In [0]: # access kaggle datasets !pip install kaggle # math operations !pip install numpy==1.15.0 !pip install eli5 !pip install shap # for machine learning !pip install catboost Requirement already satisfied: kaggle in /usr/local/lib/python3.6/dist-packa ges (1.5.3)Requirement already satisfied: urllib3<1.25,>=1.21.1 in /usr/local/lib/pytho n3.6/dist-packages (from kaggle) (1.22) Requirement already satisfied: six>=1.10 in /usr/local/lib/python3.6/dist-pa ckages (from kaggle) (1.11.0) Requirement already satisfied: certifi in /usr/local/lib/python3.6/dist-pack ages (from kaggle) (2018.11.29) Requirement already satisfied: python-dateutil in /usr/local/lib/python3.6/d ist-packages (from kaggle) (2.5.3) Requirement already satisfied: requests in /usr/local/lib/python3.6/dist-pac kages (from kaggle) (2.18.4) Requirement already satisfied: tqdm in /usr/local/lib/python3.6/dist-package s (from kaggle) (4.28.1)Requirement already satisfied: python-slugify in /usr/local/lib/python3.6/di st-packages (from kaggle) (2.0.1) Requirement already satisfied: chardet<3.1.0,>=3.0.2 in /usr/local/lib/pytho n3.6/dist-packages (from requests->kaggle) (3.0.4) Requirement already satisfied: idna<2.7,>=2.5 in /usr/local/lib/python3.6/di st-packages (from requests->kaggle) (2.6) Requirement already satisfied: Unidecode>=0.04.16 in /usr/local/lib/python3. 6/dist-packages (from python-slugify->kaggle) (1.0.23) Collecting numpy==1.15.0 Downloading https://files.pythonhosted.org/packages/88/29/f4c845648ed23264 e986cdc5fbab5f8eace1be5e62144ef69ccc7189461d/numpy-1.15.0-cp36-cp36m-manylin ux1 x86 64.whl (13.9MB) | 13.9MB 3.6MB/s 100% | torchvision 0.2.1 has requirement pillow>=4.1.1, but you'll have pillow 4.0. 0 which is incompatible. thinc 6.12.1 has requirement wrapt<1.11.0,>=1.10.0, but you'll have wrapt 1. 11.1 which is incompatible. pymc3 3.6 has requirement joblib<0.13.0, but you'll have joblib 0.13.2 which is incompatible. featuretools 0.4.1 has requirement pandas>=0.23.0, but you'll have pandas 0. 22.0 which is incompatible. albumentations 0.1.12 has requirement imgaug<0.2.7,>=0.2.5, but you'll have imgaug 0.2.8 which is incompatible. Installing collected packages: numpy Found existing installation: numpy 1.14.6 Uninstalling numpy-1.14.6: Successfully uninstalled numpy-1.14.6 Successfully installed numpy-1.15.0 Collecting eli5 Downloading https://files.pythonhosted.org/packages/8d/c8/04bed18dcce1d927 b0dd5fc3425777354b714d2e62d60ae301928b5a5bf8/eli5-0.8.1-py2.py3-none-any.whl (98kB) 100% | | 102kB 3.4MB/s Requirement already satisfied: scikit-learn>=0.18 in /usr/local/lib/python3. 6/dist-packages (from eli5) (0.20.2) Requirement already satisfied: attrs>16.0.0 in /usr/local/lib/python3.6/dist -packages (from eli5) (18.2.0) Requirement already satisfied: six in /usr/local/lib/python3.6/dist-packages (from eli5) (1.11.0) Requirement already satisfied: numpy>=1.9.0 in /usr/local/lib/python3.6/dist -packages (from eli5) (1.15.0) Requirement already satisfied: graphviz in /usr/local/lib/python3.6/dist-pac kages (from eli5) (0.10.1) Requirement already satisfied: jinja2 in /usr/local/lib/python3.6/dist-packa ges (from eli5) (2.10) Requirement already satisfied: tabulate>=0.7.7 in /usr/local/lib/python3.6/d ist-packages (from eli5) (0.8.3) Requirement already satisfied: scipy in /usr/local/lib/python3.6/dist-packag es (from eli5) (1.1.0) Requirement already satisfied: typing in /usr/local/lib/python3.6/dist-packa ges (from eli5) (3.6.6) Requirement already satisfied: MarkupSafe>=0.23 in /usr/local/lib/python3.6/ dist-packages (from jinja2->eli5) (1.1.0) Installing collected packages: eli5 Successfully installed eli5-0.8.1 Collecting shap Downloading https://files.pythonhosted.org/packages/30/b3/866b0101cbd18298 