

Подход к автоматизации решений задач ML на примере конкурса HackerEarth Machine Learning Challenge #1

Тренировка по машинному обучению

Яндекс, 08.04.2017



Чернобровов Алексей
к.ф.-м.н.

Jet  Retail

Ранее в сериале...



youtu.be/1sSr6Bs8rZo

Алексей Чернобровов

Анализ задачи предсказания выбора кредита (Tinkoff Data Science Challenge)



40+ часов изнурительного мозговой активности.
Это всё, конечно, очень увлекательно, но...

МОТИВАЦИЯ К РАБОТЕ



stasg7 Feb 28th at 11:26 AM
in #kaggle_crackers

Но если вкратце - в первой я опять тупо запустил скрипт,
немного подправив. Сами данные не смотрел

Я подумал о том, что при решении задач:

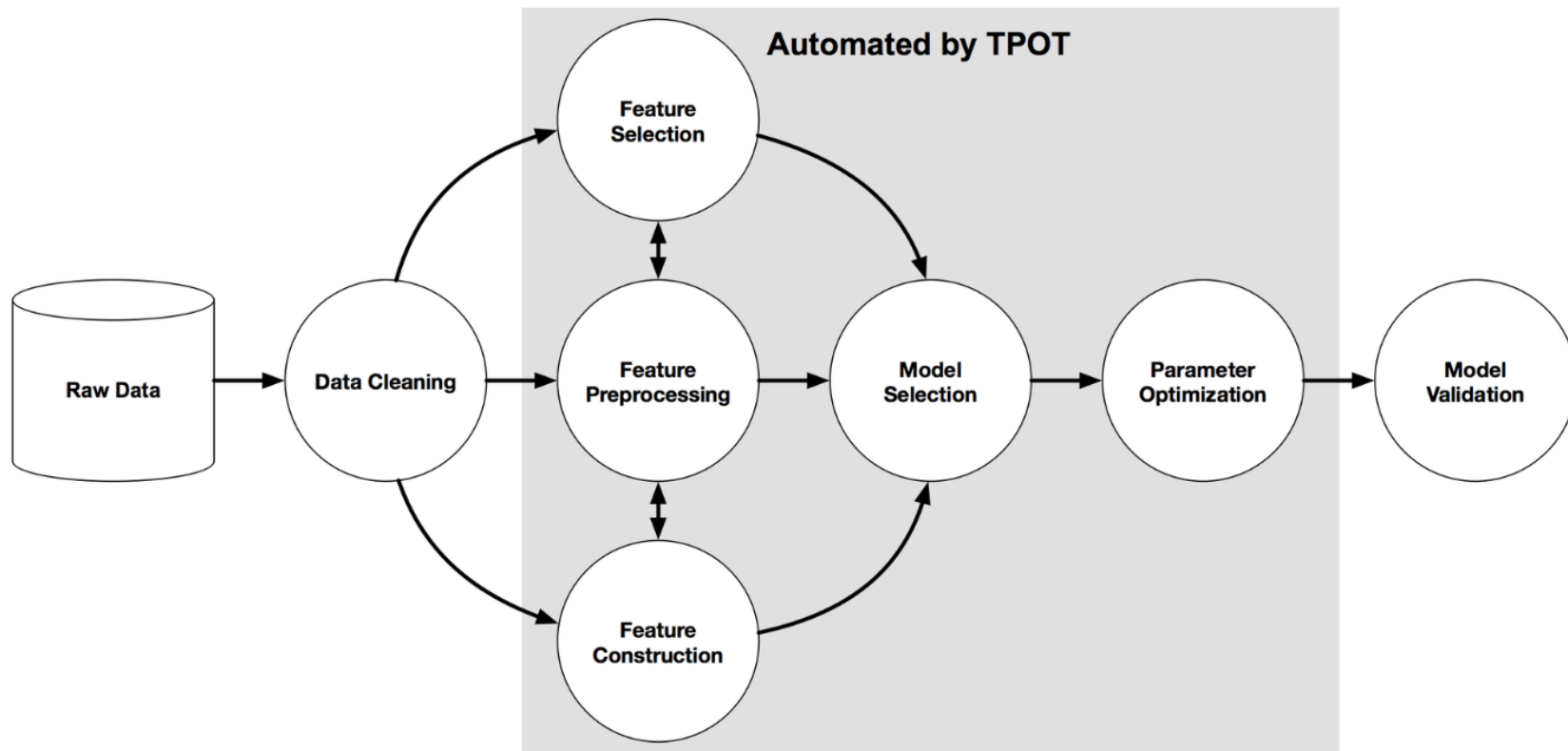
1. Очень много повторяющихся действий
2. Очень много данных, которые можно использовать повторно
3. Надо экономить время (личное и вычислительное)

Цель: написать скрипт, который бы решал задачи по ML за меня

Конечно же, я не первый, кому пришло это в голову...



