

Анализ аспектов оценки коротких историй

Данные

Данные представляют собой результаты оценки коротких историй по 9 критериям. Оставим в данных только сами критерии, удалив ненужные на данный момент колонки. Критерии оценивались по разным шкалам, что отображено в саммари.

Шкала 0 - 4

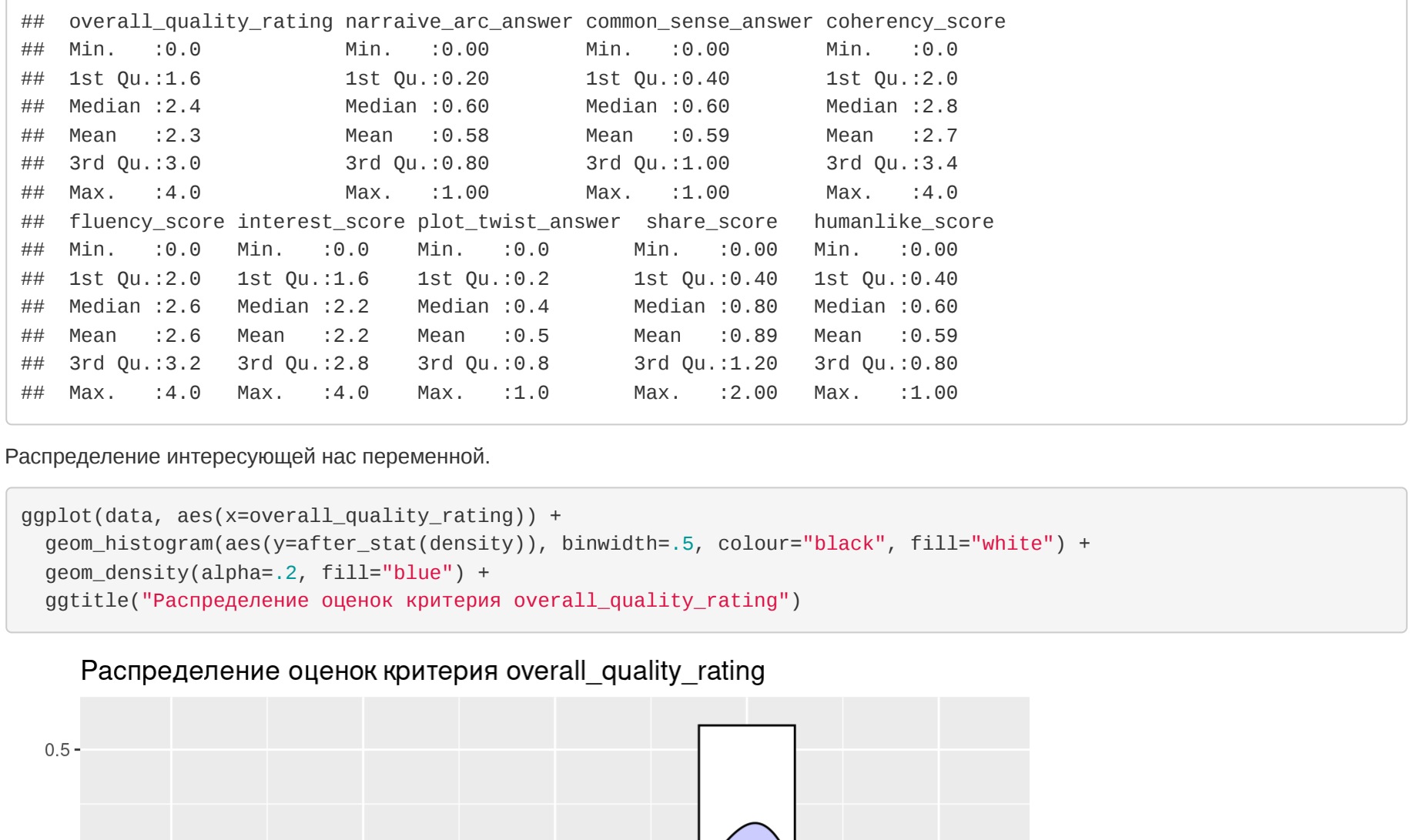
overall_quality_rating – главный критерий; общая оценка качества истории.
coherency_score – последовательность и логичность событий.
fluency_score – естественность и грамматичность языка.
interest_score – насколько история интересна.