44c35964af68c14ba522a5cce7a1e8d0f7937411d910/shap-0.28.5.tar.gz (223kB) | 225kB 9.2MB/s Requirement already satisfied: numpy in /usr/local/lib/python3.6/dist-packag es (from shap) (1.15.0) Requirement already satisfied: scipy in /usr/local/lib/python3.6/dist-packag es (from shap) (1.1.0) Requirement already satisfied: scikit-learn in /usr/local/lib/python3.6/dist -packages (from shap) (0.20.2) Requirement already satisfied: matplotlib in /usr/local/lib/python3.6/dist-p ackages (from shap) (3.0.2) Requirement already satisfied: pandas in /usr/local/lib/python3.6/dist-packa ges (from shap) (0.22.0) Requirement already satisfied: tqdm in /usr/local/lib/python3.6/dist-package s (from shap) (4.28.1) Requirement already satisfied: ipython in /usr/local/lib/python3.6/dist-pack ages (from shap) (5.5.0)Requirement already satisfied: scikit-image in /usr/local/lib/python3.6/dist -packages (from shap) (0.13.1) Requirement already satisfied: pyparsing!=2.0.4,!=2.1.2,!=2.1.6,>=2.0.1 in / usr/local/lib/python3.6/dist-packages (from matplotlib->shap) (2.3.1) Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.6/dist -packages (from matplotlib->shap) (0.10.0) Requirement already satisfied: python-dateutil>=2.1 in /usr/local/lib/python 3.6/dist-packages (from matplotlib->shap) (2.5.3) Requirement already satisfied: kiwisolver>=1.0.1 in /usr/local/lib/python3. 6/dist-packages (from matplotlib->shap) (1.0.1) Requirement already satisfied: pytz>=2011k in /usr/local/lib/python3.6/distpackages (from pandas->shap) (2018.9) Requirement already satisfied: pexpect; sys\_platform != "win32" in /usr/loca 1/lib/python3.6/dist-packages (from ipython->shap) (4.6.0) Requirement already satisfied: prompt-toolkit<2.0.0,>=1.0.4 in /usr/local/li b/python3.6/dist-packages (from ipython->shap) (1.0.15) Requirement already satisfied: pygments in /usr/local/lib/python3.6/dist-pac kages (from ipython->shap) (2.1.3) Requirement already satisfied: simplegeneric>0.8 in /usr/local/lib/python3. 6/dist-packages (from ipython->shap) (0.8.1) Requirement already satisfied: traitlets>=4.2 in /usr/local/lib/python3.6/di st-packages (from ipython->shap) (4.3.2) Requirement already satisfied: pickleshare in /usr/local/lib/python3.6/distpackages (from ipython->shap) (0.7.5) Requirement already satisfied: setuptools>=18.5 in /usr/local/lib/python3.6/ dist-packages (from ipython->shap) (40.8.0) Requirement already satisfied: decorator in /usr/local/lib/python3.6/dist-pa ckages (from ipython->shap) (4.3.2) Requirement already satisfied: six>=1.7.3 in /usr/local/lib/python3.6/dist-p ackages (from scikit-image->shap) (1.11.0) Requirement already satisfied: networkx>=1.8 in /usr/local/lib/python3.6/dis t-packages (from scikit-image->shap) (2.2) Requirement already satisfied: PyWavelets>=0.4.0 in /usr/local/lib/python3. 6/dist-packages (from scikit-image->shap) (1.0.1) Requirement already satisfied: pillow>=2.1.0 in /usr/local/lib/python3.6/dis t-packages (from scikit-image->shap) (4.0.0) Requirement already satisfied: ptyprocess>=0.5 in /usr/local/lib/python3.6/d ist-packages (from pexpect; sys platform != "win32"->ipython->shap) (0.6.0) Requirement already satisfied: wcwidth in /usr/local/lib/python3.6/dist-pack ages (from prompt-toolkit<2.0.0,>=1.0.4->ipython->shap) (0.1.7) Requirement already satisfied: ipython-genutils in /usr/local/lib/python3.6/ dist-packages (from traitlets>=4.2->ipython->shap) (0.2.0) Requirement already satisfied: olefile in /usr/local/lib/python3.6/dist-pack ages (from pillow>=2.1.0->scikit-image->shap) (0.46) Building wheels for collected packages: shap Building wheel for shap (setup.py) ... done Stored in directory: /root/.cache/pip/wheels/bf/26/bd/912db1314f1cef0171d9 b7f128dd01e8b8c92ed8d0062e632d Successfully built shap Installing collected