<https://github.com/rhiever/tpot>



Плюсы:

1. Простой интерфейс
2. Генерирует хороший код

Минусы:

1. Слабая возможность конфигурации
2. Не сохраняет промежуточные данные
3. Не стэкает
4. Часто бывает неадекватен :)

Как должен выглядеть хороший код

Data cleaning

```
df_train = pd.read_csv('train.csv', sep=',')
df_test  = pd.read_csv('test.csv', sep=',')
X, y, X_to_predict = DataClean(df_train, df_test)
```

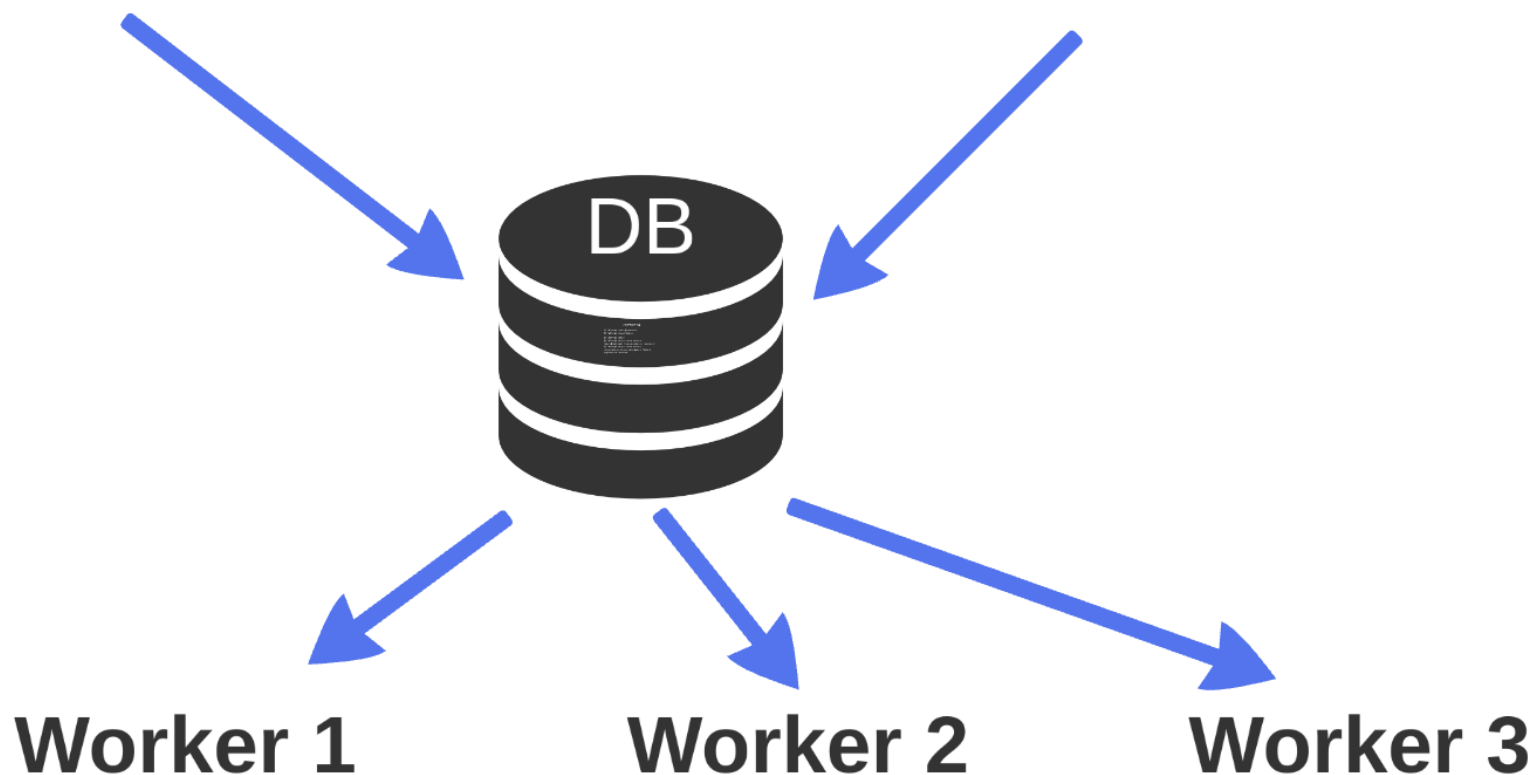
Pipelines

```
pipeline1=pipeline.make_pipeline(
    DropCols([1,2,10]),
    MultiColumnHotOneEncoder([4,5,7])
    StandardScaler(),
    LogisticRegression(penalty='l2', C=2.0)
)
pipeline2=pipeline.make_pipeline(
    DropCols([4,5,7,10,11,15]),
    RobustScaler(),
    KNeighborsClassifier(n_neighbors=100)
)
#...
clf_level2 = xgb.XGBClassifier(colsample_bytree=0.7, n_estimators=330)
y_predict=Stack(clf_level2, [pipeline1, pipeline2, ...]).fit(X, y).predict(X_to_predict)
```

