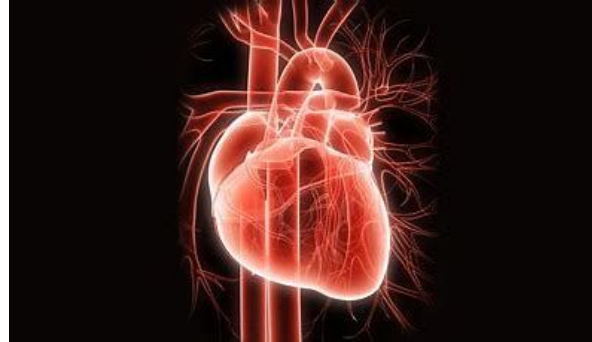


# Cardiovascular Analysis

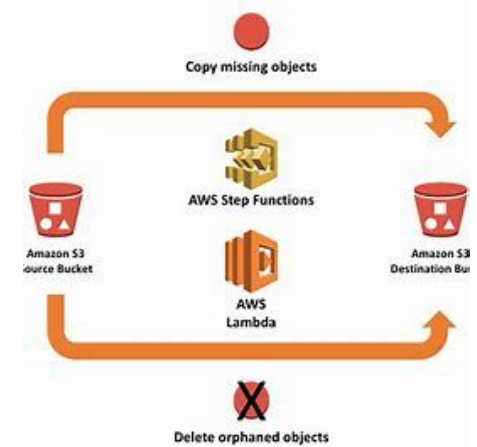
Chi - Scrum Master / Developer  
Gabby - Developer  
Alejandro - Developer



# Introduction

We aim to establish a streamlined data pipeline in a cloud environment

- Integrating Amazon S3 buckets
- Lambda functions
- RDS
- Potentially EC2



This is crucial for organizations aiming to derive actionable insights from their data assets effectively, enabling informed decisions.

# Objectives & Issues

Use dataset (Cardiovascular Disease) to run analysis on our local machine:

- Creating an S3 bucket
- Creating an RDS
- Creating a triggering Lambda function
- Pull data from the cloud environment and run analysis

**Issues:** Although we were able to successfully build this infrastructure, we ran into some issues

- Incomplete data (only about a  $\frac{1}{3}$  of our data was displayed)
- Permissions (Access denied)

# Project Tracking

- We used Jira to track our progress

## KAN board

CA

GG

TO DO 1

Discuss extensions

✓ KAN-7

CA

+ Create issue

IN PROGRESS 2

Lambda Issues due to access and permissions.

✓ KAN-15

Create database schema

✓ KAN-14

DONE 6 ✓

Checkpoint 3: Create a lambda.

✓ KAN-5

✓

Checkpoint 4: Conduct analysis.

✓ KAN-6

GG

Checkpoint 1: Set up your S3 bucket.

