МИНОБРНАУКИ РОССИИ САНКТ-ПЕТЕРБУРГСКИЙ ГОСУДАРСТВЕННЫЙ ЭЛЕКТРОТЕХНИЧЕСКИЙ УНИВЕРСИТЕТ «ЛЭТИ» ИМ. В.И. УЛЬЯНОВА (ЛЕНИНА)

Кафедра Информационной безопасности

ОТЧЕТ

по лабораторной работе № 3

по дисциплине «Распределенные системы обработки данных»

Tema: Использование PySpark для самостоятельной реализации витрины данных

Студентка гр. 1361	 Токарева У.В.
Преподаватель	Троценко В.В

Санкт-Петербург

Цель работы.

Очистить и загрузить сырые данные. Выделить категориальные переменные. Выделить бинарные признаки. Построить собственную витрину.

Основные теоретические положения.

Категориальные переменные (Categorical Variables) — это переменные, которые принимают одно из нескольких возможных значений, но для каждого наблюдения может быть указано только одно значение. Категориальные переменные удобны для анализа, когда каждая запись может быть отнесена к одной четко определенной категории, что облегчает агрегацию и создание групп по этим переменным.

Бинарные признаки (Binary Features) — это признаки, которые могут быть либо "истинными" (1), либо "ложными" (0) для каждого наблюдения. Они отражают наличие или отсутствие определенных характеристик и позволяют применять несколько признаков к одному объекту. Бинарные признаки удобны для отображения наличия или отсутствия множества характеристик в одном объекте. Они позволяют использовать несколько признаков для одного видео, что делает их гибким инструментом для меток.

Ход работы.

1. В первую очередь мы должны взять данные из предоставленного нами файла. Производим чтение CSV-файла и сохраняем их в новую таблицу:

```
csv_path = "lab3/" + NAME_DATACSV

df = spark.read.csv(csv_path, header=True,
inferSchema=True, multiLine=True)

df.createOrReplaceTempView("YT Data")
```

2. Необходимо выделить категориальные переменные. В моем случае категория: 'Трендовые персоны (Илон Маск, 21 Savage, Магнус Карлсен)'. Поэтому определяем структуру для представления категорий видео:

```
schema = StructType([
```

```
StructField("id", IntegerType(), True),
StructField("name", StringType(), True),
StructField("template", StringType(), True)

# Определяем категории видео
category = [
    (1, "Илон Маск", "%Илон Маск%"),
    (2, "21 Savage", "%21 Savage%"),
    (3, "Магнус Карлсен", "%Магнус Карлсен%")

3. Создаем статистику по общим показателям видео для каждого
месяца и сохраняем результаты во временной таблице:
view("df_total_stat", """

SELECT
```

```
SELECT

date_format(YT_data.trending_date, 'yyyy-MM')

AS YM,

SUM(YT_data.view_count) AS total_view,

SUM(YT_data.comment_count) AS total_comment,

SUM(YT_data.likes) AS total_likes,

ROUND ((SUM(YT_data.comment_count) /

SUM(YT_data.likes)), 2) AS total_CLR

FROM YT_Data AS YT_data

GROUP BY date_format(YT_data.trending_date, 'yyyy-MM')

""")
```

В данной таблице собрана основная статистика по просмотрам, комментариям, лайкам и коэффициент CLR за каждый месяц. Вывод общей статистики представлен на рисунке 1.

```
🚺 toktalk@LAPTOP-192U10JP: ~/leti-spark-course-2024
     YM|total_view|total_comment|total_likes|total_CLR|
2020-12 3524072622
                     3.0699365E7
                                    286818208
                                                    0.11
2023-08 | 5468967208 |
                     1.5162869E7
                                    305918347
                                                    0.05
2024-02 | 6897478524 |
                     1.1091736E7
                                    292969924
                                                    0.04
2021-07 4644279727
                     3.1820584E7
                                    332950747
                                                     0.1
2022-04 4681734913
                     1.6416418E7
                                    258894893
                                                    0.06
2023-12 6742291412
                     1.2821919E7
                                    364242277
                                                    0.04
2021-04 | 3557247089 |
                     2.9190677E7
                                    255985391
                                                    0.11
2023-11 6084022728
                     1.2382188E7
                                    314392817
                                                    0.04
2022-06 4735167314
                     2.6830589E7
                                    322408117
                                                    0.08
2021-06 6792900896
                     4.7180891E7
                                    425362051
                                                    0.11
2022-02 3944307350
                                                    0.08
                      1.776408E7
                                    219736901
2022-08 | 5134911223 |
                                    327417897
                                                    0.09
                     3.0337922E7
2023-07 4340800759
                                                    0.05
                     1.3167347E7
                                    280881868
2021-11 3632109978
                     1.8646073E7
                                    263235644
                                                    0.07
2023-03 4945941220
                                                    0.05
                     1.5626489E7
                                    289556726
2023-10 7298890701
                     1.3710554E7
                                    350066718
                                                    0.04
2021-03 | 3875417671 |
                     3.3547777E7
                                    303693165
                                                    0.11
2021-02 | 3111801392 |
                                    228490913
                                                     0.1
                     2.2457705E7
2022-07 4917541950
                     2.0388889E7
                                    312371852
                                                    0.07
2023-02 4664164676
                     1.1944183E7
                                    230182230
                                                    0.05
```

Рисунок 1 – Вывод общей статистики.

4. Далее создаем таблицу со всеми категориями с изначальной таблицы. Рассматриваются совпадения в названиях, описаниях и тегах. Результат обработки сохраняется в таблице:

Вывод таблицы со всеми категориями представлен на рисунке 2.

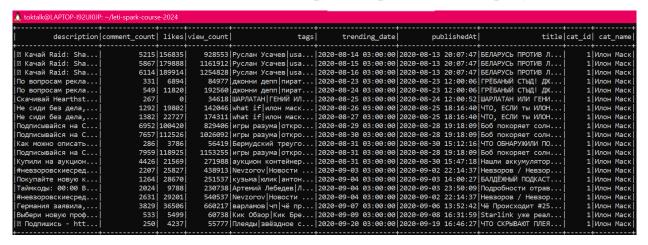


Рисунок 2 – Вывод таблицы со всеми категориями.

5. Далее создаем статистику для категорий по просмотрам, комментариям и лайкам:

Вывод таблицы со статистикой по категориям представлен на рисунке 3.

