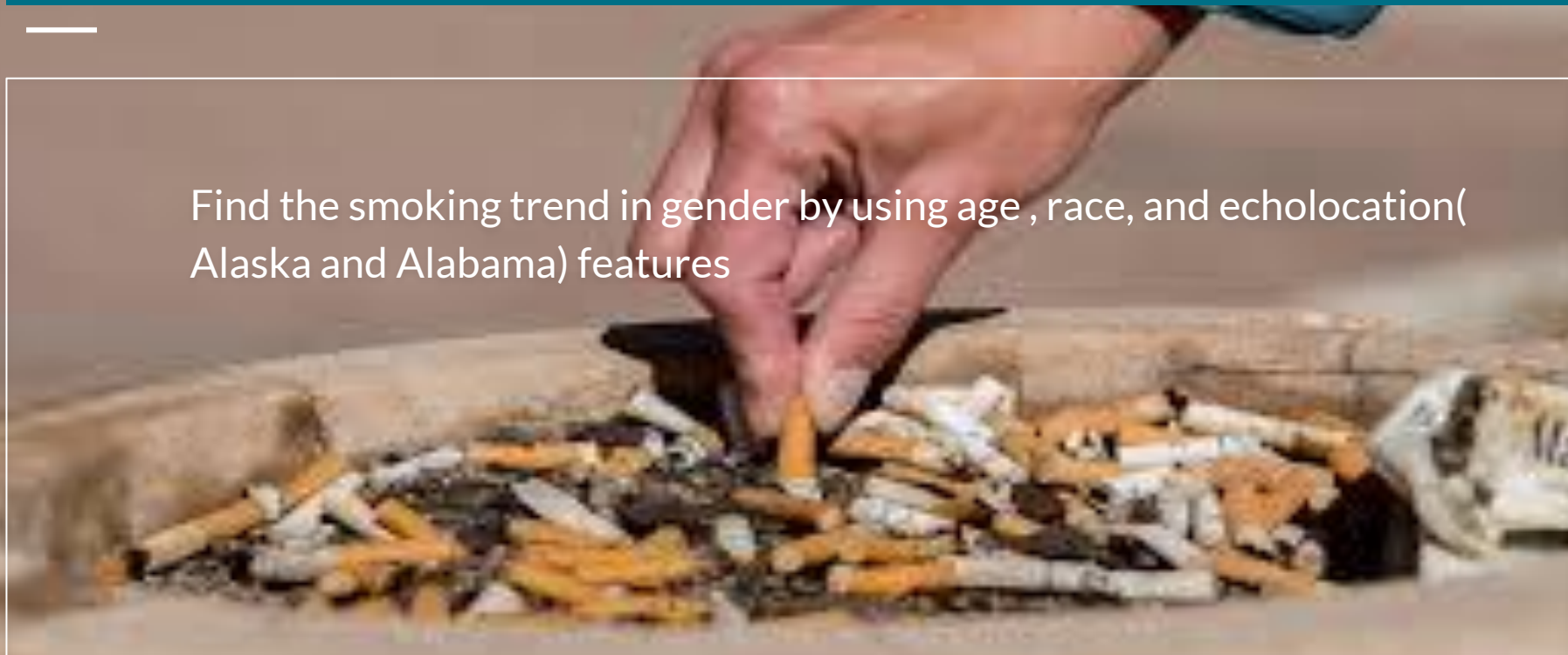# Prediction of Smoking Trend in Gender.

