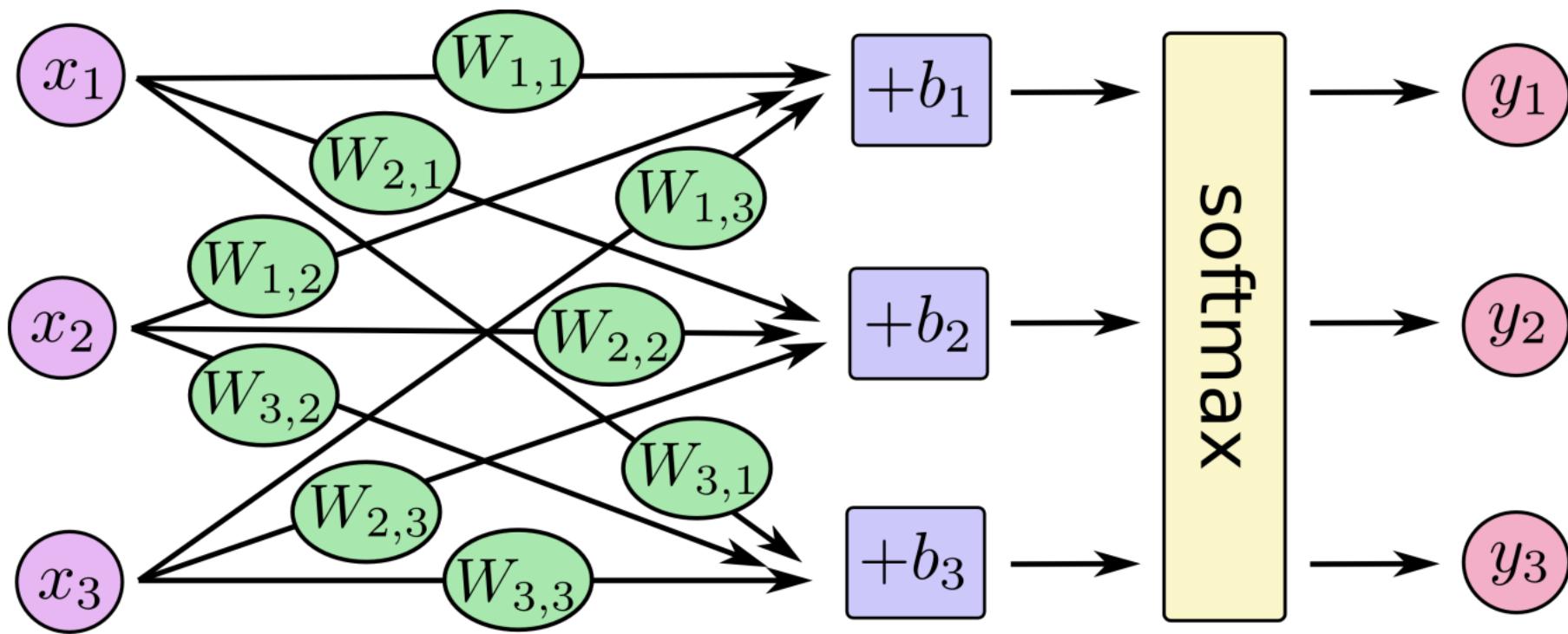
Multinomial Logistic Regression:

- An extension of logistic regression used to handle **multiclass** problems.

Multinomial Logistic Regression: used for multiclass classification
(more than 2 or K classes)

Softmax Regression - Stochastic Gradient Descent





* First, we defined 'X' as the training dataset and 'y' as the target variable

```
x = data.drop(['CLASS_Category', 'CLASS', 'Gender'], axis= 1)
y = data['CLASS_Category']
```

* Split the dataset into 80/20 training/testing split with a random state of 50

```
from sklearn.model_selection import train_test_split  
  
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=.2, random_state = 50)  
  
print(data.shape)  
print(X_train.shape)      (1000, 14)  
print(X_test.shape)      (800, 11)  
print(y_train.shape)      (800,)  
print(y_test.shape)      (200,)
```

* Trained the model

```
from sklearn.linear_model import LogisticRegression  
  
lr = LogisticRegression(multi_class='multinomial', solver='lbfgs', max_iter=2600)  
lr.fit(X_train, y_train)
```

* Predictions on the test set

```
y_pred = lr.predict(X_test)
```

```
y_pred = lr.predict(X_test)
```

* Tested the first 20 first prediction values vs real values

```
print(f"Predictions: \n{list(y_pred[:20])}\n\nReal values: \n{list(y_test[:20])}")
```

```
Predictions:  
[2, 2, 2, 0, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2]  
  
Real values:  
[2, 2, 2, 0, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2]
```



Looks promising

* Ran metrics and model evaluation

```
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
```

Accuracy score is: 93.5 %
Precision score is: 67.51 %
Recall score is: 64.06 %
f1 score is: 65.63 %



$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

$$Precision = \frac{TP}{TP + FP}$$

$$Recall = \frac{TP}{TP + FN}$$

$$F1-score = \frac{2 \times Precision \times Recall}{Precision + Recall}$$

Relatively high **accuracy score** – making correct predictions for most of the instances.

(Note: with unbalanced dataset high accuracy may be misleading)

Relatively high **accuracy score** – making correct predictions for most of the instances.

(Note: with unbalanced dataset high accuracy may be misleading)

Precision score is relatively low indicative that model may be prone to false positives i.e. identifying non-diabetic cases as diabetic.

(only 67.51% of all predicted diabetes cases were actually diabetic.)

Relatively high **accuracy score** – making correct predictions for most of the instances.

(Note: with unbalanced dataset high accuracy may be misleading)

Precision score is relatively low indicative that model may be prone to false positives i.e. identifying non-diabetic cases as diabetic.

(only 67.51% of all predicted diabetes cases were actually diabetic.)

Recall score also relatively low which indicates model wasn't great at detecting all the actual diabetic cases in the dataset (only 64.06%) i.e. missing a considerable amount of diabetic cases.

Relatively high **accuracy score** – making correct predictions for most of the instances.

(Note: with unbalanced dataset high accuracy may be misleading)

Precision score is relatively low indicative that model may be prone to false positives i.e. identifying non-diabetic cases as diabetic.

(only 67.51% of all predicted diabetes cases were actually diabetic.)

Recall score also relatively low which indicates model wasn't great at detecting all the actual diabetic cases in the dataset (only 64.06%) i.e. missing a considerable amount of diabetic cases.

F1 score shows moderate balance between correctly identifying diabetes cases and minimizing false positives (i.e. between precision and recall) HOWEVER highlights the need for further optimization.

