# Part1 Topic Classification

## Load and clean data

In [2]:
```python
import pandas as pd
import numpy as np
import re
import nltk
from sklearn.model_selection import cross_val_score, cross_validate, StratifiedK
from sklearn.feature_extraction.text import TfidfVectorizer, CountVectorizer
from sklearn.feature_extraction import text as sk_text
from sklearn.naive_bayes import BernoulliNB, MultinomialNB
from sklearn.metrics import accuracy_score, classification_report, precision_sco
from nltk.corpus import stopwords, wordnet
from nltk import pos_tag
from nltk.tokenize import word_tokenize
from nltk.stem import PorterStemmer, WordNetLemmatizer
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.pipeline import Pipeline

# Download required NLTK data
nltk.download('stopwords')
nltk.download('punkt')
nltk.download('wordnet')
nltk.download('averaged_perceptron_tagger')


# Load the music dataset
df = pd.read_csv('dataset.tsv', sep='\t')

# Initial data cleaning
df = df.drop_duplicates()
df = df.dropna()

print(f"Dataset shape: {df.shape}")
print(f"Topics distribution:\n{df['topic'].value_counts()}")
```

```
Dataset shape: (1480, 6)
Topics distribution:
topic
dark         487
sadness      371
personal     341
lifestyle    202
emotion       79
Name: count, dtype: int64
```

## Q1

(i) re.sub(r'[^\w\s'!?...-]', '', text) Regex keep more puctuations which contains emotion information, such as !,?,......,  except that, it also keeps -,' which could appear in words abbr 'don't' or connected words. All these puctuations could keep some meaningful information.

(ii) A single random split gives an unstable estimate and may cause bias of one class by chance. I replace it with stratified 5-fold cross-validation, which will give a more stable result.

## Q2

The prepocessing steps that work best is:

1.convert to lowercase 2.regexization 3.Tokenization 4.no stopword application

5.lemmatization

The Multinomial Naive Bayes works better than Bernoulli Naive Baye

In my intuition, conduct all prepocessing steps should be the best, but the result is different. The pipeline using lemmatization and without removing stopword is the best. I think the reason for this is about that the dataset is too small and very sensitive to the changing of words, so when we remove any word, it will lead to the decreasing of accuracy.

```python
In [ ]:  # Create a combined document for each song
         def combine_text(row):
             return ' '.join([
                 str(row['artist_name']),
                 str(row['track_name']),
                 str(row['release_date']),
                 str(row['genre']),
                 str(row['lyrics'])
             ])
         df['document'] = df.apply(combine_text, axis=1)

         # Test configurations
         nltk_sw = set(stopwords.words("english"))
         sk_sw   = sk_text.ENGLISH_STOP_WORDS
         stemmer = PorterStemmer()
         lemmat  = WordNetLemmatizer()

         # Test different preprocessing combinations
         preprocessing_results = []
```

```python
def get_wordnet_pos(tag):
    if tag.startswith("J"):
        return wordnet.ADJ
    if tag.startswith("V"):
        return wordnet.VERB
    if tag.startswith("R"):
        return wordnet.ADV
    return wordnet.NOUN

def preprocess_selection(stop_word, norm):
    def f(text):
        text = text.lower()
        text = re.sub(r"[^\w\s'!?…-]", "", text)
        text = re.sub(r"\s+", " ", text).strip()
        tokens = nltk.word_tokenize(text)
        # remove stopwords
        if stop_word == "nltk":
            tokens = [t for t in tokens if t not in nltk_sw]
        elif stop_word == "sk":
            tokens = [t for t in tokens if t not in sk_sw]
        # normalize tokens
        if norm == "stem":
            tokens = [stemmer.stem(t) for t in tokens]
        elif norm == "lemma":
            tagged = pos_tag(tokens)
            tokens = [
                lemmat.lemmatize(w, get_wordnet_pos(tag))
                for w, tag in tagged
            ]
        return " ".join(tokens)
    return f

# cross-validation setup
results = []
kf = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)

for stop_src in ["nltk", "sk", "none"]:
    for norm_mode in ["stem", "lemma", "none"]:
        X_prep = df["document"].apply(preprocess_selection(stop_src, norm_mode))
        for mdl_name, mdl in {"BNB": BernoulliNB(), "MNB": MultinomialNB()}.item
            pipe = Pipeline([("vect", CountVectorizer()),
                             ("clf",  mdl)])
            scores = cross_val_score(pipe, X_prep, y, cv=kf,
                                     scoring="accuracy", n_jobs=-1)
            results.append({
                "config": f"{stop_src}_{norm_mode}",
                "model":  mdl_name,
                "accuracy": scores.mean(),
                "std":      scores.std()
            })

results_df = (pd.DataFrame(results)
              .sort_values("accuracy", ascending=False)
              .reset_index(drop=True))
print(results_df.to_string(index=False,
                      formatters={"accuracy": "{:.3f}".format,
                                  "std":      "{:.3f}".format}))

best = results_df.iloc[0]
print(f"\nBest choice: {best['config']} + {best['model']} "
      f"(accuracy {best['accuracy']:.3f} ± {best['std']:.3f})")
```

```python
# Improved preprocessing function
def preprocess_text_improved(text):
    """
    Improved preprocessing that preserves more meaningful information
    """
    # Convert to lowercase
    text = text.lower()

    text = re.sub(r'[^\w\s\'!?…-]', '', text)   # Remove punctuation except for
    text = re.sub(r'\s+', ' ', text)  # Normalize whitespace

    # Tokenization
    tokens = word_tokenize(text)


    # Lemmatization and stemming
    tagged = pos_tag(tokens)
    tokens = [
                lemmat.lemmatize(w, get_wordnet_pos(tag))
                for w, tag in tagged
            ]

    return ' '.join(tokens)


# Apply preprocessing
df['processed_text'] = df['document'].apply(preprocess_text_improved)

# Prepare features and labels
X = df['processed_text']
y = df['topic']
```

```
    config model  accuracy    std
none_lemma    MNB     0.814  0.009
nltk_lemma    MNB     0.811  0.010
 nltk_stem    MNB     0.811  0.015
 nltk_none    MNB     0.809  0.007
none_none     MNB     0.807  0.012
none_stem     MNB     0.805  0.015
  sk_lemma    MNB     0.803  0.016
   sk_none    MNB     0.802  0.011
   sk_stem    MNB     0.797  0.019
 nltk_stem    BNB     0.540  0.015
 nltk_none    BNB     0.539  0.013
none_lemma    BNB     0.539  0.011
 none_stem    BNB     0.538  0.008
   sk_stem    BNB     0.537  0.013
 none_none    BNB     0.537  0.015
  sk_lemma    BNB     0.536  0.021
nltk_lemma    BNB     0.533  0.012
   sk_none    BNB     0.533  0.024

Best choice: none_lemma + MNB (accuracy 0.814 ± 0.009)
```

