

Assignment 4

October 12, 2020

*You are currently looking at **version 1.1** of this notebook. To download notebooks and datafiles, as well as get help on Jupyter notebooks in the Coursera platform, visit the [Jupyter Notebook FAQ](#) course resource.*

0.1 Assignment 4 - Understanding and Predicting Property Maintenance Fines

This assignment is based on a data challenge from the Michigan Data Science Team ([MDST](#)).

The Michigan Data Science Team ([MDST](#)) and the Michigan Student Symposium for Interdisciplinary Statistical Sciences ([MSSISS](#)) have partnered with the City of Detroit to help solve one of the most pressing problems facing Detroit - blight. [Blight violations](#) are issued by the city to individuals who allow their properties to remain in a deteriorated condition. Every year, the city of Detroit issues millions of dollars in fines to residents and every year, many of these fines remain unpaid. Enforcing unpaid blight fines is a costly and tedious process, so the city wants to know: how can we increase blight ticket compliance?

The first step in answering this question is understanding when and why a resident might fail to comply with a blight ticket. This is where predictive modeling comes in. For this assignment, your task is to predict whether a given blight ticket will be paid on time.

All data for this assignment has been provided to us through the [Detroit Open Data Portal](#). **Only the data already included in your Coursera directory can be used for training the model for this assignment.** Nonetheless, we encourage you to look into data from other Detroit datasets to help inform feature creation and model selection. We recommend taking a look at the following related datasets:

- [Building Permits](#)
- [Trades Permits](#)
- [Improve Detroit: Submitted Issues](#)
- [DPD: Citizen Complaints](#)
- [Parcel Map](#)

We provide you with two data files for use in training and validating your models: train.csv and test.csv. Each row in these two files corresponds to a single blight ticket, and includes information about when, why, and to whom each ticket was issued. The target variable is compliance, which is True if the ticket was paid early, on time, or within one month of the hearing data, False

if the ticket was paid after the hearing date or not at all, and Null if the violator was found not responsible. Compliance, as well as a handful of other variables that will not be available at test-time, are only included in train.csv.

Note: All tickets where the violators were found not responsible are not considered during evaluation. They are included in the training set as an additional source of data for visualization, and to enable unsupervised and semi-supervised approaches. However, they are not included in the test set.

File descriptions (Use only this data for training your model!)

readonly/train.csv - the training set (all tickets issued 2004-2011)

readonly/test.csv - the test set (all tickets issued 2012-2016)

readonly/addresses.csv & readonly/latlons.csv - mapping from ticket id to addresses

Note: misspelled addresses may be incorrectly geolocated.

Data fields

train.csv & test.csv

ticket_id - unique identifier for tickets

agency_name - Agency that issued the ticket

inspector_name - Name of inspector that issued the ticket

violation_code - Code of the person/organization that the ticket was issued to

violation_street_number, violation_street_name, violation_zip_code - Address where

mailing_address_str_number, mailing_address_str_name, city, state, zip_code, non_us

ticket_issued_date - Date and time the ticket was issued

hearing_date - Date and time the violator's hearing was scheduled

violation_code, violation_description - Type of violation

disposition - Judgment and judgement type

fine_amount - Violation fine amount, excluding fees

admin_fee - \$20 fee assigned to responsible judgments

state_fee - \$10 fee assigned to responsible judgments late_fee - 10% fee assigned to responsible judgments discount_amount - discount applied, if any clean_up_cost - DPW clean-up or graffiti removal cost judgment_amount - Sum of all fines and fees graffiti_status - Flag for graffiti violations

train.csv only

payment_amount - Amount paid, if any

payment_date - Date payment was made, if it was received

payment_status - Current payment status as of Feb 1 2017

balance_due - Fines and fees still owed

collection_status - Flag for payments in collections

compliance [target variable for prediction]

Null = Not responsible

0 = Responsible, non-compliant

1 = Responsible, compliant

compliance_detail - More information on why each ticket was marked compliant or non

0.2 Evaluation

Your predictions will be given as the probability that the corresponding blight ticket will be paid on time.

The evaluation metric for this assignment is the Area Under the ROC Curve (AUC).

Your grade will be based on the AUC score computed for your classifier. A model which with an AUROC of 0.7 passes this assignment, over 0.75 will receive full points. ____

For this assignment, create a function that trains a model to predict blight ticket compliance in Detroit using `readonly/train.csv`. Using this model, return a series of length 61001 with the data being the probability that each corresponding ticket from `readonly/test.csv` will be paid, and the index being the `ticket_id`.