✓ KAN-3

✓

Set up project process

✓ KAN-1

CA

Checkpoint 2: Set up an RDS

✓ KAN-4

CA

Complete and push PRD to repository

✓ KAN-2

## Timeline

CA GG

Status category ▾ Epic ▾

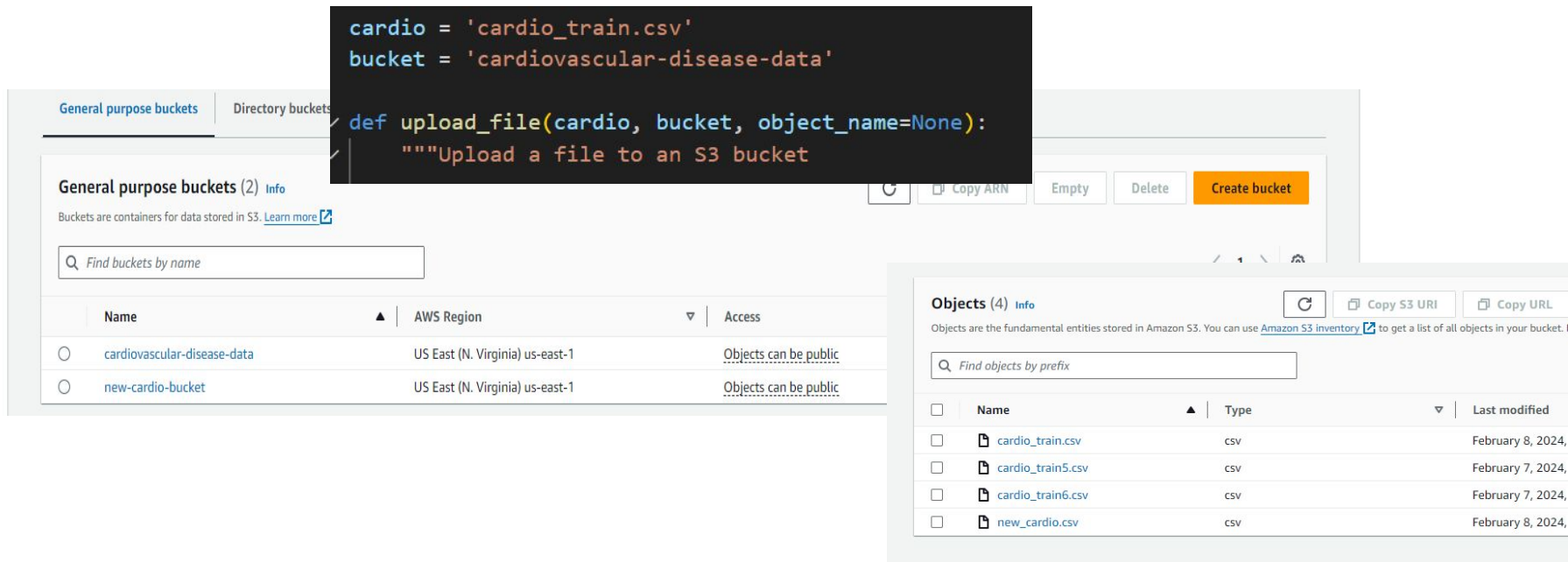
	FEB				FEB							
	1	2	3	4	5	6	7	8	9	10	11	12
⚡ KAN-8 Set up project process					█							
⚡ KAN-9 Checkpoint #1						█						
⚡ KAN-10 Checkpoint #2							█					
⚡ KAN-11 Checkpoint #3								█				
⚡ KAN-12 Checkpoint #4									█			

# Checkpoint 1: Create S3 bucket

- After specifying the bucket name and additional configuration, we were able to successfully create a bucket.
- IAM users were created and given permission to access the bucket.

```
cardio = 'cardio_train.csv'
bucket = 'cardiovascular-disease-data'

def upload_file(cardio, bucket, object_name=None):
    """Upload a file to an S3 bucket
```