Архитектура

Task generator 1

Task generator 2



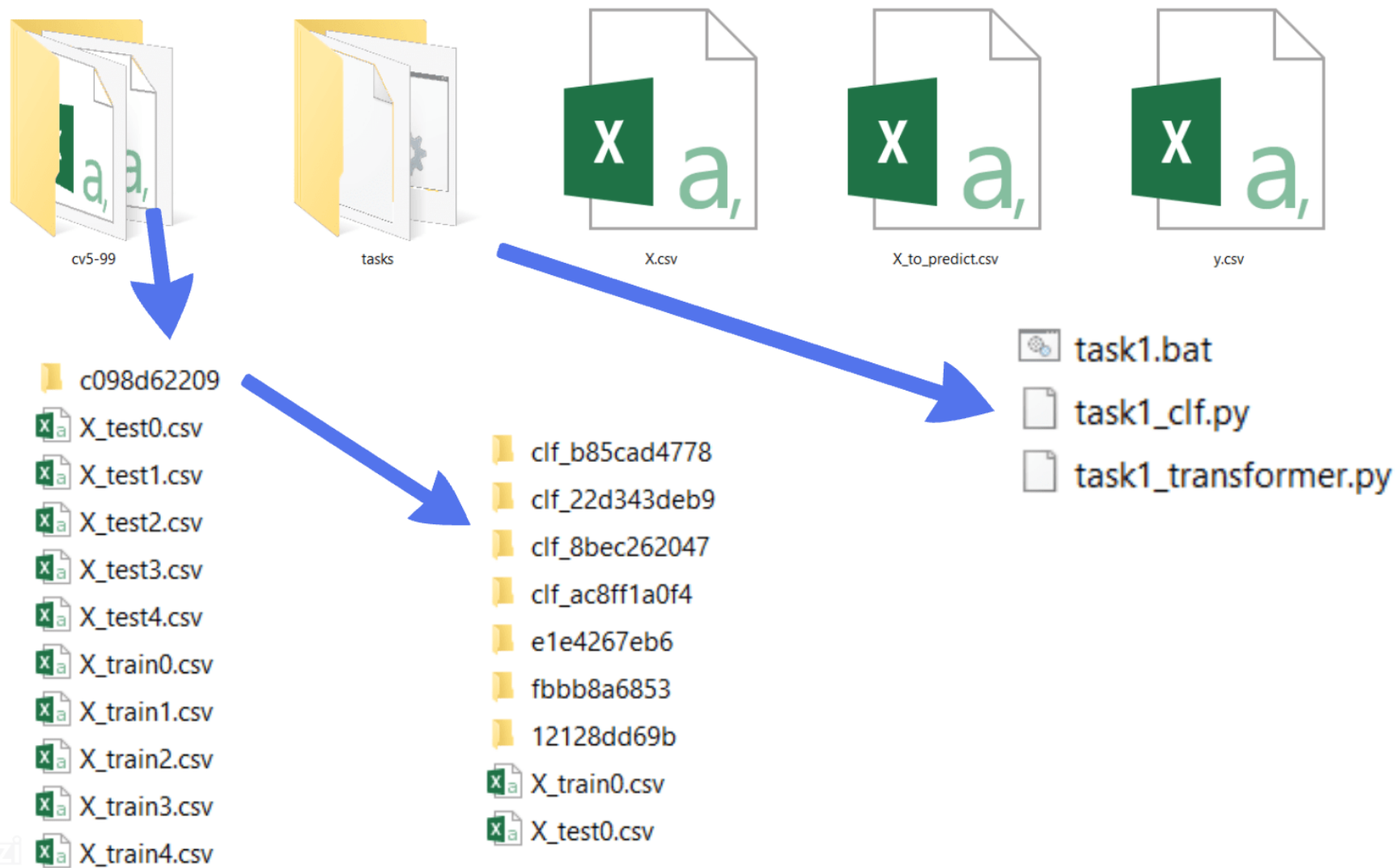
Структура БД

- 1) Таблица трансформеров
- 2) Таблица предикторов
- 3) Таблица задач
- 4) Таблица результатов работы трансформеров (число колонок, названия)
- 5) Таблица результатов работы предикторов (кросс-валидация, feature importance, intercept, ...)

