# DATA MINING PROJECT

Submitted By:

Akhil (171210005)
Ankit (171210009)
Jayprakash Kumar (171210030)

# Objective

- **Predicting whether the person would show up on their medical appointment(s).**

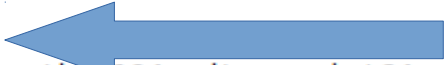**Dataset Used :** Medical Appointments No Show

**Link :** https://www.kaggle.com/joniarroba/noshowappointments

# Data Processing

## 1. Removing inconsistent values

- Negative values in age column
- Removing rows having "Appointment Day" before "Scheduled Day"



Negative values in "Age" Column
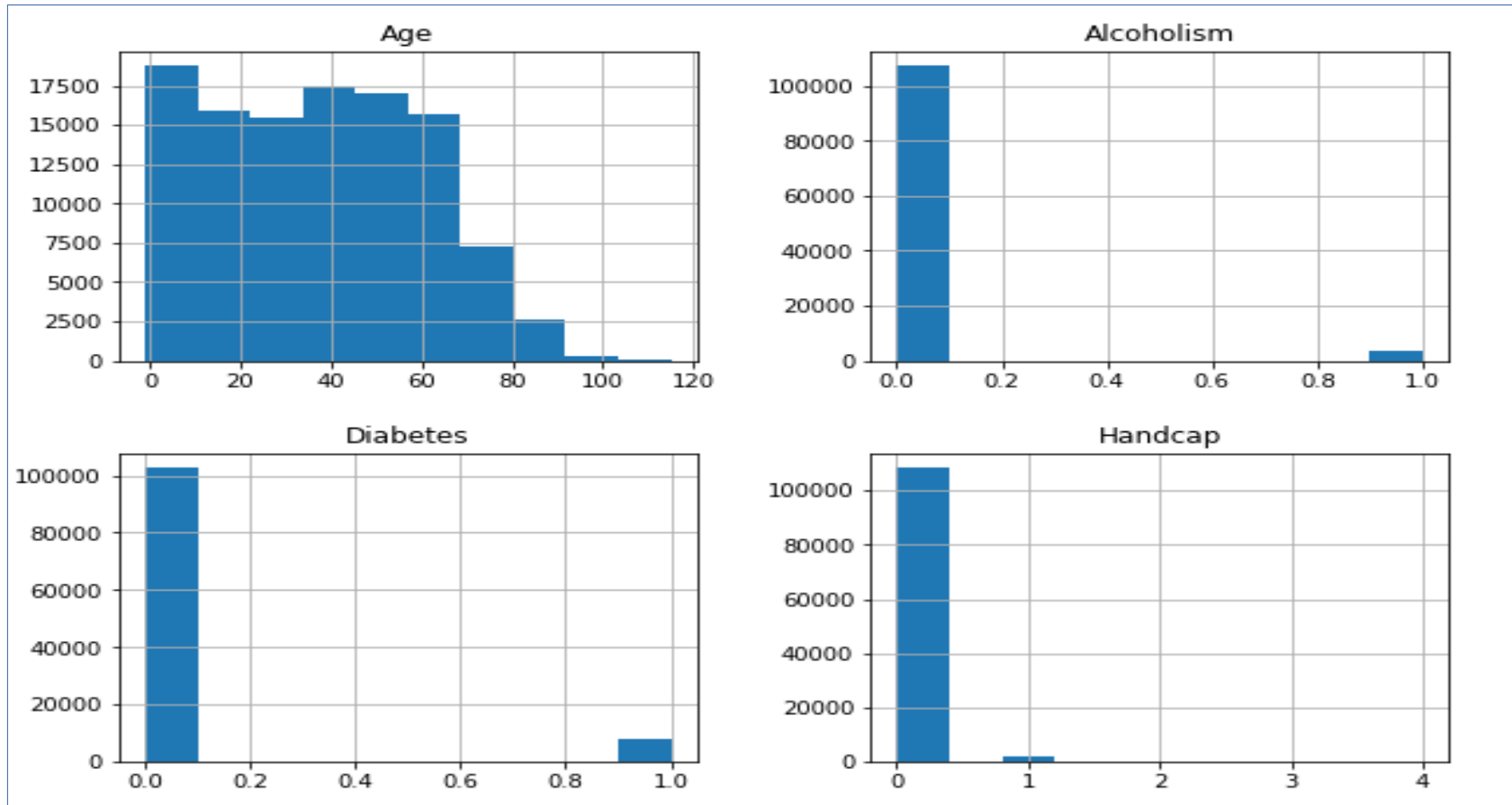
2. Encoding Genders, Neighbourhood and Days of the Week, No-Show