packages: shap Successfully installed shap-0.28.5 Collecting catboost Downloading https://files.pythonhosted.org/packages/98/03/777a0e1c12571a7f 3320a4fa6d5f123dba2dd7c0bca34f4f698a6396eb48/catboost-0.12.2-cp36-none-manyl inux1 x86 64.whl (55.5MB) | 55.5MB 818kB/s Requirement already satisfied: pandas>=0.19.1 in /usr/local/lib/python3.6/di st-packages (from catboost) (0.22.0) Requirement already satisfied: six in /usr/local/lib/python3.6/dist-packages (from catboost) (1.11.0) Requirement already satisfied: enum34 in /usr/local/lib/python3.6/dist-packa ges (from catboost) (1.1.6) Requirement already satisfied: numpy>=1.11.1 in /usr/local/lib/python3.6/dis t-packages (from catboost) (1.15.0) Requirement already satisfied: python-dateutil>=2 in /usr/local/lib/python3. 6/dist-packages (from pandas>=0.19.1->catboost) (2.5.3) Requirement already satisfied: pytz>=2011k in /usr/local/lib/python3.6/distpackages (from pandas>=0.19.1->catboost) (2018.9) Installing collected packages: catboost Successfully installed catboost-0.12.2 In [0]: # Libraries import numpy as np import pandas as pd pd.set\_option('max\_columns', None) import matplotlib.pyplot as plt import seaborn as sns %matplotlib inline import datetime import lightgbm as lgb from scipy import stats from sklearn.model selection import train test split, StratifiedKFold, KFold, cross val score, GridSearchCV, RepeatedStratifiedKFold from sklearn.preprocessing import StandardScaler import os import plotly.offline as py py.init notebook mode(connected=True) import plotly.graph objs as go import plotly.tools as tls import xgboost as xgb import lightgbm as lgb from sklearn import model selection from sklearn.metrics import accuracy\_score, roc\_auc\_score from sklearn import metrics import json import ast import time from sklearn import linear model import eli5 from eli5.sklearn import PermutationImportance import shap from tqdm import tqdm notebook from mlxtend.feature selection import SequentialFeatureSelector as SFS from mlxtend.plotting import plot sequential feature selection as plot sfs from sklearn.neighbors import NearestNeighbors from sklearn.feature selection import GenericUnivariateSelect, SelectPercentil e, SelectKBest, f classif, mutual info classif, RFE import statsmodels.api as sm import warnings warnings.filterwarnings('ignore') from catboost import CatBoostClassifier /usr/local/lib/python3.6/dist-packages/statsmodels/compat/pandas.py:56: Futu reWarning: The pandas.core.datetools module is deprecated and will be removed in a futu re version. Please use the pandas.tseries module instead. **Import Dataset from Kaggle** In [0]: # Colab's file access feature from google.colab import files # retrieve uploaded file uploaded = files.upload() # print results for fn in uploaded.keys(): print('User uploaded file "{name}" with length {length} bytes'.format( name=fn, length=len(uploaded[fn]))) # then move kaggle.json into the folder where the API expects to find it. !mkdir -p ~/.kaggle/ && mv kaggle.json ~/.kaggle/ && chmod 600 ~/.kaggle/kaggl e.json Choose Files No file chosen Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable. Saving kaggle.json to kaggle.json User uploaded file "kaggle.json" with length 62 bytes In [0]: # list files !ls sample data sample submission.csv.zip test.csv.zip train.csv.zip In [0]: # unzip training data for usage !unzip train.parquet.zip !unzip test.parquet.zip !unzip test.csv.zip !unzip train.csv.zip !unzip sample submission.csv.zip unzip: cannot find or open train.parquet.zip, train.parquet.zip.zip or trai n.parquet.zip.ZIP. unzip: cannot find or open test.parquet.zip, test.parquet.zip.zip or test.p arquet.zip.ZIP. Archive: test.csv.zip inflating: test.csv Archive: train.csv.zip inflating: train.csv Archive: sample submission.csv.zip inflating: sample submission.csv In [0]: # download santander customer transaction dataset !kaggle competitions download -c santander-customer-transaction-prediction train.csv.zip: Skipping, found more recently modified local copy (use --forc e to force download) sample submission.csv.zip: Skipping, found more recently modified local copy (use --force to force download) test.csv.zip: Skipping, found more recently modified local copy (use --force to force download) In [0]: train = pd.read csv('train.csv') test = pd.read csv('test.csv') train.shape, test.shape Out[0]: ((200000, 202), (200000, 201)) **Data exploration** In [0]: train.head() Out[0]: ID code target var\_0 var\_1 var\_2 var\_3 var\_4 var\_5 var\_6 var\_7 0 train\_0  $0 \quad 8.9255 \quad \text{-}6.7863 \quad 11.9081 \quad 5.0930 \quad 11.4607 \quad \text{-}9.2834 \quad 5.1187 \quad 18.6266 \quad \text{-}4.9200$ train 1 0 11.5006 -4.1473 13.8588 5.3890 12.3622 7.0433 5.6208 16.5338 3.1468 0 8.6093 -2.7457 12.0805 7.8928 10.5825 -9.0837 6.9427 14.6155 -4.9193 train\_2 train\_3 0 11.0604 -2.1518 8.9522 7.1957 12.5846 -1.8361 5.8428 14.9250 -5.8609 train\_4 0 9.8369 -1.4834 12.8746 6.6375 12.2772 2.4486 5.9405 19.2514 6.2654 In [0]: train[train.columns[2:]].std().plot('hist'); plt.title('Distribution of stds of all columns'); Distribution of stds of all columns 60 50 40 30 20 10 0 In [0]: train[train.columns[2:]].mean().plot('hist'); plt.title('Distribution of means of all columns'); Distribution of means of all columns 35 30 25 20 15 10 5 0 -10 0 10 20 In [0]: train.head() Out[0]: **ID\_code** target var\_0 var\_1 var\_2 var\_3 var\_4 var\_5 var\_6 var\_7 var 8 train\_0 8.9255 -6.7863 11.9081 5.0930 11.4607 -9.2834 5.1187 18.6266 -4.9200 0 11.5006 -4.1473 13.8588 5.3890 12.3622 7.0433 5.6208 16.5338 1 train\_1 3.1468 train 2 8.6093 -2.7457 12.0805 7.8928 10.5825 -9.0837 6.9427 14.6155 -4.9193 train\_3 0 11.0604 -2.1518 8.9522 7.1957 12.5846 -1.8361 5.8428 14.9250 -5.8609 train\_4 9.8369 -1.4834 12.8746 6.6375 12.2772 2.4486 5.9405 19.2514 In [0]: # we have no missing values train.isnull().any().any() Out[0]: False In [0]: print('Distributions of first 28 columns') plt.figure(figsize=(26, 24)) for i, col in enumerate(list(train.columns)[2:30]): plt.subplot(7, 4, i + 1)plt.hist(train[col]) plt.title(col) Distributions of first 28 columns In [0]: train['target'].value counts(normalize=True) Out[0]: 0 0.89951 0.10049 Name: target, dtype: float64 From this overview we can see the following things: target is binary and has disbalance: 10% of samples belong to 1 class; values in columns are more or less similar; columns have high std (up to 20) columns have a high range of means; Let's have a look at correlations now! In [0]: corrs = train.corr().abs().unstack().sort values(kind="quicksort").reset index corrs = corrs[corrs['level 0'] != corrs['level 1']] corrs.tail(30) Out[0]: level\_0 level\_1 40170 var\_80 0.057609 target 40171 var\_80 target 0.057609 **40172** var\_166 target 0.057773 40173 target var\_166 0.057773 40174 var\_99 target 0.058367 40175 var\_99 0.058367 target 40176 target var\_21 0.058483 40177 var\_21 target 0.058483 40178 target var\_22 0.060558 40179 var\_22 target 0.060558 40180 target var\_174 0.061669 40181 var\_174 target 0.061669 40182 var\_76 target 0.061917 40183 var\_76 0.061917 target 40184 var\_26 0.062422 target 40185 var\_26 target 0.062422 40186 var\_53 target 0.063399 40187 target var\_53 0.063399 40188 var\_146 target 0.063644 **40190** var\_110 target 0.064275 target var\_110 0.064275 40191 40192 target 0.066731 var\_6 40193 var\_6 0.066731 target target 0.069489 40194 var\_12 var\_12 0.069489 40195 target 40196 target var\_139 0.074080 **40197** var 139 target 0.074080 40198 var\_81 