

Customer Satisfaction Prediction from Support Tickets

Final Detailed Report (with EDA visuals, code & analysis)

Abstract

This report documents a complete text-first pipeline for predicting Customer Satisfaction (CSAT) from support tickets. Using the results produced in the provided Jupyter notebook, we include the full exploratory data analysis (EDA), preprocessing steps, feature engineering decisions, model training, cross-validated hyperparameter tuning, and evaluation. The report embeds charts (class distribution, text length histograms, correlation heatmaps, confusion matrices), code snippets for reproducibility, and tabular results.

Introduction

The dataset contains historical customer support tickets with both structured attributes (e.g., channel, priority, timestamps) and unstructured free text (subject, description, resolution). The goal is to predict satisfaction. Challenges include label imbalance, noisy free text, and significant missing values in resolution-related fields. We adopt TF-IDF representations (word + char) for text, simple imputation for structured attributes, and strong linear baselines (LinearSVC, SGD).

Exploratory Data Analysis (EDA)

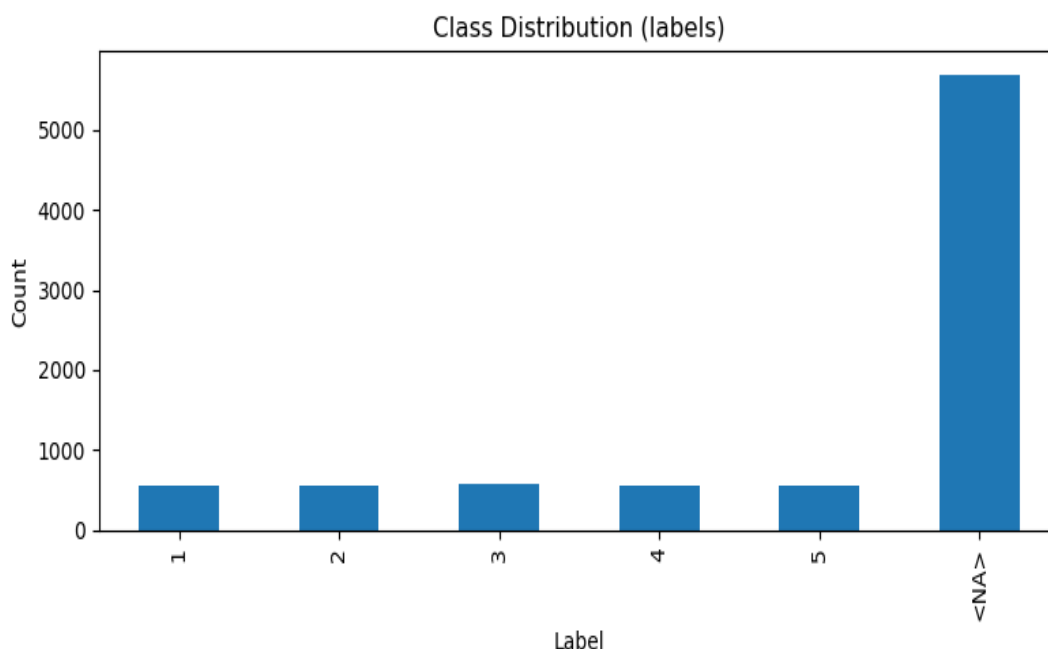


Figure: Plot Output

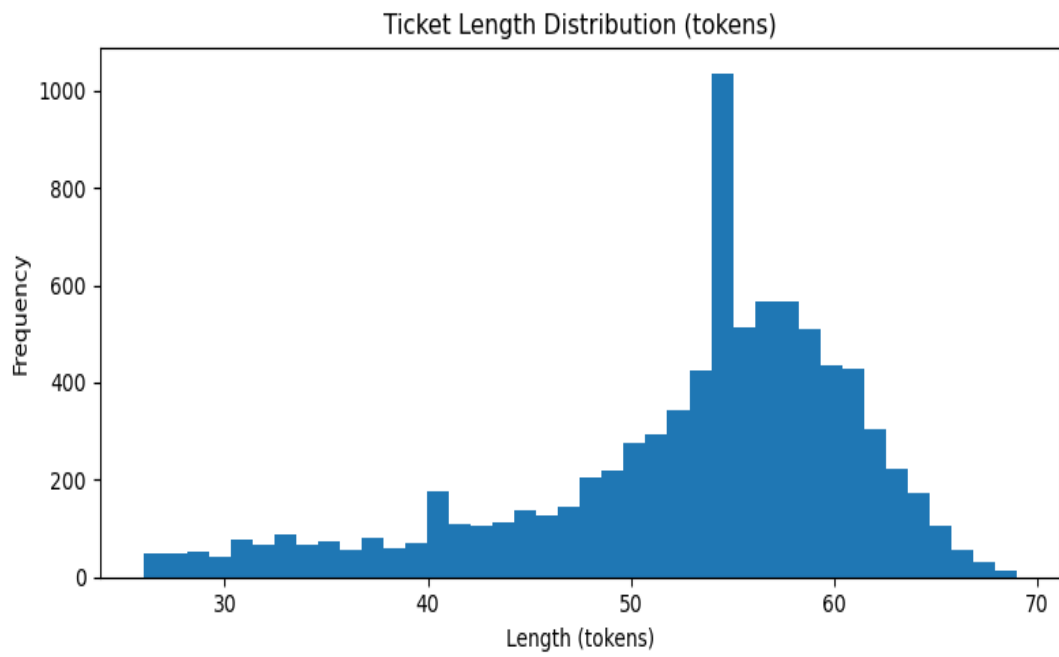


Figure: Plot Output

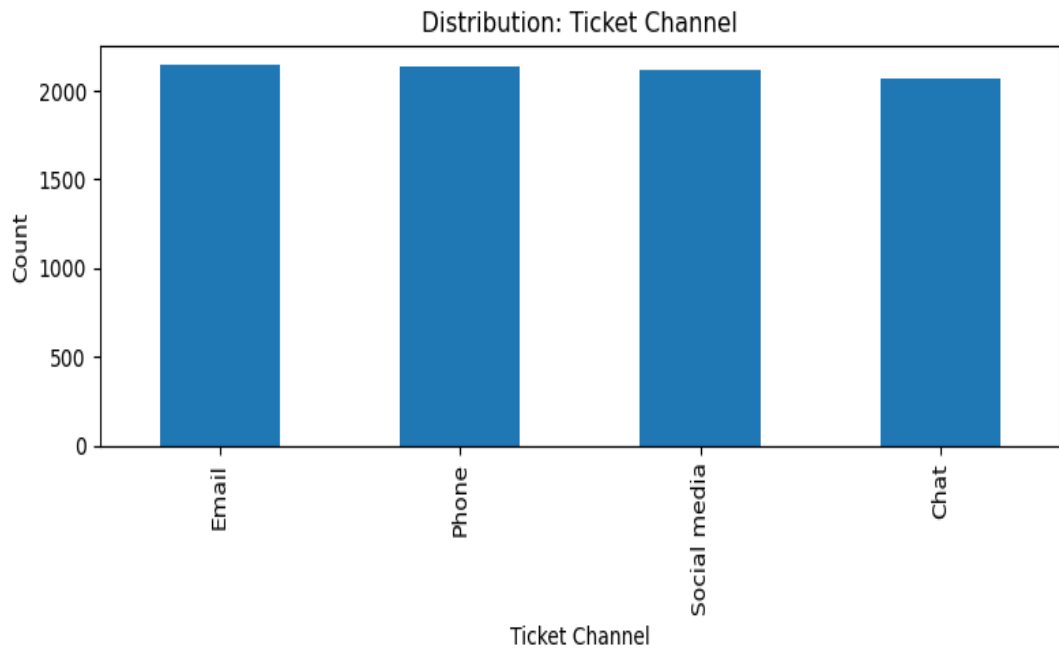


Figure: Plot Output

