## Top 10 películas (rating por promedio de reseñas)



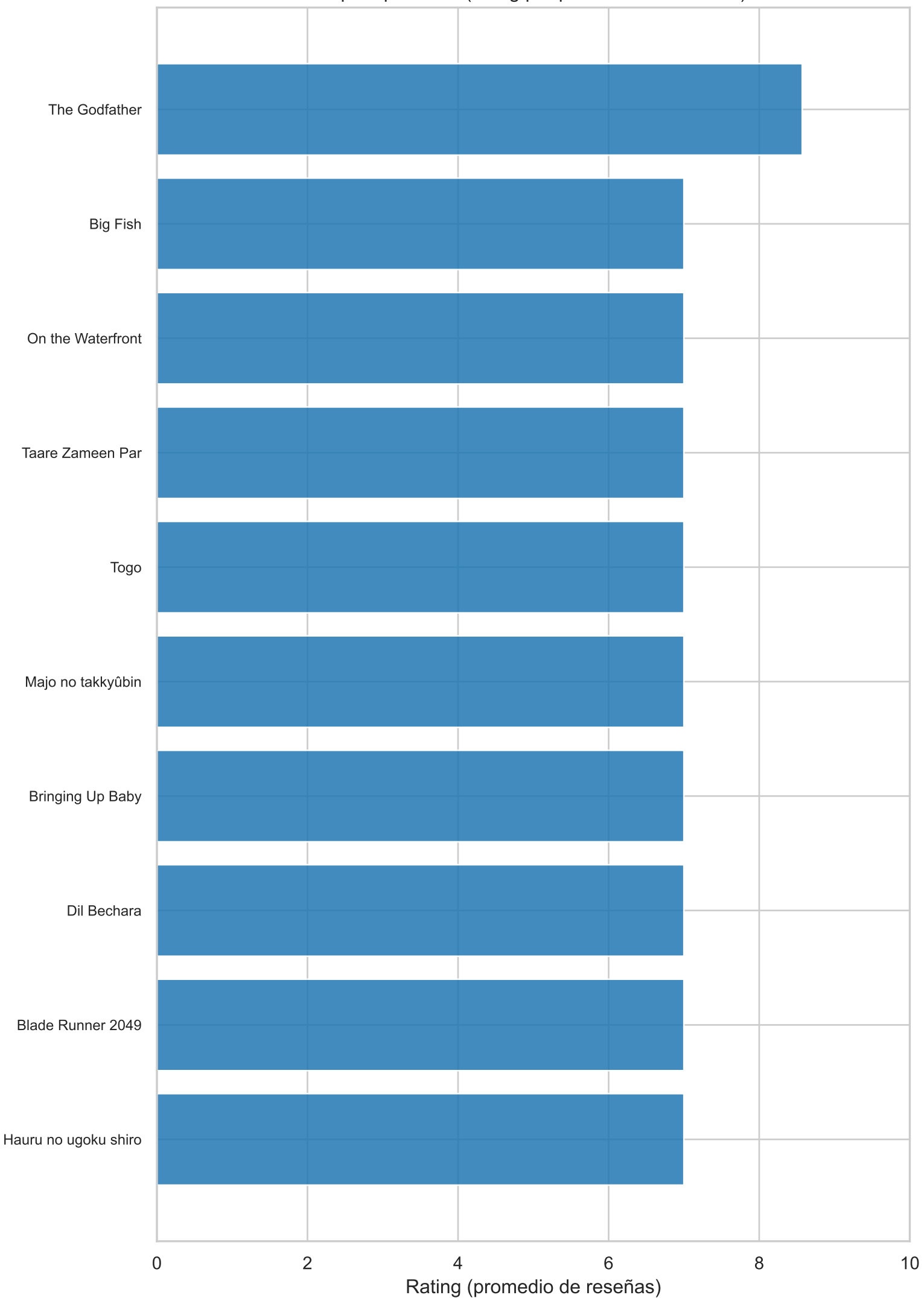| Película | Rating |
|---|---|
| The Godfather | (≈8.5) |
| Big Fish | (≈7) |
| On the Waterfront | (≈7) |
| Taare Zameen Par | (≈7) |
| Togo | (≈7) |
| Majo no takkyûbin | (≈7) |
| Bringing Up Baby | (≈7) |
| Dil Bechara | (≈7) |
| Blade Runner 2049 | (≈7) |
| Hauru no ugoku shiro | (≈7) |

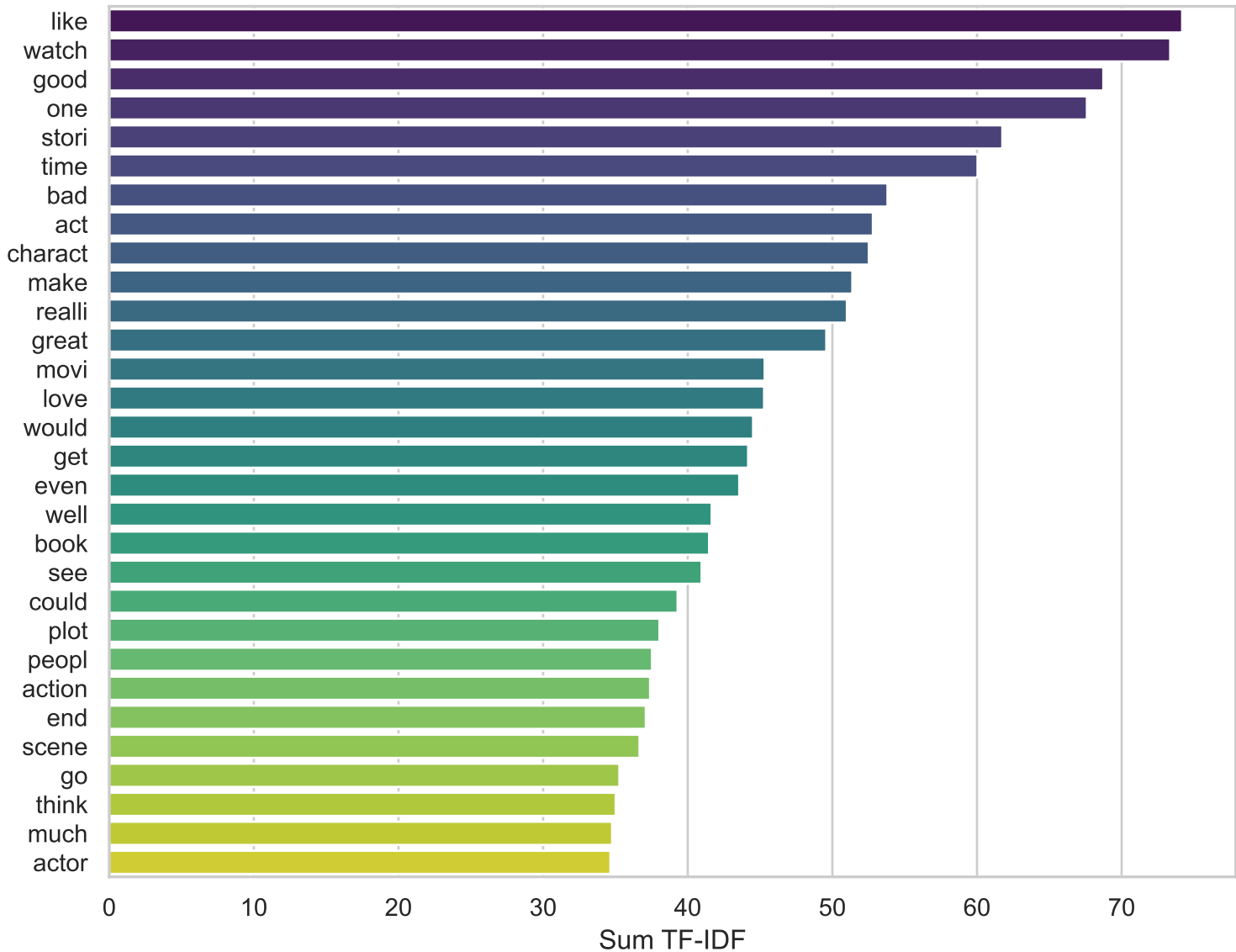Rating (promedio de reseñas)

IMDB - Model Results

Total reviews in DB: 5458
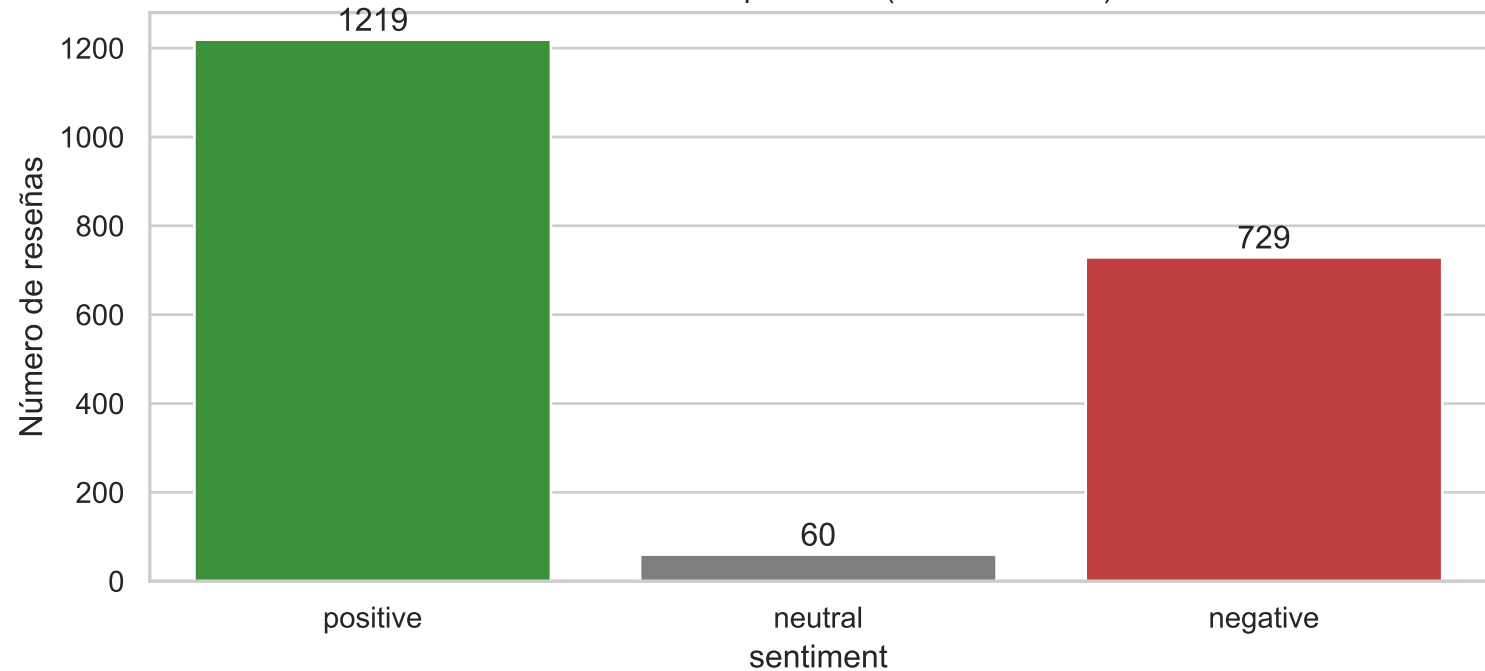Reviews used for modeling (first 4000): 3997
Reviews used for PDF visuals/evaluation (JSON): 2008

