kaggle        Search                    Q        **Competitions**    **Datasets**    **Kernels**    **Discussion**    **Learn**

**Applying LightGBM to Titanic dataset**
Python notebook using data from Titanic: Machine Learning from Disaster · 1,125
views

⌃    **8**            Fork        **12**

**Version 11**
↺ 11 commits

**Notebook**

**Data**

**Log**

**Comments**

📘
Notebook

⊞
Data

📄
Log

💬
Comments

In [1]:

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import lightgbm as lgbm
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import LabelEncoder
from sklearn.metrics import accuracy_score, recall_score
, precision_score, f1_score
%matplotlib inline
```

# An implementation of LightGBM for the Titanic problem

Load and check the data:

In [2]:

```python
train_df = pd.read_csv('../input/train.csv')
test_df = pd.read_csv('../input/test.csv')
train_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 12 columns):
PassengerId    891 non-null int64
Survived       891 non-null int64
Pclass         891 non-null int64
Name           891 non-null object
Sex            891 non-null object
Age            714 non-null float64
SibSp          891 non-null int64
Parch          891 non-null int64
Ticket         891 non-null object
Fare           891 non-null float64
Cabin          204 non-null object
Embarked       889 non-null object
dtypes: float64(2), int64(5), object(5)
memory usage: 83.6+ KB
```

Feature engineering

- Add any extra features that may be handy
- Encode and categoricals stores as strings
- Drop features that look difficult to get use out of

Want to make sure that everything I do to the training set happens to the test, so makes sense to run as a loop

In [3]:

```python
# Not sure passenger ID is useful as a feature, but need
 to save it from the test set for the submission
test_passenger_ids = test_df.pop('PassengerId')
train_df.drop(['PassengerId'], axis=1, inplace=True)

# 'Embarked' is stored as letters, so fit a label encoder
to the train set to use in the loop
embarked_encoder = LabelEncoder()
embarked_encoder.fit(train_df['Embarked'].fillna('Null'
))

# Dataframes to work on
df_list = [train_df, test_df]

for df in df_list:

    # Record anyone travelling alone
    df['Alone'] = (df['SibSp'] == 0) & (df['Parch'] == 0
)

    # Transform 'Embarked'
    df['Embarked'].fillna('Null', inplace=True)
    df['Embarked'] = embarked_encoder.transform(df['Emba
rked'])

    # Transform 'Sex'
    df.loc[df['Sex'] == 'female','Sex'] = 0
    df.loc[df['Sex'] == 'male','Sex'] = 1
    df['Sex'] = df['Sex'].astype('int8')

    # Drop features that seem unusable. Save passenger id
s if test
    df.drop(['Name', 'Ticket', 'Cabin'], axis=1, inplace
=True)
```

Prep the training set for learning

In [4]:

```python
# Separate the label
```

```python
y = train_df.pop('Survived')

# Take a hold out set randomly
X_train, X_test, y_train, y_test = train_test_split(trai
n_df, y, test_size=0.2, random_state=42)

# Create an LGBM dataset for training
categorical_features = ['Alone', 'Sex', 'Pclass', 'Embar
ked']
train_data = lgbm.Dataset(data=X_train, label=y_train, c
ategorical_feature=categorical_features, free_raw_data=F
alse)

# Create an LGBM dataset from the test
test_data = lgbm.Dataset(data=X_test, label=y_test, cate
gorical_feature=categorical_features, free_raw_data=Fals
e)

# Finally, create a dataset for the FULL training data to
give us maximum amount of data to train on after
# performance has been calibrate
final_train_set = lgbm.Dataset(data=train_df, label=y,
                               categorical_feature=categ
orical_features, free_raw_data=False)
```

Define hyperparameters for LGBM

In [5]:

```python
lgbm_params = {
    'boosting': 'dart',          # dart (drop out trees)
 often performs better
    'application': 'binary',     # Binary classification
    'learning_rate': 0.05,       # Learning rate, control
s size of a gradient descent step
    'min_data_in_leaf': 20,      # Data set is quite smal
l so reduce this a bit
    'feature_fraction': 0.7,     # Proportion of features
in each boost, controls overfitting
    'num_leaves': 41,            # Controls size of tree
 since LGBM uses leaf wise splits
    'metric': 'binary_logloss',  # Area under ROC curve a
s the evaulation metric
    'drop_rate': 0.15
            }
```

Train the model

Use the held out training data to evaluate performance and early stop

Use the held out training data to evaluate performance and early stop to control overfitting

In [6]:

```python
evaluation_results = {}
clf = lgbm.train(train_set=train_data,
                 params=lgbm_params,
                 valid_sets=[train_data, test_data],
                 valid_names=['Train', 'Test'],
                 evals_result=evaluation_results,
                 num_boost_round=500,
                 early_stopping_rounds=100,
                 verbose_eval=20
                 )
optimum_boost_rounds = clf.best_iteration
```

```
/opt/conda/lib/python3.6/site-packages/ligh
tgbm/basic.py:1036: UserWarning: Using cate
gorical_feature in Dataset.
  warnings.warn('Using categorical_feature
in Dataset.')
/opt/conda/lib/python3.6/site-packages/ligh
tgbm/basic.py:681: UserWarning: categorical
_feature in param dict is overrided.
  warnings.warn('categorical_feature in par
am dict is overrided.')


