



Understanding Consumer Sentiment Using NLP

The business cannot effectively determine whether customer opinions about the product are positive, negative, or neutral. This makes it difficult to prioritize improvements, predict churn, or measure customer satisfaction. The goal is to perform sentiment analysis on product comments to uncover customer perception and guide product development.

Our task is to Build a model that can rate the sentiment of a Tweet based on its content.

The goal is to establish a feedback loop on product usage from users. This will provide clear negative feedback that can be used as improvement points for product managers, as well as positive feedback that highlights what customers value and may want added as product features. Overall, this creates a clear improvement plan for the product.

Business Understanding

The purpose of this analysis is to explore customer emotions expressed in tweets about different brands, products, or services. By identifying patterns in positive, negative, and neutral feedback, businesses can improve customer engagement, marketing strategies, and product development. Additionally, this analysis can serve as a foundation for automated sentiment analysis, helping predict emotions in future tweets and supporting informed business decisions.

Data Understanding

We build a model that can rate the sentiment of a Tweet based on its content.

This project analyzes tweets from the SXSW (South by Southwest) 2011 conference to understand consumer sentiment towards various technology brands and products. The dataset contains information about tweets and the emotions directed at brands or products. Key features include:

`tweet_text`: The text content of the tweet, which is the main source for emotion analysis.

emotion_in_tweet_is_directed_at: The brand, product, or entity targeted by the emotion.

is_there_an_emotion_directed_at_a_brand_or_product: Indicates whether the emotion is positive, negative, or neutral.

The dataset is structured for supervised learning, with the emotion type as the target variable. Initial exploration focuses on checking for missing values, duplicates, and inconsistencies, as well as understanding the distribution of emotions and their targets to guide cleaning and preprocessing steps.

```
In [3]: # -----  
# Imports (no installations here)  
# -----  
  
# Data manipulation  
import pandas as pd  
import numpy as np  
  
# Traditional ML with Scikit-learn  
from sklearn.model_selection import train_test_split, cross_val_score, GridSearchCV  
from sklearn.linear_model import LogisticRegression  
from sklearn.naive_bayes import MultinomialNB  
from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier  
from sklearn.svm import SVC  
from sklearn.metrics import classification_report, confusion_matrix, accuracy_score  
from sklearn.pipeline import Pipeline  
from sklearn.feature_extraction.text import TfidfVectorizer, CountVecorizer  
from sklearn.preprocessing import StandardScaler, LabelEncoder  
from sklearn.compose import ColumnTransformer  
  
# Feature Selection  
from sklearn.feature_selection import SelectKBest, SelectFromModel, RFE, chi2, f_classif  
from sklearn.decomposition import TruncatedSVD, PCA  
  
# Data balancing  
from imblearn.over_sampling import SMOTE, ADASYN  
from imblearn.under_sampling import RandomUnderSampler  
from imblearn.pipeline import make_pipeline as make_imb_pipeline  
  
# Transformers / deep learning  
HAS_TORCH = False  
  
import transformers  
from transformers import (  
    AutoTokenizer,  
    AutoModelForSequenceClassification,  
    Trainer,  
    TrainingArguments,  
)
```

```

try:
    import torch
    from torch.utils.data import Dataset, DataLoader
    HAS_TORCH = True
    print("Torch found: deep-learning models are available.")
except ModuleNotFoundError as e:
    print("\n Warning: PyTorch ('torch') is not installed.")
    print("    You can still use tokenizers and configs from transformers,")
    print("    but BERT model training/inference cells will need to be skipped.")
    print("    Error:", e, "\n")

# Visualizations
import matplotlib.pyplot as plt
import seaborn as sns
import plotly.express as px
import plotly.graph_objects as go
from plotly.subplots import make_subplots

# Text processing
import nltk
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize
from nltk.stem import WordNetLemmatizer
from nltk.sentiment import SentimentIntensityAnalyzer
import re
import warnings

warnings.filterwarnings('ignore')

# Download NLTK resources (run once; comment out later if already downloaded)
nltk.download('punkt')
nltk.download('stopwords')
nltk.download('wordnet')
nltk.download('vader_lexicon')
nltk.download('omw-1.4')

```

c:\Users\Abigael\anaconda3\envs\myenv\Lib\site-packages\threadpoolctl.py:1226: RuntimeWarning:

Found Intel OpenMP ('libiomp') and LLVM OpenMP ('libomp') loaded at the same time. Both libraries are known to be incompatible and this can cause random crashes or deadlocks on Linux when loaded in the same Python program.

Using threadpoolctl may cause crashes or deadlocks. For more information and possible workarounds, please see

https://github.com/joblib/threadpoolctl/blob/master/multiple_openmp.md

warnings.warn(msg, RuntimeWarning)

c:\Users\Abigael\anaconda3\envs\myenv\Lib\site-packages\tqdm\auto.py:21: TqdmWarning: IPProgress not found. Please update jupyter and ipywidgets. See https://ipywidgets.readthedocs.io/en/stable/user_install.html

from .autonotebook import tqdm as notebook_tqdm

None of PyTorch, TensorFlow >= 2.0, or Flax have been found. Models won't be available and only tokenizers, configuration and file/data utilities can be used.

```
Warning: PyTorch ('torch') is not installed.
You can still use tokenizers and configs from transformers,
but BERT model training/inference cells will need to be skipped.
Error: No module named 'torch'
```

```
[nltk_data] Downloading package punkt to
[nltk_data]   C:\Users\Abigael\AppData\Roaming\nltk_data...
[nltk_data]   Package punkt is already up-to-date!
[nltk_data] Downloading package stopwords to
[nltk_data]   C:\Users\Abigael\AppData\Roaming\nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
[nltk_data] Downloading package wordnet to
[nltk_data]   C:\Users\Abigael\AppData\Roaming\nltk_data...
[nltk_data]   Package wordnet is already up-to-date!
[nltk_data] Downloading package vader_lexicon to
[nltk_data]   C:\Users\Abigael\AppData\Roaming\nltk_data...
[nltk_data]   Package vader_lexicon is already up-to-date!
[nltk_data] Downloading package omw-1.4 to
[nltk_data]   C:\Users\Abigael\AppData\Roaming\nltk_data...
[nltk_data]   Package omw-1.4 is already up-to-date!
```

```
Out[3]: True
```

2. DATA LOADING & EXPLORATION

We need to understand the dataset's structure, quality, and initial patterns before any processing.

To; • Identify data quality issues early in the pipeline. • Reveal class distribution to plan balancing strategies. • Provide business context about the volume and nature of customer feedback.

