

Задание

Взять страницу ВК, собрать по ней таблицу с датой постов и количеством лайков и написать SQL-запросы, которые позволят ответить на вопрос: что больше всего влияет на количество лайков: время суток публикации, день недели или промежуток между постами.

Шаги

Для написания запросов использовал MySQL. Для визуализации использовал Apache Superset и Python.

Написал скрипт на Python и собрал данные о последних постах со страницы группы [Павла Воли](#). В базу данных включил следующие поля:

ИМЯ ПОЛЯ	ПОЯСНЕНИЕ
post_id	Номер поста (с конца)
post_date	Дата и время поста
attachment_type	Тип прилагаемых файлов (фото, видео, репост, другое)
text_size	Размер текста в посте
like_count	Количество лайков
comment_count	Количество комментариев
repost_count	Количество репостов
view_count	Количество просмотров
duration	Длительность (если приложено видео, в противном случае – пустое значение)

Помимо дня недели, времени суток и интервалом между постами, хотел посмотреть на влияние других факторов на количество лайков (вид поста, размер описания, количество комментариев и т.д.)

В таблицу включил данные о последних 550 постах со страницы группы.

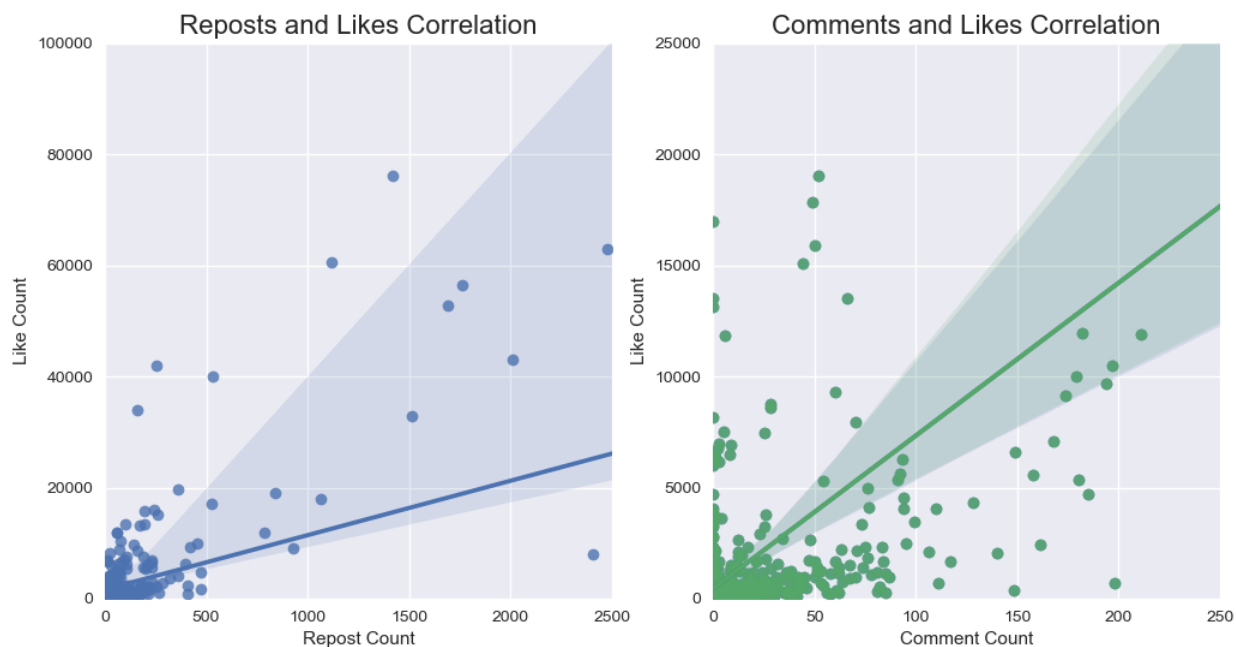
Используя библиотеку Pandas в Python, я рассчитал коэффициенты корреляции для количества лайков и следующих полей: количество репостов, комментариев, просмотров, размер текста и длительность видео. Результаты свел в таблицу:

ПОЛЕ	КОЛИЧЕСТВО ЛАЙКОВ
Количество репостов	0,705
Количество комментариев	0,659
Количество просмотров	0,404

Оказалось, что в наибольшей степени на количество лайков под постами влияют количество репостов, комментариев и просмотров (что, вообще говоря, является очевидным, поскольку с ростом

одного показателя пост попадаетея большому количеству пользователей, откуда растут и последующие показатели).

