

Recommender system test

В наведенном распределении результатов A/B теста. Регион проведения теста - Евросоюз. Задача - оценить корректность его проведения и выявить возможные недостатки.

Содержание

- Предобработка данных
- Исследовательский анализ данных

- A/B тест
- Выводы

In [1]:

```
import pandas as pd
from plotly import graph_objects as go
from plotly.subplots import make_subplots
from matplotlib import pyplot as plt
import seaborn as sns
from warnings import filterwarnings
filterwarnings("ignore")
import numpy as np
import math
from scipy import stats as st
```

Предобработка данных

In [2]:

```
# Распаковка промоакции
schedule = pd.read_csv('ab_project_marketing_events.csv')
display(schedule)
```

	name	regions	start_dt	finish_dt
0	Christmas&New Year Promo	EU, N.America	2020-12-25	2021-01-03
1	St. Valentine's Day Giveaway	EU, CIS, APAC, N.America	2020-02-14	2020-02-16
2	St. Patrick's Day Promo	EU, N.America	2020-03-17	2020-03-19
3	4th of July Promo	EU, CIS, APAC, N.America	2020-04-12	2020-04-19
4	Black Friday Ads Campaign	EU, CIS, APAC, N.America	2020-11-26	2020-12-01
5	Chinese New Year Promo	APAC	2020-01-25	2020-02-07
6	Labor Day (May 1st) Ads Campaign	EU, CIS, APAC	2020-05-01	2020-05-03
7	International Women's Day Promo	EU, CIS, APAC	2020-03-08	2020-03-10
8	Victory Day CIS (May 8th) Event	CIS	2020-05-09	2020-05-11
9	CIS New Year Gift Lottery	CIS	2020-12-30	2021-01-07
10	Dragon Boat Festival Giveaway	APAC	2020-06-25	2020-07-01
11	Single's Day Gift Promo	APAC	2020-11-11	2020-11-12
12	Chinese Moon Festival	APAC	2020-10-01	2020-10-07

In [3]:

```
# Приведем даты к корректному типу данных
schedule.start_dt = pd.to_datetime(schedule.start_dt)
schedule.finish_dt = pd.to_datetime(schedule.finish_dt)
print(schedule.info())
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 14 entries, 0 to 13
Data columns (total 4 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   name        14 non-null     object
 1   regions     14 non-null     object
 2   start_dt    14 non-null     datetime64[ns]
 3   finish_dt   14 non-null     datetime64[ns]
dtypes: datetime64[ns](2), object(2)
memory usage: 176.0+ bytes
None
```

In [4]:

```
# Выберем интересующий нас период
schedule = schedule.sort_values(by=['start_dt'], reset_index(drop=True))
actual_schedule = schedule.query('finish_dt >= "2020-01-04" and start_dt <= "2021-01-04"')
display(actual_schedule)
```

	name	regions	start_dt	finish_dt
12	Christmas&New Year Promo	EU, N.America	2020-12-25	2021-01-03
13	CIS New Year Gift Lottery	CIS	2020-12-30	2021-01-07

На интересующий нас период выглядит рождественско-новогодняя промоакция с 25 декабря по 3 января

In [5]:

```
# Новые пользователи
new_users = pd.read_csv('final_ab_new_users.csv')
display(new_users.head())
```

	user_id	first_date	region	device
0	D72A72121175D8BE	2020-12-07	EU	PC
1	F10696190F6E66	2020-12-07	N.America	Android
2	2E1BF8F4DC3EAD0F	2020-12-07	EU	PC
3	50734A223C0C83768	2020-12-07	EU	iPhone
4	E1BDDCE0DAFA2679	2020-12-07	N.America	iPhone

In [6]:

```
# Приведем даты к корректному типу
new_users.first_date = pd.to_datetime(new_users.first_date)
print(new_users.info())
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 61733 entries, 0 to 61732
Data columns (total 4 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   user_id     61733 non-null  object
 1   first_date  61733 non-null  datetime64[ns]
 2   region      61733 non-null  object
 3   device      61733 non-null  object
dtypes: datetime64[ns](1), object(3)
memory usage: 1.9+ MB
None
61733
```

In [7]:

```
# Есть ли дубликаты?
print(new_users.duplicated().sum())
print(new_users.region.unique())
```

```
0
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
```

Дубликатов нет. Интересующий нас Евросоюз обозначен только одним способом.

