

Кластеризация и визуализация каталога Spotify (200 000 треков) с целью получения алгоритма для построения плейлистов

Паремузов Матвей, 332 группа

Цели

1. Разведочный анализ данных
2. Построить непрерывную 2-D проекцию треков
3. Сгруппировать треки методом HDBSCAN
4. Интерпретировать и проверить практическую полезность

1. **Локальные соседи.** Для каждой точки x_i ищем k ближайших соседей $\mathcal{N}_k(i)$ по метрике $d(\cdot, \cdot)$ (часто *cosine* для аудио-признаков).

2. **Граф в \mathbb{R}^D .** Из расстояний делаем вероятности соседства

$$p_{ij} = \exp\left(-\frac{\max(0, d_{ij}-\rho_i)}{\sigma_i}\right), \quad j \in \mathcal{N}_k(i).$$

Параметры ρ_i (минимальный радиус) и σ_i подбираются так, чтобы $\sum_j p_{ij} = \log_2 k$ — тем самым у всех точек одинаковая «энтропия» окружения.

Симметризация (t-норма Лукаса) даёт

$$w_{ij} = p_{ij} + p_{ji} - p_{ij}p_{ji} \in [0, 1].$$

3. **Проекция в \mathbb{R}^d и оптимизация.** В низком измерении задаём

$$q_{ij} = \frac{1}{1 + a\|y_i - y_j\|^{2b}} \quad (a, b — константы из `min_dist`),$$

и минимизируем бинарную кросс-энтропию

$$\mathcal{L} = - \sum_{i < j} [w_{ij} \log q_{ij} + (1 - w_{ij}) \log(1 - q_{ij})].$$

Оптимизация — стохастический градиентный спуск с отрицательной выборкой (*negative sampling*).

2. Ключевые гиперпараметры

<code>n_neighbors</code>	k — баланс локальная \leftrightarrow глобальная структура.
<code>min_dist</code>	минимальная дистанция в проекции: $\downarrow \rightarrow$ кластеры плотнее.
<code>metric</code>	метрика в \mathbb{R}^D (<i>euclidean, cosine, ...</i>).

3. Почему UMAP подходит для музыкальных данных

- Сохраняет и форму локальных кластеров (жанры/муды), и относительное расположение «островов» — t-SNE этим жертвует.
- Линейная сложность по памяти при $k \ll N \rightarrow$ обрабатывает *200k треков* на ноутбуке.
- Параметр `min_dist` позволяет «сжимать» плотные зоны, подчёркивая жанровые и эмоциональные кластеры.

Silhouette Score: что измеряет и как считать

1. Формула для одной точки x_i

$$a(i) = \frac{1}{|C_i| - 1} \sum_{x_j \in C_i, j \neq i} d(x_i, x_j) \quad (\text{средняя дистанция до своего кластера})$$

$$b(i) = \min_{C \neq C_i} \frac{1}{|C|} \sum_{x_j \in C} d(x_i, x_j) \quad (\text{лучшая чужая «соседняя» кластер-дистанция})$$

$$s(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}} \in [-1, 1]$$

- $s(i) \approx 1$ — точка далеко от чужих, близко к своим (отлично).
- $s(i) \approx 0$ — граница между кластерами.
- $s(i) < 0$ — точка ближе к чужому кластеру, чем к своему (плохое назначение).

2. Глобальный показатель

$$\text{Silhouette}(X, \mathcal{C}) = \frac{1}{N} \sum_{i=1}^N s(i)$$

Числится по всем N точкам или по не-noise точкам, если алгоритм (HDBSCAN) выбрасывает шум.

2. Глобальный показатель

$$\text{Silhouette}(X, \mathcal{C}) = \frac{1}{N} \sum_{i=1}^N s(i)$$

3. Грубая шкала интерпретации

$s < 0$	кластеры перекрываются — пере-или недокластизация.
$0 \dots 0.20$	слабое разделение, «пятна» без чётких границ.
$0.20 \dots 0.50$	умеренное качество (норма для реальных данных).
> 0.50	высокое качество, редко встречается вне синтетических данных.

