

Feature selecture using target permutation

The notebook uses a procedure described in [this article](#).

Feature selection process using target permutation tests actual importance significance against the distribution of feature importances when fitted to noise (shuffled target).

The notebook implements the following steps :

- Create the null importances distributions : these are created fitting the model over several runs on a shuffled version of the target. This shows how the model can make sense of a feature irrespective of the target.
- Fit the model on the original target and gather the feature importances. This gives us a benchmark whose significance can be tested against the Null Importances Distribution
- for each feature test the actual importance:
 - Compute the probability of the actual importance wrt the null distribution. I will use a very simple estimation using occurrences while the article proposes to fit known distribution to the gathered data. In fact here I'll compute 1 - the proba so that things are in the right order.
 - Simply compare the actual importance to the mean and max of the null importances. This will give sort of a feature importance that allows to see major features in the dataset. Indeed the previous method may give us lots of ones.

For processing time reasons, the notebook will only cover application_train.csv but you can extend it as you wish.

Import a few packages

```
In [11]: import pandas as pd
import numpy as np

from sklearn.metrics import roc_auc_score
from sklearn.model_selection import KFold
import time
from lightgbm import LGBMClassifier
import lightgbm as lgb

import matplotlib.pyplot as plt
import matplotlib.gridspec as gridspec
import seaborn as sns
%matplotlib inline

import warnings
warnings.simplefilter('ignore', UserWarning)

import gc
gc.enable()
```

Read application_train

Read data and take care of categorical features

```
In [2]: data = pd.read_csv('../input/application_train.csv')

categorical_feats = [
    f for f in data.columns if data[f].dtype == 'object'
]

categorical_feats
for f in categorical_feats:
    data[f], _ = pd.factorize(data[f])
    # Set feature type as categorical
    data[f] = data[f].astype('category')
```

Create a scoring function

Coring function uses LightGBM in RandomForest mode fitted on the full dataset

```
In [3]: def get_feature_importances(data, shuffle, seed=None):
    # Gather real features
    train_features = [f for f in data if f not in ['TARGET', 'SK_ID_CURR']]
    # Go over fold and keep track of CV score (train and valid) and feature importances

    # Shuffle target if required
    y = data['TARGET'].copy()
    if shuffle:
        # Here you could as well use a binomial distribution
        y = data['TARGET'].copy().sample(frac=1.0)

    # Fit LightGBM in RF mode, yes it's quicker than sklearn RandomForest
    dtrain = lgb.Dataset(data[train_features], y, free_raw_data=False, silent=True)
    lgb_params = {
        'objective': 'binary',
        'boosting_type': 'rf',
        'subsample': 0.623,
        'colsample_bytree': 0.7,
        'num_leaves': 127,
        'max_depth': 8,
        'seed': seed,
        'bagging_freq': 1,
        'n_jobs': 4
    }

    # Fit the model
    clf = lgb.train(params=lgb_params, train_set=dtrain, num_boost_round=200, categorical_feature=categorical_feats)

    # Get feature importances
    imp_df = pd.DataFrame()
    imp_df['feature'] = list(train_features)
    imp_df['importance_gain'] = clf.feature_importance(importance_type='gain')
    imp_df['importance_split'] = clf.feature_importance(importance_type='split')
    imp_df['trn_score'] = roc_auc_score(y, clf.predict(data[train_features]))

    return imp_df
```

Build the benchmark for feature importance

The original paper does not talk about this but I think it makes sense to have a distribution of actual importances as well

```
In [4]: # Seed the unexpected randomness of this world
np.random.seed(123)
# Get the actual importance, i.e. without shuffling
actual_imp_df = get_feature_importances(data=data, shuffle=False)
```

```
In [5]: actual_imp_df.head()
```

Build Null Importances distribution

```
In [6]: null_imp_df = pd.DataFrame()
nb_runs = 80
import time
start = time.time()
dsp = ''
for i in range(nb_runs):
    # Get current run importances
    imp_df = get_feature_importances(data=data, shuffle=True)
    imp_df['run'] = i + 1
    # Concat the latest importances with the old ones
    null_imp_df = pd.concat([null_imp_df, imp_df], axis=0)
    # Erase previous message
    for l in range(len(dsp)):
        print('\b', end='', flush=True)
    # Display current run and time used
    spent = (time.time() - start) / 60
    dsp = 'Done with %4d of %4d (Spent %5.1f min)' % (i + 1, nb_runs, spent)
    print(dsp, end='', flush=True)
```

```
In [7]: null_imp_df.head()
```