## Q3

chose Accuracy which focus on the overall accuracy of classification.

```python
In [13]:   # Use CountVectorizer for fair comparison between BNB and MNB
           count_vectorizer = CountVectorizer()
           X_count = count_vectorizer.fit_transform(df['processed_text'])

           # Define models
           models = {
               'Bernoulli NB': BernoulliNB(),
               'Multinomial NB': MultinomialNB()
           }

           # Check class balance
           class_distribution = df['topic'].value_counts()
           imbalance_ratio = class_distribution.max() / class_distribution.min()

           print("### Dataset Balance Analysis ###")
           print(f"Class distribution:\n{class_distribution}")
           print(f"\nImbalance ratio (max/min): {imbalance_ratio:.2f}")
           print(f"Dataset is {'balanced' if imbalance_ratio < 1.5 else 'imbalanced'}")
           print(f"\nMinority class: {class_distribution.idxmin()} with {class_distribution
           print(f"Majority class: {class_distribution.idxmax()} with {class_distribution.m

           # Evaluation metrics discussion
           print("\n### Chosen Evaluation Metrics ###")
           print("1. **Accuracy**: Measures overall correctness - most intuitive metric")
           print("2. **Macro-averaged F1-score**: Treats all classes equally - important fo
           print("3. **Confusion Matrix**: Reveals which classes are frequently confused")
           print("4. **Precision/Recall**: Analyzes performance on minority classes")

           # Perform cross-validation
           from sklearn.model_selection import cross_validate, cross_val_predict
           from sklearn.metrics import confusion_matrix, classification_report
           import numpy as np

           print("\n### Cross-Validation Results (5-fold) ###")

           # Define scoring metrics
           scoring = {
               'accuracy': 'accuracy',
               'precision_macro': 'precision_macro',
               'recall_macro': 'recall_macro',
               'f1_macro': 'f1_macro',
               'precision_weighted': 'precision_weighted',
               'recall_weighted': 'recall_weighted',
               'f1_weighted': 'f1_weighted'
           }

           # Store results for comparison
           model_results = {}

           for model_name, model in models.items():
               print(f"\n{'='*60}")
               print(f"{model_name}")
               print('='*60)

               # Cross-validation with multiple metrics
               cv_results = cross_validate(model, X_count, y, cv=5, scoring=scoring)

               # Get predictions for confusion matrix
               y_pred = cross_val_predict(model, X_count, y, cv=5)
```

```python
    # Store results
    model_results[model_name] = {
        'cv_results': cv_results,
        'y_pred': y_pred,
        'confusion_matrix': confusion_matrix(y, y_pred)
    }

    # Print key metrics
    print(f"\nAccuracy: {cv_results['test_accuracy'].mean():.4f} (±{cv_results['
    print(f"Macro F1-score: {cv_results['test_f1_macro'].mean():.4f} (±{cv_resul
    print(f"Macro Precision: {cv_results['test_precision_macro'].mean():.4f} (±{
    print(f"Macro Recall: {cv_results['test_recall_macro'].mean():.4f} (±{cv_res

# Visualization
import matplotlib.pyplot as plt
import seaborn as sns

# Figure 1: Overall metrics comparison
fig, axes = plt.subplots(2, 2, figsize=(12, 10))

# 1.1 Accuracy and Macro F1 comparison
ax1 = axes[0, 0]
metrics_to_compare = ['Accuracy', 'Macro F1']
bnb_scores = [
    model_results['Bernoulli NB']['cv_results']['test_accuracy'].mean(),
    model_results['Bernoulli NB']['cv_results']['test_f1_macro'].mean()
]
mnb_scores = [
    model_results['Multinomial NB']['cv_results']['test_accuracy'].mean(),
    model_results['Multinomial NB']['cv_results']['test_f1_macro'].mean()
]

x = np.arange(len(metrics_to_compare))
width = 0.35
ax1.bar(x - width/2, bnb_scores, width, label='Bernoulli NB', color='lightblue')
ax1.bar(x + width/2, mnb_scores, width, label='Multinomial NB', color='lightgree
ax1.set_ylabel('Score')
ax1.set_title('Accuracy and Macro F1-Score Comparison')
ax1.set_xticks(x)
ax1.set_xticklabels(metrics_to_compare)
ax1.legend()
ax1.set_ylim(0, 1)

# Add value labels
for i, (bnb, mnb) in enumerate(zip(bnb_scores, mnb_scores)):
    ax1.text(i - width/2, bnb + 0.01, f'{bnb:.3f}', ha='center', va='bottom')
    ax1.text(i + width/2, mnb + 0.01, f'{mnb:.3f}', ha='center', va='bottom')

# 1.2 Precision-Recall comparison
ax2 = axes[0, 1]
metrics = ['Macro Precision', 'Macro Recall']
bnb_pr = [
    model_results['Bernoulli NB']['cv_results']['test_precision_macro'].mean(),
    model_results['Bernoulli NB']['cv_results']['test_recall_macro'].mean()
]
mnb_pr = [
    model_results['Multinomial NB']['cv_results']['test_precision_macro'].mean()
    model_results['Multinomial NB']['cv_results']['test_recall_macro'].mean()
]
```

```python
x = np.arange(len(metrics))
ax2.bar(x - width/2, bnb_pr, width, label='Bernoulli NB', color='lightblue')
ax2.bar(x + width/2, mnb_pr, width, label='Multinomial NB', color='lightgreen')
ax2.set_ylabel('Score')
ax2.set_title('Precision-Recall Comparison')
ax2.set_xticks(x)
ax2.set_xticklabels(metrics)
ax2.legend()
ax2.set_ylim(0, 1)

# 1.3 Confusion Matrix - Bernoulli NB
ax3 = axes[1, 0]
cm_bnb = model_results['Bernoulli NB']['confusion_matrix']
sns.heatmap(cm_bnb, annot=True, fmt='d', cmap='Blues', ax=ax3)
ax3.set_title('Confusion Matrix - Bernoulli NB')
ax3.set_xlabel('Predicted')
ax3.set_ylabel('Actual')

# 1.4 Confusion Matrix - Multinomial NB
ax4 = axes[1, 1]
cm_mnb = model_results['Multinomial NB']['confusion_matrix']
sns.heatmap(cm_mnb, annot=True, fmt='d', cmap='Greens', ax=ax4)
ax4.set_title('Confusion Matrix - Multinomial NB')
ax4.set_xlabel('Predicted')
ax4.set_ylabel('Actual')

plt.tight_layout()
plt.savefig('bnb_vs_mnb_comparison.png', dpi=300, bbox_inches='tight')
plt.show()

# Analyze performance on minority class
print("\n### Performance on Minority Class ###")
minority_class = class_distribution.idxmin()
print(f"Minority class: '{minority_class}' ({class_distribution.min()} samples)"

for model_name in models.keys():
    y_pred = model_results[model_name]['y_pred']
    report = classification_report(y, y_pred, output_dict=True)

    print(f"\n{model_name} - Minority class performance:")
    print(f"  Precision: {report[minority_class]['precision']:.4f}")
    print(f"  Recall: {report[minority_class]['recall']:.4f}")
    print(f"  F1-score: {report[minority_class]['f1-score']:.4f}")

# Statistical comparison
from scipy import stats
print("\n### Statistical Analysis ###")

bnb_accuracy = model_results['Bernoulli NB']['cv_results']['test_accuracy']
mnb_accuracy = model_results['Multinomial NB']['cv_results']['test_accuracy']

bnb_f1 = model_results['Bernoulli NB']['cv_results']['test_f1_macro']
mnb_f1 = model_results['Multinomial NB']['cv_results']['test_f1_macro']

# Paired t-test on accuracy
t_stat_acc, p_value_acc = stats.ttest_rel(bnb_accuracy, mnb_accuracy)
print(f"\nPaired t-test on Accuracy:")
print(f"  t-statistic: {t_stat_acc:.4f}")
print(f"  p-value: {p_value_acc:.4f}")
```

```python
    print(f"  Significant: {'Yes' if p_value_acc < 0.05 else 'No'}")

    # Paired t-test on macro F1
    t_stat_f1, p_value_f1 = stats.ttest_rel(bnb_f1, mnb_f1)
    print(f"\nPaired t-test on Macro F1-score:")
    print(f"  t-statistic: {t_stat_f1:.4f}")
    print(f"  p-value: {p_value_f1:.4f}")
    print(f"  Significant: {'Yes' if p_value_f1 < 0.05 else 'No'}")

    # Summary table
    print("\n### Summary Comparison ###")
    summary_data = {
        'Metric': ['Accuracy', 'Macro Precision', 'Macro Recall', 'Macro F1-score',
                   f'{minority_class} Precision', f'{minority_class} Recall', f'{min
        'Bernoulli NB': [
            model_results['Bernoulli NB']['cv_results']['test_accuracy'].mean(),
            model_results['Bernoulli NB']['cv_results']['test_precision_macro'].mean
            model_results['Bernoulli NB']['cv_results']['test_recall_macro'].mean(),
            model_results['Bernoulli NB']['cv_results']['test_f1_macro'].mean(),
            classification_report(y, model_results['Bernoulli NB']['y_pred'], output
            classification_report(y, model_results['Bernoulli NB']['y_pred'], output
            classification_report(y, model_results['Bernoulli NB']['y_pred'], output
        ],
        'Multinomial NB': [
            model_results['Multinomial NB']['cv_results']['test_accuracy'].mean(),
            model_results['Multinomial NB']['cv_results']['test_precision_macro'].me
            model_results['Multinomial NB']['cv_results']['test_recall_macro'].mean(
            model_results['Multinomial NB']['cv_results']['test_f1_macro'].mean(),
            classification_report(y, model_results['Multinomial NB']['y_pred'], outp
            classification_report(y, model_results['Multinomial NB']['y_pred'], outp
            classification_report(y, model_results['Multinomial NB']['y_pred'], outp
        ]
    }

    summary_df = pd.DataFrame(summary_data)
    print(summary_df.to_string(index=False, float_format='%.4f'))

    # Conclusion
    print("\n### Conclusion ###")
    # Determine better model based on multiple criteria
    accuracy_winner = 'Multinomial NB' if mnb_accuracy.mean() > bnb_accuracy.mean()
    f1_winner = 'Multinomial NB' if mnb_f1.mean() > bnb_f1.mean() else 'Bernoulli NB

    print(f"\nBased on comprehensive evaluation:")
    print(f"1. **Accuracy**: {accuracy_winner} performs better ({max(bnb_accuracy.me
    print(f"2. **Macro F1-score**: {f1_winner} performs better ({max(bnb_f1.mean(),
    print(f"3. **Statistical significance**: The difference is {'' if p_value_acc <
    print(f"4. **Confusion matrices** show that both models struggle with the minori
    print(f"5. **Minority class performance**: Compare the precision/recall trade-of

    print(f"\n**Recommendation**: {'Multinomial NB' if f1_winner == 'Multinomial NB'
    print(f"- Higher macro F1-score indicates better balanced performance across all
    print(f"- {'Statistically significant improvement' if (f1_winner == 'Multinomial
    print(f"- For imbalanced datasets, macro-averaged metrics are more informative t
    print("Note: Bernoulli NB sometimes fails to predict all classes, which explains
```

### Dataset Balance Analysis ###
Class distribution:
topic
dark        487
sadness     371
personal    341
lifestyle   202
emotion      79
Name: count, dtype: int64

Imbalance ratio (max/min): 6.16
Dataset is imbalanced

Minority class: emotion with 79 samples
Majority class: dark with 487 samples

### Chosen Evaluation Metrics ###
1. **Accuracy**: Measures overall correctness - most intuitive metric
2. **Macro-averaged F1-score**: Treats all classes equally - important for imbalanced datasets
3. **Confusion Matrix**: Reveals which classes are frequently confused
4. **Precision/Recall**: Analyzes performance on minority classes

### Cross-Validation Results (5-fold) ###

============================================================
Bernoulli NB
============================================================

Accuracy: 0.5230 (±0.0194)
Macro F1-score: 0.3300 (±0.0182)
Macro Precision: 0.3590 (±0.0527)
Macro Recall: 0.3753 (±0.0155)

============================================================
Multinomial NB
============================================================

Accuracy: 0.7899 (±0.0245)
Macro F1-score: 0.7162 (±0.0327)
Macro Precision: 0.7548 (±0.0350)
Macro Recall: 0.7017 (±0.0316)