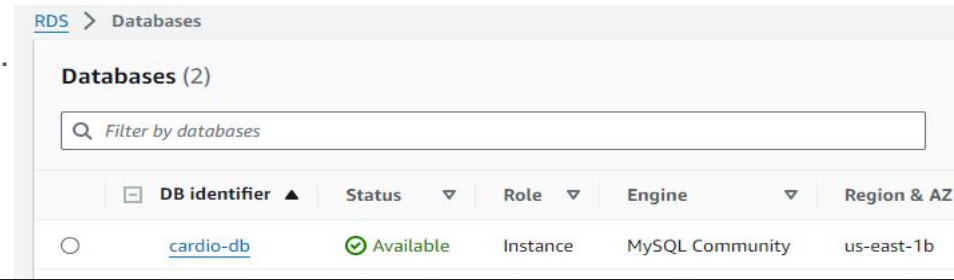
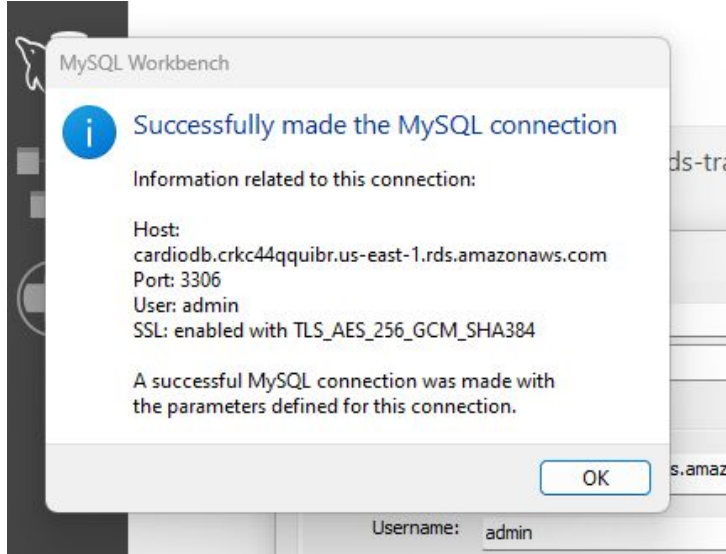
The screenshot displays the AWS Management Console interface. At the top, there are tabs for 'General purpose buckets' and 'Directory buckets'. The 'General purpose buckets' tab is active, showing a list of buckets. Below the tabs, there is a search bar labeled 'Find buckets by name'. The list of buckets includes 'cardiovascular-disease-data' and 'new-cardio-bucket', both located in 'US East (N. Virginia) us-east-1' and marked as 'Objects can be public'. To the right of the bucket list, there are buttons for 'Copy ARN', 'Empty', 'Delete', and 'Create bucket'. A dark overlay with Python code is positioned over the bucket list. Below the bucket list, there is a section for 'Objects (4)' with a search bar labeled 'Find objects by prefix'. The list of objects includes 'cardio\_train.csv', 'cardio\_train5.csv', 'cardio\_train6.csv', and 'new\_cardio.csv', all of type 'csv' and last modified on February 7, 2024, or February 8, 2024. To the right of the object list, there are buttons for 'Copy S3 URI' and 'Copy URL'.

Name	AWS Region	Access
cardiovascular-disease-data	US East (N. Virginia) us-east-1	Objects can be public
new-cardio-bucket	US East (N. Virginia) us-east-1	Objects can be public

Name	Type	Last modified
cardio_train.csv	csv	February 8, 2024,
cardio_train5.csv	csv	February 7, 2024,
cardio_train6.csv	csv	February 7, 2024,
new_cardio.csv	csv	February 8, 2024,

# Checkpoint 2: Create RDS

Used AWS console to create and configured the RDS.

A screenshot of the AWS Management Console 'RDS > Databases' page. It shows a table with two columns: 'DB identifier' and 'Status'. The first entry is 'cardio-db' with a status of 'Available'. Other columns like 'Role', 'Engine', and 'Region & AZ' are also visible but empty for this instance.

DB identifier	Status	Role	Engine	Region & AZ
cardio-db	Available	Instance	MySQL Community	us-east-1b

