# Data Science
# Modelling, concepts, techniques

**Ayyub Sheikhi**

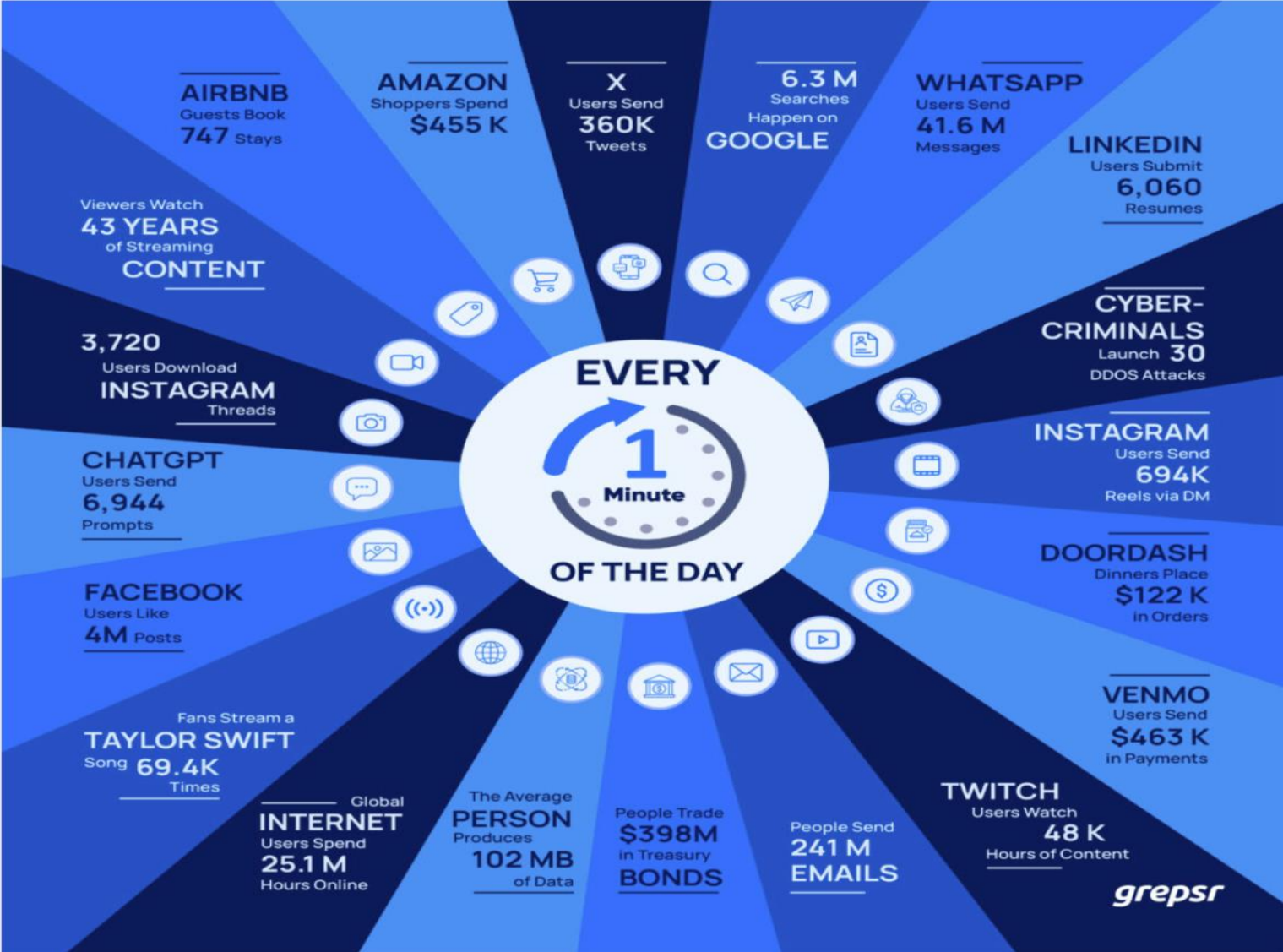**Department of Statistics**

**Shahid Bahonar university of Kerman**

# Goals

✓ Know what Data Science is and learn the basic algorithms