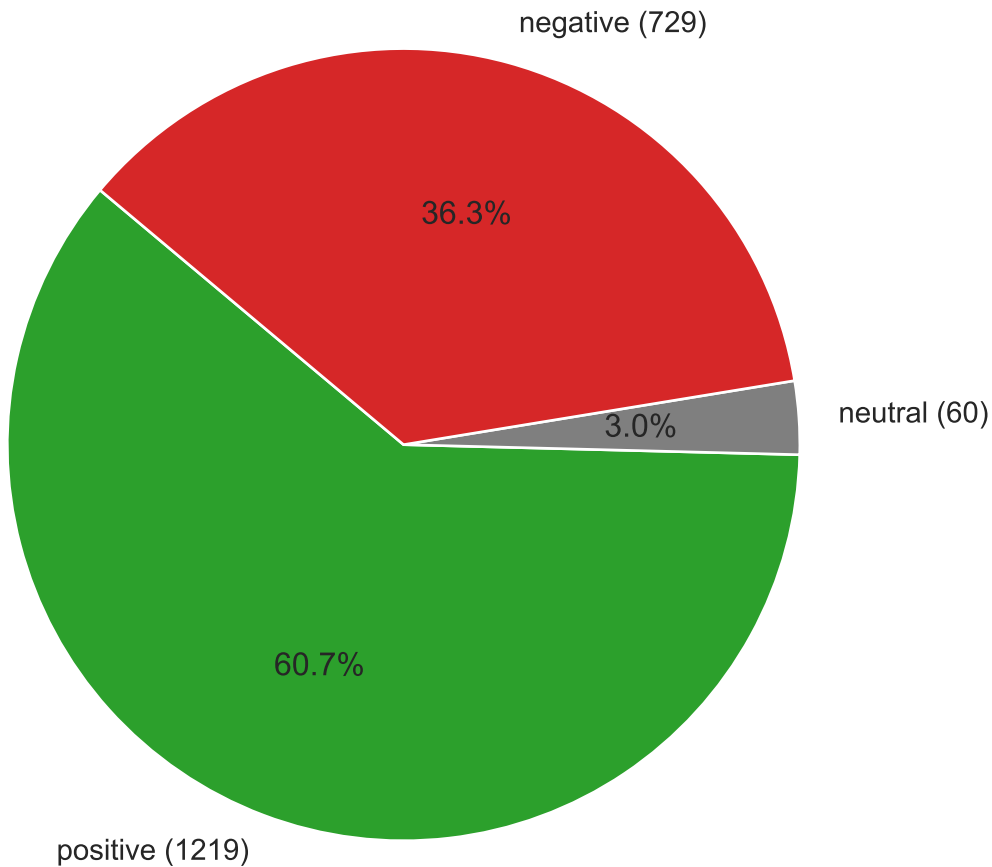Models trained: Naive Bayes, Linear SVM, Decision Tree
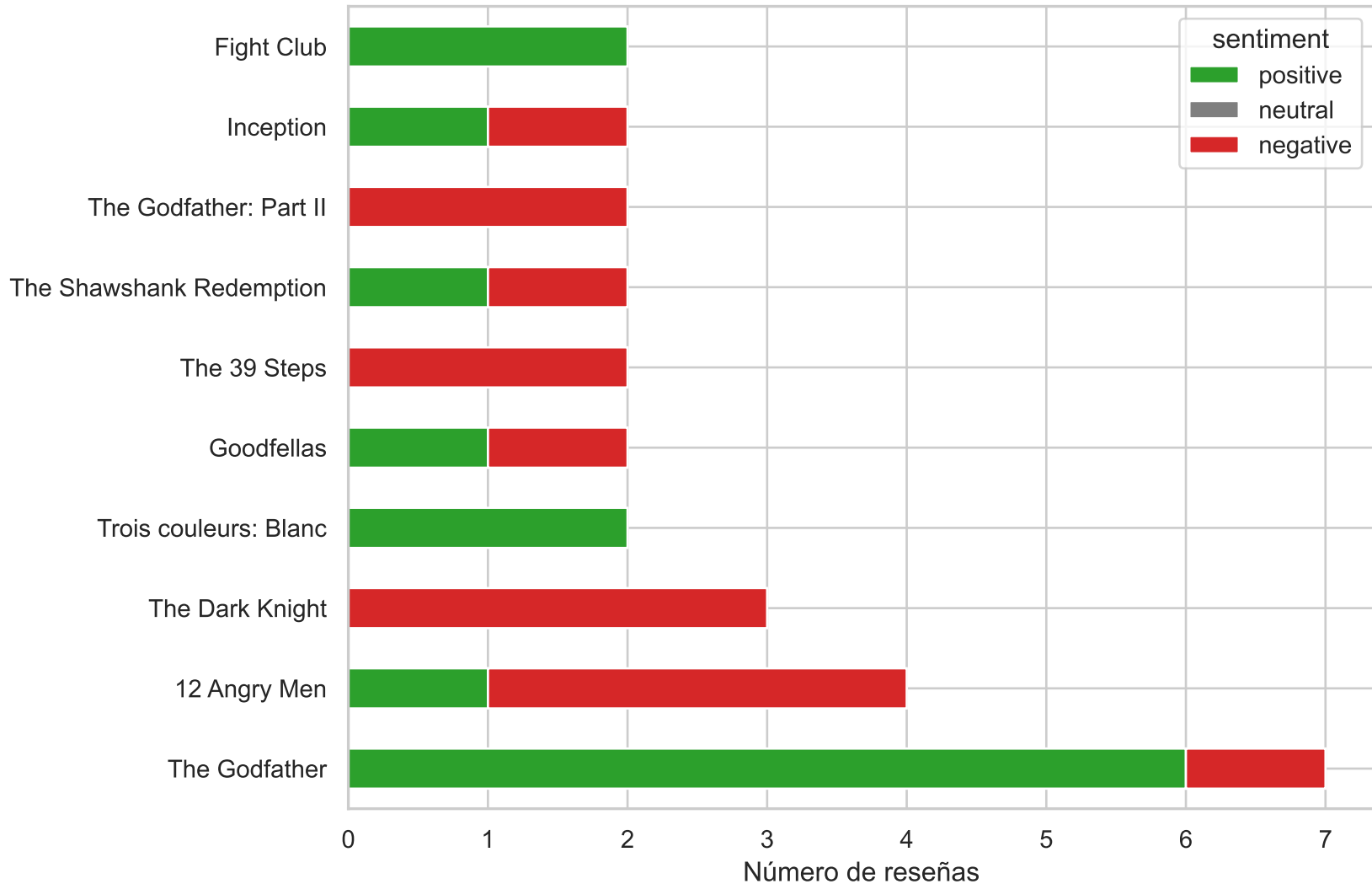
Top TF-IDF terms (reviews from JSON)

| Term | Sum TF-IDF |
|---|---|
| like | |
| watch | |
| good | |
| one | |
| stori | |
| time | |
| bad | |
| act | |
| charact | |
| make | |
| realli | |
| great | |
| movi | |
| love | |
| would | |
| get | |
| even | |
| well | |
| book | |
| see | |
| could | |
| plot | |
| peopl | |
| action | |
| end | |
| scene | |
| go | |
| think | |
| much | |
| actor | |

Distribución de polaridad (reseñas JSON)

Porcentaje de polaridad (reseñas JSON)

negative (729)

36.3%

neutral (60)

