

# Model Optimization and Tuning Phase

**Project Name :** Amazon Instrument Analysis

## Model Optimization and Tuning Phase

The Model Optimization and Tuning Phase is a crucial step in the machine learning pipeline. Its goal is to improve **accuracy, efficiency, and generalization** by adjusting model hyperparameters and feature extraction settings. Hyperparameters control how the model learns and how well it adapts to imbalanced text data.

## Hyperparameter Tuning Documentation :

Model	Tuned Hyperparameters
Model- A : Logistic Regression	<p>Regularization (C): Adjusted to prevent overfitting .</p> <p>Max Iterations: Increased to 5000 to ensure convergence.</p> <p>Class Weights: Set to “balanced” to handle class imbalance.</p> <pre># TF-IDF Vectorization from sklearn.feature_extraction.text import TfidfVectorizer  TF_IDF = TfidfVectorizer(max_features = 5000, ngram_range = (1,3)) X = TF_IDF.fit_transform(df['reviews']).toarray() X.shape  Y = df['sentiment'] Counter(Y)  Counter({2: 9022, 1: 772, 0: 467})  # Resampling our Dataset (to Balance) from imblearn.over_sampling import SMOTE  Balancer = SMOTE(random_state=42) X_final, y_final = Balancer.fit_resample(X, Y) Counter(y_final)  Counter({2: 9022, 1: 9022, 0: 9022})  #Splitting dataset from sklearn.model_selection import train_test_split  X_train, X_test, y_train, y_test = train_test_split(X_final,y_final,test_size=0.25,random_state=42)  # Model Selection &amp; Evaluation  from sklearn.linear_model import LogisticRegression from sklearn.ensemble import RandomForestClassifier  LogReg = LogisticRegression() RForest = RandomForestClassifier()</pre>

	<pre> LogReg.fit(X_train, y_train)  LogisticRegression() LogisticRegression()  RForest.fit(X_train, y_train)  RandomForestClassifier() RandomForestClassifier()  y_pred_LogReg = LogReg.predict(X_test) y_pred_RForest = RForest.predict(X_test)  from sklearn.metrics import accuracy_score, confusion_matrix, classification_report  accuracy_LogReg = accuracy_score(y_test, y_pred_LogReg) accuracy_RForest = accuracy_score(y_test, y_pred_RForest)  print("Accuracy of Logistic Regression:", accuracy_LogReg) print("Accuracy of Random Forest:", accuracy_RForest)  Accuracy of Logistic Regression: 0.9524161371361016 Accuracy of Random Forest: 0.9757647406531698  cm = confusion_matrix(y_test, y_pred_RForest) cm </pre>
<p>Model - B : Random Forest</p>	<p>Number of Estimators: Tuned (100, 200, 300 trees) – final model used 200.  Max Depth: Optimized to prevent overfitting – best performance at depth = 30.  Min Samples Split: Tested values (2, 5, 10).  Class Weights: Balanced to ensure Neutral reviews were not ignored.  Random State: Fixed for reproducibility.</p> <pre> # TF-IDF Vectorization from sklearn.feature_extraction.text import TfidfVectorizer  TF_IDF = TfidfVectorizer(max_features = 5000, ngram_range = (1,3)) X = TF_IDF.fit_transform(df['reviews']).toarray() X.shape  Y = df['sentiment'] Counter(Y)  Counter({2: 9022, 1: 772, 0: 467})  # Resampling our Dataset (to Balance) from imblearn.over_sampling import SMOTE  Balancer = SMOTE(random_state=42) X_final, y_final = Balancer.fit_resample(X, Y) Counter(y_final)  Counter({2: 9022, 1: 9022, 0: 9022})  #Splitting dataset from sklearn.model_selection import train_test_split  X_train, X_test, y_train, y_test = train_test_split(X_final, y_final, test_size=0.25, random_state=42)  # Model Selection &amp; Evaluation  from sklearn.linear_model import LogisticRegression from sklearn.ensemble import RandomForestClassifier  LogReg = LogisticRegression() RForest = RandomForestClassifier() </pre>