⚠ toktalk@	DLAPTOP-192U	JIOJP: ~/leti-spark-	course-2024			
+ YM	cats_view	cats_comment	cats_likes	cat_id	cat_	 _name +
2022-07	8786736	41075.0	534076	1	Илон	Маск
2022-12	21551069	93777.0	1013365	1	Илон	Маск
2022-11	20425282	72519.0	882849	1	Илон	Маск
2023-09	2698295	9043.0	136861	1	Илон	Маск
2020-09	2238457	12738.0	139728	1	Илон	Маск
2021-01	3075220	18998.0	221970	1	Илон	Маск
2022-02	4366644	18902.0	374159	1	Илон	Маск
2020-11	1433003	10163.0	198175	1	Илон	Маск
2021-07	3977970	26615.0	204058	1	Илон	Маск
2022-05	19643462	107779.0	991697	1	Илон	Маск
2022-06	9830552	53835.0	608053	1	Илон	Маск
2020-08	7310965	48297.0	945106	1	Илон	Маск
2023-02	4064744	25418.0	233757	1	Илон	Маск
2021-03	7017350	55633.0	513564	1	Илон	Маск
2024-01	2122229	1537.0	121483	1	Илон	Маск
2020-12	2934959	20922.0	338903	1	Илон	Маск
2021-06	16436199	77465.0	1806276	1	Илон	Маск
2023-08	3923772	32121.0	241393	1	Илон	Маск
2021-11	2467697	8890.0	255567	1	Илон	Маск
2023-12	13865888	21399.0	818103	1	Илон	Маск
+						+

Рисунок 3 – Вывод таблицы со статистикой по категориям.

6. Далее делаем сборку финальной таблицы для категорий:

Вывод финальной таблицы для категорий представлен на рисунке 4.

YM	total_comment	total_likes	total_view	total_CLR	cat_name	cats_view	cats_comment	cats_likes	cats_CLR	diff_CLR_abs	diff_CLR_relative
924- 9 3	9218794.0	351500982	 8676586316	0.03	Илон Масн	13791788	56576.0	579172	0.1	-0.07	0.3
924-92	1.1091736E7	292969924	6897478524	0.04	Илон Маск	10456901	38302.0	514000	0.07	-0.03	0.5
924-91	9150624.0	271780911	6060824937	0.03	Илон Маск	2122229	1537.0	121483	0.01	0.02	3.6
023 -1 2	1.2821919E7	364242277	6742291412	0.04	Илон Маск	13865888	21399.0	818103	0.03	0.01	1.3
023-11	1.2382188E7	314392817	6084022728	0.04	Илон Маск	1033966	3609.0	50384	0.07	-0.03	0.5
923 -1 0	1.3710554E7	350066718	7298890701	0.04	Илон Маск	8631473	50838.0	364512	0.14	-0.1	0.2
<mark>023-0</mark> 9	1.1577666E7	301046210	5963195144	0.04	Илон Масн	2698295	9043.0	136861	0.07	-0.03	0.5
<mark>923-08</mark>	1.5162869E7	305918347	5468967208	0.05	Илон Масн	3923772	32121.0	241393	0.13	-0.08	0.3
<mark>923-08</mark>	1.5162869E7	305918347	5468967208	0.05	Магнус Карлсен	451348	182.0	16339	0.01	0.04	5.
<mark>923-07</mark>	1.3167347E7	280881868	4340800759	0.05	Илон Маск	8544303	54794.0	696039	0.08	-0.03	0.6
<mark>923-06</mark>	1.4849175E7	272594274	5153586714	0.05	Илон Масн	10699341	129006.0	619005	0.21	-0.16	0.2
<mark>023-05</mark>	1.2243495E7	273610800	5281841221	0.04	Илон Маск	5414009	31273.0	348292	0.09	-0.05	0.4
<mark>023-05</mark>	1.2243495E7	273610800	5281841221	0.04	Магнус Карлсен	2217571	875.0	98808	0.01	0.03	4.
923-04	1.5712181E7	282561679	4998614296	0.06	Илон Масн	24457324	179058.0	869315	0.21	-0.15	0.2
<mark>023-04</mark>	1.5712181E7	282561679	4998614296	0.06	Магнус Карлсен	23442292	160937.0	731488	0.22	-0.16	0.2
<mark>023-03</mark>	1.5626489E7	289556726	4945941220	0.05	Илон Маск	6445182	23416.0	566580	0.04	0.01	1.2
<mark>023-02</mark>	1.1944183E7	230182230	4664164676	0.05	Илон Маск	4064744	25418.0	233757	0.11	-0.06	0.4
923- <mark>01</mark>	1.2962921E7	228015517	4394086204	0.06	Илон Масн	4131615	11976.0	152807	0.08	-0.02	0.7
22-12	1.5069554E7	286671504	5078994752	0.05	Илон Маск	21551069	93777.0	1013365	0.09	-0.04	0.5
322-11	1.3361672E7	250505617	4596294902	0.05	Илон Маск	20425282	72519.0	882849	0.08	-0.03	0.6

Рисунок 4 – Вывод финальной таблицы для категорий.

7. Создаем таблицу с бинарными признаками, отдельно рассматривая для каждого признака и той ситуации, когда они встречаются вместе:

```
,YT data.view count,YT data.tags,YT data.trending date,
YT data.publishedAt,YT data.title
    FROM YT Data AS YT data
    WHERE (YT data.description ILIKE '%#challenge%'
      OR YT data.tags ILIKE '%challenge%')
      AND (YT data.description NOT ILIKE '%#asmr%'
      OR YT data.tags NOT ILIKE '%asmr%')
""")
view("df bins asmr", """
        SELECT
YT data.description, YT data.comment count, YT data.likes
,YT data.view count,YT data.tags,YT data.trending date,
YT data.publishedAt,YT data.title
    FROM YT Data AS YT data
    WHERE (YT data.description ILIKE '%#asmr%'
    OR YT data.tags ILIKE '%asmr%')
    AND (YT data.description NOT ILIKE '%#challenge%'
        OR YT data.tags NOT ILIKE '%challenge%')
""")
view("df bins chandas", """
        SELECT
YT data.description, YT data.comment count, YT data.likes
,YT data.view count,YT data.tags,YT data.trending date,
YT data.publishedAt,YT data.title
    FROM YT Data AS YT data
    WHERE (YT data.description ILIKE '%#challenge%'
        OR YT data.tags ILIKE '%challenge%')
```

```
AND (YT_data.description ILIKE '%#asmr%'
OR YT_data.tags ILIKE '%asmr%')
""")
```

Вывод таблицы с бинарным признаком «challenge» с изначальной таблицей представлен на рисунке 5.

