

Проект

по статистическому практикуму

Зайченко Николай и Мартюшова Кира, 332 гр.

Мы объединили два набора данных:

titles.csv — содержит основную информацию о фильмах.

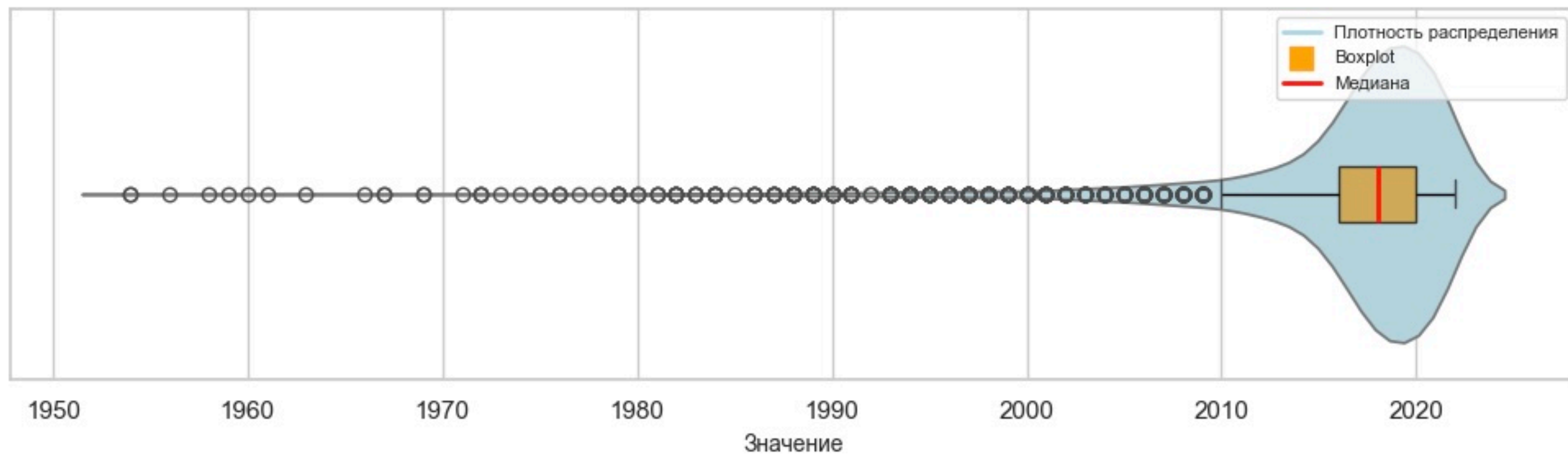
credits.csv — содержит информацию об актёрах и их участии в фильмах.

После объединения получаем следующие признаки:

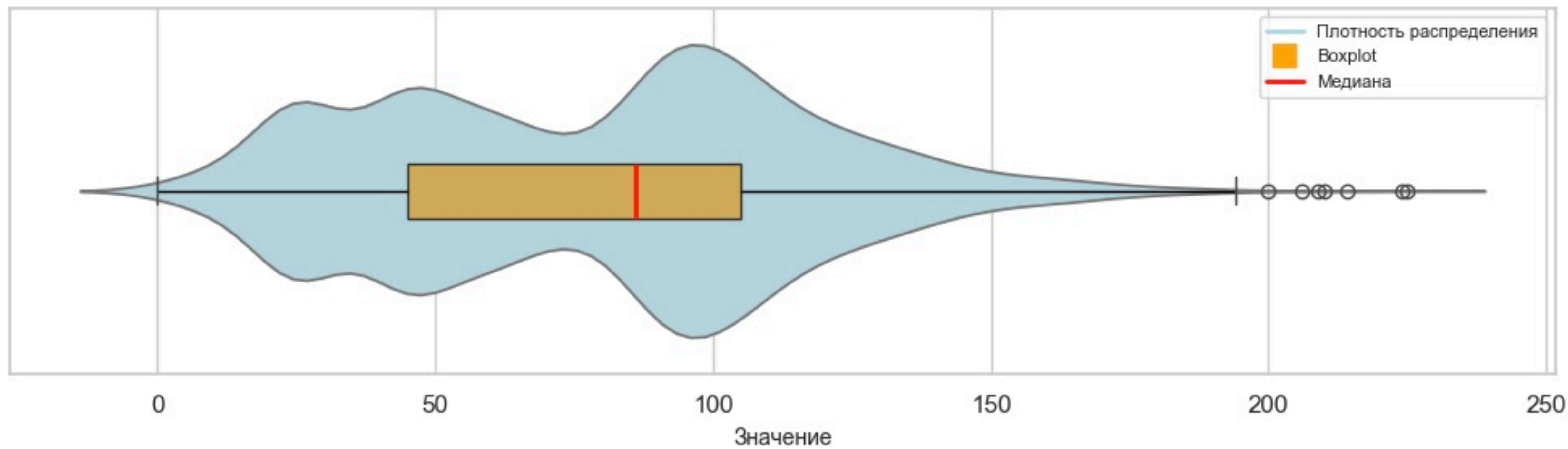
1. title — название фильма
2. type — тип (например, фильм, сериал)
3. release_year — год выпуска
4. age_certification — возрастной рейтинг
5. runtime — продолжительность фильма в минутах
6. genres — жанры (один или несколько)
7. production_countries — страна(-ы) производства
8. imdb_id — идентификатор IMDb
9. imdb_score — оценка IMDb (целевой признак)
10. imdb_votes — количество голосов на IMDb
11. actor_name — имя актёра

Анализ распределения числовых признаков: Violinplot + Boxplot

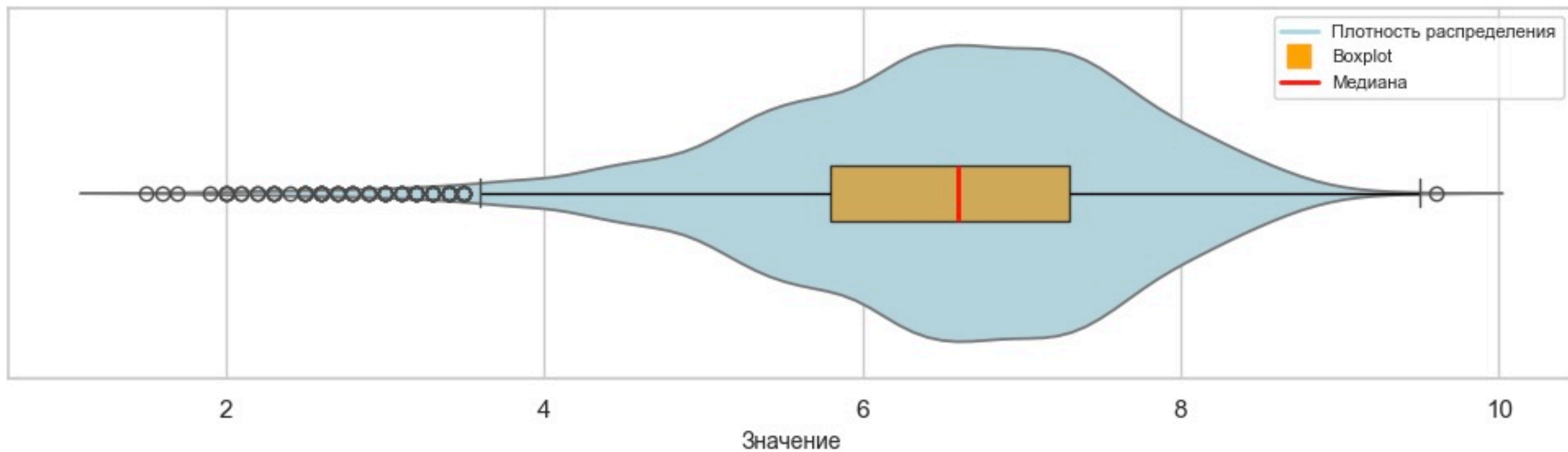
Распределение 'release_year'