0.080917 target var 81 40199 target 0.080917 corrs.head() In [0]: Out[0]: level 0 level 1 **0** var\_191 var\_75 2.703975e-08 1 var\_75 var\_191 2.703975e-08 var\_6 5.942735e-08 **2** var\_173 var\_6 var\_173 5.942735e-08 4 var\_126 var\_109 1.313947e-07 We can see that all features have a low correlation with target. So we have no highly correlated features which we could drop, on the other hand we could drop some columns with have little correlation with the target. **Basic modelling** In [0]: X = train.drop(['ID code', 'target'], axis=1) y = train['target'] X\_test = test.drop(['ID\_code'], axis=1) n fold = 5folds = StratifiedKFold(n\_splits=n\_fold, shuffle=True, random\_state=42) repeated\_folds = RepeatedStratifiedKFold(n\_splits=10, n\_repeats=20, random\_sta te=42)# scaler = StandardScaler() # X train = scaler.fit transform(X train) # X test = scaler.transform(X test) In [0]: | def train\_model(X, X\_test, y, params, folds, model\_type='lgb', plot\_feature\_im portance=False, averaging='usual', model=None): oof = np.zeros(len(X)) prediction = np.zeros(len(X\_test)) feature importance = pd.DataFrame() for fold\_n, (train\_index, valid\_index) in enumerate(folds.split(X, y)): print('Fold', fold\_n, 'started at', time.ctime()) X\_train, X\_valid = X.loc[train\_index], X.loc[valid\_index] y\_train, y\_valid = y[train\_index], y[valid\_index] if model type == 'lgb': train data = lgb.Dataset(X train, label=y train) valid\_data = lgb.Dataset(X\_valid, label=y\_valid) model = lgb.train(params, train data, num\_boost\_round=20000, valid sets = [train data, valid data], verbose eval=1000, early\_stopping\_rounds = 200) y\_pred\_valid = model.predict(X\_valid) y pred = model.predict(X test, num iteration=model.best iteration) if model\_type == 'xgb': train data = xgb.DMatrix(data=X train, label=y train, feature name s=X train.columns) valid data = xgb.DMatrix(data=X valid, label=y valid, feature name s=X train.columns) watchlist = [(train\_data, 'train'), (valid\_data, 'valid\_data')] model = xgb.train(dtrain=train\_data, num\_boost\_round=20000, evals= watchlist, early stopping rounds=200, verbose eval=500, params=params) y pred valid = model.predict(xgb.DMatrix(X valid, feature names=X train.columns), ntree\_limit=model.best\_ntree\_limit) y pred = model.predict(xgb.DMatrix(X test, feature names=X train.c olumns), ntree limit=model.best ntree limit) if model\_type == 'sklearn': model = modelmodel.fit(X train, y train) y\_pred\_valid = model.predict\_proba(X\_valid).reshape(-1,) score = roc auc score(y valid, y pred valid) # print(f'Fold {fold\_n}. AUC: {score:.4f}.') # print('') y\_pred = model.predict\_proba(X\_test)[:, 1] if model type == 'glm': model = sm.GLM(y train, X train, family=sm.families.Binomial()) model results = model.fit() model results.predict(X test) y pred valid = model results.predict(X valid).reshape(-1,) score = roc auc score(y valid, y pred valid) y pred = model results.predict(X test) if model\_type == 'cat': model = CatBoostClassifier(iterations=20000, learning rate=0.05, 1 oss\_function='Logloss', eval\_metric='AUC', \*\*params) model.fit(X\_train, y\_train, eval\_set=(X\_valid, y\_valid), cat\_featu res=[], use best model=True, verbose=False) y pred valid = model.predict proba(X valid)[:, 1] y pred = model.predict proba(X test)[:, 1] oof[valid\_index] = y\_pred\_valid.reshape(-1,) scores.append(roc auc score(y valid, y pred valid)) if averaging == 'usual': prediction += y pred elif averaging == 'rank': prediction += pd.Series(y\_pred).rank().values if model type == 'lgb': # feature importance fold\_importance = pd.DataFrame() fold importance["feature"] = X.columns fold importance["importance"] = model.feature importance() fold\_importance["fold"] = fold\_n + 1 feature\_importance = pd.concat([feature\_importance, fold\_importanc e], axis=0) prediction /= n fold print('CV mean score: {0:.4f}, std: {1:.4f}.'