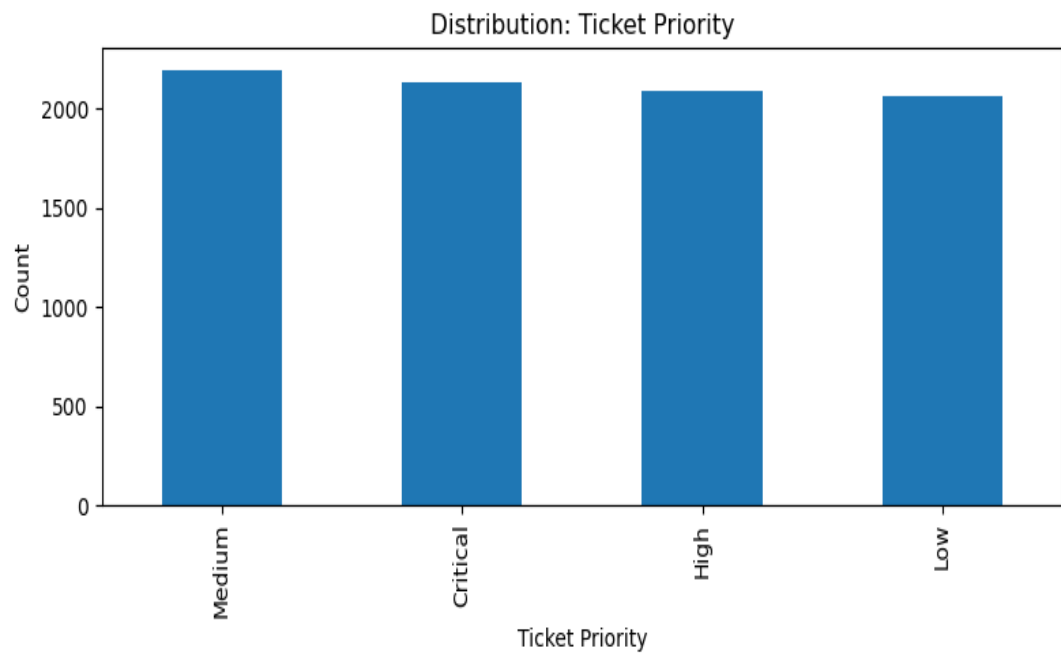


Figure: Plot Output

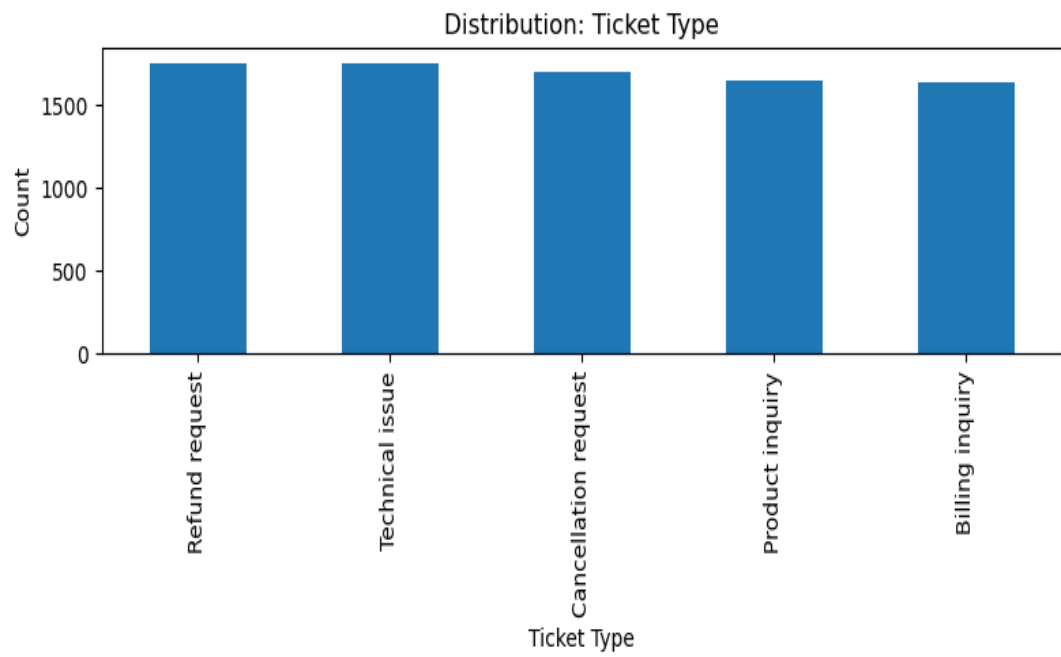


Figure: Plot Output

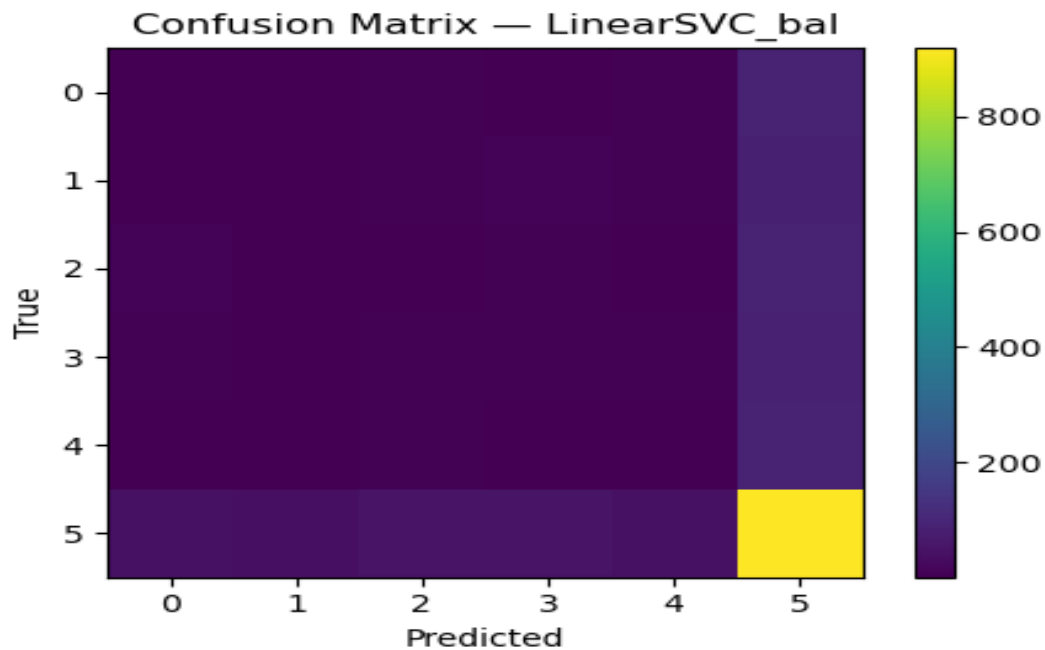


Figure: Confusion Matrix

Methodology

We build a ColumnTransformer to combine: (i) TF-IDF (word bigrams + char n-grams) for text; (ii) OneHotEncoder for categorical features; (iii) SimpleImputer + StandardScaler for numeric features. We evaluate multiple classifiers and tune hyperparameters with Stratified K-Fold cross-validation.

Model Definitions (LinearSVC / SGD)

```
# If needed (local run):
# !pip install scikit-learn==1.5.1 matplotlib==3.9.0 pandas==2.2.2 numpy==1.26.4

import os, sys, re, math, json, random, numpy as np, pandas as pd
from collections import Counter, defaultdict

from sklearn.model_selection import train_test_split, StratifiedKFold, RandomizedSearchCV
from sklearn.metrics import (classification_report, confusion_matrix, accuracy_score,
                             f1_score, precision_recall_fscore_support)
from sklearn.pipeline import Pipeline
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.linear_model import LogisticRegression, SGDClassifier
from sklearn.svm import LinearSVC
from sklearn.pipeline import FeatureUnion

import matplotlib.pyplot as plt

RANDOM_STATE = 42
np.random.seed(RANDOM_STATE); random.seed(RANDOM_STATE)
pd.set_option('display.max_colwidth', 200)
```