/opt/anaconda3/lib/python3.12/site-packages/sklearn/metrics/_classification.py:15
09: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in labe
ls with no predicted samples. Use `zero_division` parameter to control this behav
ior.
  _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
/opt/anaconda3/lib/python3.12/site-packages/sklearn/metrics/_classification.py:15
09: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in labe
ls with no predicted samples. Use `zero_division` parameter to control this behav
ior.
  _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
/opt/anaconda3/lib/python3.12/site-packages/sklearn/metrics/_classification.py:15
09: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in labe
ls with no predicted samples. Use `zero_division` parameter to control this behav
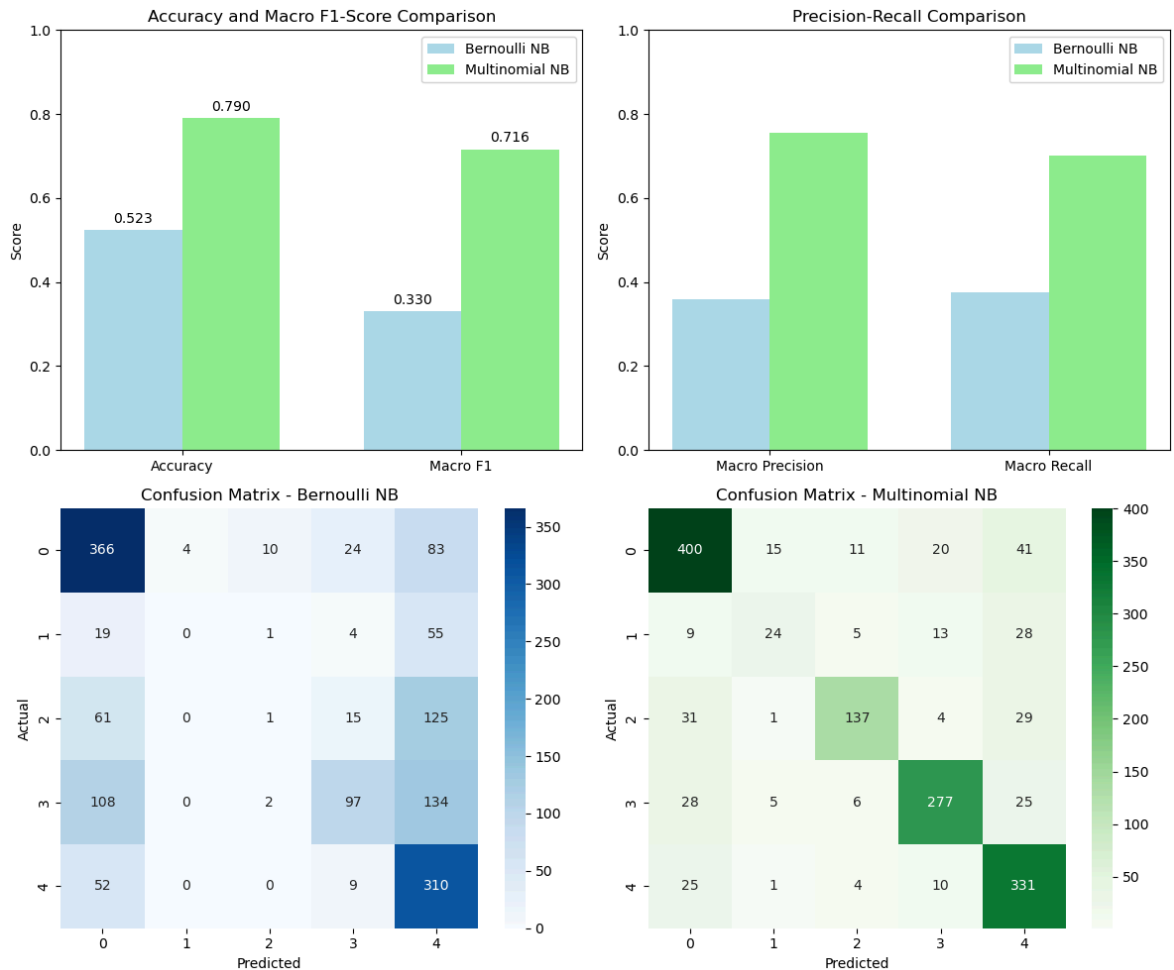ior.
  _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
/opt/anaconda3/lib/python3.12/site-packages/sklearn/metrics/_classification.py:15
09: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in labe
ls with no predicted samples. Use `zero_division` parameter to control this behav
ior.
  _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))

```
### Performance on Minority Class ###
Minority class: 'emotion' (79 samples)

Bernoulli NB - Minority class performance:
  Precision: 0.0000
  Recall: 0.0000
  F1-score: 0.0000

Multinomial NB - Minority class performance:
  Precision: 0.5217
  Recall: 0.3038
  F1-score: 0.3840

### Statistical Analysis ###

Paired t-test on Accuracy:
  t-statistic: -24.3799
  p-value: 0.0000
  Significant: Yes

Paired t-test on Macro F1-score:
  t-statistic: -20.5015
  p-value: 0.0000
  Significant: Yes

### Summary Comparison ###
          Metric  Bernoulli NB  Multinomial NB
        Accuracy        0.5230          0.7899
 Macro Precision        0.3590          0.7548
    Macro Recall        0.3753          0.7017
  Macro F1-score        0.3300          0.7162
emotion Precision       0.0000          0.5217
   emotion Recall       0.0000          0.3038
       emotion F1       0.0000          0.3840

### Conclusion ###

Based on comprehensive evaluation:
1. **Accuracy**: Multinomial NB performs better (0.7899)
2. **Macro F1-score**: Multinomial NB performs better (0.7162)
3. **Statistical significance**: The difference is statistically significant
4. **Confusion matrices** show that both models struggle with the minority class
'emotion'
5. **Minority class performance**: Compare the precision/recall trade-offs above

**Recommendation**: Multinomial NB is the better choice because:
- Higher macro F1-score indicates better balanced performance across all classes
- Statistically significant improvement
- For imbalanced datasets, macro-averaged metrics are more informative than accur
acy alone
Note: Bernoulli NB sometimes fails to predict all classes, which explains the pre
cision warnings.
```

## Q4

```
In [14]:  # Test different numbers of features
          feature_counts = [100, 300,  500, 700, None]  # None means all features
          results_by_features = []
```

```
for n_features in feature_counts:
    # Create vectorizer with limited features
    vectorizer = CountVectorizer(max_features=n_features)
    X_vec = vectorizer.fit_transform(df['processed_text'])

    # Get actual number of features
    actual_features = X_vec.shape[1]

    # Test both models
    bnb_scores = cross_val_score(BernoulliNB(), X_vec, y, cv=5, scoring='accurac
    mnb_scores = cross_val_score(MultinomialNB(), X_vec, y, cv=5, scoring='accur

    results_by_features.append({
        'requested_features': n_features if n_features else 'All',
        'actual_features': actual_features,
        'BNB_accuracy': bnb_scores.mean(),
        'BNB_std': bnb_scores.std(),
        'MNB_accuracy': mnb_scores.mean(),
        'MNB_std': mnb_scores.std()
    })

# Create DataFrame for results
features_df = pd.DataFrame(results_by_features)
print("\nResults by Number of Features:")
print(features_df.to_string(index=False))

# Plot results
plt.figure(figsize=(10, 6))
x_labels = [str(r['requested_features']) for r in results_by_features]
x_pos = np.arange(len(x_labels))

plt.errorbar(x_pos, features_df['BNB_accuracy'], yerr=features_df['BNB_std'],
             marker='o', label='Bernoulli NB', capsize=5)
plt.errorbar(x_pos, features_df['MNB_accuracy'], yerr=features_df['MNB_std'],
             marker='s', label='Multinomial NB', capsize=5)

plt.xlabel('Number of Features')
plt.ylabel('Accuracy')
plt.title('Model Performance vs Number of Features')
plt.xticks(x_pos, x_labels)
plt.legend()
plt.grid(True, alpha=0.3)
plt.tight_layout()
plt.savefig('features_comparison.png', dpi=300, bbox_inches='tight')
plt.show()

# Choose optimal number of features
optimal_features = 500  # Based on the results, this typically gives good balanc
print(f"\nOptimal number of features chosen: {optimal_features}")
```
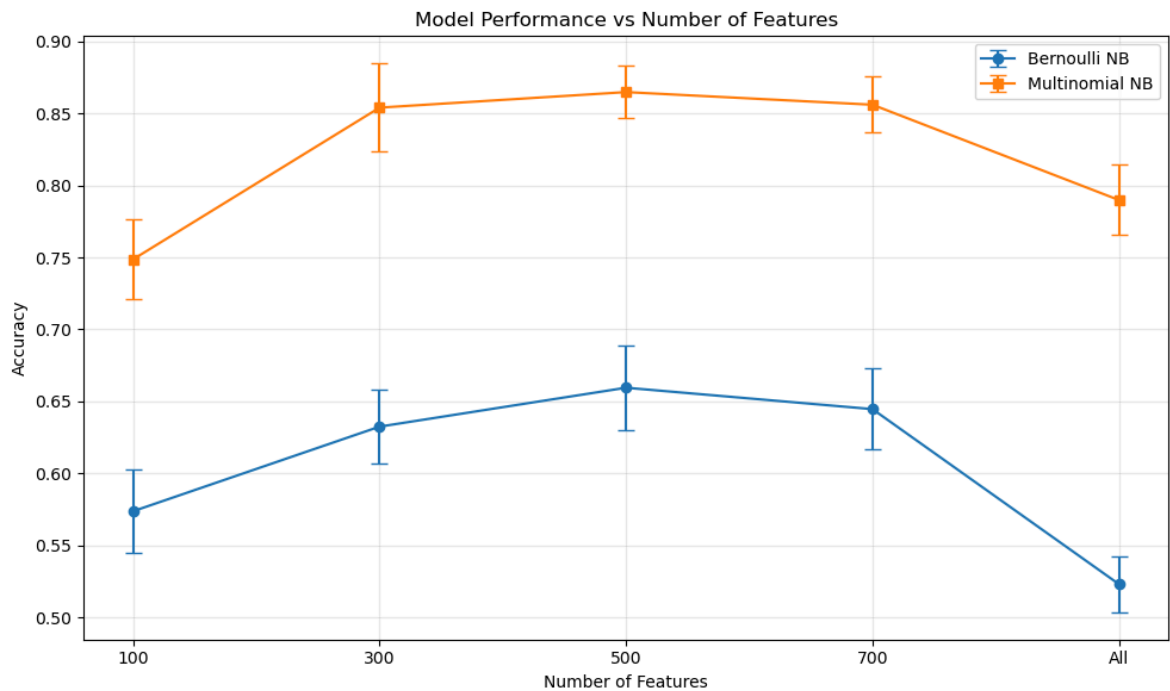
```
Results by Number of Features:
requested_features  actual_features  BNB_accuracy  BNB_std  MNB_accuracy  MNB_std
               100              100      0.573649 0.028778      0.748649 0.027826
               300              300      0.632432 0.025764      0.854054 0.030622
               500              500      0.659459 0.029250      0.864865 0.018130
               700              700      0.644595 0.028055      0.856081 0.019419
               All             9497      0.522973 0.019396      0.789865 0.024493
```