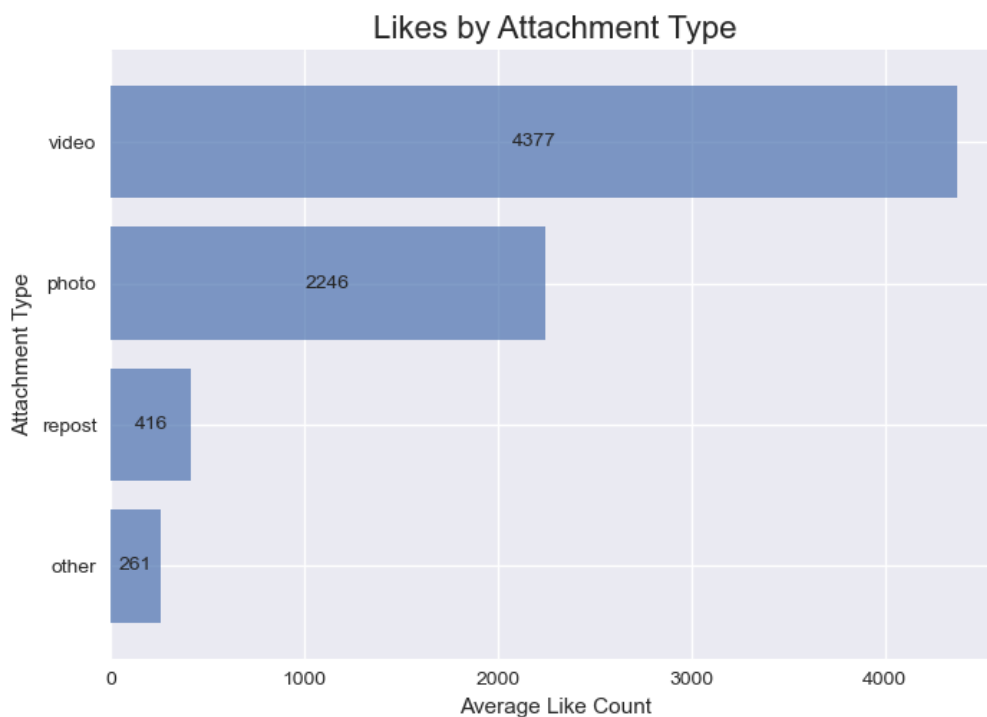
Более детально данные зависимости можно посмотреть на графиках ниже:



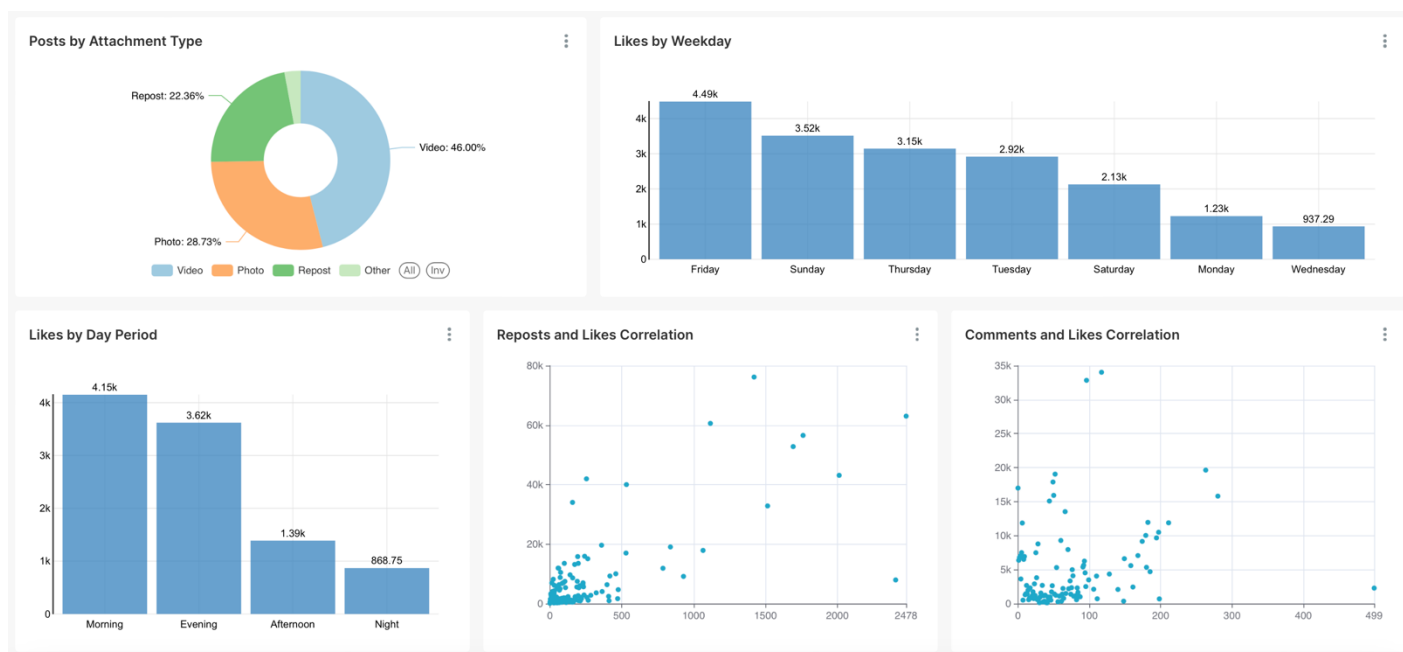
Слева – зависимость количества репостов от лайков, справа – количество комментариев от лайков.

