



Machine Learning Assignment

PROJECT REPORT

<TEAM ID : 16>

<PROJECT TITLE: Insincere Question Classification on Quora>

| Name | SRN |
|------------------|---------------|
| CHETAN NADICHAGI | PES2UG23CS149 |
| DIVYA J | PES2UG24CS810 |

Problem Statement

Online platforms such as Quora allow users to ask and answer questions freely. However, some questions are insincere or toxic, intended to provoke, mislead, or offend rather than seek genuine information.

These toxic questions negatively affect the user experience and platform moderation.

Our Solution

Our solution uses Natural Language Processing (NLP) and Machine Learning to automatically detect toxic or insincere questions on Quora.

The system cleans and processes text, converts it into numerical features using TF-IDF, and classifies questions using models like Logistic Regression and Random Forest.

This helps identify and filter harmful content to maintain a respectful platform environment.

Objective / Aim

Primary Objective:

To develop a text classification model capable of identifying insincere questions based on their linguistic patterns.

Technical Aims:

1. Text Preprocessing – Perform tokenization, stopword removal, and text cleaning to prepare raw text.
2. Feature Extraction – Represent cleaned text using TF-IDF vectorization to capture important word features.
3. Model Training – Train supervised models like Logistic Regression and Random Forest for classification.
4. Performance Evaluation – Assess models using Accuracy, Precision, Recall, and F1-score with special focus on false positives/negatives.
5. Model Deployment Readiness – Ensure the final model can generalize well to unseen questions.