4.1 Logistic Regression - Classification

4.2 - Normalization of the data - MinMaxScaler

```
* We then normalized our data using MinMaxScaler to see if this would improve the performance of our Logistic
Regression model (in our case normalization may help in particular with scaling of features with different ranges)
from sklearn.preprocessing import MinMaxScaler

* Created a MinMaxScaler object
scaler = MinMaxScaler()

* Fit and transform the data - columns to be normalized selected, transformation scales the columns, and result of
transformation stored in new dataframe
x_norm = pd.DataFrame(scaler.fit_transform([[ 'Age', ... 'Cr', ... 'Bal', ... 'Bal', ... 'M', ... 'L', ... 'Y']]), 
columns=[['Age', ... 'Cr', ... 'Bal', ... 'Bal', ... 'M', ... 'L', ... 'Y']])
x_norm = pd.concat([x_norm, df['Age', ... 'M', ... 'Gender_Catogory']], axis=1)

* Training/testing split on the X_norm dataset
x_norm_train, x_norm_test, y_norm_train, y_norm_test = train_test_split(x_norm, y, test_size=0.2, random_state=42)

* Trained the model
from sklearn.linear_model import LogisticRegression
lr_norm = LogisticRegression(multi_class='multinomial', solver='liblinear', max_iter=200)
lr_norm.fit(x_norm_train, y_norm_train)

* Predictions on the test set
y_norm_pred = lr_norm.predict(x_norm_test)

* Tested the first 20 first prediction values vs real values
print(f'Predictions:\n{list(y_norm_pred[:20])}\nReal values: \n{list(y_test[:20])}')
print(f'')
* Ran metrics and model evaluation
print(f'Accuracy score: {lr_norm.score(x_norm_test, y_norm_test)}')
print(f'Precision score: {precision_score(y_norm_test, y_norm_pred, average="weighted")}')
print(f'Recall score: {recall_score(y_norm_test, y_norm_pred, average="weighted")}')
print(f'F1 score: {f1_score(y_norm_test, y_norm_pred, average="weighted")}'
```

* We then normalized our data using **MinMaxScaler** to see if this would improve the performance of our Logistic Regression model (in our case normalization may help in particular with scaling of features with different ranges)

```
from sklearn.preprocessing import MinMaxScaler
```

* Created a MinMaxScaler object

```
scaler = MinMaxScaler()
```

* Fit and transform the data – columns to be normalized selected, transformation scales the columns, and result of transformation stored in new dataframe

```
X_norm = pd.DataFrame(scaler.fit_transform(X[['Urea', 'Cr', 'HbA1c', 'Chol', 'TG', 'HDL', 'LDL', 'VLDL']]),  
columns=X[['Urea', 'Cr', 'HbA1c', 'Chol', 'TG', 'HDL', 'LDL', 'VLDL']].columns)  
  
X_norm = pd.concat([X_norm, X[['AGE', 'BMI', 'Gender_Category']]], axis=1)
```

* Training/testing split on the X_norm dataset

```
X_norm_train, X_norm_test, y_norm_train, y_norm_test = train_test_split(X_norm, y, test_size=.2, random_state = 42)
```

* Trained the model

```
from sklearn.linear_model import LogisticRegression

lr_norm = LogisticRegression(multi_class='multinomial', solver='lbfgs', max_iter=700)
lr_norm.fit(X_norm_train, y_norm_train)
```

* Predictions on the test set

```
y_norm_pred = lr_norm.predict(X_norm_test)
```

*Tested the first 20 first prediction values vs real values

```
print(f"Predictions: \n{list(y_norm_pred[:20])}\n\nReal values: \n{list(y_test[:20])}")
```

```
Predictions:  
[2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 0, 2, 0, 2, 2, 2, 2, 2, 2]  
  
Real values:  
[2, 2, 2, 0, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2]
```



Some errors

- * Ran metrics and model evaluation

```
Accuracy score is: 92.0 %
Precision score is: 69.9 %
Recall score is: 65.92 %
f1 score is: 64.5 %
```



- Similarly to before, relatively high **accuracy score** – making correct predictions for most of the instances.

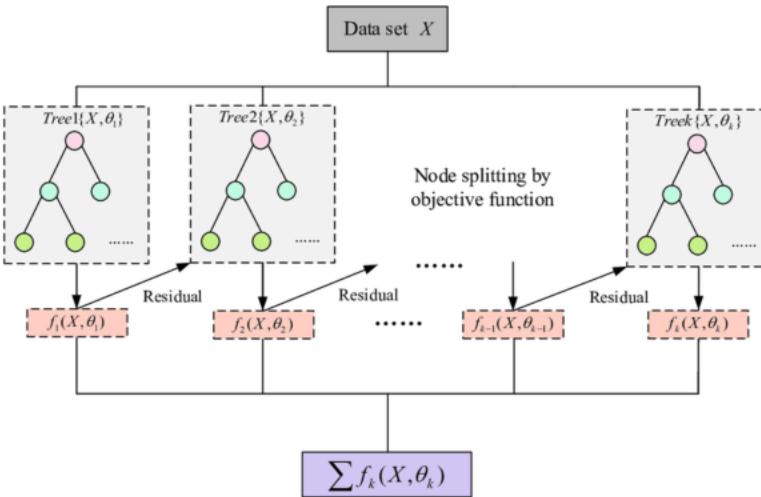
- Similarly to before, relatively high **accuracy score** – making correct predictions for most of the instances.
- **Precision score** is relatively low (slightly higher after normalization) indicative that model may be prone to false positives i.e. identifying non-diabetic cases as diabetic.

- Similarly to before, relatively high **accuracy score** – making correct predictions for most of the instances.
- **Precision score** is relatively low (slightly higher after normalization) indicative that model may be prone to false positives i.e. identifying non-diabetic cases as diabetic.
- **Recall score** also relatively low (very slightly higher after normalization) – indicates model wasn't great at detecting all the actual diabetic cases i.e. missing a considerable amount of diabetic cases.

- Similarly to before, relatively high **accuracy score** – making correct predictions for most of the instances.
- **Precision score** is relatively low (slightly higher after normalization) indicative that model may be prone to false positives i.e. identifying non-diabetic cases as diabetic.
- **Recall score** also relatively low (very slightly higher after normalization) – indicates model wasn't great at detecting all the actual diabetic cases i.e. missing a considerable amount of diabetic cases.
- **F1 score** is slightly lower than before normalization and shows moderate balance between correctly identifying diabetes cases and minimizing false positives, HOWEVER, once again, this highlights the need for further optimization and a new strategy.