Training until validation scores don't impr
ove for 100 rounds.
[20]    Train's binary_logloss: 0.507372
Test's binary_logloss: 0.54198
[40]    Train's binary_logloss: 0.433745
Test's binary_logloss: 0.479584
[60]    Train's binary_logloss: 0.413204
Test's binary_logloss: 0.461935
[80]    Train's binary_logloss: 0.393179
Test's binary_logloss: 0.444887
[100]   Train's binary_logloss: 0.385953
Test's binary_logloss: 0.442896
[120]   Train's binary_logloss: 0.373897
Test's binary_logloss: 0.438083
[140]   Train's binary_logloss: 0.36198 Tes
t's binary_logloss: 0.431068
[160]   Train's binary_logloss: 0.365415
Test's binary_logloss: 0.433233
[180]   Train's binary_logloss: 0.347998
Test's binary_logloss: 0.420761
[200]   Train's binary_logloss: 0.333363
Test's binary_logloss: 0.415733
[220]   Train's binary_logloss: 0.326614
Test's binary_logloss: 0.409778
[240]   Train's binary_logloss: 0.314509
Test's binary_logloss: 0.404197
```

```
[260]   Train's binary_logloss: 0.307755
Test's binary_logloss: 0.401085
[280]   Train's binary_logloss: 0.296964
Test's binary_logloss: 0.399032
[300]   Train's binary_logloss: 0.292589
Test's binary_logloss: 0.399077
[320]   Train's binary_logloss: 0.286967
Test's binary_logloss: 0.399104
[340]   Train's binary_logloss: 0.284326
Test's binary_logloss: 0.398559
[360]   Train's binary_logloss: 0.287643
Test's binary_logloss: 0.400735
[380]   Train's binary_logloss: 0.282073
Test's binary_logloss: 0.399788
[400]   Train's binary_logloss: 0.276178
Test's binary_logloss: 0.399436
Early stopping, best iteration is:
[301]   Train's binary_logloss: 0.291464
Test's binary_logloss: 0.397995
```

Visualise training performance

In [7]:

```python
fig, axs = plt.subplots(1, 2, figsize=[15, 4])

# Plot the log loss during training
axs[0].plot(evaluation_results['Train']['binary_logloss'
], label='Train')
axs[0].plot(evaluation_results['Test']['binary_logloss'
], label='Test')
axs[0].set_ylabel('Log loss')
axs[0].set_xlabel('Boosting round')
axs[0].set_title('Training performance')
axs[0].legend()

# Plot feature importance
importances = pd.DataFrame({'features': clf.feature_name
(),
                            'importance': clf.feature_im
portance()}).sort_values('importance', ascending=False)
axs[1].bar(x=np.arange(len(importances)), height=importa
nces['importance'])
axs[1].set_xticks(np.arange(len(importances)))
axs[1].set_xticklabels(importances['features'])
axs[1].set_ylabel('Feature importance (# times used to s
plit)')
axs[1].set_title('Feature importance')

plt.show()
```

Examine model performance

Accuracy score can often be misleading for classifiers, so have a look at the precision and recall too

In [8]:

```
preds = np.round(clf.predict(X_test))
print('Accuracy score = \t {}'.format(accuracy_score(y_test, preds)))
print('Precision score = \t {}'.format(precision_score(y_test, preds)))
print('Recall score =   \t {}'.format(recall_score(y_test, preds)))
print('F1 score =       \t {}'.format(f1_score(y_test, preds)))
```

```
Accuracy score =        0.8324022346368715
```

Data

Data Sources

⌄  🏆  Titanic: Machi...

   ▦  ...  418 x 2

   ▦  te  418 x 11

   ▦  t  891 x 12

## Titanic: Machine Learning from Disaster

**Start here! Predict survival on the Titanic and get familiar with ML basics**

Last Updated: 7 years ago

### About this Competition

### Overview

The data has been split into two groups:

- training set (train.csv)
- test set (test.csv)

**The training set** should be used to build your machine learning models. For the training set, we provide the outcome (also known as the "ground truth") for each passenger. Your model will be based on "features" like passengers'

gender and class. You can also use feature engineering to create new features.

**The test set** should be used to see how well your model performs on unseen data. For the test set, we do not provide the ground truth for each passenger. It is your job to predict these outcomes. For each passenger in the test set, use the model you trained to predict whether or not they survived the sinking of the Titanic.

We also include **gender_submission.csv**, a set of predictions that assume all and only female passengers survive, as an example of what a submission file should look like.

## Run Info

| | | | |
|---|---|---|---|
| Succeeded | **True** | Run Time | **24.1 seconds** |
| Exit Code | **0** | Queue Time | **0 seconds** |
| Docker Image Name | kaggle/python(Dockerfile) | Output Size | **0** |
| Timeout Exceeded | **False** | Used All Space | **False** |
| Failure Message | | | |

## Log

**Download Log**