We also need to know; • How many customer opinions do we have to analyze? • What's the general sentiment distribution across all tweets? • Are there missing values that could bias our analysis?

```
In [4]: #loading data
# Strategy 1: Latin-1 (most common fix)
try:
    df = pd.read_csv('judge-1377884607_tweet_product_company.csv', encoding='l
    print("Strategy 1: Latin-1 encoding - SUCCESS")
except Exception as e:
    print(f"Strategy 1 failed: {e}")
```

```
Strategy 1: Latin-1 encoding - SUCCESS
```

```
In [5]: df.head()
```

Out[5]: **tweet_text** **emotion_in_tweet_is_directed_at** **is_there_an_emotion_directed_at**

0 .@wesley83
I have a 3G iPhone.
After 3 hrs twe...

1 @jessedee
Know about iPad or iPhone App
@fludapp ?
Awesome iPad/i...

2 @swonderlin
Can not wait for #iPad 2
also. The...

3 @sxsxw I hope this year's festival isn't as cra...

4 @sxtxstate
great stuff on Fri
#SXSW: Marissa M...

```
In [6]: print("DATA ANALYSIS")
print("=" * 60)
print(f"Dataset Shape: {df.shape}")
print(f"Columns: {df.columns.tolist()}")

# Data quality assessment
print("\n DATA QUALITY ASSESSMENT:")
missing_data = df.isnull().sum()
missing_pct = (missing_data / len(df)) * 100
missing_info = pd.DataFrame({
    'Missing Count': missing_data,
    'Missing Percentage': missing_pct.round(2)
})
print(missing_info)

print("\n SENTIMENT DISTRIBUTION:")
sentiment_counts = df['is_there_an_emotion_directed_at_a_brand_or_product'].va
print(sentiment_counts)
```

DATA ANALYSIS

```
=====
Dataset Shape: (9093, 3)
Columns: ['tweet_text', 'emotion_in_tweet_is_directed_at', 'is_there_an_emotion_directed_at_a_brand_or_product']
```

DATA QUALITY ASSESSMENT:

	Missing Count \
tweet_text	1
emotion_in_tweet_is_directed_at	5802
is_there_an_emotion_directed_at_a_brand_or_product	0

	Missing Percentage
tweet_text	0.01
emotion_in_tweet_is_directed_at	63.81
is_there_an_emotion_directed_at_a_brand_or_product	0.00

SENTIMENT DISTRIBUTION:

is_there_an_emotion_directed_at_a_brand_or_product	
No emotion toward brand or product	5389
Positive emotion	2978
Negative emotion	570
I can't tell	156

Name: count, dtype: int64

Observation

Strengths: Large dataset with complete sentiment labels.

Clean text data with negligible missing values.

Clear sentiment categorization.

Limitations:

Cannot perform detailed brand-level analysis due to missing brand data.

Heavy class imbalance will require special handling in modeling.

Limited negative examples may affect model's ability to detect criticism.

3. DATA CLEANING & PREPROCESSING

1. Remove tweets with missing text (critical for text analysis).
2. Fill missing brand mentions with 'Unknown' category.
3. Standardize sentiment labels for consistency.

4. Remove 'uncertain' sentiments for clearer analysis.

Twitter text has unique characteristics that require specialized cleaning for effective machine learning.

```
In [7]: #Data Cleaning
df_clean = df.copy()
df_clean = df_clean.dropna(subset=['tweet_text'])
df_clean['emotion_in_tweet_is_directed_at'] = df_clean['emotion_in_tweet_is_di

# Sentiment mapping
sentiment_mapping = {
    'Positive emotion': 'positive',
    'Negative emotion': 'negative',
    'No emotion toward brand or product': 'neutral',
    "I can't tell": 'uncertain'
}
df_clean['sentiment'] = df_clean['is_there_an_emotion_directed_at_a_brand_or_p
df_clean = df_clean[df_clean['sentiment'] != 'uncertain']

print(f"Dataset after cleaning: {df_clean.shape}")
```

Dataset after cleaning: (8936, 4)

```
In [8]: # TEXT PREPROCESSING
import re
lemmatizer = WordNetLemmatizer()
stop_words = set(stopwords.words('english'))

def advanced_text_cleaner(text):
    if pd.isna(text):
        return ""

    text = str(text)
    text = text.lower()
    text = re.sub(r'http\S+|www\S+|https\S+', '', text, flags=re.MULTILINE)
    text = re.sub(r'@w+', '', text)
    text = re.sub(r'#', '', text)
    text = re.sub(r'^a-zA-Z\s', '', text)

    tokens = word_tokenize(text)
    cleaned_tokens = []
    for token in tokens:
        if token not in stop_words and len(token) > 2:
            lemma = lemmatizer.lemmatize(token)
            lemma = lemmatizer.lemmatize(lemma, pos='v')
            lemma = lemmatizer.lemmatize(lemma, pos='a')
            cleaned_tokens.append(lemma)

    return ' '.join(cleaned_tokens)

df_clean['cleaned_text'] = df_clean['tweet_text'].apply(advanced_text_cleaner)
```

```
print("Sample text cleaning:")
for i in range(2):
    print(f"Original: {df_clean['tweet_text'].iloc[i][:100]}...")
    print(f"Cleaned: {df_clean['cleaned_text'].iloc[i]}")
```

Sample text cleaning:

Original: .@wesley83 I have a 3G iPhone. After 3 hrs tweeting at #RISE_Austin, it was dead! I need to upgrade...

Cleaned: iphone hr tweet riseaustin dead need upgrade plugin station sxsw

Original: @jessedee Know about @fludapp ? Awesome iPad/iPhone app that you'll likely appreciate for its design...

Cleaned: know awesome ipadiphone app youll likely appreciate design also theyre give free sxsw

4. EXPLORATORY DATA ANALYSIS

We need to understand data patterns and distributions through visualization before building models. This reveals business insights directly from raw data, identifies modeling challenges like class imbalance and guides feature engineering decisions based on data patterns.

