# **OVERVIEW OF THE FINAL PROCESS**

### **BUSINESS PROBLEM**

Classification of tickets and assignment to their appropriate resolver groups based on available ticket data and features

### **OUR APPROACH**

Division of data into two groups, one for groups with low number of records and the other for groups with high number of records

### **SELECTED STRATEGY ELEMENTS**

#### **Features of Data:**

- Ticket descriptions are Multilanguage
- Data is highly imbalanced
- · Data has text as well as numeric features

# Translation used:

- G-translate loops with backout
- Translation runs before and after data cleaning

### **Our preprocessing steps:**

- Extract numeric features
- Cleaning data, encodings using FTFY
- · Translating records, Augmenting the data

### **Augmentation used:**

- Random augmentation of upto 10 words in a record
- Random word dropping at ~15%

### Algorithms used:

- Logistic Regression, SVM, XGB
- Random Forest, LightGBM, CatBoost
- GRU, LSTM, Bidirectional LST

### **Techniques used:**

- Skip connections in GRU
- Joining LSTM text and Numeric models
- Running data through multiple models to maximize precision and recall

# **SOLUTION ROADMAP: STEP BY STEP**

|        | ACTION  | FINDING  | DETERMINATION   |
|--------|---|--|---|
| Step 1 | Perform EDA on data to discover data challenges                                 | Imbalanced target group<br>distribution, garbage<br>words & multilanguage<br>source data | Need for Data, cleaning<br>translation &<br>augmentation  |
| Step 2 | Perform Data cleaning<br>on the dataset to<br>remove fluff words and<br>garbage | Data encoded in UTF-8  | Requires conversion into corresponding language codes through FTFY  |
| Step 3 | Initiate Language     Translation through     google translate libraries        | Partial translations on first attempt  | <ul> <li>Greater translation<br/>requires removal of<br/>additional markers such<br/>as <chinese text="">,34323</chinese></li> <li>→ detected as English</li> </ul> |

# **SOLUTION ROADMAP: STEP BY STEP**

| ACTION |  | FINDING   | DETERMINATION   |  |
|--------|--|---|---|--|
| Step 4 | Augment data to create     a larger dataset  | Standard augmentation<br>with synonyms fails to<br>capture latent<br>relationships  | Inclusion of drop words augmentation  |  |
| Step 5 | Demarcate datasets into<br>rule-based and Al-based<br>and perform<br>vectorization | <ul> <li>Increasing features<br/>increases accuracy for<br/>TFIDF, additionally,<br/>some feature<br/>engineering also helps</li> </ul> | Improved feature<br>extraction required, and<br>further combination of<br>models needed |  |
| Step 6 | Perform machine learning and deep learning on data                                 | Accuracy lacking in some target groups  | Improve feature<br>extraction, translation<br>and data cleaning                         |  |

# **MODEL EVALUATION**

#### **FINAL MODEL**

✓ Random Forest

#### **BEST PARAMETERS**

- √ 'criterion': 'entropy'
- √ 'max\_depth': None
- √ 'max features': 'log2'
- √ 'n\_estimators': 50

#### **PROMINENT PARAMETERS**

✓ Increasing n\_estimators directly reduces overfitting of the model and increases test accuracy



- High overall accuracy on test set
- High individual accuracy on target classes
- > Strong balance between precision and recall



Evaluation of Success

- > Performance on test data
- > Performance on validation data
- > Performance on re-augmented test set
- Individual class-wise accuracy/precision/recall above threshold

# **INTERIM PERFORMANCE OF MODELS**

# Logistic Regression

Fast training and scales well

#### **Performance**

- >84% (Single Model)
- ≥90% (Tuned Model)

## GRU

Fast training and scales well

### **Performance**

≥91% (Base Model)

#### **XGB**

Known for accuracy and stability of results

#### **Performance**

- > 84% (Single Model)
- > 80% (Tuned Model)

### **LSTM**

- Offers increased accuracy over GRU
- ➤ Marginally slower

#### **Performance**

➤ 89.3% (Base Model)

#### SVC

Works well with text data

#### **Performance**

>84% (Single Model)