4.3 XGboost

XGboost (Extreme Gradient Boosting) is known as a “**black box**” algorithm and is an advanced boosting algorithm that **creates an ensemble of decision trees in a sequential manner** – correcting errors made by previous trees to build a more accurate and robust predictive model.



```
from xgboost import XGBClassifier
* Trained the model and created the model instance
bst = XGBClassifier(n_estimators=3, max_depth=3, learning_rate=1, objective='multiple:logistic')
```

* Fitted the model

```
bst.fit(X_train, y_train)
```

* Made predictions on test set

```
y_XG_preds = bst.predict(X_test)
```

```
*Tested the first 20 first prediction values vs real values
print(f"Predictions:\n{list(y_XG_preds[:20])}\n\nReal values: \n{list(y_test[:20])}")
```

Predictions: [2, 2, 2, 0, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2]	→	Looks good
Real values: [2, 2, 2, 0, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2]		

Model Evaluation

```
Accuracy score is: 100.0 %
Precision score is: 100.0 %
Recall score is: 100.0 %
f1 score is: 100.0 %
```

These scores suggest that this model is performing **exceptionally well** on the dataset

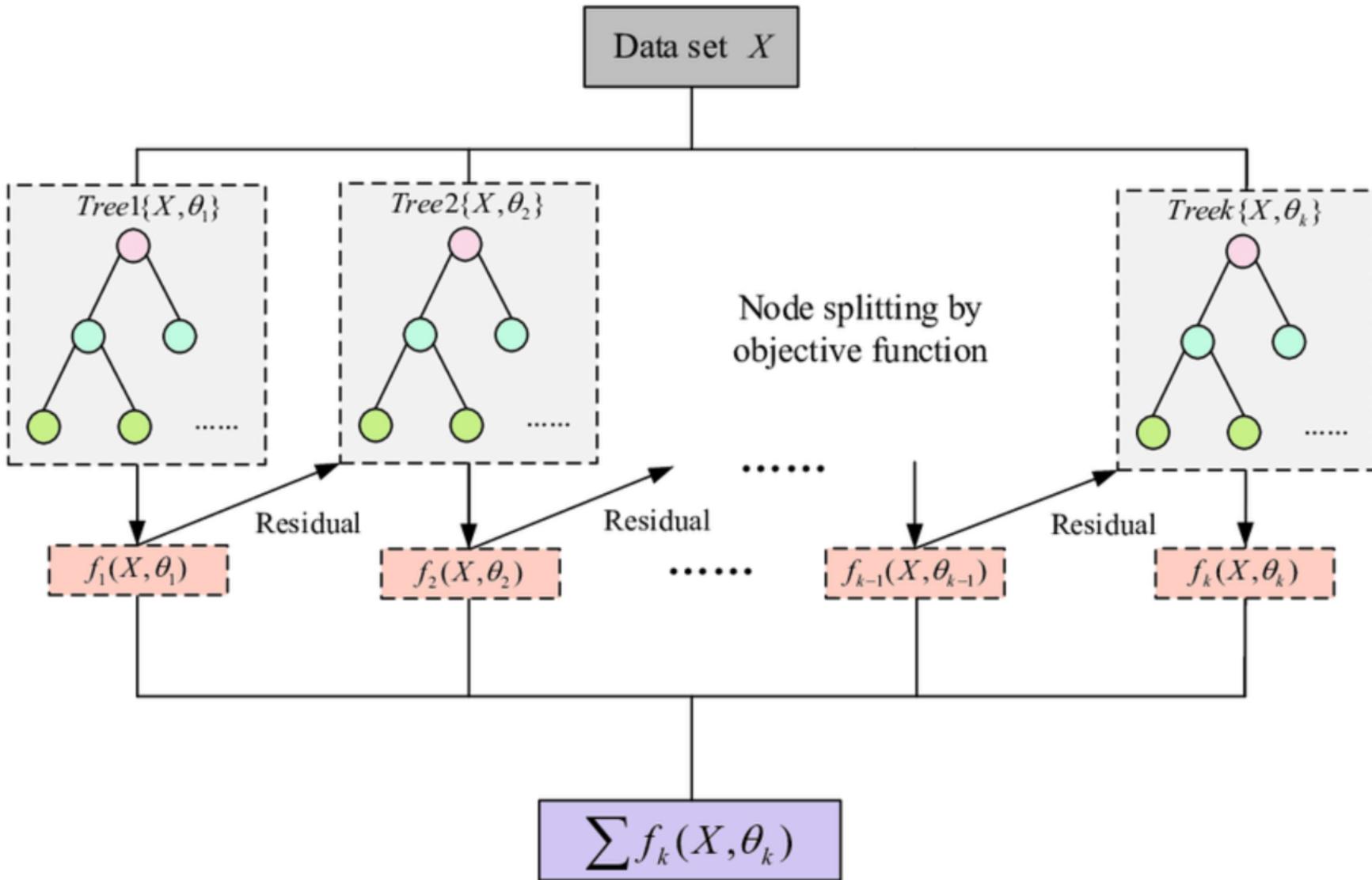
It is identifying all the 3 possible results, without making any false positive or negative predictions.

XGboost (Extreme Gradient Boosting) is known as a “**black box**” algorithm and is an advanced boosting algorithm that **creates an ensemble of decision trees in a sequential manner** – correcting errors made by previous trees to build a more accurate and robust predictive model.

XGboost (Extreme Gradient Boosting) is known as a “**black box**” algorithm and is an advanced boosting algorithm that **creates an ensemble of decision trees in a sequential manner** – correcting errors made by previous trees to build a more accurate and robust predictive model.

It uses a gradient boosting algorithm to

- **minimize the residual/loss** when adding new models.
- It is able to **handle complex data, reduce overfitting**, and
- provide **high prediction performance**.



```
from xgboost import XGBClassifier
```

* Trained the model and created the model instance

```
bst = XGBClassifier(n_estimators=3, max_depth=3, learning_rate=1, objective='multiple:logistic')
```

* Fitted the model

```
bst.fit(X_train, y_train)
```

* Fitted the model

```
bst.fit(X_train, y_train)
```

* Made predictions on test set

```
y_XG_preds = bst.predict(X_test)
```

*Tested the first 20 first prediction values vs real values

```
print(f"Predictions: \n{list(y_XG_preds[:20])}\n\nReal values: \n{list(y_test[:20])}")
```

```
Predictions:  
[2, 2, 2, 0, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2]  
  
Real values:  
[2, 2, 2, 0, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2]
```



Looks good

Model Evaluation

Accuracy score is: 100.0 %

Precision score is: 100.0 %

Recall score is: 100.0 %

f1 score is: 100.0 %

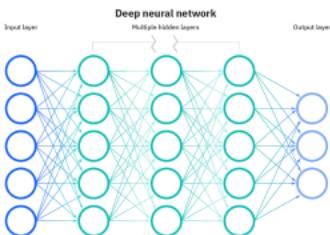
These scores suggest that this model is performing **exceptionally well** on the dataset

It is identifying all the 3 possible results, without making any false positive or negative predictions.