+			+	+	+	++	
description	comment_count	likes	view_count	tags	trending_date	publishedAt	title
Самый сложный чел	1057	10272	72790	Как похудеть как	2020-08-13 03:00:00	2020-08-12 14:49:38	КАК Я ПОХУДЕЛА ЗА
ТО по-настоящему	13184	281865	2228152	КТО ВЫЖИВЕТ НА ОТ	2020-08-14 03:00:00	2020-08-13 16:00:11	КТО ВЫЖИВЕТ НА ОТ
оня и разносчица	0	46496	585215	Непета Непета Стр	2020-08-14 03:00:00	2020-08-14 10:00:11	СИРЕНОГОЛОВЫЕ ПРО
кадре — команда	1322	29505	256457	Реакция NAVI CSGO	2020-08-14 03:00:00	2020-08-13 16:05:16	Реакция NAVI CSGO
оня и разносчица	0			Непета Непета Стр			
ТО по-настоящему	14508	317794		КТО ВЫЖИВЕТ НА ОТ			
ы готовы к новом	97	3105		123go челлендж 12			
одписывайся на м	6471	101163	1112021	Купить ПОНТОРЕЗКУ	2020-08-16 03:00:00	2020-08-15 19:29:14	Купить ПОНТОРЕЗКУ
оня и разносчица	0	59450	1140056	Непета Непета Стр	2020-08-16 03:00:00	2020-08-14 10:00:11	СИРЕНОГОЛОВЫЕ ПРО
ТО по-настоящему	15187	337087		КТО ВЫЖИВЕТ НА ОТ			
(еребьевка турнир	69	5742	33197	FIFA FIFA 17 FIFA	2020-08-17 03:00:00	2020-08-17 11:07:51	МЫСЛИТЬ КАК ПОДПИ
одписывайся на м	7237	115705	1377461	Купить ПОНТОРЕЗКУ	2020-08-17 03:00:00	2020-08-15 19:29:14	Купить ПОНТОРЕЗКУ
оня и разносчица	0	62133	1223877	Непета Непета Стр	2020-08-17 03:00:00	2020-08-14 10:00:11	СИРЕНОГОЛОВЫЕ ПРО
ream Team Family	4372	64492	450557	дрим тим хаус дри	2020-08-18 03:00:00	2020-08-17 15:00:06	РЕАКЦИЯ на УЖАСНЫ
незапно захотело	134	3484		123go челлендж 12			
OMAROY https://w	2460	11053		ФУТБОЛ FOOTBALL Г			
незапно захотело	200	4592	338084	123go челлендж 12	2020-08-20 03:00:00	2020-08-19 06:00:05	ЧЕЛЛЕНДЖ С СИНЕЙ
лава:Инстаграм:	135			КТО ВЫЖИВЕТ НА ОТ			
бязательно подпи	4630	50996	278395	ДЕНЬГИ ВЗЛОМАЕШЬ	2020-08-21 03:00:00	2020-08-21 11:30:11	ДЕНЬГИ или ВЗЛОМА
OMAROY https://w	4337	13051	129866	ФУТБОЛ FOOTBALL Г	2020-08-21 03:00:00	2020-08-20 10:59:17	ЛУЧШИЙ МАТЧ КУБКА

Рисунок 5 — Вывод таблицы с бинарным признаком «challenge» с изначальной таблицей.

Вывод таблицы с бинарным признаком «asmr» с изначальной таблицей представлен на рисунке 6.

	++					+	
description	comment_count	likes	view_count	tags	trending_date	publishedAt	title
Кочешь знать, что	1125	9057	103667	асмр асмр для сна	2020-08-12 03:00:00	2020-08-11 19:30:03	АСМР КАК ИСПЫТАТЬ
Спасибо за просмо	999	3817	56286	ACMP ASMR Расслаб	2020-08-14 03:00:00	2020-08-13 19:56:02	ASMR LICKING MARV
Спасибо за просмо	1236	5060	86494	ACMP ASMR Расслаб	2020-08-15 03:00:00	2020-08-13 19:56:02	ASMR LICKING MARV
В этом видео я от	1157	14291		Вопрос Ответ АСМР	2020-08-15 03:00:00		
Thanks for watchi	44	2587	317726	Kluna Tik Kluna T	2020-08-18 03:00:00	2020-08-17 16:00:12	KLUNA TIK Destroy
Внезапно захотело	134	3484	202515	123go челлендж 12			
незапно захотело				123go челлендж 12	2020-08-20 03:00:00		
Во время просмотр				асмр asmr мурашки	2020-08-22 03:00:00		
Во время просмотр		7190		асмр asmr мурашки	2020-08-23 03:00:00		
№ КУПИТЬ МОЙ МЕРЧ	605	2616		асмр asmr relax w			
КУПИТЬ МОЙ МЕРЧ	788	3546		асмр asmr relax w	2020-08-26 03:00:00	2020-08-24 19:00:00	АСМР НЕДОВОЛЬНЫЙ
Сегодня тебя ждет	1277	15751		ACMP для сна ASMR	2020-09-01 03:00:00		
Сегодня тебя ждет				ACMP для сна ASMR			
[ASMR] Во время п		3598		асмр asmr асмр дл	2020-09-05 03:00:00		
Hello friends! Ap	15789	149720	1049185	I tried acrylic p	2020-09-06 03:00:00	2020-09-06 04:53:34	I Tried Acrylic P
[ASMR] Во время п				aсмр asmr aсмр дл	2020-09-06 03:00:00	2020-09-04 19:31:58	ACMP [ASMR] 100%
Hello friends! Ap		231120		I tried acrylic p			
(ак приготовить с	1435	31218	353950	rfr как приготови	2020-09-08 03:00:00	2020-09-08 12:00:47	ГОВЯЖЬЕ СЕРДЦЕ в
(ак приготовить с	2359	51291	901422	rfr как приготови	2020-09-09 03:00:00	2020-09-08 12:00:47	ГОВЯЖЬЕ СЕРДЦЕ в
Как приготовить с	2674	59065	1178107	rfr как приготови	2020-09-10 03:00:00	2020-09-08 12:00:47	ГОВЯЖЬЕ СЕРДЦЕ в

Рисунок 6 – Вывод таблицы с бинарным признаком «asmr» с изначальной таблицей.

Вывод таблицы с бинарными признаками с изначальной таблицей представлен на рисунке 7.

	+		+	+			.
description	comment_count	likes	view_count	tags	trending_date	publishedAt	titl
Внезапно захотело	134	3484	202515	 123go челлендж 12	2020-08-19 03:00:00	2020-08-19 06:00:05	ЧЕЛЛЕНДЖ С СИНЕЙ .
Внезапно захотело	200	4592	338084	123go челлендж 12	2020-08-20 03:00:00	2020-08-19 06:00:05	челлендж с синей .
Черная пятница в	8525	33602	296364	Самый острый соли	2020-11-29 03:00:00	2020-11-28 16:48:53	САМЫЙ ОСТРЫЙ ЧИПС.
Черная пятница в	9464	38527		Самый острый соли			
2 2 2 2 ! 2 2 2 2 2	74	3031		MOOMOOSTUDIO MOOM			
2 2 2 2 ! 2 2 2 2 2	115			MOOMOOSTUDIO MOOM			
2 2 2 2 ! 2 2 2 2 2	139			MOOMOOSTUDIO MOOM			
2 2 2 2 ! 2 2 2 2 2	152			MOOMOOSTUDIO MOOM			
2 2 2 2 ! 2 2 2 2 2	183			MOOMOOSTUDIO MOOM			
Some of you have		500451		hairstyle kpop bo			
A fun viral tikto		152519		Pop it Popit Popi			
A fun viral tikto		223132		Pop it Popit Popi			
A fun viral TikTo		254614		Pop it Popit Popi			
A fun viral TikTo		525517		Pop it Popit Popi			
ТОП-5 наушников Р				ИРП сухой паёк су			
ТОП-5 наушников Р	2243			ИРП сухой паёк су			
ТОП-5 наушников Р				ИРП сухой паёк су			
Did you like this				Pop it Popit Fidg			
Compilation asmr		5763		dessert bottle je			
Compilation asmr	99	6817	1236820	dessert bottle je	2023-01-28 03:00:00	2023-01-26 10:00:10	ASMR MUKBANG Supe