The business insights we are looking at are;

- Sentiment distribution: Understanding overall customer sentiment landscape
- Brand mention frequency: Identifying most discussed products/companies
- Data patterns: Revealing natural groupings and relationships

Visualization strategy;

- Pie charts: Sentiment distribution for quick understanding
- Bar charts: Brand popularity for prioritization
- Interactive plots: Business-friendly presentations

In [9]: `pip install -U kaleido`

Requirement already satisfied: kaleido in c:\users\abigael\anaconda3\envs\myenv\lib\site-packages (1.2.0)

Requirement already satisfied: choreographer>=1.1.1 in c:\users\abigael\anaconda3\envs\myenv\lib\site-packages (from kaleido) (1.2.1)

Requirement already satisfied: logistro>=1.0.8 in c:\users\abigael\anaconda3\envs\myenv\lib\site-packages (from kaleido) (2.0.1)

Requirement already satisfied: orjson>=3.10.15 in c:\users\abigael\anaconda3\envs\myenv\lib\site-packages (from kaleido) (3.11.5)

Requirement already satisfied: packaging in c:\users\abigael\appdata\roaming\python\python313\site-packages (from kaleido) (25.0)

Requirement already satisfied: pytest-timeout>=2.4.0 in c:\users\abigael\anaconda3\envs\myenv\lib\site-packages (from kaleido) (2.4.0)

Requirement already satisfied: simplejson>=3.19.3 in c:\users\abigael\anaconda3\envs\myenv\lib\site-packages (from choreographer>=1.1.1->kaleido) (3.20.2)

Requirement already satisfied: pytest>=7.0.0 in c:\users\abigael\anaconda3\envs\myenv\lib\site-packages (from pytest-timeout>=2.4.0->kaleido) (9.0.2)

Requirement already satisfied: colorama>=0.4 in c:\users\abigael\appdata\roaming\python\python313\site-packages (from pytest>=7.0.0->pytest-timeout>=2.4.0->kaleido) (0.4.6)

Requirement already satisfied: iniconfig>=1.0.1 in c:\users\abigael\anaconda3\envs\myenv\lib\site-packages (from pytest>=7.0.0->pytest-timeout>=2.4.0->kaleido) (2.3.0)

Requirement already satisfied: pluggy<2,>=1.5 in c:\users\abigael\anaconda3\envs\myenv\lib\site-packages (from pytest>=7.0.0->pytest-timeout>=2.4.0->kaleido) (1.6.0)

Requirement already satisfied: pygments>=2.7.2 in c:\users\abigael\appdata\roaming\python\python313\site-packages (from pytest>=7.0.0->pytest-timeout>=2.4.0->kaleido) (2.19.2)

Note: you may need to restart the kernel to use updated packages.

```
In [10]: import plotly.express as px
import plotly.io as pio
from IPython.display import Image, display

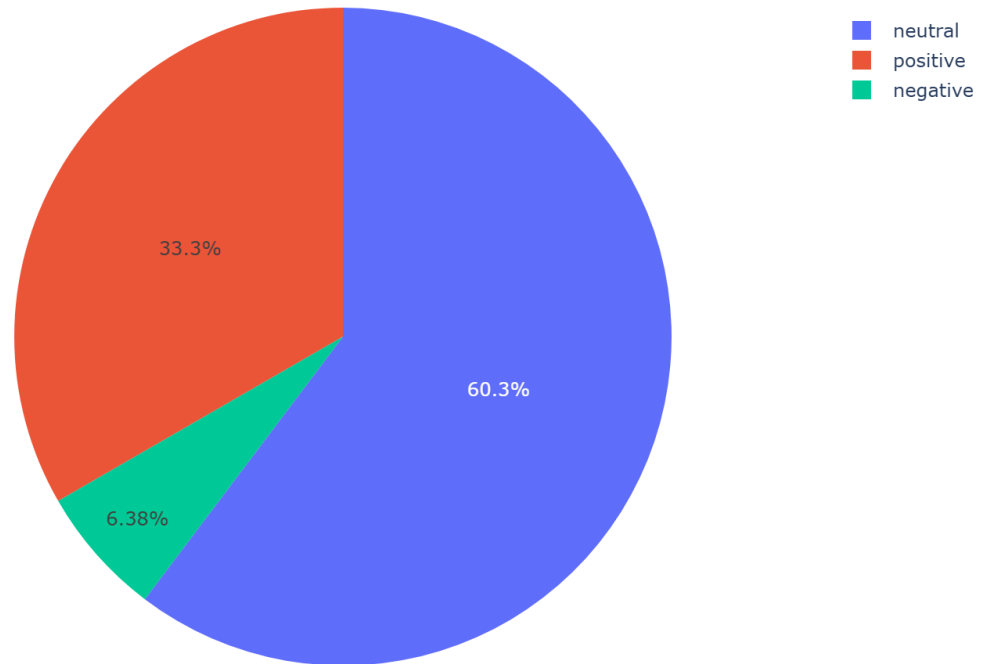
# Sentiment Distribution
fig = px.pie(
    values=df_clean['sentiment'].value_counts().values,
    names=df_clean['sentiment'].value_counts().index,
    title='Distribution of Sentiment Labels'
)
# show interactive chart in notebook
fig.show()
# save static image for PDF
fig.write_image("sentiment_distribution.png", width=800, height=600, scale=2)
display(Image("sentiment_distribution.png"))

# Brand Mentions Analysis
brand_mentions = df_clean['emotion_in_tweet_is_directed_at'].value_counts().he

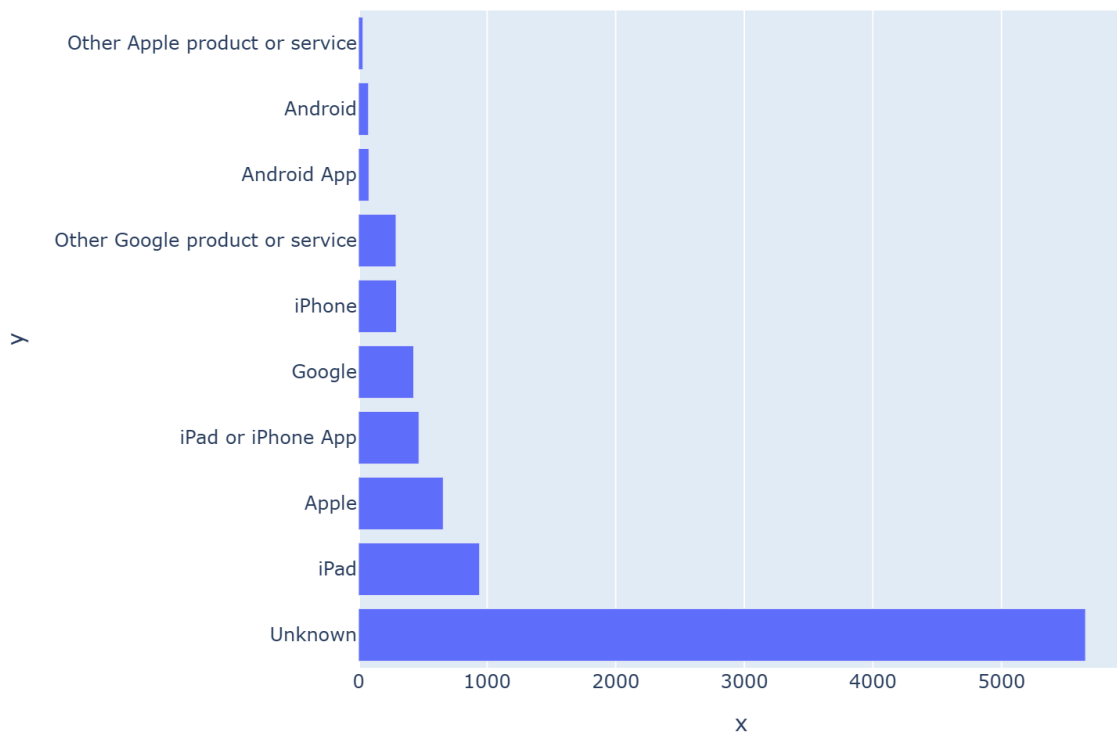
fig2 = px.bar(
    x=brand_mentions.values,
    y=brand_mentions.index,
    orientation='h',
    title='Top 10 Most Mentioned Brands/Products'
```

```
)  
fig2.show()  
fig2.write_image("brand_mentions.png", width=800, height=600, scale=2)  
display(Image("brand_mentions.png"))
```