4.4 NEURAL NETWORKS

Neural Networks (NN)

- A class of **machine learning** models inspired by the structure and function of the **human brain**.
 - They consist of interconnected nodes or "**neurons**" that receive input data, process it, and then pass the output to other nodes.



```
from keras.models import Sequential  
from keras.layers import Dense  
from keras.utils import to_categorical
```

- A type of neural network
 - Designed to process input data in a sequential manner
 - Multiple layers of nodes are organized in a linear, sequential way forming a feedforward network.

* Converted labels in the target variable (`y_train`) using one-hot encoding

```
y_NN_train = to_categorical(y_train)
```

```
print(y_train[:3])
```

```
886    2  
488    2  
265    2  
Name: CLASS_Categor...  
  
print(y_NN_train[:3])  
  
[[0. 0. 1.]  
 [0. 0. 1.]  
 [0. 0. 1.]]
```

* Trained the keras model

```
model = Sequential()
model.add(Dense(200, activation='relu', input_dim=100))
model.add(Dense(200, activation='relu'))
model.add(Dense(3, activation='softmax'))
```

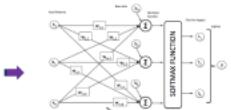
```

model.summary()

Model: "sequential"
_________________________________________________________________
Layer (type)                 Output Shape              Param #
dense (Dense)                (None, 200)             24000
dense_1 (Dense)               (None, 200)             40000
dense_2 (Dense)               (None, 3)               603
_________________________________________________________________

Total params: 43,283
Trainable params: 43,283
Non-trainable params: 0

```



```
* Compiled the keras model

model.compile(
    loss='categorical_crossentropy',
    optimizer='adam',
    metrics=['accuracy'])

# as model on the dataset



```

Dataset = 2000
Batch_size = 32
X_M_train, X_M_val, y_M_train, y_M_val = train_test_split(X_M, y_M, test_size=0.2)
X_train, X_val, y_train, y_val = train_test_split(X_M_train, y_M_train, test_size=0.2)

x_train = np.reshape(X_train, (X_train.shape[0], 28, 28, 1))
x_val = np.reshape(X_val, (X_val.shape[0], 28, 28, 1))

y_train = np_utils.to_categorical(y_train, num_classes)
y_val = np_utils.to_categorical(y_val, num_classes)

model.fit(x_train, y_train,
 batch_size=32,
 epochs=10,
 validation_data=(x_val, y_val),
 verbose=1)

Model evaluation
loss, accuracy = model.evaluate(X_M_train, y_M_train)
print("Accuracy: {}% \nLoss: {}".format(accuracy * 100, loss))

Class predictions with the model
predictions = model.predict(X_M_val)
predictions[:5]
array([[0., 0., 1.],
 [0., 0., 1.],
 [0., 0., 1.],
 [0., 0., 1.],
 [0., 0., 1.]], dtype=float32)

Class predictions with the model but round them into the closest integer
predictions = predictions.round()
predictions[:5]
array([[1., 0., 0.],
 [1., 0., 0.],
 [1., 0., 0.],
 [1., 0., 0.],
 [1., 0., 0.]], dtype=float32)

* First 10 predictions vs real values

Predictions: Real values:
[[0., 0., 1.]] [[0., 0., 1.]]
[[0., 0., 1.]] [[0., 0., 1.]]
[[0., 0., 1.]] [[0., 0., 1.]]
[[1., 0., 0.]] [[1., 0., 0.]]
[[0., 0., 1.]] [[0., 0., 1.]]
[[0., 0., 1.]] [[0., 0., 1.]]
[[0., 0., 1.]] [[0., 0., 1.]]
[[0., 0., 1.]] [[0., 0., 1.]]
[[0., 0., 1.]] [[0., 0., 1.]]
[[0., 0., 1.]] [[0., 0., 1.]]
```



* Changed categorical values into integers as it was in the y_test



```
y_M_pred = []
for prediction in predictions:
 max_index = np.argmax(prediction)
 y_M_pred.append(max_index)
y_M_pred[:5]
```



* First 20 prediction values vs real values



```
Predictions:
[2, 2, 2, 0, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2]

Real values:
[2, 2, 2, 0, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2]
```



* Metrics and model evaluation



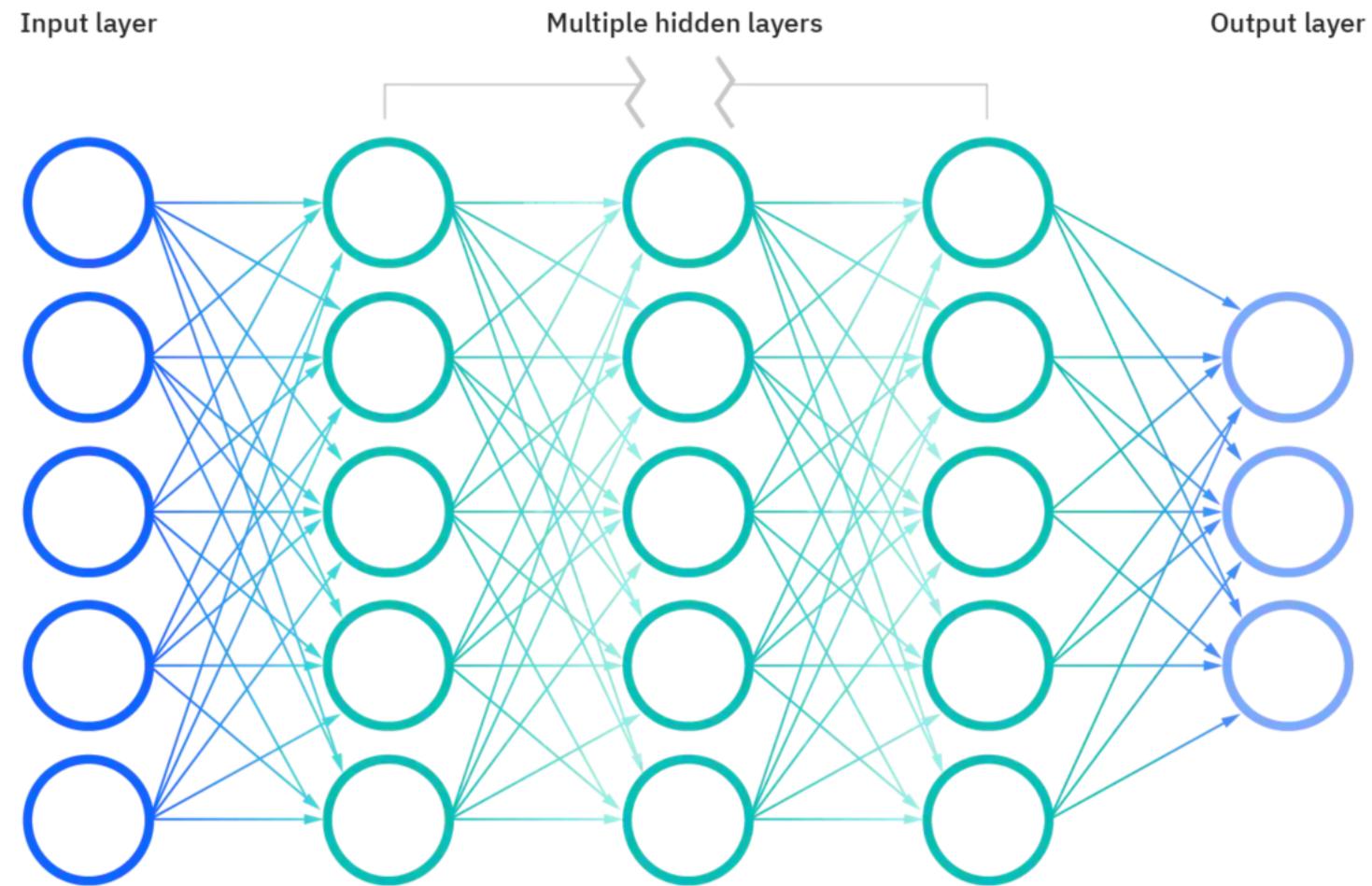
```
Accuracy score is: 95.5 %
Precision score is: 96.61 %
Recall score is: 95.5 %
F1 score is: 95.7 %
```