By Collin Tully, Sahrish Afzal, Eric Cacdac,Rucha Soni
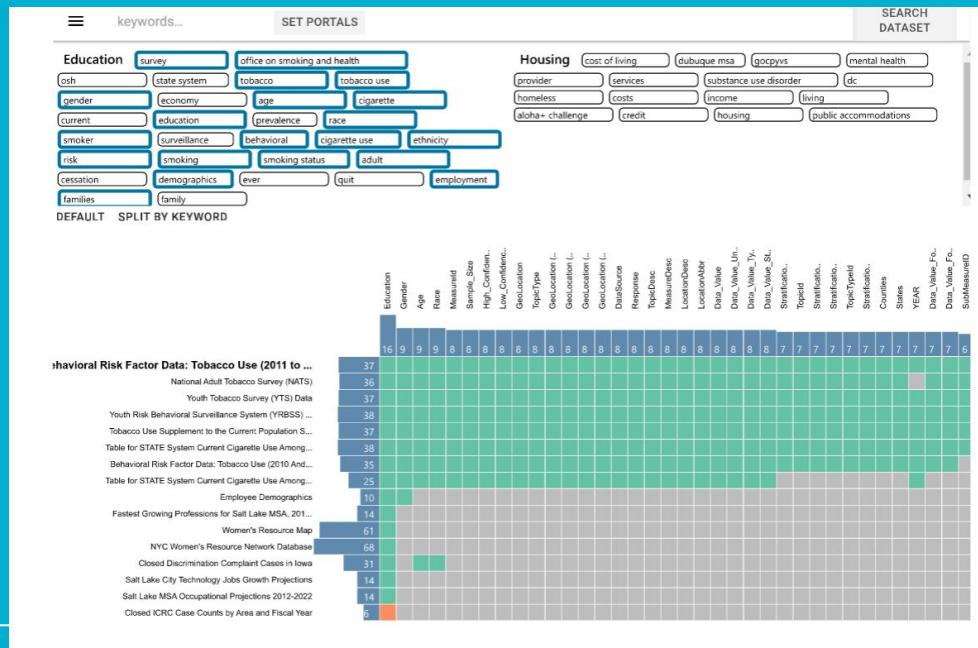
# Objective of the Analysis

Find the smoking trend in gender by using age , race, and echolocation( Alaska and Alabama) features

# The Dataset

## Used data set:

- WE have used the Behavioral Risk Factor Data by using shown portals
- That gave us data with 30 features
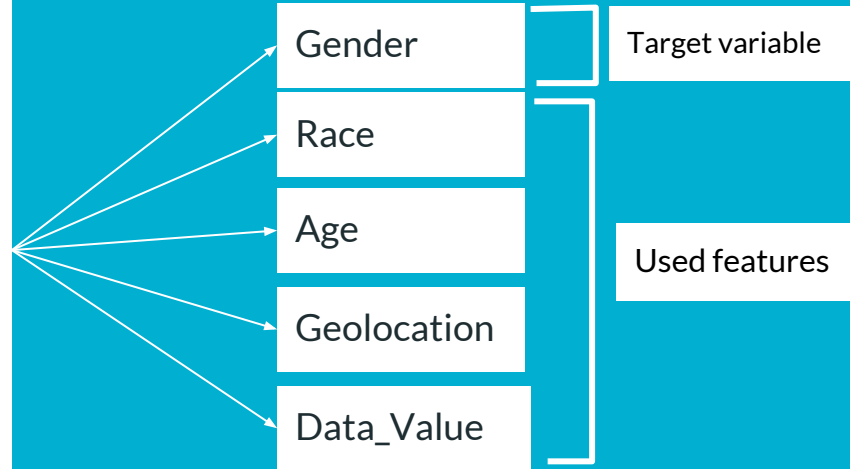
# Features of the Dataset

## Original Data

| | |
|---|---|
| 1 | LOCATIONABBR |
| 2 | MEASUREDESC |
| 3 | RESPONSE |
| 4 | DATA_VALUE_UNIT |
| 5 | DATA_VALUE_TYPE |
| 6 | DATA_VALUE_FOOTNOTE_SYMBOL |
| 7 | DATA_VALUE_FOOTNOTE |
| 8 | DATA_VALUE_STD_ERR |
| 9 | LOW_CONFIDENCE_LIMIT |
| 10 | HIGH_CONFIDENCE_LIMIT |
| 11 | TOPICTYPEID |
| 12 | STRATIFICATIONID2 |
| 13 | STRATIFICATIONID3 |
| 14 | STRATIFICATIONID4 |
| 15 | SUBMEASUREID |
| 16 | YEAR |
| 17 | LOCATIONDESC |
| 18 | TOPICTYPE |
| 19 | TOPICDESC |
| 20 | DATASOURCE |
| 21 | DATA_VALUE |
| 22 | SAMPLE_SIZE |
| 23 | GENDER |
| 24 | RACE |
| 25 | AGE |
| 26 | EDUCATION |
| 27 | GEOLOCATION |
| 28 | TOPICID |
| 29 | MEASUREID |
| 30 | DISPLAYORDER |