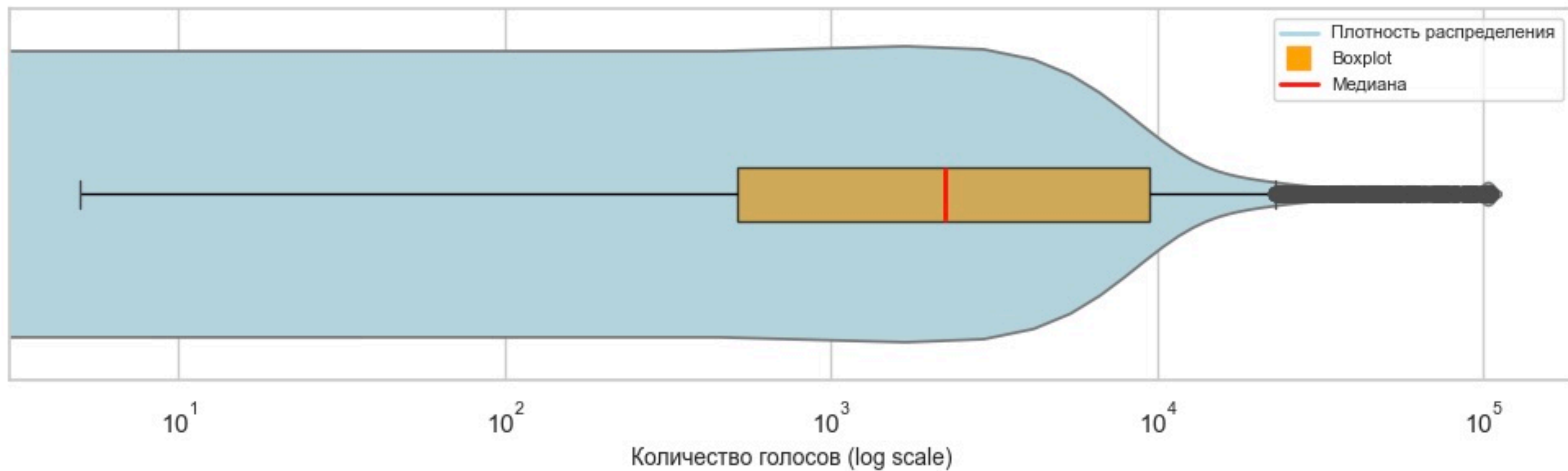
Распределение 'runtime'



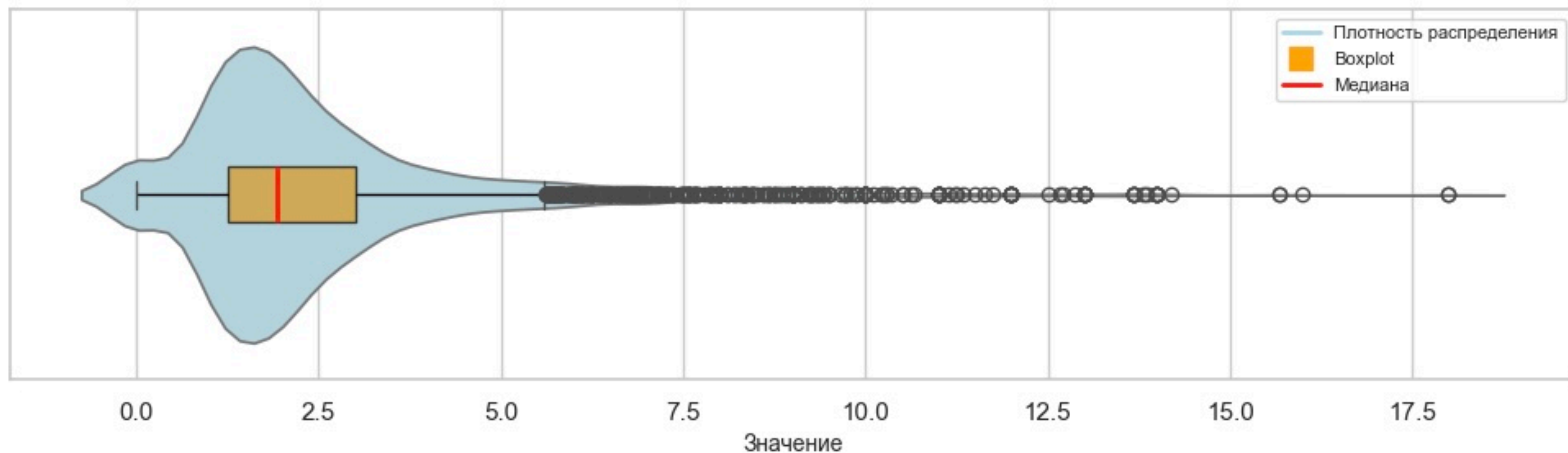
Распределение 'imdb_score'