✓ Perform Data Science techniques
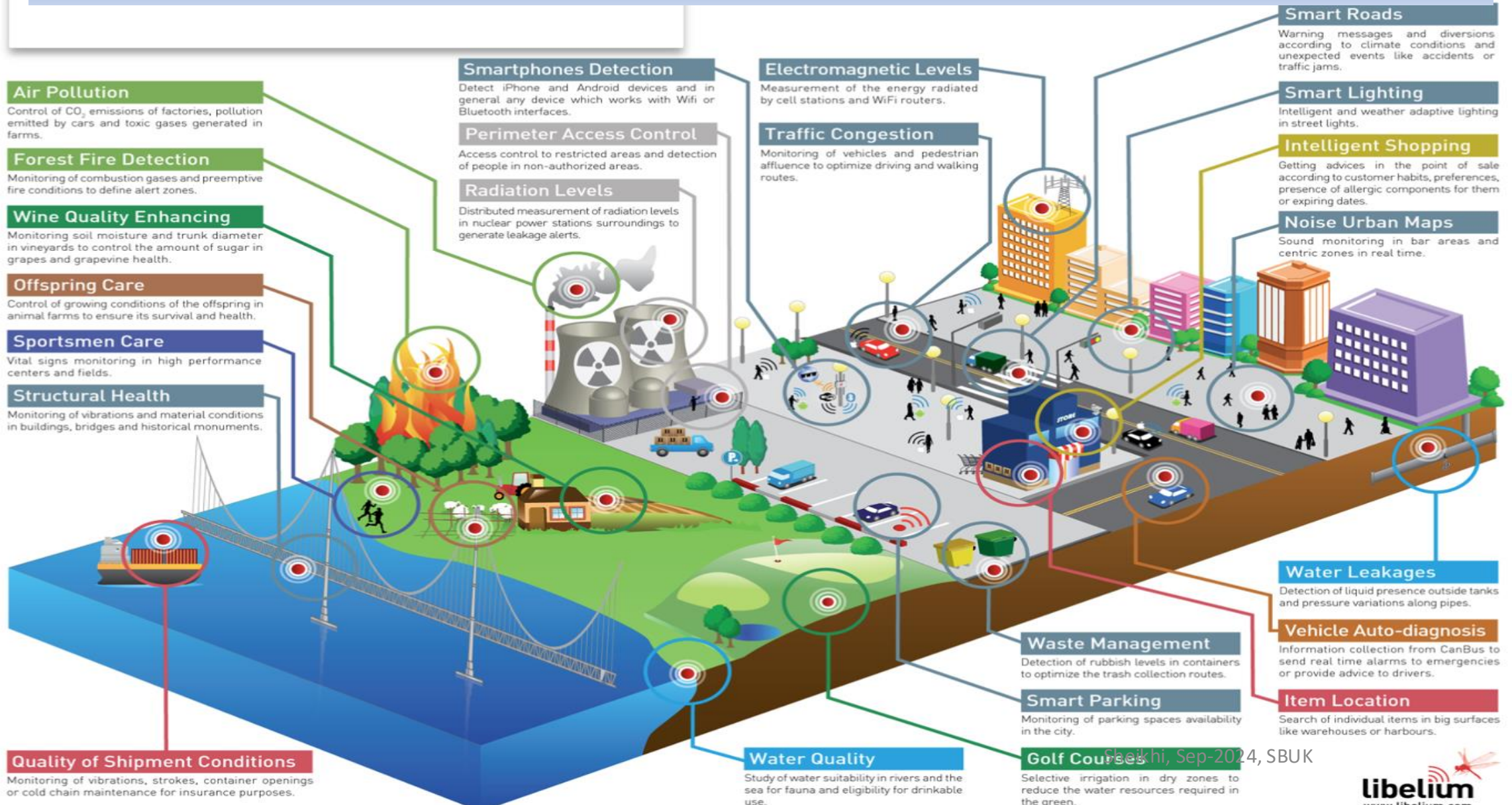
✓ Know how to apply algorithms to real-world applications

## Chapter 1.  Introduction

➢ Why Data Science ?