Display distribution examples

A few plots are better than any words

```
In [25]: def display_distributions(actual_imp_df, null_imp_df, feature_):
    plt.figure(figsize=(13, 6))
    gs = gridspec.GridSpec(1, 2)
    # Plot Split importances
    ax = plt.subplot(gs[0, 0])
    a = ax.hist(null_imp_df.loc[actual_imp_df['feature'] == feature_, 'importance_split'].values, label='Null importances')
    ax.vlines(x=actual_imp_df.loc[actual_imp_df['feature'] == feature_, 'importance_split'].mean(),
              ymin=0, ymax=np.max(a[0]), color='r', linewidth=10, label='Real Target')
    ax.legend()
    ax.set_title('Split Importance of %s' % feature_.upper(), fontweight='bold')
    plt.xlabel('Null Importance (split) Distribution for %s ' % feature_.upper())
    # Plot Gain importances
    ax = plt.subplot(gs[0, 1])
    a = ax.hist(null_imp_df.loc[actual_imp_df['feature'] == feature_, 'importance_gain'].values, label='Null importances')
    ax.vlines(x=actual_imp_df.loc[actual_imp_df['feature'] == feature_, 'importance_gain'].mean(),
              ymin=0, ymax=np.max(a[0]), color='r', linewidth=10, label='Real Target')
    ax.legend()
    ax.set_title('Gain Importance of %s' % feature_.upper(), fontweight='bold')
    plt.xlabel('Null Importance (gain) Distribution for %s ' % feature_.upper())

In [26]: display_distributions(actual_imp_df=actual_imp_df, null_imp_df=null_imp_df, feature_='LIVINGAPARTMENT_S_AVG')

In [27]: display_distributions(actual_imp_df=actual_imp_df, null_imp_df=null_imp_df, feature_='CODE_GENDER')

In [28]: display_distributions(actual_imp_df=actual_imp_df, null_imp_df=null_imp_df, feature_='EXT_SOURCE_1')

In [29]: display_distributions(actual_imp_df=actual_imp_df, null_imp_df=null_imp_df, feature_='EXT_SOURCE_2')

In [30]: display_distributions(actual_imp_df=actual_imp_df, null_imp_df=null_imp_df, feature_='EXT_SOURCE_3')
```

From the above plot I believe the power of the exposed feature selection method is demonstrated. In particular it is well known that :

- Any feature sufficient variance can be used and made sense of by tree models. You can always find splits that help scoring better
- Correlated features have decaying importances once one of them is used by the model. The chosen feature will have strong importance and its correlated suite will have decaying importances

The current method allows to :

- Drop high variance features if they are not really related to the target
- Remove the decaying factor on correlated features, showing their real importance (or unbiased importance)

Score features

There are several ways to score features :

- Compute the number of samples in the actual importances that are away from the null importances recorded distribution.
- Compute ratios like Actual / Null Max, Actual / Null Mean, Actual Mean / Null Max

In a first step I will use the log actual feature importance divided by the 75 percentile of null distribution.

```
In [41]: feature_scores = []
for f in actual_imp_df['feature'].unique():
    f_nullimps_gain = null_imp_df.loc[actual_imp_df['feature'] == f, 'importance_gain'].values
    f_actimps_gain = actual_imp_df.loc[actual_imp_df['feature'] == f, 'importance_gain'].values
    gain_score = 100 * (f_nullimps < np.percentile(f_actimps, 25)).sum() / f_nullimps.size
    f_nullimps_split = null_imp_df.loc[actual_imp_df['feature'] == f, 'importance_split'].values
    f_actimps_split = actual_imp_df.loc[actual_imp_df['feature'] == f, 'importance_split'].values
    split_score = 100 * (f_nullimps < np.percentile(f_actimps, 25)).sum() / f_nullimps.size
    correlation_scores.append((f, split_score, gain_score))

scores_df = pd.DataFrame(feature_scores, columns=['feature', 'split_score', 'gain_score'])

plt.figure(figsize=(16, 16))
gs = gridspec.GridSpec(1, 2)
# Plot Split importances
ax = plt.subplot(gs[0, 0])
sns.barplot(x='split_score', y='feature', data=scores_df.sort_values('split_score', ascending=False).iloc[0:70], ax=ax)
ax.set_title('Feature scores wrt split importances', fontweight='bold', fontsize=14)
# Plot Gain importances
ax = plt.subplot(gs[0, 1])
sns.barplot(x='gain_score', y='feature', data=scores_df.sort_values('gain_score', ascending=False).iloc[0:70], ax=ax)
ax.set_title('Feature scores wrt gain importances', fontweight='bold', fontsize=14)
plt.tight_layout()
```