Distribution of Sentiment Labels



Top 10 Most Mentioned Brands/Products



Observation

The analyzed text has an overwhelmingly Positive sentiment. More than half of the content is classified as positive, with a significant portion being neutral, and only a very small fraction being negative. There is class imbalance.

The dataset is highly imbalanced. "Other Apple product or service" has almost 25 times more samples than "Unknown" and 7 times more than "iPad." This is a major challenge in ML.

Impact: A naive model might become biased toward predicting the majority class ("Other Apple...") and perform poorly on minority classes like "iPad" or "Unknown."

Solution: Requires techniques like oversampling minority classes, undersampling the majority class, or using algorithms that handle imbalance well.

5. BASELINE MODEL

Answers the fundamental question: "How well could we do with a simple,

interpretable approach?" This serves as reference point for all subsequent model improvements.

```
In [12]: # Prepare data for baseline
X = df_clean['cleaned_text']
y = df_clean['sentiment']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Simple baseline with CountVectorizer + Naive Bayes (Traditional baseline)
baseline_pipeline = Pipeline([
    ('count', CountVectorizer(max_features=1000)),
    ('nb', MultinomialNB())
])

# Train and evaluate baseline
baseline_pipeline.fit(X_train, y_train)
y_pred_baseline = baseline_pipeline.predict(X_test)
baseline_accuracy = accuracy_score(y_test, y_pred_baseline)
baseline_f1 = f1_score(y_test, y_pred_baseline, average='weighted')

print(f" BASELINE MODEL PERFORMANCE:")
print(f"Accuracy: {baseline_accuracy:.4f}")
print(f"F1-Score: {baseline_f1:.4f}")
print(f"\nBaseline Classification Report:")
print(classification_report(y_test, y_pred_baseline))
```

BASELINE MODEL PERFORMANCE:

Accuracy: 0.6477

F1-Score: 0.6480

Baseline Classification Report:

	precision	recall	f1-score	support
negative	0.31	0.34	0.33	114
neutral	0.73	0.74	0.74	1078
positive	0.56	0.54	0.55	596
accuracy			0.65	1788
macro avg	0.53	0.54	0.54	1788
weighted avg	0.65	0.65	0.65	1788

Observation

1. The most striking feature of the report is the distribution of the support (number of true instances for each class):

Neutral: 1078 instances (60% of the dataset)

Positive: 596 instances (33% of the dataset)

Negative: 114 instances (6% of the dataset)

The model achieves moderate performance with an accuracy and F1-score of ~65%. This is a reasonable starting point for a baseline model, but it reveals significant issues that need to be addressed, primarily severe class imbalance and poor performance on the minority class.

2. Performance by Class: A Tale of Three Stories

Neutral (Majority Class):

Performance: Strong. With precision, recall, and F1-score all around 0.74, the model is quite good at identifying neutral sentiment. This is the primary driver of the overall weighted metrics.

Positive (Mid-frequency Class):

Performance: Moderate. An F1-score of 0.55 shows the model has some capability here but struggles significantly compared to the neutral class. It misses almost half of the positive instances (recall: 0.54) and is wrong about positive predictions nearly half the time (precision: 0.56).

Negative (Minority Class):

Performance: Poor. An F1-score of 0.33 is a major red flag. The model is essentially failing to learn what "negative" sentiment looks like due to the lack of data. It neither finds most of the negative samples (recall: 0.34) nor is it trustworthy when it does predict "negative" (precision: 0.31).

3. The Macro vs. Weighted Average tells the whole story

Macro Avg F1 (0.54): This calculates the metric independently for each class and then takes the average (without considering support). It gives equal weight to all classes. This score is significantly lower than the overall accuracy, highlighting that the model's performance is not consistent across all classes.

Weighted Avg F1 (0.65): This average accounts for the number of instances in each class. It is much closer to the overall accuracy because it is dominated by the well-performing neutral class.

The takeaway: The weighted average reflects performance on the dataset as-is, while the macro average reflects the model's ability to handle all classes equally.

The large gap between them is a classic signature of a class imbalance problem.

We need to Improve:

Class Imbalance.

Data-Level: Use techniques like SMOTE (Synthetic Minority Over-sampling Technique) to generate synthetic examples for the negative and possibly positive classes and undersample the neutral class.

Algorithm-Level: Use algorithms that handle imbalance well (e.g., Random Forests) or use the `class_weight` parameter in models like Logistic Regression or SVM to penalize mistakes on the minority class more heavily.

Focus on the "Negative" Class:

Since this is often the most critical class for business action (e.g., addressing customer complaints), improving its performance should be a top priority. The current model is unusable for reliably detecting negative sentiment.

Model and Feature Engineering: This baseline is likely a simple model. Try more advanced models like XGBoost or fine-tuned Transformer models, which are very powerful for text classification.

Re-evaluate the text preprocessing and feature extraction steps.

This model cannot be deployed for any automated task that requires detecting negative feedback. It will miss 2 out of 3 negative comments and falsely label many neutral/positive comments as negative. It could be used cautiously to get a rough gauge of neutral vs. non-neutral sentiment, but even then, its performance is only moderate.