Рисунок 7 – Вывод таблицы с бинарными признаками с изначальной таблицей.

8. Создаем статистику для бинарных признаков:

```
view("df bins stat challenge", """
    SELECT SUM(a.likes) AS bins likes,
         SUM(a.comment count) AS bins comment,
         SUM(a.view count) AS bins view,
         date format(a.trending date, 'yyyy-MM') AS YM
    FROM df bins challenge AS a
    GROUP BY date format(a.trending date, 'yyyy-MM')
""")
view("df bins stat asmr", """
    SELECT SUM(a.likes) AS bins likes,
         SUM(a.comment count) AS bins comment,
         SUM(a.view count) AS bins view,
         date format(a.trending date, 'yyyy-MM') AS YM
    FROM df bins asmr AS a
    GROUP BY date format(a.trending date, 'yyyy-MM')
""")
view("df bins stat chandas", """
```

```
SELECT SUM(a.likes) AS bins_likes,

SUM(a.comment_count) AS bins_comment,

SUM(a.view_count) AS bins_view,

date_format(a.trending_date, 'yyyy-MM')AS YM

FROM df_bins_chandas AS a

GROUP BY date_format(a.trending_date, 'yyyy-MM')
```

Вывод таблицы статистики для бинарного признака «challenge» представлен на рисунке 8.

🍌 toktalk@LAF	PTOP-192UI0JP: ~/le	ti-spark-cour	se-2024
+	+-		+
bins_likes	bins_comment	bins_view	YM
+	+-		+
10462419	714363.0	104271164	2020-12
2846149	46318.0	57730238	2023-08
7519272	116776.0	329973126	2024-02
21622837	451671.0	246460367	2021-07
6834511	170810.0	150237564	2022-04
10945070	90271.0	399940626	2023-12
4111545	246566.0	79805061	2021-04
5807445	145414.0	184340548	2023-11
3031466	92283.0	93863097	2022-06
15150929	309776.0	479480365	2021-06
7333679	596288.0	99873263	2022-02
7516061	117750.0	148262084	2022-08
7178076	156710.0	131998799	2023-07
10021541	313714.0	186404722	2021-11
7806492	153599.0	243654860	2023-03
6627626	68053.0	230588443	2023-10
14480102	724373.0	128532182	2021-03
7659101	467662.0	66012398	2021-02
10136750	190375.0	308453523	2022-07
5059505	278507.0	122884379	2023-02
+			+

Рисунок 8 – Вывод статистики для «challenge».

Вывод таблицы статистики для бинарного признака «asmr» представлен на рисунке 9.