Анализ данных:

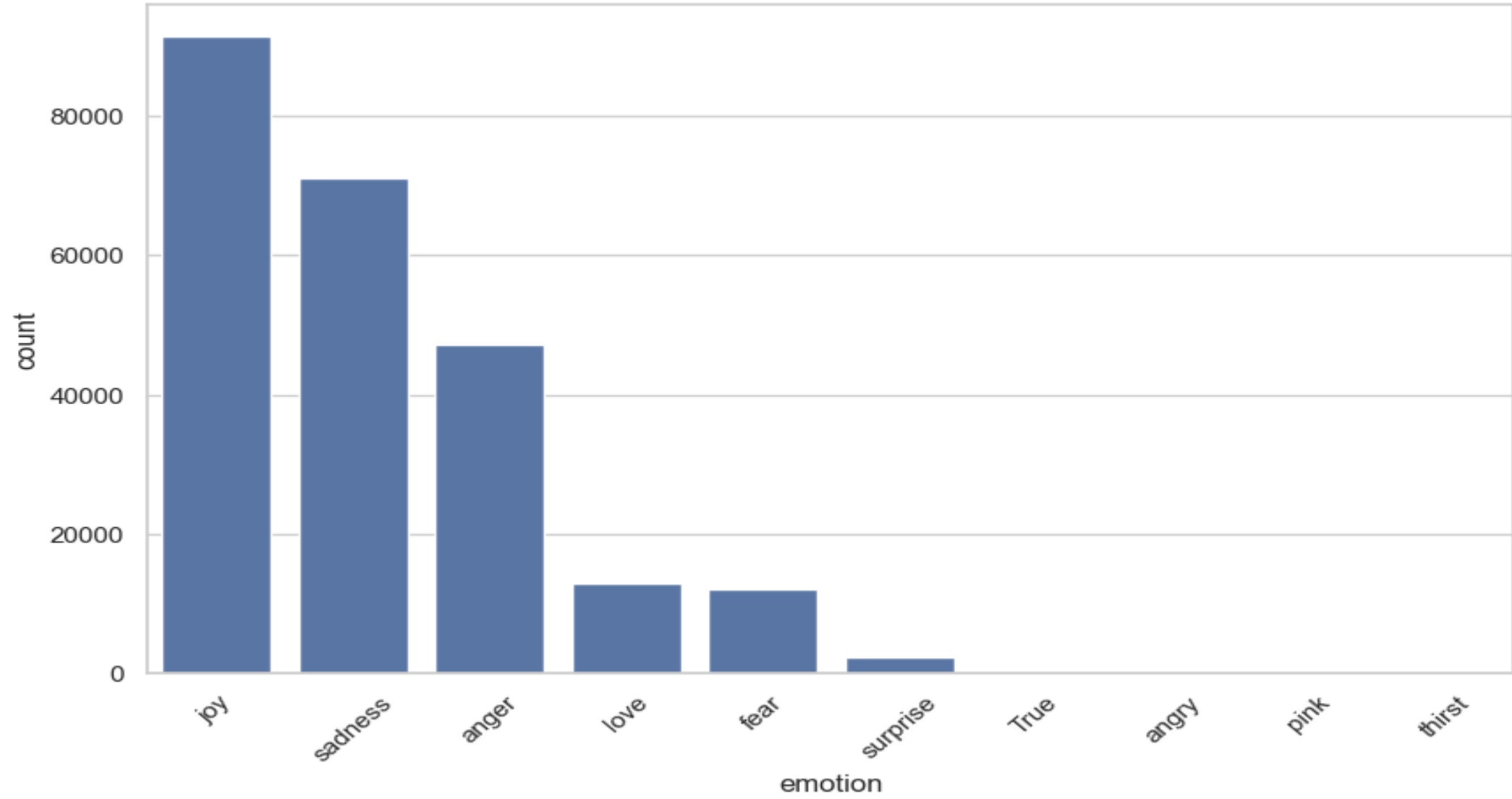
Признаки по трекам:

- Исполнитель
- Название трека
- Мера изменчивости звуковых характеристик песни
- Жанр
- Год выпуска
- Тональность
- Темп (в Уд\мин)
- Средняя громкость песни
- Наличие ненормативной лексики