In [8]:

```
# Новые события
new_events = pd.read_csv('final_ab_events.csv')
display(new_events.head())
```

	user_id	event_dt	event_name	details
0	E1BDDCE0DAFA2679	2020-12-07 09:22:53	purchase	99.99
1	7B6452F0B61FA904	2020-12-07 09:22:53	purchase	9.99
2	9CD9F34540FD254C	2020-12-07 12:59:29	purchase	4.99
3	96F27A054B191457	2020-12-07 04:02:40	purchase	4.99
4	1FD7680FDF94CAF1F	2020-12-07 10:15:09	purchase	4.99

In [9]:

```
# Приведем даты к корректному типу
new_events.event_dt = pd.to_datetime(new_events.event_dt)
print(new_events.info())
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 62417 entries, 0 to 62416
Data columns (total 4 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   user_id     62417 non-null  object
 1   event_dt    62417 non-null  datetime64[ns]
 2   event_name  62417 non-null  object
 3   details     62410 non-null  float64
dtypes: datetime64[ns](1), float64(1), object(2)
memory usage: 13.4+ MB
None
```

In [10]:

```
# Есть ли дубликаты?
print(new_events.duplicated().sum())
print(new_events.details.value_counts())
```

```
purchase    62410
income      62410
None
(62410, 4)
```

Дубликатов нет. Число покупок равно числу "деталей". Значит, переименуем "детали" в "доход". Пропуски оставим, здесь они не мешают.

In [11]:

```
new_events = new_events.rename(columns = {'details': 'income'})
```

In [12]:

```
# Тестовые тестовые
test_groups = pd.read_csv('final_ab_participants.csv')
display(test_groups.head())
```

	user_id	group	ab_test
0	D1AB43E28B7B6A73	A	recommender_system_test
1	A74364BD6242119	A	recommender_system_test
2	D4BC14FDDAFD028E	A	recommender_system_test
3	04986C5DF18632E	A	recommender_system_test
4	482F14783456021B	B	recommender_system_test

In [13]:

```
print(test_groups.info())
print()
print(test_groups.duplicated().sum())
print()
print(test_groups.ab_test.unique())
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 18268 entries, 0 to 18267
Data columns (total 4 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   user_id     18268 non-null  object
 1   group       18268 non-null  object
 2   ab_test     18268 non-null  object
dtypes: object(3)
memory usage: 428.3+ KB
None
0
```

Дубликатов нет, пропусков нет. Тестов - два вида.

Исследовательский анализ данных

Корректность набора участников теста

In [14]:

```
# Сравним сред по количеству нас тесту
recommender_system_test = test_groups.query('ab_test == "recommender_system_test"').reset_index(drop=True)
eu_users = new_users.query('region == "EU"').reset_index(drop=True)
```

```
# Проверим, есть ли совпадения между двумя наборами данных
eu_users_tests = eu_users.merge(recommender_system_test, on='user_id', how='left')
eu_users_tests.dropna(inplace=True)
eu_users_tests.reset_index(drop=True, inplace=True)
eu_users_tests.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 6701 entries, 0 to 6700
Data columns (total 4 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   user_id     6701 non-null  object
 1   group       6701 non-null  object
 2   ab_test     6701 non-null  object
dtypes: object(3)
memory usage: 13.4+ MB
None
6701
```

Пользователей, попавших и в контрольную и в тестовую группы, нет: число уникальных пользователей в таблице recommender_system_test равно числу строк в ней.

