

## Рассмотрение признаков, выводы из визуального анализа:

### float data

**1. page\_likes\_num** (done: убрано, чтобы ускорить алгоритм, этот признак почти не влияет на результат)

- можно попробовать прологарифмировать (стало хуже на тысячные)
- можно попробовать убрать выбросы лучше после логарифмирования (стало хуже на тысячные)
- корреляция с целевой переменной слаба (0.06 по Пирсону), можно попробовать убрать (стало лучше на тысячные)

**2. page\_checkins** (done: убрано, чтобы ускорить алгоритм, этот признак почти не влияет на результат)

- слишком много нулей (больше половины: 25 тысяч из 40)
- может быть стоит убрать из рассмотрения из-за нулей
- может стоит перекодировать (0 против всех)
- корреляция с целевой переменной слаба

**3. page\_talking\_about** (done: логарифм от +0,05)

- может стоит прологарифмировать, но при этом стоит обработать -inf, например  $\text{np.log}(\text{data}[\text{'page\_talking\_about'}]+0.05)$  (стало лучше на сотые)
- возможно, стоит попробовать убрать этот признак, корреляция с целевой переменной слаба

**4. page\_statistics' 5-29 cols** (done: убрано 10, 15, 16, 17, 18, 19, 21, 23)

- 'page\_statistics\_10' и 'page\_statistics\_23' очень слабо коррелируют с целевой переменной (их можно убрать)
- между некоторыми значениями наблюдаются сильные корреляции (например 0 и 15 - по модулю больше 0,99), скорее всего стоит убрать: 10, 15, 16, 17, 18, 19, 21, 23 (более 0,99) (стало на несколько тысячных хуже) и возможно стоит убрать: 3, 7, 12, 14, 24 (более 0,95) (оставить)
- признаки 0, 5, 10 и 15 имеют слишком много нулей (больше 37 тысяч из 40). Их лучше убрать, либо один против всех (10 и 15 убраны, если убрать 0 и 5 - оценки уменьшаются на тысячные, если перекодировать - средняя оценка такая же, отклонение стало меньше, лучше будет убрать, так как дисбаланс очень большой)
- признак 23 имеет половину 0, можно сделать один против всех (признак убран)
- для 1,2,11 можно попробовать корень или логарифм распределения (оценки стали немного лучше)

**5. 'comments\_'** (done: убран признак comments\_num\_in\_first\_24\_hours, один против всех comments\_num\_in\_last\_48\_to\_24\_hours)

- 'comments\_num\_before\_base\_time' сильно коррелирует с 'comments\_num\_in\_first\_24\_hours', можно попробовать без одного из этих полей (убрать: comments\_num\_in\_first\_24\_hours)
- во всех наблюдается большое количество нулей, может стоит кодировать (один против всех, там где 15-20 тысяч и другое кодирование для других, но из-за больших корреляций с целевой переменной может быть не очень удачно)

- можно попробовать логарифм (не сильно повлияло на распределение)

#### **6. character\_num\_in\_post** (done: убрать)

- Целочисленное значение, которое можно попробовать закодировать (до 100 знаков и тд)
- Много нулей: чуть меньше 5 тысяч из 40
- Наблюдается закономерность с целевой переменной: чем больше символов, тем меньше лайков чаще всего (а корреляция очень слаба)

#### **7. share\_num** (done: оставить как есть)

- Нулей нет, много единиц: более 8 тысяч из 40, можно попробовать перекодировать по диапазонам (нет сильного влияния на оценку)
- Логарифм не сильно помогает из-за большого числа единиц

### **cat data**

#### **8. 'page\_cat'** (done: убрать)

- 81 категория
- каких-то значений много 7491, каких-то мало 1, возможно, признак не стоит рассматривать, либо стоит закодировать по частоте (при пороге 1500 оценки стали, при 2000 оценки стали незначительно хуже на тысячные)
- в pdf файле есть текстовые названия, названия содержат однокоренные слова (sport, art), можно попробовать рассмотреть как тексты (стало хуже на тысячные)
- можно закодировать one-hot (улучшилось, но не намного)
- можно и тексты, и one-hot (стало хуже на тысячные)