Шкала 0 - 1

narrative_arc_answer – наличие нарративной арки.
common_sense_answer – соответствие истории здравому смыслу.
plot_twist_answer – наличие неожиданной смены событий.
humanlike_score – насколько история похожа на написанную человеком.

Шкала 0 - 2

share_score – желание поделиться историей с другими (да/не уверен/нет).



Гипотеза

Создателя истории (как сгенерированных, так и написанных) важно понимать, какие истории нравятся читателям и почему. В этом исследовании изучим факторы, составляющие общую оценку и выясним, какой вклад вносит каждый из них.

Гипотеза – общий критерий *overall_quality_rating* независим от остальных критериев, которые рассматриваются в этой работе.
Альтернативная гипотеза – общая оценка истории зависит от остальных критериев.

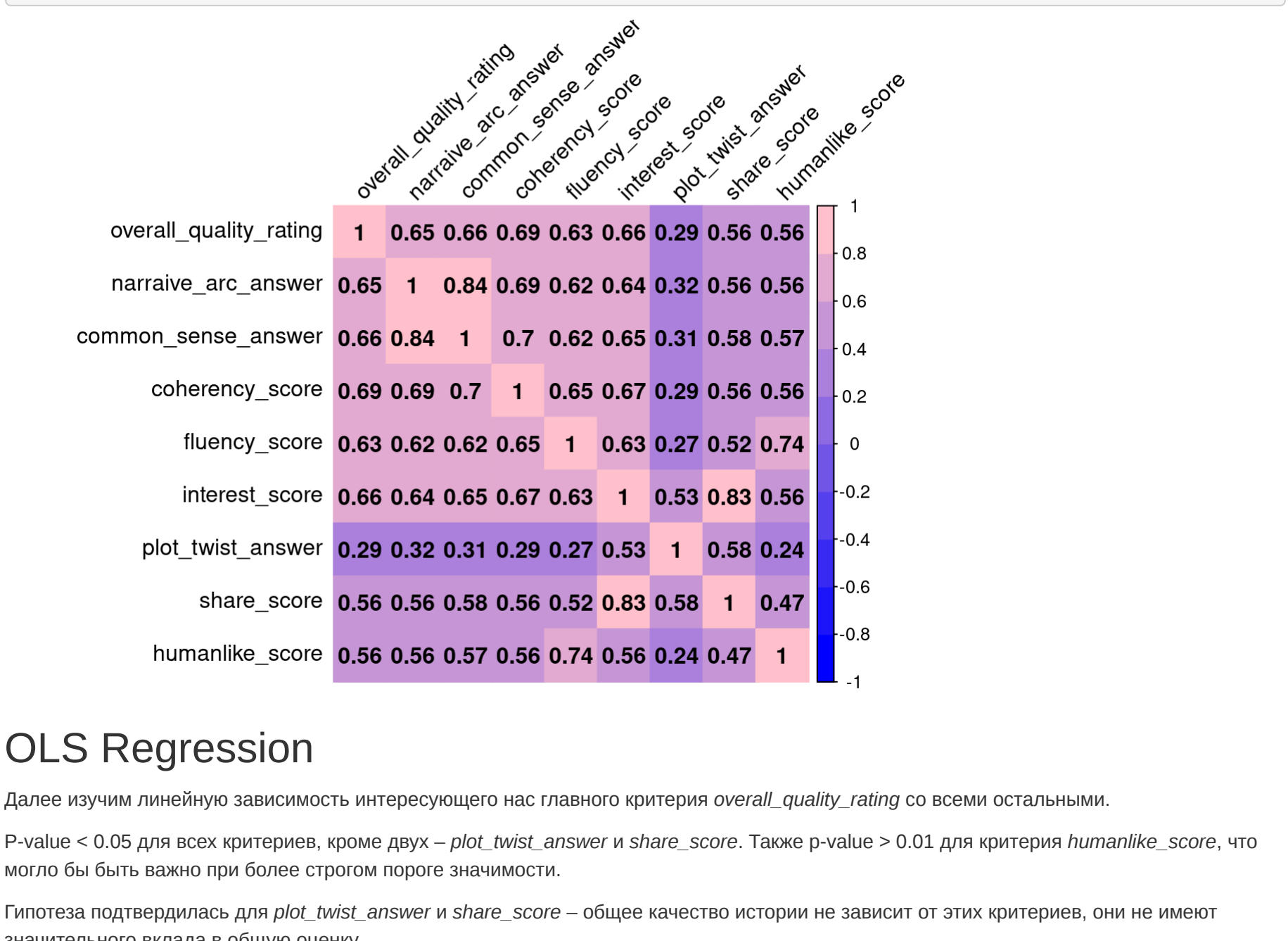
Корреляция

Сначала посчитаем корреляцию Пирсона между всеми критериями. Примечательные заметки по графику:

- Наиболее высокая корреляция для критерия *overall_quality_rating* наблюдается с критерием последовательности, а также с критериями интересности и соответствия здравому смыслу.
- Критерий *plot_twist_answer* имеет самую низкую корреляцию с остальными факторами, но при этом заметно коррелирует с фактором *interest_score*.

