

```
In [1]: #####
# Purpose: Identify a claim can be fast-tracked using Machine Learning.
# Created by: Suriya Mohan
# Created on: 12-Nov-2016
#####

# Import python - scikit-learn machine learning packages.
import numpy as np;
import pandas as pd;
from pandas import Series, DataFrame;
from sklearn.preprocessing import LabelEncoder
from sklearn.preprocessing import OneHotEncoder
from sklearn.preprocessing import StandardScaler
from sklearn.preprocessing import Imputer
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.svm import SVC
from sklearn.ensemble import RandomForestClassifier
from sklearn.ensemble import GradientBoostingClassifier
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import learning_curve
from sklearn.model_selection import validation_curve
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import confusion_matrix
from sklearn.pipeline import Pipeline
%matplotlib inline
```

```
In [2]: # Load the data in CSV file into panda dataframe.
#df_claims = pd.read_csv('C:\Users\gbu4moh\Desktop\ML - FAST TRACK\owning_adjuste
df_claims = pd.read_csv('C:\Users\gbu4moh\Desktop\ML - FAST TRACK\ML_FASTTRACK_DA
```

```
In [3]: def claim_labelEncode(df, column_name):
    ...
    Purpose: Function used to convert the nominal features (fields with string va
    Input: Dataframe and column name to encode
    Output: Modify the passed dataframe by adding encoded field to end and drop t
    ...
    df[column_name].fillna(value=df[column_name].value_counts().idxmax(), inplace

    le = LabelEncoder()
    encoder = le.fit_transform(df[column_name].values)
    df.insert(0, 'code' + column_name, encoder.astype(int))
    df.drop(column_name, axis=1, inplace=True)
```

```
In [4]: def claim_oneHotEncode(df,column_array):
        ...
        Purpose: Function used to hot encode fields - example field having 3 distinct
        Input: Dataframe and column array
        Output: Returns the data frame with encoded values.
        ...
        ohe = OneHotEncoder(categorical_features=column_array,sparse=False)
        df_ohe = ohe.fit_transform(df)
        return df_ohe
```

```
In [5]: def claim_standardScaler(df):
        ...
        Purpose: Normalize the dataframe so the features are on the same scale.
        Input: Dataframe
        Output: Returns the data frame with standarized values.
        ...
        sc = StandardScaler()
        df_std = sc.fit_transform(df)
        return df_std
```

```
In [6]: def claim_Imputer(df):
        ...
        Purpose: Normalize the dataframe so the features are on the same scale.
        Input: Dataframe
        Output: Returns the data frame with standarized values.
        ...
        imr = Imputer(missing_values=np.NaN, strategy='most_frequent',axis=0)
        imr = imr.fit_transform(df)
        return imr
```

```
In [7]: # Display all the fields in the input dataframe.
        df_claims.columns
```

```
Out[7]: Index([u'CLAIM_NUM', u'CLAIM_FRAUD_CD_DESC', u'CASUALTY_UNIT_IND',
               u'FATALITY_IND', u'SUBROGATION_IND', u'HIGH_VALUE_CLAIM_IND',
               u'ACCIDENT_FAULT_IND', u'LEGAL_ENTITY_NM', u'WEATHER_CD_DESC',
               u'CLAIM_JURISDICTION_ST_CD_DESC', u'LOSS_TYPE_CD_DESC',
               u'CLAIM_COMPLEXITY_CD_DESC', u'CLAIM_FILE_TYPE_DESC',
               u'CLAIM_TIER_CD_DESC', u'FAULT_RATING_CD_DESC', u'CATASTROPHE_IND',
               u'LOSS_SUB_TYPE_DESC', u'LIABILITY_CD_DESC', u'LIABILITY_PERCENTAGE',
               u'INCIDENT_ONLY_IND', u'GLASS_ONLY_IND', u'NON_CHARGEABLE_CLAIM_IND',
               u'TOTAL_LOSS_IND', u'PRODUCT_TYPE_DESC', u'ATTORNEY_ON_CLAIM_IND',
               u'VEHICLE_TOW_STORAGE_IND', u'COVERAGE_TYPE_CD_DESC',
               u'TOW_TYPE_CD_DESC', u'INCIDENT_TYPE_CD_DESC',
               u'COVERAGE_SUB_TYPE_CD_DESC', u'CLAIM_TYPE_DESC', u'COVERAGE_NM',
               u'PROPERTY_DESC', u'DAMAGE_DESC', u'CYCLE_TIME'],
              dtype='object')
```

```
In [30]: # Check value and count of each value in all the fields in the dataframe.
# for col in df_claims.columns:
#     if col in ['CLAIM_NUM']:
#         continue

#     print col
#     value_cnt = df_claims[col].value_counts()
#     print value_cnt
#     print '*****'
```

```
In [9]: # Move the target field to array.['CLAIM_NUM','CLAIM_ADJUSTER_GROUP_TYPE']
# df_target = df_claims[['CLAIM_NUM','CLAIM_ADJUSTER_GROUP_TYPE']]
df_claims = df_claims[(df_claims.CYCLE_TIME != 0)]

df_target = df_claims[['CLAIM_NUM','CYCLE_TIME']]

df_target.groupby('CYCLE_TIME').count()

df_target['Fast_track'] = np.where(df_target['CYCLE_TIME']<= 10, 1,0)
df_claims['Fast_track'] = np.where(df_claims['CYCLE_TIME']<= 10, 1,0)
```