3.0%

60.7%

positive (1219)
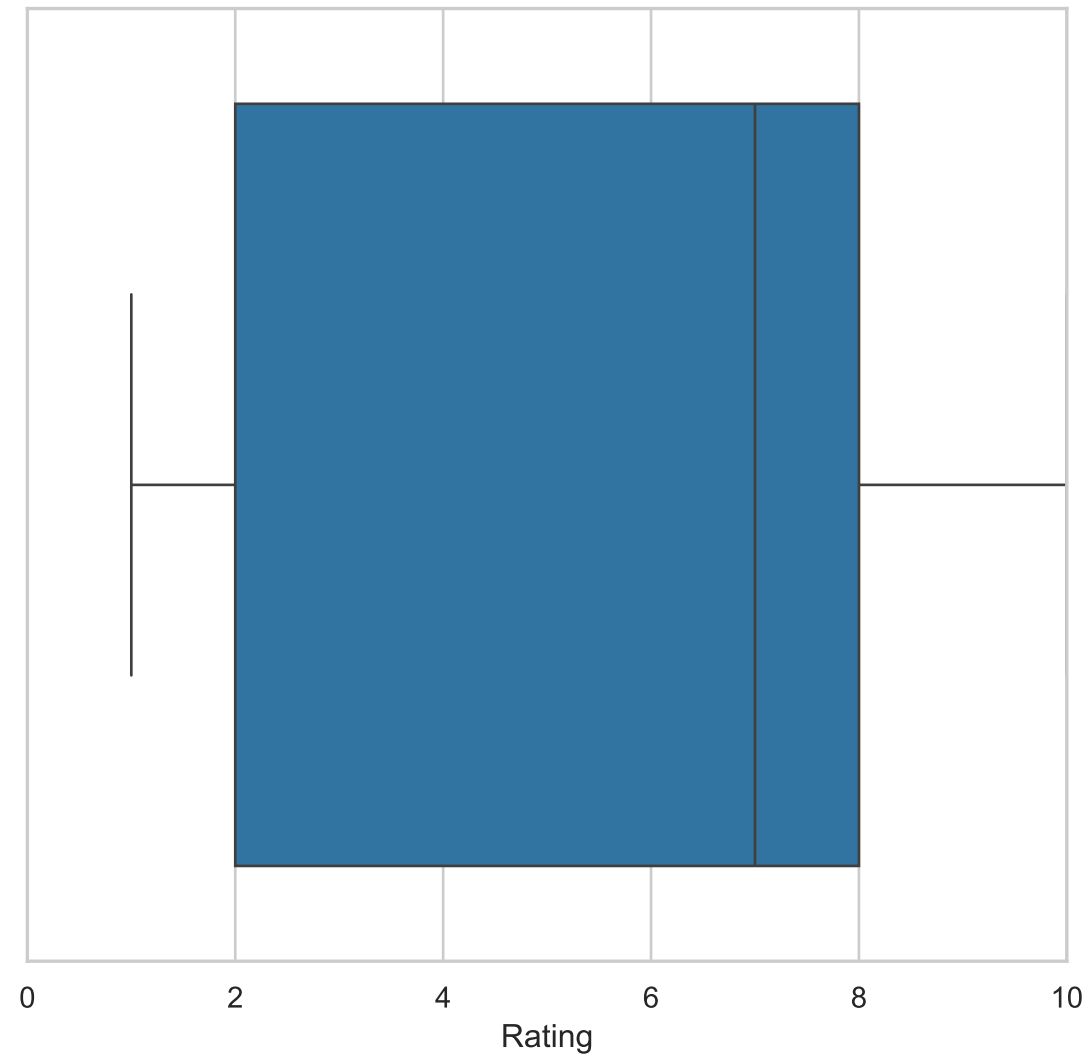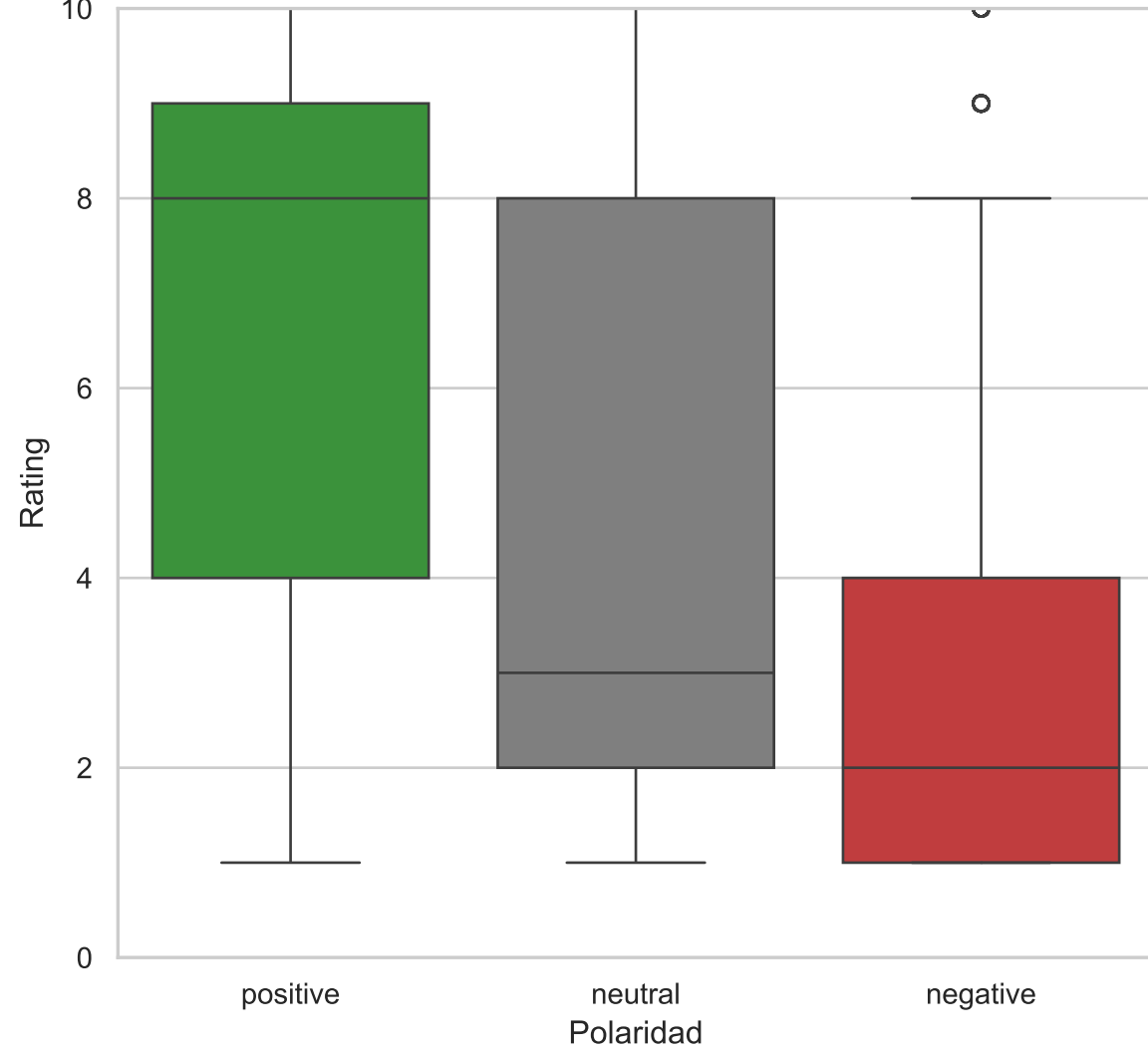
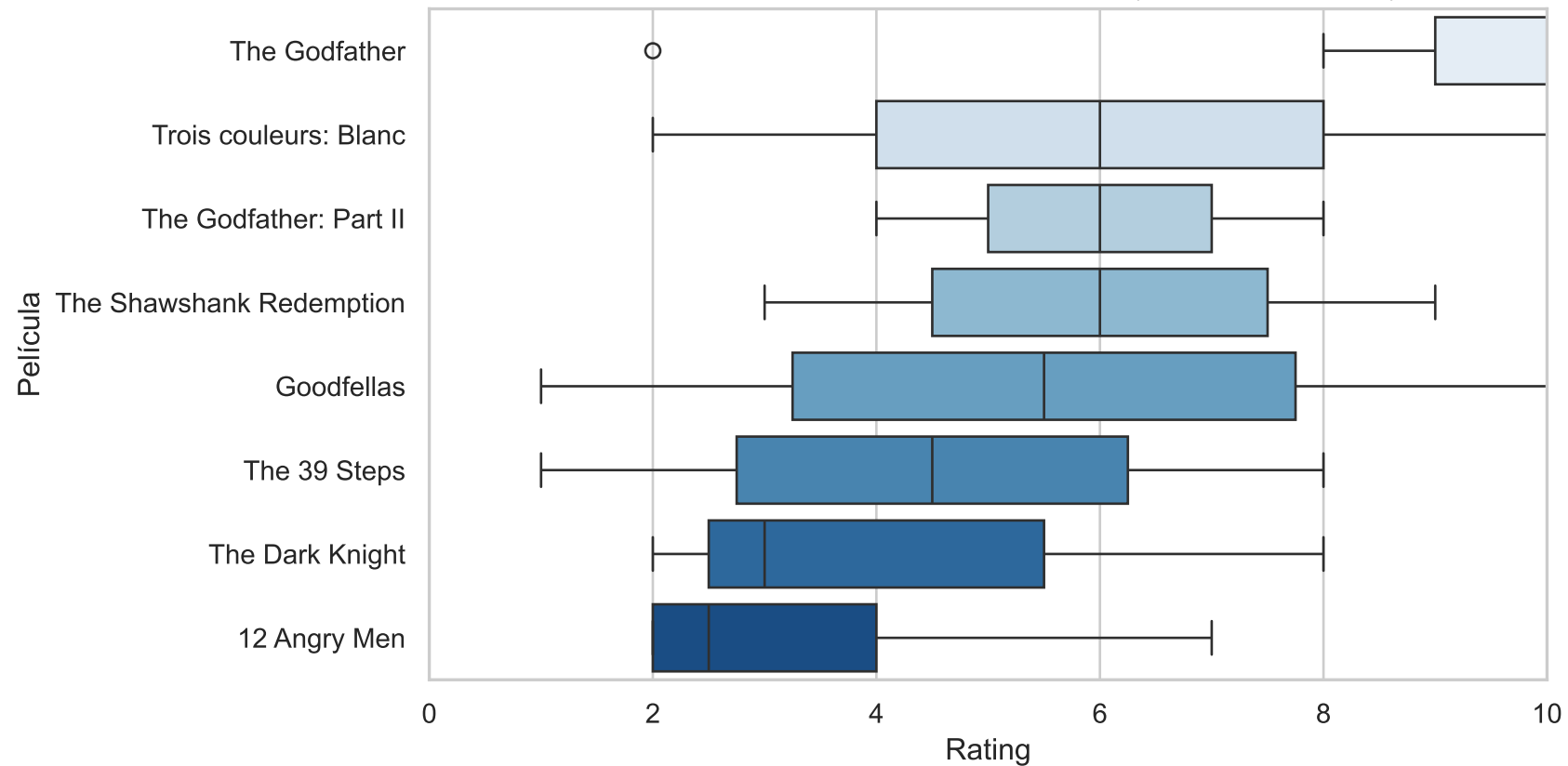Polaridad por película (top 10 por #reseñas en JSON)