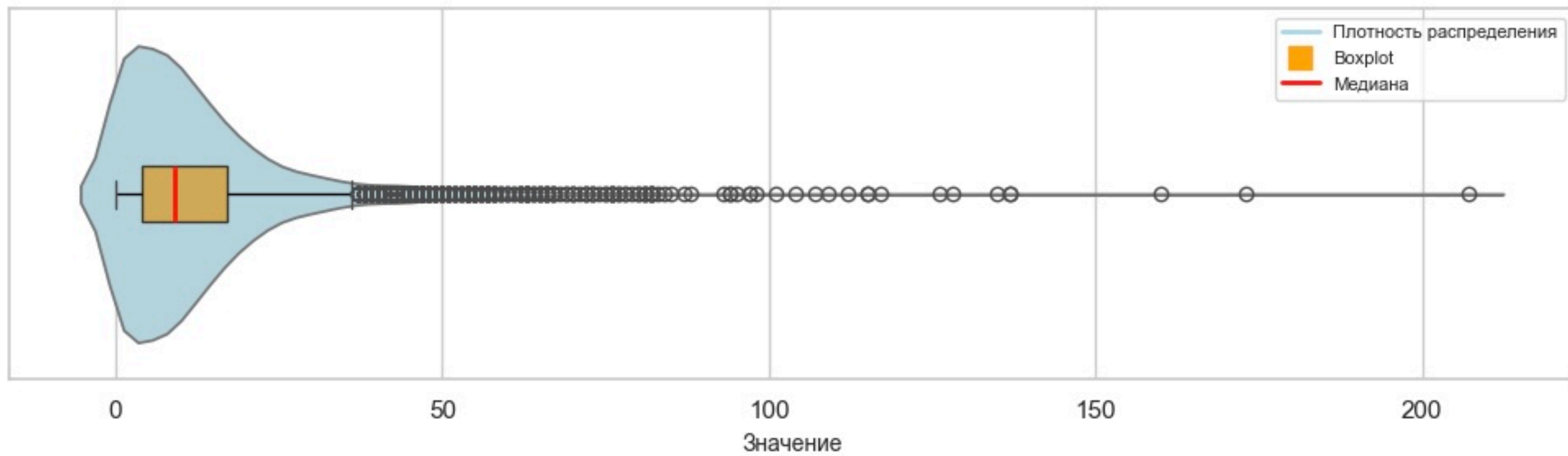
Распределение 'imdb_votes' (min=5)



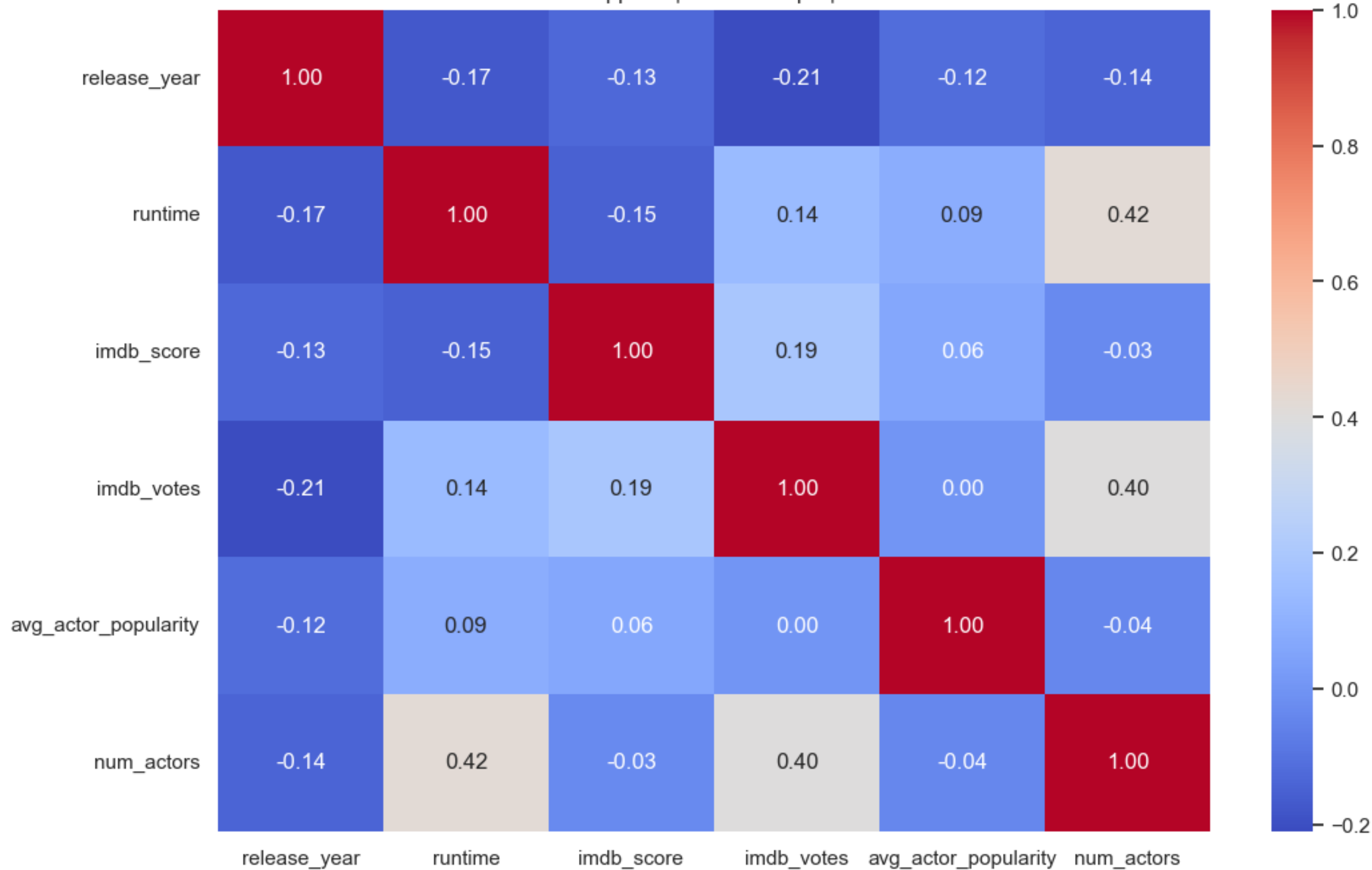
Распределение 'avg_actor_popularity'