In [15]:

```
# Создадим сред пользователей из ЕС
eu_users = new_users.query('region == "EU"').reset_index(drop=True)
eu_users_tests = eu_users_tests.reset_index(drop=True)
eu_users_tests.reset_index(drop=True, inplace=True)
eu_users_tests.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 6701 entries, 0 to 6700
Data columns (total 4 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   user_id     6701 non-null  object
 1   group       6701 non-null  object
 2   ab_test     6701 non-null  object
dtypes: object(3)
memory usage: 13.4+ MB
None
6701
```

Несколько сотен участников теста были и не из ЕС. Теперь остались только нужные, их чуть больше необходимых 6000

In [16]:

```
# Проверим, есть ли общие пользователи в параллельном тесте
interface_eu_test = test_groups.query('ab_test == "interface_eu_test"').reset_index(drop=True)
total = 0
for i in eu_users_tests.user_id.unique():
    if i in interface_eu_test.user_id.unique():
        total += 1
print(total)
```

1602

1602 пользователя принимают участие в обоих тестах. То есть примерно каждый четвертый пользователь из нашего теста также есть в другом, это ставит под сомнение корректность проведения теста, потому что участие в другом тесте может влиять на поведение этих пользователей.

Продуктовая воронка

Для воронки нам понадобится вычислить, сколько уникальных пользователей приходится на события каждого типа в каждой группе, а также общее число пользователей, включая, не совершивших ни одного действия - тоже для обеих групп.

In [17]:

```
eu_users_tests_a = eu_users_tests.query('group == "A"')
eu_users_tests_b = eu_users_tests.query('group == "B"')
eu_users_tests_a.reset_index(drop=True, inplace=True)
eu_users_tests_b.reset_index(drop=True, inplace=True)
eu_users_tests_a.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 3634 entries, 0 to 3633
Data columns (total 4 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   user_id     3634 non-null  object
 1   group       3634 non-null  object
 2   ab_test     3634 non-null  object
dtypes: object(3)
memory usage: 13.4+ MB
None
3634
2717
```

In [18]:

```
# Теперь найдем пользователей, которым совершить хотя бы одно действие
eu_users_and_events = pd.merge(eu_users, new_events, on='user_id')
eu_users_and_events.reset_index(drop=True, inplace=True)
eu_users_and_events.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 6701 entries, 0 to 6700
Data columns (total 4 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   user_id     6701 non-null  object
 1   group       6701 non-null  object
 2   ab_test     6701 non-null  object
dtypes: object(3)
memory usage: 13.4+ MB
None
6701
```

	user_id	first_date	device	event_dt	event_name	income	group
0	D72A72121175D8BE	2020-12-07	PC	2020-12-07 21:52:07	product_page	NaN	A
1	D72A72121175D8BE	2020-12-07	PC	2020-12-07 21:52:07	login	NaN	A
2	DD4352CDDFC8CD57	2020-12-07	Android	2020-12-07 15:32:54	product_page	NaN	B
3	DD4352CDDFC8CD57	2020-12-07	Android	2020-12-08 08:29:31	product_page	NaN	B
4	DD4352CDDFC8CD57	2020-12-07	Android	2020-12-10 18:18:27	product_page	NaN	B

In [19]:

```
# В обеих группах посчитаем число уникальных пользователей для каждого типа событий.
eu_users_tests_a.reset_index(drop=True, inplace=True)
eu_users_tests_b.reset_index(drop=True, inplace=True)
eu_users_tests_a.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 3634 entries, 0 to 3633
Data columns (total 4 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   user_id     3634 non-null  object
 1   group       3634 non-null  object
 2   ab_test     3634 non-null  object
dtypes: object(3)
memory usage: 13.4+ MB
None
3634
2717
```

In [20]:

```
fig = make_subplots(rows=1, cols=2)
fig.add_trace(go.Funnel(name="Group A", y = funnel_a.event_name, x=funnel_a.num_users, row=1, col=1))
fig.add_trace(go.Funnel(name="Group B", y = funnel_b.event_name, x=funnel_b.num_users, row=1, col=2))
fig.show()
```



