# Data analysis of diabetes and chronic kidney disease
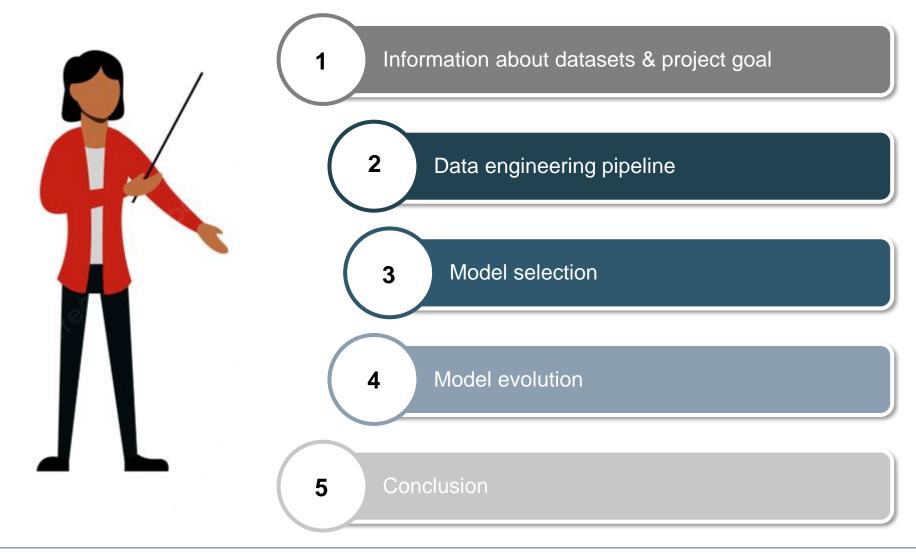
METHODS OF ADVANCED DATA ENGINEERING

Zeinab Aliakbari Mamaghani

Student ID: 23028078

# Outline
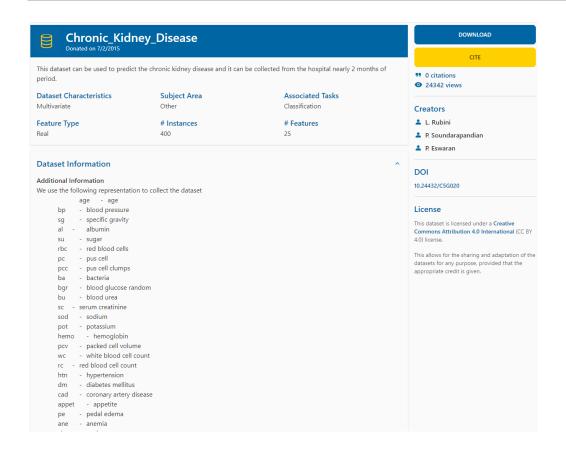
**FAU**

# How can I choose two suitable datasets?

- Open access

- Being interesting to me

- Accessible in common formats

- Possibility of relationship between two datasets



Thought bubbles: Weather data, Traffic Data, Public Transportation data, Crime data, Medical data, Shopping data

# Datasets

## Chronic_Kidney_Disease
Donated on 7/2/2015

This dataset can be used to predict the chronic kidney disease and it can be collected from the hospital nearly 2 months of period.

| Dataset Characteristics | Subject Area | Associated Tasks |
|---|---|---|
| Multivariate | Other | Classification |

| Feature Type | # Instances | # Features |
|---|---|---|
| Real | 400 | 25 |

### Dataset Information

**Additional Information**

We use the following representation to collect the dataset
```
        age    - age
bp      - blood pressure
sg      - specific gravity
al      -  albumin
su      - sugar
rbc     - red blood cells
pc      - pus cell
pcc     - pus cell clumps
ba      - bacteria
bgr     - blood glucose random
bu      - blood urea
sc      - serum creatinine
sod     - sodium
pot     - potassium
hemo     - hemoglobin
pcv     - packed cell volume
wc      - white blood cell count
rc      - red blood cell count
htn     - hypertension
dm      - diabetes mellitus
cad     - coronary artery disease
appet    - appetite
pe      - pedal edema
ane     - anemia
```

**DOWNLOAD**

**CITE**

❝ 0 citations
👁 24342 views

**Creators**
- L. Rubini
- P. Soundarapandian
- P. Eswaran

**DOI**
10.24432/C5G020

**License**
This dataset is licensed under a **Creative Commons Attribution 4.0 International** (CC BY 4.0) license.

This allows for the sharing and adaptation of the datasets for any purpose, provided that the appropriate credit is given.

---

## Pima Indians Diabetes Database

Predict the onset of diabetes based on diagnostic measures

Data Card    Code (2770)    Discussion (49)

### About Dataset

**Context**

This dataset is originally from the National Institute of Diabetes and Digestive and Kidney Diseases. The objective of the dataset is to diagnostically predict whether or not a patient has diabetes, based on certain diagnostic measurements included in the dataset. Several constraints were placed on the selection of these instances from a larger database. In particular, all patients here are females at least 21 years old of Pima Indian heritage.

**Content**

The datasets consists of several medical predictor variables and one target variable, `Outcome`. Predictor variables includes the number of pregnancies the patient has had, their BMI, insulin level, age, and so on.