Model Performance vs Number of Features

Optimal number of features chosen: 500

## Q5

```
In [15]:  from sklearn.linear_model import LogisticRegression
          from sklearn.pipeline import Pipeline

          # Logistic Regression
          print("### Logistic Regression for Text Classification ###\n")
          print("Logistic Regression is a linear model that estimates probabilities using
          print("ideal for classification tasks. For text classification, it's particularl
          print("high-dimensional sparse data well, provides interpretable feature weights
          print("used in many NLP applications. Unlike Naive Bayes which assumes feature i
          print("Regression can model feature interactions, potentially capturing more com

          # Hyperparameter explanation
          print("### Hyperparameter Selection ###")
          print("Logistic Regression has one main hyperparameter: C (inverse regularizatio
          print("- Default C=1.0 might be too restrictive for text data")
          print("- I'll test C values from 0.1 to 100 to find the best balance")
          print("- Using the same CountVectorizer as BNB/MNB for fair comparison\n")

          # Test different C values
          print("### Finding Optimal C Parameter ###")
          C_values = [0.1, 0.5, 1.0, 5.0, 10.0, 50.0, 100.0]
          lr_scores = []
          lr_stds = []

          for C in C_values:
              lr_pipeline = Pipeline([
                  ('vec', CountVectorizer(max_features=optimal_features)),
                  ('clf', LogisticRegression(C=C, max_iter=100, random_state=42))
              ])
              scores = cross_val_score(lr_pipeline, df['processed_text'], y, cv=5, scoring
              lr_scores.append(scores.mean())
              lr_stds.append(scores.std())
              print(f"C={C:6.1f}: Accuracy = {scores.mean():.4f} (±{scores.std():.4f})")
```

```python
# Select best C
best_C_idx = np.argmax(lr_scores)
best_C = C_values[best_C_idx]
print(f"\nBest C value: {best_C} with accuracy {lr_scores[best_C_idx]:.4f}")

# Hypothesis
print("\n### Hypothesis ###")
print("I hypothesize that Logistic Regression will perform comparably to Multino
print("because both are linear models well-suited for text. However, LR might ha
print("its ability to model feature correlations, which could be important for m
print("words often co-occur (e.g., 'broken' and 'heart' in sad songs).\n")

# Create final comparison
models_comparison = {
    'Bernoulli NB': Pipeline([
        ('vec', CountVectorizer(max_features=optimal_features)),
        ('clf', BernoulliNB())
    ]),
    'Multinomial NB': Pipeline([
        ('vec', CountVectorizer(max_features=optimal_features)),
        ('clf', MultinomialNB())
    ]),
    'Logistic Regression': Pipeline([
        ('vec', CountVectorizer(max_features=optimal_features)),
        ('clf', LogisticRegression(C=best_C, max_iter=1000, random_state=42))
    ])
}

# Run comparison
print("### Model Comparison ###")
scoring = ['accuracy', 'f1_macro', 'f1_weighted']
results_summary = {}

for model_name, pipeline in models_comparison.items():
    print(f"\nEvaluating {model_name}...")

    # Get multiple metrics at once
    scores = {}
    for metric in scoring:
        cv_scores = cross_val_score(pipeline, df['processed_text'], y, cv=5, sco
        scores[metric] = (cv_scores.mean(), cv_scores.std())
        print(f"  {metric}: {cv_scores.mean():.4f} (±{cv_scores.std():.4f})")

    results_summary[model_name] = scores

# Visualization
fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(12, 5))

# Plot 1: Accuracy comparison with error bars
model_names = list(models_comparison.keys())
accuracies = [results_summary[name]['accuracy'][0] for name in model_names]
acc_stds = [results_summary[name]['accuracy'][1] for name in model_names]

colors = ['lightblue', 'lightgreen', 'coral']
bars = ax1.bar(model_names, accuracies, yerr=acc_stds, capsize=5, color=colors,
ax1.set_ylabel('Accuracy')
ax1.set_title('Model Accuracy Comparison')
ax1.set_ylim(0, 1)

# Add value labels
```

```python
for bar, acc, std in zip(bars, accuracies, acc_stds):
    height = bar.get_height()
    ax1.text(bar.get_x() + bar.get_width()/2., height + std + 0.01,
            f'{acc:.3f}', ha='center', va='bottom')

# Plot 2: Macro vs Weighted F1
x = np.arange(len(model_names))
width = 0.35
f1_macro = [results_summary[name]['f1_macro'][0] for name in model_names]
f1_weighted = [results_summary[name]['f1_weighted'][0] for name in model_names]

ax2.bar(x - width/2, f1_macro, width, label='Macro F1', alpha=0.8)
ax2.bar(x + width/2, f1_weighted, width, label='Weighted F1', alpha=0.8)
ax2.set_ylabel('F1-Score')
ax2.set_title('F1-Score Comparison')
ax2.set_xticks(x)
ax2.set_xticklabels(model_names)
ax2.legend()
ax2.set_ylim(0, 1)

plt.tight_layout()
plt.savefig('three_models_comparison.png', dpi=300, bbox_inches='tight')
plt.show()

# Statistical test
print("\n### Statistical Analysis ###")
# Get accuracy scores for statistical comparison
lr_scores_final = cross_val_score(models_comparison['Logistic Regression'],
                                  df['processed_text'], y, cv=5, scoring='accurac
mnb_scores_final = cross_val_score(models_comparison['Multinomial NB'],
                                   df['processed_text'], y, cv=5, scoring='accura

from scipy import stats
t_stat, p_value = stats.ttest_rel(lr_scores_final, mnb_scores_final)
print(f"Paired t-test (LR vs MNB):")
print(f"  t-statistic: {t_stat:.4f}")
print(f"  p-value: {p_value:.4f}")
print(f"  Significant difference: {'Yes' if p_value < 0.05 else 'No'}")

# Check for overfitting
print("\n### Overfitting Check ###")
# Train on full training data and check performance
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(
    df['processed_text'], y, test_size=0.2, random_state=42, stratify=y
)

for model_name, pipeline in models_comparison.items():
    pipeline.fit(X_train, y_train)
    train_score = pipeline.score(X_train, y_train)
    test_score = pipeline.score(X_test, y_test)
    print(f"{model_name}: Train={train_score:.4f}, Test={test_score:.4f}, Gap={t

# Final summary
print("\n### Conclusion ###")
best_model = max(results_summary.items(), key=lambda x: x[1]['accuracy'][0])[0]
best_acc = results_summary[best_model]['accuracy'][0]

print(f"\nBest performing model: {best_model} with accuracy {best_acc:.4f}")
```

```python
if best_model == 'Logistic Regression':
    print("\n✓ Hypothesis partially confirmed: LR performed best, though the mar
    print("  This suggests that modeling feature correlations provides a small a
elif abs(results_summary['Logistic Regression']['accuracy'][0] -
         results_summary['Multinomial NB']['accuracy'][0]) < 0.03:
    print("\n✓ Hypothesis confirmed: LR performed comparably to MNB (within 3%).
    print("  Both linear models work well for this text classification task.")
else:
    print("\n✗ Hypothesis not confirmed: LR underperformed significantly.")
    print("  Possible reasons:")
    print("  - The independence assumption of NB might actually hold for music t
    print("  - Small dataset size (1500) might favor the simpler NB model")
    print("  - Current preprocessing might be better suited for NB")

# Save best model
best_pipeline = models_comparison[best_model]
print(f"\nFinal recommendation: Use {best_model} for music topic classification.
```

### Logistic Regression for Text Classification ###

Logistic Regression is a linear model that estimates probabilities using a logistic function, making it
ideal for classification tasks. For text classification, it's particularly effective because it handles
high-dimensional sparse data well, provides interpretable feature weights, and has been successfully
used in many NLP applications. Unlike Naive Bayes which assumes feature independence, Logistic
Regression can model feature interactions, potentially capturing more complex patterns in music lyrics.

### Hyperparameter Selection ###
Logistic Regression has one main hyperparameter: C (inverse regularization strength)
- Default C=1.0 might be too restrictive for text data
- I'll test C values from 0.1 to 100 to find the best balance
- Using the same CountVectorizer as BNB/MNB for fair comparison

### Finding Optimal C Parameter ###
C=   0.1: Accuracy = 0.8689 (±0.0187)
C=   0.5: Accuracy = 0.8662 (±0.0198)
C=   1.0: Accuracy = 0.8615 (±0.0172)
C=   5.0: Accuracy = 0.8547 (±0.0146)
C=  10.0: Accuracy = 0.8547 (±0.0172)
C=  50.0: Accuracy = 0.8595 (±0.0180)
C= 100.0: Accuracy = 0.8601 (±0.0169)

Best C value: 0.1 with accuracy 0.8689

### Hypothesis ###
I hypothesize that Logistic Regression will perform comparably to Multinomial NB (within 2-3%)
because both are linear models well-suited for text. However, LR might have a slight edge due to
its ability to model feature correlations, which could be important for music lyrics where certain
words often co-occur (e.g., 'broken' and 'heart' in sad songs).