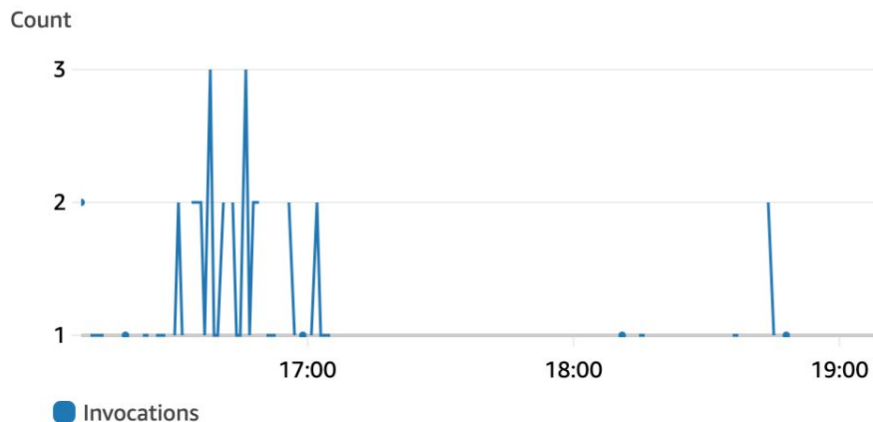
```
# RDS configuration
rds_endpoint = 'cardio-db.crkc44quibr.us-east-1.rds.amazonaws.com'
rds_port = 3306
rds_username = 'admin'
rds_password = 'rootroot'
```

```
def connect_to_rds():
    try:
        connection = mysql.connector.connect(
            host=rds_endpoint,
            port=rds_port,
            user=rds_username,
            password=rds_password,
        )
        print("Connected to RDS successfully.")
        return connection
    except Exception as e:
        print("An error occurred while connecting to RDS:", e)
        return None
```

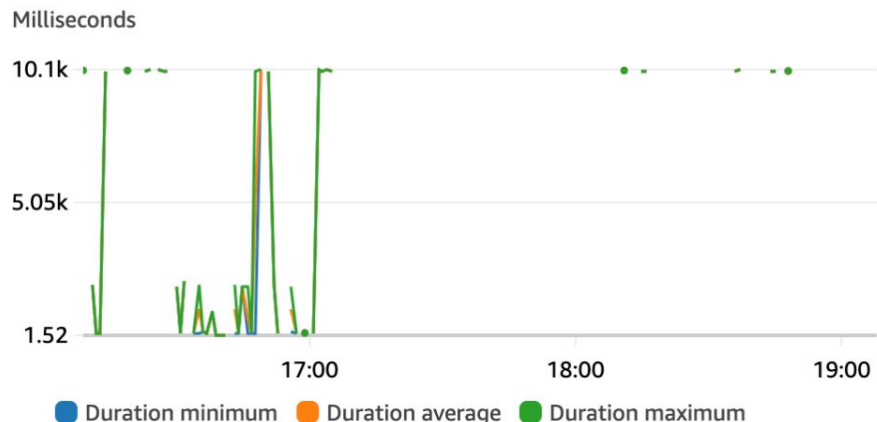
# Checkpoint 3: Create a Lambda function

- This function will act as a triggering event
  - Listening for any activity in the s3 bucket
  - Cleaning the data and registering it to the rds database.
- Creating the Lambda from the rds.
  - Creating a direct connection to the rds proxy

## Invocations



## Duration



# Issues

## Lambda Issues:

- Permissions
  - Getting the correct access for the correct actions
- Timeouts
  - Understanding where and why the data was not being processed successfully
- Troubleshooting
  - Using standard debugging technique and pair programing to pinpoint a solution





▶	2024-02-08T13:44:46.600-05:00	START RequestId: 80923608-6bf8-4215-9142-34953122feaa Version: \$LATEST
▶	2024-02-08T13:44:48.191-05:00	succesfully pulled data from bucket
▶	2024-02-08T13:44:48.191-05:00	bucket: cardiovascular-disease-data
▶	2024-02-08T13:44:48.191-05:00	key: cardio_train.csv
▼	2024-02-08T13:44:56.652-05:00	2024-02-08T18:44:56.652Z 80923608-6bf8-4215-9142-34953122feaa Task timed out after 10.05 seconds
		2024-02-08T18:44:56.652Z 80923608-6bf8-4215-9142-34953122feaa Task timed out after 10.05 seconds



# Database

- All the information needed (ex. patient information and details) - was stored under the same table.

```
1 • use cardio_train;  
2  
3 • select * from cardio_train order by id asc;  
4  
5
```