**Acknowledgements**

Smith, J.W., Everhart, J.E., Dickson, W.C., Knowler, W.C., & Johannes, R.S. (1988). Using the ADAP learning algorithm to forecast the onset of diabetes mellitus. *In Proceedings of the Symposium on Computer Applications and Medical Care* (pp. 261--265). IEEE Computer Society Press.

**Inspiration**

Can you build a machine learning model to accurately predict whether or not the patients in the dataset have diabetes or not?

**Usability** ⓘ
8.82

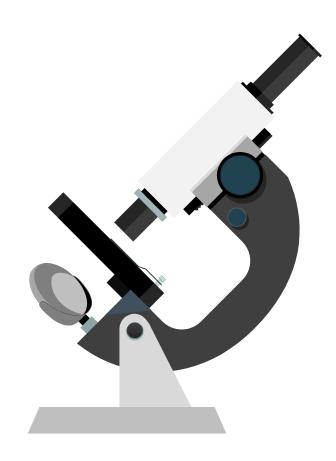**License**
CC0: Public Domain

**Expected update frequency**
Not specified

**Tags**
Earth and Nature    Health
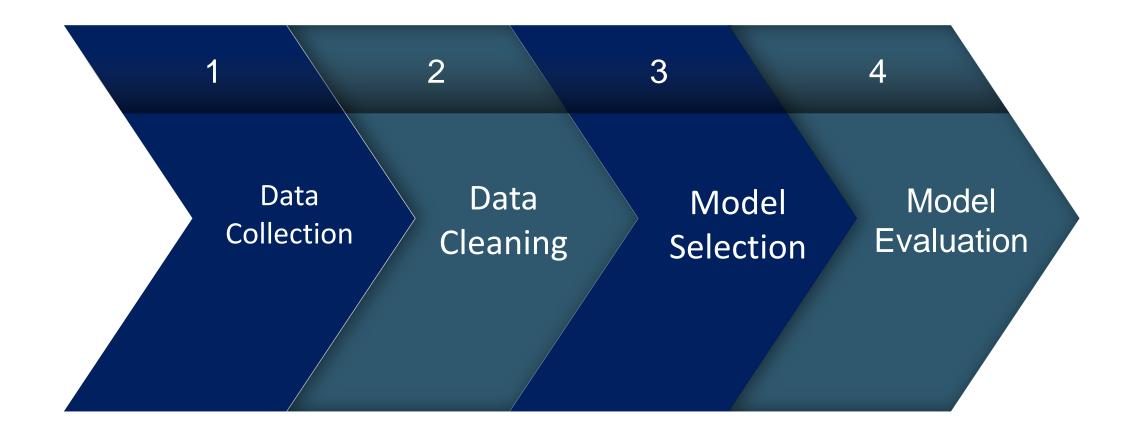Diabetes    India
Healthcare

**What factors affect diabetes and chronic kidney disease? And is there a relationship between these diseases?**

# Methods & Processes

# A closer look at the dataset : Data cleaning



**Chronic Kidney Disease.csv**

**What I did in Pipeline :**
- Abbreviation Replacement
- Categorical to Numerical
- Handling Null/Zero Values

**Pima Indians Diabetes.csv**

# A closer look at the dataset : Data cleaning

**What I did in Pipeline :**

- Remove outliers

**Kidney Disease Dataset :**
- ❖ Z-Score Transformation with a threshold of 3

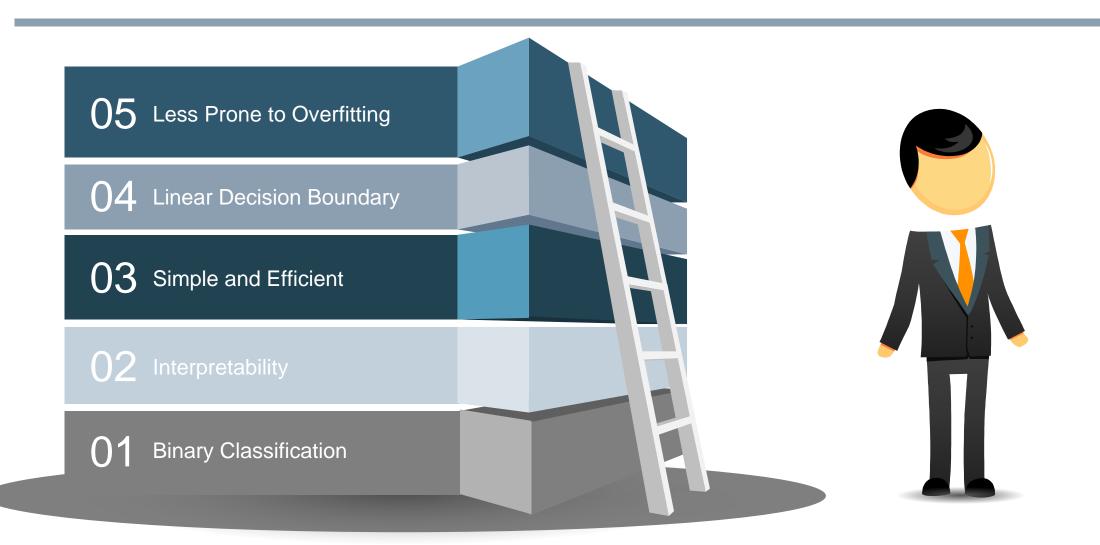**Diabetes Dataset:**
- ❖ Z-Score Transformation for features such as Glucose, Blood Pressure, SkinThickness, BMI with a threshold of 3-5.
- ❖ Interquartile Range (IQR) transformation for features such as Pregnancies, Insulin, DiabetesPedigreeFunction and Age with a threshold of 1.5-3.

