CSE558 **Data Science**

# Final Project Presentation

Aditya Girdhar (**2021005**)
Bhavya Narnoli (**2021316**)
Hardik Singh (**2021390**)
Dhruv Sood (**2021387**)

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

**Built-in features**

1. drugName *(categorical)*: name of drug

2. condition *(categorical)*: name of condition

3. review *(text)*: patient review

4. rating *(numerical)*: 10-star patient rating

5. date *(date)*: date of review entry

6. usefulCount *(numerical)*: no. of users who found review useful

## Engineered Features

1. sentimentPolarity *(numerical)*: measured using NLTK

2. usefulScore *(numerical)*: normalised usefulCount*rating, essentially a weighted average of rating scaled down from 0 - 1 . ( usefulCount

# Recap (2/4): Problem Statement

We aim to cluster users based on their review patterns and preferences and recommend drugs to users who have reviewed drugs for similar conditions or have given similar ratings to drugs.

1. **Problem A**: Can we predict user's condition from their review?
2. **Problem B**: Can we recommend medicines on the basis of given condition and review?

# Recap (3/4): Preprocessing

1. Made all reviews **lower-case**

2. Cleaned the text by removing any HTML artefacts

3. Removed numbers (since they often mess with language processing)

4. Expanded all contractions (<u>you've</u> is a contraction of <u>you have</u>)

# Recap (4/4): Empirical Evaluations

1. **Hypothesis Testing**:
   a. Conducted chi-square tests for rating independence.
   b. Performed t-tests to assess correlation between ratings and sentiment  scores.

2. **Insights Gained**:
   a. Identified key patterns and trends in patient reviews.
   b. Discovered correlations between review sentiment and drug ratings.

3. **Challenges Overcome**:
   a. Addressed data quality issues including unformatted HTML codes and  missing values.
   b. Managed complexities in data standardization and feature extraction.