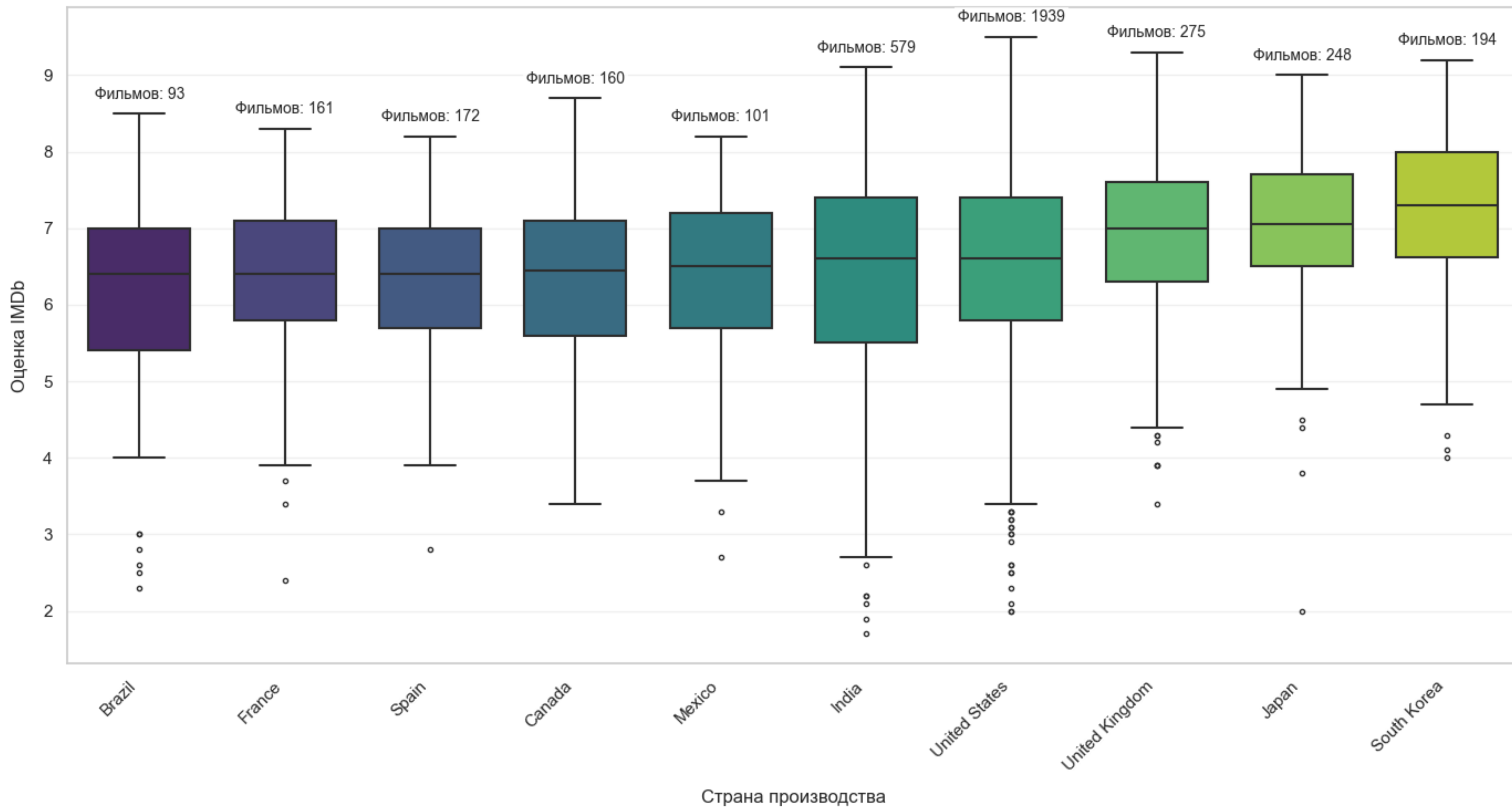
Распределение 'num_actors'