## Present Data

| | |
|---|---|
| 16 | YEAR |
| 17 | LOCATIONDESC |
| 18 | TOPICTYPE |
| 19 | TOPICDESC |
| 20 | DATASOURCE |
| 21 | DATA_VALUE |
| 22 | SAMPLE_SIZE |
| 23 | GENDER |
| 24 | RACE |
| 25 | AGE |
| 26 | EDUCATION |
| 27 | GEOLOCATION |
| 28 | TOPICID |
| 29 | MEASUREID |
| 30 | DISPLAYORDER |

## Focus Data

Gender — Target variable

Race

Age

Geolocation

Data_Value

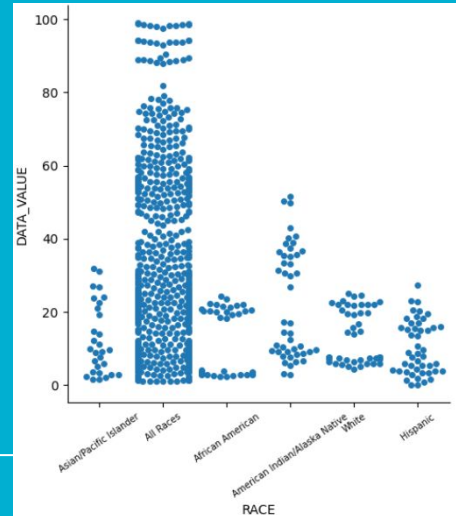Used features

# Graph: Feature Vs Target

# First Model: Linear Regression

Purpose: To determine if there is a relationship between the features and the target variable.
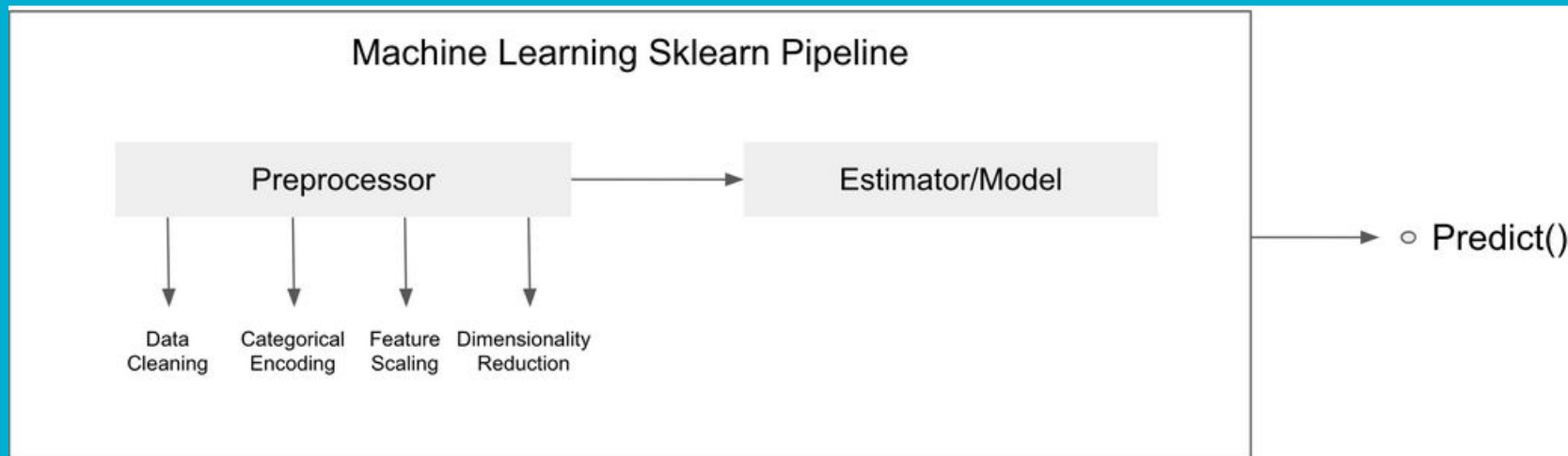
Features Used:  Race, Age, Location, Probability of Smoking.
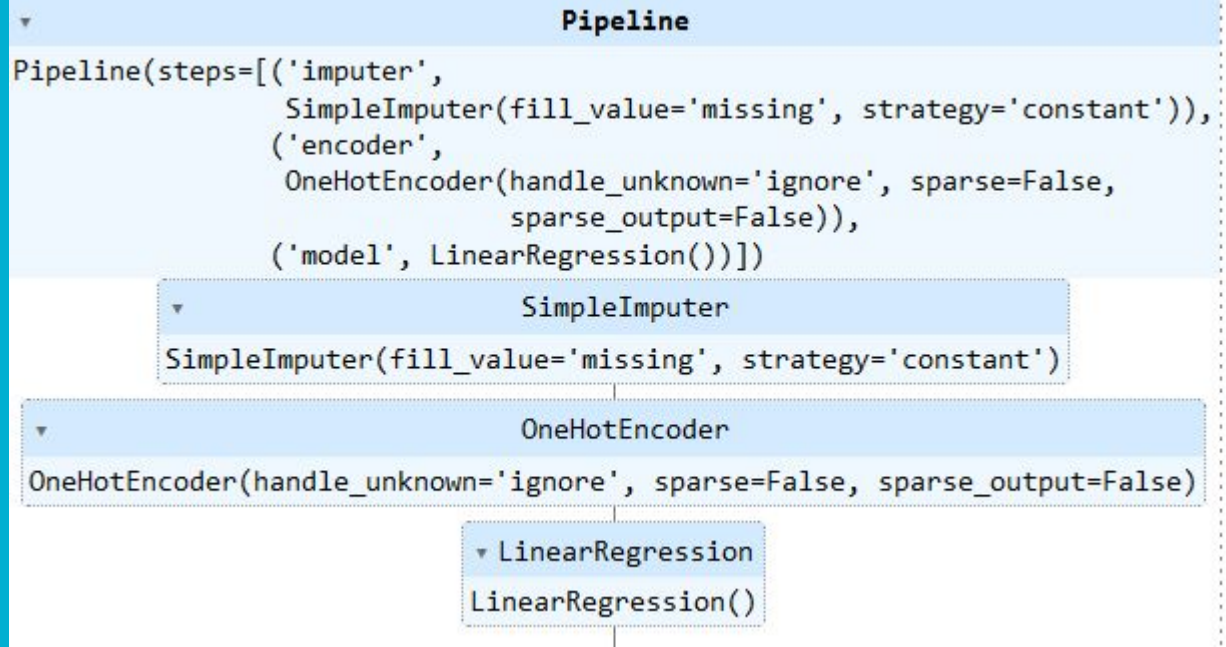
Target Variable: Sex