# First
## Goal

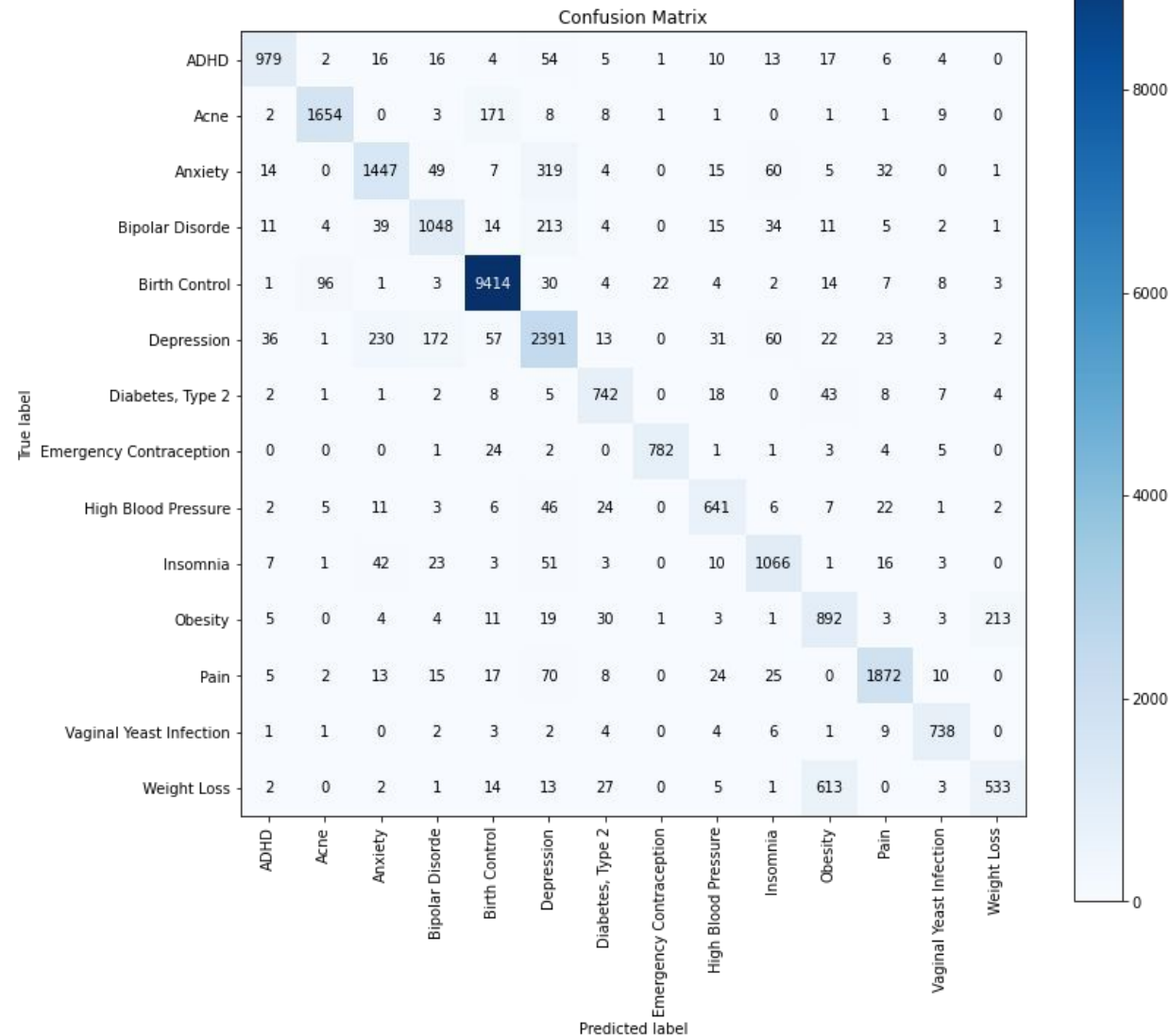Predict Condition from Review

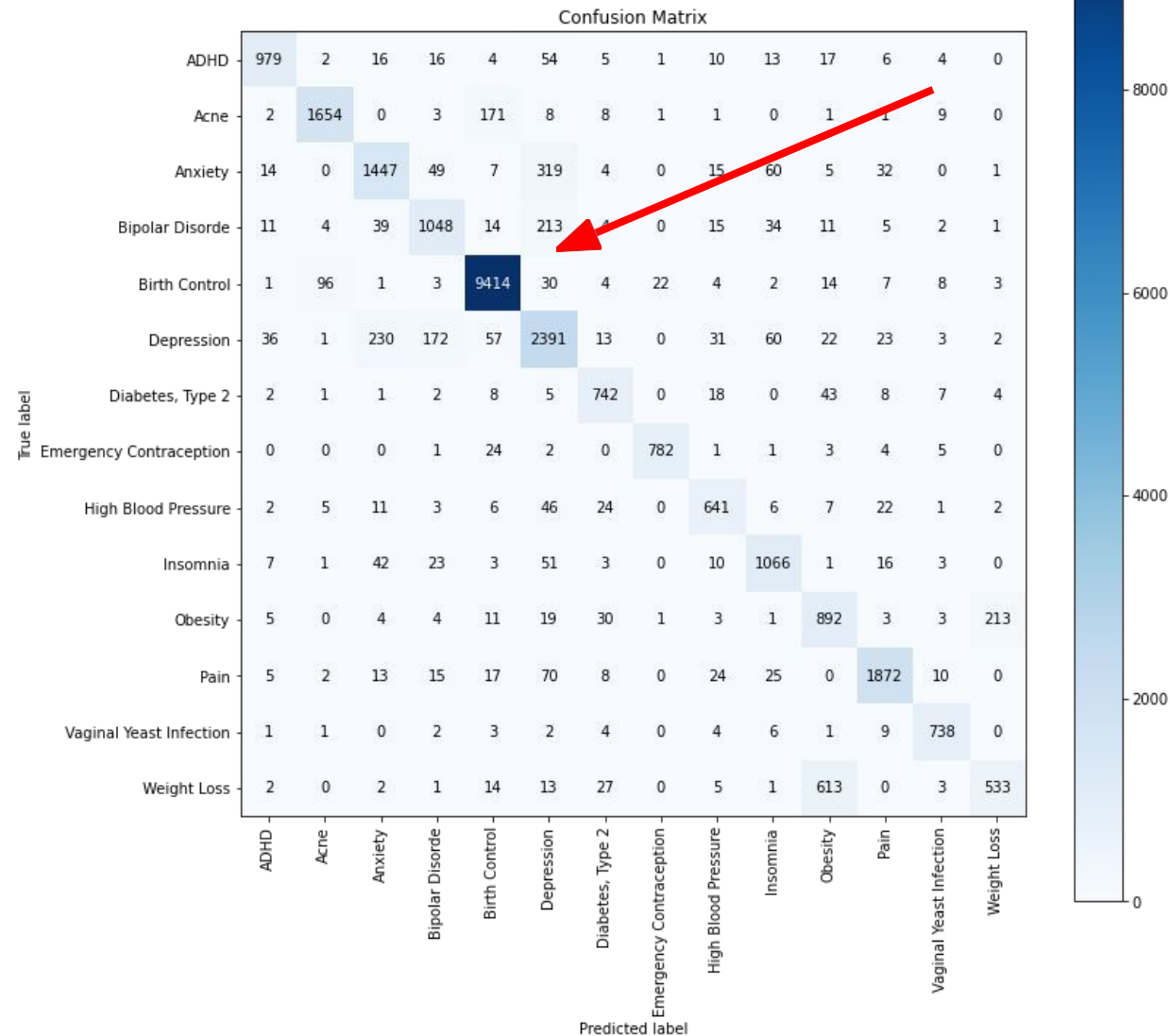# **Model**: Long Short-Term Memory Network

**Architecture Overview**

1. **Embedding Layer**: Uses nn.Embedding for token-to-vector conversion with fixed pre-trained word vectors (embedding_matrix).