Конфиг Task Generators

- Objective
- Cross-validation (folds, random_state)
- Metrics
- Genetic algorithm (selection, mutation, ...)
- Features type (numeric, categorical, binary)
- Transformers
- Predictors
- *Zoo

Структура папок



В итоге получился soft:

1. Отказоустойчивый
2. Масштабируемый
3. Расширяемый
4. *Почти* без проблем "холодного старта"
5. С широкими возможностями для конфигурации
6. Стэкает из коробки :)

HackerEarth Machine Learning Challenge #1

```
X = pd.read_csv('train_indessa.csv', sep=',')
X_test = pd.read_csv('test_indessa.csv', sep=',')
df = pd.concat([X, X_test])
```

```
df['term']=df['term'].apply(lambda x: int(x.replace(' months', "")))
```

```
grade_dict={'A' : 1, 'B' : 2, 'C' : 3, 'D' : 4, 'E' : 5, 'F' : 6, 'G' : 7}
for i in grade_dict:
    df['sub_grade']=df['sub_grade'].apply(lambda x: (x.replace(i, grade_dict[i])))
```

```
df['emp_length']=df['emp_length'].apply(lambda x: (x.replace(' years', "")))
df['emp_length']=df['emp_length'].apply(lambda x: (x.replace('10+', "11")))
df['emp_length']=df['emp_length'].apply(lambda x: (x.replace('< 1 year', "0")))
df['emp_length']=df['emp_length'].apply(lambda x: (x.replace('1 year', "1")))
df['emp_length']=df['emp_length'].apply(lambda x: int(x.replace('n/a', "-1")))
```

```
cat_features0 = ['pymnt_plan', 'purpose', 'initial_list_status', 'application_type', 'addr_state']
for col in cat_features0:
    df[col], _ = pd.factorize(df[col])
```

```
df['zip_code1']=df['zip_code'].apply(lambda x: int(x[0]))
df['zip_code2']=df['zip_code'].apply(lambda x: int(x[0:2]))
df['zip_code']=df['zip_code'].apply(lambda x: int(x.replace('xx', "")))
```

```
df['last_week_pay']=df['last_week_pay'].apply(lambda x: (x.replace('th week', "")))
df['last_week_pay']=df['last_week_pay'].apply(lambda x: (x.replace('NA', "-1")))
```

```
df['desc_isnan']=pd.isnull(df['desc'])
df['emp_title_isnan']=pd.isnull(df['emp_title'])
```

```
cat_to_del=['desc', 'funded_amnt', 'grade', 'loan_amnt', 'batch_enrolled']
```

Это весь
Data
cleaning

**Я просто запустил свой скрипт
И...
НАЧАЛ ПРАВИТЬ БАГИ :)**

Далеко не полный перечень проблем:

1. Нужно быть очень аккуратным с параметрами алгоритмов
2. Любой баг в "шаблоне" приводит к полному отсутствию предиктора
3. Нужно следить за именами колонок.
4. Делать PCA 3 раза подряд - не самая лучшая идея.
5. Такая же история с PolynomialFeatures... :)
6. У некоторых классификаторов нет `predict_proba`, а только `predict`

...

В итоге:

1. Финальная модель:

Стэкинг из 5 lgbm, 3 xgb, 2 adabost, 1 rf, 1 et.

Для большинства из них были разные трансформеры.

И xgb на втором уровне.

2. Время работы.

Мое: 4 часа (2 часа data cleaning)

Вычислительное: примерно 60-70 часов.

3. Результат: 9 место из 3000+

Полезные ссылки

1. Яндекс изнутри:
как машинное обучение делает поиск умнее

 youtu.be/HYKdyhB1Kcs

2. Юрий Кашницкий. 16 ядер и 30 Гб под капотом Вашего Jupyter за \$0.25 в час

 habrahabr.ru/post/280562/

Спасибо за внимание!

Вопросы?

Подход к автоматизации решений задач ML
на примере «HackerRank»
HackerRank Machine Learning Challenge #1

Технологии на которых реализовано решение

Python 3.6.5, Jupyter Notebook, TensorFlow 1.12.0



Конфиг Task Generators

- Описание
- Ссылка на данные (URL, локальный диск)
- Настройки
- Описание задачи (описание, название, метрика)
- Описание задачи (описание, название, метрика)
- Настройки
- Описание
- Ссылка



В итоге получили soft:

1. Описание задачи
2. Описание задачи
3. Описание задачи
4. Описание задачи
5. Описание задачи
6. Описание задачи

HackerRank Machine Learning Challenge #1

Эта часть

Task generator	Worker	Score
1	1	0.000000
1	2	0.000000
1	3	0.000000
2	1	0.000000
2	2	0.000000
2	3	0.000000

Я просто запустил свой скрипт
И...
НАЧАЛ ПРАВИТЬ БАГИ :)

- В итоге:
1. Внесены правки.
 2. Внесены правки.
 3. Внесены правки.

Полезные ссылки

- Описание задачи
- Описание задачи
- Описание задачи



vk.com/chernobrovov



OpenDataScience @alex4er

Алексей Чернобровов