

Supervised Machine Learning: Classification

01. Main Objective

The main objective is to predict the type of emotion expressed in tweets. The model will focus on accurate prediction to help businesses understand customer sentiment in real-time.

02. Dataset Description

```
In [1]: import pandas as pd

# Load the dataset
df = pd.read_csv('Users/nahthiyyashar/Downloads/tweet_emotions.csv')
df.head()
```

```
Out[1]:
```

	tweet_id	sentiment	content
0	1956967341	empty	@iffanyue i know i was listenin to bad hab...
1	1956967666	sadness	Layin n bed with a headache ughhhh..watin o...
2	1956967696	sadness	Funeral ceremony..gloomy friday...
3	1956967789	enthusiasm	wants to hang out with friends SOON!
4	1956968416	neutral	@dannygastilo We want to trade with someone w...

```
In [2]: import warnings
warnings.filterwarnings('ignore')
```

```
In [3]: print(f"Dataset contains {df.shape[0]} rows and {df.shape[1]} columns.")
print("Columns:", df.columns.tolist())
print("Sentiment Classes:", df['sentiment'].unique())
print("Class Distribution:\n", df['sentiment'].value_counts())

Dataset contains 40008 rows and 3 columns.
Columns: ['tweet_id', 'sentiment', 'content']
Sentiment Classes: ['empty' 'sadness' 'enthusiasm' 'neutral' 'worry' 'surprise' 'love' 'fun'
'hate' 'happiness' 'boredom' 'relief' 'anger']
Class Distribution:
sentiment
neutral    9539
worry      8459
happiness  5289
sadness    5195
love       3842
surprise   3187
fun        3208
hate       3223
empty      827
enthusiasm 759
boredom    179
anger      179
Name: count, dtype: int64
```

We can see the details of the dataset above. The number of rows, columns, the unique sentiment classes, and the class distribution for each sentiment class can be observed.

03. Data Exploration and Cleaning

```
In [4]: # Checking for missing values
print(df.isnull().sum())

tweet_id    0
sentiment    0
content      0
dtype: int64
```

```
In [5]: # Checking for duplicates
print(f"Number of duplicate rows: {df.duplicated().sum()}")

Number of duplicate rows: 8
```

There are no missing values or duplicate values, so there isn't much cleaning to do.

Now we will preprocess the text:

```
In [6]: import re
import nltk
from nltk.corpus import stopwords
from sklearn.model_selection import train_test_split
from sklearn.feature_extraction.text import TfidfVectorizer

nltk.download('stopwords')
stop_words = set(stopwords.words('english'))

def preprocess_text(text):
    # Lowercase the text
    text = text.lower()
    # Remove URLs, mentions, hashtags, and special characters
    text = re.sub('http\S+|@\S+|#\S+|/\S+', '', text)
    # Remove stopwords
    text = ' '.join(word for word in text.split() if word not in stop_words)
    return text

# Apply preprocessing to the content column
df['cleaned_content'] = df['content'].apply(preprocess_text)
```

```
[nltk_data] Downloading package stopwords to
[nltk_data]   http://www.nltk.org/nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
```

Next, we will do some feature engineering to convert the text into numerical features:

```
In [7]: from sklearn.model_selection import train_test_split
from sklearn.feature_extraction.text import TfidfVectorizer

# Convert text to numerical features using TF-IDF
vectorizer = TfidfVectorizer(max_features=5000) # Limiting to 5000 features for simplicity
x = vectorizer.fit_transform(df['cleaned_content'])
y = df['sentiment']
```

```
In [8]: # Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(x, y, test_size=0.2, random_state=42)
```

04. Model Training

Baseline Model: Logistic Regression

```
In [9]: from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix

# Train a Logistic Regression model
lr_model = LogisticRegression(max_iter=1000, random_state=42)
lr_model.fit(X_train, y_train)

y_pred_lr = lr_model.predict(X_test)

print("Logistic Regression Accuracy:", accuracy_score(y_test, y_pred_lr))
print("Classification Report:\n", classification_report(y_test, y_pred_lr))

Logistic Regression Accuracy: 0.347
Classification Report:
```

	precision	recall	f1-score	support
anger	0.80	0.80	0.80	19
boredom	0.89	0.89	0.89	31
empty	0.33	0.81	0.61	162
enthusiasm	0.89	0.89	0.89	163
fun	0.32	0.81	0.63	338
happiness	0.34	0.37	0.35	1828
hate	0.51	0.16	0.25	268
love	0.51	0.27	0.33	762
neutral	0.33	0.57	0.42	1740
relief	0.37	0.82	0.64	352
sadness	0.34	0.24	0.28	1846
surprise	0.33	0.85	0.69	425
worry	0.33	0.46	0.39	1666
accuracy			0.35	8900
macro avg	0.27	0.18	0.18	8800
weighted avg	0.34	0.35	0.31	8800

