

CTR prediction

- [План](#)
 - [1. xgb](#)
 - [2. DeepCTR](#)
 - [3. xlearn](#)
 - [4. Ансамбль \(самая интересная часть\)](#)
-

План

Изначально планировал получить три решения:

1. xgb на всем датасете с hyperopt
 2. Модель из [DeepCTR](#)
 3. Модель из [xlearn](#)
-

1. xgb

Ключевые моменты:

- пробовал делать one hot encoding для всех категориальных фичей, и это даже считалось, но одна итерация обучения бустинга занимала около 30 минут, что очень долго для подбора параметров
- поэтому решил использовать другой подход для категориальных фичей: а именно [catboost encoder](#) из [categorical encoding](#) (правда этот пакет не умеет работать с ruspark, поэтому пришлось сначала сделать все в pandas Dataframe, а затем перевести в spark — что могло бы являться проблемой, если данных было бы сильно больше)
- Для подбора параметров:
 - не просто `hp.choice`, так как зачастую есть порядок в гиперпараметрах (`num_round`, `eta`, ...), а `hp.quniform`, `hp.uniform` ([hyperopt parameter expresion](#))
 - решил подбирать все параметры *сразу*, а не поочередно

- Общее пространство поиска параметров:

```
space = {
    # Optimize
    'num_round': hp.quniform('num_round', 10, 100, 10),
    'eta': hp.uniform('eta', 0.025, 0.8),

    'max_depth': hp.quniform('max_depth', 3, 9, 1),
    'min_child_weight': hp.quniform('min_child_weight', 0, 100, 15),

    'gamma': hp.loguniform('gamma', -3, 0),

    'subsample': hp.uniform('subsample', 0.3, 1.),
    'colsample_bytree': hp.uniform('colsample_bytree', 0.3, 1.),

    **static_params
}
```

- выставил 30 попыток, что, возможно, маловато
- поиск занял 1 час 20 минут
- лучшие параметры нашлись при последней (*sic!*) попытке
- Подбор параметров дал **1.4** прироста *ate* по *roc-auc* относительно модели со стандартными параметрами
- На kaggle с этим решением получил следующий результат:

full_dataset_mte_xgb.csv	0.77948	0.77880	✓
9 days ago by Yury Belousov			
full train, mte, xgb, tuned			

- (возможно), если бы отправил нетюненную модель, то все равно стал бы вторым
- Обучение на всем датасете дало **0.045** прироста по *roc-auc* в сравнении с обучением на случайной половине с подбором параметров (разница с моделью на половине без подбора - **0.056**), что достаточно много, и в очередной раз показывает, что обучение на подвыборке не самый лучший подход.
- код доступен в [01.ctr_xgb.ipynb](#)

2. DeepCTR

- Понравилась библиотека, достаточно легко с ней работать.
- Использовал самую стандартную модель [DeepFM](#), легко завести и работать (единственное что tensorflow почему-то не может строить граф с переменными, которые содержат нижнее подчеркивание, поэтому нужно обязательно переименовать колонки: `train.rename(columns=lambda x: x.replace('_', ''))`).
- и нужно обязательно dense фичи скастить к `np.float32`, а sparse к `string`
- Для категориальных фичей использовал [DeepCTR feature hashing](#)
- На kaggle вполне хороший результат:

deepfm_full_dataset.csv	0.77862	0.77674	✓
9 days ago by Yury Belousov			
full train, hash, deepfm			

- И это при том, что я не особо экспериментировал и подбирал параметры
- код доступен в [02.1.ctr_deepctr_deepfm.ipynb](#)
- После окончания конкурса попробовал другую модель: [FGCNN](#), и она тоже показала себя неплохо (правда обучалась в 4 раза дольше):

FGCNN_full_dataset.csv	0.76246	0.76169	<input type="checkbox"/>
9 days ago by Yury Belousov			
full train, hash, FGCNN			

- код доступен в [02.2.ctr_deepctr_fgcn.ipynb](#)