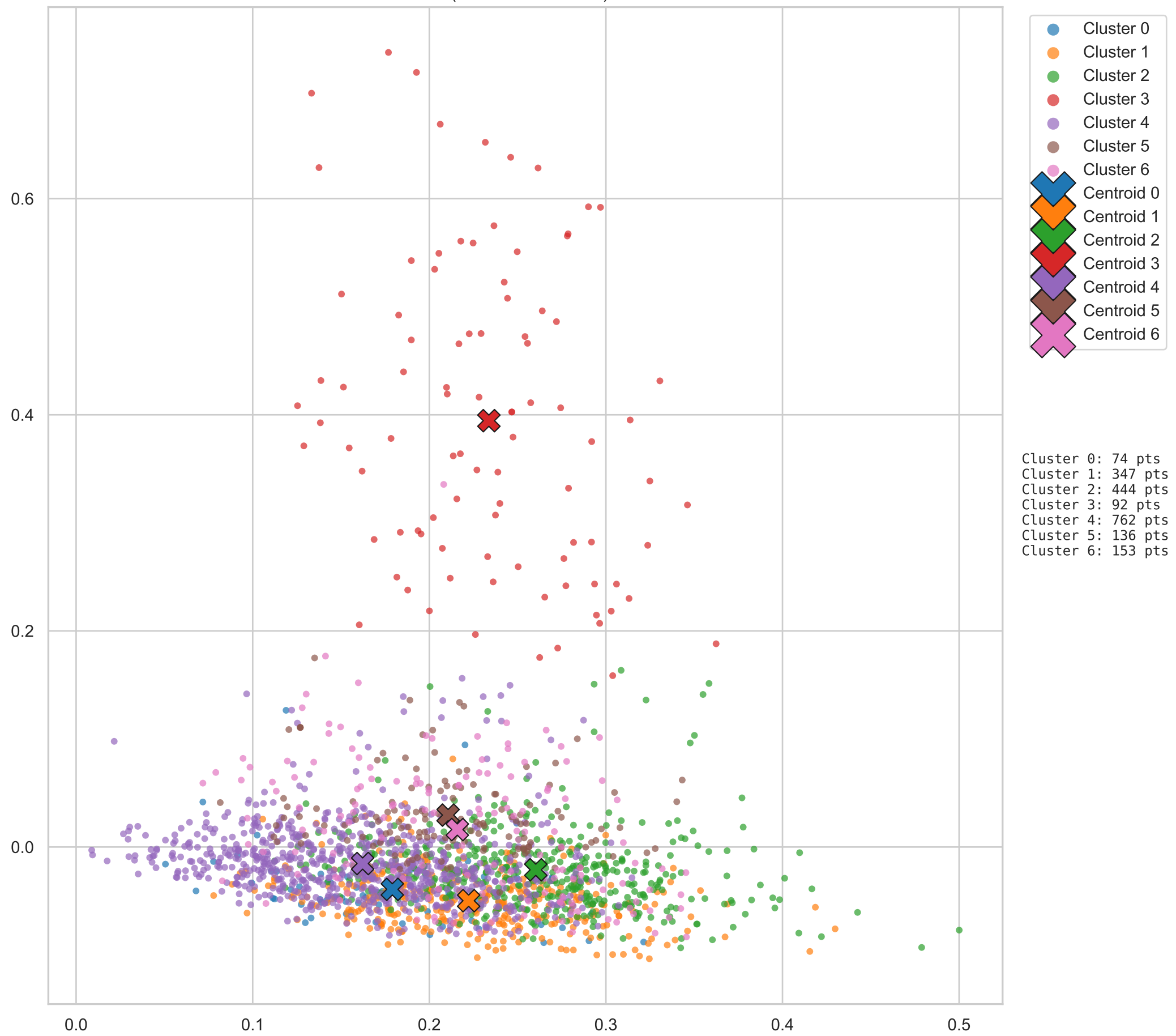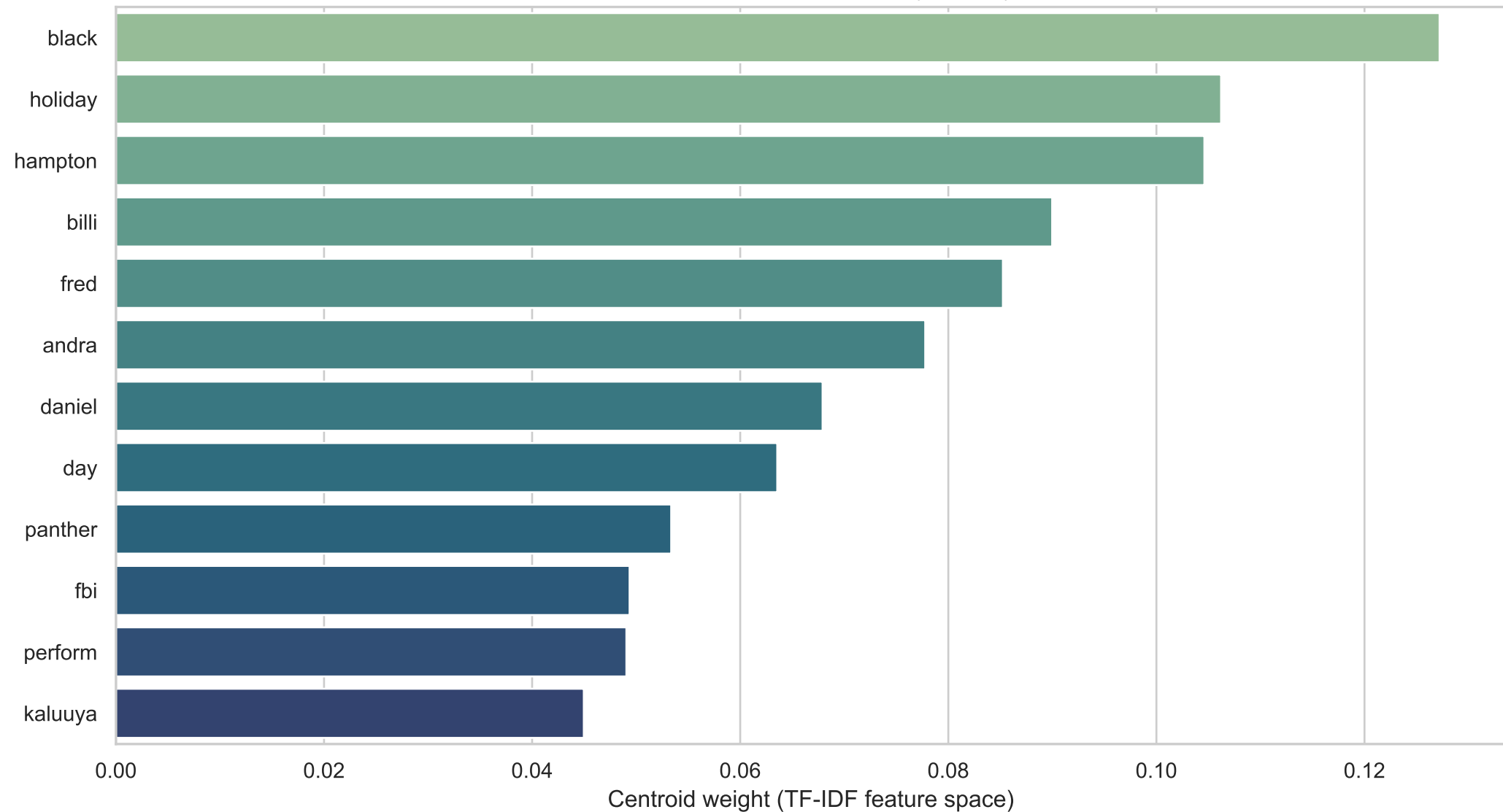**Distribución de calificaciones (boxplot) - global** y **Distribución de calificaciones por polaridad (boxplot)**
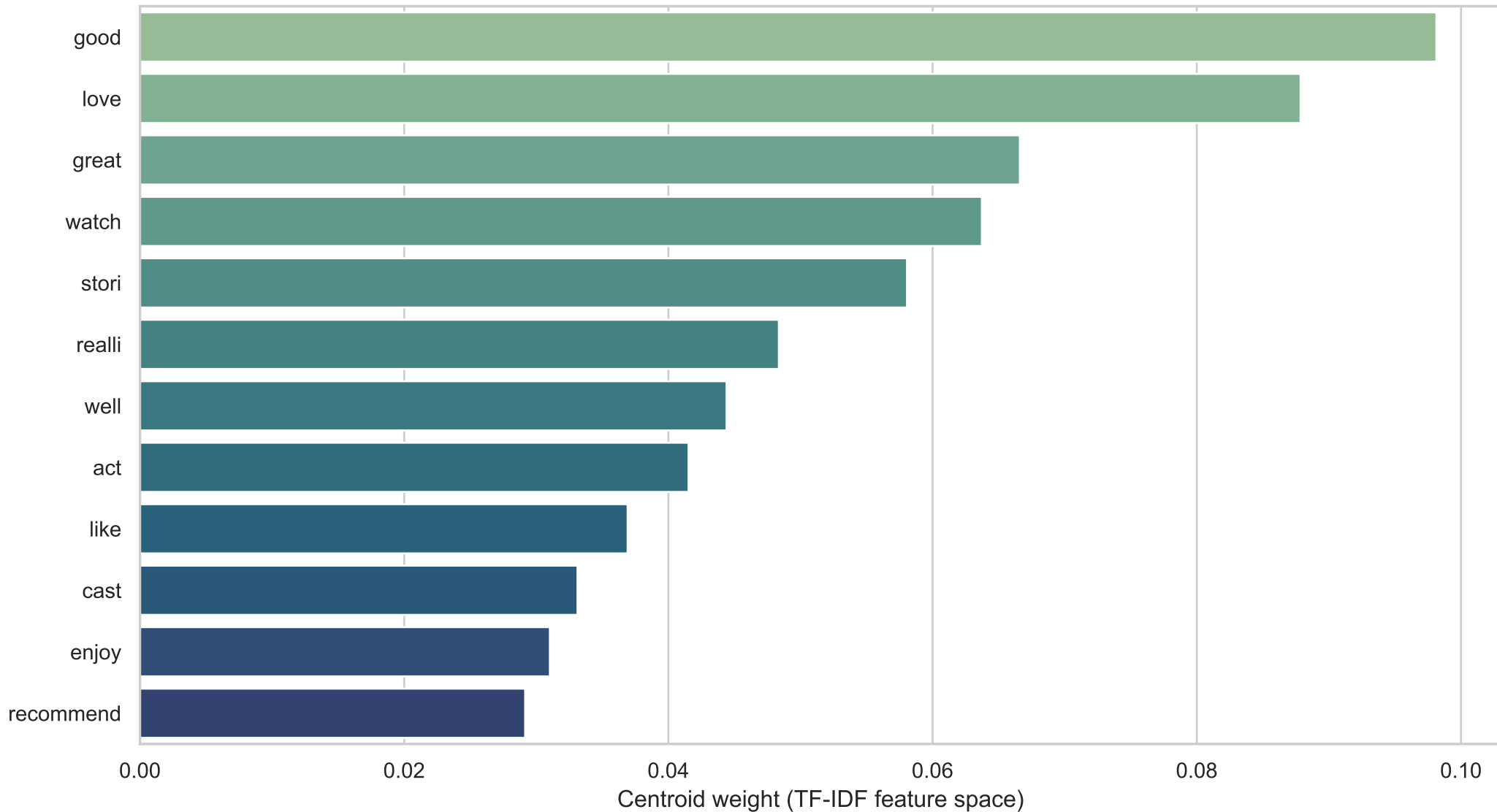
Distribución de calificaciones por película (top 8 por #reseñas)

KMeans clusters (TruncatedSVD 2D) - JSON reviews

Cluster 0
Cluster 1
Cluster 2
Cluster 3
Cluster 4
Cluster 5
Cluster 6
Centroid 0
Centroid 1
Centroid 2
Centroid 3
Centroid 4
Centroid 5
Centroid 6