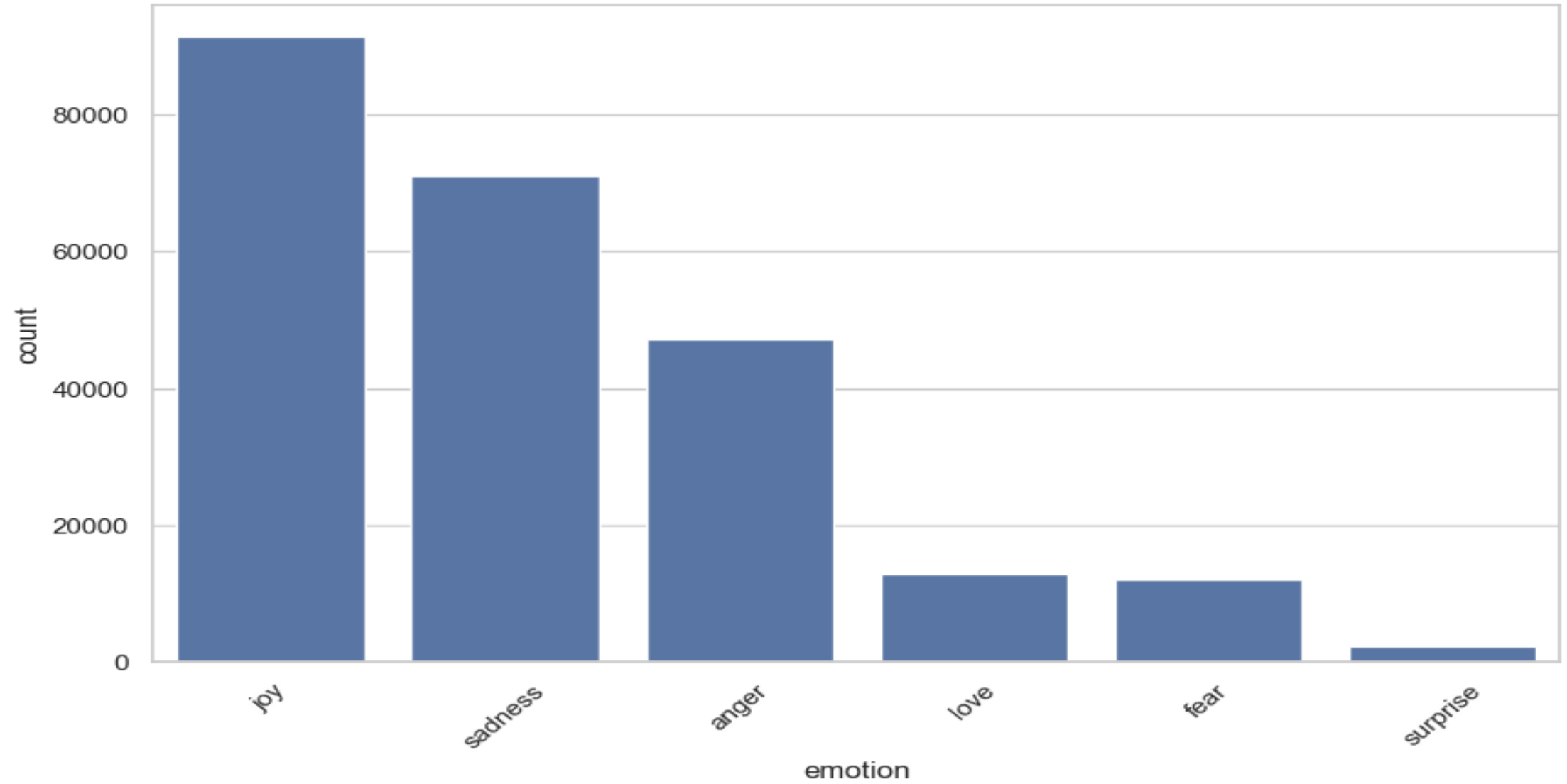
Также в данном датасете присутствуют признаки, полученные с помощью нейросети T5 от компании Google. Вот эти признаки:

- Основная эмоция песни
- Насколько вероятно трек является танцевальным (0-100)
- Энергичность трека
- Позитивность
- Речетативность
- Вероятность записи трека быть живым выступлением
- Акустическое качество
- Вероятность того, что трек инструментальный