3. xlearn

- Захотел использовать алгоритм FFM, который использовался в выигрышном решении оригинального соревнования. Но оригинальная реализация ([LIBFFM](#)) достаточно древняя, поэтому решил использовать более новую и быструю (по заявлению авторов) реализацию [xlearn](#)
- Однако этот интерфейс я тоже не до конца понял, а именно что за [libffm](#) формат (конкретно как составляется [field](#)).
- Поэтому использовал просто factorization machine (FM), так как нашел понятный [FM пример](#) с использованием scikit learn api.
- Но результат совсем не очень:

FMM_full_dataset.csv	0.65948	0.65786	<input type="checkbox"/>
9 days ago by Yury Belousov			
full train, FMM			

- код доступен в [03.ctr_xl_fm.ipynb](#)

4. Ансамбль (самая интересная часть)

- Как известно, зачастую в соревнованиях побеждают комбинации различных алгоритмов (а формально — теорема Кондорсе о присяжных).
- Поэтому я захотел объединить результаты первого и второго подхода.
- Быстро написав и отправив код, я очень удивился:

combined.csv	0.70540	0.70417	<input type="checkbox"/>
9 days ago by Yury Belousov			
mean of deepfm and xgb			
deepfm_full_dataset.csv	0.77862	0.77674	✓
9 days ago by Yury Belousov			
full train, hash, deepfm			
full_dataset_mte_xgb.csv	0.77948	0.77880	✓
9 days ago by Yury Belousov			
full train, mte, xgb, tuned			

- Так как результат объединения **сильно** хуже каждого из подхода по отдельности.
- Это была последняя попытка в последний день, поэтому я ничего не мог уже исправить

- Лишь на следующий день, когда я посмотрел на результаты первых двух подходов:

deepfm_full_dataset.csv	full_dataset_mte_xgb.csv
1 id,proba	1 id,proba
2 566935904713,0.4102064	2 566936179499,0.15740602
3 566935904715,0.3677692	3 584116096827,0.044796906
4 566935904727,0.40551883	4 575525985092,0.102862604
5 566935904737,0.5109164	5 575525778555,0.5752256
6 566935904741,0.38027614	6 584115725080,0.3822638
7 566935904745,0.08481385	7 566936098987,0.40038332
8 566935904752,0.32745838	8 575526158169,0.34860474
9 566935904759,0.012269236	9 584115768000,0.3369564
10 566935904767,0.5162818	10 575525703619,0.6717163
11 566935904768,0.077426665	11 575525995952,0.117264174
12 566935904778,0.7996467	12 566935961087,0.13080785
13 566935904785,0.06712084	13 575525855671,0.4796365
14 566935904789,0.7352005	14 584115565705,0.08831968
15 566935904804,0.25856566	15 584115652121,0.55467355
16 566935904810,0.4860264	16 584115683934,0.08740573
17 566935904817,0.10084507	17 575525681971,0.112382025
18 566935904819,0.28937745	18 584115674185,0.30997187
19 566935904830,0.4274971	19 592705576855,0.4577617
20 566935904847,0.27776772	20 575526053659,0.052949328
21 566935904848,0.24683109	21 592705495073,0.4116129
22 566935904849,0.04271594	22 584115934718,0.14455374
23 566935904853,0.24116129	23 584116016277,0.31964454
24 566935904867,0.13399896	24 566936160190,0.18716021

я заметил, что в csv от русpark `id` идут группами по 5 (возможно это количество партиций).

- а вот как у меня было написано объединение:

```
probs = []

submission_filenames = ['deepfm_full_dataset.csv', 'full_dataset_mte_xgb.csv']

for submission_filename in submission_filenames:
    submission = pd.read_csv(submission_filename)
    probs.append(submission['proba'])

proba = np.array(probs).mean(axis=0)
```

- то есть у меня никак не учитывался порядок `id` 😞😞😞
- Для исправления достаточно было добавить `sort_values` или, ещё лучше, написать через `merge`.
- В итоге, исправив, я получил следующий результат:

Submission and Description	Private Score	Public Score	Use for Final Score
combined.csv 8 days ago by Yury Belousov deemfm & xgb correct combining	0.78771	0.78639	<input type="checkbox"/>

Что почти на `0.01` больше чем у первого места 😞

- неправильный код доступен в [04.CombineSolutions.ipynb](#)