```

Neural Networks (NN)

- A class of **machine learning** models inspired by the structure and function of the **human brain**.
- They consist of interconnected nodes or “**neurons**” that receive input data, process it, and then pass the output to other nodes.

Deep neural network



```
from keras.models import Sequential  
from keras.layers import Dense  
from keras.utils import to_categorical
```



Sequential Model

- A type of neural network
- Designed to process input data in a sequential manner
- Multiple layers of nodes are organized in a linear, sequential way forming a feedforward network.

* Converted labels in the target variable (`y_train`) using one-hot encoding

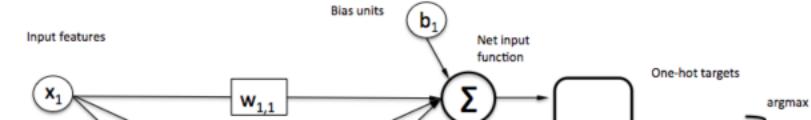
```
y_NN_train = to_categorical(y_train)
```

```
print(y_train[:3])
```

```
886    2  
488    2  
265    2  
Name: CLASS_Category, dtype: int64
```

```
print(y_NN_train[:3])
```

```
[[0. 0. 1.]  
 [0. 0. 1.]  
 [0. 0. 1.]]
```



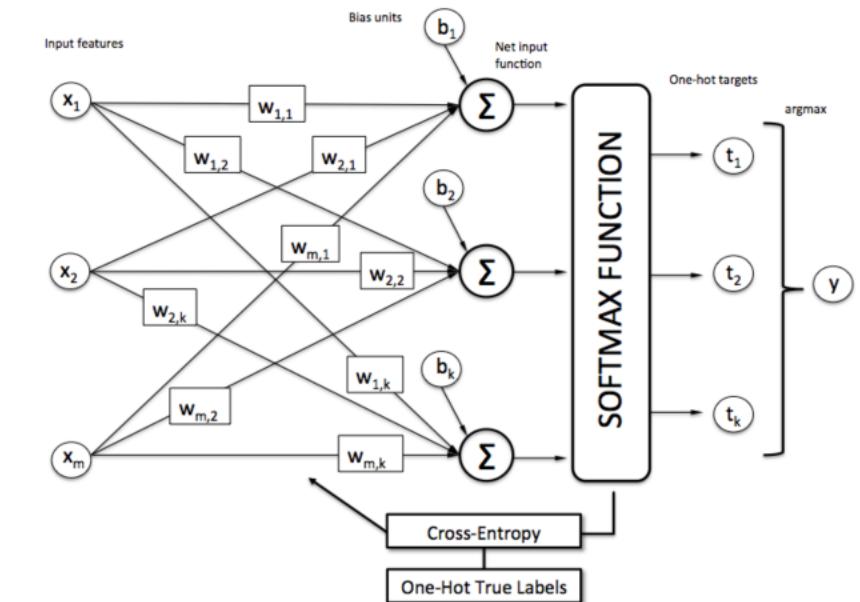
Name: CLASS_Category, dtype: int64

```
print(y_NN_train[:3])
```

```
[[0. 0. 1.]  
 [0. 0. 1.]  
 [0. 0. 1.]]
```