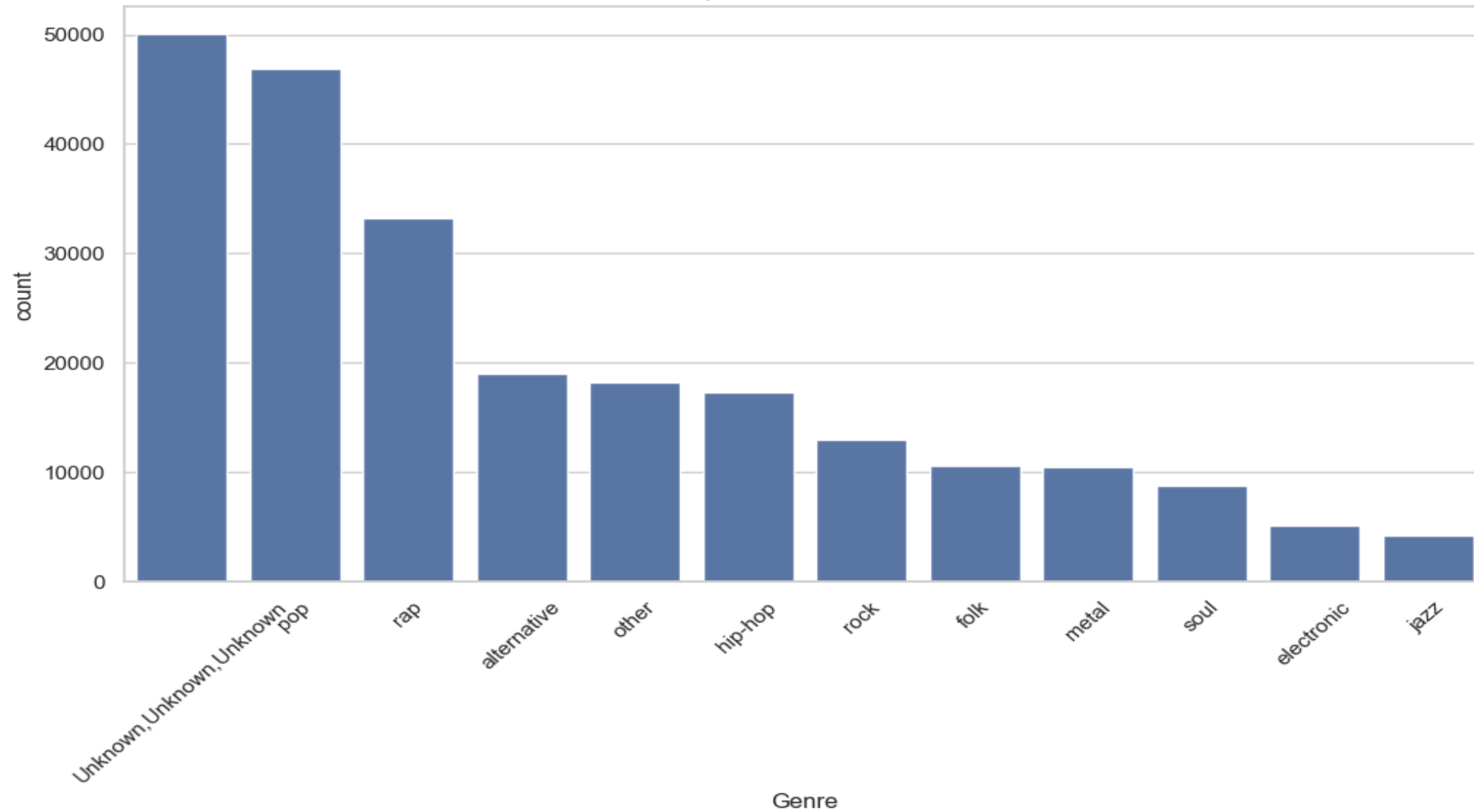
Distribution of Song Emotions



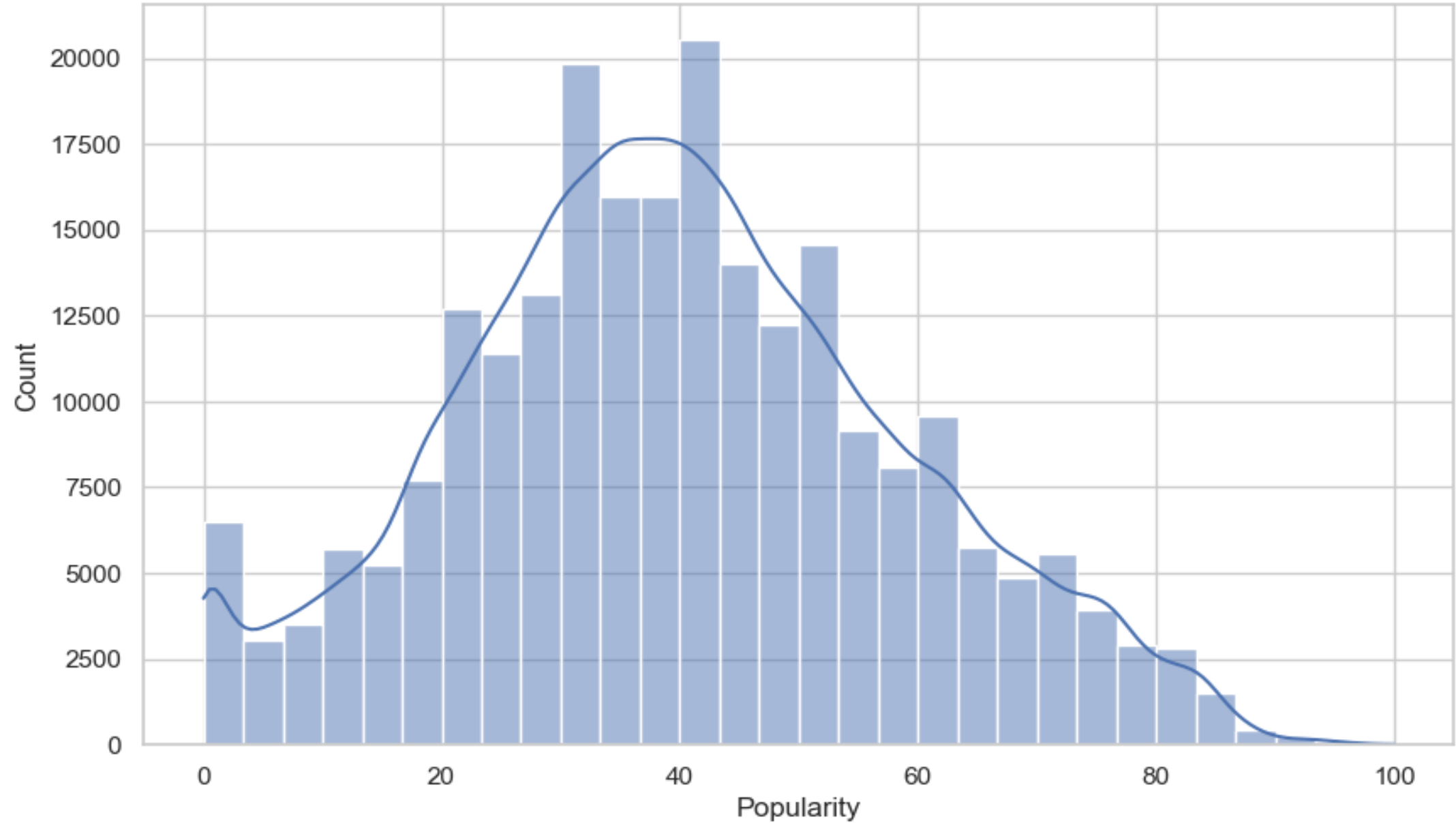