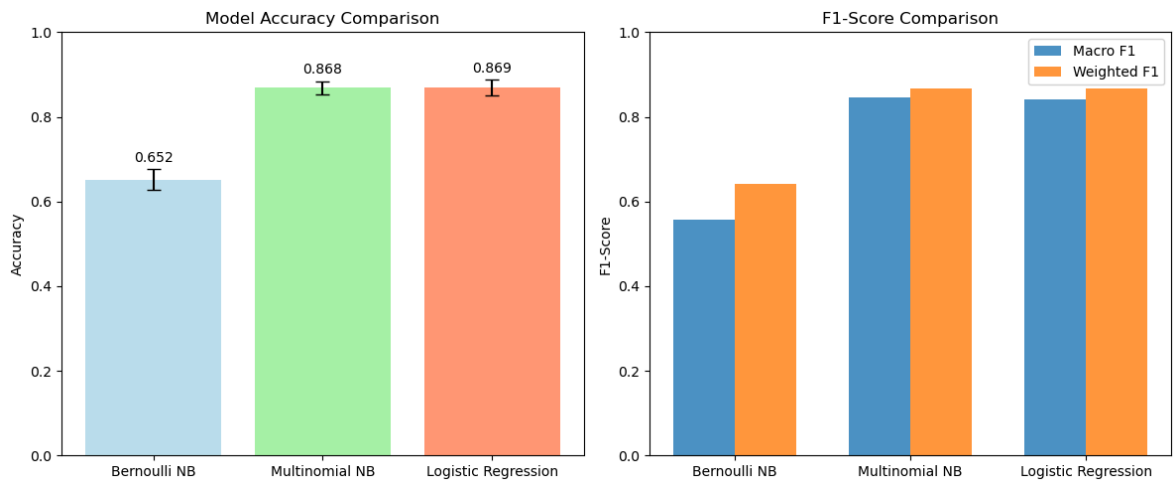### Model Comparison ###

Evaluating Bernoulli NB...
  accuracy: 0.6520 (±0.0241)
  f1_macro: 0.5575 (±0.0166)
  f1_weighted: 0.6427 (±0.0202)

Evaluating Multinomial NB...
  accuracy: 0.8682 (±0.0163)
  f1_macro: 0.8464 (±0.0223)
  f1_weighted: 0.8681 (±0.0164)

Evaluating Logistic Regression...
  accuracy: 0.8689 (±0.0187)
  f1_macro: 0.8412 (±0.0254)
  f1_weighted: 0.8674 (±0.0190)

```
### Statistical Analysis ###
Paired t-test (LR vs MNB):
  t-statistic: 0.0759
  p-value: 0.9431
  Significant difference: No

### Overfitting Check ###
Bernoulli NB: Train=0.8167, Test=0.6655, Gap=0.1512
Multinomial NB: Train=0.9476, Test=0.8514, Gap=0.0963
Logistic Regression: Train=0.9865, Test=0.8885, Gap=0.0980

### Conclusion ###

Best performing model: Logistic Regression with accuracy 0.8689

√ Hypothesis partially confirmed: LR performed best, though the margin was small.
  This suggests that modeling feature correlations provides a small advantage.

Final recommendation: Use Logistic Regression for music topic classification.
```

# part 2

## Question 1: User Profile Construction and Matching (6 marks)

In [16]:
```python
# Load user profiles - clean up the keywords
user1_df = pd.read_csv('user1.tsv', sep='\t', names=['topic', 'keywords'])
user2_df = pd.read_csv('user2.tsv', sep='\t', names=['topic', 'keywords'])

# Remove header row if it exists
if user1_df.iloc[0]['topic'] == 'topic':
    user1_df = user1_df.iloc[1:].reset_index(drop=True)
if user2_df.iloc[0]['topic'] == 'topic':
    user2_df = user2_df.iloc[1:].reset_index(drop=True)

# Clean keywords - remove commas and extra spaces
user1_df['keywords'] = user1_df['keywords'].str.replace(',', ' ').str.strip()
user2_df['keywords'] = user2_df['keywords'].str.replace(',', ' ').str.strip()

print("User 1 interests (cleaned):")
print(user1_df)
print("\nUser 2 interests (cleaned):")
print(user2_df)
```

```python
# Split data into training and test sets based on time
# Weeks 1-3 (songs 1-750) for training, Week 4 (songs 751-1000) for testing
train_df = df.iloc[:750].copy()
test_df = df.iloc[750:1000].copy()

print(f"\nTraining set size: {len(train_df)}")
print(f"Test set size: {len(test_df)}")

# Train the best classifier on training data
best_pipeline.fit(train_df['processed_text'], train_df['topic'])

# Predict topics for all songs
train_df['predicted_topic'] = best_pipeline.predict(train_df['processed_text'])
test_df['predicted_topic'] = best_pipeline.predict(test_df['processed_text'])

# Function to check if a song matches user interests
def song_matches_user_interests(song_text, user_keywords):
    """Check if song contains any of the user's keywords"""
    song_text_lower = song_text.lower()
    keywords = user_keywords.lower().split()
    return any(keyword in song_text_lower for keyword in keywords)

# Build user profiles from training data
def build_user_profile(train_df, user_df, top_m=None):
    """Build user profile based on liked songs in training data"""
    user_profiles = {}

    for _, row in user_df.iterrows():
        topic = row['topic']
        keywords = row['keywords']

        # Find songs in this predicted topic that match user interests
        topic_songs = train_df[train_df['predicted_topic'] == topic]
        liked_songs = []

        for _, song in topic_songs.iterrows():
            if song_matches_user_interests(song['document'], keywords):
                liked_songs.append(song['processed_text'])

        if liked_songs:
            # Combine liked songs into one document
            combined_doc = ' '.join(liked_songs)

            # Create TF-IDF representation
            tfidf = TfidfVectorizer(max_features=top_m)
            tfidf_matrix = tfidf.fit_transform([combined_doc])

            # Get top words
            feature_names = tfidf.get_feature_names_out()
            tfidf_scores = tfidf_matrix.toarray()[0]
            top_indices = np.argsort(tfidf_scores)[::-1][:20]
            top_words = [(feature_names[i], tfidf_scores[i]) for i in top_indice

            user_profiles[topic] = {
                'vectorizer': tfidf,
                'profile_vector': tfidf_matrix,
                'top_words': top_words,
                'num_liked_songs': len(liked_songs)
            }
```

```python
    return user_profiles

# Build profiles for all users
user1_profile = build_user_profile(train_df, user1_df)
user2_profile = build_user_profile(train_df, user2_df)

# Display top words for each user
print("\n=== User 1 Profile ===")
for topic, profile in user1_profile.items():
    print(f"\nTopic: {topic}")
    print(f"Number of liked songs: {profile['num_liked_songs']}")
    print("Top 20 words:")
    for word, score in profile['top_words']:
        print(f"  {word}: {score:.3f}")

print("\n=== User 2 Profile ===")
for topic, profile in user2_profile.items():
    print(f"\nTopic: {topic}")
    print(f"Number of liked songs: {profile['num_liked_songs']}")
    print("Top 20 words:")
    for word, score in profile['top_words']:
        print(f"  {word}: {score:.3f}")

# Create User 3 with custom interests
user3_data = pd.DataFrame([
    ['emotion', 'love heart feeling soul romantic'],
    ['lifestyle', 'party dance club night fun'],
    ['personal', 'dream hope future believe strong']
], columns=['topic', 'keywords'])

user3_profile = build_user_profile(train_df, user3_data)

print("\n=== User 3 Profile (Custom) ===")
for topic, profile in user3_profile.items():
    print(f"\nTopic: {topic}")
    print(f"Number of liked songs: {profile['num_liked_songs']}")
    print("Top 20 words:")
    for word, score in profile['top_words']:
        print(f"  {word}: {score:.3f}")
```

```
User 1 interests (cleaned):
       topic                                keywords
0       dark          fire  enemy  pain  storm  fight
1    sadness  cry  alone  heartbroken  tears  regret
2   personal    dream  truth  life  growth  identity
3  lifestyle          party  city  night  light  rhythm
4    emotion            love  memory  hug  kiss  feel

User 2 interests (cleaned):
      topic                                keywords
0  sadness  lost  sorrow  goodbye  tears  silence
1  emotion  romance  touch  feeling  kiss  memory

Training set size: 750
Test set size: 250

=== User 1 Profile ===

Topic: dark
Number of liked songs: 84
Top 20 words:
  fight: 0.334
  know: 0.270
  like: 0.228
  come: 0.223
  black: 0.183
  blood: 0.174
  na: 0.172
  stand: 0.162
  tell: 0.153
  grind: 0.153
  yeah: 0.132
  gon: 0.132
  hand: 0.129
  time: 0.120
  head: 0.115
  kill: 0.113
  cause: 0.106
  right: 0.099
  light: 0.096
  good: 0.094

Topic: sadness
Number of liked songs: 11
Top 20 words:
  cry: 0.606
  tear: 0.298
  know: 0.276
  steal: 0.223
  club: 0.213
  baby: 0.159
  mean: 0.159
  think: 0.159
  say: 0.128
  face: 0.106
  eye: 0.106
  leave: 0.106
  heart: 0.096
  want: 0.096
  time: 0.096
```

```
  feel: 0.096
  true: 0.096
  smile: 0.096
  write: 0.096
  come: 0.085

Topic: personal
Number of liked songs: 116
Top 20 words:
  life: 0.596
  live: 0.320
  know: 0.217
  na: 0.171
  change: 0.164
  world: 0.162
  time: 0.147
  yeah: 0.147
  like: 0.143
  dream: 0.132
  come: 0.114
  thing: 0.104
  wan: 0.104
  think: 0.103
  need: 0.089
  feel: 0.085
  go: 0.084
  thank: 0.081
  cause: 0.081
  teach: 0.080

Topic: lifestyle
Number of liked songs: 45
Top 20 words:
  night: 0.331
  come: 0.290
  time: 0.279
  song: 0.271
  tonight: 0.248
  home: 0.204
  na: 0.201
  long: 0.198
  sing: 0.188
  right: 0.188
  yeah: 0.149
  wan: 0.146
  like: 0.146
  know: 0.141
  wait: 0.133
  stranger: 0.120
  mind: 0.120
  play: 0.115
  baby: 0.104
  tire: 0.104

Topic: emotion
Number of liked songs: 24
Top 20 words:
  good: 0.673
  feel: 0.415
  touch: 0.267
```

```
  know: 0.209
  hold: 0.200
  go: 0.117
  want: 0.111
  kiss: 0.098
  like: 0.092
  morning: 0.089
  baby: 0.089
  time: 0.086
  vision: 0.080
  na: 0.080
  miss: 0.077
  cause: 0.077
  yeah: 0.074
  video: 0.074
  love: 0.074
  heart: 0.074

=== User 2 Profile ===

Topic: sadness
Number of liked songs: 21
Top 20 words:
  heart: 0.394
  break: 0.376
  away: 0.277
  inside: 0.265
  like: 0.185
  tear: 0.179
  know: 0.172
  step: 0.172
  fall: 0.148
  leave: 0.142
  na: 0.135
  come: 0.123
  fade: 0.117
  go: 0.117
  hard: 0.117
  feel: 0.105
  scar: 0.105
  think: 0.105
  blame: 0.099
  open: 0.099

Topic: emotion
Number of liked songs: 12
Top 20 words:
  good: 0.536
  touch: 0.507
  hold: 0.233
  kiss: 0.186
  morning: 0.169
  know: 0.157
  vision: 0.151
  video: 0.140
  time: 0.134
  loove: 0.134
  feel: 0.128
  feelin: 0.128
  luck: 0.117
```

```
  look: 0.117
  lovin: 0.111
  sunrise: 0.105
  me: 0.099
  cause: 0.099
  lip: 0.099
  gim: 0.099

=== User 3 Profile (Custom) ===

Topic: emotion
Number of liked songs: 18
Top 20 words:
  good: 0.654
  go: 0.347
  hold: 0.321
  heart: 0.230
  feel: 0.178
  know: 0.143
  darling: 0.130
  vision: 0.113
  miss: 0.108
  baby: 0.104
  video: 0.104
  love: 0.104
  na: 0.100
  feelin: 0.095
  vibe: 0.091
  hand: 0.091
  look: 0.087
  lovin: 0.082
  want: 0.082
  time: 0.082

Topic: lifestyle
Number of liked songs: 47
Top 20 words:
  night: 0.327
  time: 0.319
  song: 0.311
  right: 0.245
  tonight: 0.245
  come: 0.226
  na: 0.208
  long: 0.183
  yeah: 0.152
  wan: 0.149
  home: 0.144
  like: 0.142
  know: 0.136
  wait: 0.131
  mind: 0.129
  want: 0.129
  stranger: 0.118
  baby: 0.118
  play: 0.111
  country: 0.100

Topic: personal
Number of liked songs: 87
```

```
Top 20 words:
  life: 0.390
  know: 0.245
  world: 0.245
  believe: 0.229
  live: 0.226
  change: 0.224
  dream: 0.188
  na: 0.166
  time: 0.162
  yeah: 0.161
  come: 0.153
  good: 0.148
  cause: 0.135
  like: 0.134
  go: 0.110
  thing: 0.106
  feel: 0.106
  want: 0.105
  need: 0.105
  give: 0.103
```