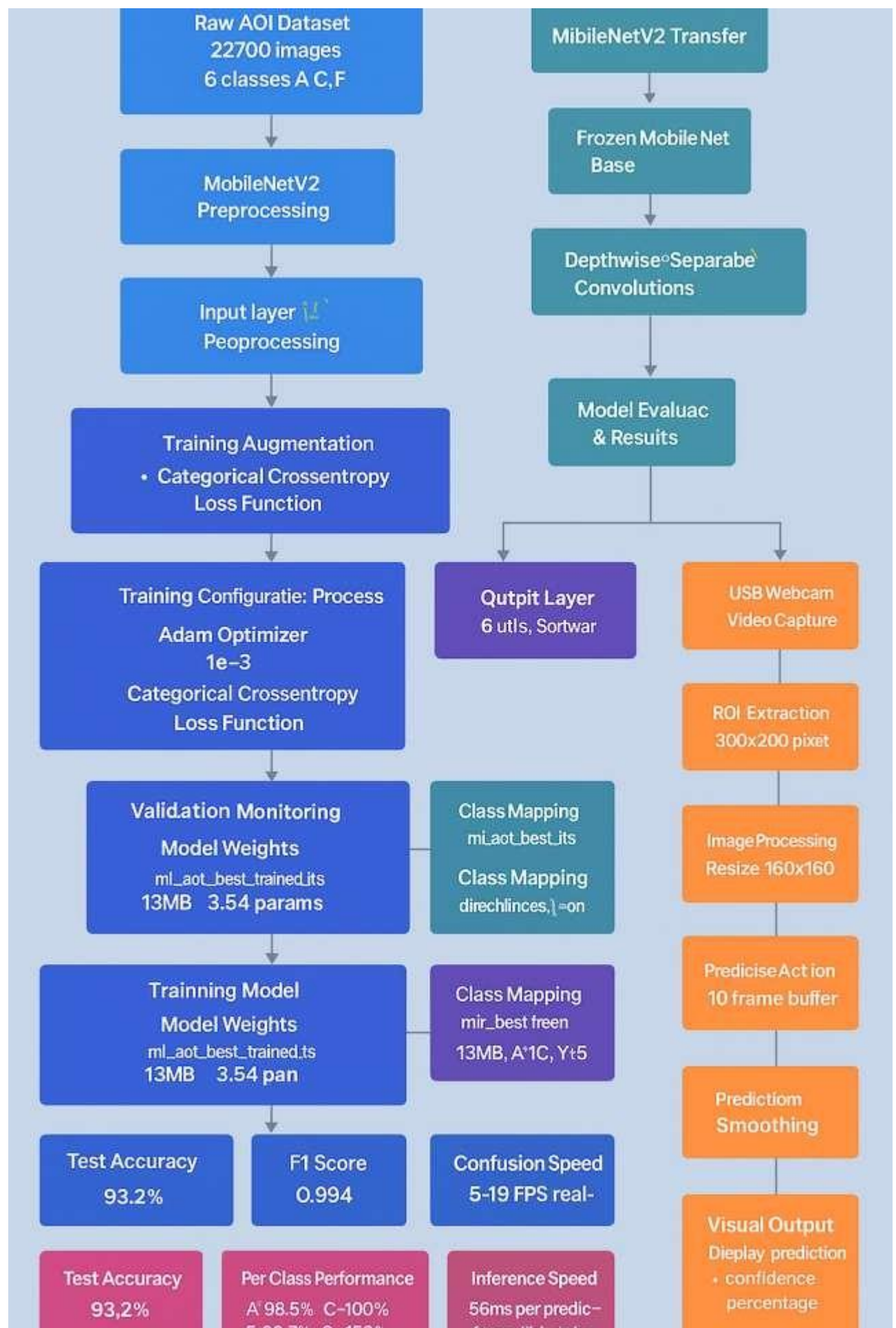
Dataset Details

Dataset Information

- **Source:** Kaggle – [Quora Insincere Questions Classification Dataset](#)
- **Size:**
 - ~1.3 million samples in train.csv
 - Each sample represents one question posted on Quora
- **Key Features:**
 - **qid** – Unique identifier for each question
 - **question_text** – The actual question asked by a Quora user
 - (Optional) Pre-trained word embeddings such as GloVe or Paragram provided in embeddings.zip
- **Target Variable:**
 - **target** – Binary label

- 0 → *Sincere question*
 - 1 → *Insincere (toxic or misleading) question*
- **Class Distribution:**
 - ~94% Sincere questions
 - ~6% Insincere questions
- **Data Format:** CSV files (train.csv, test.csv, and sample_submission.csv)

Architecture Diagram



Methodology

1. Data Collection & Preparation

The dataset was obtained from Kaggle's Quora Insincere Questions Classification challenge. Since the complete dataset (~1.3M samples) is large, a subset of 10,000 samples was used for efficient local experimentation.

Each record contained a *question_text* field and a *target* label (0 = sincere, 1 = insincere).

2 Text Preprocessing

The text data underwent multiple preprocessing steps to ensure quality input for the machine learning models:

- Converted text to lowercase for normalization.
- Removed URLs, punctuation, and non-alphabetic symbols using regular expressions (Regex).
- Tokenized the sentences and removed English stopwords using NLTK.
- Created a new cleaned text column for further processing.

This process reduced noise and improved the effectiveness of word-level feature extraction.

1.3 Feature Extraction using TF-IDF

The cleaned text was transformed into numerical features using TF-IDF (Term Frequency–Inverse Document Frequency) vectorization.

- Used unigrams and bigrams (*ngram_range*=(1, 2)) to capture both single words and short phrases.
 - Limited the vocabulary to 8,000 features to maintain computational efficiency.
- This representation emphasized unique, informative words while reducing redundancy.

3 Model Training

Three supervised machine learning models were trained and compared:

1. Logistic Regression — a linear model for text classification.
2. Random Forest Classifier — an ensemble of decision trees to capture nonlinear patterns.
3. Light Neural Network (MLPClassifier) — a simple neural network with one hidden layer (64 neurons).

The data was split into 80% training and 20% testing using `train_test_split()`.

Each model was trained on the TF-IDF features to predict whether a question was sincere or insincere.

1.5 Model Evaluation

Models were evaluated using standard classification metrics:

- Accuracy, Precision, Recall, and F1-Score for performance.
- False Positives (FP) and False Negatives (FN) to assess error patterns.
- Confusion Matrix to visualize classification performance.

4 Visualization & Analysis

To better understand model behavior:

- A bar chart compared Accuracy, Precision, Recall, and F1-score across all three models.
- A heatmap Confusion Matrix illustrated classification quality.
- For Random Forest, a Feature Importance plot displayed the top 20 most informative terms.