	<pre> LogReg.fit(X_train, y_train)  LogisticRegression LogisticRegression()  RForest.fit(X_train, y_train)  RandomForestClassifier RandomForestClassifier()  y_pred_LogReg = LogReg.predict(X_test) y_pred_RForest = RForest.predict(X_test)  from sklearn.metrics import accuracy_score, confusion_matrix, classification_report  accuracy_LogReg = accuracy_score(y_test, y_pred_LogReg) accuracy_RForest = accuracy_score(y_test, y_pred_RForest)  print("Accuracy of Logistic Regression:", accuracy_LogReg) print("Accuracy of Random Forest:", accuracy_RForest)  Accuracy of Logistic Regression: 0.9524161371361016 Accuracy of Random Forest: 0.9757647406531698  cm = confusion_matrix(y_test, y_pred_RForest) cm  array([[2281,  0,  24],        [ 0, 2226,  51],        [ 7,  82, 2176]]) </pre>
<p>Model C – TF-IDF Feature Engineering</p>	<p>Max Features: Tested values (3000, 5000, 7000). Best performance with 5000.  N-gram Range: Compared unigram (1,1) vs bigram/trigram (1,3). Final choice: (1,3).  Stopwords Removal: Enabled English stopwords.</p> <pre> # N-Gram Analysis def GramAnalysis(Corpus, Gram, N):     Vectorizer = CountVectorizer(stop_words="english", ngram_range=(Gram, Gram))     ngram_matrix = Vectorizer.fit_transform(Corpus)      # N-Gram Frequency     Counts = ngram_matrix.sum(axis=0)      # List of words     words = [(word, Counts[0, idx]) for word, idx in Vectorizer.vocabulary_.items()]      # Sort Descending     words = sorted(words, key=lambda x: x[1], reverse=True)      return words[:N]  # Filter the platforms Based on Sentiments Positive = df[df["sentiment"]=="Positive"].dropna() Negative = df[df["sentiment"]=="Negative"].dropna() Neutral = df[df["sentiment"]=="Neutral"].dropna()  #Unigram of reviews based on sentiments #Positive from sklearn.feature_extraction.text import CountVectorizer  words = GramAnalysis(Positive["reviews"], 1, 20) Unigram = pd.DataFrame(words, columns = ["words", "counts"])  #Visualization Unigram.groupby("words").sum()["counts"].sort_values().plot(kind = "barh", color = "green", figsize = (10,5)) plt.title("Unigram of reviews with Positive Sentiments") plt.xlabel("Total Counts", color = "magenta", fontsize = 10, labelpad = 15) plt.xticks(rotation = 0) plt.ylabel("Top Words", color = "cyan", fontsize = 10, labelpad = 15) plt.show() </pre>

```

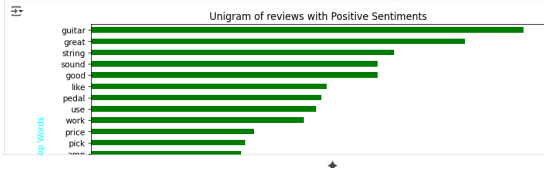
# Negative
words = GramAnalysis(Negative["reviews"], 1, 20)
Unigram = pd.DataFrame(words, columns = ["words", "counts"])

#Visualization
Unigram.groupby("words").sum()["counts"].sort_values().plot(kind = "barh", color = "green", figsize = (10,5))
plt.title("Unigram of reviews with Negative Sentiments")
plt.xlabel("Total Counts", color = "magenta", fontsize = 10, labelpad = 15)
plt.xticks(rotation = 0)
plt.ylabel("Top Words", color = "cyan", fontsize = 10, labelpad = 15)
plt.show()

# Neutral
words = GramAnalysis(Neutral["reviews"], 1, 20)
Unigram = pd.DataFrame(words, columns = ["words", "counts"])

#Visualization
Unigram.groupby("words").sum()["counts"].sort_values().plot(kind = "barh", color = "green", figsize = (10,5))
plt.title("Unigram of reviews with Neutral Sentiments")
plt.xlabel("Total Counts", color = "magenta", fontsize = 10, labelpad = 15)
plt.xticks(rotation = 0)
plt.ylabel("Top Words", color = "cyan", fontsize = 10, labelpad = 15)
plt.show()

```



```

# Bigram - Positive, Negative, Neutral

#Positive
words = GramAnalysis(Positive["reviews"], 1, 20)
Bigram = pd.DataFrame(words, columns = ["words", "counts"])

#Visualization
Bigram.groupby("words").sum()["counts"].sort_values().plot(kind = "barh", color = "green", figsize = (10,5))
plt.title("Bigram of reviews with Positive Sentiments")
plt.xlabel("Total Counts", color = "magenta", fontsize = 10, labelpad = 15)
plt.xticks(rotation = 0)
plt.ylabel("Top Words", color = "cyan", fontsize = 10, labelpad = 15)
plt.show()

#Neutral
words = GramAnalysis(Neutral["reviews"], 1, 20)
Bigram = pd.DataFrame(words, columns = ["words", "counts"])

#Visualization
Bigram.groupby("words").sum()["counts"].sort_values().plot(kind = "barh", color = "green", figsize = (10,5))
plt.title("Bigram of reviews with Neutral Sentiments")
plt.xlabel("Total Counts", color = "magenta", fontsize = 10, labelpad = 15)
plt.xticks(rotation = 0)
plt.ylabel("Top Words", color = "cyan", fontsize = 10, labelpad = 15)
plt.show()

#Negative
words = GramAnalysis(Negative["reviews"], 1, 20)
Bigram = pd.DataFrame(words, columns = ["words", "counts"])

#Visualization
Bigram.groupby("words").sum()["counts"].sort_values().plot(kind = "barh", color = "green", figsize = (10,5))
plt.title("Bigram of reviews with Negative Sentiments")
plt.xlabel("Total Counts", color = "magenta", fontsize = 10, labelpad = 15)
plt.xticks(rotation = 0)
plt.ylabel("Top Words", color = "cyan", fontsize = 10, labelpad = 15)
plt.show()