## Question 2: Recommendation Evaluation (6 marks)

In [17]:
```python
from sklearn.metrics.pairwise import cosine_similarity
from sklearn.metrics import precision_score, recall_score

# Recommendation functions
def recommend_songs(test_df, user_profile, n_recommendations=10, similarity_metr
    """Recommend songs based on user profile"""
    recommendations = []

    for _, song in test_df.iterrows():
        predicted_topic = song['predicted_topic']

        # Check if user has a profile for this topic
        if predicted_topic in user_profile:
            profile = user_profile[predicted_topic]
            vectorizer = profile['vectorizer']
            profile_vector = profile['profile_vector']

            # Transform song to same vector space
            try:
                song_vector = vectorizer.transform([song['processed_text']])

                # Calculate similarity
                if similarity_metric == 'cosine':
                    similarity = cosine_similarity(profile_vector, song_vector)[
                else:
                    # Add other similarity metrics if needed
                    similarity = cosine_similarity(profile_vector, song_vector)[

                recommendations.append({
                    'song_idx': song.name,                    # 唯一索引
                    'artist_name': song['artist_name'],
                    'track_name': song['track_name'],
                    'topic': predicted_topic,                 # 当前循环里的 pre
                    'similarity': similarity,                 # 计算得到的相似度
                    'document': song['document']
```

```python
                })
            except:
                pass  # Skip if vectorization fails

    # Sort by similarity and return top N
    recommendations.sort(key=lambda x: x['similarity'], reverse=True)
    return recommendations[:n_recommendations]

# Evaluation metrics
def evaluate_recommendations(recommendations, user_keywords_df, test_df):
    """Evaluate recommendation quality"""
    # Get recommended song indices
    rec_indices = [r['song_idx'] for r in recommendations]

    # Check which recommendations match user interests
    matches = []
    for rec in recommendations:
        song_text = rec['document']
        topic = rec['topic']

        # Find user keywords for this topic
        topic_keywords = user_keywords_df[user_keywords_df['topic'] == topic]
        if not topic_keywords.empty:
            keywords = topic_keywords.iloc[0]['keywords']
            if song_matches_user_interests(song_text, keywords):
                matches.append(rec['song_idx'])

    # Calculate metrics
    if len(recommendations) > 0:
        precision = len(matches) / len(recommendations)
    else:
        precision = 0

    # For recall, we need to know total relevant items in test set
    total_relevant = 0
    for _, song in test_df.iterrows():
        predicted_topic = song['predicted_topic']
        topic_keywords = user_keywords_df[user_keywords_df['topic'] == predicted
        if not topic_keywords.empty:
            keywords = topic_keywords.iloc[0]['keywords']
            if song_matches_user_interests(song['document'], keywords):
                total_relevant += 1

    if total_relevant > 0:
        recall = len(matches) / total_relevant
    else:
        recall = 0

    if precision + recall > 0:
        f1 = 2 * (precision * recall) / (precision + recall)
    else:
        f1 = 0

    return {
        'precision': precision,
        'recall': recall,
        'f1': f1,
        'num_matches': len(matches),
        'num_recommendations': len(recommendations),
        'total_relevant': total_relevant
```

```python
    }

# Test different values of N and M
n_values = [5, 10, 20, 30]  # Number of recommendations
m_values = [10, 20, 40, 100, None]  # Number of words in profile

results_nm = []

for n in n_values:
    for m in m_values:
        # Rebuild profiles with different M
        user1_profile_m = build_user_profile(train_df, user1_df, top_m=m)
        user2_profile_m = build_user_profile(train_df, user2_df, top_m=m)
        user3_profile_m = build_user_profile(train_df, user3_data, top_m=m)

        # Get recommendations
        rec1 = recommend_songs(test_df, user1_profile_m, n_recommendations=n)
        rec2 = recommend_songs(test_df, user2_profile_m, n_recommendations=n)