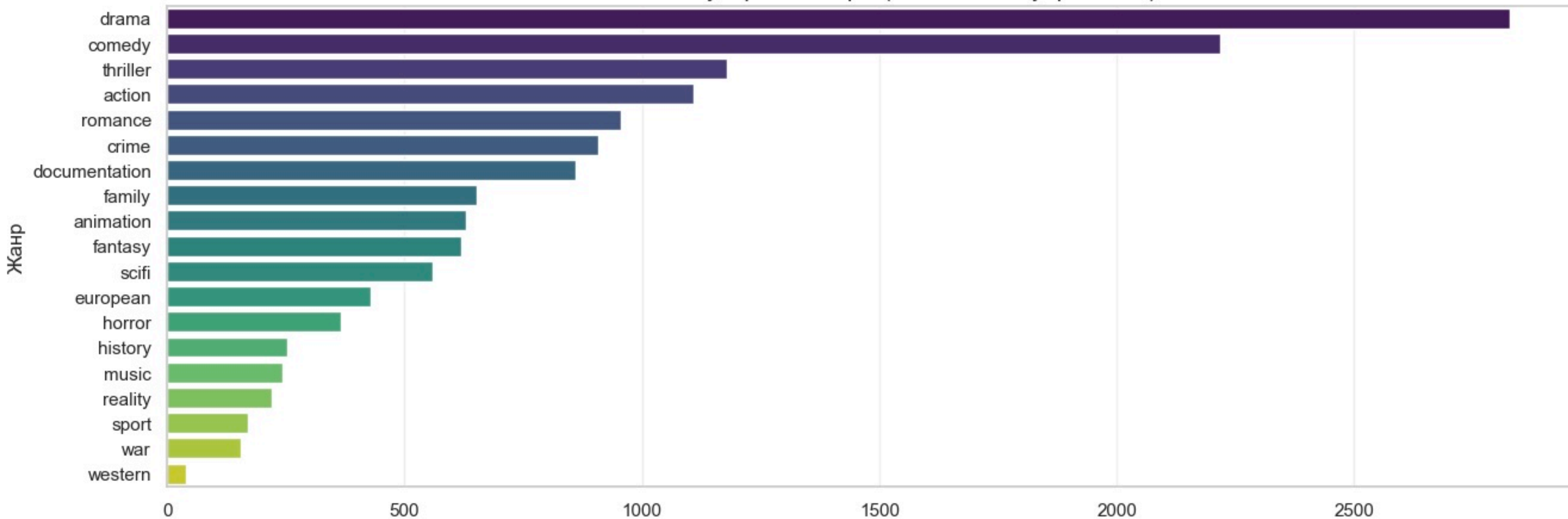
Корреляционная матрица



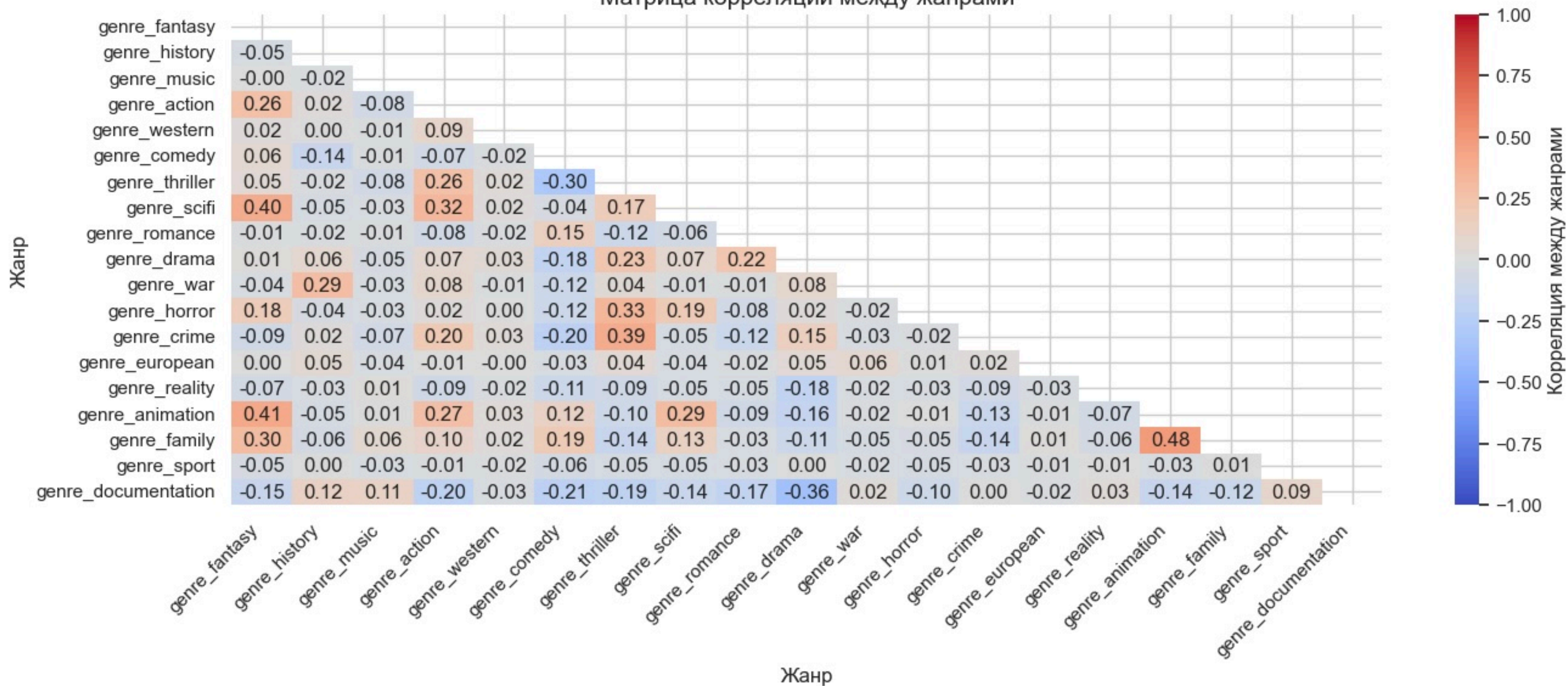
Распределение оценок IMDb по топ-10 странам производства



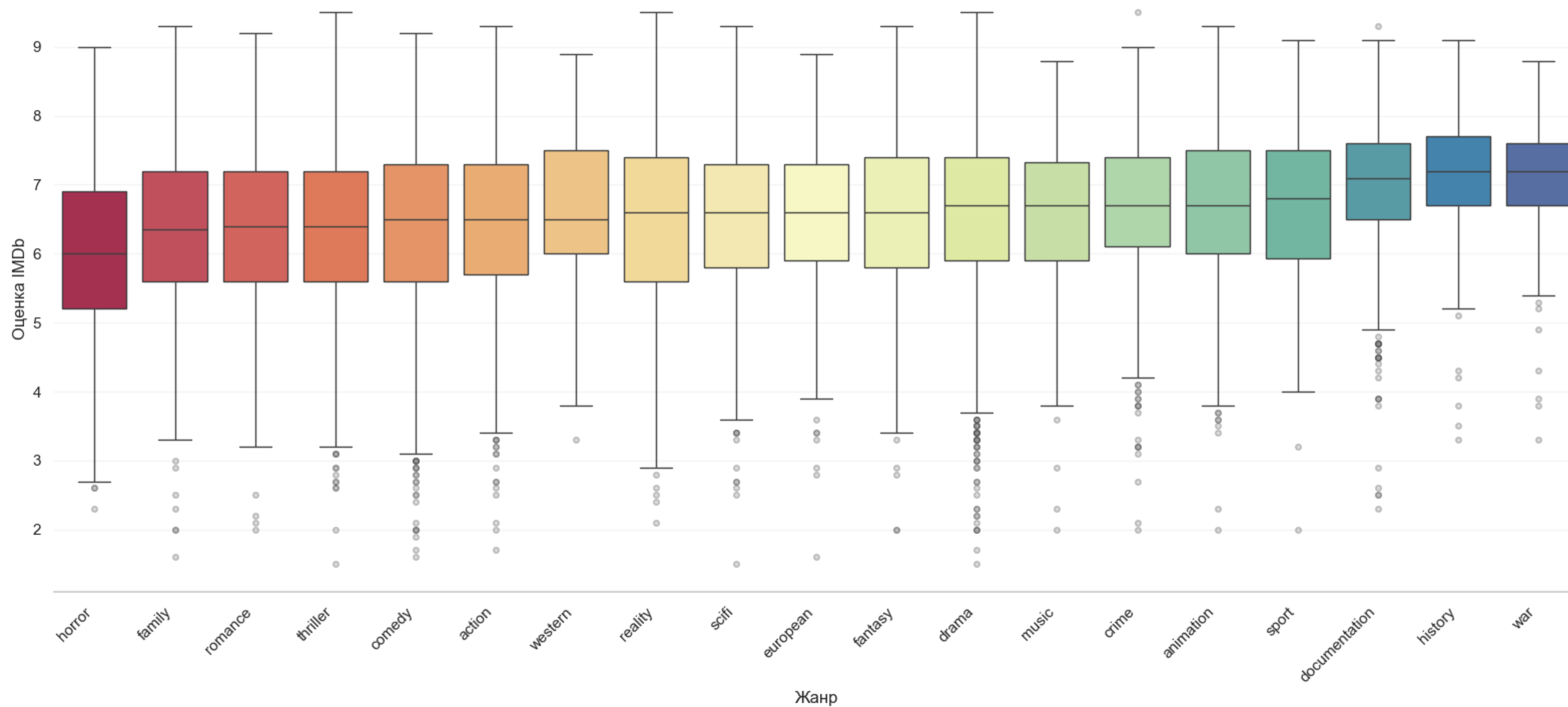
Самые популярные жанры (по количеству фильмов)