Distribution of Song Emotions

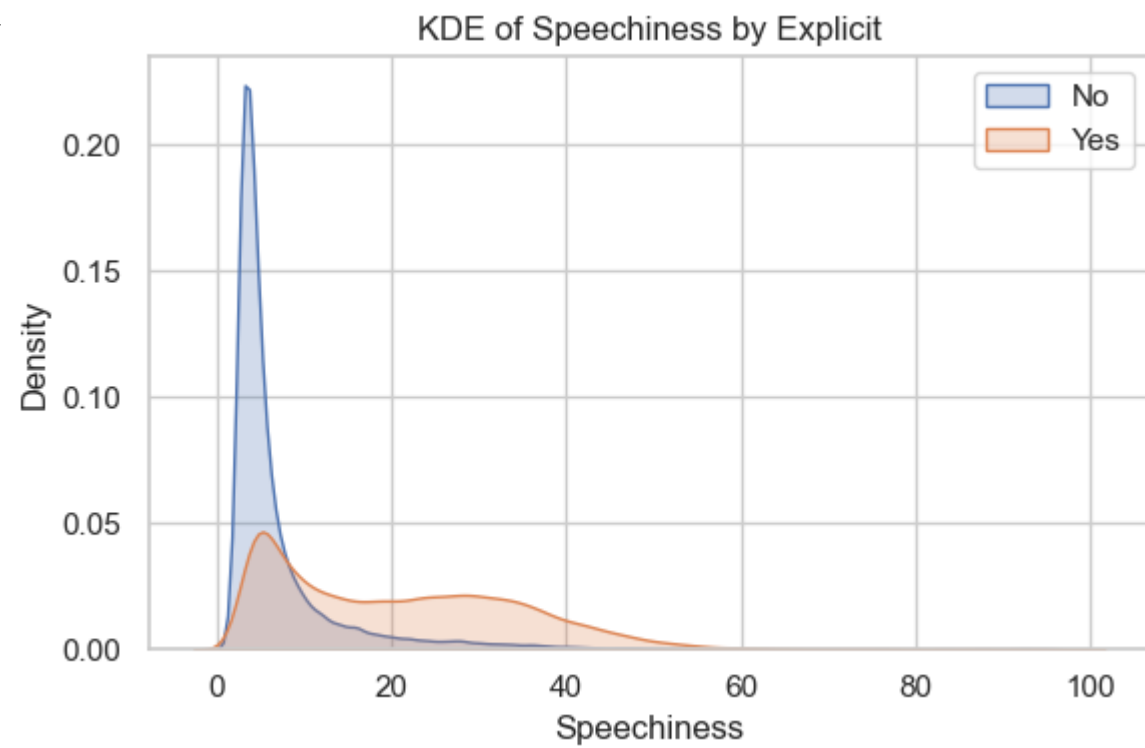
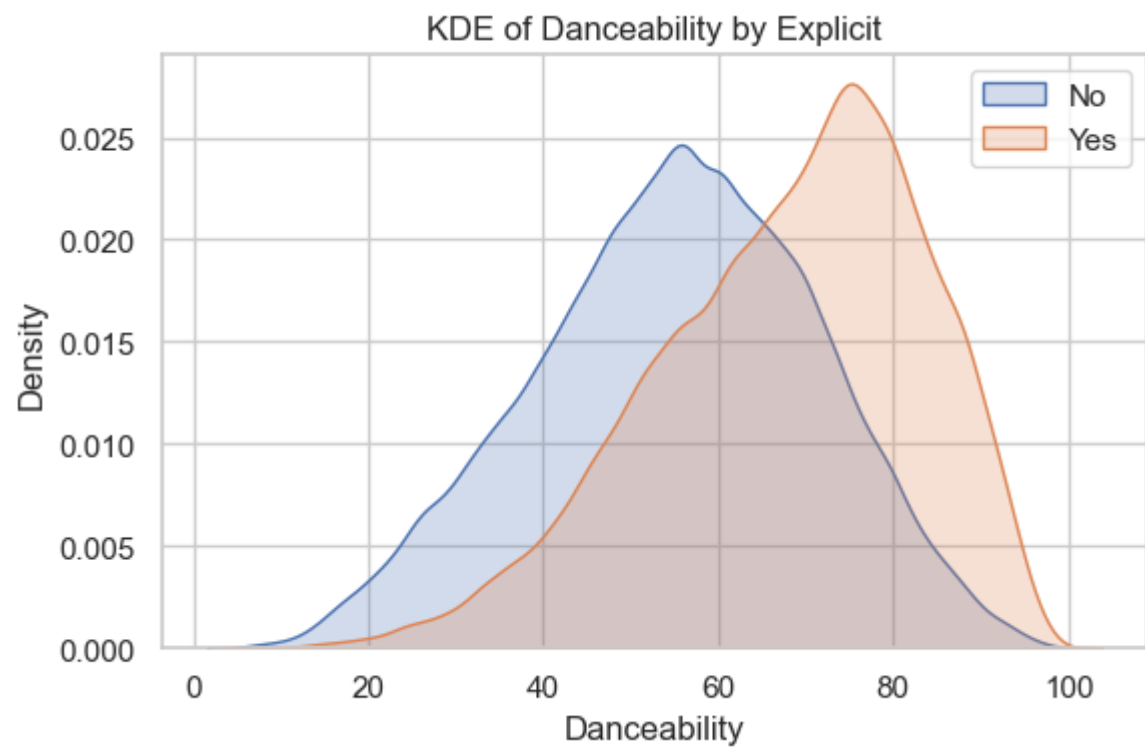


Top 10 Genres + Other



Distributuion of Song Popularity





UMAP projection with HDBSCAN clusters