```
Time    Line #  Log Message
5.1s        1   [NbConvertApp] Converting notebook
                script.ipynb to html
5.1s        2   [NbConvertApp] Executing notebook with kernel:
                python3
9.0s        3   [LightGBM] [Info] Number of positive: 268,
                number of negative: 444
9.0s        4   [LightGBM] [Info] Total Bins 199
                [LightGBM] [Info] Number of data: 712, number
                of used features: 8
                [LightGBM] [Warning] No further splits with
                positive gain, best gain: -inf
                [LightGBM] [Warning] No further splits with
                positive gain, best gain: -inf
                [LightGBM] [Warning] No further splits with
                positive gain, best gain: -inf
                [LightGBM] [Warning] No further splits with
                positive gain, best gain: -inf
                [LightGBM] [Warning] No further splits with
                positive gain, best gain: -inf
9.1s        5   [LightGBM] [Warning] No further splits with
                positive gain, best gain: -inf
9.1s        6   [LightGBM] [Warning] No further splits with
                positive gain, best gain: -inf
                [LightGBM] [Warning] No further splits with
                positive gain, best gain: -inf
                [LightGBM] [Warning] No further splits with
                positive gain, best gain: -inf
                [LightGBM] [Warning] No further splits with
                positive gain, best gain: -inf
                [LightGBM] [Warning] No further splits with
                positive gain, best gain: -inf
                [LightGBM] [Warning] No further splits with
                positive gain, best gain: -inf
                [LightGBM] [Warning] No further splits with
                positive gain, best gain: -inf
                [LightGBM] [Warning] No further splits with
                positive gain, best gain: -inf
                [LightGBM] [Warning] No further splits with
                positive gain, best gain: -inf
```

```
                              [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
                              [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
                              [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
                              [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
                              [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
                              [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
        9.2s      7           [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
        9.2s      8           [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
                              [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
                              [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
                              [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
                              [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
                              [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
                              [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
                              [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
                              [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
                              [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
        9.3s      9           [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
        9.3s     10           [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
                              [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
                              [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
                              [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
                              [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
                              [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
                              [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
                              [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
        9.4s     11           [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
        9.4s     12           [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
                              [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
                              [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
                              [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
                              [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
                              [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
                              [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
                              [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
        9.5s     13           [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
        9.5s     14           [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
                              [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
                              [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
                              [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
                              [LightGBM] [Warning] No further splits with
                              positive gain, best gain: -inf
```

|       |    |                                                                                     |
|-------|----|-------------------------------------------------------------------------------------|
|       |    | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |
|       |    | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |
|       |    | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |
| 9.6s  | 15 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |
| 9.6s  | 16 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |
|       |    | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |
|       |    | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |
|       |    | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |
|       |    | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |
|       |    | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |
|       |    | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |
| 9.7s  | 17 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |
| 9.7s  | 18 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |
|       |    | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |
|       |    | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |
|       |    | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |
|       |    | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |
|       |    | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |
| 9.8s  | 19 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |
| 9.8s  | 20 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |
|       |    | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |
|       |    | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |
|       |    | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |
|       |    | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |
| 9.9s  | 21 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |
| 9.9s  | 22 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |
|       |    | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |
|       |    | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |
|       |    | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |
|       |    | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |
|       |    | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |
|       |    | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |
| 10.0s | 23 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |
| 10.0s | 24 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |
|       |    | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |
|       |    | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |
|       |    | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |
|       |    | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |
|       |    | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |
|       |    | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf           |

| | | |
|---|---|---|
| 10.1s | 25 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf |
| 10.1s | 26 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf |
| 10.2s | 27 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf |
| 10.2s | 28 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf |
| 10.3s | 29 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf |
| 10.3s | 30 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf |
| 10.4s | 31 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf |
| 10.4s | 32 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf |
| 10.5s | 33 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf |
| 10.5s | 34 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf |
| 10.6s | 35 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf |
| 10.6s | 36 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf |
| 10.7s | 37 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf |
| 10.7s | 38 | [LightGBM] [Warning] No further splits with |