806898 195182.0 25479054 2020-12 972978 9812.0 21934161 2023-08 372812 3710.0 16948984 2024-02 139422 9782.0 2793246 2021-07 19753 1223.0 197798 2022-04 981541 8305.0 40940254 2023-12 185301 11669.0 3683672 2021-04 3050175 8195.0 97983460 2023-11 145550 7799.0 4030143 2022-06 102886 6443.0 5274694 2021-06 790544 13393.0 16592853 2022-02 109548 4806.0 3001010 2022-08 1376859 13366.0 24240074 2023-07 1143205 30478.0 24852924 2021-11 4612315 28594.0 79240618 2023-03	toktalk@LAI	PTOP-192U10JP: ~/	leti-spark-cou	rse-2024
806898 195182.0 25479054 2020-12 972978 9812.0 21934161 2023-08 372812 3710.0 16948984 2024-02 139422 9782.0 2793246 2021-07 19753 1223.0 197798 2022-04 981541 8305.0 40940254 2023-12 185301 11669.0 3683672 2021-04 3050175 8195.0 97983460 2023-11 145550 7799.0 4030143 2022-06 102886 6443.0 5274694 2021-06 790544 13393.0 16592853 2022-02 109548 4806.0 3001010 2022-08 1376859 13366.0 24240074 2023-07 1143205 30478.0 24852924 2021-11 4612315 28594.0 79240618 2023-03				
806898 195182.0 25479054 2020-12 972978 9812.0 21934161 2023-08 372812 3710.0 16948984 2024-02 139422 9782.0 2793246 2021-07 19753 1223.0 197798 2022-04 981541 8305.0 40940254 2023-12 185301 11669.0 3683672 2021-04 3050175 8195.0 97983460 2023-11 145550 7799.0 4030143 2022-06 102886 6443.0 5274694 2021-06 790544 13393.0 16592853 2022-02 109548 4806.0 3001010 2022-08 1376859 13366.0 24240074 2023-07 1143205 30478.0 24852924 2021-11 4612315 28594.0 79240618 2023-03		L-2		++ \alpha
972978 9812.0 21934161 2023-08 372812 3710.0 16948984 2024-02 139422 9782.0 2793246 2021-07 19753 1223.0 197798 2022-04 981541 8305.0 40940254 2023-12 185301 11669.0 3683672 2021-04 3050175 8195.0 97983460 2023-11 145550 7799.0 4030143 2022-06 102886 6443.0 5274694 2021-06 790544 13393.0 16592853 2022-02 109548 4806.0 3001010 2022-08 1376859 13366.0 24240074 2023-07 1143205 30478.0 24852924 2021-11 4612315 28594.0 79240618 2023-03	bins_likes	bins_comment	bins_view	YM
972978 9812.0 21934161 2023-08 372812 3710.0 16948984 2024-02 139422 9782.0 2793246 2021-07 19753 1223.0 197798 2022-04 981541 8305.0 40940254 2023-12 185301 11669.0 3683672 2021-04 3050175 8195.0 97983460 2023-11 145550 7799.0 4030143 2022-06 102886 6443.0 5274694 2021-06 790544 13393.0 16592853 2022-02 109548 4806.0 3001010 2022-08 1376859 13366.0 24240074 2023-07 1143205 30478.0 24852924 2021-11 4612315 28594.0 79240618 2023-03	806808	105182 0	25479954	2020-12
372812 3710.0 16948984 2024-02 139422 9782.0 2793246 2021-07 19753 1223.0 197798 2022-04 981541 8305.0 40940254 2023-12 185301 11669.0 3683672 2021-04 3050175 8195.0 97983460 2023-11 145550 7799.0 4030143 2022-06 102886 6443.0 5274694 2021-06 790544 13393.0 16592853 2022-02 109548 4806.0 3001010 2022-08 1376859 13366.0 24240074 2023-07 1143205 30478.0 24852924 2021-11 4612315 28594.0 79240618 2023-03			:	: :
139422 9782.0 2793246 2021-07 19753 1223.0 197798 2022-04 981541 8305.0 40940254 2023-12 185301 11669.0 3683672 2021-04 3050175 8195.0 97983460 2023-11 145550 7799.0 4030143 2022-06 102886 6443.0 5274694 2021-06 790544 13393.0 16592853 2022-02 109548 4806.0 3001010 2022-08 1376859 13366.0 24240074 2023-07 1143205 30478.0 24852924 2021-11 4612315 28594.0 79240618 2023-03			!	: :
19753 1223.0 197798 2022-04 981541 8305.0 40940254 2023-12 185301 11669.0 3683672 2021-04 3050175 8195.0 97983460 2023-11 145550 7799.0 4030143 2022-06 102886 6443.0 5274694 2021-06 790544 13393.0 16592853 2022-02 109548 4806.0 3001010 2022-08 1376859 13366.0 24240074 2023-07 1143205 30478.0 24852924 2021-11 4612315 28594.0 79240618 2023-03				
981541 8305.0 40940254 2023-12 185301 11669.0 3683672 2021-04 3050175 8195.0 97983460 2023-11 145550 7799.0 4030143 2022-06 102886 6443.0 5274694 2021-06 790544 13393.0 16592853 2022-02 109548 4806.0 3001010 2022-08 1376859 13366.0 24240074 2023-07 1143205 30478.0 24852924 2021-11 4612315 28594.0 79240618 2023-03			:	: :
185301 11669.0 3683672 2021-04 3050175 8195.0 97983460 2023-11 145550 7799.0 4030143 2022-06 102886 6443.0 5274694 2021-06 790544 13393.0 16592853 2022-02 109548 4806.0 3001010 2022-08 1376859 13366.0 24240074 2023-07 1143205 30478.0 24852924 2021-11 4612315 28594.0 79240618 2023-03			:	: :
3050175 8195.0 97983460 2023-11 145550 7799.0 4030143 2022-06 102886 6443.0 5274694 2021-06 790544 13393.0 16592853 2022-02 109548 4806.0 3001010 2022-08 1376859 13366.0 24240074 2023-07 1143205 30478.0 24852924 2021-11 4612315 28594.0 79240618 2023-03			!	: :
145550 7799.0 4030143 2022-06 102886 6443.0 5274694 2021-06 790544 13393.0 16592853 2022-02 109548 4806.0 3001010 2022-08 1376859 13366.0 24240074 2023-07 1143205 30478.0 24852924 2021-11 4612315 28594.0 79240618 2023-03				: :
102886 6443.0 5274694 2021-06 790544 13393.0 16592853 2022-02 109548 4806.0 3001010 2022-08 1376859 13366.0 24240074 2023-07 1143205 30478.0 24852924 2021-11 4612315 28594.0 79240618 2023-03				: :
790544				:
109548				: :
1376859 13366.0 24240074 2023-07 1143205 30478.0 24852924 2021-11 4612315 28594.0 79240618 2023-03				: :
1143205 30478.0 24852924 2021-11 4612315 28594.0 79240618 2023-03			:	: :
4612315 28594.0 79240618 2023-03				: :
27/16006	3746006		:	: :
3746006 22475.0 150644595 2023-10 171147 14773.0 3023245 2021-03				:
			!	: :
			:	: :
114734 5675.0 2710350 2022-07 7980218 191878.0 111873187 2023-02			:	: :
7980218 191878.0 111873187 2023-02	7980218	1918/8.0	1119/318/	2023-02

Рисунок 9 – Вывод статистики для «asmr».

Вывод таблицы статистики для бинарных признаков представлен на рисунке 10.