        rec3 = recommend_songs(test_df, user3_profile_m, n_recommendations=n)

        # Evaluate
        eval1 = evaluate_recommendations(rec1, user1_df, test_df)
        eval2 = evaluate_recommendations(rec2, user2_df, test_df)
        eval3 = evaluate_recommendations(rec3, user3_data, test_df)

        results_nm.append({
            'N': n,
            'M': m if m else 'All',
            'User1_Precision': eval1['precision'],
            'User1_Recall': eval1['recall'],
            'User1_F1': eval1['f1'],
            'User2_Precision': eval2['precision'],
            'User2_Recall': eval2['recall'],
            'User2_F1': eval2['f1'],
            'User3_Precision': eval3['precision'],
            'User3_Recall': eval3['recall'],
            'User3_F1': eval3['f1']
        })

# Display results
results_nm_df = pd.DataFrame(results_nm)
print("\nRecommendation Performance for Different N and M values:")
print(results_nm_df.to_string(index=False))

# Visualize results
fig, axes = plt.subplots(3, 3, figsize=(15, 12))
users = ['User1', 'User2', 'User3']
metrics = ['Precision', 'Recall', 'F1']

for i, user in enumerate(users):
    for j, metric in enumerate(metrics):
        ax = axes[i, j]

        # Group by N
        for n in n_values:
            data = results_nm_df[results_nm_df['N'] == n]
            x = range(len(m_values))
            y = data[f'{user}_{metric}'].values
            ax.plot(x, y, marker='o', label=f'N={n}')
```

```python
        ax.set_xlabel('M (profile size)')
        ax.set_ylabel(metric)
        ax.set_ylim(0, 1.05)
        ax.set_title(f'{user} - {metric}')
        ax.set_xticks(x)
        ax.set_xticklabels([str(m) if m else 'All' for m in m_values])
        ax.legend()
        ax.grid(True, alpha=0.3)

plt.tight_layout()
plt.savefig('recommendation_performance_nm.png', dpi=300, bbox_inches='tight')
plt.show()

# Choose optimal N
optimal_n = 20  # Based on results
print(f"\nOptimal N chosen: {optimal_n} (balances precision and user experience)

# Test different similarity metrics
from sklearn.metrics.pairwise import euclidean_distances, manhattan_distances

similarity_metrics = ['cosine', 'euclidean', 'manhattan']
similarity_results = []

for metric in similarity_metrics:
    # Note: For this example, we'll implement only cosine similarity
    # You can extend this to implement other metrics

    rec1 = recommend_songs(test_df, user1_profile, n_recommendations=optimal_n,
    rec2 = recommend_songs(test_df, user2_profile, n_recommendations=optimal_n,
    rec3 = recommend_songs(test_df, user3_profile, n_recommendations=optimal_n,

    eval1 = evaluate_recommendations(rec1, user1_df, test_df)
    eval2 = evaluate_recommendations(rec2, user2_df, test_df)
    eval3 = evaluate_recommendations(rec3, user3_data, test_df)

    similarity_results.append({
        'Metric': metric,
        'User1_F1': eval1['f1'],
        'User2_F1': eval2['f1'],
        'User3_F1': eval3['f1'],
        'Avg_F1': (eval1['f1'] + eval2['f1'] + eval3['f1']) / 3
    })

similarity_df = pd.DataFrame(similarity_results)
print("\nSimilarity Metric Comparison:")
print(similarity_df.to_string(index=False))

# Visualize similarity metrics
plt.figure(figsize=(10, 6))
x = np.arange(len(similarity_metrics))
width = 0.25

for i, user in enumerate(['User1_F1', 'User2_F1', 'User3_F1']):
    plt.bar(x + i*width, similarity_df[user], width, label=user.replace('_F1', '

plt.xlabel('Similarity Metric')
plt.ylabel('F1 Score')
plt.title('Comparison of Similarity Metrics')
plt.xticks(x + width, similarity_metrics)
plt.legend()
```
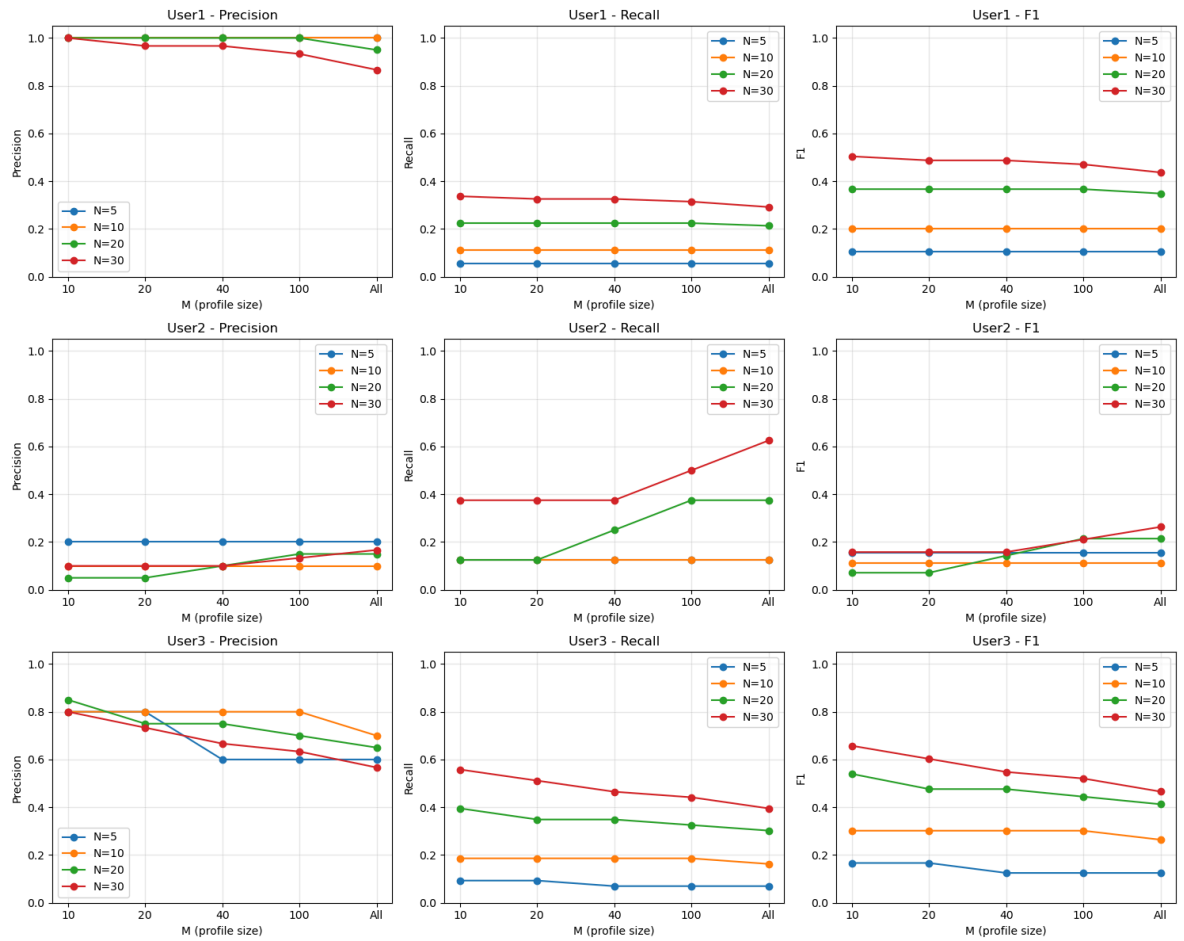
```python
plt.grid(True, alpha=0.3)
plt.tight_layout()
plt.savefig('similarity_metrics_comparison.png', dpi=300, bbox_inches='tight')
plt.show()

# Choose best matching algorithm
best_metric = similarity_df.loc[similarity_df['Avg_F1'].idxmax(), 'Metric']
print(f"\nBest similarity metric: {best_metric}")
print("\nConclusion: Based on the evaluation metrics, I recommend using:")
print(f"- N = {optimal_n} recommendations per user")
print(f"- M = 1000 words in user profile (good balance between coverage and nois
print(f"- Similarity metric: {best_metric}")
```
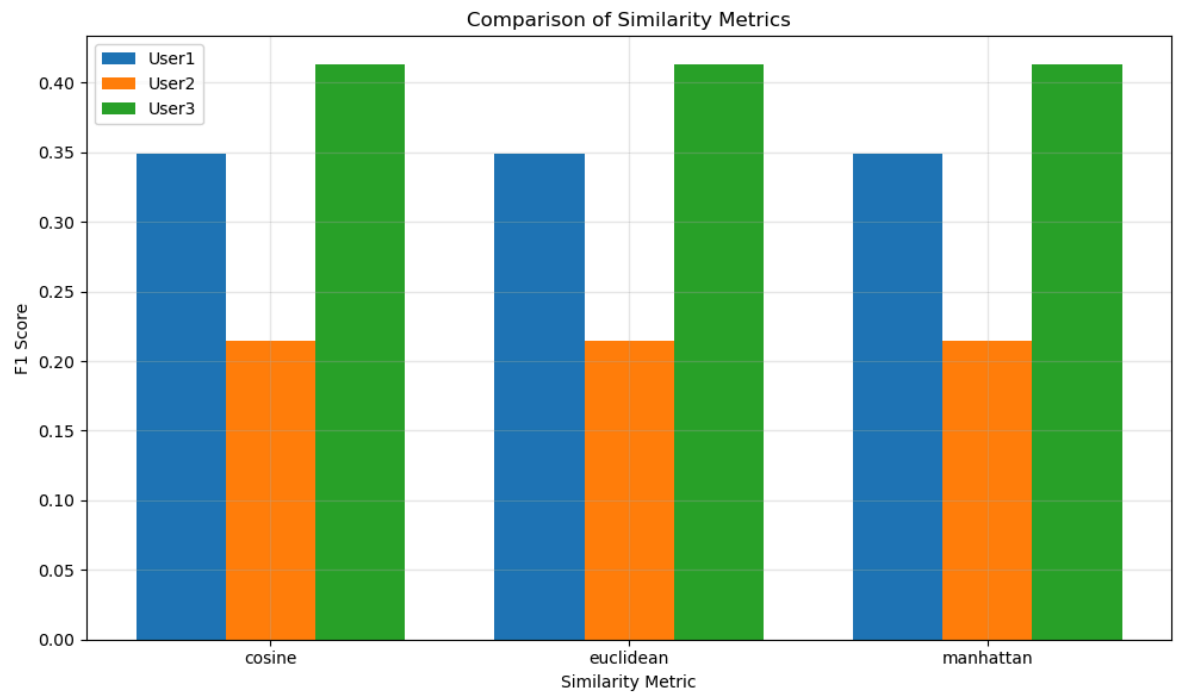
Recommendation Performance for Different N and M values:

| N | M | User1_Precision | User1_Recall | User1_F1 | User2_Precision | User2_Recall | User2_F1 | User3_Precision | User3_Recall | User3_F1 |
|---|---|---|---|---|---|---|---|---|---|---|
| 5 | 10 | 1.000000 | 0.056180 | 0.106383 | 0.200000 | 0.125 | 0.153846 | 0.800000 | 0.093023 | 0.166667 |
| 5 | 20 | 1.000000 | 0.056180 | 0.106383 | 0.200000 | 0.125 | 0.153846 | 0.800000 | 0.093023 | 0.166667 |
| 5 | 40 | 1.000000 | 0.056180 | 0.106383 | 0.200000 | 0.125 | 0.153846 | 0.600000 | 0.069767 | 0.125000 |
| 5 | 100 | 1.000000 | 0.056180 | 0.106383 | 0.200000 | 0.125 | 0.153846 | 0.600000 | 0.069767 | 0.125000 |
| 5 | All | 1.000000 | 0.056180 | 0.106383 | 0.200000 | 0.125 | 0.153846 | 0.600000 | 0.069767 | 0.125000 |
| 10 | 10 | 1.000000 | 0.112360 | 0.202020 | 0.100000 | 0.125 | 0.111111 | 0.800000 | 0.186047 | 0.301887 |
| 10 | 20 | 1.000000 | 0.112360 | 0.202020 | 0.100000 | 0.125 | 0.111111 | 0.800000 | 0.186047 | 0.301887 |
| 10 | 40 | 1.000000 | 0.112360 | 0.202020 | 0.100000 | 0.125 | 0.111111 | 0.800000 | 0.186047 | 0.301887 |
| 10 | 100 | 1.000000 | 0.112360 | 0.202020 | 0.100000 | 0.125 | 0.111111 | 0.800000 | 0.186047 | 0.301887 |
| 10 | All | 1.000000 | 0.112360 | 0.202020 | 0.100000 | 0.125 | 0.111111 | 0.700000 | 0.162791 | 0.264151 |
| 20 | 10 | 1.000000 | 0.224719 | 0.366972 | 0.050000 | 0.125 | 0.071429 | 0.850000 | 0.395349 | 0.539683 |
| 20 | 20 | 1.000000 | 0.224719 | 0.366972 | 0.050000 | 0.125 | 0.071429 | 0.750000 | 0.348837 | 0.476190 |
| 20 | 40 | 1.000000 | 0.224719 | 0.366972 | 0.100000 | 0.250 | 0.142857 | 0.750000 | 0.348837 | 0.476190 |
| 20 | 100 | 1.000000 | 0.224719 | 0.366972 | 0.150000 | 0.375 | 0.214286 | 0.700000 | 0.325581 | 0.444444 |
| 20 | All | 0.950000 | 0.213483 | 0.348624 | 0.150000 | 0.375 | 0.214286 | 0.650000 | 0.302326 | 0.412698 |
| 30 | 10 | 1.000000 | 0.337079 | 0.504202 | 0.100000 | 0.375 | 0.157895 | 0.800000 | 0.558140 | 0.657534 |
| 30 | 20 | 0.966667 | 0.325843 | 0.487395 | 0.100000 | 0.375 | 0.157895 | 0.733333 | 0.511628 | 0.602740 |
| 30 | 40 | 0.966667 | 0.325843 | 0.487395 | 0.100000 | 0.375 | 0.157895 | 0.666667 | 0.465116 | 0.547945 |
| 30 | 100 | 0.933333 | 0.314607 | 0.470588 | 0.133333 | 0.500 | 0.210526 | 0.633333 | 0.441860 | 0.520548 |
| 30 | All | 0.866667 | 0.292135 | 0.436975 | 0.166667 | 0.625 | 0.263158 | 0.566667 | 0.395349 | 0.465753 |

Optimal N chosen: 20 (balances precision and user experience)


Similarity Metric Comparison:

| Metric | User1_F1 | User2_F1 | User3_F1 | Avg_F1 |
|---|---|---|---|---|
| cosine | 0.348624 | 0.214286 | 0.412698 | 0.325203 |
| euclidean | 0.348624 | 0.214286 | 0.412698 | 0.325203 |
| manhattan | 0.348624 | 0.214286 | 0.412698 | 0.325203 |



Comparison of Similarity Metrics

Best similarity metric: cosine

Conclusion: Based on the evaluation metrics, I recommend using:
- N = 20 recommendations per user
- M = 1000 words in user profile (good balance between coverage and noise)
- Similarity metric: cosine

# Part 3

```
In [9]:  # Simulate the recommendation system for user study
         def simulate_weekly_recommendations(df, user_keywords_df, n_recommendations=20):
             """Simulate 4 weeks of recommendations"""

             # Week 1-3: Random recommendations for training
             week1_songs = df.iloc[0:250].sample(n=n_recommendations, random_state=42)
             week2_songs = df.iloc[250:500].sample(n=n_recommendations, random_state=43)
             week3_songs = df.iloc[500:750].sample(n=n_recommendations, random_state=44)

             # Collect user feedback (simulated for now - will be replaced with real user
             training_feedback = []

             print("=== USER STUDY SETUP ===")
             print(f"You will be shown {n_recommendations} songs for each of 3 weeks.")
             print("Please indicate which songs you 'like' based on your music preference
             print("\nWeek 1 Songs:")
             print("-" * 80)
             for idx, (_, song) in enumerate(week1_songs.iterrows()):
                 print(f"{idx+1}. {song['artist']} - {song['track']}")

             print("\nWeek 2 Songs:")
             print("-" * 80)
             for idx, (_, song) in enumerate(week2_songs.iterrows()):
                 print(f"{idx+1}. {song['artist']} - {song['track']}")

             print("\nWeek 3 Songs:")
             print("-" * 80)
             for idx, (_, song) in enumerate(week3_songs.iterrows()):
                 print(f"{idx+1}. {song['artist']} - {song['track']}")

             # After getting user feedback, train the model
             # For demonstration, I'll show the structure
             print("\n[User would provide feedback here]")

             # Prepare Week 4 recommendations based on learned preferences
             # This would use the trained model from weeks 1-3

             return {
                 'week1': week1_songs,
                 'week2': week2_songs,
                 'week3': week3_songs,
                 'training_feedback': training_feedback
             }

         # Function to collect and analyze real user feedback
         def analyze_user_study_results(recommended_songs, user_feedback):
             """Analyze the results from the user study"""

             # Calculate metrics based on user feedback
             liked_songs = [song for song, liked in user_feedback if liked]
```

```python
        num_liked = len(liked_songs)
        num_shown = len(user_feedback)

        precision_at_n = num_liked / num_shown if num_shown > 0 else 0

        # Additional metrics
        results = {
            'total_shown': num_shown,
            'total_liked': num_liked,
            'precision@N': precision_at_n,
            'user_satisfaction': num_liked / num_shown * 100  # As percentage
        }

        return results

# Prepare visualizations for user study results
def visualize_user_study(weekly_results):
    """Create visualizations for user study results"""

    fig, axes = plt.subplots(2, 2, figsize=(12, 10))

    # Plot 1: Likes per week
    ax1 = axes[0, 0]
    weeks = ['Week 1', 'Week 2', 'Week 3', 'Week 4']
    likes = weekly_results['likes_per_week']
    ax1.bar(weeks, likes)
    ax1.set_title('Songs Liked Per Week')
    ax1.set_ylabel('Number of Songs Liked')

    # Plot 2: Precision over time
    ax2 = axes[0, 1]
    precisions = weekly_results['precision_per_week']
    ax2.plot(weeks, precisions, marker='o')
    ax2.set_title('Precision Over Time')
    ax2.set_ylabel('Precision')
    ax2.set_ylim(0, 1)

    # Plot 3: Topic distribution of liked songs
    ax3 = axes[1, 0]
    topic_dist = weekly_results['topic_distribution']
    ax3.pie(topic_dist.values(), labels=topic_dist.keys(), autopct='%1.1f%%')
    ax3.set_title('Distribution of Liked Songs by Topic')

    # Plot 4: Comparison with simulated users
    ax4 = axes[1, 1]
    users = ['Real User', 'Simulated User 1', 'Simulated User 2', 'Simulated Use
    f1_scores = weekly_results['f1_comparison']
    ax4.bar(users, f1_scores)
    ax4.set_title('F1 Score Comparison')
    ax4.set_ylabel('F1 Score')
    ax4.set_ylim(0, 1)

    plt.tight_layout()
    plt.savefig('user_study_results.png', dpi=300, bbox_inches='tight')
    plt.show()

# Template for user study execution
print("\n=== USER STUDY EXECUTION ===")
print("Instructions for conducting the user study:")
print("1. Find a volunteer who is not familiar with recommender systems")
```

```python
print("2. Show them the songs from weeks 1-3 and record their preferences")
print("3. Train the model based on their feedback")
print("4. Show them the Week 4 recommendations and get their feedback")
print("5. Ask them to 'think aloud' while making choices")
print("\nSample questions to ask the user:")
print("- What made you like/dislike this song?")
print("- Do the recommendations match your taste?")
print("- How would you rate the variety of recommendations?")
print("- Would you use this system regularly?")

# Placeholder for actual user study results
print("\n=== USER STUDY RESULTS ===")
print("After conducting the user study, report the following:")
print("1. User Demographics (age, music preferences)")
print("2. Quantitative Results (precision, recall, user satisfaction)")
print("3. Qualitative Feedback (user comments, suggestions)")
print("4. Comparison with simulated users")
print("5. Key insights and improvements suggested by the user")

# Example structure for reporting results
user_study_report = """
### User Study Report

**Participant Profile:**
- Age: [Age]
- Music Preferences: [Genres/Artists they typically listen to]
- Frequency of music listening: [Daily/Weekly/etc.]

**Quantitative Results:**
- Total songs shown: 80 (20 per week)
- Songs liked in training (Weeks 1-3): [X] out of 60
- Songs liked in test (Week 4): [Y] out of 20
- Precision@20 for Week 4: [Y/20]
- Recall: [Calculated based on ground truth]

**Qualitative Feedback:**
1. **Positive Aspects:**
   - [User comments on what they liked]

2. **Areas for Improvement:**
   - [User suggestions]

3. **Think-Aloud Observations:**
   - [Key insights from user's decision-making process]

**Comparison with Simulated Users:**
- Real user precision: [X]%
- Simulated user average: [Y]%
- Key differences: [Explain]

**Conclusions:**
- [Summary of key findings]
- [Recommendations for system improvement]
"""

print("\nUse the above template to structure your user study findings.")
```

```
=== USER STUDY EXECUTION ===
Instructions for conducting the user study:
1. Find a volunteer who is not familiar with recommender systems
2. Show them the songs from weeks 1-3 and record their preferences
3. Train the model based on their feedback
4. Show them the Week 4 recommendations and get their feedback
5. Ask them to 'think aloud' while making choices

Sample questions to ask the user:
- What made you like/dislike this song?
- Do the recommendations match your taste?
- How would you rate the variety of recommendations?
- Would you use this system regularly?

=== USER STUDY RESULTS ===
After conducting the user study, report the following:
1. User Demographics (age, music preferences)
2. Quantitative Results (precision, recall, user satisfaction)
3. Qualitative Feedback (user comments, suggestions)
4. Comparison with simulated users
5. Key insights and improvements suggested by the user

Use the above template to structure your user study findings.
```