```
                        positive gain, best gain: -inf
                        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
                        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
                        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
                        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
                        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
                        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
                        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
      10.8s   39        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
      10.8s   40        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
      10.9s   41        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
      10.9s   42        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
                        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
                        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
                        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
      11.0s   43        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
      11.0s   44        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
                        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
                        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
                        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
      11.1s   45        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
      11.1s   46        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
                        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
                        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
                        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
                        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
                        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
                        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
      11.2s   47        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
      11.2s   48        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
                        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
                        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
                        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
                        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
      11.3s   49        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
      11.3s   50        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
                        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
                        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
      11.4s   51        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
      11.4s   52        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
                        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
                        [LightGBM] [Warning] No further splits with
                        positive gain, best gain: -inf
                        [LightGBM] [Warning] No further splits with
```

```
                               positive gain, best gain: -inf
                               [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
                               [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
        11.5s    53            [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
        11.5s    54            [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
        11.6s    55            [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
        11.6s    56            [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
                               [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
                               [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
        11.7s    57            [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
        11.7s    58            [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
                               [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
                               [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
        11.8s    59            [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
        11.8s    60            [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
                               [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
        11.9s    61            [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
        11.9s    62            [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
                               [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
        12.0s    63            [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
        12.0s    64            [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
                               [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
                               [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
        12.1s    65            [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
        12.1s    66            [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
                               [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
                               [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
        12.2s    67            [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
        12.2s    68            [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
                               [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
                               [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
                               [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
        12.3s    69            [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
        12.3s    70            [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
                               [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
                               [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
        12.4s    71            [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
        12.4s    72            [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
                               [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
                               [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
        12.5s    73            [LightGBM] [Warning] No further splits with
                               positive gain, best gain: -inf
```

| | | |
|---|---|---|
| 12.5s | 74 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf |
| 12.6s | 75 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf |
| 12.6s | 76 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf |
| 12.7s | 77 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf |
| 12.7s | 78 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf |
| 12.8s | 79 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf |
| 12.8s | 80 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf |
| 12.9s | 81 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf |
| 12.9s | 82 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf |
| 13.0s | 83 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf |
| 13.0s | 84 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf |
| 13.1s | 85 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf |
| 13.1s | 86 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf |
| 13.2s | 87 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf |
| 13.2s | 88 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf |
| 13.3s | 89 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf |
| 13.3s | 90 | [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf [LightGBM] [Warning] No further splits with positive gain, best gain: -inf |

```
                    [LightGBM] [Warning] No further splits with
                    positive gain, best gain: -inf
                    [LightGBM] [Warning] No further splits with
                    positive gain, best gain: -inf
    13.4s      91   [LightGBM] [Warning] No further splits with
                    positive gain, best gain: -inf
    13.4s      92   [LightGBM] [Warning] No further splits with
                    positive gain, best gain: -inf
    13.5s      93   [LightGBM] [Warning] No further splits with
                    positive gain, best gain: -inf
    13.5s      94   [LightGBM] [Warning] No further splits with
                    positive gain, best gain: -inf
                    [LightGBM] [Warning] No further splits with
                    positive gain, best gain: -inf
    13.6s      95   [LightGBM] [Warning] No further splits with
                    positive gain, best gain: -inf
    13.6s      96   [LightGBM] [Warning] No further splits with
                    positive gain, best gain: -inf
                    [LightGBM] [Warning] No further splits with
                    positive gain, best gain: -inf
    13.7s      97   [LightGBM] [Warning] No further splits with
                    positive gain, best gain: -inf
    13.7s      98   [LightGBM] [Warning] No further splits with
                    positive gain, best gain: -inf
    13.8s      99   [LightGBM] [Warning] No further splits with
                    positive gain, best gain: -inf
    13.8s     100   [LightGBM] [Warning] No further splits with
                    positive gain, best gain: -inf
                    [LightGBM] [Warning] No further splits with
                    positive gain, best gain: -inf
                    [LightGBM] [Warning] No further splits with
                    positive gain, best gain: -inf
                    [LightGBM] [Warning] No further splits with
                    positive gain, best gain: -inf
    23.6s     259
    23.6s     260   ...
    23.6s     261   Complete. Exited with code 0.
```

## Comments (3)

Sort by

All Comments ⌄          Hotness ⌄

Click here to enter a comment...

楼上膜一下 • Posted on Latest Version • a year ago • Options • Reply  ⌃ 1 ⌄

It's a good kernel to help me learning usage of LightGBM and sklearn, much appreciate!

Ilya Khris... • Posted on Latest Version • 6 months ago • Options • Reply  ⌃ 0 ⌄

A good kernel to play with LGBM! Thank you!

Adriano ... • Posted on Latest Version • 10 days ago • Options • Reply  ⌃ 0 ⌄

Thanks. That's what I am looking for.