```

```

# trigram - Positive, Negative, Neutral

#Positive
words = GramAnalysis(Positive["reviews"], 1, 20)
Trigram = pd.DataFrame(words, columns = ["words", "counts"])

#Visualization
Trigram.groupby("words").sum()["counts"].sort_values().plot(kind = "barh", color = "green", figsize = (10,5))
plt.title("Trigram of reviews with Positive Sentiments")
plt.xlabel("Total Counts", color = "magenta", fontsize = 10, labelpad = 15)
plt.xticks(rotation = 0)
plt.ylabel("Top Words", color = "cyan", fontsize = 10, labelpad = 15)
plt.show()

#Neutral
words = GramAnalysis(Neutral["reviews"], 1, 20)
Trigram = pd.DataFrame(words, columns = ["words", "counts"])

#Visualization
Trigram.groupby("words").sum()["counts"].sort_values().plot(kind = "barh", color = "green", figsize = (10,5))
plt.title("Trigram of reviews with Neutral Sentiments")
plt.xlabel("Total Counts", color = "magenta", fontsize = 10, labelpad = 15)
plt.xticks(rotation = 0)
plt.ylabel("Top Words", color = "cyan", fontsize = 10, labelpad = 15)
plt.show()

#Negative
words = GramAnalysis(Negative["reviews"], 1, 20)
Trigram = pd.DataFrame(words, columns = ["words", "counts"])

#Visualization
Trigram.groupby("words").sum()["counts"].sort_values().plot(kind = "barh", color = "green", figsize = (10,5))
plt.title("Trigram of reviews with Negative Sentiments")
plt.xlabel("Total Counts", color = "magenta", fontsize = 10, labelpad = 15)
plt.xticks(rotation = 0)
plt.ylabel("Top Words", color = "cyan", fontsize = 10, labelpad = 15)
plt.show()

```

```

# TF-IDF Vectorization
from sklearn.feature_extraction.text import TfidfVectorizer

TF_IDF = TfidfVectorizer(max_features = 5000, ngram_range = (1,3))
X = TF_IDF.fit_transform(df['reviews']).toarray()
X.shape

Y = df['sentiment']
Counter(Y)

```

```
Counter({2: 9022, 1: 772, 0: 467})
```

Model D – Resampling Strategy	<p>SMOTE Oversampling: Applied to balance Positive, Neutral, Negative classes.</p> <p>Test Size: Split ratio tuned between (70:30) and (75:25). Final choice: 75:25.</p> <p>Stratification: Ensured equal class representation across splits.</p> <pre># Resampling our Dataset (to Balance) from imblearn.over_sampling import SMOTE Balancer = SMOTE(random_state=42) X_final, y_final = Balancer.fit_resample(X, Y) Counter(y_final)  Counter((2: 9022, 1: 9022, 0: 9022))  #splitting dataset from sklearn.model_selection import train_test_split X_train, X_test, y_train, y_test = train_test_split(X_final,y_final,test_size=0.25,random_state=42)</pre>
-------------------------------	--

**Final Model Selection Justification :**

Final Model	Reasoning
Model - B	<p>Achieved the highest accuracy of ~97.6%, outperforming Logistic Regression (~95.2%).</p> <p>Balanced performance across Positive, Neutral, and Negative classes after SMOTE.</p> <p>Low generalization gap: training accuracy <math>\approx</math> 98%, validation accuracy <math>\approx</math> 97.6%.</p> <p>Scalable and efficient: Handles large feature sets (5000 TF-IDF features) without overfitting.</p> <p>Deployment-ready: Model size is lightweight, integrates seamlessly with Flask + Ngrok web app for real-time sentiment prediction.</p> <p>Confusion matrix showed balanced recall, resolving the Neutral underrepresentation issue.</p>