➢What is Data Science ? How big is Data Science

➢Characteristics of Data Science

➢Top most popular Data Science algorithms

➢Major issues in Data Science

# Why Data Science

# Why Data Science

**Some links of Data Science**

https://everysecond.io

https://www.worldometers.info/coronavirus

http://irsc.ut.ac.ir

https://finance.yahoo.com/quote/BTC-USD/history

## Why Data Science ?

- The Explosive Growth of Data: from terabytes to petabytes

  - Business: Web, e-commerce, transactions, stocks, …

  - Science: Remote sensing, bioinformatics, scientific simulation, …

  - Society and everyone: news, digital cameras, YouTube

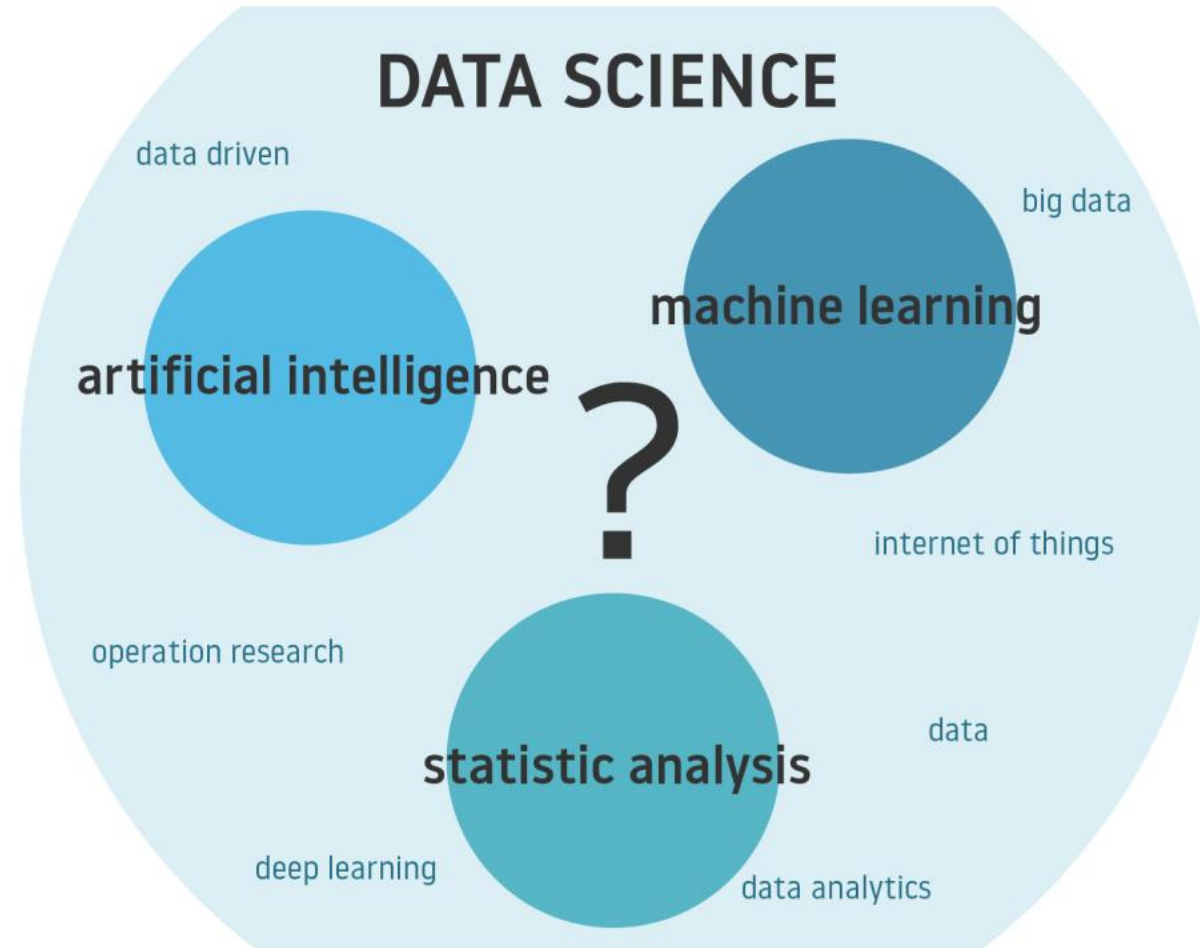  - Healthcares, recording patient symptom using online monitoring

## Evolution of Sciences (from empirical sciences to data science)

- Before 1600, **empirical science**

- 1600-1950s, **theoretical science**
  - Each discipline has grown a *theoretical* component. Theoretical models often motivate experiments and generalize our understanding.