Также выделяются следующие факты:

- Критерий *share_score* положительно коррелирует с *interest_score*.
- Критерий *fluency_score* положительно коррелирует с *humanlike_score*.
- Критерий *common_sense_answer* положительно коррелирует с *narrative_arc_answer*, что довольно сложно объяснить.

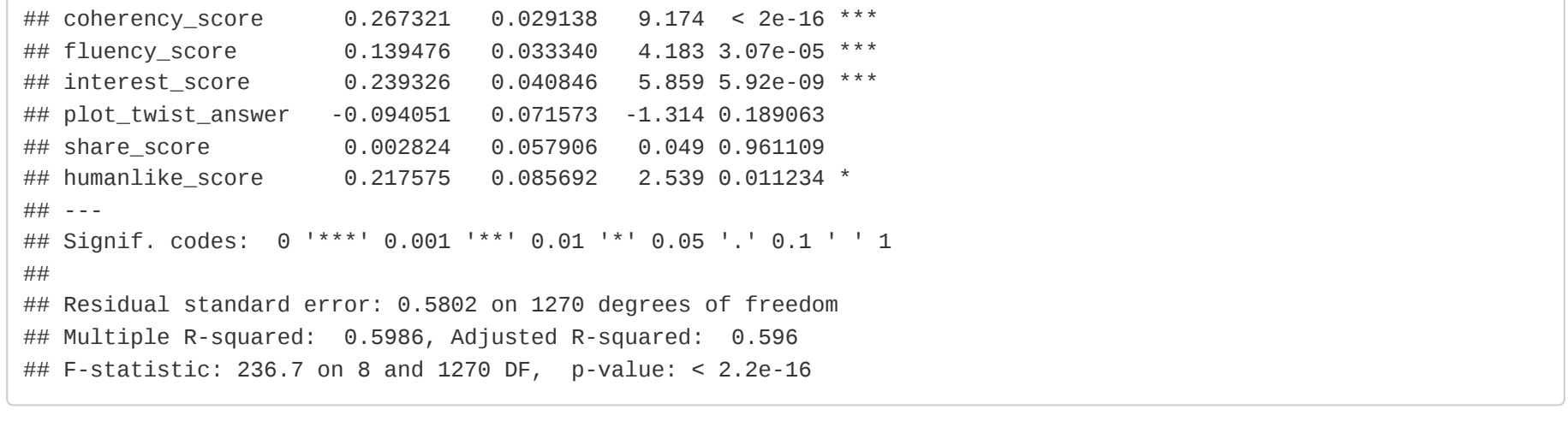


OLS Regression

Далее изучим линейную зависимость интересующего нас главного критерия *overall_quality_rating* со всеми остальными. R-value < 0.05 для всех критериев, кроме двух – *plot_twist_answer* и *share_score*. Также p-value > 0.01 для критерия *humanlike_score*, что могло бы быть важно при более строгом пороге значимости.