# First Mode: Linear Regression Encoding

- Binary Encoding
- Encoding using Sklearn Pipeline

# Evaluation of First Model: Linear Regression

R2 = 0.728

# Issues With the First Model: Linear Regression

- Encoding of categorical target variable. Some of the data was mixed with male and female.
- Issue with collinearity
- VIF: 1 lack of collinearity,
- VIF: >5 correlation between predictor
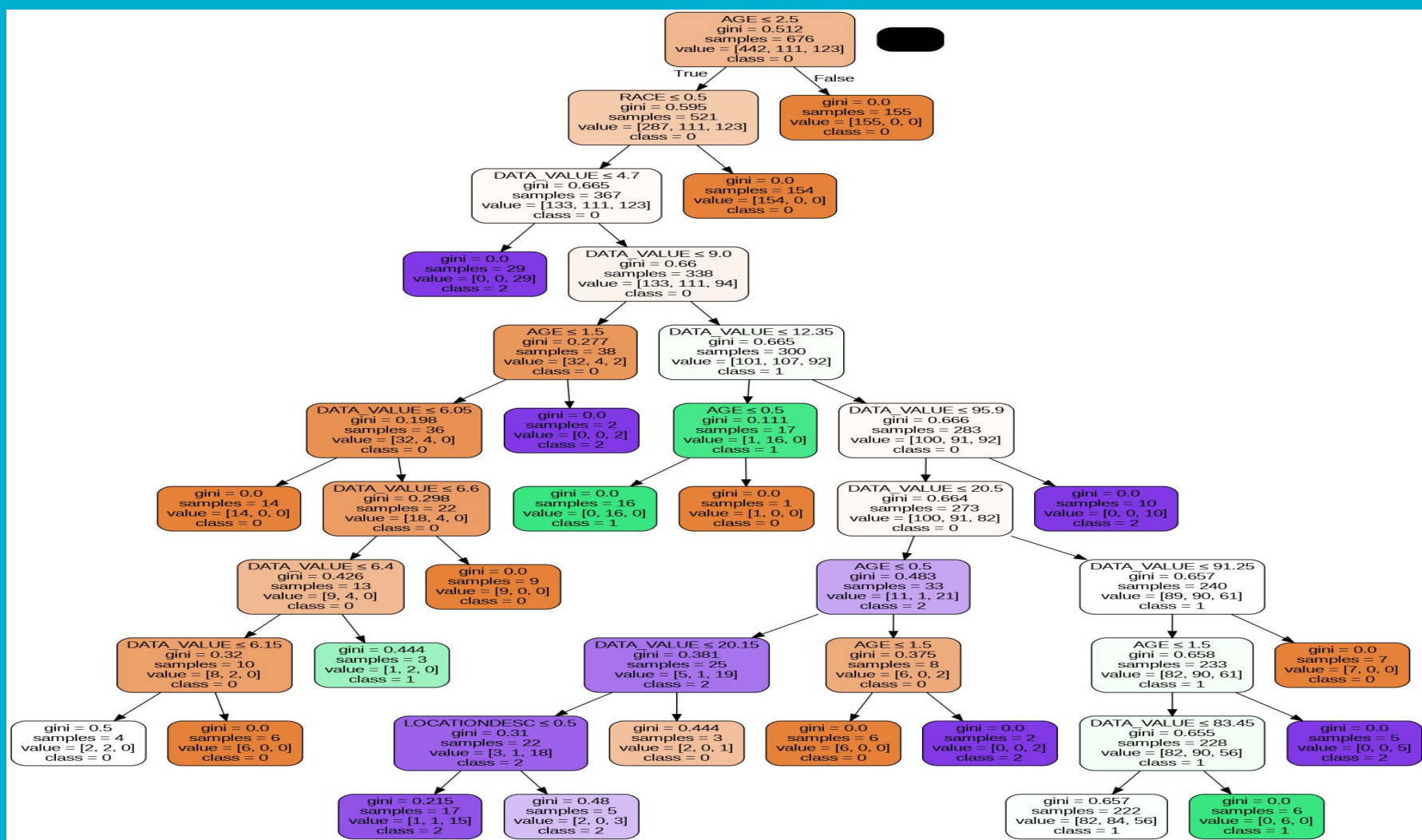- Variance Inflation Factor