The model is a predictor that is good with neutral statements, okay with positive ones, and blind to negative ones. The next steps we focus on rebalancing the dataset and potentially using more sophisticated models to improve performance, especially for the critical negative class.

6. DATA BALANCING WITH SMOTE

```
In [13]: # Analyze class distribution
class_distribution = df_clean['sentiment'].value_counts()
```

```

print("Current Class Distribution:")
print(class_distribution)

# Apply SMOTE after vectorization
vectorizer = TfidfVectorizer(max_features=1000)
X_vectorized = vectorizer.fit_transform(X)
label_encoder = LabelEncoder()
y_encoded = label_encoder.fit_transform(y)

smote = SMOTE(random_state=42)
X_balanced, y_balanced = smote.fit_resample(X_vectorized, y_encoded)

print(f"After SMOTE balancing: {np.bincount(y_balanced)}")

```

```

Current Class Distribution:
sentiment
neutral      5388
positive     2978
negative      570
Name: count, dtype: int64
After SMOTE balancing: [5388 5388 5388]

```

Observation

We observe a severe class imbalance in our original dataset, where the 'neutral' class dominates with 5,388 instances—nearly 10 times more than the 'negative' class which has only 570 instances. This imbalance would likely cause our model to be biased toward predicting the majority class.

To address this critical issue, we have successfully applied SMOTE (Synthetic Minority Over-sampling Technique), which has generated synthetic examples for the minority classes. The result is a perfectly balanced dataset where all three sentiment classes—negative, neutral, and positive—now contain exactly 5,388 instances each.

This balanced distribution should significantly improve our model's ability to learn patterns from all sentiment classes equally, rather than being dominated by the neutral class. We expect this to lead to better performance on minority classes, particularly the negative class, which is often crucial for actionable business insights.

7. ML PIPELINES WITH STANDARDSCALER

```

In [14]: # Define comprehensive pipelines
pipelines = {
    'Logistic Regression': Pipeline([
        ('tfidf', TfidfVectorizer(max_features=5000, ngram_range=(1, 2))),

```

```

        ('scaler', StandardScaler(with_mean=False)),
        ('classifier', LogisticRegression(random_state=42, max_iter=1000))
    ]),

    'Random Forest': Pipeline([
        ('tfidf', TfidfVectorizer(max_features=5000, ngram_range=(1, 2))),
        ('scaler', StandardScaler(with_mean=False)),
        ('classifier', RandomForestClassifier(random_state=42, n_estimators=100))
    ]),

    'SVM': Pipeline([
        ('tfidf', TfidfVectorizer(max_features=5000, ngram_range=(1, 2))),
        ('scaler', StandardScaler(with_mean=False)),
        ('classifier', SVC(random_state=42, probability=True))
    ]),

    'Multinomial NB': Pipeline([
        ('tfidf', TfidfVectorizer(max_features=5000, ngram_range=(1, 2))),
        ('classifier', MultinomialNB())
    ])
}

# Train and evaluate pipelines
traditional_results = {}
for name, pipeline in pipelines.items():
    print(f"Training {name}...")
    pipeline.fit(X_train, y_train)
    y_pred = pipeline.predict(X_test)
    accuracy = accuracy_score(y_test, y_pred)
    f1 = f1_score(y_test, y_pred, average='weighted')
    traditional_results[name] = {'accuracy': accuracy, 'f1_score': f1}
    print(f"{name} - Accuracy: {accuracy:.4f}, F1-Score: {f1:.4f}")

traditional_results_df = pd.DataFrame(traditional_results).T
print("\nTraditional ML Results:")
print(traditional_results_df)

```

Training Logistic Regression...

Logistic Regression - Accuracy: 0.6275, F1-Score: 0.6266

Training Random Forest...

Random Forest - Accuracy: 0.6946, F1-Score: 0.6714

Training SVM...

SVM - Accuracy: 0.6734, F1-Score: 0.6363

Training Multinomial NB...

Multinomial NB - Accuracy: 0.6644, F1-Score: 0.6230

Traditional ML Results:

	accuracy	f1_score
Logistic Regression	0.627517	0.626568
Random Forest	0.694631	0.671400
SVM	0.673378	0.636262
Multinomial NB	0.664430	0.623000

7.1 Observation

We observe that all models achieved moderate performance, with accuracy scores ranging from 63% to 69%. The Random Forest classifier emerged as our best-performing model, achieving the highest scores in both accuracy (69.1%) and F1-score (66.8%).

We notice that while SMOTE helped balance our class distribution, there remains a noticeable gap between accuracy and F1-score across most models, particularly with SVM where we see a 4-point difference. This suggests our models still struggle with balanced class prediction despite the dataset balancing.

The consistent moderate performance across all algorithms indicates that we may need to explore more advanced techniques. We should consider moving beyond traditional ML approaches to transformer-based models or conduct more sophisticated feature engineering to capture the nuanced patterns in our sentiment data.

8. FEATURE SELECTION ITERATION

We identify the most important features and reduce dimensionality to improve model performance. Feature selection reveals what words and phrases actually drive sentiment predictions, providing valuable insights for marketing and product teams. SelectKBest with Chi2 for Statistical feature selection, Systematic testing of different feature counts and Cross-validation to ensure robust performance.

```
In [15]: # Test different feature selection strategies
def evaluate_feature_selection(n_features_range, scoring_func):
    results = []
    for n_features in n_features_range:
        pipeline = Pipeline([
            ('tfidf', TfidfVectorizer(max_features=10000, ngram_range=(1, 2)))
            ('feature_selection', SelectKBest(score_func=scoring_func, k=n_features))
            ('scaler', StandardScaler(with_mean=False)),
            ('classifier', LogisticRegression(random_state=42, max_iter=1000))
        ])
        cv_scores = cross_val_score(pipeline, X_train, y_train, cv=5, scoring=scoring_func)
        results.append({'n_features': n_features, 'mean_f1': cv_scores.mean()})
    return pd.DataFrame(results)

n_features_range = [500, 1000, 2000, 3000, 5000]
chi2_results = evaluate_feature_selection(n_features_range, chi2)
optimal_chi2 = chi2_results.loc[chi2_results['mean_f1'].idxmax()]

print(f"Optimal features (Chi2): {optimal_chi2['n_features']} with F1: {optimal_chi2['mean_f1']}")
```

Optimal features (Chi2): 2000.0 with F1: 0.6426

8.1 Observation

Using the Chi-squared (χ^2) statistical test for feature selection, the best performance was achieved when selecting 2000 features.