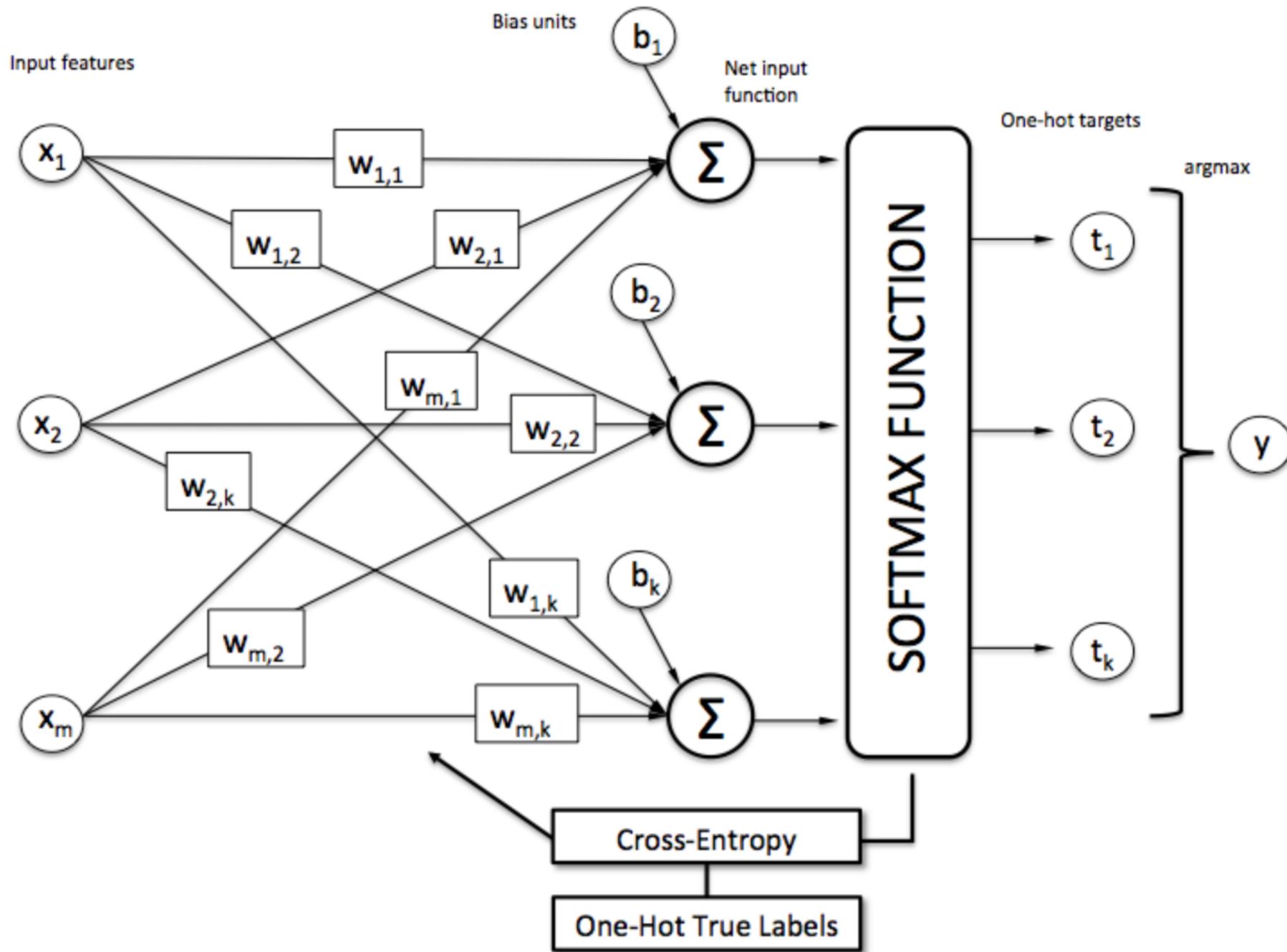
* Trained the keras model

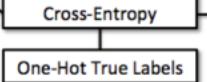
```
model = Sequential()  
model.add(Dense(200, activation='relu', input_shape=(11,)))  
model.add(Dense(200, activation='relu'))  
model.add(Dense(3, activation='softmax'))
```



```
model.summary()
```

Model: "sequential"





```
model.summary()
```

Model: "sequential"

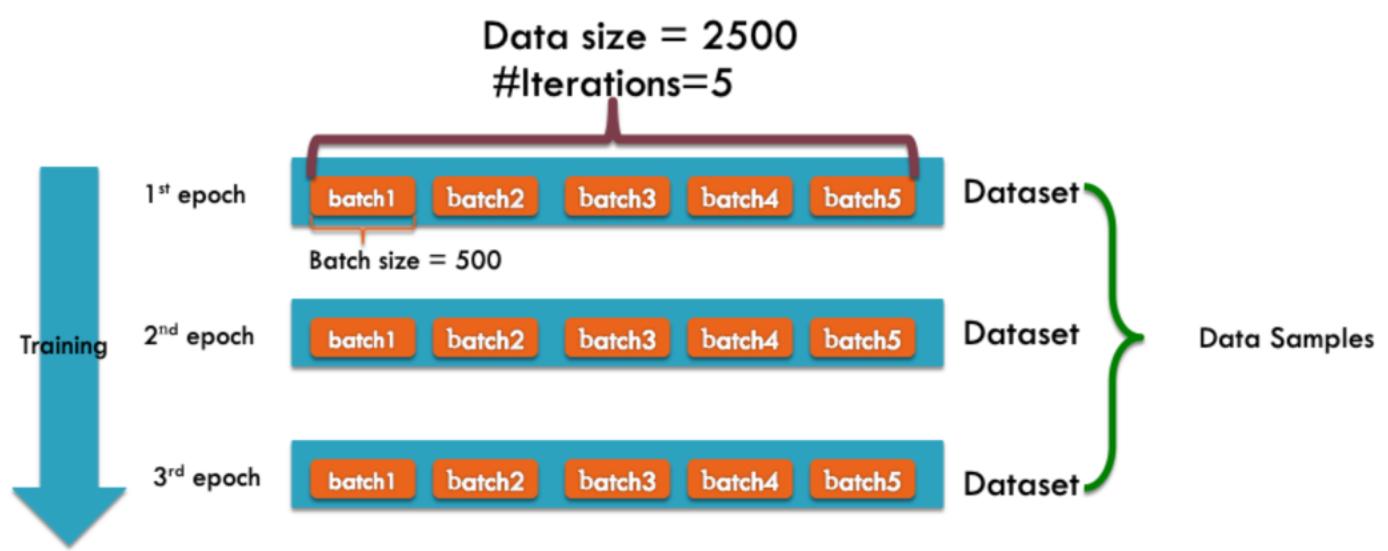
Layer (type)	Output Shape	Param #
=====		
dense (Dense)	(None, 200)	2400
dense_1 (Dense)	(None, 200)	40200
dense_2 (Dense)	(None, 3)	603
=====		
Total params: 43,203		
Trainable params: 43,203		
Non-trainable params: 0		

* Compiled the keras model

```
model.compile(  
    loss='categorical_crossentropy',  
    optimizer='adam',  
    metrics=['accuracy'])
```

* Fitted the keras model on the dataset

```
model.fit(  
    X_train,  
    y_NN_train,  
    epochs=500,  
    batch_size=30  
)
```



```
Epoch 494/500
27/27 [=====] - 0s 4ms/step - loss: 0.0036 - accuracy: 1.0000
Epoch 495/500
27/27 [=====] - 0s 3ms/step - loss: 0.0043 - accuracy: 1.0000
Epoch 496/500
27/27 [=====] - 0s 3ms/step - loss: 0.0033 - accuracy: 1.0000
Epoch 497/500
27/27 [=====] - 0s 3ms/step - loss: 0.0028 - accuracy: 1.0000
Epoch 498/500
27/27 [=====] - 0s 3ms/step - loss: 0.0027 - accuracy: 1.0000
Epoch 499/500
27/27 [=====] - 0s 3ms/step - loss: 0.0029 - accuracy: 1.0000
Epoch 500/500
27/27 [=====] - 0s 3ms/step - loss: 0.0025 - accuracy: 1.0000
<keras.callbacks.History at 0x7d9fc11e19c0>
```

* Evaluation of the keras model

```
loss, accuracy = model.evaluate(X_train, y_NN_train)
```

Accuracy:100.00%
Loss: 0.0015

* Class predictions with the model

```
predictions = model.predict(X_test)  
predictions[:5]
```

```
7/7 [=====] - 0s 2ms/step  
array([[1.0368682e-10, 2.9687493e-05, 9.9997032e-01],  
       [1.9128512e-16, 2.3144937e-06, 9.9999774e-01],  
       [1.6650246e-19, 2.4478529e-28, 1.0000000e+00],  
       [9.9999857e-01, 1.0152958e-06, 3.1146052e-07],  
       [4.9838639e-10, 2.4898136e-05, 9.9997509e-01]], dtype=float32)
```

* Class predictions with the model but round them into the closest integer

```
predictions = predictions.round()  
predictions[:5]
```

```
array([[0., 0., 1.],  
       [0., 0., 1.],  
       [0., 0., 1.],  
       [1., 0., 0.],  
       [0., 0., 1.]], dtype=float32)
```

* First 10 predictions vs real values

* Changed categorical values into integers as it was in the y_test

```
y_NN_pred = []

for prediction in predictions:
    max_index = np.argmax(prediction)      [2, 2, 2, 0, 2]
    y_NN_pred.append(max_index)

y_NN_pred[:5]
```

* First 20 prediction values vs real values

Predictions:

```
[2, 2, 2, 0, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2]
```

Real values:

```
[2, 2, 2, 0, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2]
```

* Metrics and model evaluation

Accuracy score is: 95.5 %

Precision score is: 96.61 %

Recall score is: 95.5 %

F1 score is: 95.7 %

- High **accuracy score** – making correct predictions for most of the instances.