Здесь более детально видна зависимость между двумя показателями. Она в особенности заметна на правом графике.

Также важно отметить, что в среднем количество лайков больше на постах с видео. Вероятнее всего, что данный показатель зависит от типа контента на определенной странице. В данном случае полученный результат имеет смысл, учитывая то, какой контент размещается на выбранной мною странице.

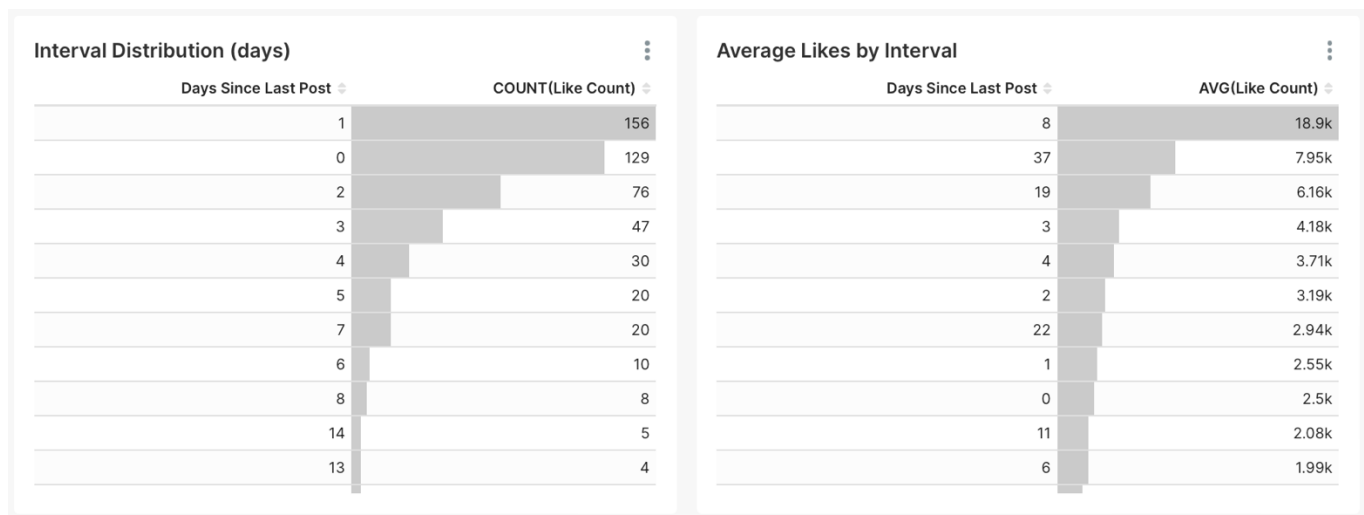


Среднее число лайков по типу постов



Для визуализации остальных показателей я построил дашборд в Superset. Из него можно сделать следующие выводы:

- Виды постов: видео (46%), фото (28,73%), репосты (22,36%), другие (2,91%)
- В среднем больше лайков набирают посты, выложенные в пятницу и воскресенье утром (с 6 до 11) и вечером (с 18 до 23)
- Также включены графики зависимости количества комментариев и репостов от лайков.