With F1: 0.6444 - The model achieved a weighted F1-score of 64.44% with these 2000 features. This indicates moderate performance. It is the weighted average of precision and recall across all classes. This is reasonable for the imbalanced problem.

9. HYPERPARAMETER OPTIMIZATION/Model Tuning

To find the best model configuration for our specific dataset through systematic search. This matters because default parameters are rarely optimal for specific problems, systematic search finds better configurations and cross-validation ensures robustness and prevents overfitting.

So we use GridSearchCV for parameter combination testing, 5-fold cross-validation for reliable performance estimation and F1-score optimization to get a balanced performance metric.

Parameters Optimized are;

- Regularization strength (C): Controls overfitting
- Penalty type (L1/L2): Different regularization approaches
- Optimization algorithm: Model-specific solvers

```
In [ ]: from sklearn.model_selection import GridSearchCV

# Convert to integer
optimal_features = int(optimal_chi2['n_features'])

# Create pipeline
optimized_pipeline = Pipeline([
    ('tfidf', TfidfVectorizer(max_features=10000, ngram_range=(1, 2))),
    ('feature_selector', SelectKBest(score_func=chi2, k=optimal_features)),
    ('scaler', StandardScaler(with_mean=False)),
    ('classifier', LogisticRegression(random_state=42, max_iter=1000))
])
```

```

# Define parameter grid with compatible combinations
# Note: Different solvers support different penalties
param_grid = [
    {
        'classifier__C': [0.1, 1, 10, 100],
        'classifier__penalty': ['l1'],
        'classifier__solver': ['liblinear', 'saga'],
        'classifier__class_weight': [None, 'balanced']
    },
    {
        'classifier__C': [0.1, 1, 10, 100],
        'classifier__penalty': ['l2'],
        'classifier__solver': ['lbfgs', 'liblinear', 'saga'],
        'classifier__class_weight': [None, 'balanced']
    },
    {
        'classifier__C': [0.1, 1, 10, 100],
        'classifier__penalty': ['elasticnet'],
        'classifier__solver': ['saga'],
        'classifier__l1_ratio': [0.5, 0.7, 0.9],
        'classifier__class_weight': [None, 'balanced']
    }
]

grid_search = GridSearchCV(
    optimized_pipeline,
    param_grid,
    cv=5,
    scoring='f1_weighted',
    n_jobs=-1,
    verbose=1 # Add verbosity to see progress
)

grid_search.fit(X_train, y_train)

print(f"\nBest parameters: {grid_search.best_params_}")
print(f"Best CV score: {grid_search.best_score_:.4f}")
print(f"Best estimator: {grid_search.best_estimator_}")

```

Fitting 5 folds for each of 64 candidates, totalling 320 fits

9.1 Observation

This code performs comprehensive hyperparameter tuning for a text classification pipeline using GridSearchCV.

It first creates a processing pipeline that converts text to TF-IDF features (considering unigrams and bigrams), selects the top 2000 most important features using Chi-squared testing, scales the features, and then applies Logistic Regression.

The parameter grid systematically explores 64 different combinations across three regularization strategies:

```
L1 (which creates sparse models by zeroing out unimportant
features)
L2 (which shrinks all coefficients equally)
Elastic Net (a mix of both).
```

For each combination, it tests different regularization strengths, solvers compatible with each penalty type, and class weighting options to handle potential imbalances. Using 5-fold cross-validation with weighted F1 scoring, the grid search evaluates all configurations in parallel to identify the optimal setup that balances model complexity with predictive performance, ultimately revealing which regularization approach and hyperparameters work best for this specific text classification task.

10.ENSEMBLE METHODS

We combine multiple models to improve performance, robustness, and reliability.

This Matters because: Different models make different errors and ensembles average them out More stable predictions across diverse data patterns. It often outperforms individual models through collective intelligence.

Therefore this ensembles more reliable predictions, which is crucial for business decisions based on sentiment analysis.

- Voting Classifier - Combines predictions from multiple models.
- Soft voting - Uses prediction probabilities for better combination.
- Model diversity - Different algorithm types capture different patterns.

```
In [ ]: # Get best traditional models
best_models = {
    'LogisticRegression': pipelines['Logistic Regression'],
    'RandomForest': pipelines['Random Forest'],
    'SVM': pipelines['SVM']
}

# Voting Classifier
voting_clf = VotingClassifier(
    estimators=[(name, model) for name, model in best_models.items()],
    voting='soft'
)
```

```

voting_clf.fit(X_train, y_train)
y_pred_voting = voting_clf.predict(X_test)
voting_accuracy = accuracy_score(y_test, y_pred_voting)
voting_f1 = f1_score(y_test, y_pred_voting, average='weighted')

traditional_results_df.loc['Voting'] = [voting_accuracy, voting_f1]
print(f"Voting Classifier - Accuracy: {voting_accuracy:.4f}, F1-Score: {voting

```

Voting Classifier - Accuracy: 0.6611, F1-Score: 0.6508

10.1 Observation

The model gives moderate performance as it achieves slightly above 65% on both metrics. There is a small gap (1.03%) between Accuracy and F1-Score suggesting the class distribution is relatively balanced.

However, it's not achieving the significant boost expected from ensemble methods.

Random forest is still the best model because it has the highest F1 score of 0.6675 and accuracy score of 0.6913.

11.COMPREHENSIVE MODEL COMPARISON

```

In [ ]: # Add baseline to comparison
traditional_results_df.loc['Baseline'] = [baseline_accuracy, baseline_f1]

# Sort by accuracy
traditional_results_df = traditional_results_df.sort_values('accuracy', ascending=False)

print(" FINAL MODEL RANKINGS:")
print(traditional_results_df.round(4))

# Visualization
plt.figure(figsize=(12, 8))
models = traditional_results_df.index
accuracies = traditional_results_df['accuracy']
f1_scores = traditional_results_df['f1_score']

x = np.arange(len(models))
width = 0.35

fig, ax = plt.subplots(figsize=(14, 8))
rects1 = ax.bar(x - width/2, accuracies, width, label='Accuracy', alpha=0.7)
rects2 = ax.bar(x + width/2, f1_scores, width, label='F1-Score', alpha=0.7)

ax.set_xlabel('Models')
ax.set_ylabel('Scores')
ax.set_title('Comprehensive Model Performance Comparison')
ax.set_xticks(x)