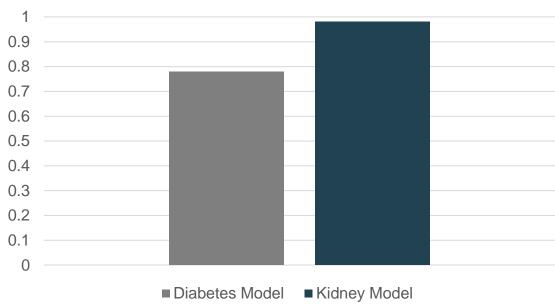# Selecting a prediction Model : Which model and why?

**05** Less Prone to Overfitting

**04** Linear Decision Boundary

**03** Simple and Efficient

**02** Interpretability

**01** Binary Classification

# Are the results satisfactory?

**Selected Model : Logistic Regression**

## Comparison of model accuracy



**kidney model accuracy: 0.98**

**Diabetes model accuracy : 0.78**

# What about the importance of features?



Kidney Feature Importance



Diabetes Feature Importance
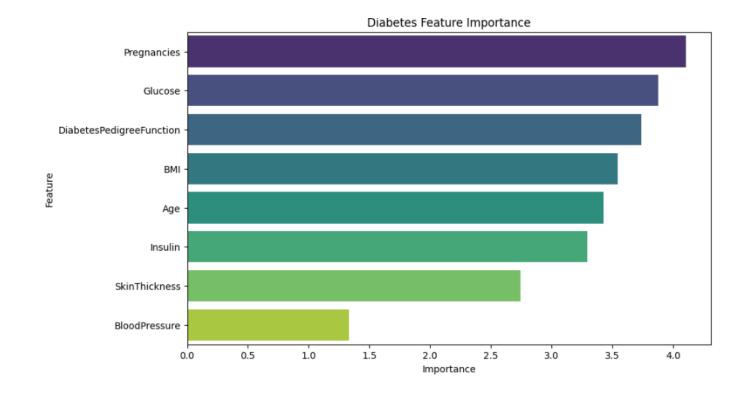
**Ensemble modeling
Or
Neural Network**

# New Model for Diabetes dataset

## Neural Network

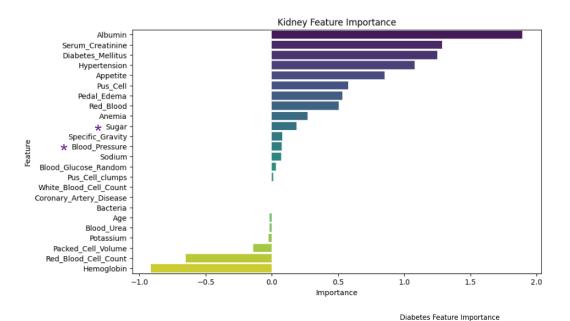Two hidden layers, the first with 8 neurons and the second with 4 neurons
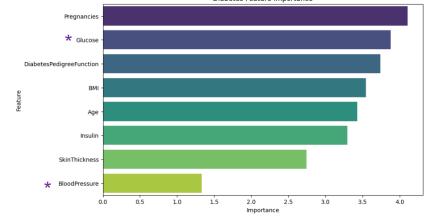
**Diabetes model accuracy:** 0.84



Diabetes Feature Importance

# Discussion & Conclusion

| Highlights | Diabetes Model | Kidney Model |
|---|---|---|
| **Selected Model** | Neural Network | Logistic Regression |
| **Accuracy** | 0.84 | 0.98 |
| **Important Features** | Number of Pregnancies* , Glucose | Albumin, Serum_Creatinine and *Diabetes_Mellitus |
| **Common Features** | Blood Sugar (Glucose) and Blood Pressure | |



Kidney Feature Importance



Diabetes Feature Importance

- The diabetes dataset only contains information about women.
- The kidney disease dataset don't have the gender feature.

* https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7153959
* https://www.niddk.nih.gov/health-information/diabetes/overview/preventing-problems/diabetic-kidney-disease
* https://www.cdc.gov/diabetes/managing/diabetes-kidney-disease.html

# References

- https://archive.ics.uci.edu/dataset/336/chronic+kidney+disease

- https://www.kaggle.com/datasets/uciml/pima-indians-diabetes-database/data

- https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7153959

- https://www.niddk.nih.gov/health-information/diabetes/overview/preventing-problems/diabetic-kidney-disease

- https://www.cdc.gov/diabetes/managing/diabetes-kidney-disease.html