Cluster 0: 74 pts
Cluster 1: 347 pts
Cluster 2: 444 pts
Cluster 3: 92 pts
Cluster 4: 762 pts
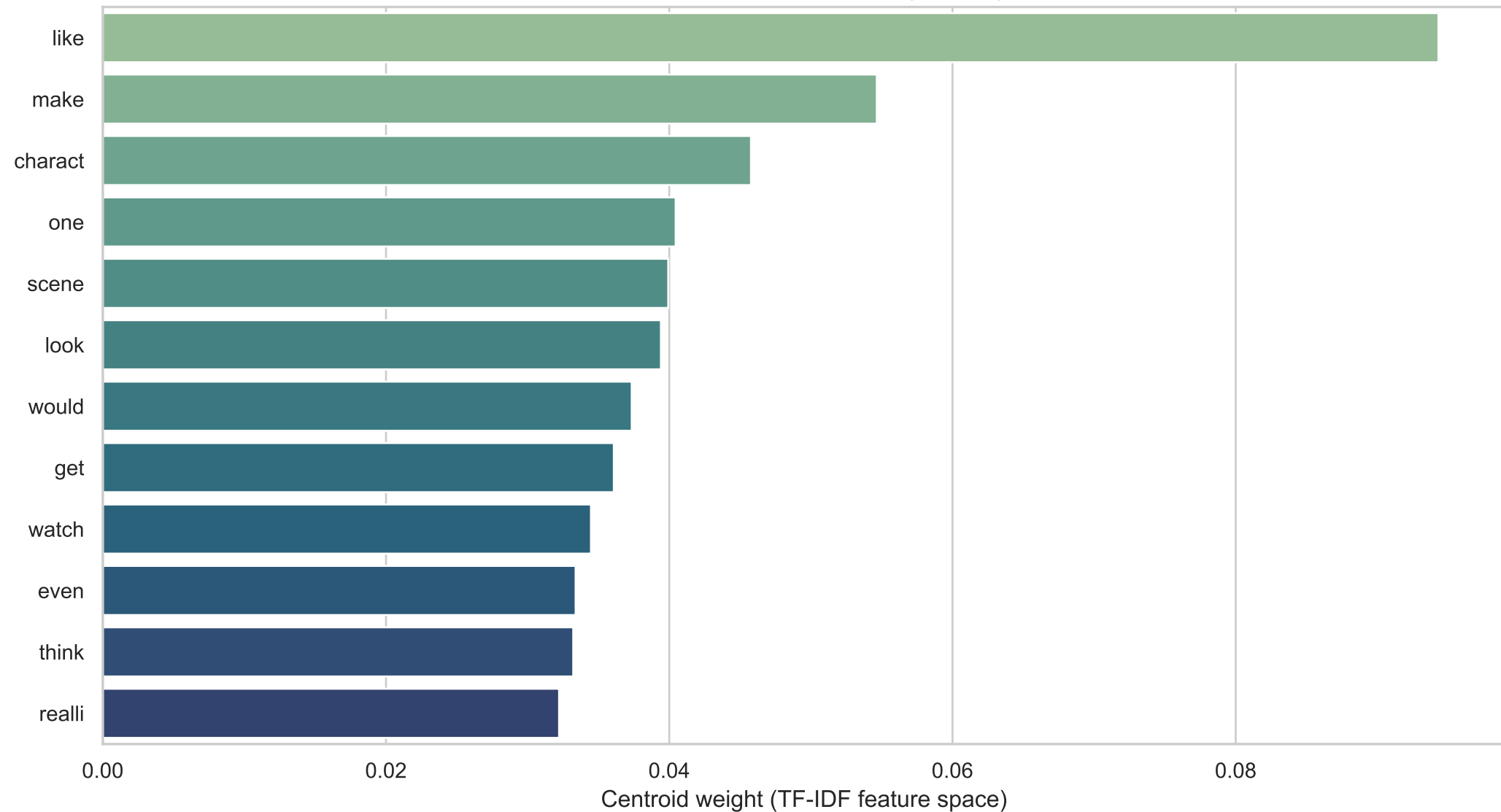Cluster 5: 136 pts
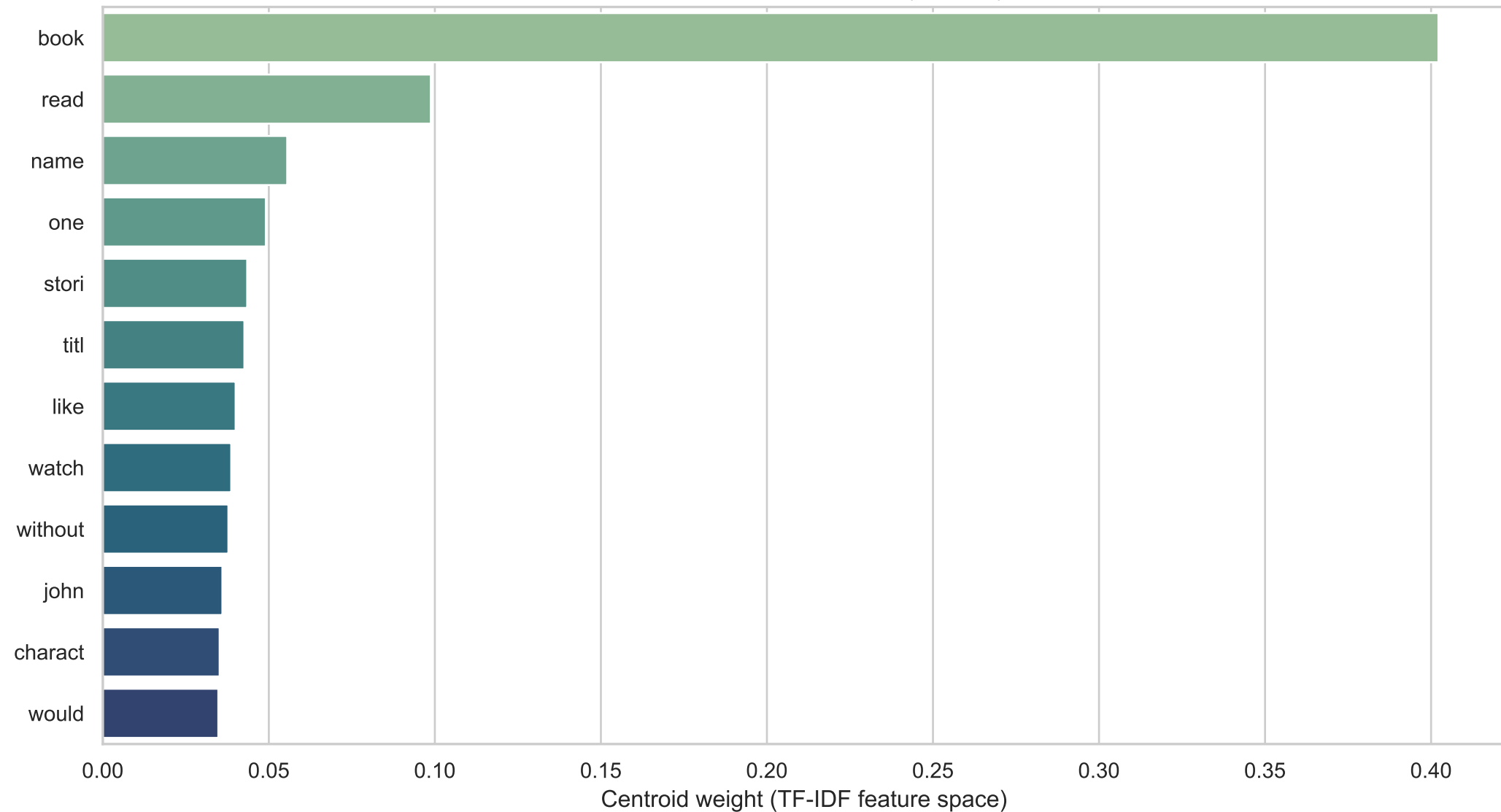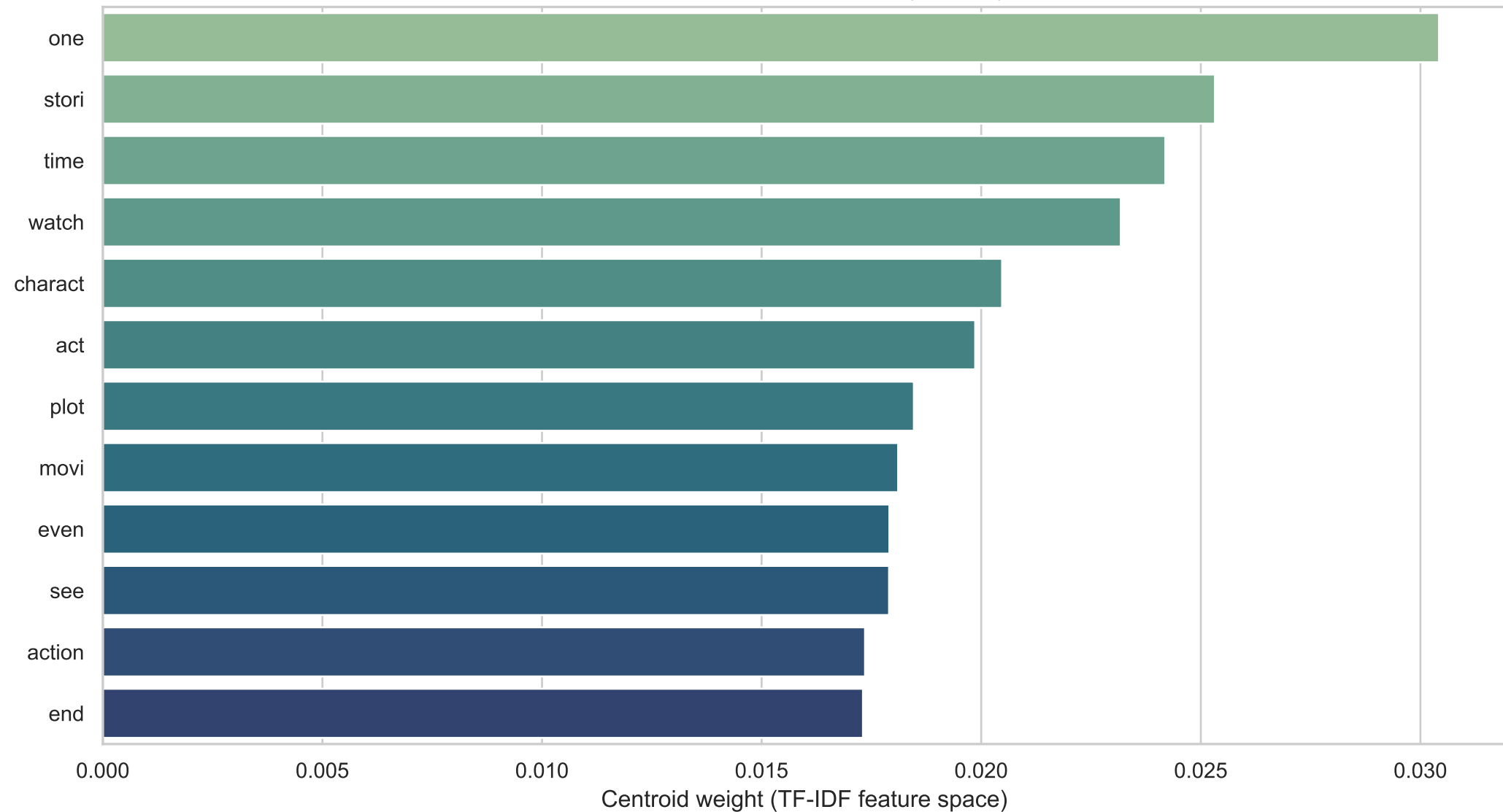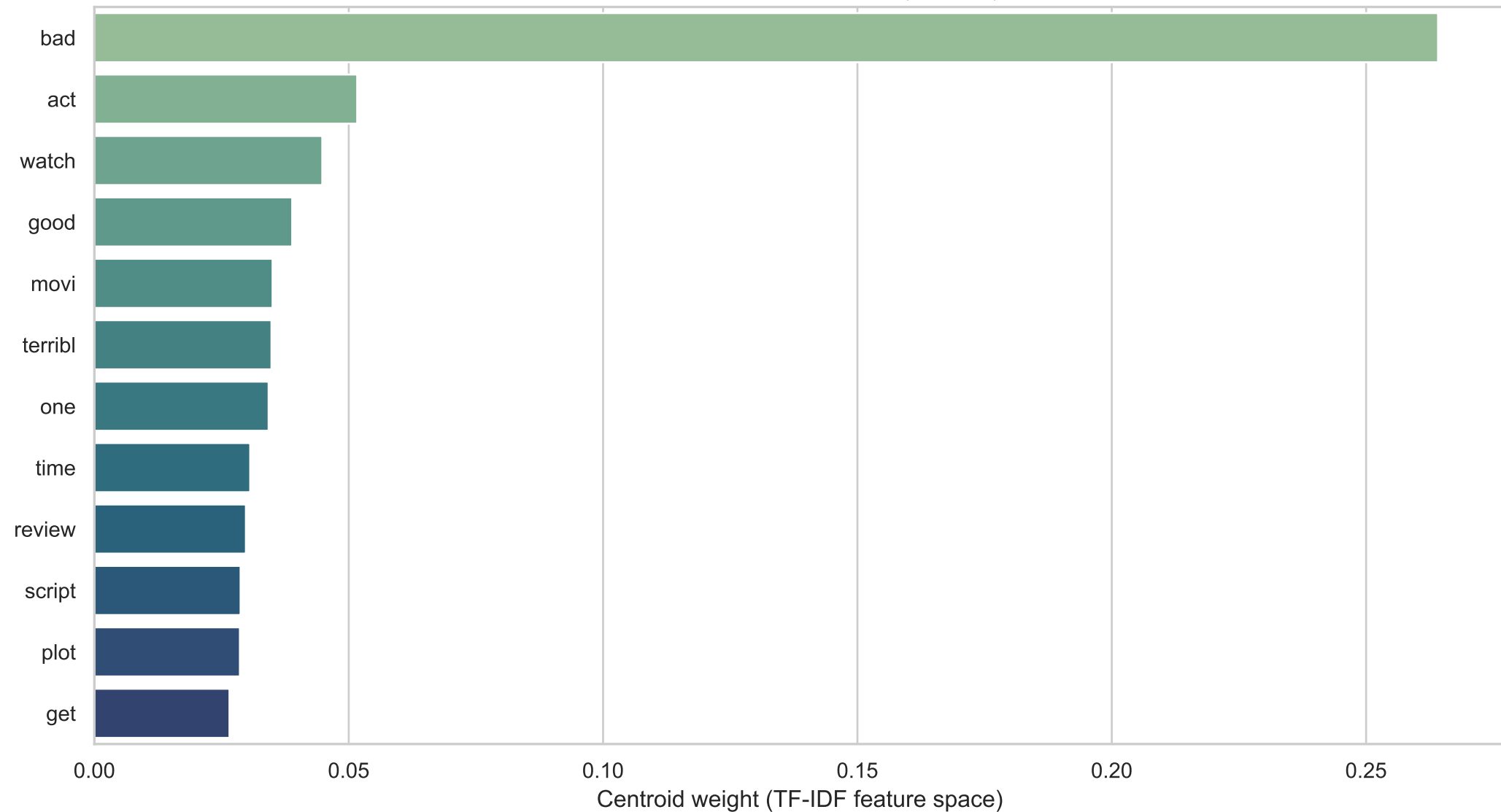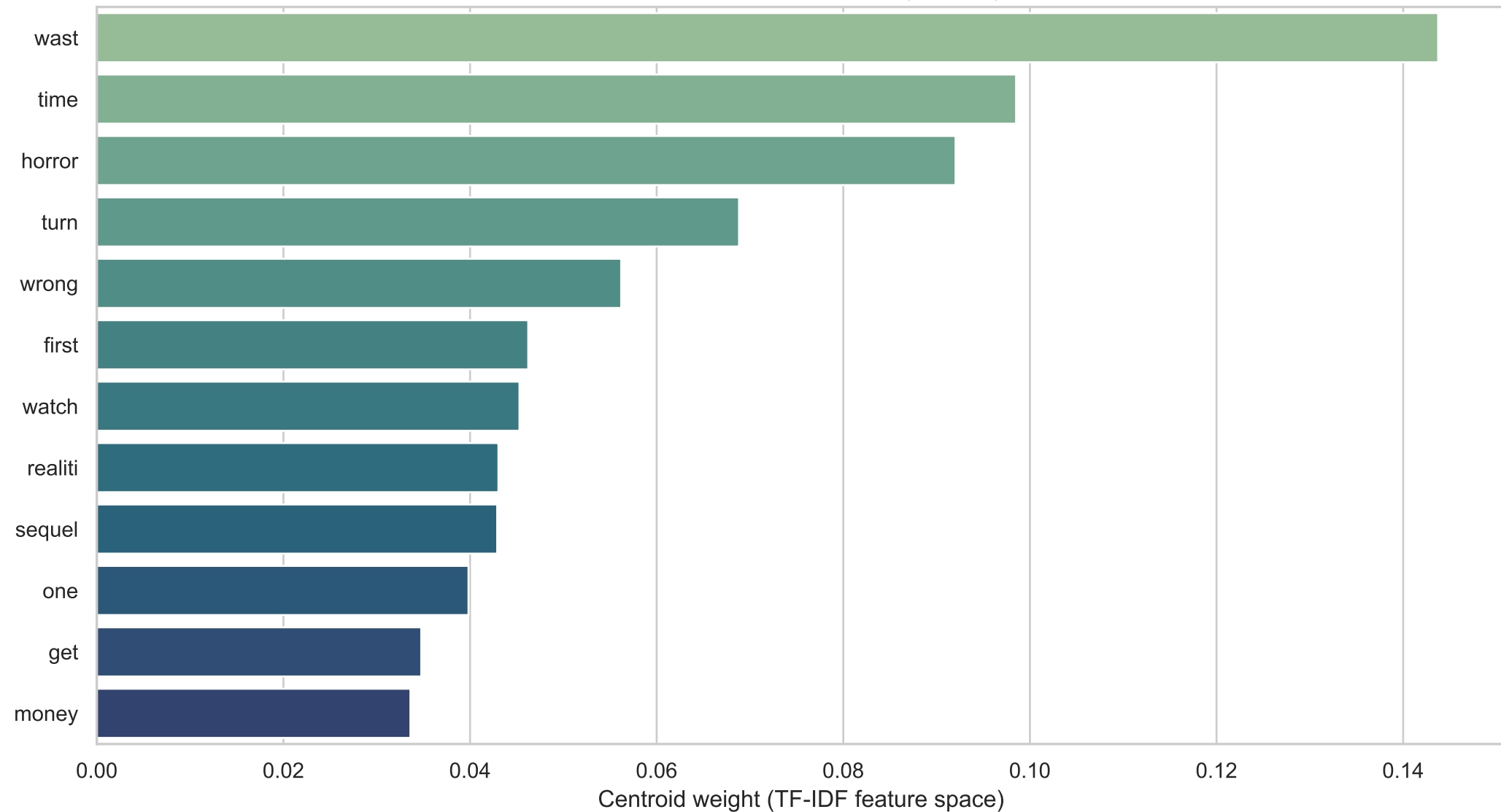Cluster 6: 153 pts

