



PROJECT

Identify Fraud from Enron Email

A part of the Data Analyst Nanodegree Program

PROJECT REVIEW

CODE REVIEW 7

NOTES

SHARE YOUR ACCOMPLISHMENT!  

Meets Specifications

Dear Student,

Great work!. Yours is a strong project!. You demonstrate a good understanding of Machine Learning and your report is written in explanatory terms allowing your audience to follow and understand the work done.

Congratulations on passing your exam and stay Udacious! 

please keep the tips and references coming on what and how to improve further. decided to put the code refactoring to another time :p

I left some tips on the review that I hope you find interesting!, in case you feel something is missed, please reach us, we'll be glad to help!!

Quality of Code

Code reflects the description in the answers to questions in the writeup. i.e. code performs the functions documented in the writeup and the writeup clearly specifies the final analysis strategy.

Your written response perfectly describes your strategy, includes the required level of detail when required and your code includes all the processes mentioned. Machine learning is not just about building good models, it is also about communicating results to your audience in a clear and direct way. You achieve both goals. Well done!

As a suggestion, this project is a great opportunity for you to create a new repository in Github that becomes part of your online portfolio and allow potential employers to review your work, in case you are not familiar with Github, this is a [great post](#) and an [Udacity Course](#) for deeper understanding. This report defines your credentials, so it is important that you put special attention not just to the technical side of the project but also the communications side since this is a critical characteristic for any data scientist. For your reference, check this [Kaggle post](#) for further reference, as you can see this is really a hot topic in the data science world! 🤖

`poi_id.py` can be run to export the dataset, list of features and algorithm, so that the final algorithm can be checked easily using `tester.py`.

All required `.pk1` files are included and `poi_id.py` worked without problems.

Understanding the Dataset and Question

Student response addresses the most important characteristics of the dataset and uses these characteristics to inform their analysis. Important characteristics include:

- total number of data points
- allocation across classes (POI/non-POI)
- number of features used
- are there features with many missing values? etc.

Your report includes a description of the main characteristics of the dataset.

Note these numbers are particularly important since they describe the main dataset characteristics:

1. The small data set is why the tester.py file uses StratifiedShuffleSplit instead of a simpler cross-validation method such as TrainTestSplit. StratifiedShuffleSplit will make randomly chosen training and test sets multiple times and average the results over all the tests.
2. The data is unbalanced with many more non-POIs than POIs. StratifiedShuffleSplit also makes sure that the ratio of non-POI:POI is the same in the training and test sets as it was in the larger data set.
3. The unbalanced data is also why we use precision and recall instead of accuracy as our evaluation metric.

For your reference, some [techniques](#) to handle unbalanced datasets and this [repo](#) with plenty of tools.

Student response identifies outlier(s) in the financial data, and explains how they are removed or otherwise handled.

Well done by identifying TOTAL, LOCKHART EUGENE E and THE TRAVEL AGENCY IN THE PARK. These are the three most prominent outliers in the dataset.

As a suggestion, the dataset comes in a json format that makes it difficult to handle and explore. My suggestion is to use [Pandas](#):

```
data_dict = pickle.load(open("final_project_dataset.pkl", "r") )
###creating dataframe from dictionary - pandas
df = pandas.DataFrame.from_dict(data_dict, orient='index', dtype=np.float)
print df.describe().loc[:, ['salary', 'bonus']]
```

Check the [docs](#) for more information on how pandas can read data from different sources.

Optimize Feature Selection/Engineering

At least one new feature is implemented. Justification for that feature is provided in the written response. The effect of that feature on final algorithm performance is tested or its strength is compared to other features in feature selection. The student is not required to include their new feature in their final feature set.

Good work engineering your features, including your reasons and testing their impact over your classifier.