#### **9. base\_time** (done: one-hot)

- В документации указано, что это закодированное значение в диапазоне от 0 до 71, на практике от 0 до 72
- Разница в объеме данных для каждого значения небольшая (минимум 489 записей, максимум 625 записей)
- Распределение целевой переменной для некоторых значений могут сильно различаться
- Можно попробовать развернуть в one-hot encoding (оценки значительно улучшились), но тогда возможно будет потеряна связь между значениями признака, а зависимость целевой переменной от признака к признаку заметно: чем больше признак, тем меньше разброс значений переменной.

#### **10. h\_local** (done: one-hot)

- значения от 1 до 24, очень плохо сбалансировано, почти все значения равны 24. Можно попробовать убрать или сделать один против всех или one-hot (стало немного лучше)
- зависимость целевой переменной между соседними признаками не наблюдаема, корреляция с целевой переменной мала, но это категориальный признак

#### **11. post\_published\_weekday\_** (done: убрано)

- для каждого столбца 5-6 тысяч ненулевых значений

- для всех столбцов boxplot целевой переменной для значений 0 и 1 выглядит одинаково, но очень большое количество данных отсеиваются как выбросы, поэтому стоит посмотреть на матрицу корреляций

- из матрицы корреляций видно, что признаки очень слабо коррелируют с целевой переменной, можно попробовать без них

## **12. base\_datetime\_weekday\_ (done: убрано)**

- для каждого столбца 5-6 тысяч ненулевых значений

- для всех столбцов boxplot, кроме пятницы и субботы целевой переменной для значений 0 и 1 выглядит одинаково, но очень большое количество данных отсеиваются как выбросы, поэтому стоит посмотреть на матрицу корреляций

- из матрицы корреляций видно, что признаки очень слабо коррелируют с целевой переменной, можно попробовать без них

## **target**

- очень много нулей: 22 тысячи из 40

## **Возможные дополнительные признаки:**

1. Из признаков post\_published\_weekday\_ и base\_datetime\_weekday\_ вывести бинарный признак выходного дня
2. Отношение share\_num к comments\_num\_before\_base\_time
3. Отношение page\_talking\_about к page\_likes\_num #доля репостов к лайкам
4. Отношение comments\_num\_in\_last\_24\_hours к comments\_num\_before\_base\_time
5. Отношение comments\_num\_in\_last\_24\_hours к comments\_num\_in\_last\_48\_to\_24\_hours #скорость падения интереса
6. Отношение 'comments\_num\_in\_last\_24\_hours' к 'comments\_num\_in\_first\_24\_hours'
7. Попробовать квадратичные признаки

## **Выводы для новых признаков:**

1. **is\_post\_published\_in\_holiday** и **is\_base\_datetime\_holiday** - очень много нулей: 30 тысяч из 40. Из scatter plot можно увидеть небольшую зависимость целевой переменной, стоит посмотреть на корреляцию. Корреляции слабые. (результаты ухудшились на тысячные) (done: не добавлены)
2. **share\_part\_in\_comm** - можно попробовать взять логарифм, корреляция слаба (результаты ухудшились на тысячные) (done: не добавлен)
3. **page\_share\_part\_in\_likes** - стоит попробовать рассмотреть корень 4 степени. С этим признаком скорее всего можно работать. (результаты ухудшились на тысячные) (done: не добавлен)
4. **comm\_part\_1** - 15 тысяч нулей из 40, получилось мультимодальное распределение -> стоит разделить на 2 части бинарно: является 0 или нет (done: добавлен comm\_part\_1)
5. **comm\_part\_2** - 15 тысяч нулей из 40, может стоит сделать бинарным. Из выведенных признаков сильнее всего коррелирует с целевой переменной, скорее всего из-за числителя. (стало хуже на тысячные) (done: не добавлен)

6. **comm\_part\_3** - 15 тысяч нулей из 40, может стоит сделать бинарным. (done: не добавлен)