# **Bidirectional LSTM**

- ➤ Has not shown significant gain over LSTM
- ➤ Much slower

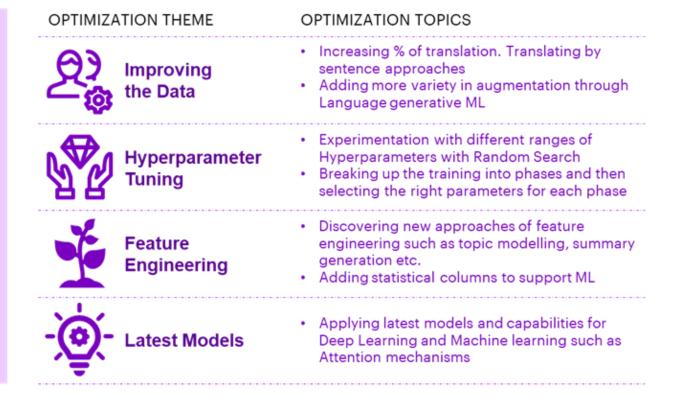
### **Performance**

>85% (Base Model)

We were unable to hyperparameterize XGB and SVC due to the significant time they were consuming. For deep learning, we will start the tuning in the next phase. The low performance of LSTM and Bidirectional LSTM is solely due to the need for greater training time.

# **OUR INTERIM PLAN ON FUTURE OPTIMIZATIONS**

Improving our outcomes can be done with 4 broad themes



| OUR FINAL MODELS |  |
|------------------|--|
|                  |  |
|                  |  |

| Classifica   |
|--|
| Olassilica   |
| Logistic   |
| ogistic Regression using Rand  |
| Logistic Gr  |
| XGBoos   |
| XGBoos   |
| XGBoost Classifier   |
| Support Vector Cl  |
| upport Vector Classification (C  |
| Random Fore  |
| Random Forest Classifier (C  |
|  |
|  |
| LightGBM (C  |
| CatBoo   |
|  |
| CatBoo   |
| CatBoo<br>CatBoost Cla   |
| CatBoo<br>CatBoost Cla<br>CatBoost Cla   |
| CatBoo<br>CatBoost Cla<br>CatBoost Cla<br>LSTM Me  |
| CatBoo<br>CatBoost Cla<br>CatBoost Cla<br>LSTM Me<br>Bi Direc  |
| CatBoo<br>CatBoost Cla<br>CatBoost Cla<br>LSTM Me  |
| CatBoo<br>CatBoost Cla<br>CatBoost Cla<br>LSTM Me<br>Bi Direc  |
| CatBoo<br>CatBoost Cla<br>CatBoost Cla<br>LSTM Me<br>Bi Direc<br>GRU Skip                              |
| CatBoo<br>CatBoost Cla<br>CatBoost Cla<br>LSTM Me<br>Bi Direc<br>GRU Skip (<br>Random Forest Classifie |

| Classification Model                      |  | Training Accuracy Test Acc |       |  |
|---|--|----------------------------|-------|--|
| Logistic Regression                       |  | 90.74                      | 89.62 |  |
| ogistic Regression using Random Search    |  | 94.32                      | 93.57 |  |
| Logistic Grid Classifier                  |  | 94.41                      | 93.64 |  |
| XGBoost Classifier1                       |  | 91.68                      | 90.08 |  |
| XGBoost Classifier2                       |  | 95.78                      | 94.42 |  |
| XGBoost Classifier Grid Search            |  | 96.60                      | 95.52 |  |
| Support Vector Classification             |  | 94.03                      | 93.45 |  |
| pport Vector Classification (Grid Search) |  | 96.98                      | 96.91 |  |
| Random Forest Classifier                  |  | 98.41                      | 97.78 |  |
| Random Forest Classifier (Grid Search)    |  | 98.46                      | 98.31 |  |
| LightGBM                                  |  | 98.28                      | 97.52 |  |
| LightGBM (Grid Search)                    |  | 98.35                      | 97.62 |  |
| CatBoost Classifier                       |  | 89.58                      | 88.10 |  |
| CatBoost Classifier Grid                  |  | 82.34                      | 80.54 |  |
| CatBoost Classifier Best                  |  | 92.41                      | 90.96 |  |
| LSTM Merged Model                         |  | 94.99                      | 94.26 |  |
| LSTM                                      |  | 94.89                      | 94.42 |  |
| GRU                                       |  | 95.34                      | 94.85 |  |
| Bi Directional LSTM                       |  | 95.39                      | 94.84 |  |
| GRU Skip Connection                       |  | 95.63                      | 95.28 |  |
| Random Forest Classifier Base Data        |  | 98.46                      | 65.61 |  |
| LightGBM Base Data                        |  | 97.89                      | 68.18 |  |
| Logistic Regression Base Data             |  | 82.18                      | 60.29 |  |
| Models selected                           |  | Models on base data        |       |  |