3. Adding a column Waiting Time(in days)

where Waiting Time =  Appointement Day – Scheduled Day

4. Removed unnecessary columns like "Patient ID", "Appointment ID", "Scheduled Day" and "Appointment Day"

# Visualising Data

To know the dependency of target variable on different features



Histograms for different features

# Choosing Optimal Algorithm

- As this was a Yes/No based classification problem
- It was best suited for a Decision Tree based classification.
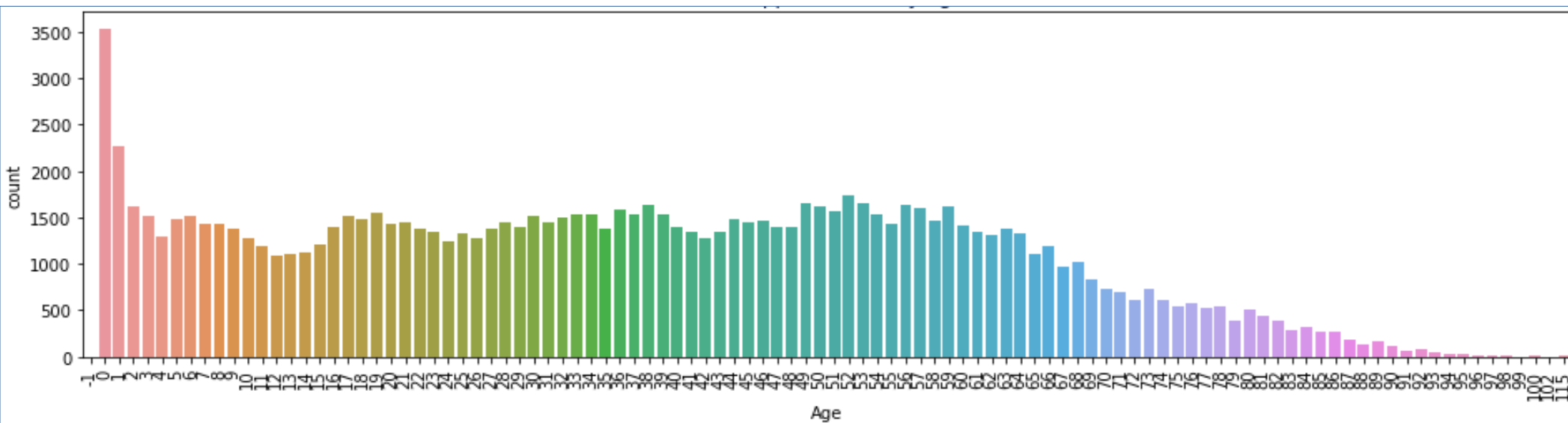
```
Feature ranking:
   1.   feature: Age (0.23031279358109608)
   2.   feature: Scholarship (0.11661695625956661)
   3.   feature: Hypertension (0.03158117856512685)
   4.   feature: Diabetes (0.0196827458256893 9)
   5.   feature: Alcoholism (0.01753557615127744)
   6.   feature: SMS_received (0.01702952934783208 4)
   7.   feature: WaitingTime (0.01028657055282368)
   8.   feature: AppointmentDayOfWeek (0.009488809077866552)
   9.   feature: Handicap_0 (0.009204400671004856)
   10.   feature: Handicap_1 (0.00892060077705996)
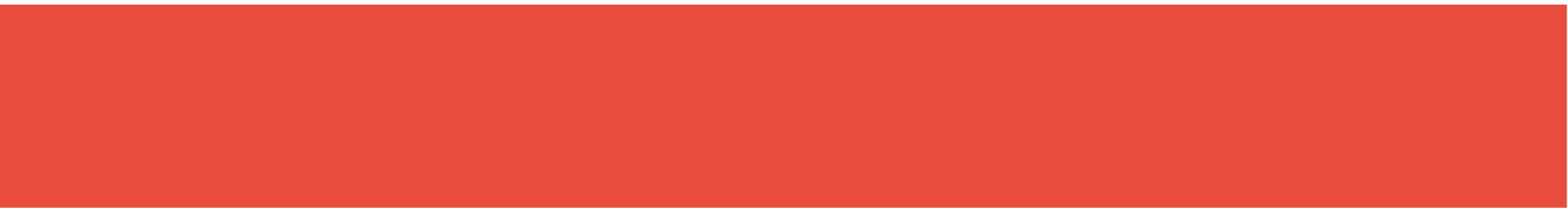```

Ranking of features based on their importances

# Conclusion

After obtaining the Decision Tree , we came to following conclusions:

- Showing Up of person depends a lot on their "Age"
- And it depends the least on the place where they live
  i.e Neighbourhood



Histogram Plot of Age of Persons

# THANK YOU