K-Nearest Neighbors

```
In [10]: from sklearn.neighbors import KNeighborsClassifier

knn_model = KNeighborsClassifier(n_neighbors=5)

knn_model.fit(X_train, y_train)
y_pred_knn = knn_model.predict(X_test)

print("KNN Accuracy:", accuracy_score(y_test, y_pred_knn))
print("KNN Classification Report:\n", classification_report(y_test, y_pred_knn))

KNN Accuracy: 0.2
KNN Classification Report:
```

	precision	recall	f1-score	support
anger	0.80	0.80	0.80	19
boredom	0.89	0.89	0.89	31
empty	0.83	0.81	0.82	162
enthusiasm	0.89	0.89	0.89	163
fun	0.89	0.82	0.85	338
happiness	0.27	0.19	0.14	1828
hate	0.35	0.17	0.11	268
love	0.44	0.16	0.23	762
neutral	0.31	0.27	0.29	1740
relief	0.13	0.81	0.82	352
sadness	0.14	0.61	0.23	1846
surprise	0.15	0.81	0.82	425
worry	0.25	0.14	0.18	1666
accuracy			0.20	8800
macro avg	0.17	0.10	0.10	8800
weighted avg	0.24	0.20	0.18	8800

Decision Tree

```
In [11]: from sklearn.tree import DecisionTreeClassifier

dt_model = DecisionTreeClassifier(random_state=42)
dt_model.fit(X_train, y_train)

y_pred_dt = dt_model.predict(X_test)

print("Decision Tree Accuracy:", accuracy_score(y_test, y_pred_dt))
print("Classification Report:\n", classification_report(y_test, y_pred_dt))

Decision Tree Accuracy: 0.262875
Classification Report:
```

	precision	recall	f1-score	support
anger	0.80	0.80	0.80	19
boredom	0.89	0.89	0.89	31
empty	0.87	0.84	0.85	162
enthusiasm	0.84	0.82	0.83	163
fun	0.13	0.89	0.11	338
happiness	0.25	0.27	0.26	1828
hate	0.25	0.19	0.21	268
love	0.22	0.22	0.22	762
neutral	0.32	0.39	0.35	1740
relief	0.19	0.85	0.87	352
sadness	0.23	0.22	0.23	1846
surprise	0.12	0.88	0.10	425
worry	0.28	0.11	0.16	1666
accuracy			0.26	8900
macro avg	0.16	0.15	0.16	8800
weighted avg	0.25	0.26	0.25	8800

Random Forest

```
In [12]: from sklearn.ensemble import RandomForestClassifier

rf_model = RandomForestClassifier(n_estimators=100, random_state=42)
rf_model.fit(X_train, y_train)

y_pred_rf = rf_model.predict(X_test)

print("Random Forest Accuracy:", accuracy_score(y_test, y_pred_rf))
print("Classification Report:\n", classification_report(y_test, y_pred_rf))

Random Forest Accuracy: 0.328625
Classification Report:
```

	precision	recall	f1-score	support
anger	0.89	0.89	0.89	19
boredom	0.89	0.89	0.89	31
empty	0.89	0.81	0.82	162
enthusiasm	0.89	0.89	0.89	163
fun	0.14	0.83	0.85	338
happiness	0.30	0.35	0.32	1828
hate	0.39	0.20	0.27	268
love	0.46	0.40	0.43	762
neutral	0.34	0.51	0.41	1740
relief	0.26	0.83	0.85	352
sadness	0.31	0.21	0.25	1846
surprise	0.18	0.84	0.85	425
worry	0.31	0.46	0.37	1666
accuracy			0.33	8800
macro avg	0.21	0.17	0.17	8800
weighted avg	0.30	0.33	0.30	8800