toktalk@LAPTOP-I92UI0JP: ~/leti-spark-course-2024 bins_likes bins_comment bins_view YM				
109887	♣ toktalk@LAI	PTOP-192U10JP: ~/l	eti-spark-cou	rse-2024
109887				
109887	+	+		++
360038	bins_likes	bins_comment	bins_view	YM
360038	+	+		++
1026287 1570.0 40799257 2023-11 117149 6387.0 1445523 2022-06 24804 663.0 3677081 2021-06 40748 32.0 2716941 2023-07 780131 401.0 13885801 2021-11 320574 2774.0 23666501 2023-03 156619 143.0 7209875 2023-10 375651 320.0 7203295 2021-10 72129 17989.0 661959 2020-11 84853 485.0 3046824 2023-09 80524 99.0 7407704 2023-04 104972 95.0 6544252 2023-05 47124 193.0 3776090 2023-01 479495 1470.0 29808958 2024-03 105407 281.0 7103180 2023-06 8076 334.0 540599 2020-08 936580 2361.0 45754954 2024-04	109887	1791.0	1413359	2024-02
117149 6387.0 1445523 2022-06 24804 663.0 3677081 2021-06 40748 32.0 2716941 2023-07 780131 401.0 13885801 2021-11 320574 2774.0 23666501 2023-03 156619 143.0 7209875 2023-10 375651 320.0 7203295 2021-10 72129 17989.0 661959 2020-11 84853 485.0 3046824 2023-09 80524 99.0 7407704 2023-04 104972 95.0 6544252 2023-05 47124 193.0 3776090 2023-01 479495 1470.0 29808958 2024-03 105407 281.0 7103180 2023-06 8076 334.0 540599 2020-08 936580 2361.0 45754954 2024-04	360038	2405.0	21008692	2023-12
24804 663.0 3677081 2021-06 40748 32.0 2716941 2023-07 780131 401.0 13885801 2021-11 320574 2774.0 23666501 2023-03 156619 143.0 7209875 2023-10 375651 320.0 7203295 2021-10 72129 17989.0 661959 2020-11 84853 485.0 3046824 2023-09 80524 99.0 7407704 2023-04 104972 95.0 6544252 2023-05 47124 193.0 3776090 2023-01 479495 1470.0 29808958 2024-03 105407 281.0 7103180 2023-06 8076 334.0 540599 2020-08 936580 2361.0 45754954 2024-04	1026287	1570.0	40799257	2023-11
40748 32.0 2716941 2023-07 780131 401.0 13885801 2021-11 320574 2774.0 23666501 2023-03 156619 143.0 7209875 2023-10 375651 320.0 7203295 2021-10 72129 17989.0 661959 2020-11 84853 485.0 3046824 2023-09 80524 99.0 7407704 2023-04 104972 95.0 6544252 2023-05 47124 193.0 3776090 2023-01 479495 1470.0 29808958 2024-03 105407 281.0 7103180 2023-06 8076 334.0 540599 2020-08 936580 2361.0 45754954 2024-04	117149	6387.0	1445523	2022-06
780131 401.0 13885801 2021-11 320574 2774.0 23666501 2023-03 156619 143.0 7209875 2023-10 375651 320.0 7203295 2021-10 72129 17989.0 661959 2020-11 84853 485.0 3046824 2023-09 80524 99.0 7407704 2023-04 104972 95.0 6544252 2023-05 47124 193.0 3776090 2023-01 479495 1470.0 29808958 2024-03 105407 281.0 7103180 2023-06 8076 334.0 540599 2020-08 936580 2361.0 45754954 2024-04	24804	663.0	3677081	2021-06
320574 2774.0 23666501 2023-03 156619 143.0 7209875 2023-10 375651 320.0 7203295 2021-10 72129 17989.0 661959 2020-11 84853 485.0 3046824 2023-09 80524 99.0 7407704 2023-04 104972 95.0 6544252 2023-05 47124 193.0 3776090 2023-01 479495 1470.0 29808958 2024-03 105407 281.0 7103180 2023-06 8076 334.0 540599 2020-08 936580 2361.0 45754954 2024-04	40748	32.0	2716941	2023-07
156619 143.0 7209875 2023-10 375651 320.0 7203295 2021-10 72129 17989.0 661959 2020-11 84853 485.0 3046824 2023-09 80524 99.0 7407704 2023-04 104972 95.0 6544252 2023-05 47124 193.0 3776090 2023-01 479495 1470.0 29808958 2024-03 105407 281.0 7103180 2023-06 8076 334.0 540599 2020-08 936580 2361.0 45754954 2024-04	780131	401.0	13885801	2021-11
375651	320574	2774.0	23666501	2023-03
72129 17989.0 661959 2020-11 84853 485.0 3046824 2023-09 80524 99.0 7407704 2023-04 104972 95.0 6544252 2023-05 47124 193.0 3776090 2023-01 479495 1470.0 29808958 2024-03 105407 281.0 7103180 2023-06 8076 334.0 540599 2020-08 936580 2361.0 45754954 2024-04	156619	143.0	7209875	2023-10
84853 485.0 3046824 2023-09 80524 99.0 7407704 2023-04 104972 95.0 6544252 2023-05 47124 193.0 3776090 2023-01 479495 1470.0 29808958 2024-03 105407 281.0 7103180 2023-06 8076 334.0 540599 2020-08 936580 2361.0 45754954 2024-04	375651	320.0	7203295	2021-10
80524 99.0 7407704 2023-04 104972 95.0 6544252 2023-05 47124 193.0 3776090 2023-01 479495 1470.0 29808958 2024-03 105407 281.0 7103180 2023-06 8076 334.0 540599 2020-08 936580 2361.0 45754954 2024-04	72129	17989.0	661959	2020-11
104972 95.0 6544252 2023-05 47124 193.0 3776090 2023-01 479495 1470.0 29808958 2024-03 105407 281.0 7103180 2023-06 8076 334.0 540599 2020-08 936580 2361.0 45754954 2024-04	84853	485.0	3046824	2023-09
47124 193.0 3776090 2023-01 479495 1470.0 29808958 2024-03 105407 281.0 7103180 2023-06 8076 334.0 540599 2020-08 936580 2361.0 45754954 2024-04	80524	99.0	7407704	2023-04
479495 1470.0 29808958 2024-03 105407 281.0 7103180 2023-06 8076 334.0 540599 2020-08 936580 2361.0 45754954 2024-04	104972	95.0	6544252	2023-05
105407 281.0 7103180 2023-06 8076 334.0 540599 2020-08 936580 2361.0 45754954 2024-04	47124	193.0	3776090	2023-01
8076 334.0 540599 2020-08 936580 2361.0 45754954 2024-04	479495	1470.0	29808958	2024-03
936580 2361.0 45754954 2024-04	105407	281.0	7103180	2023-06
	8076	334.0	540599	2020-08
500451 2286.0 7088648 2021-09 +	936580	2361.0	45754954	2024-04
+	500451	2286.0	7088648	2021-09
	+			++

Рисунок 10 – Вывод статистики для бинарных признаков.

```
9.
        Делаем сборку статистики для бинарных признаков:
view("df total bins stat", """
    SELECT chandas.YM,
        challenge.bins likes AS challenge likes,
        challenge.bins comment AS challenge comment,
        challenge.bins view AS challenge view,
        ROUND
                      ((challenge.bins comment
challenge.bins likes ), 2) AS challenge CLR,
        asmr.bins likes AS asmr likes,
        asmr.bins comment AS asmr comment,
        asmr.bins view AS asmr view,
        ROUND
              ((asmr.bins comment / asmr.bins likes),
2) AS asmr CLR,
        chandas.bins likes AS chandas likes,
        chandas.bins comment AS chandas comment,
        chandas.bins view AS chandas view,
        ROUND
                       ((chandas.bins comment
chandas.bins likes), 2) AS chandas CLR
    FROM df bins stat chandas AS chandas
         df bins stat asmr AS asmr ON chandas.YM =
asmr.YM
          df bins stat challenge AS
    JOIN
                                          challenge
                                                       ON
chandas.YM = challenge.YM
""")
```

Вывод сборки статистики для бинарных признаков представлен на рисунке 11.

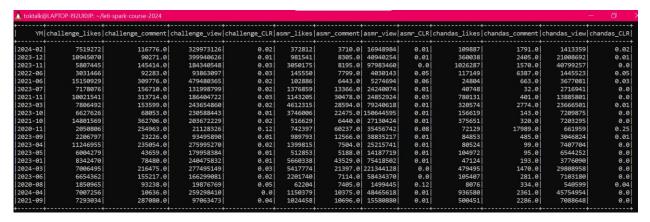


Рисунок 11 – Вывод сборки статистики для бинарных признаков.

