# Insightful Identity Analysis: Detecting Age, Gender, and Ethnicity

Team Members:
Sarvagya Kaushik
Kunal Sharma
Vansh

**IIID**

INDRAPRASTHA INSTITUTE *of*
INFORMATION TECHNOLOGY
**DELHI**

# Motivation

1. For those working in the fields of computer vision and facial recognition, the UTKFace dataset is a helpful resource.
2. The unique characteristics of this dataset and its potential applications across numerous domains serve as the inspiration for training on it.

Reasons UTKFace is revolutionary in the field of face data sets:
1. Age Diversity
2. Real-World Challenges
3. Annotations for Gender and Ethnicity
4. Multi-Task Learning

Benchmark:
1. Contribute to the advancement of state-of-the-art techniques in facial analysis.
2. Compare our algorithm against existing state-of-the-art methods using the UTKFace dataset.
3. Understand how your solutions measure up in terms of accuracy, efficiency, and robustness.

# Literature Survey

## 1. GRA_Net:

- Consists of multiple layers, each containing an attention block.

- Each attention block combines features from the previous layer with attention weights to produce refined feature representation.

- Gating mechanism dynamically controls the influence of attention on the feature at each layer.

- GRA_Net is trained using standard deep learning techniques, such as backpropagation and gradient descent.

The classification accuracies achieved by the proposed GRA_Net model for UTKFace datasets was found to be 99.2%.



Thorough comparison with various alternative models, GRA_Net has unequivocally demonstrated its supremacy by consistently yielding superior results.

| Model | Gender(%) | Age(%) |
|---|---|---|
| Facenet | 91.2 | 56.9 |
| Finetuned Facanet (FFNet) | 96.1 | 64 |
| MTCNN | 98.23 | 70.1 |
| RAN (Wang et al. (2017)) | 97.5 | 85.4 |
| **Proposed model** | **99.2** | **93.7** |

# Literature Survey

## 2. Feature Extraction based Face Recognition, Gender and Age Classification algorithm

- The algorithm yields good results with small training data.

- Steps involved:
  - Preprocessing:
    - Color Conversion
    - Noise Reduction
    - Edge detection

  - Feature Extraction:
    - Computation: Ratios* are calculated.
    - Gender Classification: Naive Bayes

  - Training on the dataset
    - Artificial Neural Network(ANN): carried out in two parts:
      - Feed-forward path
      - Feedback path
    - Back Propagation

**Performance Analysis**

| Gender | Sample size | Correctly Labeled(CL) | Correct Rate(CR) | Total CR |
|--------|-------------|-----------------------|------------------|----------|
| Male   | 40          | 38                    | 95%              | 94.82%   |
| Female | 18          | 17                    | 94.44%           |          |

| Algorithm | AG | Sample size | CL | CR | Total CR |
|-----------|----|-------------|----|----|----------|
| FEBFRGAC  | Y  | 28          | 25 | 89.3% | 89.65% |
|           | M  | 20          | 18 | 90%   |        |
|           | O  | 10          | 09 | 90%   |        |
| CAGBFF    | Y  | 44          | 37 | 84.4% | 78.49% |
|           | M  | 32          | 25 | 78.1% |        |
|           | O  | 17          | 11 | 64.7% |        |

*Ratios that were taken into account were left-to-right eye distance upon eye-to-nose distance, left-to-right eye distance upon eye-to-lip distance, eye-to-nose distance upon eye to chin distance, and eye-to-nose distance upon eye-to-lip distance.

# Dataset Description

The UTKFace dataset is a fairly large dataset with over 20,000 face images with annotations of age, gender and ethinicity.

The subjects covered in the dataset consisted of people ranging from the age of 0 to 116 years old, over 4 ethnicities.

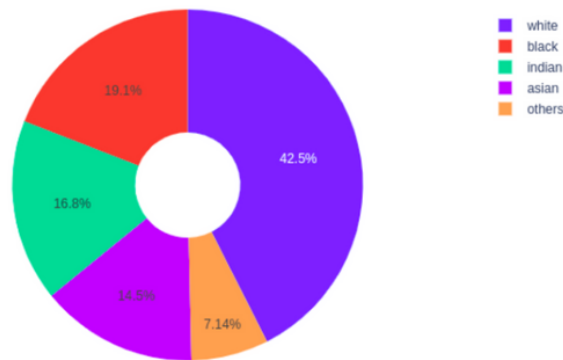Use Cases: Face detection, Age estimation, Age progression/regression,Landmark localization, etc.