In [21]:

```
# Число пользователей, совершивших хотя бы один шаг, совпадает с числом логинов (см. воронку)
eu_users_tests_a.reset_index(drop=True, inplace=True)
eu_users_tests_b.reset_index(drop=True, inplace=True)
eu_users_tests_a.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 3634 entries, 0 to 3633
Data columns (total 4 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   user_id     3634 non-null  object
 1   group       3634 non-null  object
 2   ab_test     3634 non-null  object
dtypes: object(3)
memory usage: 13.4+ MB
None
3634
2717
```

Очевидно, обязательным действием, без которого не возможны другие шаги, является вход (login). И действительно, число логинов совпадает с количеством уникальных пользователей в таблице пользователей, совершивших хотя бы одно действие. Из этого можно заключить, что 1840 пользователей из группы В не совершили ни одного действия, даже не зашли на сайт. Из-за этого 1840 случаев переход к следующим шагам могут быть суммами шагов.

Далше, при переходе к следующим шагам могут быть суммами. Например, переход в корзину - не обязательный шаг, судя по тому, что переходов в корзину меньше, чем совершивших покупку.

Тестовая группа демонстрирует ярко выраженную отрицательную тенденцию по сравнению с контрольной группой. Конверсия в покупку 9% против 28% в группе А. Выход на тестовую страницу - 46 % в контрольной группе, в тестовой - 18%. Метрики ухудшились в два с половиной раза.

Распределение событий по пользователям

In [22]:

```
# Количество событий по пользователям
eu_events_by_user = eu_events.groupby('user_id', as_index=False).agg({'event_dt': 'count'})
eu_events_by_user.reset_index(drop=True, inplace=True)
eu_events_by_user.info()
```

	user_id	num_events
0	0001719F4D0B1D1B	6
1	00019F18E7A4E568	6
2	000249B6327275C7	9
3	0002CE61FF2C4011	12
4	000456437D0E7F1E	4

In [23]:

```
# Объединим таблицу всех участников теста с таблицей количества событий на пользователя,
eu_events_by_user.reset_index(drop=True, inplace=True)
eu_events_by_user.info()
```

group	A	B
event_day		
2020-12-07	318.0	356.0
2020-12-08	313.0	238.0
2020-12-09	371.0	338.0
2020-12-10	331.0	249.0


```
In [37]: # Просмотр корзины

# Нулевая гипотеза: доли просмотревших корзину в обеих выборках равны.
# Альтернативная гипотеза: между долями есть статистически значимая разница.

print("Гипотеза о равенстве долей просмотревших корзину:")
print()
calculate_som_e_stats(funnel_a['num_users'][0], funnel_a['num_users'][4], funnel_b['num_users'][0], funnel_b['num_users'][4])
```

Гипотеза о равенстве долей просмотревших корзину:

p-значение 0.0

Отвергаем нулевую гипотезу, разница между долями статистически значима

```
In [38]: # Совершившие покупку

# Нулевая гипотеза: доли совершивших покупку в обеих выборках равны.
# Альтернативная гипотеза: между долями есть статистически значимая разница.

print("Гипотеза о равенстве долей совершивших покупку:")
print()
calculate_som_e_stats(funnel_a['num_users'][0], funnel_a['num_users'][3], funnel_b['num_users'][0], funnel_b['num_users'][3])

Гипотеза о равенстве долей совершивших покупку:
```

p-значение 0.0

Отвергаем нулевую гипотезу, разница между долями статистически значима

Различия между выборками по всем трем показателям статистически значимы, но при этом мы получили отрицательный результат: метрики в группе B значительно хуже, чем в группе A.

Выводы

Исследовательский анализ данных показал, и A/B тест подтвердил, что внедрение новой рекомендательной системы не пошло на пользу ключевым метрикам. Конверсия сократилась по всем показателям в два с половиной раза.

В то же время корректность проведения теста вызывает сомнения: во-первых, из-за 1602 общих участников с параллельным тестом, во-вторых из-за влияния рождественской промоакции, в-третьих, из-за прекращения регистрации новых участников 21-го декабря.

```
In [ ]:
```