Result Grid    Filter Rows: <input type="text"/>   Export:    Wrap Cell Content:    Fetch rows: 													
	id	age	gender	height	weight	ap_hi	ap_lo	cholesterol	gluc	smoke	alco	active	cardio
▶	0	18393	2	168	62	110	80	1	1	0	0	1	0
	1	20228	1	156	85	140	90	3	1	0	0	1	1
	2	18857	1	165	64	130	70	3	1	0	0	0	1
	3	17623	2	169	82	150	100	1	1	0	0	1	1
	4	17474	1	156	56	100	60	1	1	0	0	0	0
	8	21914	1	151	67	120	80	2	2	0	0	0	0
	9	22113	1	157	93	130	80	3	1	0	0	1	0
	12	22584	2	178	95	130	90	3	3	0	0	1	1
	13	17668	1	158	71	110	70	1	1	0	0	1	0

cardio_train ▼	
◇	id INT
◇	age INT
◇	gender INT
◇	height INT
◇	weight DOUBLE
◇	ap_hi INT
◇	ap_lo INT
◇	cholesterol INT
◇	gluc INT
◇	smoke INT
◇	alco INT
◇	active INT
◇	cardio INT

# Checkpoint 4: Data Analysis

Checkpoint 4- Pull Data from RDS to create meaningful insight

To uncover patterns, correlations and potential risk factors related to cardiovascular disease we used:

- Pandas
- Matplotlib.pyplot
- SQLAlchemy to create engine

```
from sqlalchemy import create_engine
from config import rds_port, rds_username, rds_password, rds_endpoint

hostname = rds_endpoint
port = rds_port
username = rds_username
password = rds_password
database_name = 'cardio_train'

engine = create_engine(f'mysql+mysqlconnector://{username}:{password}@{hostname}:{port}/{database_name}')
connection = engine.connect()
```

✓ 0.7s

```
import pandas as pd
import matplotlib.pyplot as plt
from sqlalchemy import create_engine
```

```
cardio_csv = 'cardio_train2.csv'
```

```
cardio_df = pd.read_csv(cardio_csv)
```

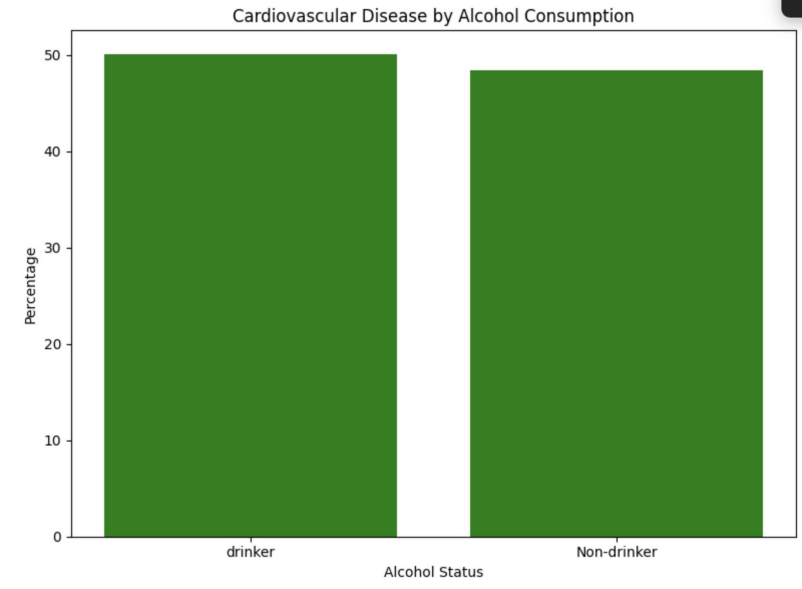
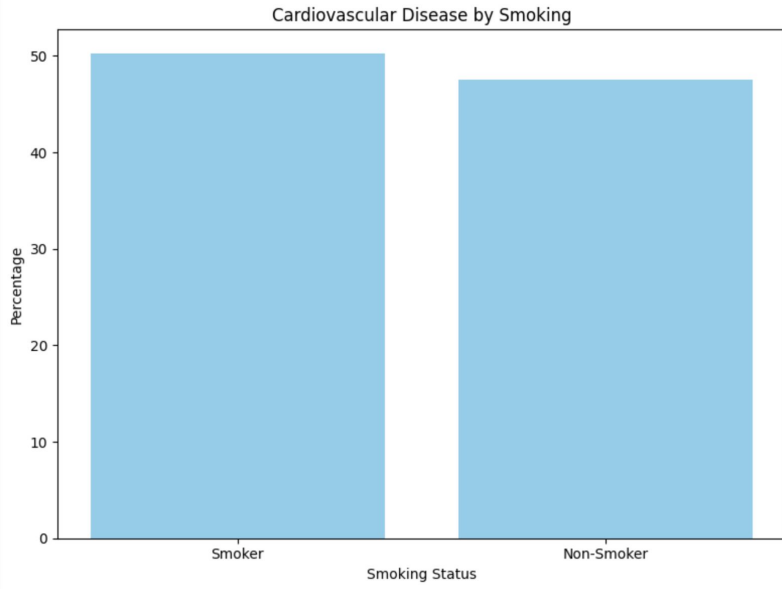
✓ 0.0s

```
engine = create_engine('sqlite:///cardio.db', echo=False)
cardio_df.to_sql('cardio_table', engine, if_exists='replace', index=False)

print("Database created successfully.")
```