Figure 12. Race Distribution

- white
- black
- indian
- asian
- others

19.1%
42.5%
16.8%
14.5%
7.14%

- Data was collected from a wide number of sources across the internet
-  Model may have slightly higher bias with context to race
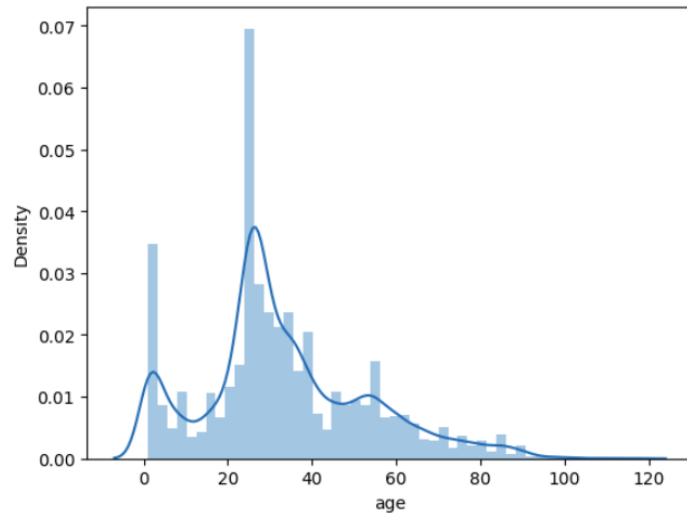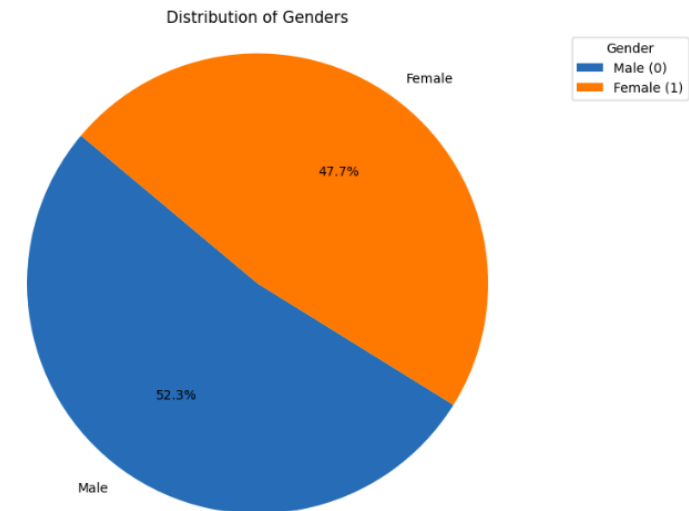
# Dataset Description


Figure 14. Age distribution plot

- The plot shows us that the data is normally distributed.
- Surface Observation shows us that the data is skewed to left ie most people in hf the dataset are less than 40 years old


Distribution of Genders

Gender
- Male (0)
- Female (1)

47.7%

52.3%

- Data spread was fairly symmetrical with regard to gender
- Not from IIIT D !

# Pre-processing

- **Effective Preprocessing is critical for ML tasks.**

- **Extracted features from image path name.**

- **Resizing is essential for overall quality and adaptability: Resized from 128 x 128 pixels to 28 x 28 pixels**

- **Converted from RGB Scale to grayscale to reduce data complexity and processing resource requirements**

- **Normalized the pixel values.**

# Methodology

## Models Used:-

- Logistic Regression for gender prediction

- K-Nearest Neighbours for gender prediction

# Logistic Regression

- **Splitting the dataset** into training and test set. For example, we can have 70% data for training, 20% for validation and the rest for testing.

- Creation of model

- $P(Y = 1) = 1 / 1 + e-(\mathbf{\beta}0+\mathbf{\beta}1X1+\mathbf{\beta}2X2+...+\mathbf{\beta}pXp)$, where $P(Y = 1)$ is the probability of the event occurring and $\mathbf{\beta}0, \mathbf{\beta}1, \mathbf{\beta}2, . . . , \mathbf{\beta}p$ are coefficients that represent the relationship between the independent variables X1, X2, . . . , Xp and the probability of the event.

