# DA421 Assignment 1

Name: Varun Nagpal
Roll No: 210150020

# Implementation

1. Data Preprocessing
   a. Missing value imputation
   b. Categorical feature encoding, Standardization and min-max normalization
   c. Train test splits for preventionN algorithm
   d. Data sampling for census income

2. Custom Distance function for KNN
   a. Interval-Scaled Features
   b. Nominal Features
   c. Ordinal Features

3. DiscoveryN Algorithm
4. PreventionN Algorithm

# Experimental Setup

- Protected Attributes
  - Adult Dataset: Race (non-whites)
  - Census-Income Dataset: Race (non-whites) and Marital Status (Divorced, Separated and Widowed).

- DiscoveryN algorithm
  - Decision Tree classifier for 'disc' label on the t-labelled dataset
  - Analyzed performance metrics such as accuracy, precision, recall, F1-score.                    (8
    split on the t-labelled dataset).

- PreventionN algorithm
  - Tested on original and t-corrected data (t = 0.1) using the original 66.66:33.33 splits..
  - Classifiers used: Decision Tree, Naive Bayes, Logistic Regression
  - Metrics reported: Accuracy, t = 0.10 discrimination, Classifier discrimination on predictions.
  - Additional experiment on 0.05-corrected data: Accuracy and Classifier discrimination measured.

# Results

## Discrimination Discovery

| | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|
| Adult (Race) | 89.12% | 48.73% | 51.33% | 50.00% |
| Census-Income (Race) | 92.09% | 32.39% | 29.87% | 31.08% |
| Census-Income (Marital Status) | 89.63% | 22.05% | 20.54% | 21.27% |

# Results

## Research Paper Results (Adult Dataset)

| | No preprocessing | | 0.1 correction | |
|---|---|---|---|---|
| Classifier | Accuracy | 0.1 disc | Accuracy | 0.1 disc |
| Decision Tree | 85.60% | 4.24% | 84.94% | 1.07% |
| Naive Bayes | 82.46% | 4.06% | 82.33% | 2.23% |
| Logistic Regression | 85.28% | 6.61% | 84.70% | 0.61% |

# Results

## Our Implementation (Adult Dataset)

| Classifier | No preprocessing | | | 0.1 correction | | | 0.05 correction | |
|---|---|---|---|---|---|---|---|---|
| | Accuracy | Classifier disc | 0.1 disc | Accuracy | Classifier disc | 0.1 disc | Accuracy | Classifier disc |
| Decision Tree | 81.90% | 7.82% | 4.72% | 81.21% | 6.82% | 1.92% | 79.75% | 3.85% |
| Naive Bayes | 79.31% | 15.34% | 6.39% | 77.32% | 12.58% | 3.38% | 75.15% | 7.59% |
| Logistic Regression | 85.20% | 8.46% | 5.26% | 85.03% | 6.65% | 1.16% | 84.71% | 4.20% |

# Results

## Our Implementation (Census Dataset, Race Sensitive Attribute)

| | No preprocessing | | | 0.1 correction | | | 0.05 correction | |
|---|---|---|---|---|---|---|---|---|
| Classifier | Accuracy | Classifier disc | 0.1 disc | Accuracy | Classifier disc | 0.1 disc | Accuracy | Classifier disc |
| Decision Tree | 92.53% | 4.06% | 1.06% | 91.70% | 1.72% | 0.32% | 90.07% | -1.07% |
| Naive Bayes | 88.82% | 11.71% | 1.41% | 87.90% | 9.46% | 0.86% | 83.44% | 7.56% |
| Logistic Regression | 94.86% | 2.08% | 1.04% | 93.46% | -1.79% | 0.79% | 93.44% | -9.62% |

# Results

## Our Implementation (Census Dataset, Marital Status Sensitive Attribute)

| Classifier | No preprocessing | | | 0.1 correction | | | 0.05 correction | |
|---|---|---|---|---|---|---|---|---|
| | Accuracy | Classifier disc | 0.1 disc | Accuracy | Classifier disc | 0.1 disc | Accuracy | Classifier disc |
| Decision Tree | 92.46% | 3.07% | 1.16% | 91.13% | 1.86% | 0.36% | 90.69% | -3.85% |
| Naive Bayes | 88.82% | 8.34% | 1.24% | 87.30% | 7.58% | 0.28% | 82.72% | 6.59% |
| Logistic Regression | 94.86% | 0.93% | 1.23% | 93.70% | -1.65% | 0.03% | 90.53% | -3.20% |