Example:

```
ticket_id
284932    0.531842
285362    0.401958
285361    0.105928
285338    0.018572
...
376499    0.208567
376500    0.818759
369851    0.018528
Name: compliance, dtype: float32
```

0.2.1 Hints

- Make sure your code is working before submitting it to the autograder.
- Print out your result to see whether there is anything weird (e.g., all probabilities are the same).
- Generally the total runtime should be less than 10 mins. You should NOT use Neural Network related classifiers (e.g., `MLPClassifier`) in this question.
- Try to avoid global variables. If you have other functions besides `blight_model`, you should move those functions inside the scope of `blight_model`.
- Refer to the pinned threads in Week 4's discussion forum when there is something you could not figure it out.

```
In [35]: import pandas as pd
import numpy as np

def blight_model():

    from sklearn.preprocessing import LabelEncoder
    from sklearn.ensemble import RandomForestRegressor
    from sklearn.model_selection import train_test_split, GridSearchCV
    from sklearn.metrics import roc_auc_score
```

```

# load data
train = pd.read_csv('train.csv', encoding='cp1252')
test = pd.read_csv('test.csv')
addresses = pd.read_csv('addresses.csv')
latlons = pd.read_csv('latlons.csv')
# pre-processing

# drop all rows with Null compliance
train=train[train['compliance'].notnull()]

# drop all rows not in the U.S
train = train[train.country == 'USA']
test = test[test.country == 'USA']
# merge latlons and addresses with data
train = pd.merge(train, pd.merge(addresses, latlons, on='address'), on='address')
test = pd.merge(test, pd.merge(addresses, latlons, on='address'), on='address')
# drop all unnecessary columns
train.drop(['agency_name', 'inspector_name', 'violation_name', 'non_us_
            'grafitti_status', 'state_fee', 'admin_fee', 'ticket_issue
            # columns not available in test
            'payment_amount', 'balance_due', 'payment_date', 'payment
            'collection_status', 'compliance_detail',
            # address related columns
            'violation_zip_code', 'country', 'address', 'violation_str
            'violation_street_name', 'mailing_address_str_number', 'ma
            'city', 'state', 'zip_code', 'address'], axis=1, inplace=True)
# discretizing relevant columns
le = LabelEncoder()
train['disposition'] = le.fit_transform(train['disposition'])
test['disposition'] = le.fit_transform(test['disposition'])
le = LabelEncoder()
train['violation_code'] = le.fit_transform(train['violation_code'])
test['violation_code'] = le.fit_transform(test['violation_code'])
train['lat'] = train['lat'].fillna(train['lat'].mean())
train['lon'] = train['lon'].fillna(train['lon'].mean())
test['lat'] = test['lat'].fillna(test['lat'].mean())
test['lon'] = test['lon'].fillna(test['lon'].mean())
X=train.ix[:,train.columns!='compliance']
y=train['compliance']
test=test[X.columns]

# train the model
X_train, X_test, y_train, y_test = train_test_split(X,y)
regr_rf = RandomForestRegressor()
grid_values = {'n_estimators': [10, 100], 'max_depth': [None, 30]}
grid_clf_auc = GridSearchCV(regr_rf, param_grid=grid_values, scoring='roc_auc')
grid_clf_auc.fit(X_train, y_train)
#print('Grid best parameter (max. AUC): ', grid_clf_auc.best_params_)

```

```
#print('Grid best score (AUC): ', grid_clf_auc.best_score_)
```

```
return pd.DataFrame({'test':test['ticket_id'],'prediction':grid_clf_auc.predict(test['ticket_id'])})
```

```
In [32]: blight_model()
```

```
/opt/conda/lib/python3.6/site-packages/IPython/core/interactiveshell.py:2827: DtypeWarning:
  if self.run_code(code, result):
```

```
Grid best parameter (max. AUC): {'max_depth': 30, 'n_estimators': 100}
```

```
Grid best score (AUC): 0.810423621896
```

```
Out[32]:
```

	prediction
--	------------

test	
284932	0.58
285362	0.15
285361	0.23
285338	0.27
285346	0.30
285345	0.27
285347	0.33
285342	0.97
285530	0.35
284989	0.50
285344	0.39
285343	0.28
285340	0.38
285341	0.31
285349	0.29
285348	0.26
284991	0.49
285532	0.41
285406	0.30
285001	0.33
285006	0.34
285405	0.15
285337	0.13
285496	0.51
285497	0.55
285378	0.19
285589	0.19
285585	0.26
285501	0.58
285581	0.15
...	...
376367	0.21
376366	0.26

376362	0.41
376363	0.39
376365	0.21
376364	0.26
376228	0.32
376265	0.28
376286	0.98
376320	0.62
376314	0.67
376327	0.98
376385	0.96
376435	0.83
376370	0.78
376434	0.23
376459	0.52
376478	0.11
376473	0.58
376484	0.35
376482	0.16
376480	0.17
376479	0.17
376481	0.17
376483	0.45
376496	0.14
376497	0.14
376499	0.58
376500	0.58
369851	0.91

[61001 rows x 1 columns]