Матрица корреляций между жанрами

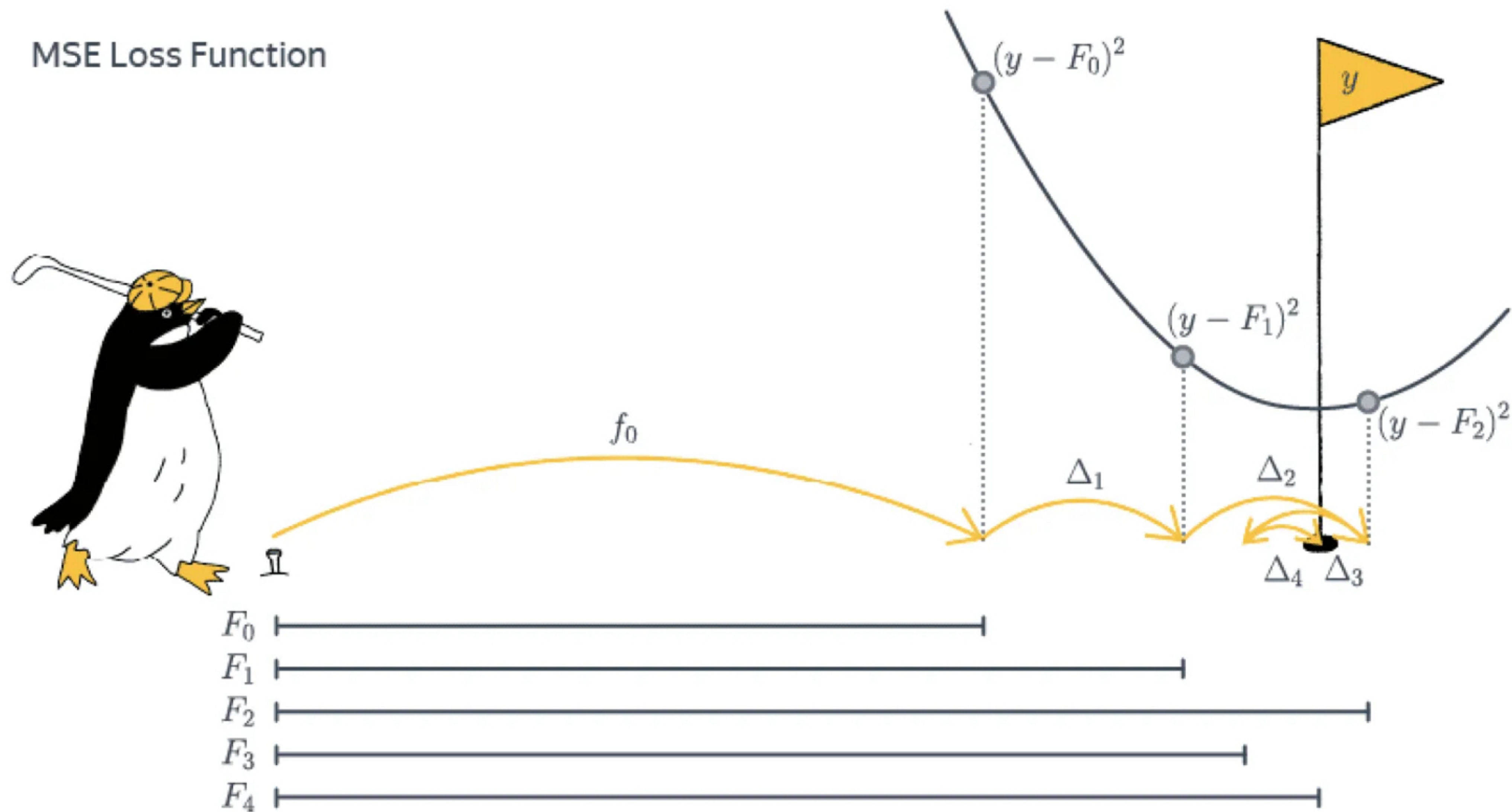


Распределение оценок IMDb по жанрам



Градиентный бустинг

MSE Loss Function



Градиентный бустинг строит взвешенную композицию алгоритмов:

$$a(x) = \sum_{i=1}^T \alpha_i b_i(x),$$

где $b_i(x)$ — базовые алгоритмы, а α_i — их веса.

$$a(x) = \sum_{i=1}^T \alpha_i b_i(x)$$

$$Q(\alpha, b) = \sum_{i=1}^l L(a(x_i), y_i) = \sum_{i=1}^l L \left(\underbrace{\sum_{t=1}^{T-1} \alpha_t b_t(x_i)}_{f_{T-1,j}} + \alpha b(x_i), y_i \right) \rightarrow \min_{\alpha, b}$$

$$Q(\alpha, b) = \sum_{i=1}^l L(a(x_i), y_i) = \sum_{i=1}^l L \left(\underbrace{\sum_{t=1}^{T-1} \alpha_t b_t(x_i)}_{f_{T-1,j}} + \alpha b(x_i), y_i \right) \rightarrow \min_{\alpha, b}$$

$$f_{T-1} = [f_{T-1,1}, f_{T-1,2}, \dots, f_{T-1,l}]^T$$

$$f_T = [f_{T,1}, f_{T,2}, \dots, f_{T,l}]^T$$

$$Q(\alpha, b) = \sum_{i=1}^l L(a(x_i), y_i) = \sum_{i=1}^l L \left(\underbrace{\sum_{t=1}^{T-1} \alpha_t b_t(x_i)}_{f_{T-1,j}} + \alpha b(x_i), y_i \right) \rightarrow \min_{\alpha, b}$$

$$f_{T-1} = [f_{T-1,1}, f_{T-1,2}, \dots, f_{T-1,l}]^T$$

$$f_T = [f_{T,1}, f_{T,2}, \dots, f_{T,l}]^T$$

$$g_i = \frac{\partial L(f_{T-1,i}, y_i)}{\partial f_{T-1,i}} = L'(f_{T-1,i}, y_i), \quad i = 1, \dots, l$$

$$Q(\alpha, b) = \sum_{i=1}^l L(a(x_i), y_i) = \sum_{i=1}^l L \left(\underbrace{\sum_{t=1}^{T-1} \alpha_t b_t(x_i)}_{f_{T-1,j}} + \underbrace{\alpha b(x_i)}_{f_{T,j}}, y_i \right) \rightarrow \min_{\alpha, b}$$

$$f_{T-1} = [f_{T-1,1}, f_{T-1,2}, \dots, f_{T-1,l}]^T \qquad f_{T,i} = f_{T-1,i} - \alpha g_i, \quad i = 1, \dots, l$$

$$f_T = [f_{T,1}, f_{T,2}, \dots, f_{T,l}]^T \qquad f_{T,i} = f_{T-1,i} + \alpha b(x_i), \quad i = 1, \dots, l$$

$$g_i = \frac{\partial L(f_{T-1,i}, y_i)}{\partial f_{T-1,i}} = L'(f_{T-1,i}, y_i), \quad i = 1, \dots, l$$

$$b_T = \arg \min_b \sum_{i=1}^l (b(x_i) + g_i)^2$$

$$b_T = \arg \min_b \sum_{i=1}^l (b(x_i) + g_i)^2$$

$$\alpha_T = \arg \min_{\alpha > 0} \sum_{i=1}^l L(f_{T-1,i} + \alpha b_T(x_i), y_i)$$

В качестве входных признаков (**X**) используются числовые характеристики фильмов:

- Длительность фильма (runtime);
- Год выпуска (release_year);
- Число актеров (num_actors)
- Популярность актера (avg_actor_popularity);
- Есть ли в фильме топ-актер (has_top_actor);
- Жанр фильма (genres);
- Страна производства (country).

В качестве целевой переменной (**Y**) берётся рейтинг IMDb (imdb_score), который необходимо предсказать.