Naive Bayes

```
In [13]: from sklearn.naive_bayes import MultinomialNB

# Train a Naive Bayes model
nb_model = MultinomialNB()
nb_model.fit(X_train, y_train)

# Predict on the test set
y_pred_nb = nb_model.predict(X_test)

# Evaluate the model
print("Naive Bayes Accuracy:", accuracy_score(y_test, y_pred_nb))
print("Classification Report:\n", classification_report(y_test, y_pred_nb))

Naive Bayes Accuracy: 0.21825
Classification Report:
```

	precision	recall	f1-score	support
anger	0.80	0.80	0.80	19
boredom	0.89	0.89	0.89	31
empty	0.80	0.80	0.80	162
enthusiasm	0.89	0.89	0.89	163
fun	0.89	0.89	0.89	338
happiness	0.34	0.28	0.30	1828
hate	0.80	0.80	0.80	268
love	0.51	0.27	0.35	762
neutral	0.30	0.55	0.39	1740
relief	0.80	0.80	0.80	352
sadness	0.25	0.12	0.18	1846
surprise	1.00	0.81	0.81	425
worry	0.29	0.57	0.39	1666
accuracy			0.32	8800
macro avg	0.22	0.14	0.13	8800
weighted avg	0.32	0.32	0.26	8800

Gradient Boosting

```
In [14]: from sklearn.ensemble import GradientBoostingClassifier

# Train a Gradient Boosting model
gb_model = GradientBoostingClassifier(n_estimators=100, random_state=42)
gb_model.fit(X_train, y_train)

# Predict on the test set
y_pred_gb = gb_model.predict(X_test)

# Evaluate the model
print("Gradient Boosting Accuracy:", accuracy_score(y_test, y_pred_gb))
print("Classification Report:\n", classification_report(y_test, y_pred_gb))

Gradient Boosting Accuracy: 0.320875
Classification Report:
```

	precision	recall	f1-score	support
anger	0.89	0.89	0.89	19
boredom	0.89	0.89	0.89	31
empty	0.86	0.82	0.83	162
enthusiasm	0.89	0.89	0.89	163
fun	0.88	0.82	0.85	338
happiness	0.34	0.25	0.29	1828
hate	0.34	0.22	0.27	268
love	0.49	0.36	0.41	762
neutral	0.39	0.66	0.51	1740
relief	0.29	0.84	0.87	352
sadness	0.42	0.19	0.26	1846
surprise	0.21	0.85	0.88	425
worry	0.34	0.34	0.34	1666
accuracy			0.32	8800
macro avg	0.22	0.17	0.17	8800
weighted avg	0.32	0.32	0.29	8800

05. Model Recommendation

Let's first do cross validation and find out the Cross Validation Accuracies for all the models:

```
In [15]: from sklearn.model_selection import cross_val_score

models = [
    'Logistic Regression', LogisticRegression(max_iter=1000, random_state=42),
    'KNN', KNeighborsClassifier(n_neighbors=5),
    'Decision Tree', DecisionTreeClassifier(random_state=42),
    'Random Forest', RandomForestClassifier(n_estimators=100, random_state=42),
    'Gradient Boosting', GradientBoostingClassifier(n_estimators=100, random_state=42),
    'Naive Bayes', MultinomialNB()
]

# Function to perform cross-validation
def evaluate_model(model, X, y):
    scores = cross_val_score(model, X, y, cv=5, scoring='accuracy')
    return scores.mean(), scores.std()
```

Evaluate each model using cross-validation

for name, model in models.items():	mean_score, std_dev = evaluate_model(model, X, y)	print(f'{name} - Cross-Validation Accuracy: {mean_score:.4f} ± {std_dev:.4f}')
Logistic Regression	Cross-Validation Accuracy: 0.3328 ± 0.8259	
KNN	Cross-Validation Accuracy: 0.2401 ± 0.0884	
Decision Tree	Cross-Validation Accuracy: 0.2648 ± 0.8118	
Random Forest	Cross-Validation Accuracy: 0.3282 ± 0.8159	
Gradient Boosting	Cross-Validation Accuracy: 0.3163 ± 0.8190	
Naive Bayes	Cross-Validation Accuracy: 0.2088 ± 0.8249	

Next, let's find out all the Performance Metrics:

```
In [16]: from sklearn.metrics import ConfusionMatrixDisplay

def evaluate_metrics(y_true, y_pred):
    accuracy = accuracy_score(y_true, y_pred)
    precision = precision_score(y_true, y_pred, average='weighted')
    recall = recall_score(y_true, y_pred, average='weighted')
    f1 = f1_score(y_true, y_pred, average='weighted')
    return accuracy, precision, recall, f1

# Create a dictionary to store model names and their predictions
model_predictions = {
    'Logistic Regression': y_pred_lr,
    'KNN': y_pred_knn,
    'Decision Tree': y_pred_dt,
    'Random Forest': y_pred_rf,
    'Gradient Boosting': y_pred_gb,
    'Naive Bayes': y_pred_nb,
}

# Evaluate each model
for model_name, y_pred in model_predictions.items():
    accuracy, precision, recall, f1 = evaluate_metrics(y_test, y_pred)
    print(f'{model_name} -')
    print(f'Accuracy: {accuracy:.4f}')
    print(f'Precision: {precision:.4f}')
    print(f'Recall: {recall:.4f}')
    print(f'F1-Score: {f1:.4f}')
    print()
```