- Model Performance

- Tuning hyperparameters and regularization

- Gender prediction
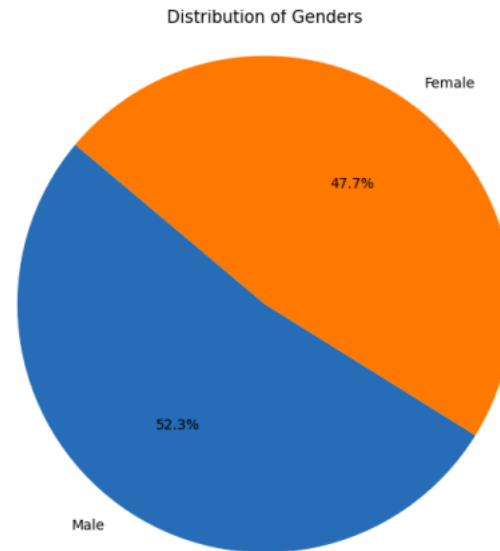
# K-Nearest Neighbours

- **Splitting the dataset** into training and test set. For example, we can have 70% data for training, 20% for validation and the rest for testing.

- **Creation of a classifier model** based on the KNN algorithm which is a non-parametric, instance-based algorithm that classifies data points based on their similarity to the k-nearest neighbors in the training data.

- Model training

- Model performance

- Tuning hyperparameters and regularization

- Gender prediction

# Results and Analysis

## Gender Distribution



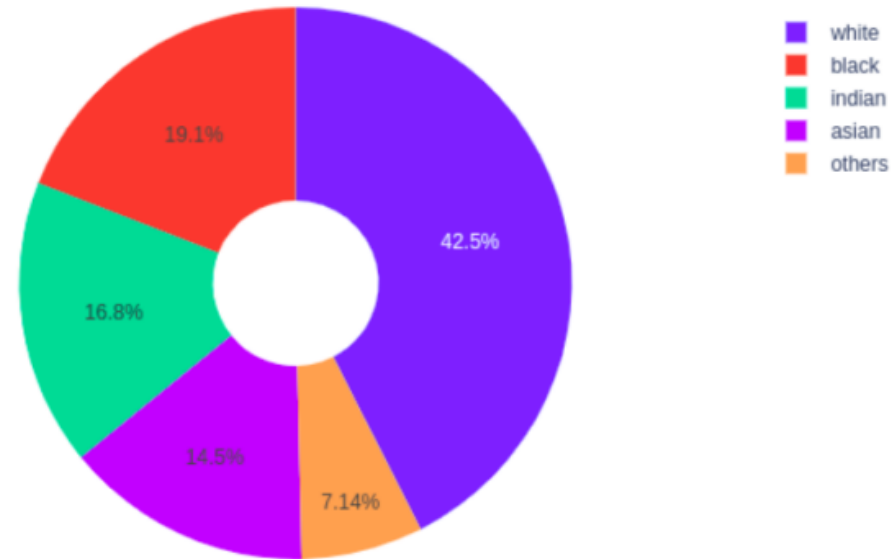Distribution of Genders

Female 47.7%

Male 52.3%

 The above figure gives a visualization of gender distribution. We can see that the percentage of the male population is slightly greater than females but the difference is minor. It's not capable of creating high bias.
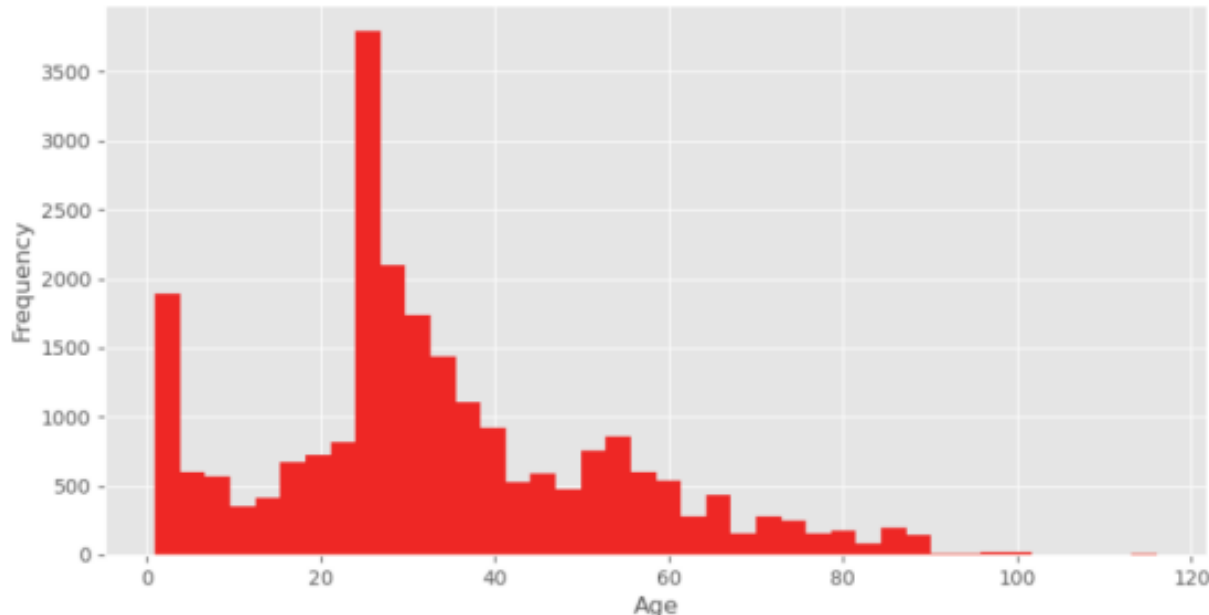
# Results and Analysis

