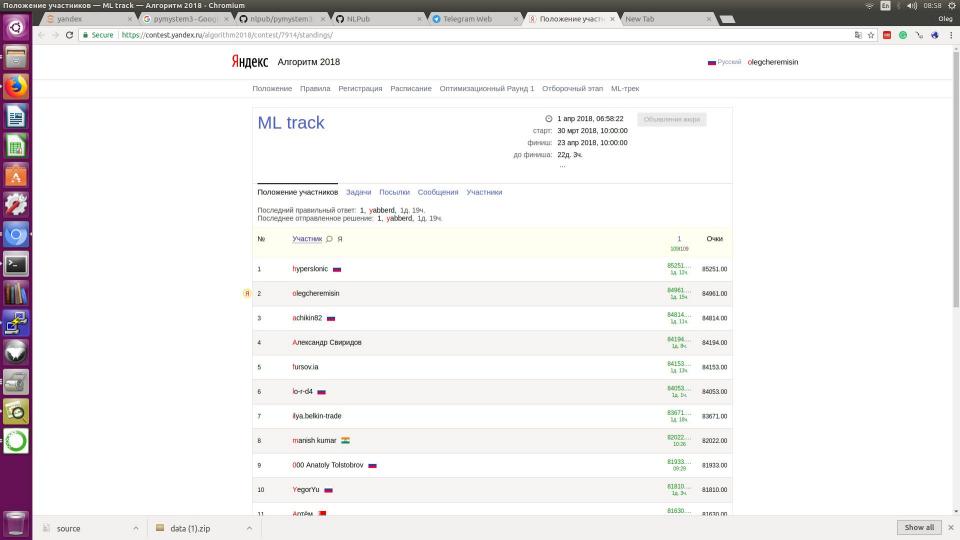
Яндекс Алгоритм 2018 ML track

Побывать в топе пока не подтянулись те, у кого модель училась две недели :)



Ooops

```
In [1]: import pandas as pd
        import numpy as np
In [3]: train = pd.read csv('datasets/yandex/train.tsv', sep='\t', header=None)
                    def read(self, nrows=None):
           1746
           1747
                        trv:
                            data = self. reader.read(nrows)
        -> 1748
                        except StopIteration:
           1749
                            if self. first chunk:
           1750
        pandas/ libs/parsers.pyx in pandas. libs.parsers.TextReader.read (pandas/ libs/parsers.c:10862)()
        pandas/ libs/parsers.pyx in pandas. libs.parsers.TextReader. read low memory (pandas/ libs/parsers.c:11138)()
        pandas/ libs/parsers.pyx in pandas. libs.parsers.TextReader. read rows (pandas/ libs/parsers.c:11884)()
        pandas/ libs/parsers.pyx in pandas. libs.parsers.TextReader. tokenize rows (pandas/ libs/parsers.c:11755)()
        pandas/ libs/parsers.pyx in pandas. libs.parsers.raise parser error (pandas/ libs/parsers.c:28765)()
        ParserError: Error tokenizing data. C error: Expected 8 fields in line 134, saw 10
```

Данные

Out[12]:									
000[12].	0	1		2	3	4	5	6	7
	0 22579918886	кликни на меня а потом на надпись " видео - зв	о, я те	ебя вижу .	ладно , повесь трубку .	0	не могу .	good	0.875352
	1 22579918886	кликни на меня а потом на надпись " видео - зв	о, я те	ебя вижу .	ладно , повесь трубку .	1	нет , звонить буду я .	neutral	0.900968
	2 22579918886	кликни на меня а потом на надпись " видео - зв	о,яте	ебя вижу .	ладно , повесь трубку .	2	слушай , я не мог уйти .	bad	0.884320
	3 ₂₂₅₇₉₉₁₈₈₈₆	кликни на меня а потом на надпись " видео - зв	о, я те	ебя вижу .	ладно , повесь трубку .	3	я не прекращу звон <mark>ить</mark> .	good	0.982530
	4 22579918886	кликни на меня а потом на надпись " видео - зв	о,яте	ебя вижу .	ладно , повесь трубку .	4	я звоню им .	good	0.838054
	5 22579918886	кликни на меня а потом на надпись " видео - зв	о, я те	ебя вижу .	ладно , повесь трубку .	5	просто повесь трубку .	bad	0.955718
	6 50117448291	бывало и получше .	слушайте , мы с женой . забл	совсем я пы удились .	ытаюсь добраться до юджина , но , кажется ,	0	едете ?	bad	0.909115
in [13]:	test.head(7)								
Out[13]:	0		1	2	1		3 4		5
	0 138920940977	знаешь , я иногда подумываю , что	тебе надо пр	не - а	*/		нет? 0		неа
	1 138920940977	знаешь , я иногда подумываю , что	тебе надо пр	не-а	5 3		нет? 1	нет	, <mark>не хоч</mark> у
	2 138920940977	знаешь , я иногда подумываю , что	тебе надо пр	не - а	•/		нет ? 2		нет.
	3	знаешь , я иногда подумываю , что	тебе нало						

Метрика

Задача участников – возвратить ранжирование реплик-кандидатов представленных в порядке убывания скоров, выданных моделями участников. Метрика для оценивания этих ранжирований – NDCG.

$$ext{DCG}_{ ext{p}} = \sum_{i=1}^p rac{rel_i}{\log_2(i+1)} \quad ext{nDCG}_{ ext{p}} = rac{DCG_p}{IDCG_p}$$

 rel_i принимает три возможных значения - 2, 1 и 0 - для меток **good**, **neutral** и **bad** соответственно.