.format(np.mean(scores), np.s td(scores))) if model type == 'lgb': feature\_importance["importance"] /= n\_fold if plot\_feature\_importance: cols = feature\_importance[["feature", "importance"]].groupby("feat ure").mean().sort\_values( by="importance", ascending=False)[:50].index best feature = feature importance.loc[feature\_importance.feature. isin(cols)] plt.figure(figsize=(16, 12)); sns.barplot(x="importance", y="feature", data=best\_features.sort\_v alues(by="importance", ascending=False)); plt.title('LGB Features (avg over folds)'); return oof, prediction, feature\_importance return oof, prediction, scores else: return oof, prediction, scores In [0]: # %%time # model = linear model.LogisticRegression(class weight='balanced', penalty='l 2', C=0.1) # oof lr, prediction lr, scores = train model(X, X\_test, y, params=None, folds =folds, model type='sklearn', model=model) In [0]: params = {'num\_leaves': 8, 'min\_data\_in\_leaf': 42, 'objective': 'binary', 'max depth': 16, 'learning rate': 0.0123, 'boosting': 'gbdt', 'bagging freq': 5, 'feature\_fraction': 0.8201, 'bagging seed': 11, 'reg\_alpha': 1.728910519108444, 'reg\_lambda': 4.9847051755586085, 'random\_state': 42, 'metric': 'auc', 'verbosity': -1, 'subsample': 0.81, 'min\_gain\_to\_split': 0.01077313523861969, 'min child weight': 19.428902804238373, 'num threads': 4} oof\_lgb, prediction\_lgb, scores = train\_model(X, X\_test, y, params=params, fol ds=folds, model\_type='lgb', plot\_feature\_importance=True) Fold 0 started at Wed Feb 27 07:27:43 2019 Training until validation scores don't improve for 200 rounds. [1000] training's auc: 0.885014 valid\_1's auc: 0.864492 [2000] training's auc: 0.909721 valid\_1's auc: 0.882867 [3000] training's auc: 0.921202 valid\_1's auc: 0.89042 [4000] training's auc: 0.928088 valid\_1's auc: 0.89445 [5000] training's auc: 0.932895 valid\_1's auc: 0.896647 [6000] training's auc: 0.936977 valid\_1's auc: 0.897681 In [0]: sub = pd.read csv('sample submission.csv') sub['target'] = prediction lgb sub.to csv('lgb.csv', index=False) ELI5 In [0]: | model = lgb.LGBMClassifier(\*\*params, n estimators = 20000, n jobs = -1) X train, X valid, y train, y valid = train test split(X, y, test size=0.2, str atify=y) model.fit(X\_train, y\_train, eval\_set=[(X\_train, y\_train), (X\_valid, y\_valid)], verbose=1000, early\_stopping\_rounds=200) In [0]: eli5.show weights (model, targets=[0, 1], feature names=list(X train.columns), top=40, feature filter=lambda x: x != '<BIAS>') ELI5 didn't help up to eliminate features, but let's at least try to take top-100 and see how it helps. In [0]: top features = [i for i in eli5.formatters.as dataframe.explain weights df(mod el).feature if 'BIAS' not in i][:100]  $X1 = X[top_features]$ X train, X valid, y train, y valid = train test split(X1, y, test size=0.2, st ratify=y) model.fit(X\_train, y\_train, eval\_set=[(X\_train, y\_train), (X\_valid, y\_valid)], verbose=1000, early stopping rounds=200) In [0]: def calculate metrics (model, X train: pd.DataFrame() = None, y train: pd.DataF rame() = None, X valid: pd.DataFrame() = None, y valid: pd.DataFrame() = None, columns: list = []) -> p d.DataFrame(): columns = columns if len(columns) > 0 else list(X train.columns) train pred = model.predict proba(X train[columns]) valid pred = model.predict proba(X valid[columns]) f1 = 0best t = 0for t in np.arange(0.1, 1, 0.05): valid pr = (valid pred[:, 1] > t).astype(int) valid f1 = metrics.f1 score(y valid, valid pr) if valid f1 > f1: f1 = valid f1best t = tt = best ttrain pr = (train pred[:, 1] > t).astype(int) valid\_pr = (valid\_pred[:, 1] > t).astype(int) train\_f1 = metrics.f1\_score(y\_train, train\_pr) valid\_f1 = metrics.f1\_score(y\_valid, valid\_pr) score df = []print(f'Best threshold: {t:.2f}. Train f1: {train\_f1:.4f}. Valid f1: {vali d f1:.4f}.') score df.append(['F1', np.round(train f1, 4), np.round(valid f1, 4)]) train\_r = metrics.recall\_score(y\_train, train\_pr) valid r = metrics.recall score(y valid, valid pr) score df.append(['Recall', np.round(train r, 4), np.round(valid r, 4)]) train\_p = metrics.precision\_score(y\_train, train\_pr) valid p = metrics.precision score(y valid, valid pr) score df.append(['Precision', np.round(train p, 4), np.round(valid p, 4)]) train\_roc = metrics.roc\_auc\_score(y\_train, train\_pred[:, 1]) valid roc = metrics.roc\_auc\_score(y\_valid, valid\_pred[:, 1]) score df.append(['ROCAUC', np.round(train roc, 4), np.round(valid roc, 4 )]) train apc = metrics.average precision score(y train, train pred[:, 1]) valid apc = metrics.average precision score(y valid, valid pred[:, 1]) score\_df.append(['APC', np.round(train\_apc, 4), np.round(valid\_apc, 4)]) print(metrics.confusion\_matrix(y\_valid, valid\_pr)) score df = pd.DataFrame(score df, columns=['Metric', 'Train', 'Valid']) print(score df) return score df, t In [0]: = calculate metrics(model, X train, y train, X valid, y valid) Feature generation **Feature interaction** In [0]: X = train.drop(['ID code', 'target'], axis=1) X test = test.drop(['ID code'], axis=1) columns = top features = [i for i in eli5.formatters.as dataframe.explain weig hts df(model).feature if 'BIAS' not in i][:20] for coll in tqdm notebook(columns): for col2 in columns: X[col1 + ' ' + col2] = X[col1] \* X[col2]X test[col1 + ' ' + col2] = X test[col1] \* X test[col2] In [0]: | # oof lgb, prediction lgb inter, scores = train model(X, X test, y, params=par ams, folds=folds, model type='lgb', plot feature importance=True) In [0]: # sub = pd.read csv('sample submission.csv') # sub['target'] = prediction lgb inter # sub.to csv('lgb inter.csv', index=False) Scaling ! Notice scaling severely decreases score In [0]: X = train.drop(['ID code', 'target'], axis=1) X test = test.drop(['ID code'], axis=1) scaler = StandardScaler() X train[X train.columns] = scaler.fit transform(X train[X train.columns]) X test[X train.columns] = scaler.transform(X test[X train.columns]) # oof\_lgb, prediction\_lgb\_scaled, scores = train\_model(X, X\_test, y, params=pa rams, folds=folds, model\_type='lgb', plot\_feature\_importance=True) # sub = pd.read csv('sample submission.csv') # sub['target'] = prediction lgb scaled # sub.to csv('lgb scaled.csv', index=False) **Statistics** In [0]: # X = train.drop(['ID\_code', 'target'], axis=1) # X test = test.drop(['ID code'], axis=1) # X['std'] = X.std(1)# X\_test['std'] = X\_test.std(1) # X['mean'] = X.mean(1) # X test['mean'] = X test.mean(1) # oof\_lgb, prediction\_lgb\_stats, scores = train\_model(X, X\_test, y, params=par ams, folds=folds, model\_type='lgb', plot\_feature\_importance=True) # sub = pd.read csv('sample submission.csv') # sub['target'] = prediction lgb stats # sub.to\_csv('lgb\_stats.csv', index=False) Training with these features gives the same score on LB: 0.899 NN features In [0]: # %%time # X = train.drop(['ID\_code', 'target'], axis=1) # X\_test = test.drop(['ID\_code'], axis=1) # neigh = NearestNeighbors(3, n jobs=-1) # neigh.fit(X) # dists, \_ = neigh.kneighbors(X, n\_neighbors=3) # mean dist = dists.mean(axis=1) # max\_dist = dists.max(axis=1) # min dist = dists.min(axis=1) # X['mean dist'] = mean dist # X['max dist'] = max dist # X['min\_dist'] = min\_dist # test dists, \_ = neigh.kneighbors(X\_test, n\_neighbors=3) # test\_mean\_dist = test\_dists.mean(axis=1) # test\_max\_dist = test\_dists.max(axis=1) # test\_min\_dist = test\_dists.min(axis=1)