- High **accuracy score** – making correct predictions for most of the instances.
- High **precision score** is indicative that model is making very few false positive predictions i.e. very few non-diabetic cases as diabetic (96.61% of all predicted diabetes cases were actually diabetic).

- High **accuracy score** – making correct predictions for most of the instances.
- High **precision score** is indicative that model is making very few false positive predictions i.e. very few non-diabetic cases as diabetic (96.61% of all predicted diabetes cases were actually diabetic).
- **Recall score** also high – indicates model very good at detecting all the actual diabetic cases in the dataset (only 95.5%) i.e. missing very few diabetic cases.

- High **accuracy score** – making correct predictions for most of the instances.
- High **precision score** is indicative that model is making very few false positive predictions i.e. very few non-diabetic cases as diabetic (96.61% of all predicted diabetes cases were actually diabetic).
- **Recall score** also high – indicates model very good at detecting all the actual diabetic cases in the dataset (only 95.5%) i.e. missing very few diabetic cases.
- **F1 score** shows good balance between correctly identifying diabetes cases and minimizing false positives

Diabetes Prediction Model

Using Machine Learning Regression Algorithms and Neural Networks



Index

Presentation by Nastaran Ashoori & Daniella Vuletić
CodeOp DSPT10 - July 2023



1- Introduction

**2- Data Pre-
Processing & Feature
Engineering**

**3- Exploratory
Data Analysis
(EDA)**

**4- Machine
Learning (ML)**

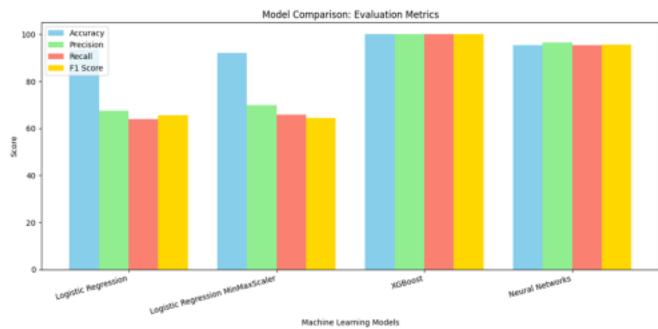
5- Conclusion

5 -CONCLUSIONS

Quick backtrack to study objective...

We aimed to design a **predictive diabetes model** that could estimate the likelihood of diabetes development based on the provided factors.

SUMMARY OF RESULTS



REFERENCES

1. <https://www-who.int/news-room/detail/5-august-2018-diabetes>
2. <http://www.cdc.gov/diabetes/managing/insulinuse/>
3. DiMeglio LA, Evans-Molina C, Orman RA. Type 1 diabetes. *Lancet*. 2018 Jun 16;391(10138):2449-2462. doi: 10.1016/S0140-6736(18)31320-5. PMID: 29916386; PMCID: PMC601118
4. Risérus U, Whitt WC, Hu FB. Dietary fats and prevention of type 2 diabetes. *Prog Lipid Res*. 2009 Jan;48(1):44-51. doi: 10.1016/j.plipres.2008.10.002. Epub 2008 Dec 7. PMID: 19052963; PMCID: PMC2654180
5. Taliash SA. Is Diabetes Becoming the Biggest Epidemic of the Twenty-first Century? *Int J Health Sci (Qassim)*. 2007 Jul;1(2):V-VIII. PMID: 21675425; PMCID: PMC2068646
6. Taliash AG, Hender C, Rutherford W, Brunner EL, Kivimaki M. Prediabetes: a high-risk state for diabetes development. *Lancet*. 2012 Jun 16;379(9833):2279-90. doi: 10.1016/S0140-6736(12)60283-9. Epub 2012 Jun 9. PMID: 22681126; PMCID: PMC3891203
7. Carris NW, Magness RR, Labovitz AJ. Prevention of Diabetes Mellitus in Patients With Prediabetes. *Am J Cardiol*. 2019 Feb 1;123(3):507-512. doi: 10.1016/j.amjcard.2018.10.032. Epub 2018 Nov 6. PMID: 30528418; PMCID: PMC6395098
8. <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3891203/>
9. <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6395098/>
10. <https://data.mendeley.com/datasets/cu9h7vndzg-2/>
11. <https://data.mendeley.com/datasets/cu9h7vndzg-2/>

Big thanks to FitLipa and Shandy!