10. Делаем сборку общей статистики для бинарных признаков:

```
view("df final bins stat",
WITH sub AS (
    SELECT t.YM,
           t.total comment,
           t.total likes,
           t.total view,
           ROUND ((t.total comment / t.total likes), 2)
AS total CLR,
           b.challenge likes,
           b.challenge comment,
           b.challenge view,
           b.challenge CLR,
           b.asmr likes,
           b.asmr comment,
           b.asmr view,
           b.asmr CLR,
           b.chandas_likes,
           b.chandas comment,
           b.chandas view,
```

b.chandas CLR

```
FROM df total stat AS t
   JOIN df total bins stat AS b
    ON t.YM = b.YM
   ORDER BY t.YM DESC)
SELECT *,
      ROUND((total CLR - challenge CLR),
                                             2)
                                                  AS
diff CLR abs,
      ROUND((total CLR - asmr CLR), 2)
                                                  AS
diff CLR abs,
      ROUND((total CLR - chandas CLR),
                                            2)
                                                  AS
diff CLR abs,
      ROUND((total CLR / challenge CLR),
                                             2)
                                                  AS
diff CLR relative,
      ROUND((total CLR / asmr CLR), 2)
                                                  AS
diff CLR relative,
      ROUND((total CLR / chandas CLR), 2)
                                                  AS
diff CLR relative
FROM sub
""")
```

Вывод общей статистики для бинарных признаков представлен на рисунке 12.



Рисунок 12 – Вывод общей статистики для бинарных признаков.

Выводы.

В работе предоставленные данные были обработаны по категориальным переменным и бинарным признакам, а также посчитан CLR.

В результате были получены 2 таблицы: для категориальных переменных и для бинарных признаков.

Исходный код программы.

```
from pyspark.sql import DataFrame, SparkSession
from pyspark.sql.functions import *
from pyspark.sql.types import *
from pyspark.shell import spark
```

Определяем переменную для названия таблицы CSV-файла

NAME_DATACSV = "RU_youtube_trending_data.csv" #

Название таблицы CSV-файла

def view(name: str, target: DataFrame | str) ->
DataFrame:

"""Создает или заменяет временную представительскую таблицу."""

df = target if isinstance(target, DataFrame) else
spark.sql(target)

df.createOrReplaceTempView(name)

return

Создаем экземпляр SparkSession spark = SparkSession.builder.appName("MyApp").getOrCreate()

Указываем путь к CSV-файлу и читаем его в Dataframe csv_path = "lab3/" + NAME_DATACSV df = spark.read.csv(csv_path, header=True, inferSchema=True, multiLine=True)

```
df.createOrReplaceTempView("YT Data")
```

```
# Бинарный признак (надо выбрать 2 признака из группы):
'Популярные теги
                    (#challenge, #asmr)';
# Категория: 'Трендовые персоны (Илон Маск, 21 Savage,
                                       Карлсен) ';
Магнус
# Коэффициент 'CLR - Процент комментариев от общего
                                          лайков'
числа
    Определяем схему
                        для категорий видео
schema
                                      StructType([
   StructField("id",
                    IntegerType(),
                                          True),
   StructField("name", StringType(),
                                          True),
   StructField("template", StringType(),
                                            True)
])
         Определяем
                     категории
                                           видео
category
                                               Mack",
                              "%Илон
                                       Macκ%"),
          "Илон
   (1,
                              "%21
          "21
                  Savage",
                                       Savage%"),
   (2,
   (3,
         "Магнус Карлсен", "%Магнус
                                       Карлсен%")
1
#
 Создаем DataFrame для категорий видео
categorical df = spark.createDataFrame(category,
schema)
categorical df.createOrReplaceTempView("category")
# ----- создание общей статистики ------
```

```
** ** **
view("df total stat",
    SELECT
         date format(YT data.trending date, 'yyyy-MM')
AS
                                                      YM,
         SUM (YT data.view count)
                                      AS total view,
         SUM(YT data.comment count) AS
                                          total comment,
         SUM(YT data.likes)
                                   AS
                                             total likes,
                     ((SUM(YT data.comment count)
         ROUND
SUM(YT data.likes)),
                            2)
                                      AS
                                                total CLR
    FROM
                   YT Data
                                     AS
                                                  YT data
    GROUP BY date format (YT data.trending date, 'yyyy-
MM')
""")
 ----- создание таблицы со всеми категориями с
изначальной
                     таблицой
                                                      11 11 11
view("df categoryes",
    SELECT
YT data.description, YT data.comment count, YT data.likes
,YT data.view count,YT data.tags,YT data.trending date,
YT data.publishedAt,YT data.title,
           c.id
                               AS
                                                  cat id,
                                AS
                                                 cat name
           c.name
    FROM
                   YT Data
                                     AS
                                                  YT data
    JOIN
                    category
                                         AS
                                                        С
           YT data.description
    ON
                                    ILIKE
                                               c.template
     OR
              YT data.title
                                  ILIKE
                                               c.template
              YT data.tags
                                               c.template
     OR
                                  ILIKE
""")
```

```
----- создание статистики для категорий ---
                                                     11 11 11
view("df cats stat",
    SELECT
         date format(a.trending date, 'yyyy-MM') AS YM,
         SUM(a.view count)
                                   AS
                                              cats view,
         SUM(a.comment count) AS cats comment,
         SUM(a.likes)
                               AS
                                            cats likes,
         a.cat id,
         a.cat name
                                          AS
   FROM
                  df categoryes
                                                       а
                 ΒY
                         a.cat id,
   GROUP
                                            a.cat name,
                                              'yyyy-MM')
date format(a.trending date,
""")
# ----- сборка таблицы для категорий ------
                                                     11 11 11
view("df final cats stat",
   WITH
                                       AS
                      sub
                                                       (
   SELECT
                                                   t.YM,
         t.total comment,
         t.total likes,
         t.total view,
         t.total CLR,
         c.cat name,
         c.cats view,
         c.cats comment,
         c.cats likes,
         ROUND ((c.cats comment / c.cats likes), 2) AS
```

```
cats CLR
   FROM
                  df total stat
                                         AS
                                                      t
   JOIN
                  df cats stat
                                         AS
                                                      С
                   t.YM
   ON
                                                   c.YM
             ΒY
   ORDER
                    t.YM
                             DESC, c.cat id
                                                    ASC
    )
    SELECT
                                                     *,
         ROUND((total CLR - cats CLR),
                                               2)
                                                     AS
diff CLR abs,
         CASE
                         cats CLR
                                                      0
           WHEN
                                          ! =
                   ROUND((total CLR / cats CLR),
              THEN
                                                     2)
              ELSE
                                                   NULL
                                      diff CLR relative
              END
                           AS
    FROM
                                                    sub
    """)
    Показываем содержимое конечного DataFrame
                                                      С
категориями
                                 df total stat").show()
spark.sql("SELECT
                         FROM
spark.sql("SELECT
                    *
                         FROM
                                df categoryes").show()
spark.sql("SELECT
                 *
                                 df cats stat").show()
                          FROM
                   * FROM df final cats stat").show()
spark.sql("SELECT
# Бинарный признак (надо выбрать 2 признака из группы):
'Популярные
                 теги
                           (#challenge,
                                               #asmr)';
# ----- создание таблицы с бинарными признаками с
изначальной
                     таблицы
                                                    11 11 11
view ("df bins challenge",
```