Для модели градиентного бустинга наиболее популярной функцией потерь является среднеквадратичная ошибка (MSE):

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2.$$

Для модели градиентного бустинга наиболее популярной функцией потерь является среднеквадратичная ошибка (MSE):

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2.$$

Средняя абсолютная ошибка (MAE):

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|;$$

Коэффициент детерминации (R^2):

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}.$$

Результаты

Модель	MSE	MAE	R^2
Градиентный бустинг	0,9885	0,7575	0,2914

Кластеризация: OPTICS и спектральный подход

OPTICS (Ordering Points To Identify the Clustering Structure)

- Похоже на DBSCAN, но не требует фиксированного радиуса
- Строит упорядоченный список точек по плотности
- Обнаруживает кластеры разной плотности и шум

Параметры OPTICS

- **min_samples:** определяет минимальный размер кластера
- **cluster_method:** метод кластеризации (в нашем случае это метод “xi”)
- **xi:** чувствительность к изменению плотности

Ключевые понятия OPTICS

Core Distance — радиус, включающий min_samples

Reachability Distance(q,p) = $\max(\text{core_dist}(p), \text{dist}(p,q))$

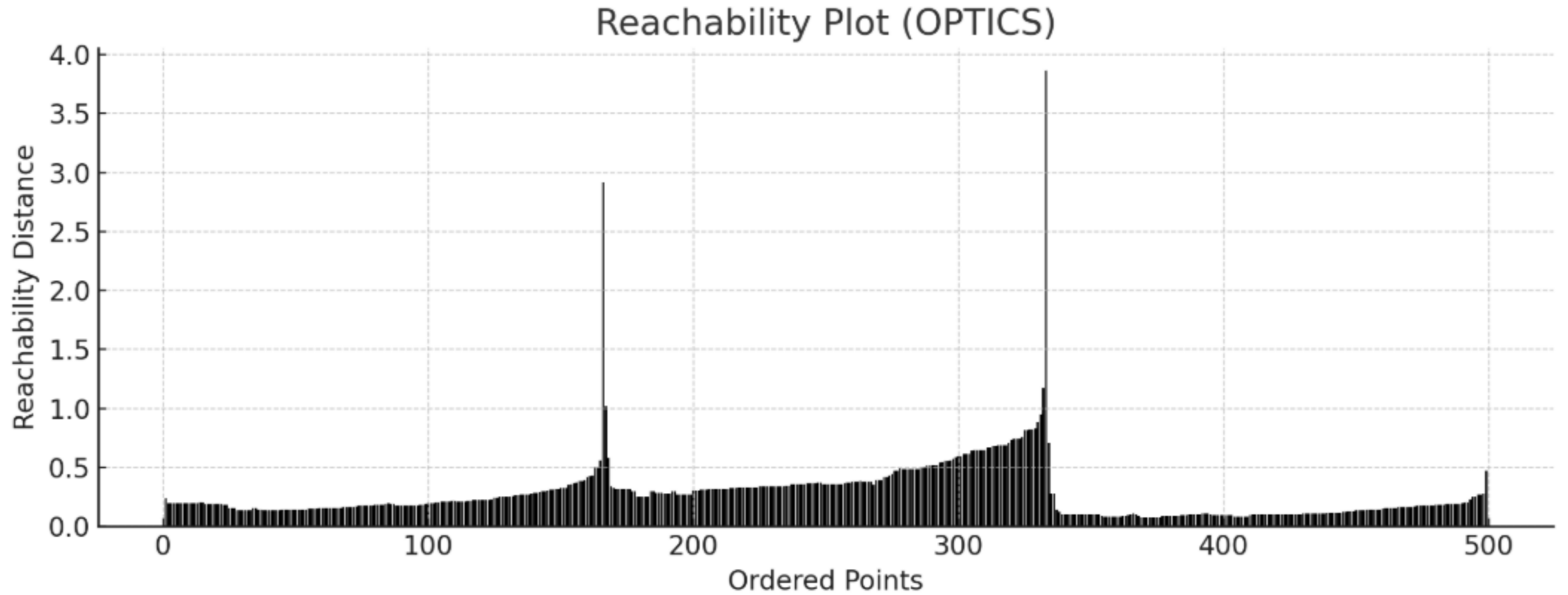
Метод ξ в OPTICS

- Кластеры определяются как области **стабильной плотности**, между которыми резкие изменения.
- Используется относительное изменение расстояния достижимости:

$$\frac{\text{reachability_dist}_{i+1} - \text{reachability_dist}_i}{\text{reachability_dist}_i} > \xi$$

- Параметр $\xi \in (0, 1)$ — чувствительность к границам кластеров (обычно 0.03–0.05).

Пример графика достижимости



Спектральная кластеризация

1. Строим граф похожести
2. Вычисляем лапласиан графа
3. Находим собственные вектора
4. Применяем k-means

Матрица смежности и лапласиан

- Матрица A : $A_{ij} = \begin{cases} 1, & \text{если } j \in kNN(i) \\ 0, & \text{иначе} \end{cases}$
- Матрица степеней: $D_{ii} = \sum_j A_{ij}$
- Нормализованный лапласиан: $L = I - D^{-1/2}AD^{-1/2}$

Кластеризация фильмов с помощью комбинации данных методов:

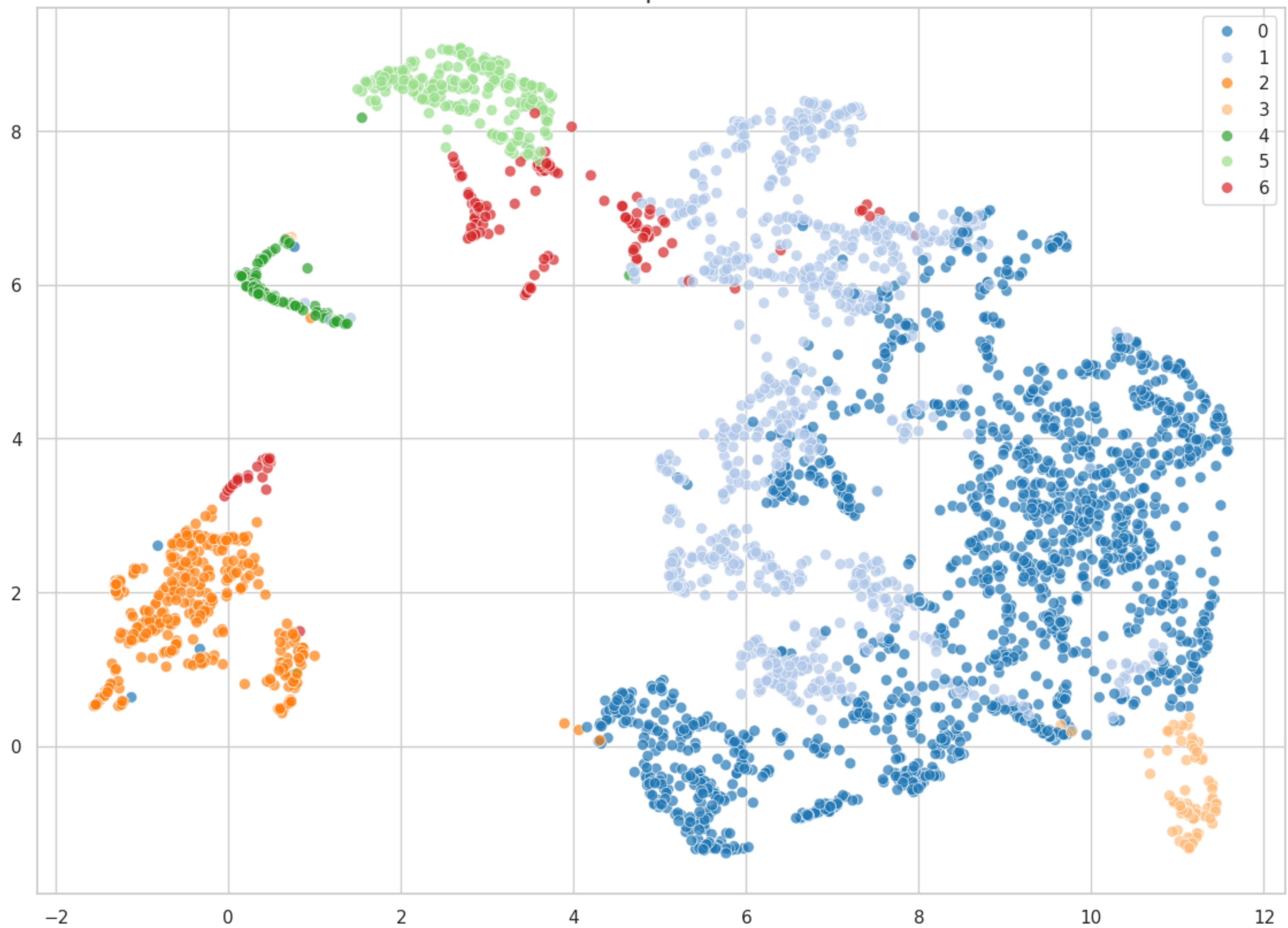
Цель: Мы хотим автоматически разделить фильмы на группы (кластеры), основываясь на количественных характеристиках.

Это позволяет понять, какие типы фильмов существуют в датасете, какие у них общие свойства, и как можно использовать эту информацию, например, для рекомендаций или анализа рынка.

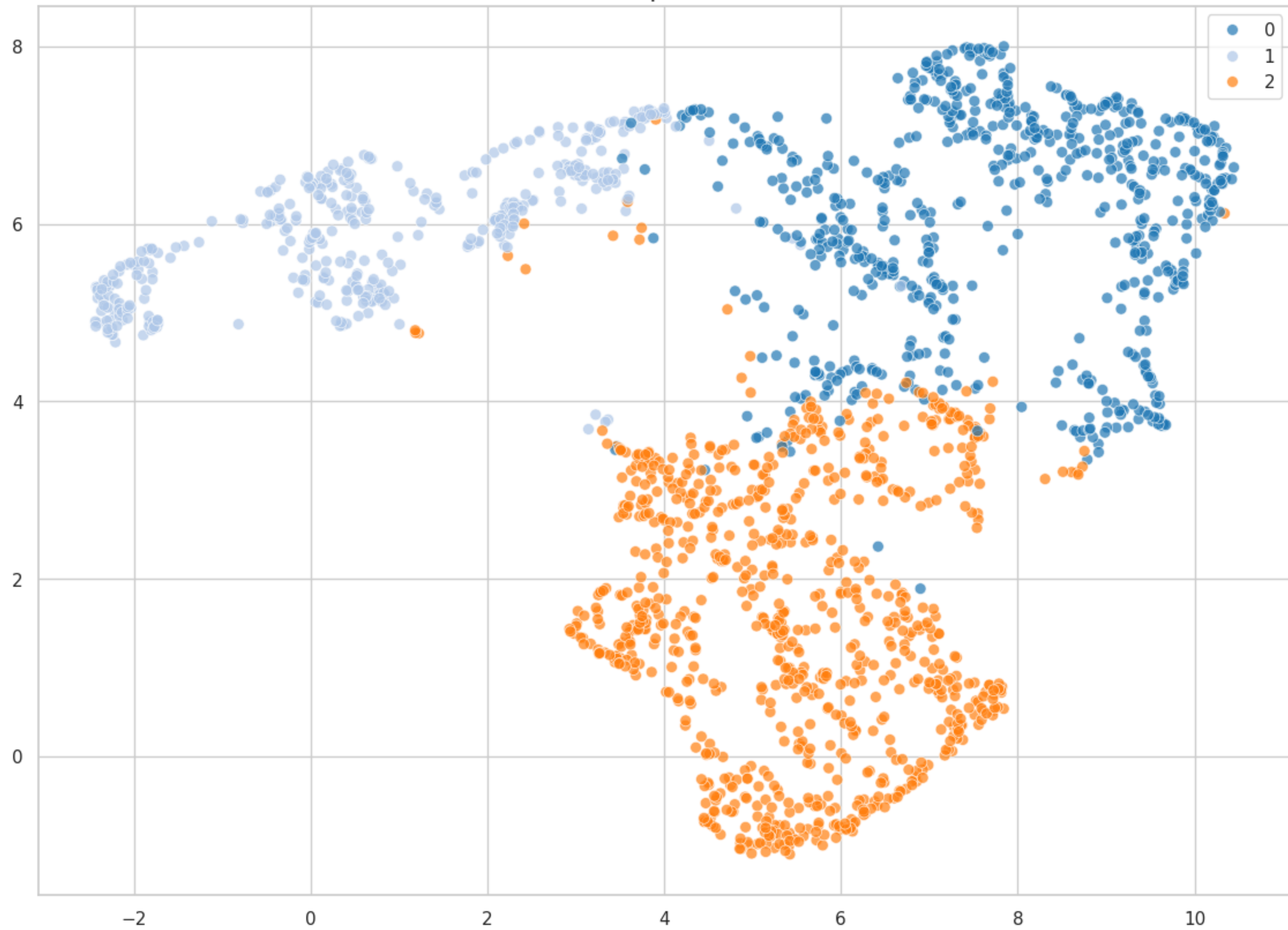
Используемые признаки:

- runtime — длительность фильма;
- release_year — год выхода;
- imdb_score — рейтинг IMDb;
- imdb_votes — количество голосов на IMDb;
- жанры.

Кластеры MOVIE



Кластеры SHOW



Кластеризация фильмов и сериалов