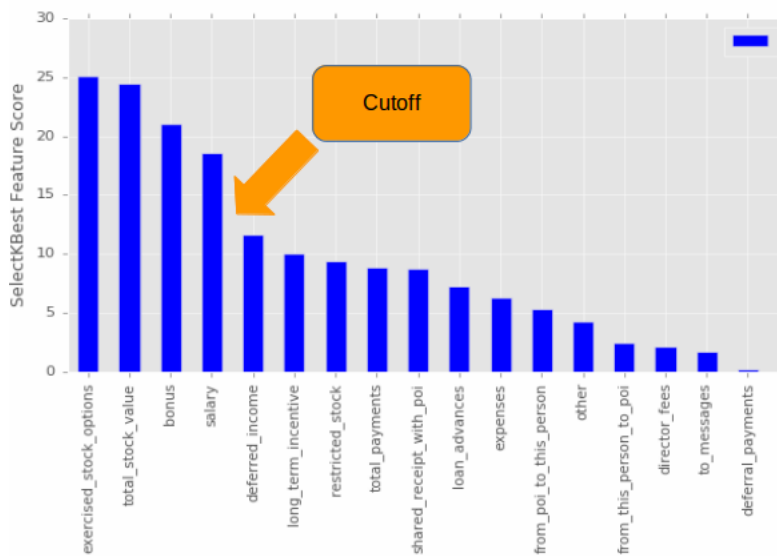
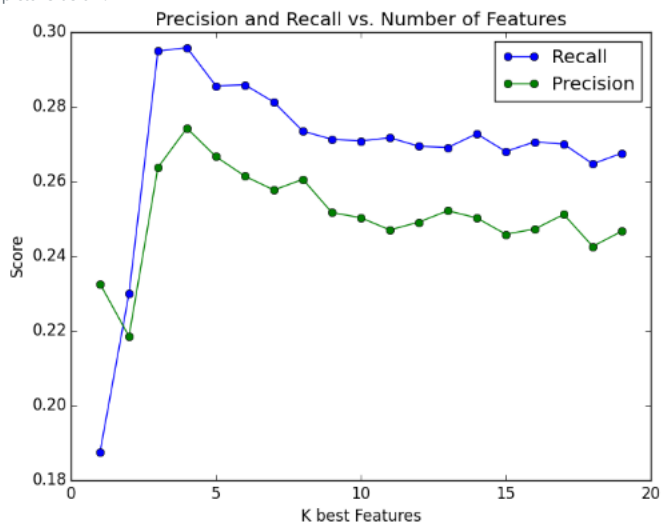
Univariate or recursive feature selection is deployed, or features are selected by hand (different combinations of features are attempted, and the performance is documented for each one). Features that are selected are reported and the number of features selected is justified. For an algorithm that supports getting the feature importances (e.g. decision tree) or feature scores (e.g. SelectKBest), those are documented as well.

Good work in this section. Note feature selection process is key in Machine Learning problems, the idea behind it is that you want to have the minimum number of features that capture trends and patterns in your data. Your machine learning algorithm is just going to be as good as the features you put into it. For that reason, this is definitely a critical step into any ML problem and the methodology deployed must be scientific and exhaustive without room for intuition.

Having said so, your methodology is scientific and exhaustive as required, SelectKBest ranks the level of correlation between features and the target label, and by testing the different feature sets generated for each value of K in its range, your methodology is exploratory. Now, since you repeated this

process for different classifiers, you identified which is the best performing classifier for this dataset, that goes beyond the project expectations, but definitely, it is worth to do in the real world. Well done!

As a suggestion, when working with larger datasets, since training and testing your models for each feature selection is computationally expensive, a good decision is to use the feature scores (ordered by ranking) to determine where the feature scores drop and so determine the cutoff as explained in the picture below.



An alternative (when working with small/medium dataset sizes) to identify the optimal value of `k` is to use Pipelines and SearchGridCV. [Here](#) is a good example where multiple feature selection methods are combined in a single pipeline that is later tuned in a GridCV object. Note this approximation brings two drawbacks, firstly it requires a lot of computational resources and secondly, it does not give you a deep understanding of how features relate with the label. For your reference:

NOTE: This process is computationally expensive!

```
clf_params= {
    'kbest_k' : range(1,len(feature_list)),
    'clf_C': [1e-5, 1, 10, 1e2, 1e5],
    'clf_gamma': [0.0],
    'clf_kernel': ['linear', 'poly', 'rbf'],
    'clf_tol': [1e-1, 1e-2, 1e-5],
    'clf_class_weight': [{True: 12, False: 1},
                        {True: 10, False: 1},
                        'balanced', None]
```

```

    }

#For this Pipeline:
pipe = Pipeline(steps=[('minmaxer', MinMaxScaler()), ('kbest', SelectKBest()), ('clf', SVC())])
cv = cross_validation.StratifiedShuffleSplit(n_splits = 50, random_state = 42)
a_grid_search = GridSearchCV(pipe, param_grid = clf_params, cv = cv, scoring = 'recall')
a_grid_search.fit(features, labels)

# pick a winner
best_clf = a_grid_search.best_estimator_
print best_clf

```

If algorithm calls for scaled features, feature scaling is deployed.

Good understanding of when scaled features are required. 🍌

Pick and Tune an Algorithm

At least two different algorithms are attempted and their performance is compared, with the best performing one used in the final analysis.

Good work testing several algorithms and comparing their performance in terms of precision&recall.

Response addresses what it means to perform parameter tuning and why it is important.

Good understanding of tuning.