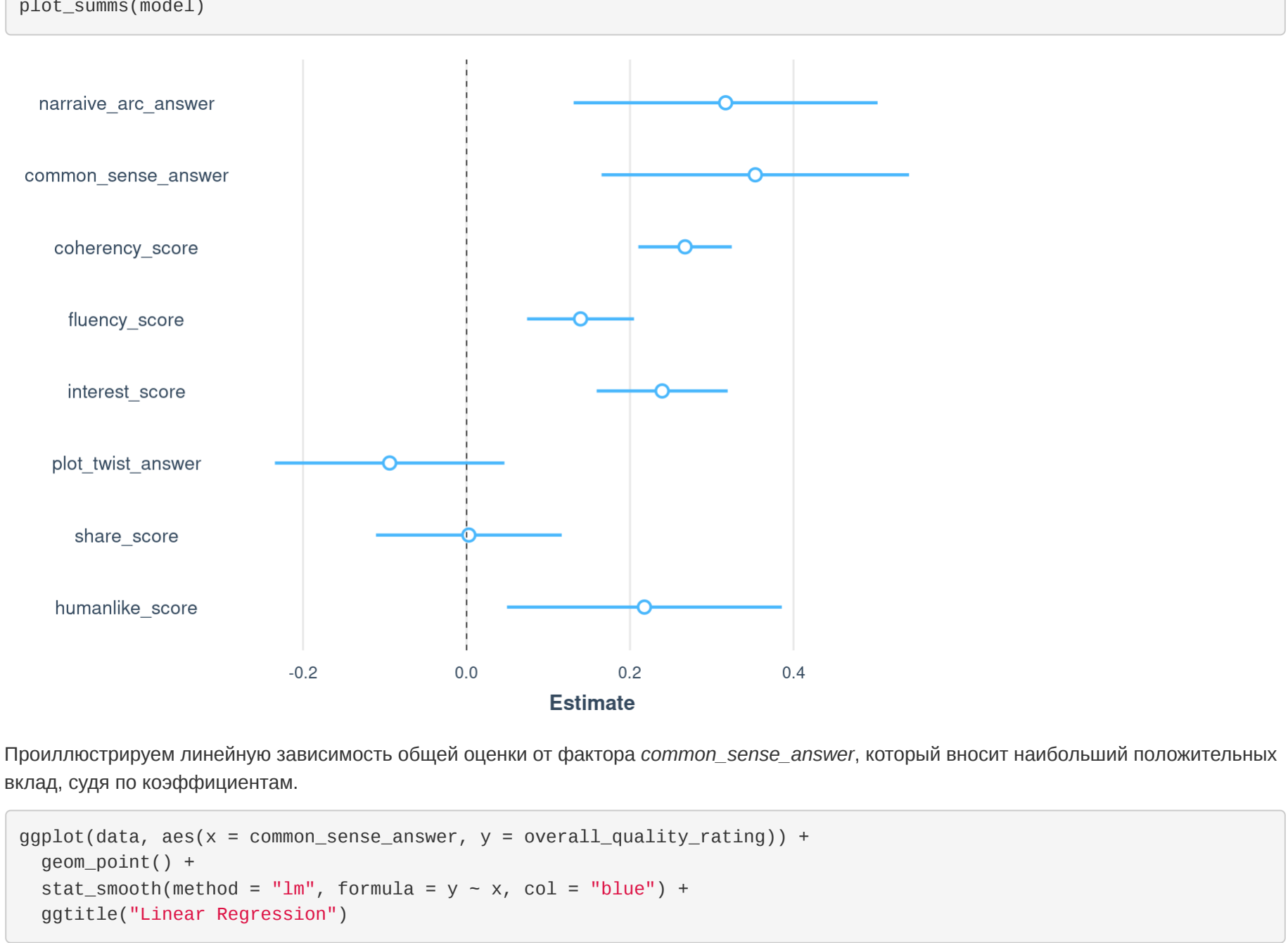
Гипотеза подтвердилась для *plot_twist_answer* и *share_score* – общее качество истории не зависит от этих критериев, они не имеют значительного вклада в общую оценку.

Для остальных критериев нулевая гипотеза была отвергнута и принята альтернативная. То есть общая оценка истории статистически значимо зависит от оставшихся факторов. Все обнаруженные зависимости положительные, значит, при улучшении отдельных факторов ожидается рост общего качества.