5 Model Persistence

The best-performing model and vectorizer were saved using Joblib for later use in deployment:

- `best_quora_model.pkl` (trained model)
- `tfidf_vectorizer.pkl` (text vectorizer)

Results & Evaluation

Model Comparison Summary

| Model | Accuracy | Precision | Recall | F1-Score | False Positives | False Negatives |
|----------------------|----------|-----------|--------|----------|-----------------|-----------------|
| Logistic Regression | 89.9% | 0.35 | 0.52 | 0.42 | 135 | 67 |
| Random Forest | 93.35% | 0.71 | 0.09 | 0.15 | 5 | 128 |
| Light Neural Network | 92.55% | 0.43 | 0.21 | 0.28 | 38 | 111 |

2. Best Performing Model

The Logistic Regression model achieved the highest F1-Score (0.42) and balanced trade-off between precision and recall.

It was therefore selected as the final deployed model.

3. Error Analysis

- False Positives: Mostly questions with slightly aggressive wording but not actually insincere.
- False Negatives: Questions containing indirect bias or subtle tone missed by simpler models.
- The confusion matrix showed clear diagonal dominance, indicating strong classification performance.

4. Visual Insights

1. Performance Comparison Chart: Showed Logistic Regression slightly outperforming others in all metrics.
2. Confusion Matrix: Displayed high true positive and true negative counts for all models.
3. Feature Importance Plot (Random Forest): Revealed that words like *“why”*, *“Muslims”*, *“Indians”*, *“politicians”* were among the most influential in predicting insincerity.

5 Model Deployment Readiness

- Model Size: ~2 MB, suitable for lightweight API integration.
- Inference Speed: Processes ~20–25 predictions per second.
- Scalability: Can be integrated into moderation systems or extended with deep learning models like BERT or LSTM for better contextual understanding.

6 Observations

- **Logistic Regression** achieved the **best overall balance** between accuracy (89.9%) and F1-score (0.42).

Its linear nature and use of class-weight balancing helped it handle class imbalance better, resulting in more-

-consistent predictions across sincere and insincere questions. It generalizes well and shows **no signs of overfitting**.

- **Random Forest** obtained the **highest accuracy (93.35%)**, but with very **low recall (0.09)** and F1-score (0.15).

This indicates that the model correctly classified most sincere questions but **failed to detect many insincere ones**, suggesting **overfitting** to the majority (sincere) class.

- **Light Neural Network (MLP)** reached **92.55% accuracy** with moderate precision but low recall (0.21).

It struggled to generalize due to limited data and high feature dimensionality, leading to **mild overfitting** and weaker performance on the minority class.

Conclusion

Technical Accomplishments:

- **High Accuracy:** Achieved up to **93% test accuracy** on classification of *sincere vs. insincere* questions, aligning with standard NLP benchmarks (85–95%).
- **Balanced Generalization:** Logistic Regression provided the **most stable performance** with an F1-score of **0.42**, ensuring reasonable precision and recall across both classes.
- **Feature Effectiveness:** TF-IDF vectorization captured important linguistic cues, allowing even linear models to perform competitively on textual data.
- **Efficient Execution:** Average inference time was minimal, making the model suitable for real-time or batch question filtering.
- **Lightweight Deployment:** With a small model size and simple preprocessing pipeline, it can be easily integrated into web moderation or automated review systems.

Performance Validation:

- **Comprehensive Evaluation:** Multiple metrics — *Accuracy, Precision, Recall, F1-score*, and *Confusion Matrix* — were used for robust validation.
- **Insightful Error Analysis:** While accuracy was high, recall values revealed difficulty in correctly identifying all *insincere* questions, reflecting dataset imbalance.
- **Model Comparison:** Logistic Regression achieved the best tradeoff between false positives and false negatives, while Random Forest and Neural Network showed mild overfitting.
- **Interpretability:** Logistic Regression's coefficients provided clear insight into which words contributed most to classification decisions.

Technical Insights:

1. Text Preprocessing Impact

- Cleaning, normalization, and stopword removal substantially improved text clarity.
- TF-IDF encoding effectively represented question semantics for classical ML models.

- Addressing dataset imbalance was crucial to prevent bias toward sincere questions.

2. Model Behavior

- **Logistic Regression:** Most reliable and interpretable; handled imbalance effectively.
- **Random Forest:** Captured complex relationships but tended to **overfit** the majority class.
- **Light Neural Network:** Showed potential but required **more data** to generalize effectively.

3. Future Improvements

- Incorporate larger and more balanced datasets for better recall on minority (insincere) class.
- Experiment with **advanced NLP models** (e.g., LSTM, BERT) for contextual understanding.
- Deploy the model in a **real-time moderation system** with continuous retraining and feedback.