Слева – распределение постов по интервалам, справа – среднее количество лайков под постом, опубликованный с определенным интервалом

При исследовании зависимости интервалов между постами и количеством лайков можно сделать следующие выводы:

- За все время большинство постов выкладывались с интервалом 0 (несколько постов в один и тот же день) – 4 дня
- Больше всего лайков набрали посты с интервалом в 8 дней.

Основываясь на полученных данных, сложно сказать о зависимости между частотой публикаций и количеством лайков. Тем более в данную выборку с интервалом 8 дней входят только 8 постов. Возможно, для более точного результата понадобилась бы выборка большего размера, однако вероятность того, что результат будет другим – мала.

ВЫВОД

Для анализа использовалась выборка с последними постами со страницы группы Павла Воли. В итоге в таблице собралось 550 записей.

На основе полученных данных можно сделать следующие выводы:

1. количество лайков напрямую зависит от количества репостов, комментариев и просмотров. С увеличением данных показателей растет популярность поста и, соответственно, количество лайков под ним;
2. большую активность пользователи проявляют в пятницу и воскресенье в утреннее и вечернее время;
3. интервал между постами, как оказалось, не играет большой роли в количестве лайков под постом.

SQL-ЗАПРОСЫ:

```
1  /*
2   Основной SQL-запрос
3   Использовал оператор CASE, чтобы из всех типов постов
4   выбрать наиболее популярные (фото, видео, репост), вывел их
5   с заглавной буквы
6   Оставшиеся пометил как "другие"
7  */
8
9  SELECT
10     post_id,
11     post_date,
12     CASE
13     WHEN attachment_type IN ('photo', 'video', 'repost') THEN
14         CONCAT(UPPER(SUBSTR(attachment_type, 1, 1)), SUBSTR(attachment_type, 2))
15     ELSE 'Other'
16     END AS 'attachment_type',
17     text_size,
18     like_count,
19     comment_count,
20     repost_count,
21     view_count,
22     duration
23 FROM posts;
```

```

1 /*
2  Данный запрос выводит информацию
3  о среднем количестве лайков в разное время суток
4  Группируется по времени суток и сортируется по
5  среднему количеству лайков в порядке убывания
6 */
7
8 SELECT
9     CASE
10     WHEN HOUR(post_date) between 0 AND 5 THEN 'Night'
11     WHEN HOUR(post_date) between 6 AND 11 THEN 'Morning'
12     WHEN HOUR(post_date) between 12 AND 17 THEN 'Afternoon'
13     WHEN HOUR(post_date) between 18 AND 23 THEN 'Evening'
14     END AS 'post_hour',
15     ROUND(AVG(like_count), 2) as avg_likes
16 FROM
17     posts
18 GROUP BY
19     post_hour
20 ORDER BY
21     avg_likes DESC

```

```

1 /*
2  Данный запрос выводит информацию
3  о среднем количестве лайков в разные дни недели
4  Группируется по дню недели и сортируется по
5  среднему количеству лайков в порядке убывания
6 */
7
8 SELECT
9     DAYNAME(post_date) as day_of_week,
10     ROUND(AVG(like_count), 2) as avg_likes
11 FROM
12     posts
13 GROUP BY
14     day_of_week
15 ORDER BY
16     avg_likes DESC;

```



```
1 /*
2  Данный запрос выводит информацию о постах и количестве
3  дней, прошедших с момента публикации предыдущего поста.
4  Чтобы это реализовать, я использовал CROSS JOIN. Я
5  соединил таблицы таким образом, чтобы каждый пост в таблице
6  "а" сопоставлялся с предыдущим по счету постом в таблице
7  "б". Это условие я реализовал в WHERE (поскольку посты в
8  таблице находятся в порядке убывания по дате,
9  ID таблицы "а" должно быть на 1 меньше ID таблицы "б").
10 В запрос включено количество постов для каждого интервала
11 (post_count) и среднее количество лайков для каждого
12 интервала (avg_likes)
13 */
14
15 SELECT
16     DATEDIFF(a.post_date, b.post_date) AS days_since_last_post,
17     COUNT(a.like_count) AS post_count,
18     AVG(a.like_count) as avg_likes
19 FROM
20     posts a, posts b
21 WHERE
22     a.post_id = b.post_id - 1
23 GROUP BY
24     days_since_last_post;
```