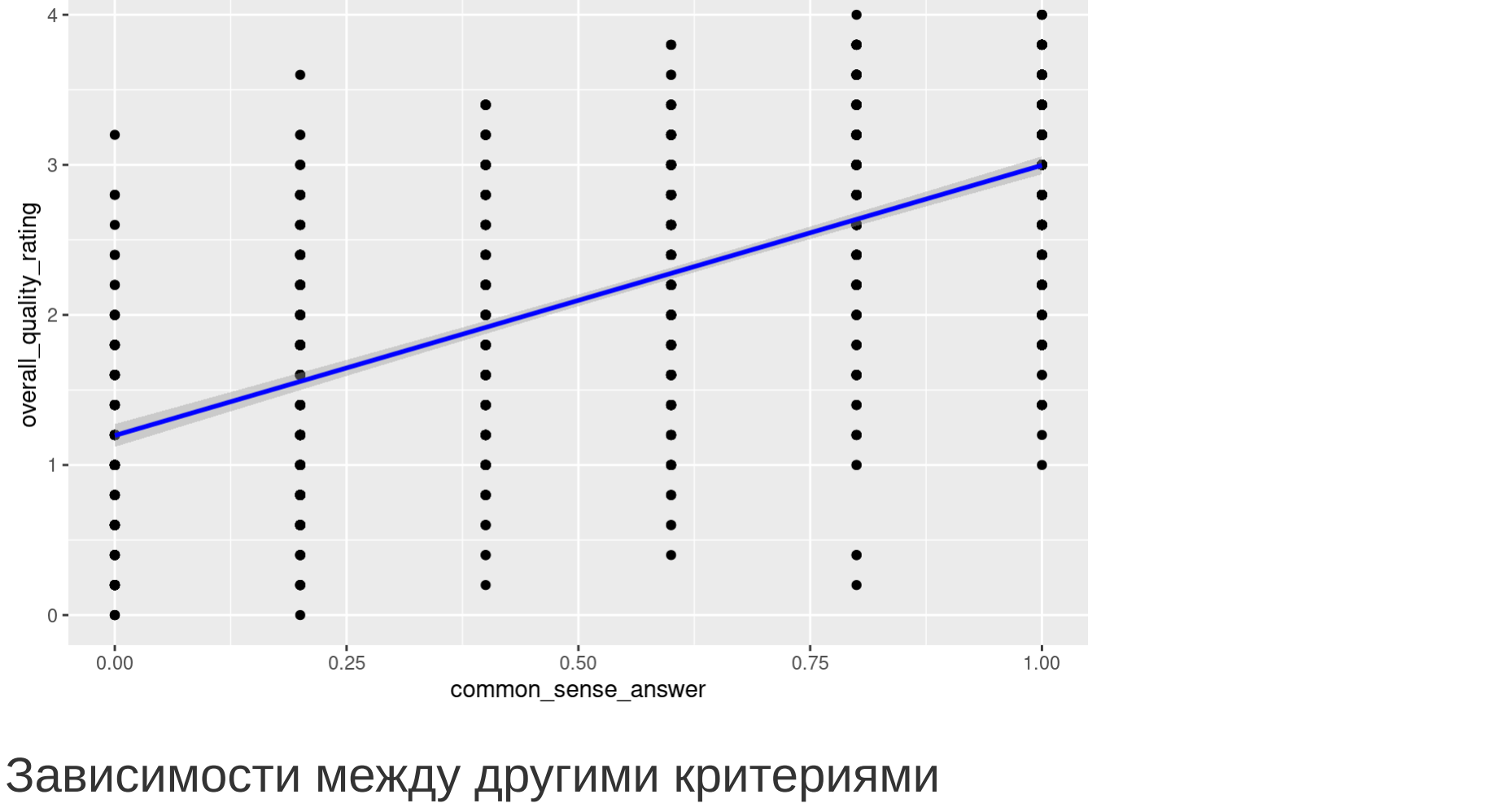


Построим график, на котором наглядно показаны коэффициенты при каждом факторе в линейной зависимости.

Хотя мы выяснили, что общее качество не зависит от наличия неожиданного перехода (plot twist), интересно, что коэффициент этого фактора отрицательный.



Проиллюстрируем линейную зависимость общей оценки от фактора *common_sense_answer*, который вносит наибольший положительных вклад, судя по коэффициентам.



Зависимости между другими критериями

Мы выяснили зависимость общего качества истории от остальных критериев. Интересно изучить зависимости между другими критериями, понять, как каждый из них объясним с помощью других.

Выводы по этой табличке дублируют те, что были приведены в анализе таблички корреляции, поэтому я их опускаю. Можно добавить, что последовательность (*coherency_score*) вносит большой вклад в большую часть критериев, как и *overall_quality_rating*.



Оценка сгенерированных и написанных людьми историй

В данных имеется информация о происхождении истории. Она могла быть написана человеком или могла быть сгенерированная языковой моделью.

Поставим нулевую гипотезу, что написанные людьми не отличаются в общей оценке от тех, что были сгенерированы. Проведем тест Манна-Уитни, потому что данные распределены не нормально.

Тест отверг нулевую гипотезу, так как p-value < 0.05. Получается, что написанные людьми истории оцениваются значительно выше (в среднем 2.9/4), а сгенерированные ниже (1.9/4). Интересно, что написанные истории тоже не идеальны.