# **COMPARISON TO BENCHMARK**

## WE DELIVERED OVER 7% IMPROVEMENT TO OUR 91% BENCHMARK IN THE INTERIM









Leveraging Skip Connections, Merged deep learning models, we were able to improve deep learning accuracy past 95%

Compared to our interim efforts, we extracted multiple new features of value to the classification We improved our augmentation outcomes and variety, increased the degree of translation to almost 100% and reduced the data loss from cleaning

We brought our hyperparameters closer to optimal ranges using coarse and fine gradations from the default values

Our final accuracy was above 98% using random forest models, this is a 30% improvement over the translated data without augmentation

# IMPLICATIONS FOR BUSINESS



# **ASSIGNMENT TO GROUPS**



# ACCURACY OF ASSIGNMENT





The solution proceeds through the test group (19744 records) in around 10 minutes for 3 models. This means that assignment is being performed at 2000 records per minute. **Compared to human** assignment, this is several hundred times faster

The best model of the lot is random forests, which provides 98.31% accuracy. **Combining this with other** deep learning models, we can move towards nearly 99% accuracy which would be a significant improvement over human assignment

We recommend the below

- 1. Dissolution of groups which have less than 1% records from the total
- 2. Auto resolution of several tickets that relate to SID's or Hosts or batches that have failed. They can be automatically restarted based on detection

Based on the set of values we obtain from our best models; we can estimate the accuracy to be within 96.3% -98.61%. 99% of the times.

# LIMITATIONS OF THE MODEL

# LIMITATIONS, REAL WORLD DEFICIENCIES, AND FUTURE ENHANCEMENTS



#### **LIMITATION:**

#### **CHINESE LANGUAGE**

Currently, the caller is a key differentiating factor in the tickets raised in Chinese language. However, as tickets increase, there could be some overlap between the users in different groups causing accuracy reduction.

Additionally, tickets raised automatically in groups 5,6,8,9 are challenging for the algorithm to classify

#### **DEFICIENCY:**

#### **LIMITED DATA**

Real world data can be vastly different from the data provided. In production scale ticketing systems, organizations can get hundreds of thousands of tickets per year. This means that we would not have captured the entire vocabulary for real time assignment, nor planned for its full impact on performance & future training

#### **ENHANCEMENT:**

# CONTINUED LEARNING FROM CLOSED TICKETS

We need to configure the solution in a way that allows models to continuously be trained on new data as it arrives

The solution also needs to generate new patterns of language and vocabulary to create a deeper coverage of what a caller may raise in a ticket

# **OBSERVATION**

## LEARNING FOR FUTURE



LOW ACCURACY IN 6% OF TARGET GROUPS DISCOVERED POST MODEL BUILDING

GROUP BY GROUP RECORD REVIEW
IN AS PART OF INITIAL EDA





ACCURACY REDUCTION DUE TO DATA LOSS IN CLEANING DATA

DETERMINATION OF VALUE OF EVERY FEATURE/COLUMN PRIOR TO CLEANING





LARGE DEVIATION FROM BASE HYPERPARAMETERS DURING TUNING CAUSING LOSS OF EFFORT GRADUAL GRADATION OF HYPERPARAMETERS FROM DEFAULT SETTINGS





A SINGLE MODEL DOES NOT SUFFICE FOR ALL INFORMATION

STACKING OR MERGING MODELS MAY IMPROVE ACCURACY