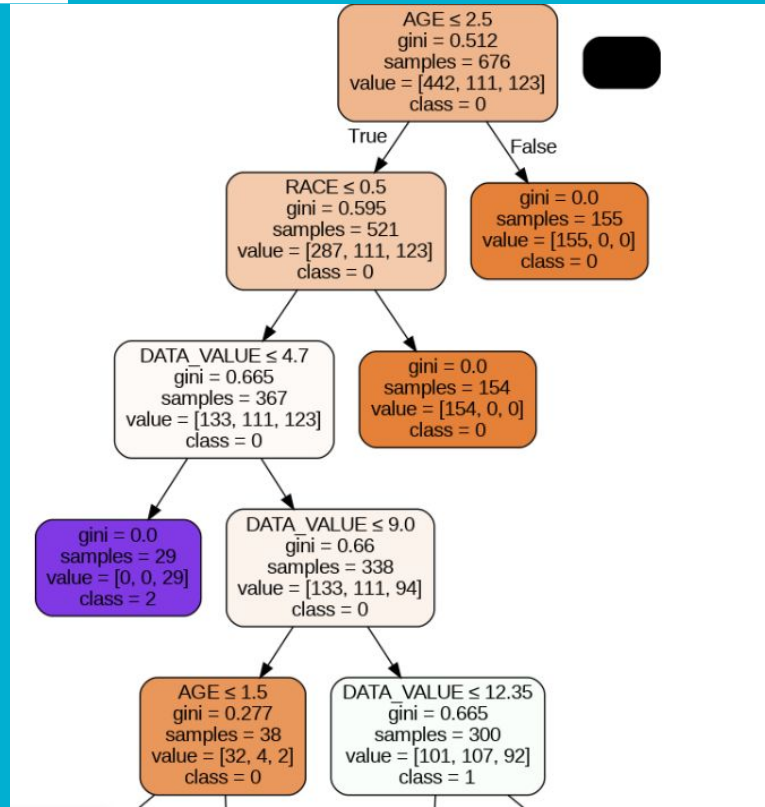|    | variable | VIF |
|----|----------|-----|
| 0 | Intercept | 63.039294 |
| 1 | RACE[T.All Races] | 4.655352 |
| 2 | RACE[T.American Indian/Alaska Native] | 1.993132 |
| 3 | RACE[T.Asian/Pacific Islander] | 1.649183 |
| 4 | RACE[T.Hispanic] | 2.018633 |
| 5 | RACE[T.White] | 2.040517 |
| 6 | LOCATIONDESC[T.Alaska] | 1.006983 |
| 7 | AGE[T.18 to 44 Years] | 1.946569 |
| 8 | AGE[T.25 to 44 Years] | 1.911454 |
| 9 | AGE[T.45 to 64 Years] | 1.910976 |
| 10 | AGE[T.65 Years and Older] | 1.913885 |
| 11 | AGE[T.Age 20 and Older] | 3.584818 |
| 12 | AGE[T.Age 25 and Older] | 3.584532 |
| 13 | AGE[T.All Ages] | 8.405993 |
| 14 | DATA_VALUE | 1.469572 |

# Second Model: Decision Tree

Objective: Predict the target variable (gender) based on the independent variables (features: race, age, location, data value)

Assigned all feature values (non-numerical) to a numerical value:

- Gender: Male (0), Female (1), Overall (2)
- Race: All Races (0), Asian/Pacific Islander (1), American Indian/Alaskan Native(2), African American (3), White (4), Hispanic (5)
- Age: All Ages (0), 18 to 24 (1), 18 to 44 (2), 20+ (3), 25+ (4), 25 to 44 (5), 45 to 64 (6), 65+ (7)
- Location: Alabama (0), Alaska (1)

# Evaluation of Second Model: Decision Tree



- Accuracy score: 0.742268
- Each node is assigned a class (gender) based on whether they fit the criteria of the feature
- Gini is the impurity score (rate of samples that do not match the assigned class)
- R^2 score: 0.482

# Comparisons Between the Models
Decision Tree R^2 = 0.482    Linear Regression R^2=0.728

## Advantages of Decision Tree

- Decision Tree can capture complex relationships in the data
- It does not require much data preprocessing for feature engineering
- Decision tree can handle both categorical and numerical data

## Disadvantages of Decision Tree

- Prone to overfitting, especially with complex trees
- Decision tree is computationally expensive with large datasets
- It can be unstable and sensitive to small changes in the data

## Advantages of Linear Regression

- Linear Regression can provide insight into the relationship between the independent and dependent variables
- It is Computationally efficient and can handle large datasets
- Linear Regression model is simple and easy to understand model

## Disadvantages of Linear Regression

- Assumes a linear relationship between independent and dependent variables
- Can be affected by multicollinearity among independent variables
- This model cannot capture complex relationships in the data

# Possible Flaws With Data Itself

- Encoding with target variable with data that was unspecific. Some of the surveys marked male and female.
- Age bounds of some surveys overlapped.
- Difference in R2 Models