Logistic Regression:
Accuracy: 0.3478
Precision: 0.3989
Recall: 0.3478
F1-Score: 0.3118

KNN:
Accuracy: 0.2098
Precision: 0.2434
Recall: 0.2098
F1-Score: 0.1780

Decision Tree:
Accuracy: 0.2629
Precision: 0.2475
Recall: 0.2629
F1-Score: 0.2524

Random Forest:
Accuracy: 0.3286
Precision: 0.3941
Recall: 0.3286
F1-Score: 0.2983

Gradient Boosting:
Accuracy: 0.3209
Precision: 0.3152
Recall: 0.3209
F1-Score: 0.2991

Naive Bayes:
Accuracy: 0.2182
Precision: 0.2185
Recall: 0.2182
F1-Score: 0.2026

The Logistic Regression Model has the highest mean cross validation accuracy, followed closely by the Random Forest Model and the Gradient Boosting Model.

In terms of the Performance Metrics, the Logistic Regression has the highest accuracy and the highest F1 score. It also has a fairly good precision and recall. Random Forest and the Gradient Boosting Model follow but they lag behind the Logistic Regression Model in terms of the F1 score and the precision. Decision Tree and the Naive Bayes are in the middle but don't outperform the Logistic Regression Model.

The KNN Model performs badly in the Performance Metrics and the cross-validation accuracy score.

The Logistic Regression Model is the best model considering both the cross-validation accuracy and test set performance metrics. It consistently performs well across different evaluation criteria.

06. Key Findings and Insights

Feature Importance

```
In [17]: import numpy as np

feature_names = vectorizer.get_feature_names_out()
coefs = lr_model.coef_

for i, class_label in enumerate(lr_model.classes_):
    top5 = np.argsort(coefs[i])[0:5]
    print(f"Top words for {class_label}: ('', '.join([feature_names[j] for j in top5]))")

Top words for anger: mcfly, hard, respond, people, damned, sheesh, annoying, know, right, fault
Top words for boredom: getting, cleaning, stuck, ix, hell, sitting, minutes, tired, boring, bored
Top words for empty: lights, sleep, messages, stinks, vista, something, inside, boredom, change, bored
Top words for enthusiasm: awesome, great, first, drank, wants, hang, want, hope, hey, wait
Top words for fun: night, playing, hahaha, wait, look, fans, haha, funny, lol, fun
Top words for happiness: hahaha, thanks, hana, fun, woohoo, awesome, excellent, great, excited, happy
Top words for hate: angry, fucking, damn, fuck, suck, shit, stupid, hates, sucks, hate
Top words for love: best, mothers, amazing, happy, loves, cute, lovely, loving, loved, love
Top words for neutral: suggest, hog, grass, download, lmo, how, ask, now, weeks, gut
Top words for relief: thanx, best, relax, thanx, done, sleep, thanks, fmg, glad, finally
Top words for sadness: depressed, hurt, disappointed, missed, missing, sorry, cry, sadly, miss, sad
Top words for surprise: wow, thought, wtf, surprised, surprisingly, oh, believe, surprise, omg, wow
Top words for worry: hoping, hands, scared, sick, hurts, hurt, poor, sorry, worried, sad
```

The above shows the most common words used to predict a certain emotion.

```
In [42]: from sklearn.metrics import ConfusionMatrixDisplay
import matplotlib.pyplot as plt

fig = plt.figure(figsize=(12, 10))
cm = ConfusionMatrixDisplay.from_estimator(lr_model, X_test, y_test, cmap='viridis')
plt.title('Confusion Matrix - Logistic Regression')
plt.xlabel('Predicted')
plt.ylabel('True')
plt.show()
```

<Figure size 1200x900 with 8 Axes>

Using the above we can see how many true labels the model was able to predict.

```
In [43]: import seaborn as sns
from sklearn.metrics import confusion_matrix

cm = confusion_matrix(y_test, y_pred_lr)
sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', xticklabels=emotion_classes, yticklabels=emotion_classes)
plt.title('Confusion Matrix Heatmap - Logistic Regression')
plt.xlabel('Predicted')
plt.ylabel('True')
plt.show()
```