Cluster 0 - top TF-IDF words (centroid)

Cluster 1 - top TF-IDF words (centroid)

Cluster 2 - top TF-IDF words (centroid)

Centroid weight (TF-IDF feature space)

# Cluster 3 - top TF-IDF words (centroid)



| Word | Centroid weight (TF-IDF feature space) |
|---|---|
| book | 0.40 |
| read | 0.10 |
| name | 0.055 |
| one | 0.05 |
| stori | 0.043 |
| titl | 0.042 |
| like | 0.04 |
| watch | 0.038 |
| without | 0.037 |
| john | 0.036 |
| charact | 0.035 |
| would | 0.035 |

Cluster 4 - top TF-IDF words (centroid)

# Cluster 5 - top TF-IDF words (centroid)



| | |
|---|---|
| bad | |
| act | |
| watch | |
| good | |
| movi | |
| terribl | |
| one | |
| time | |
| review | |
| script | |
| plot | |
| get | |

Centroid weight (TF-IDF feature space)

0.00    0.05    0.10    0.15    0.20    0.25

Cluster 6 - top TF-IDF words (centroid)

Naive Bayes - Confusion Matrix

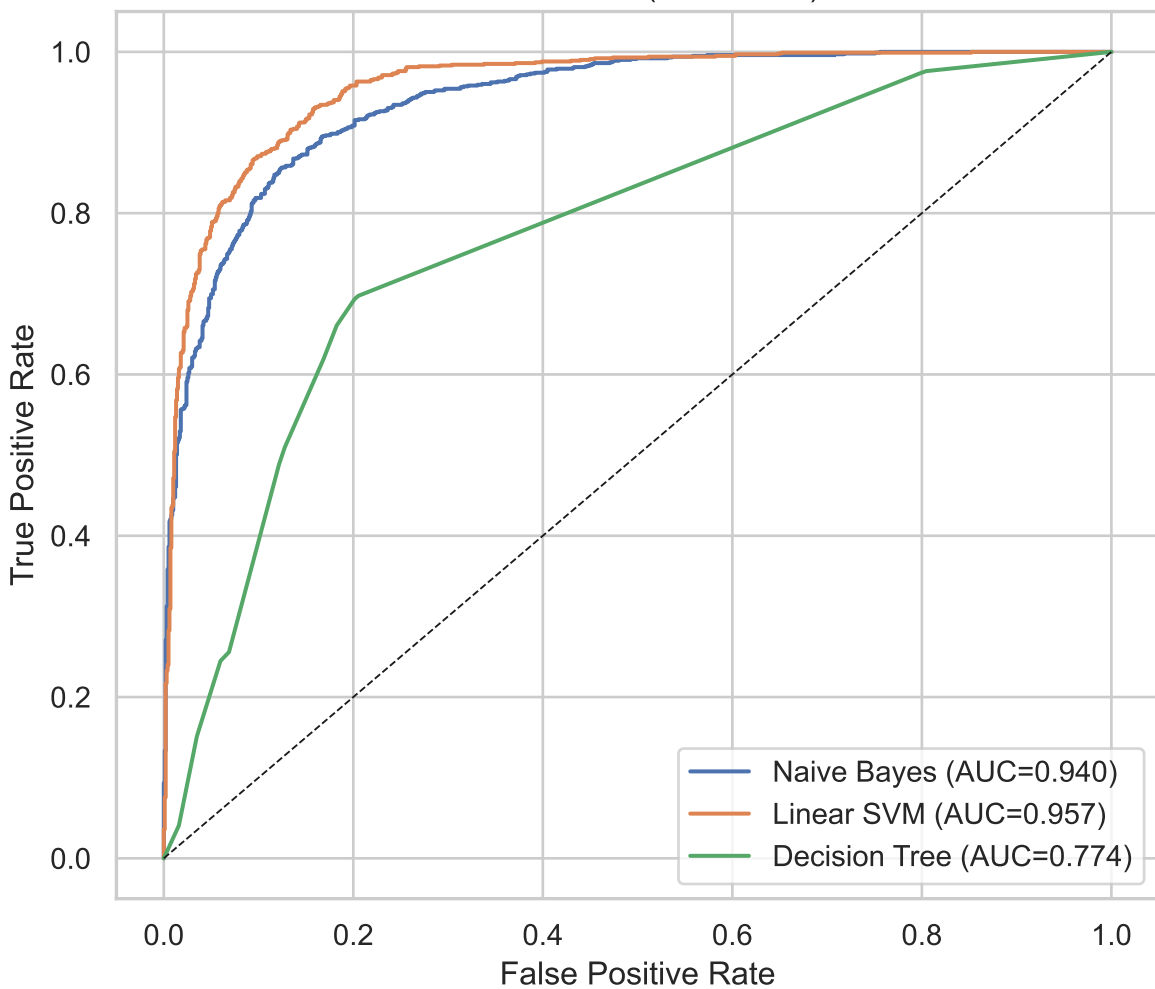|                | Predicted 0 | Predicted 1 |
|----------------|-------------|-------------|
| **True 0**     | 735         | 268         |
| **True 1**     | 56          | 949         |

Linear SVM - Confusion Matrix

Decision Tree - Confusion Matrix

ROC Curves (JSON eval)