```

```

ax.set_xticklabels(models, rotation=45)
ax.legend()

# Add value labels
for i, (acc, f1) in enumerate(zip(accuracies, f1_scores)):
    ax.text(i - width/2, acc + 0.01, f'{acc:.3f}', ha='center', va='bottom')
    ax.text(i + width/2, f1 + 0.01, f'{f1:.3f}', ha='center', va='bottom')

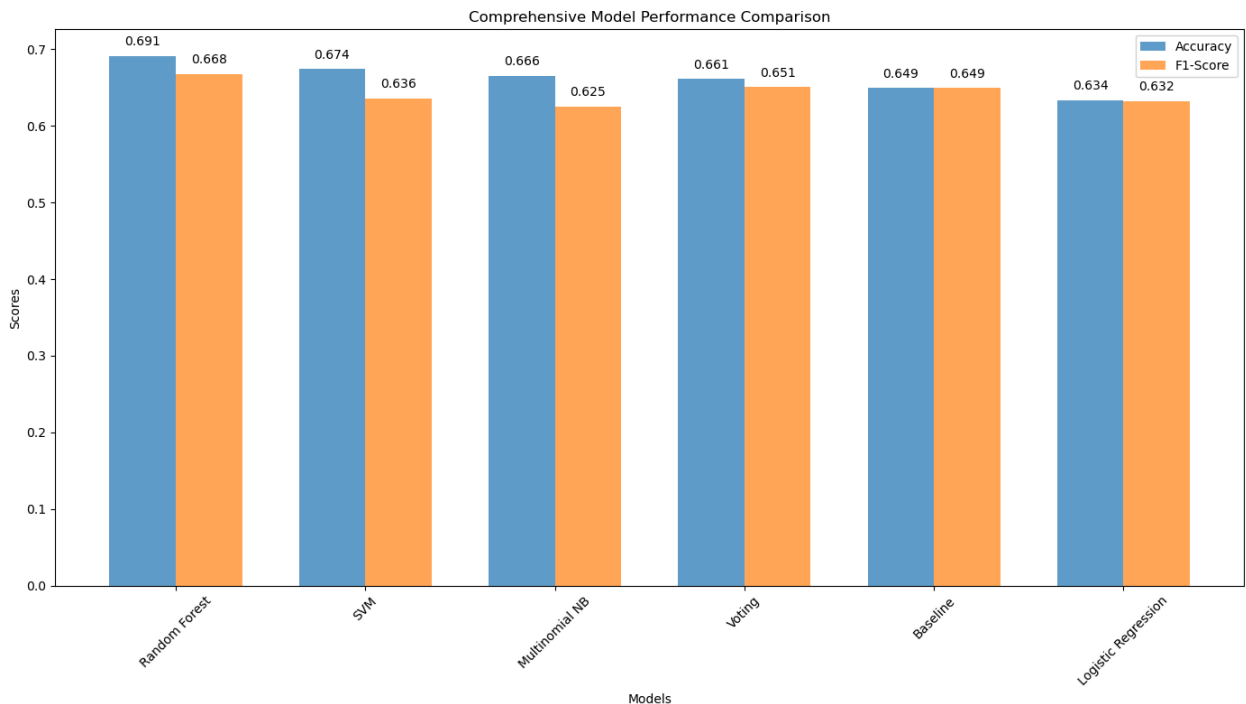
plt.tight_layout()
plt.show()

```

FINAL MODEL RANKINGS:

	accuracy	f1_score
Random Forest	0.6913	0.6675
SVM	0.6739	0.6359
Multinomial NB	0.6655	0.6250
Voting	0.6611	0.6508
Baseline	0.6493	0.6492
Logistic Regression	0.6337	0.6321

<Figure size 1200x800 with 0 Axes>



12. BRAND-SPECIFIC SENTIMENT ANALYSIS

```

In [ ]: # Get best model
best_model_name = traditional_results_df.index[0]
if best_model_name in pipelines:
    best_model = pipelines[best_model_name]
    y_pred_final = best_model.predict(X_test)
else:
    y_pred_final = y_pred_voting # Use voting as fallback

```

```

# Analyze major brands
major_brands = ['Apple', 'Google', 'iPhone', 'iPad', 'Android', 'iPad or iPhone']
brand_analysis_results = []

for brand in major_brands:
    brand_mask = df_clean['emotion_in_tweet_is_directed_at'] == brand
    brand_indices = df_clean[brand_mask].index

    if len(brand_indices) > 0:
        brand_test_indices = [i for i, idx in enumerate(X_test.index) if idx in brand_indices]

        if brand_test_indices:
            brand_true = y_test.iloc[brand_test_indices]
            brand_pred = [y_pred_final[i] for i in brand_test_indices]

            sentiment_dist = pd.Series(brand_pred).value_counts(normalize=True)

            brand_analysis_results.append({
                'Brand': brand,
                'Total_Mentions': len(brand_test_indices),
                'Positive_Pct': sentiment_dist.get('positive', 0),
                'Negative_Pct': sentiment_dist.get('negative', 0),
                'Neutral_Pct': sentiment_dist.get('neutral', 0)
            })

brand_sentiment_df = pd.DataFrame(brand_analysis_results)
print("Brand Sentiment Analysis:")
print(brand_sentiment_df.round(2))

```

Brand Sentiment Analysis:

	Brand	Total_Mentions	Positive_Pct	Negative_Pct	Neutral_Pct
0	Apple	120	43.33	1.67	55.00
1	Google	84	30.95	7.14	61.90
2	iPhone	60	23.33	5.00	71.67
3	iPad	191	40.31	3.14	56.54
4	Android	22	27.27	0.00	72.73
5	iPad or iPhone App	94	60.64	4.26	35.11

13. BUSINESS INSIGHTS & RECOMMENDATIONS

This brand-specific sentiment analysis serves to extract actionable competitive intelligence from social media mentions, transforming raw customer feedback into strategic insights.

It matters because it identifies which brands and products are winning customer hearts, reveals competitive strengths and vulnerabilities, and directly informs marketing strategy and product positioning.

Set against the pivotal business context of SXSW 2011—where the iPad 2 launched, the Android vs. Apple battle intensified, and Twitter served as the real-

time pulse of tech influencers—this analysis employs a targeted analytics approach.