C:\Users\gbu4moh\AppData\Local\Enthought\Canopy\User\lib\site-packages\ipykernel\\_main\_.py:9: SettingWithCopyWarning:

A value is trying to be set on a copy of a slice from a DataFrame.

Try using .loc[row\_indexer,col\_indexer] = value instead

See the caveats in the documentation: <http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing-view-versus-copy> (<http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing-view-versus-copy>)

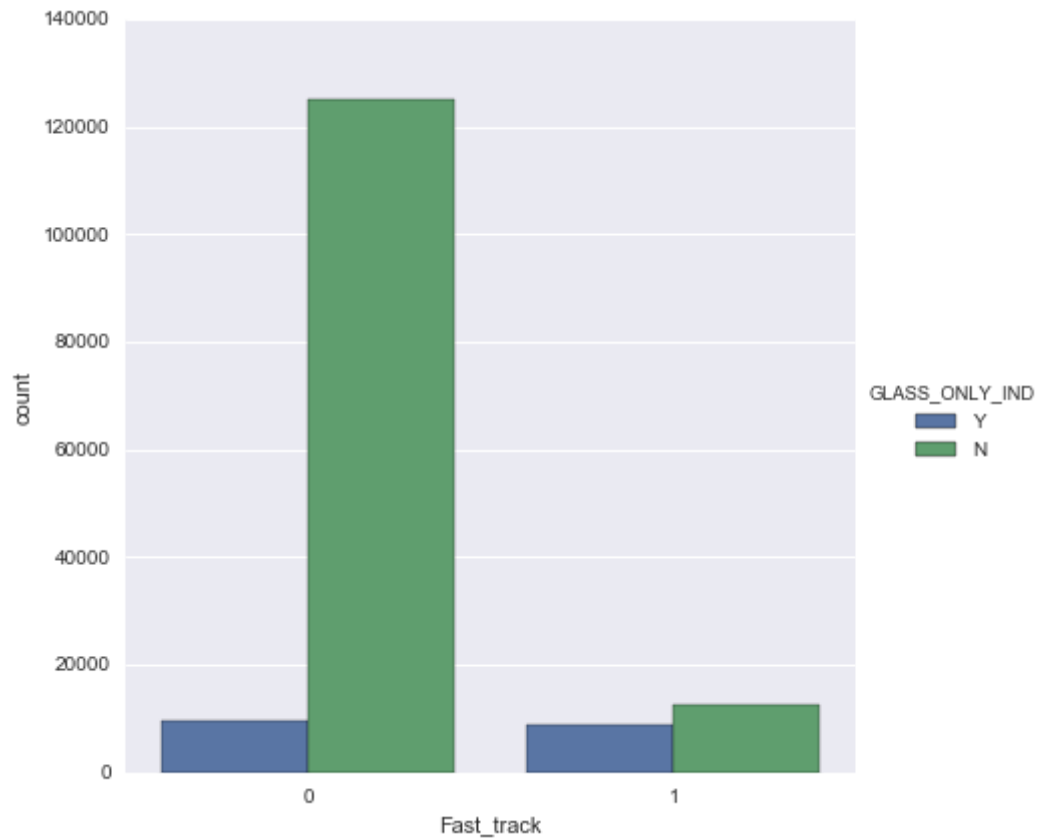
```
In [10]: df_target.describe()
```

```
Out[10]:
```