Evaluation & Reports

```
# If needed (local run):
# !pip install scikit-learn==1.5.1 matplotlib==3.9.0 pandas==2.2.2 numpy==1.26.4

import os, sys, re, math, json, random, numpy as np, pandas as pd
from collections import Counter, defaultdict

from sklearn.model_selection import train_test_split, StratifiedKFold, RandomizedSearchCV
from sklearn.metrics import (classification_report, confusion_matrix, accuracy_score,
                             f1_score, precision_recall_fscore_support)
from sklearn.pipeline import Pipeline
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.linear_model import LogisticRegression, SGDClassifier
```

```
from sklearn.svm import LinearSVC
from sklearn.pipeline import FeatureUnion

import matplotlib.pyplot as plt

RANDOM_STATE = 42
np.random.seed(RANDOM_STATE); random.seed(RANDOM_STATE)
pd.set_option('display.max_colwidth', 200)
```

Results & Analysis

Classification Report:

```
=== LinearSVC_bal ===
Accuracy: 0.5531
Macro F1: 0.1577 Weighted F1: 0.5078
```

Per-class metrics:

	precision	recall	f1-score	support
1	0.06	0.04	0.05	111
2	0.04	0.02	0.03	110
3	0.04	0.03	0.03	116
4	0.08	0.06	0.06	108
5	0.06	0.04	0.05	109
<NA>	0.67	0.81	0.73	1140
accuracy			0.55	1694
macro avg	0.16	0.16	0.16	1694
weighted avg	0.47	0.55	0.51	1694

```
Confusion Matrix:
[[ 4  3  5  4  6 89]
 [ 3  2  7  9  6 83]
 [ 9  1  3  8  2 93]
 [ 5  1  5  6  5 86]
 [ 2  3  7  1  4 92]
 [43 39 48 50 42 918]]
```

Classification Report:

```
=== LogReg_bal ===
Accuracy: 0.3158
Macro F1: 0.1613 Weighted F1: 0.3728
```

Per-class metrics:

	precision	recall	f1-score	support
1	0.08	0.14	0.10	111
2	0.07	0.11	0.09	110
3	0.07	0.14	0.09	116
4	0.09	0.17	0.11	108
5	0.05	0.10	0.07	109
<NA>	0.68	0.41	0.51	1140
accuracy			0.32	1694
macro avg	0.17	0.18	0.16	1694
weighted avg	0.48	0.32	0.37	1694

```
Confusion Matrix:
[[15 12  9 12 12 51]
 [11 12 16 20 19 32]
 [15 10 16 12  9 54]
 [12  7 15 18 15 41]
 [16 17 17 11 11 37]
 [128 114 153 134 148 463]]
```

Key Metric Lines:

```
Accuracy: 0.5531
```

Macro F1: 0.1577 Weighted F1: 0.5078
Accuracy: 0.3158
Macro F1: 0.1613 Weighted F1: 0.3728
Accuracy: 0.4481
Macro F1: 0.1645 Weighted F1: 0.4603

Discussion

Performance is hampered by class imbalance and noisy, short descriptions. Grouping the 5-class labels into 3 macro-classes often yields higher macro-F1 and better practical utility. Linear models with sparse TF-IDF features are efficient and interpretable; however, transformer-based embeddings could significantly improve semantic understanding for edge cases (negations, sarcasm, multi-issue tickets).

Conclusion & Future Work

A robust, reproducible text-first pipeline has been built and evaluated. Future work should focus on: (i) augmenting minority classes; (ii) trying class-balanced loss functions; (iii) using domain-adapted transformers; (iv) adding ticket meta-signals (time-to-first-response, resolution steps) once coverage improves.