At least one important parameter tuned with at least 3 settings investigated systematically, or any of the following are true:

- GridSearchCV used for parameter tuning
- Several parameters tuned
- Parameter tuning incorporated into algorithm selection (i.e. parameters tuned for more than one algorithm, and best algorithm-tune combination selected for final analysis).

Good work using [SearchGridCV](#), and as a suggestion to enhance its performance, use the `cv` parameter you can pass a cross-validation object to validate your search results that best adapt to your dataset characteristics.

For example:

```

# Set up cross validator (will be used for tuning all classifiers)
cv = cross_validation.StratifiedShuffleSplit(n_splits=100, random_state = 42)
a_grid_search = GridSearchCV(clf, param_grid = clf_params, cv = cv, scoring = 'recall')

```

This GridCV object validates algorithm performance using a cross-validation object that best adapts to the dataset characteristics and searches for those parameters that maximize recall.

For your reference, note you can also [visualize your grid results](#).

Validate and Evaluate

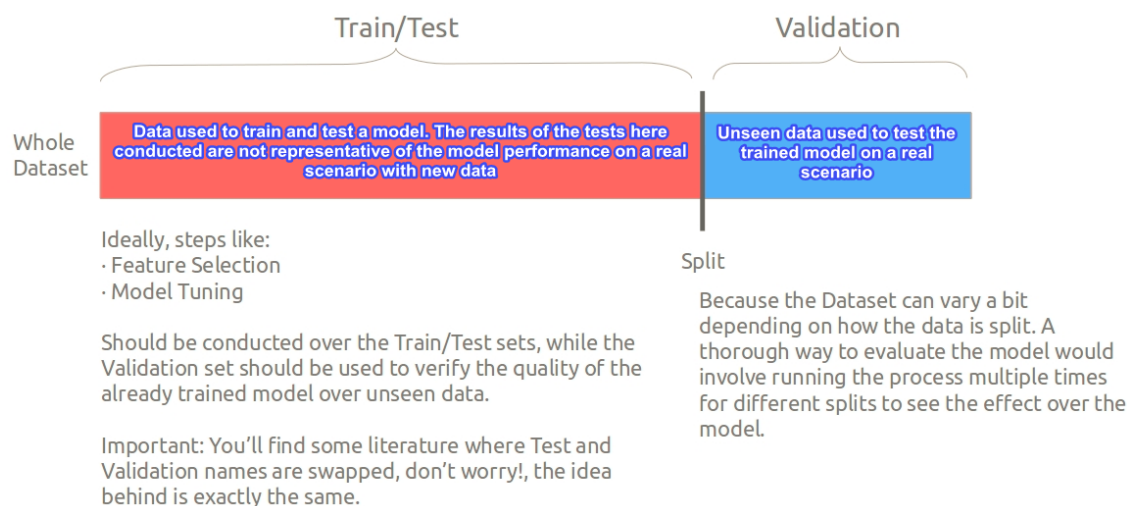
At least two appropriate metrics are used to evaluate algorithm performance (e.g. precision and recall), and the student articulates what those metrics measure in context of the project task.

Good work using precision&recall to evaluate your classifier. Also, good definition of both in terms of POI detection.

As a side comment: Note accuracy is not a good score in this case since the dataset in this project is very small and the ratio of negatives to positives is highly skewed (18 POIs out of 146 total cases), a classifier that predicted only non-POIs as output, would obtain an accuracy score of 87.4%. In this regard, accuracy is not a strong validation metric for our classifier. For your reference, [this interesting post](#).

Response addresses what validation is and why it is important.

Good understanding of validation. As a side comment, since this dataset is really really small, validation is performed just using the test set. In the real world, you would split your dataset in train/test/validation in order to validate your trained&tuned model against completely unseen data. For further reference, check [this excellent answer](#).



Performance of the final algorithm selected is assessed by splitting the data into training and testing sets or through the use of cross validation, noting the specific type of validation performed.

For this particular problem, a stratified shuffle split is preferred for a number of reasons. Since the dataset is small, a shuffle split will randomly assign entries to test and training sets in a fold. However, the stratification preserves the percentage in the target class as in the complete set. This is particularly important Validation for our investigation as we have such a small percentage of Poles - if we did not stratify we could easily have a test or train set which contained no entries that were Poles. If we had no entries that were Poles in either a test or train set then the model would perform very badly.

When tester.py is used to evaluate performance, precision and recall are both at least 0.3.

Well done!. Your classifier is over 0.3 for precision&recall scores. 🍌

[↓ DOWNLOAD PROJECT](#)

7 CODE REVIEW COMMENTS



RETURN TO PATH

Rate this review

[Student FAQ](#)