Особо отмечаем, что информация об уверенности меток (доступная только для тренировочных данных) никак не учитывается в метрике.

Скор участников отображаемый в контесте - это среднее NDCG для всех context_id тестовых данных, умноженное на 100 000.

```
# Neutral (1)

# Bad (0)

In [143]: def rank2num(st):

    if st == 'good':

        return 2

    else:

        if st == 'neutral':

            return 1

        else:

        return 0
```

train['rank'] = train[6].apply(rank2num)
<pre>train['target'] = train['rank'] * train[7]</pre>
train.head()

Good (2)

In [142]: # расставляем значения по формуле хорошесть * confidence

	train.neau()												
Out[144]:		0	1	L	2		3	4	5	6	7	rank	target
	0	22579918886	кликни на меня а потом на надпись " видео - зв		о,я тебя вижу	ладн	ю , повесь трубку .	0	не могу .	good	0.875352	2	1.750703
	1	22579918886	кликни на меня а потом на надпись " видео - зв		о,я тебя вижу	ладн	ю , повесь трубку .	1	нет , звонить буду я .	neutral	0.900968	1	0.900968
	2	22579918886	кликни на меня а потом на надпись " видео - зв		о , я тебя вижу	ладн	ю , повесь трубку .	2	слушай , я не мог уйти	bad	0.884320	0	0.000000
	3	22579918886	кликни на меня а потом на надпись " видео - зв		о,я тебя вижу	ладн	ю , повесь трубку .	3	я не прекращу звонить .	good	0.982530	2	1.965061
	4	22579918886	кликни на меня а потом на надпись " видео - зв		о,я тебя вижу	ладн	ю , повесь трубку .	4	я звоню им .	good	0.838054	2	1.676107

Quick'n'dirty валидация

```
In [18]: X train.shape, X test.shape
Out[18]: ((97533, 505874), (9968, 505874))
In [19]: y train = train['target']
In [20]: train part size = int(0.75 * train['target'].shape[0])
         X train part = X train[:train part size, :]
         v train part = v train[:train part size]
         X valid = X train[train part size:, :]
         y valid = y train[train part size:]
In [21]: X train part.shape
Out[21]: (73149, 505874)
In [22]: from sklearn.linear model import LinearRegression, Ridge, Lasso
         reg = Ridge()
         reg.fit(X train part, y train part)
Out[22]: Ridge(alpha=1.0, copy X=True, fit intercept=True, max iter=None,
            normalize=False, random state=None, solver='auto', tol=0.001)
In [23]: req pred = req.predict(X valid)
In [24]: %%time
         from sklearn.metrics import mean absolute error
         valid mae = mean absolute error(y valid, reg pred)
         print(valid mae)
         0.72641238803
         CPU times: user 214 µs, sys: 3.92 ms, total: 4.13 ms
         Wall time: 2.85 ms
```

Submission

```
In [202]: sub['context_id'] = test[0]
           sub['reply id'] = test[4]
           sub['rank'] = - y test
In [203]: sub.head()
Out[203]:
                context_id reply_id
                                     rank
           0 138920940977
                               0 -0.659569
           1 138920940977
                               1 -0.922145
           2 138920940977
                               2 -0.356104
           3 138920940977
                              3 -0.730093
           4 138920940977
                               4 -0.852768
In [205]: submission = sub.sort values(by=['context id', 'rank'])
In [206]: del submission['rank']
In [207]: submission.head()
Out[207]:
                context_id reply_id
           1 138920940977
           4 138920940977
           3 138920940977
                               0
           0 138920940977
           5 138920940977
In [208]: test.shape, sub.shape
Out[208]: ((9968, 6), (9968, 3))
In [209]: submission.to csv('yandex-ml-ridge.tsv',header=None, index=False, sep=' ')
```

Дальнейшие планы

- Реализовать метрику из постановки задачи для валидации
- Попробовать стемминг или работу с морфологией (к сожалению, pymystem3 оказался слишком медленным)
- Попробовать другое ранжирование (делать классификации Good / Normal / Bad)
- Попробовать что-нибудь новое модное ранее не пробованное из упомянутого на следующем слайде

Материалы по теме

- Запись трансляции тренировки ML 31.03.2018 | Kaggle Mercari, Recruit, Whatever Hack в начале обсуждается конкурс, Алиса и подходы к чатикам на темы "ни о чем" https://youtu.be/N2WMVUKUt81
- Sberbank Data Science Contest: определение релевантности вопроса параграфу <u>https://youtu.be/eN_jxfofBOk</u>
- Sberbank Data Science Journey: базовые подходы и идеи решения https://youtu.be/c7EdABKz2SE
- DrQA, a PyTorch implementation of the DrQA system described in the ACL 2017 paper <u>Reading Wikipedia to Answer Open-Domain Questions</u>
 https://github.com/facebookresearch/DrQA
- Another Twitter sentiment analysis with Python—Part 6 (Doc2Vec)
 https://towardsdatascience.com/another-twitter-sentiment-analysis-with-python-part-6-doc2vec-603f11832504

Предложенное на семинаре

- SyntaxNet
- SVD/LSA
- Сбалансировать пропорции good/normal/bad для каждого контекста (например, если мало bad - добавить строку из случайного места)
- Добавить CountVectorizer к TFIDF