## Ethnicity Breakdown



 The above figure gives visualization for ethnicity distribution in our data set. Our data set majorly consists of images of white ethnicity with 42.5 percent. It is followed by black at 19.1 percent, Indian with 16.8 percent and Asian with 14.5 percent. Rest of the population are categorized by others.

# Results and Analysis

## Age Distribution



 The above figure gives visualization for age distribution in our data set. From surface observation, we can see that the data is skewed to the left. Thus our data set mainly consists of a population less than 40 years. From the figure we can also see that the data is normally distributed.

# Results and Analysis

## Performance Metrics

The model was trained using a batch size of 32, binary cross-entropy as the loss function, and stochastic gradient descent (SGD) as the optimization algorithm. After 10 epochs of training, we achieved an accuracy of 80

The following statistics summarize the model's performance:

- **Training Loss: 0.3654**
- **Test Loss: 0.3598**
- **Test Accuracy: 84.41**

These results indicate that our logistic regression model performs well in classifying images into male and female categories. The relatively low training and test losses suggest that the model effectively minimized the classification error, and the test accuracy of 84.41 percent demonstrates its ability to correctly classify the gender of previously unseen images

# Results and Analysis

## Performance Metrics

The model was trained by flattening the image dimensions into one dimension and setting the k parameter to 20.

The following statistics and the classification report summarize the model's performance:

- **Accuracy: 0.7344**

| Class | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| 0 (Male) | 0.70 | 0.85 | 0.77 | 2468 |
| 1 (Female) | 0.79 | 0.61 | 0.69 | 2273 |
| Accuracy | | | 0.73 | 4741 |
| Macro Avg | 0.75 | 0.73 | 0.73 | 4741 |
| Weighted Avg | 0.74 | 0.73 | 0.73 | 4741 |

1. True Positive: 2095
2. False Negative: 373
3. False Positive: 886
4. True Negative: 1387

# Model Performance

## Logistic Regression

- The logistic regression model achieved an accuracy of 80 percent after 10 epochs of training. Its relatively low training and test losses, along with an accuracy of 84.41 percent, demonstrated its effectiveness in classifying gender from images. However, further analysis and potential fine-tuning are expected to enhance the model's performance.

## K-Nearest Neighbours

- The k-NN model, with 'k' set to 20, provided an accuracy of 73.44 percent. This model showed promise in gender classification, with distinct strengths in precision and recall for different gender categories. Future work may involve optimizing the choice of 'k' and exploring additional feature engineering techniques.

# Conclusion

The study explored gender classification using image data with logistic regression and k-Nearest Neighbors (k-NN). Logistic regression achieved 80% accuracy with potential for improvement, while k-NN achieved 73.44% accuracy, showing promise with precision and recall.

Our study and literature survey suggests further research into models like CNNs, SVMs, and decision trees for better performance. Overall, while logistic regression and k-NN have potential, ongoing research is needed to enhance gender classification from images, address biases, and ensure fairness in practical applications.

# Timeline Progress

| August 30 – September 30 | September 30 – October 15 | October 15 – November 15 | November 15 – End sems |
|---|---|---|---|
| **Data Analysis** | **Data Cleaning** | **Model Preparation** | **Optimization and Finalization** |
| • Visualizations<br>• Insights | • RGB to Grayscale<br>• Rescaling<br>• Normalisation | • Logistic Regression<br>• K-Nearest Neighbours<br>• More Models | • Code cleanup<br>• Optimization<br>• Final adjustments |

# Individual Contribution

- Sarvagya Kaushik : Training and Testing Models, Data Visualisation, Exploratory Data Analysis, Report preparation
- Kunal Sharma : Literature Survey, Preprocessing, Hyperparameter tuning, Report preparation
- Vansh : Training and Testing Models, Presentation preparation, Preprocessing, Report preparation