- 1950s-1990s, **computational science**
  - Over the last 50 years, most disciplines have grown a third, *computational* branch (e.g. empirical, theoretical, and computational ecology, or physics, or linguistics.)
  - Computational Science traditionally meant simulation. It grew out of our inability to find closed-form solutions for complex mathematical models.

- 1990-now, **data science/ Big Science**
  - The flood of data from new scientific instruments and simulations
  - The ability to economically store and manage petabytes of data online
  - The Internet and computing Grid that makes all these archives universally accessible
  - Scientific info. management, acquisition, organization, query, and visualization tasks scale almost linearly with data volumes.  Data Science is a major new challenge!
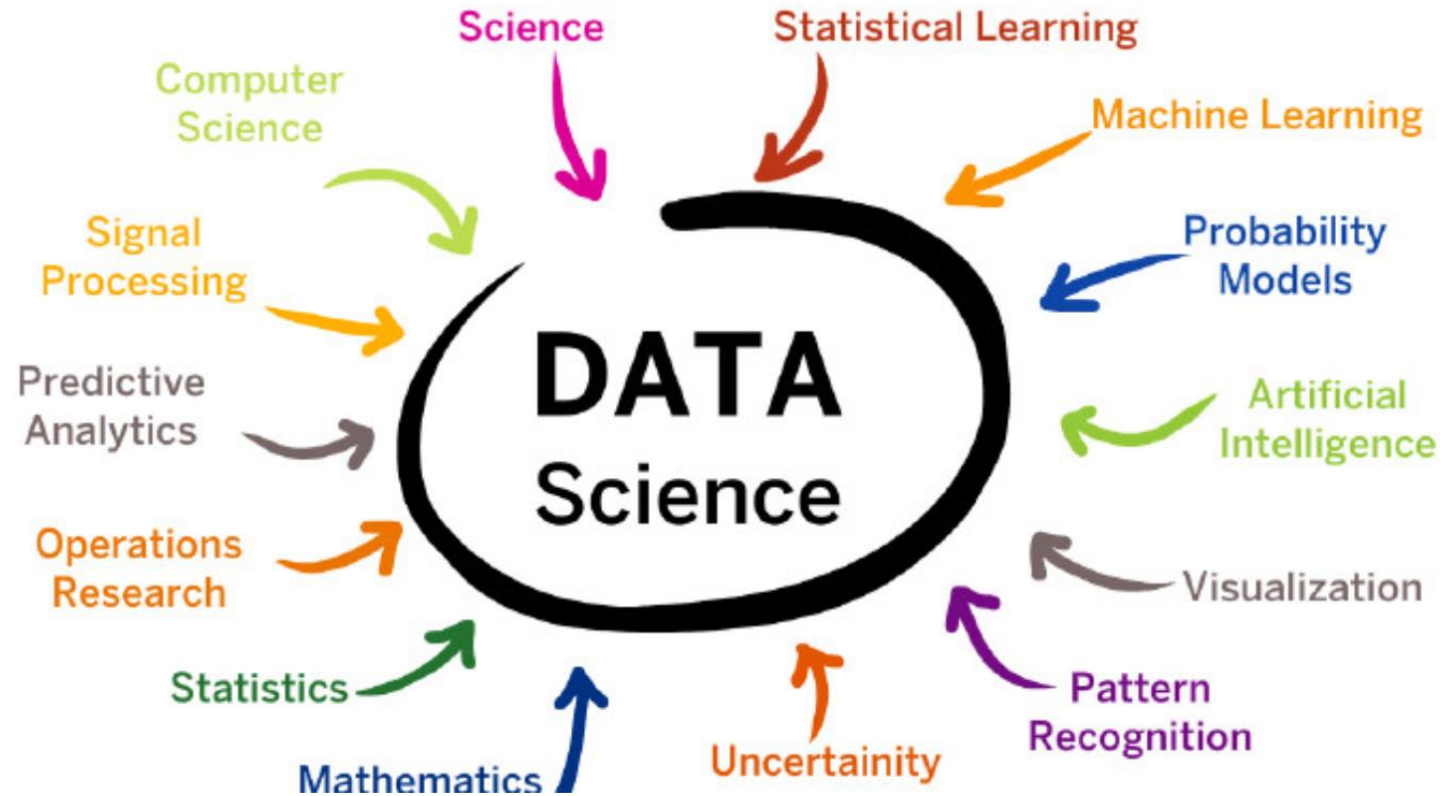
# Characteristics of Data Science in big data ?