# X\_test['mean\_dist'] = test\_mean\_dist
# X\_test['max\_dist'] = test\_max\_dist
# X\_test['min\_dist'] = test\_min\_dist

# sub['target'] = prediction\_lgb\_dist
# sub.to\_csv('lgb\_dist.csv', index=False)

**Blend** 

e': 0.9,

# sub = pd.read\_csv('sample\_submission.csv')

# oof\_lgb, prediction\_lgb\_dist, scores = train model(X, X test, y, params=para

'objective': 'binary:logistic', 'eval\_metric': 'auc', 'silent': True

ms, folds=folds, model\_type='lgb', plot\_feature\_importance=True)

In [0]: xgb params = {'eta': 0.05, 'max depth': 3, 'subsample': 0.9, 'colsample bytre

**Customer Transaction Prediction** 

In Santander Customer Transaction Prediction competition we have a binary classification task. Train and test data have 200k samples each and we have 200 anonimyzed numerical columns. It

features. In fact this competition seems to be similar to another current competition: don't overfit! In

would be interesting to try good models without overfitting and knowing the meaning of the

• Trying various approaches to feature selection including taking top features from eli5;

RIHAD VARIAWA

this kernel I'll write the following things:

• Feature generation;

· Other things

· Hyperparameter optimization for models;

• EDA on the features and trying to get some insights;

• Using permutation importance to select most impactful features;

· Comparing various models: linear models, tree based models and others;

26-02-2019

**Overview**