2. **Bidirectional LSTM**: Employs nn.LSTM bidirectionally to capture past/future context (Input: embed_size, Hidden: self.hidden_size).

3. **Linear Layer & Activation**: Reduces LSTM output, applying ReLU for non-linearity.

4. **Pooling & Concatenation**: Uses average and max pooling to capture temporal information, concatenating these for feature extraction.

5. **Dropout**: Helps prevent overfitting by deactivating 10% of units during training (drp = 0.1) using self.dropout.

Model Accuracy: **86.8%**

# LSTM Confusion
# Matrix



Confusion Matrix

# LSTM Confusion Matrix


Confusion Matrix

# LSTM Confusion
# Matrix



Confusion Matrix

# **Model**: Convolutional Neural Network

**Architecture Overview**

1. **Embedding Layer**: Converts tokens to dense vectors using pre-trained  embeddings (fixed weights).

2. **Convolutional Layers**: Multiple layers capture n-grams with varied kernel  sizes. Extracts features using num_filters filters.

3. **Activation and Pooling**: Applies ReLU for non-linearity and max pooling  along the temporal dimension.

4. **Dropout**: Mitigates overfitting by deactivating 10% units randomly.

5. **Fully Connected Layer**: Flattens pooled output for predictions with  neuron count determined by filter_sizes and num_filters.

Model Accuracy: **88.4%**

# CNN Confusion Matrix



Confusion Matrix

# Why does **CNN** perform better than LSTM?

1.  **Local Feature Extraction**: CNNs excel at capturing local patterns and features within data. In text data, for instance, they can identify and learn representations of specific word sequences or n-grams efficiently. For tasks where local structures are crucial (like in image analysis or certain types of text classification), CNNs might outperform LSTMs.

2.  **Translation Invariance**: CNNs inherently possess translation invariance, which means they can recognize patterns regardless of their location in the input. This property is particularly useful in tasks where the position of features doesn't significantly impact the final prediction, such as in image classification.

# Second Goal

Recommend Medicines

# Term Frequency: Inverse Document Frequency

**TF-IDF** : Weighting strategy in which words are offered with weight, not count.

It gives low importance to the terms that often appear in the dataset, which implies TF-IDF estimates relevance, not a recurrence. Term frequency (TF) can be called the likelihood of locating a word in a document.

$$tf(t, d) = \log(1 + freq(t, d))$$

Inverse document frequency (IDF) is the opposite of the number of times a specific term showed up in the whole corpus. It catches how a specific term is document specific.

$$tfidf(t, d, D) = tf(t, d) \cdot idf(t, D)$$

TF-IDF is the multiplication of TF with IDF, suggesting how vital and relevant a word is in the document. The selected n-gram range for TF-IDF in this is ([1],[2]). [ unigram and bigram ]

Source: link

# Features and Metric for recommendation

Metric :