Naive Bayes (AUC=0.940)
Linear SVM (AUC=0.957)
Decision Tree (AUC=0.774)
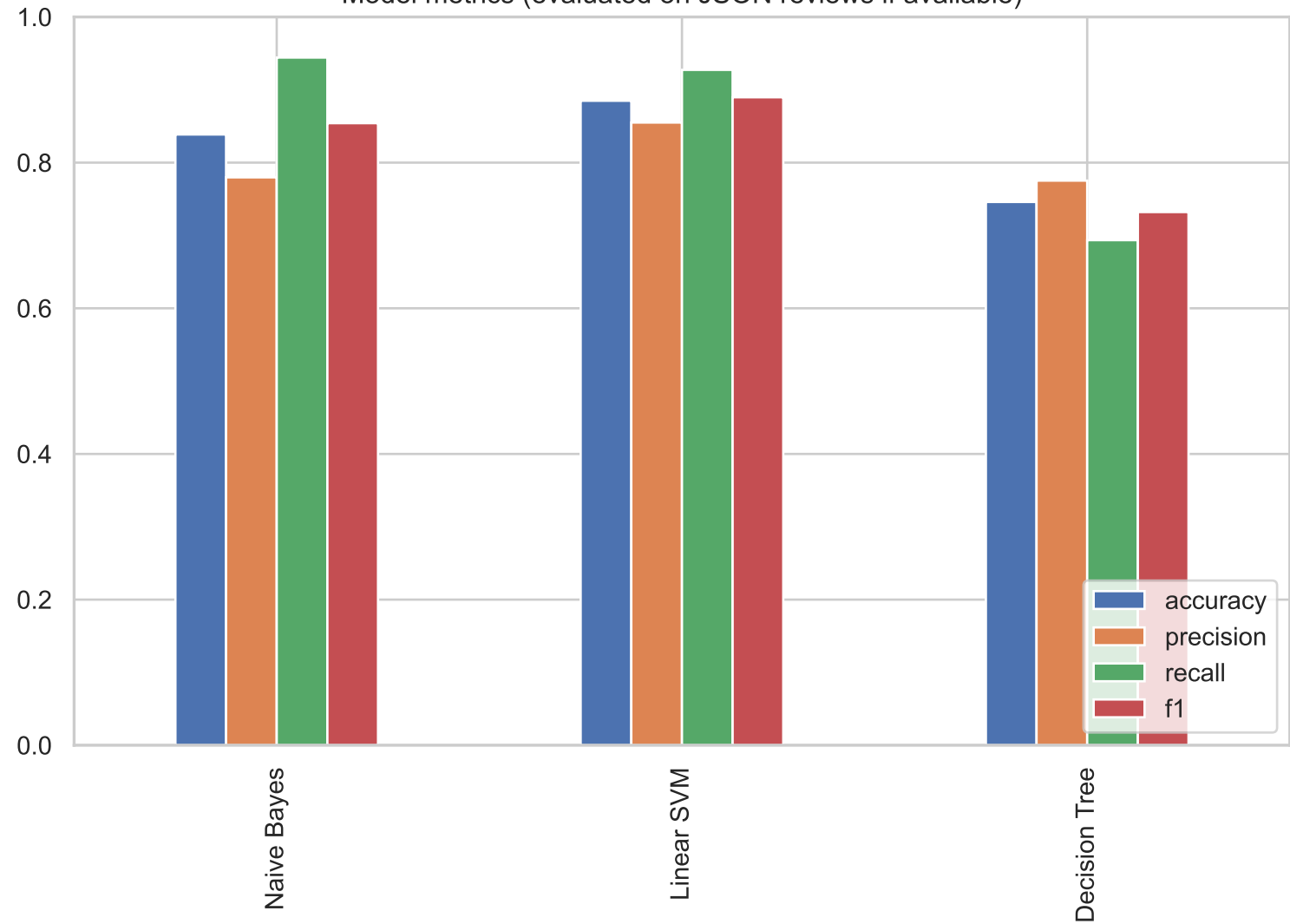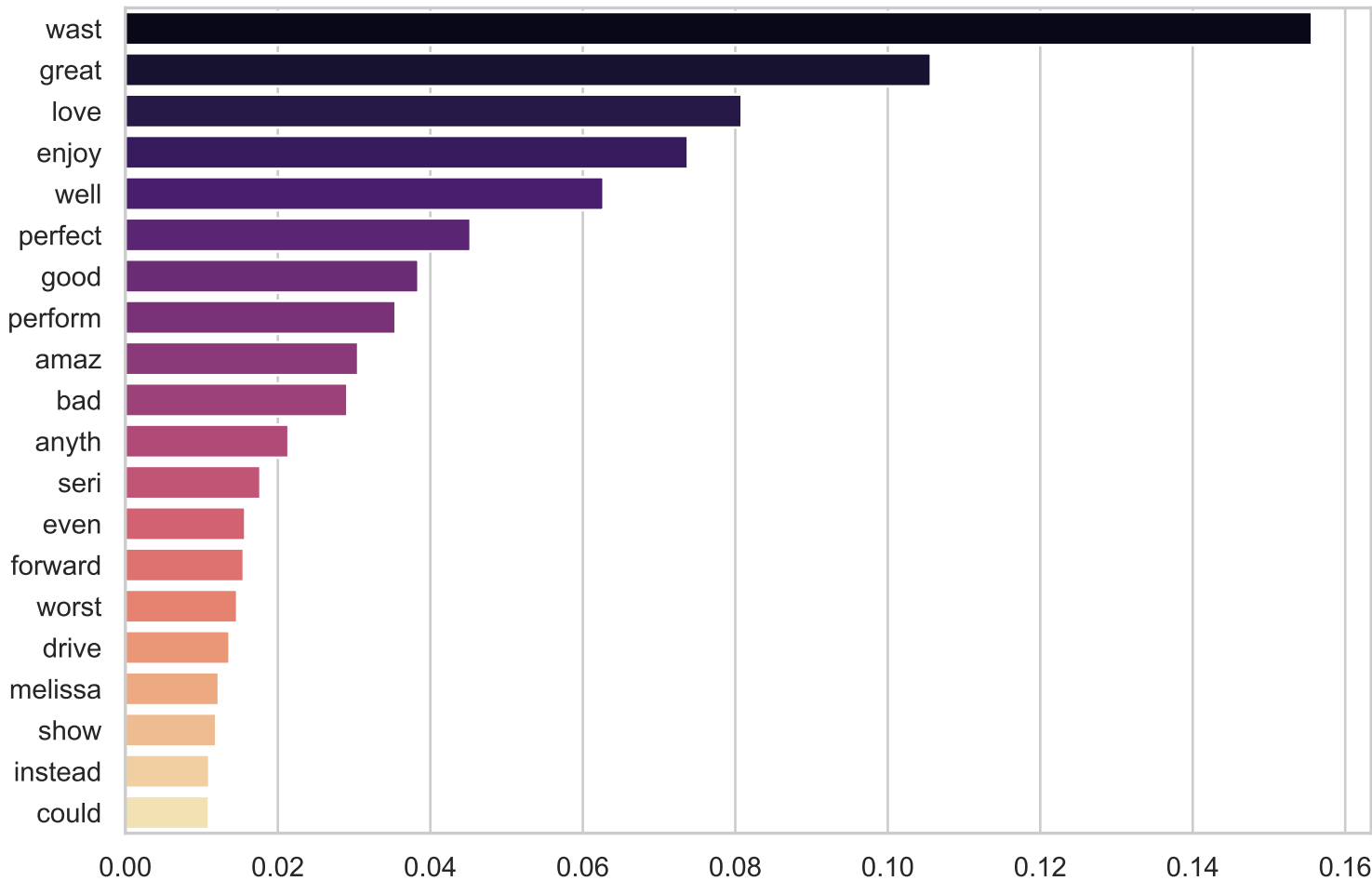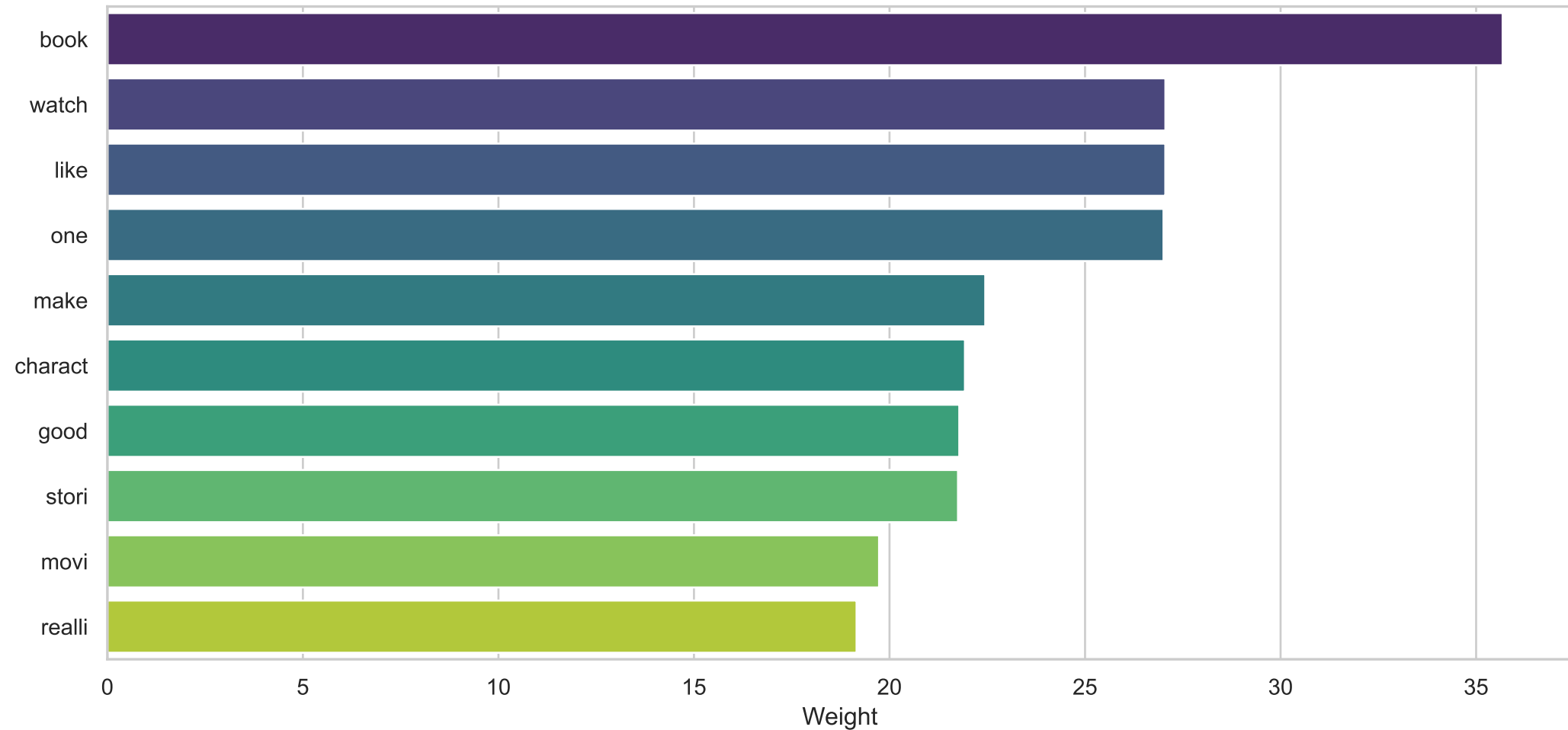
Precision-Recall Curves (JSON eval)

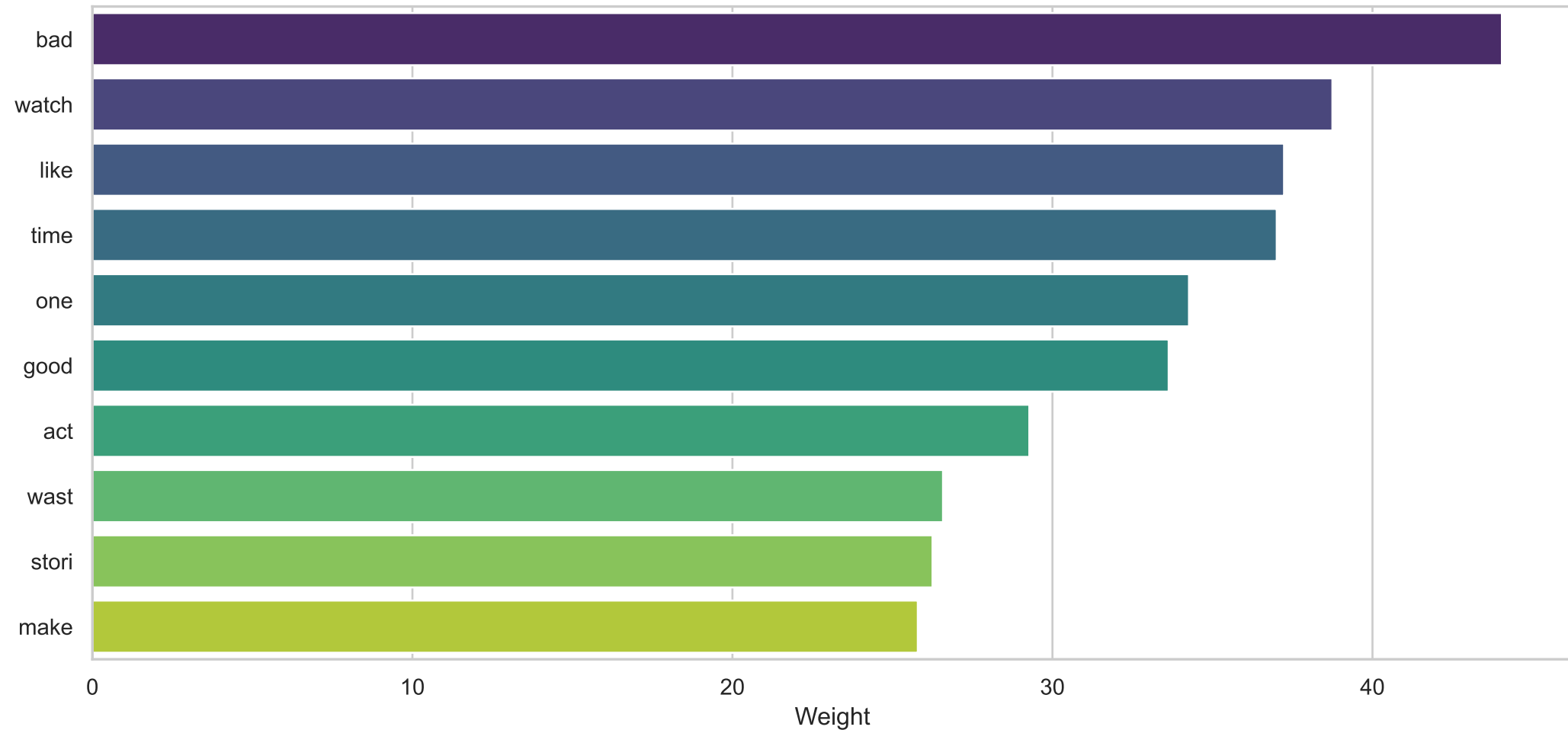Model metrics (evaluated on JSON reviews if available)