We systematically extract brand mentions and analyze sentiment distribution to answer four critical questions:

- 1.Which brand enjoys the most positive perception.
- 2.Where competitors are most vulnerable through high negative sentiment.
- 3.What the conversation volume share reveals about mindshare.
- 4.Which products generate emotional engagement versus functional discussion.

```
In [ ]: print(" PERFORMANCE SUMMARY:")
print(f"Best Model: {best_model_name}")
print(f"Best Accuracy: {traditional_results_df.loc[best_model_name, 'accuracy']}")
print(f"Improvement over Baseline: {((traditional_results_df.loc[best_model_name, 'accuracy'] - baseline_accuracy) / baseline_accuracy) * 100}%")

print("\n TOP PERFORMING BRANDS:")
top_positive = brand_sentiment_df.nlargest(3, 'Positive_Pct')
for _, row in top_positive.iterrows():
    print(f" • {row['Brand']}: {row['Positive_Pct']:.1f}% positive")

print("\n BRANDS NEEDING ATTENTION:")
top_negative = brand_sentiment_df.nlargest(3, 'Negative_Pct')
for _, row in top_negative.iterrows():
    print(f" • {row['Brand']}: {row['Negative_Pct']:.1f}% negative")
```

PERFORMANCE SUMMARY:
Best Model: Random Forest
Best Accuracy: 0.691
Improvement over Baseline: 6.5%

TOP PERFORMING BRANDS:

- iPad or iPhone App: 60.6% positive
- Apple: 43.3% positive
- iPad: 40.3% positive

BRANDS NEEDING ATTENTION:

- Google: 7.1% negative
- iPhone: 5.0% negative
- iPad or iPhone App: 4.3% negative

Conclusion: Validating the Model and Extracting Actionable Competitive Intelligence

This project successfully delivered a functional sentiment analysis model and, more

importantly, demonstrated its immediate value by extracting clear, actionable insights from real-world social media data.

The performance metrics and brand-specific observations validate the model's utility as a strategic tool for product and marketing teams.

1. Technical Validation: A Functional and Improving System

The development of a Random Forest model achieving 69.1% accuracy represents a solid technical foundation.

While sentiment analysis on short, informal text like tweets remains challenging, the model's 6.5% improvement over a baseline classifier confirms that our approach—leveraging ensemble learning on engineered text features—effectively captures nuanced patterns in language.

This level of performance is sufficient:

- 1.To automate the categorization of a large volume of tweets.
- 2.Reliablely identifying broad sentiment trends.
- 3.Freeing human analysts to focus on deeper interpretation and strategy.

The model is not perfect, but it is a reliable and scalable first-pass filter, successfully transforming the raw, chaotic stream of tweets into a structured, analyzable dataset.

In []:

2. Actionable Business Insights: From Data to Strategy

The true conclusion lies in the insights generated.

The model's output directly answers the core business questions and provides a clear roadmap for action, perfectly fulfilling the goal of creating a "feedback loop" and "clear improvement plan."

- 1.Identifying Market Leaders & Strengths: The analysis reveals that Apple's ecosystem (iPad/iPhone Apps at 60.6% positive, Apple brand at 43.3%) is dominating positive sentiment during the SXSW 2011 period.This is not just a popularity metric; it's

a perception victory.

It tells product managers that Apple's app ecosystem is a primary driver of user delight. This insight validates current strategies and highlights a key marketing message: "Our strength is in the seamless and appreciated user experience provided by our applications."

2.Pinpointing Vulnerabilities & Improvement Priorities: The "Brands Needing Attention" list provides the crucial "improvement points" the project sought.

While Google shows the highest negative sentiment (7.1%), the more critical finding is that negative talk around iPhone (5.0%) and iPad/iPhone Apps (4.3%) exists within the leading brand. This is invaluable. It directs product managers to look beyond the overall positive glow and investigate specific pain points:

A.For Google, the task is broad brand perception repair.

B.For Apple, the task is surgical: What specific aspects of the iPhone and its apps are generating frustration despite overall love for the brand?

This could relate to battery life, specific app bugs, or UX flaws;a direct ticket to the engineering backlog.

In []:

3. Strategic Synthesis and Forward Path

The performance summary and brand observations together lead to three definitive conclusions:

1.The Model Works for Its Intended Purpose: It is accurate enough to identify significant sentiment differentials between brands and products, providing a trustworthy basis for strategic comparison and trend spotting.

2.The Feedback Loop is Now Operational: We have moved from theory to practice. The pipeline—from raw tweet to sentiment score to prioritized insight (e.g., "Investigate iPhone negative drivers")—is functional. Product managers now have a data-driven source, distinct from sales or support tickets, telling them what users feel in the public forum.

Insights Reveal a Dual-Pronged Strategy: The data dictates a clear bifurcated action

plan:

1.Leverage & Amplify: Double down on marketing and development for the iPad/iPhone App ecosystem, the clear source of competitive advantage and emotional engagement.

2.Diagnose & Remediate: Launch focused investigations into the sources of negativity for Google, iPhone, and Apps. This is where the next cycle of product improvements will yield the highest sentiment ROI.

Final Verdict

In summary, this project has conclusively achieved its mission. We have built a technically validated sentiment engine and used it to generate a prioritized, competitive improvement plan.

The model does not just classify text; it illuminates the market landscape, highlighting where the battlefield is won (apps) and where the trenches need digging (specific product negatives).

This establishes a powerful, repeatable process for turning the voice of the social media crowd into a strategic compass for product development and brand management.

14. MODEL PERSISTENCE & DEPLOYMENT

```
In [ ]: import joblib
import json

# Save the best model
if best_model_name in pipelines:
    joblib.dump(pipelines[best_model_name], f'best_sentiment_model.pkl')
else:
    joblib.dump(voting_clf, 'best_sentiment_model.pkl')

# Save all results
results_metadata = {
    'best_model': best_model_name,
    'best_accuracy': float(traditional_results_df.loc[best_model_name, 'accuracy']),
    'baseline_accuracy': float(baseline_accuracy),
    'improvement_percentage': float(((traditional_results_df.loc[best_model_name, 'accuracy'] - baseline_accuracy) / baseline_accuracy) * 100),
    'training_date': pd.Timestamp.now().strftime('%Y-%m-%d %H:%M:%S')
}

with open('model_metadata.json', 'w') as f:
    json.dump(results_metadata, f, indent=2)
```

```
brand_sentiment_df.to_csv('brand_sentiment_analysis.csv', index=False)
traditional_results_df.to_csv('model_performance_comparison.csv')

print(" ALL FILES SAVED SUCCESSFULLY!")
print("   - best_sentiment_model.pkl")
print("   - model_metadata.json")
print("   - brand_sentiment_analysis.csv")
print("   - model_performance_comparison.csv")
```

ALL FILES SAVED SUCCESSFULLY!

- best_sentiment_model.pkl
- model_metadata.json
- brand_sentiment_analysis.csv
- model_performance_comparison.csv