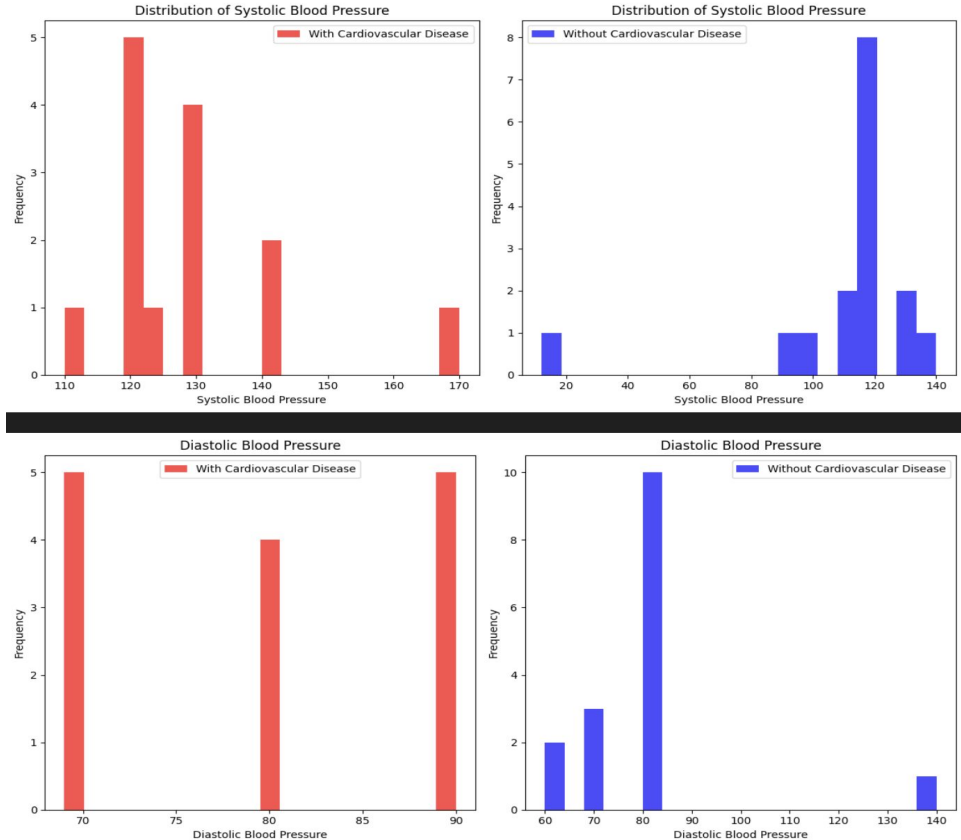
# Impact to Lifestyle

Both Alcohol consumption and Smoking should increase the risk of cardiovascular disease. However from the data set everything was almost split evenly.



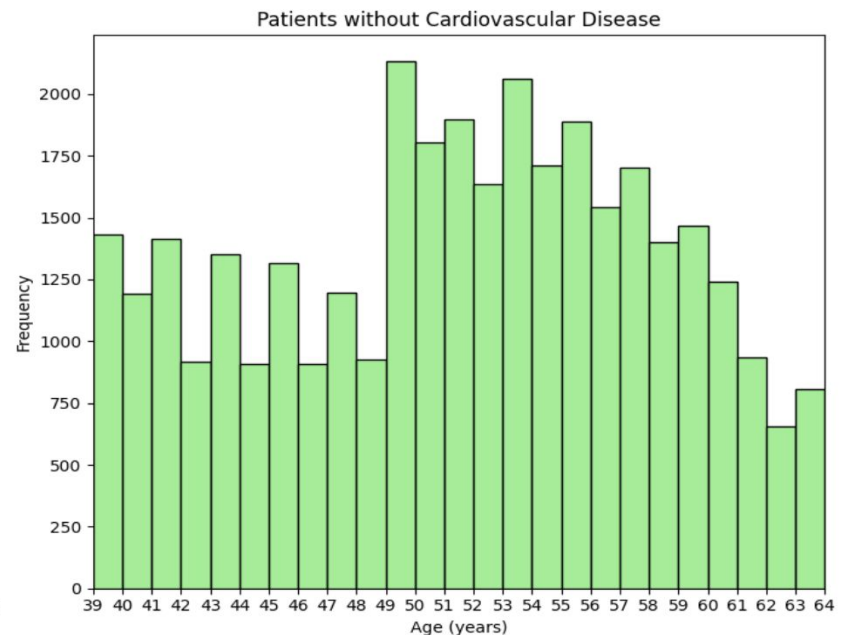
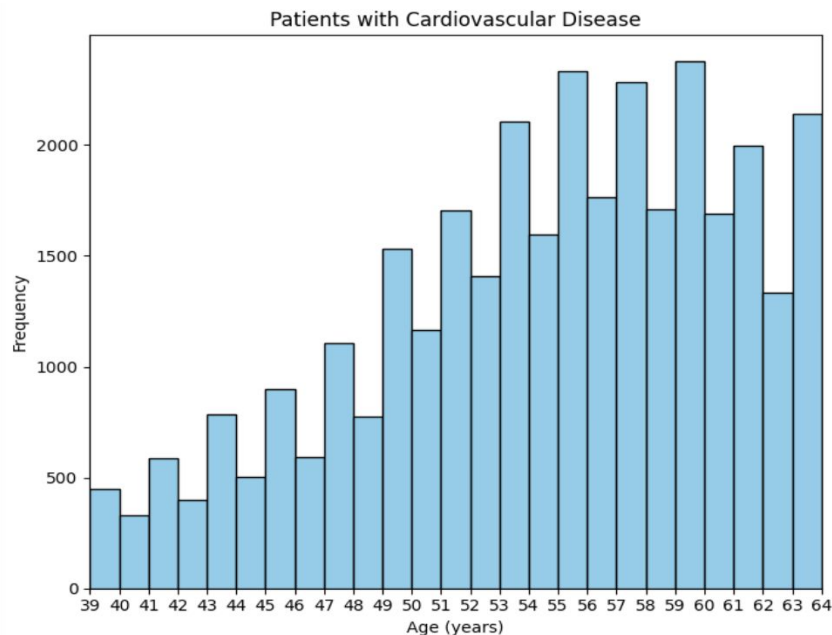
# Blood Pressure and Cardiovascular Disease

Normal blood pressure for most adults is defined as a systolic pressure of less than 120 and a diastolic pressure of less than 80. Elevated blood pressure is defined as a systolic pressure between 120 and 129 with a diastolic pressure of less than 80.