usefulScore was normalised if usefulScore > 0.5 : Recommendable Drug

else if usefulScore < 0.5 : Not Recommendable Drug

Training features : TF-IDF matrix, along with condition, sentiment_score

Output Binary Prediction : useful Score

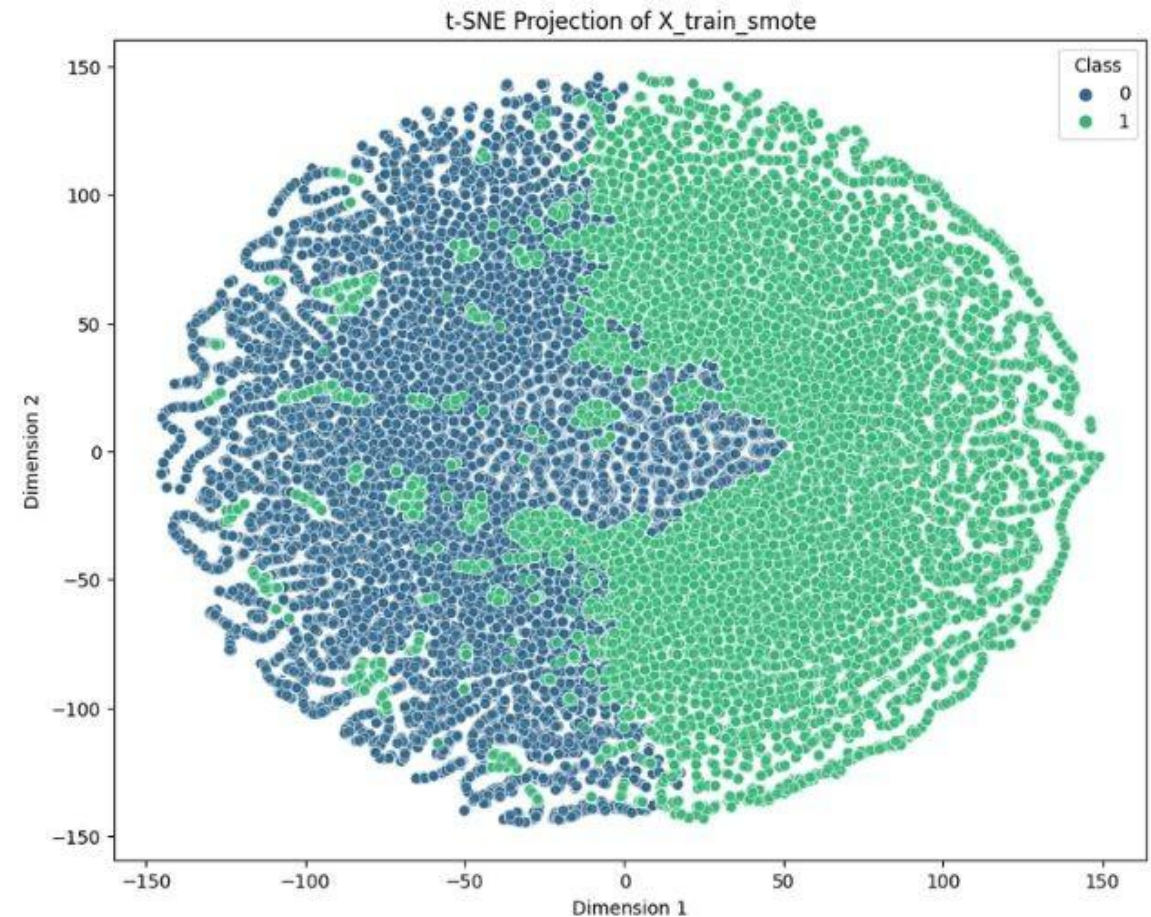Distribution of the data

0 : 206755

1 : 5923

This had class imbalance issue, which was overcome by using SMOTE oversampling which uses  KNN in the feature space to oversample the minor class, by generating similar data

# T-sne of above trainable features

The linear separability could indicate that these clusters are easily distinguishable from each other.

# Models (Preliminary)

**Perceptron:** A perceptron is a basic artificial neuron that takes binary inputs, applies weights, sums them up, passes through an activation function, and outputs a binary result. The same above inputs and outputs were considered in building the model. It demonstrates in an overall accuracy of **51%**.

**Logistic Regression** is a binary classification algorithm that models the probability of an instance belonging to a particular class, providing precision, recall, and F1-score metrics, with an overall accuracy of **62%**.The same above features (target and input were used in the model)

**Random Forest** is an ensemble model known for its robustness and ability to handle complex relationships.It was used to predict that for given (sentiment score, TF-IDF vectors, condition) the medicine should be recommended or not. The model showed high precision, recall, and F1-scores for both positive and negative classes, resulting in an overall accuracy of **98%**.

# Best Medicine for a given condition

The code iterates through each unique medical condition in the training set. For each condition, it identifies the data point with the highest confidence score (distance from the decision boundary) as predicted by the logistic regression model.

Output file : <u>link</u>

# **Model:** Best-K Review Scores

We devised a composite score to gauge product/service performance  based on textual reviews. The formula is as follows:

**review_score** = 0.6*rating + 0.3*useful_count + 0.1*sentiment_score

*Top 5 and Bottom 5 Drugs for* a condition *Guides doctors' decisions, showcasing public sentiment. Valuable for research, identifying concerns with poorly-rated drugs, fostering proactive healthcare management.*

# **Model:** Best-K Review Scores

```
Top 5 Drugs for Depression:
       drugName    reviewScore
0     Sertraline    393.39772
1         Zoloft    393.39772
2         Zoloft    290.03185
3     Sertraline    290.03185
4     Citalopram    236.02773


Bottom 5 Drugs for Depression:
          drugName    reviewScore
0        Fluoxetine      0.50150
1   Levomilnacipran      0.50248
2        Vilazodone      0.51090
3           Viibryd      0.51090
4        Citalopram      0.51212
```

# Thank you!

## Team Members

- Aditya Girdhar (**2021005**)

- Bhavya Narnoli (**2021316**)

- Hardik Singh (**2021390**)

- Sanskar Ranjan (**2021096**)

All members equally contributed to this project.