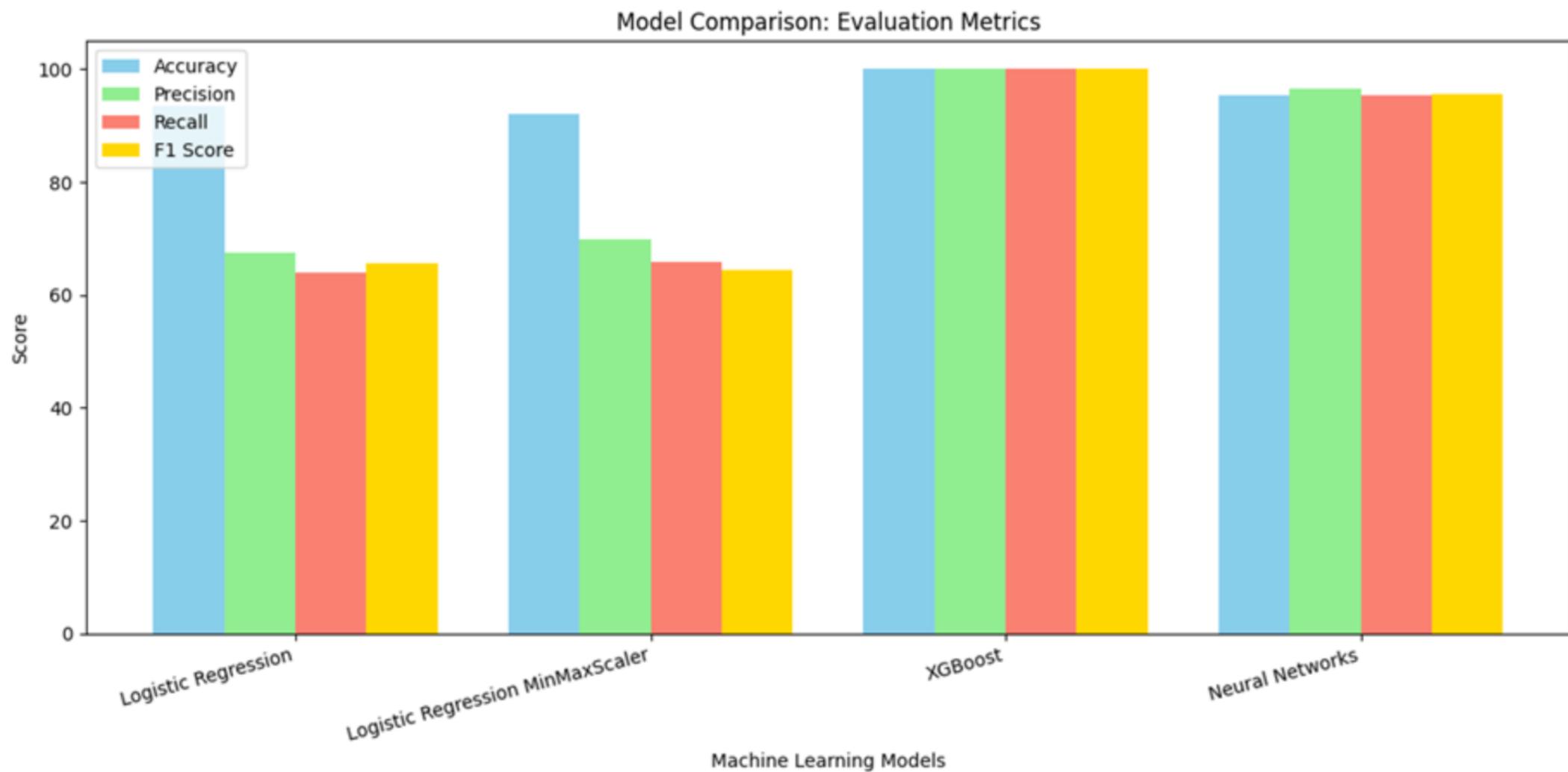


Contact details:
danielaluvetic21@gmail.com
n.ashoori1989@gmail.com
Nastaran Ashoori on LinkedIn

Quick backtrack to study objective...

We aimed to design a **predictive diabetes model** that could estimate the likelihood of diabetes development based on the provided factors.

SUMMARY OF RESULTS



- According to our dataset and analysis, the **variables** with the strongest correlation, and therefore most likely to put people at **risk** for the development of **diabetes**, include an increased body mass index (**BMI**), **HbA1c** (blood glucose) and **age**.

- According to our dataset and analysis, the **variables** with the strongest correlation, and therefore most likely to put people at **risk** for the development of **diabetes**, include an increased body mass index (**BMI**), **HbA1c** (blood glucose) and **age**.

- Our machine learning models were able to successfully predict whether a patient belonged to the non-diabetic, pre-diabetic or diabetic class, based on the given variables. The **XGboost** model performing the best.

- According to our dataset and analysis, the **variables** with the strongest correlation, and therefore most likely to put people at **risk** for the development of **diabetes**, include an increased body mass index (**BMI**), **HbA1c** (blood glucose) and **age**.

- Our machine learning models were able to successfully predict whether a patient belonged to the non-diabetic, pre-diabetic or diabetic class, based on the given variables. The **XGboost** model performing the best.

- Further studies would benefit from an increased **sample size** and perhaps a more **balanced dataset**. A more **globally representative** sample would also add value.

REFERENCES

1. <https://www.who.int/news-room/fact-sheets/detail/diabetes>
2. <https://www.cdc.gov/diabetes/managing/problems>
3. DiMeglio LA, Evans-Molina C, Oram RA. Type 1 diabetes. Lancet. 2018 Jun 16;391(10138):2449-2462. doi: 10.1016/S0140-6736(18)31320-5. PMID: 29916386; PMCID: PMC6661119.
4. Risérus U, Willett WC, Hu FB. Dietary fats and prevention of type 2 diabetes. Prog Lipid Res. 2009 Jan;48(1):44-51. doi: 10.1016/j.plipres.2008.10.002. Epub 2008 Nov 7. PMID: 19032965; PMCID: PMC2654180.
5. <https://bmcmedicine.biomedcentral.com/articles/10.1186/s12916-017-0901x>
6. Tabish SA. Is Diabetes Becoming the Biggest Epidemic of the Twenty-first Century? Int J Health Sci (Qassim). 2007 Jul;1(2):V-VIII. PMID: 21475425; PMCID: PMC3068646
7. <https://diabetesatlas.org>
8. Tabák AG, Herder C, Rathmann W, Brunner EJ, Kivimäki M. Prediabetes: a high-risk state for diabetes development. Lancet. 2012 Jun 16;379(9833):2279-90. doi: 10.1016/S0140-6736(12)60283-9. Epub 2012 Jun 9. PMID: 22683128; PMCID: PMC3891203.
9. Carris NW, Magness RR, Labovitz AJ. Prevention of Diabetes Mellitus in Patients With Prediabetes. Am J Cardiol. 2019 Feb 1;123(3):507-512. doi: 10.1016/j.amjcard.2018.10.032. Epub 2018 Nov 6. PMID: 30528418; PMCID: PMC6350898.
10. <https://www.cdc.gov/diabetes/basics/prediabetes>
11. <https://data.mendeley.com/datasets/wj9rwkp9c2/1>

Big thanks to FilLipa and Shandy!



Contact details:
daniellavuletic218@gmail.com
n.ashoori1989@gmail.com
Nastaran Ashoori on [Linkedin](#)



Diabetes Prediction Model

Using Machine Learning Regression Algorithms and Neural Networks



Index

Presentation by Nastaran Ashoori & Daniella Vuletić
CodeOp DSPT10 - July 2023



1- Introduction

**2- Data Pre-
Processing & Feature
Engineering**

**3- Exploratory
Data Analysis
(EDA)**

**4- Machine
Learning (ML)**

5- Conclusion

Diabetes Prediction Model

Using Machine Learning Regression Algorithms and Neural Networks



Index

Presentation by Nastaran Ashoori & Daniella Vuletić
CodeOp DSPT10 - July 2023



1- Introduction

**2- Data Pre-
Processing & Feature
Engineering**

**3- Exploratory
Data Analysis
(EDA)**

**4- Machine
Learning (ML)**

5- Conclusion