SELECT

YT_data.description, YT_data.comment_count, YT_data.likes
,YT_data.view_count, YT_data.tags, YT_data.trending_date,
YT_data.publishedAt, YT_data.title

FROM YT Data YT data AS '%#challenge%' WHERE (YT data.description ILIKE OR YT data.tags '%challenge%') ILIKE '%#asmr%' (YT data.description AND NOT ILIKE YT data.tags '%asmr%') OR NOTILIKE

11 11 11

11 11 11

""")

view("df_bins_asmr",

SELECT

YT_data.description, YT_data.comment_count, YT_data.likes
,YT_data.view_count, YT_data.tags, YT_data.trending_date,
YT_data.publishedAt, YT_data.title

YT Data FROM AS YT data '%#asmr%' (YT data.description WHERE ILIKE '%asmr%') YT data.tags OR ILIKE AND (YT data.description NOT ILIKE '%#challenge%'

OR YT_data.tags NOT ILIKE '%challenge%')

view("df_bins_chandas",

SELECT

YT_data.description, YT_data.comment_count, YT_data.likes
,YT_data.view_count, YT_data.tags, YT_data.trending_date,
YT_data.publishedAt, YT_data.title

```
(YT data.description ILIKE
                                         '%#challenge%'
   WHERE
              YT data.tags ILIKE
                                         '%challenge%')
       OR
              (YT data.description ILIKE
                                             '%#asmr%'
       AND
                YT data.tags
                                             '%asmr%')
       OR
                                  ILIKE
""")
                                               бинарных
    ----- создание
                            статистики
                                         ДЛЯ
признаков
                                                    11 11 11
view ("df bins stat challenge",
   SELECT
                SUM(a.likes)
                                  AS
                                            bins likes,
        SUM(a.comment count)
                                 AS
                                          bins comment,
        SUM(a.view count)
                                             bins view,
                                  AS
        date format(a.trending date, 'yyyy-MM') AS YM
   FROM
                df bins challenge
                                          AS
          BY date format(a.trending_date, 'yyyy-MM')
   GROUP
""")
                                                    11 11 11
view("df bins stat asmr",
   SELECT
                SUM(a.likes)
                                  AS
                                            bins likes,
        SUM (a.comment count) AS
                                         bins comment,
        SUM(a.view count)
                                             bins view,
                                  AS
        date format(a.trending date, 'yyyy-MM')AS
   FROM
                  df bins asmr
                                         AS
                                                      а
   GROUP BY date format(a.trending date, 'yyyy-MM')
""")
                                                    11 11 11
view("df bins stat chandas",
                                            bins likes,
   SELECT
                SUM(a.likes)
                                   AS
```

YT Data

AS

YT data

FROM

```
SUM(a.comment count)
                                    AS
                                            bins comment,
         SUM(a.view count)
                                    AS
                                               bins view,
         date format(a.trending date, 'yyyy-MM')AS
                  df bins chandas
    FROM
                                            AS
                                                         а
                date format(a.trending date, 'yyyy-MM')
    GROUP
           ΒY
""")
                     Сборка
                                                  бинарных
                              статистики
                                            ДЛЯ
признаков
                                                       11 11 11
view ("df total bins stat",
    SELECT
                                              chandas.YM,
        challenge.bins likes
                                         challenge likes,
                                  AS
        challenge.bins comment
                                       challenge comment,
                                  AS
        challenge.bins view
                                          challenge view,
                                  AS
        ROUND
                       ((challenge.bins comment
challenge.bins likes ),
                               2)
                                     AS
                                           challenge CLR,
        asmr.bins likes
                                  AS
                                              asmr likes,
        asmr.bins comment
                                            asmr comment,
                                  AS
        asmr.bins view
                                  AS
                                               asmr view,
               ((asmr.bins comment / asmr.bins likes),
        ROUND
2)
                        AS
                                                 asmr CLR,
        chandas.bins likes
                                           chandas likes,
                                  AS
        chandas.bins comment
                                  AS
                                         chandas comment,
        chandas.bins view
                                            chandas view,
                                  AS
                        ((chandas.bins comment
        ROUND
```

```
df bins stat chandas
                                        AS
    FROM
                                                 chandas
          df bins stat asmr AS asmr
    JOIN
                                       ON chandas.YM =
asmr.YM
           df bins stat challenge AS challenge
    JOIN
                                                      ON
chandas.YM
                                            challenge.YM
""")
  ----- Сборка общей статистики для бинарных
признаков
                                                      11 11 11
view ("df final bins stat",
WITH
                   sub
                                      AS
                                                        (
    SELECT
                                                   t.YM,
           t.total comment,
           t.total likes,
           t.total view,
           ROUND ((t.total comment / t.total likes), 2)
AS
                                              total CLR,
           b.challenge likes,
           b.challenge comment,
           b.challenge view,
           b.challenge CLR,
           b.asmr likes,
           b.asmr comment,
           b.asmr view,
           b.asmr CLR,
           b.chandas likes,
           b.chandas comment,
           b.chandas view,
```

2)

AS

chandas CLR

chandas.bins likes),

b.chandas CLR

```
df total stat
    FROM
                                           AS
                                                        t
                 df total bins stat
    JOIN
                                            AS
                                                       b
                    t.YM
                                                    b.YM
     ON
    ORDER
                     ΒY
                                                    DESC)
                                   t. YM
SELECT
                                                       *,
       ROUND((total CLR - challenge CLR),
                                                 2)
                                                       AS
diff CLR abs ch,
       ROUND ((total CLR
                                 asmr CLR),
                                                2)
                                                       AS
diff CLR abs as,
       ROUND((total CLR
                                chandas CLR),
                                                 2)
                                                       AS
diff CLR abs chas,
       ROUND((total CLR
                          / challenge CLR),
                                                 2)
                                                       AS
diff CLR relative ch,
       ROUND((total CLR /
                                 asmr CLR),
                                                2)
                                                       AS
diff CLR relative as,
       ROUND((total CLR
                         /
                                chandas CLR),
                                                 2)
                                                       AS
diff CLR relative chas
FROM
                                                      sub
""")
# Показываем содержимое конечного DataFrame с бинарными
```

Показываем содержимое конечного DataFrame с бинарными признаками

```
spark.sql("SELECT * FROM df_bins_challenge").show()
spark.sql("SELECT * FROM df_bins_asmr").show()
spark.sql("SELECT * FROM df_bins_chandas").show()
spark.sql("SELECT * FROM
df_bins_stat_challenge").show()
spark.sql("SELECT * FROM df bins stat asmr").show()
```

```
spark.sql("SELECT * FROM df_bins_stat_chandas").show()
spark.sql("SELECT * FROM df_total_bins_stat").show()
spark.sql("SELECT * FROM df_final_bins_stat").show()

# Завершаем работу с SparkSession
spark.stop()
```