Save data

```
In [29]: null_imp_df.to_csv('null_importances_distribution_rf.csv')
actual_imp_df.to_csv('actual_importances_distribution_rf.csv')
```

Check the impact of removing uncorrelated features

Here I'll use a different metric to asses correlation to the target

```
In [49]: correlation_scores = []
for f in actual_imp_df['feature'].unique():
    f_nullimps = null_imp_df.loc[actual_imp_df['feature'] == f, 'importance_gain'].values
    f_actimps = actual_imp_df.loc[actual_imp_df['feature'] == f, 'importance_gain'].values
    gain_score = 100 * (f_nullimps < np.percentile(f_actimps, 25)).sum() / f_nullimps.size
    f_nullimps_split = null_imp_df.loc[actual_imp_df['feature'] == f, 'importance_split'].values
    f_actimps_split = actual_imp_df.loc[actual_imp_df['feature'] == f, 'importance_split'].values
    split_score = 100 * (f_nullimps < np.percentile(f_actimps, 25)).sum() / f_nullimps.size
    correlation_scores.append((f, split_score, gain_score))

corr_scores_df = pd.DataFrame(correlation_scores, columns=['feature', 'split_score', 'gain_score'])

fig = plt.figure(figsize=(16, 16))
gs = gridspec.GridSpec(1, 2)
# Plot Split importances
ax = plt.subplot(gs[0, 0])
sns.barplot(x='split_score', y='feature', data=corr_scores_df.sort_values('split_score', ascending=False).iloc[0:70], ax=ax)
ax.set_title('Feature scores wrt split importances', fontweight='bold', fontsize=14)
# Plot Gain importances
ax = plt.subplot(gs[0, 1])
sns.barplot(x='gain_score', y='feature', data=corr_scores_df.sort_values('gain_score', ascending=False).iloc[0:70], ax=ax)
ax.set_title('Feature scores wrt gain importances', fontweight='bold', fontsize=14)
plt.tight_layout()
plt.suptitle('Features' split and gain scores', fontweight='bold', fontsize=16)
fig.subplots_adjust(top=0.93)
```

Score feature removal for different thresholds

```
In [61]: def score_feature_selection(df=None, train_features=None, cat_feats=None, target=None):
    # Fit LightGBM
    dtrain = lgb.Dataset(df[train_features], target, free_raw_data=False, silent=True)
    lgb_params = {
        'objective': 'binary',
        'boosting_type': 'gbdt',
        'learning_rate': .1,
        'subsample': 0.8,
        'colsample_bytree': 0.8,
        'num_leaves': 31,
        'max_depth': -1,
        'seed': 13,
        'n_jobs': 4,
        'min_split_gain': .00001,
        'reg_alpha': .00001,
        'reg_lambda': .00001,
        'metric': 'auc'
    }

    # Fit the model
    hist = lgb.cv(
        params=lgb_params,
        train_set=dtrain,
        num_boost_round=2000,
        categorical_feature=cat_feats,
        nfold=5,
        stratified=True,
        shuffle=True,
        early_stopping_rounds=50,
        verbose_eval=0,
        seed=17
    )

    # Return the last mean / std values
    return hist['auc-mean'][-1], hist['auc-stdv'][-1]

# features = [f for f in data.columns if f not in ['SK_ID_CURR', 'TARGET']]
# score_feature_selection(df=data[features], train_features=features, target=data['TARGET'])

for threshold in [0, 10, 20, 30, 40, 50, 60, 70, 80, 90, 95, 99]:
    split_feats = [f for f, _score, _ in correlation_scores if _score >= threshold]
    split_cat_feats = [f for f, _score, _ in correlation_scores if (_score >= threshold) & (f in categorical_feats)]
    gain_feats = [f for f, _score in correlation_scores if _score >= threshold]
    gain_cat_feats = [f for f, _score in correlation_scores if (_score >= threshold) & (f in categorical_feats)]

    print('Results for threshold %3d' % threshold)
    split_results = score_feature_selection(df=data, train_features=split_feats, cat_feats=split_cat_feats, target=data['TARGET'])
    print('\t\t SPLIT : %6f +/- %6f' % (split_results[0], split_results[1]))
    gain_results = score_feature_selection(df=data, train_features=gain_feats, cat_feats=gain_cat_feats, target=data['TARGET'])
    print('\t\t GAIN : %6f +/- %6f' % (gain_results[0], gain_results[1]))
```