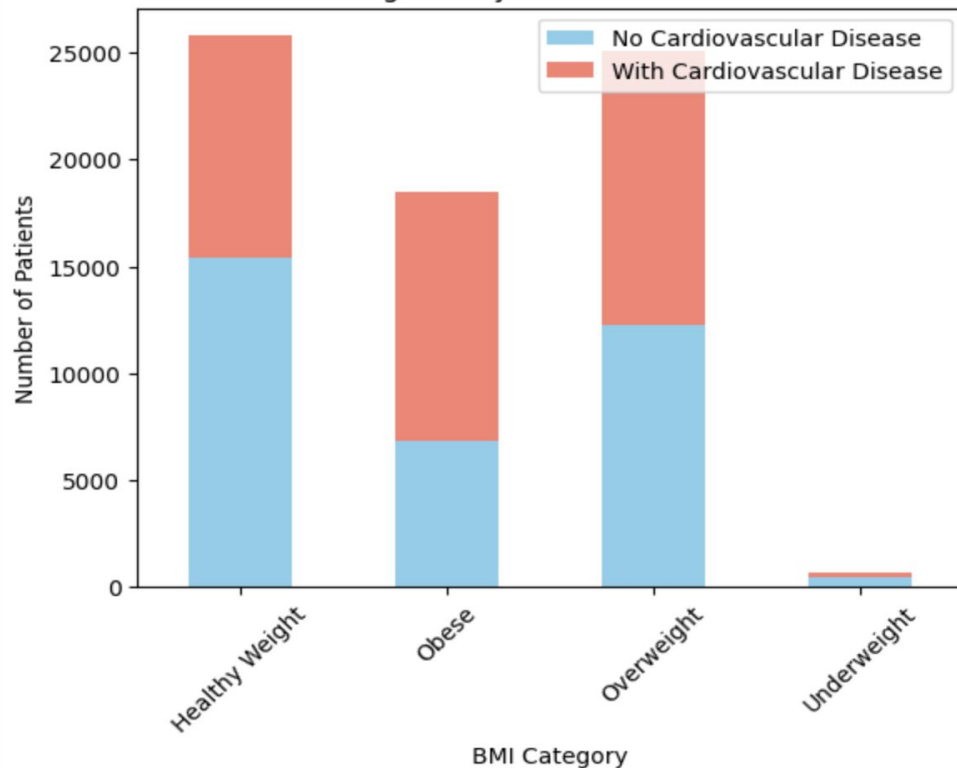
# Age Distribution

The first graph is the main one used for analysis. We can see with an increase of age the risk of cardiovascular disease.



# BMI

BMI Categories by Cardiovascular Disease



Finding the BMI gave clearer insight to which patients were more at risk for cardiovascular disease.

```
print("Numerical values for each BMI category with and without cardiovascular disease:")
print(bmi_cardio_grouped)
```

✓ 0.0s

```
Numerical values for each BMI category with and without cardiovascular disease:
cardio      0      1
bmi_category
Healthy Weight  15440  10350
Obese          6823  11651
Overweight     12290  12796
Underweight     468   182
```

BMI is calculated as  $BMI = \text{weight(kg)} / \text{height(m)}^2$  the dataset had height in cm so we did a conversion.

# Cardiovascular Disease by Gender

We wanted to see what the correlation between Gender and cardiovascular disease was.