# Data Science: Confluence of Multiple Disciplines

# Data Science: Confluence of Multiple Disciplines

**Types of Variables**

- *Nominal* : Name only--Gender, hair color, ethnicity

- *Ordinal* : Nominal categories with an implied order-- Low, medium, high.

- *Discrete*:  Reflects a number obtained by counting— no decimal.

- *Continuous*: Reflects a measurement; the number of decimal places depends on the precision of the measuring device.

## Types of Variables

There are two basic types of variables:

*categorical* and *numerical*.

*Categorical Variables*: variables defined by the classes or categories into which an individual member falls.

*Numerical Variables*: variables to which a number is assigned as a quantitative value.

## Data and Variables

**Data** are often discussed in terms of variables, where a **variable** is:

Any characteristic that *varies* from one member of a population to another.

A simple example is height in centimeters, which varies from person to person.

## Definition of Variables in a data Matrix

- AGE: Age in years

- BMI: Body mass index, weight/height$^2$ in kg/m$^2$

- FFNUM: The average number of times eating "fast food" in a week

- TEMP: High temperature for the day

- GENDER: 1- Female  0- Male

- EXERCISE LEVEL: 1- Low 2- Medium 3- High

- QUESTION: Compared to others, what is your satisfaction rating of the National Practitioner Data Bank?

  1- Very Satisfied  2- Somewhat Satisfied  3- Neutral

  4- Somewhat dissatisfied 5- Dissatisfied

# Data table/Data matrix

| OBS | AGE | BMI | FFNUM | TEMP($^0$F) | GENDER | EXERCISE LEVEL | QUESTION |
|-----|-----|-----|-------|-------------|--------|----------------|----------|
| 1 | 26 | 23.2 | 0 | 61.0 | 0 | 1 | 1 |
| 2 | 30 | 30.2 | 9 | 65.5 | 1 | 3 | 2 |
| 3 | 32 | 28.9 | 17 | 59.6 | 1 | 3 | 4 |
| 4 | 37 | 22.4 | 1 | 68.4 | 1 | 2 | 3 |
| 5 | 33 | 25.5 | 7 | 64.5 | 0 | 3 | 5 |
| 6 | 29 | 22.3 | 1 | 70.2 | 0 | 2 | 2 |
| 7 | 32 | 23.0 | 0 | 67.3 | 0 | 1 | 1 |
| 8 | 33 | 26.3 | 1 | 72.8 | 0 | 3 | 1 |
| 9 | 32 | 22.2 | 3 | 71.5 | 0 | 1 | 4 |
| 10 | 33 | 29.1 | 5 | 63.2 | 1 | 1 | 4 |
| 11 | 26 | 20.8 | 2 | 69.1 | 0 | 1 | 3 |
| 12 | 34 | 20.9 | 4 | 73.6 | 0 | 2 | 3 |
| 13 | 31 | 36.3 | 1 | 66.3 | 0 | 2 | 5 |
| 14 | 31 | 36.4 | 0 | 66.9 | 1 | 1 | 5 |
| 15 | 27 | 28.6 | 2 | 70.2 | 1 | 2 | 2 |
| 16 | 36 | 27.5 | 2 | 68.5 | 1 | 3 | 3 |
| 17 | 35 | 25.6 | 143 | 67.8 | 1 | 3 | 4 |

# Co to the following address

https://raw.githubusercontent.com/amrrs/sample_revenue_dashboard_shiny/master/recommendation.csv

Account,Product,Region,Revenue--------------------→name of attributes

Axis Bank,FBB,North,2000

HSBC,FBB,South,30000

SBI,FBB,East,1000

ICICI,FBB,West,1000

Bandhan Bank,FBB,West,200

Axis Bank,SIMO,North,200

HSBC,SIMO,South,300

SBI,SIMO,East,100

ICICI,SIMO,West,100

Thanks for your attention