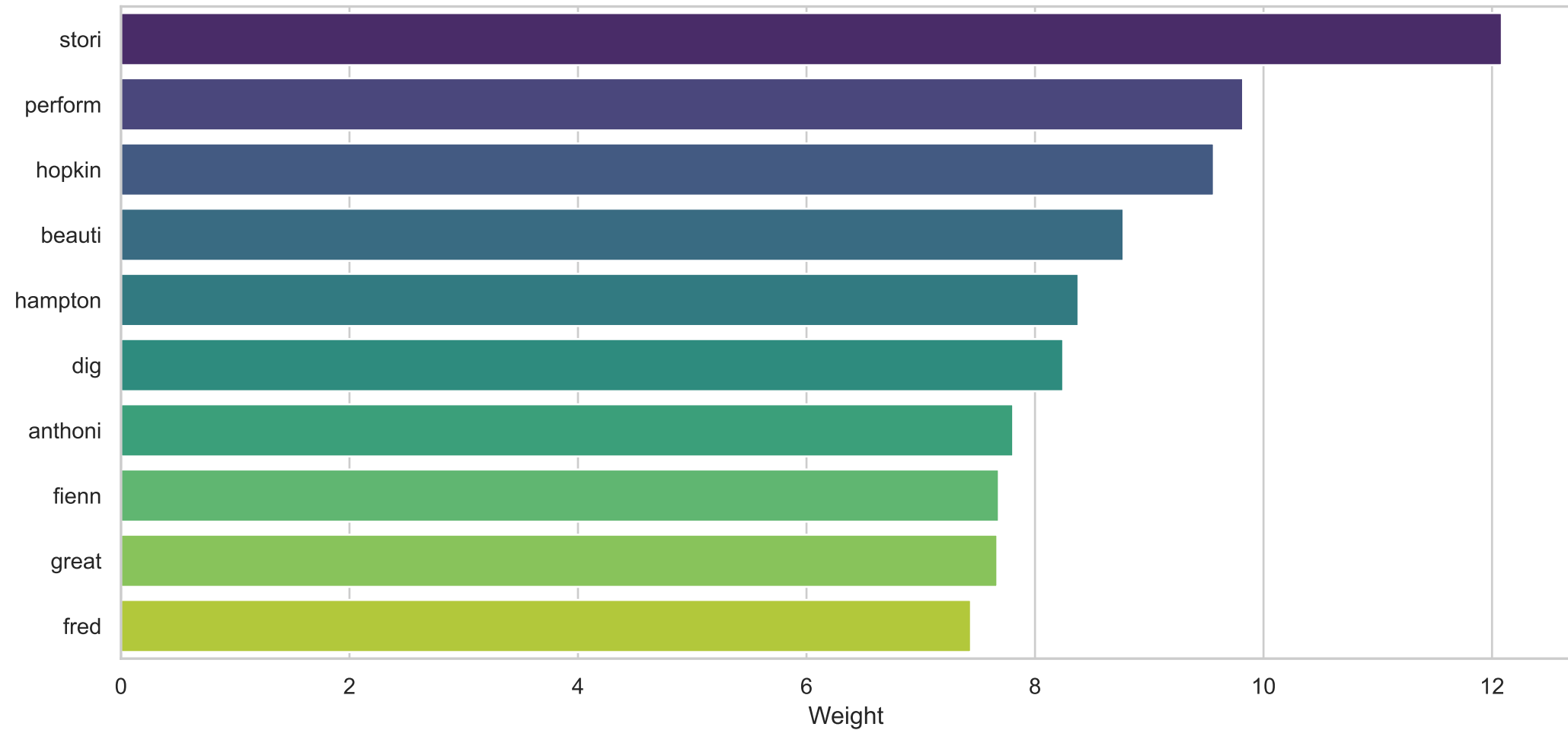Decision Tree - top feature importances (training)
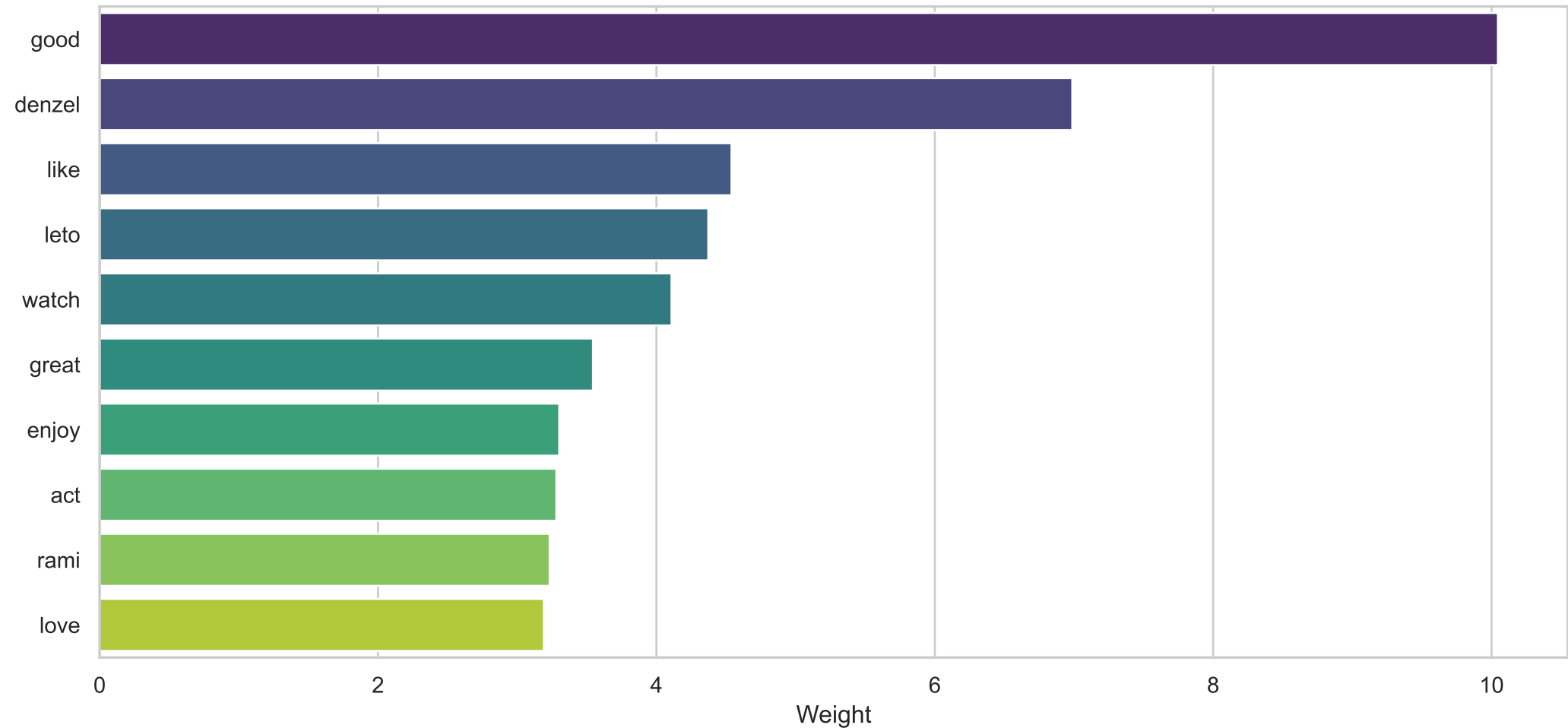
LDA Topic 0 - top words (LDA)
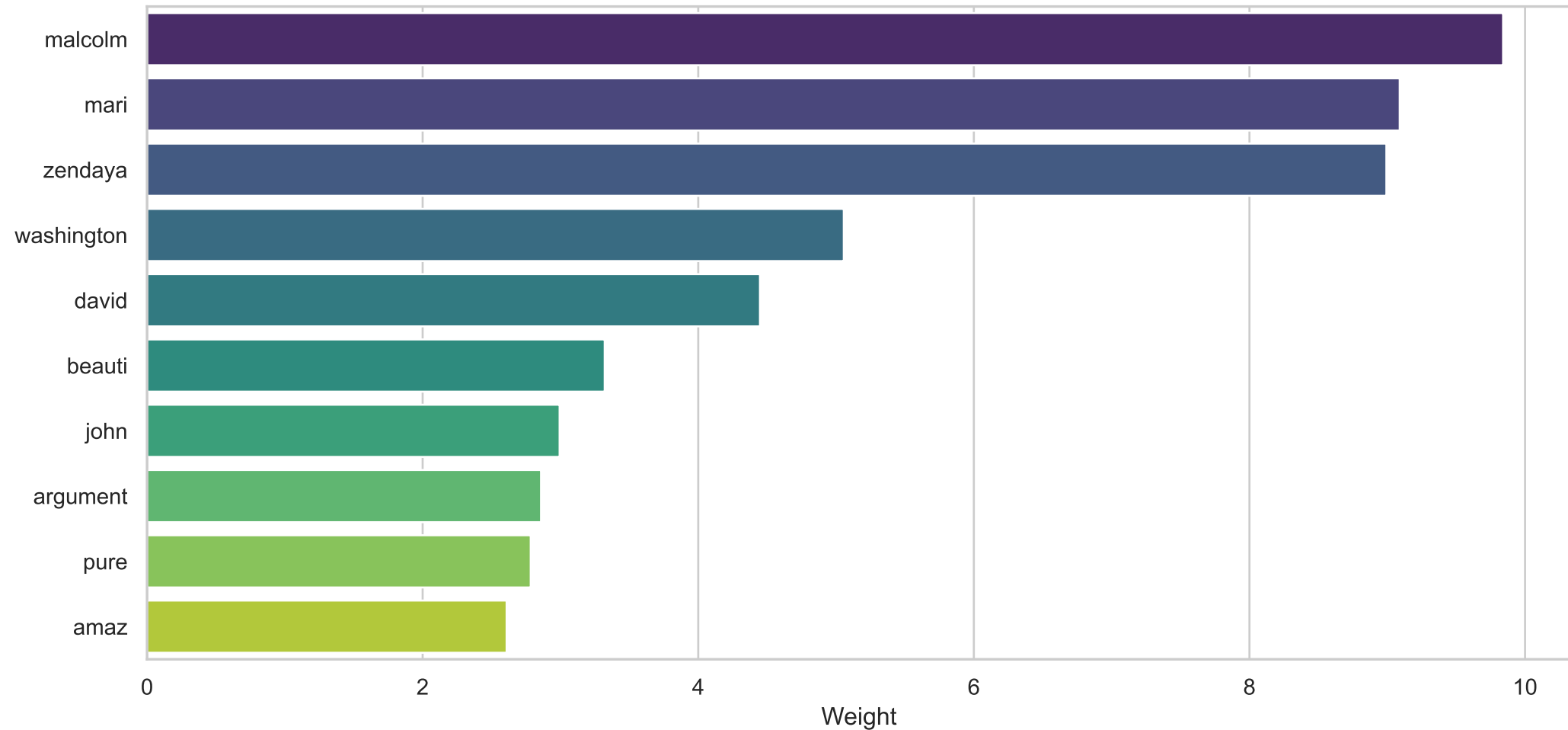
LDA Topic 1 - top words (LDA)

LDA Topic 2 - top words (LDA)

LDA Topic 3 - top words (LDA)

LDA Topic 4 - top words (LDA)

```
Classification reports (evaluated on JSON reviews if available)

Model: Naive Bayes
              precision    recall  f1-score   support

         NEG       0.93      0.73      0.82      1003
         POS       0.78      0.94      0.85      1005

    accuracy                           0.84      2008
   macro avg       0.85      0.84      0.84      2008
weighted avg       0.85      0.84      0.84      2008


Model: Linear SVM
              precision    recall  f1-score   support

         NEG       0.92      0.84      0.88      1003
         POS       0.86      0.93      0.89      1005

    accuracy                           0.88      2008
   macro avg       0.89      0.88      0.88      2008
weighted avg       0.89      0.88      0.88      2008


Model: Decision Tree
              precision    recall  f1-score   support

         NEG       0.72      0.80      0.76      1003
         POS       0.78      0.69      0.73      1005

    accuracy                           0.75      2008
   macro avg       0.75      0.75      0.75      2008
weighted avg       0.75      0.75      0.75      2008
```