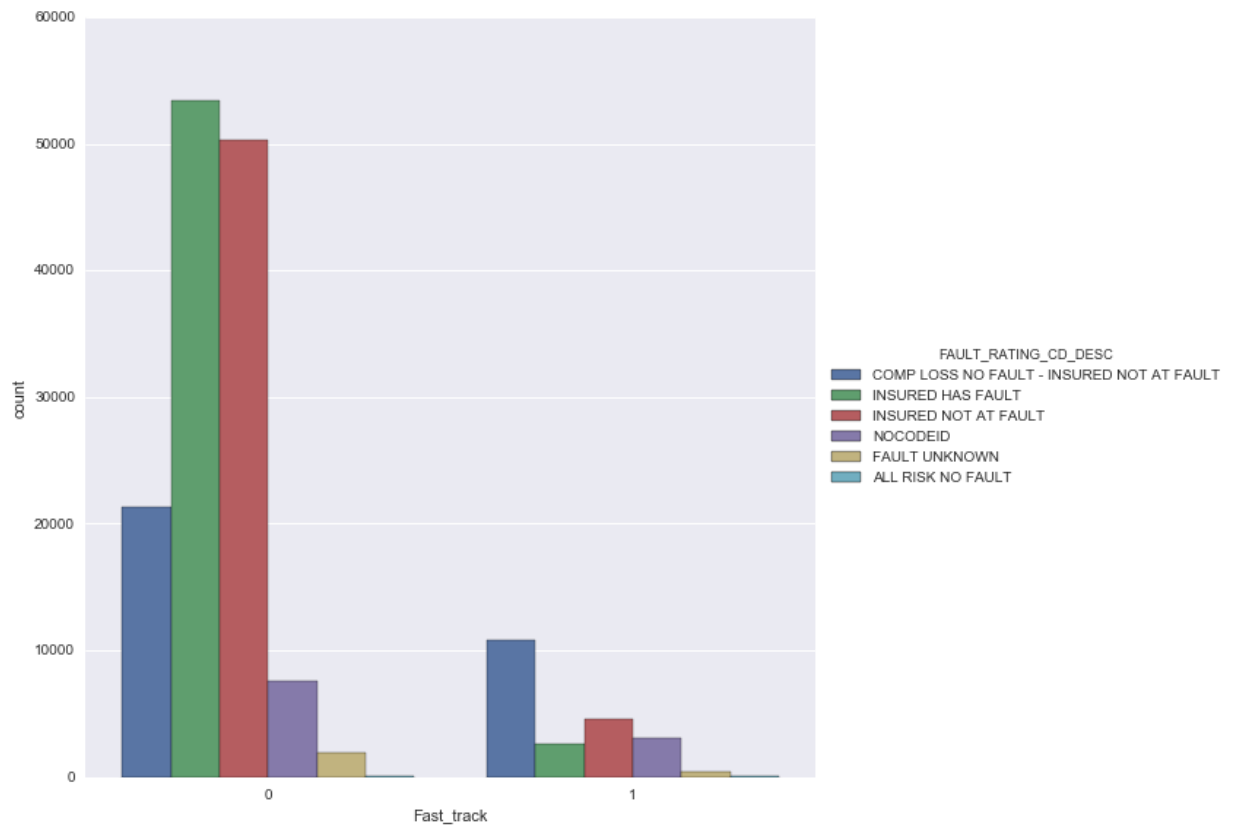
	CYCLE_TIME	Fast_track
count	130939.000000	156053.000000
mean	53.727132	0.138350
std	46.047580	0.345269
min	1.000000	0.000000
25%	17.000000	0.000000
50%	42.000000	0.000000
75%	76.000000	0.000000
max	230.000000	1.000000

```
In [11]: # Plot the glass only indicator and claim adjuster group.  
# = df_claims[df_claims['CLAIM_ADJUSTER_GROUP_TYPE'] != 'FNOL']  
  
sns.factorplot('Fast_track', data=df_claims, hue='GLASS_ONLY_IND' , kind='count', siz
```

Out[11]: <seaborn.axisgrid.FacetGrid at 0xd6acb38>



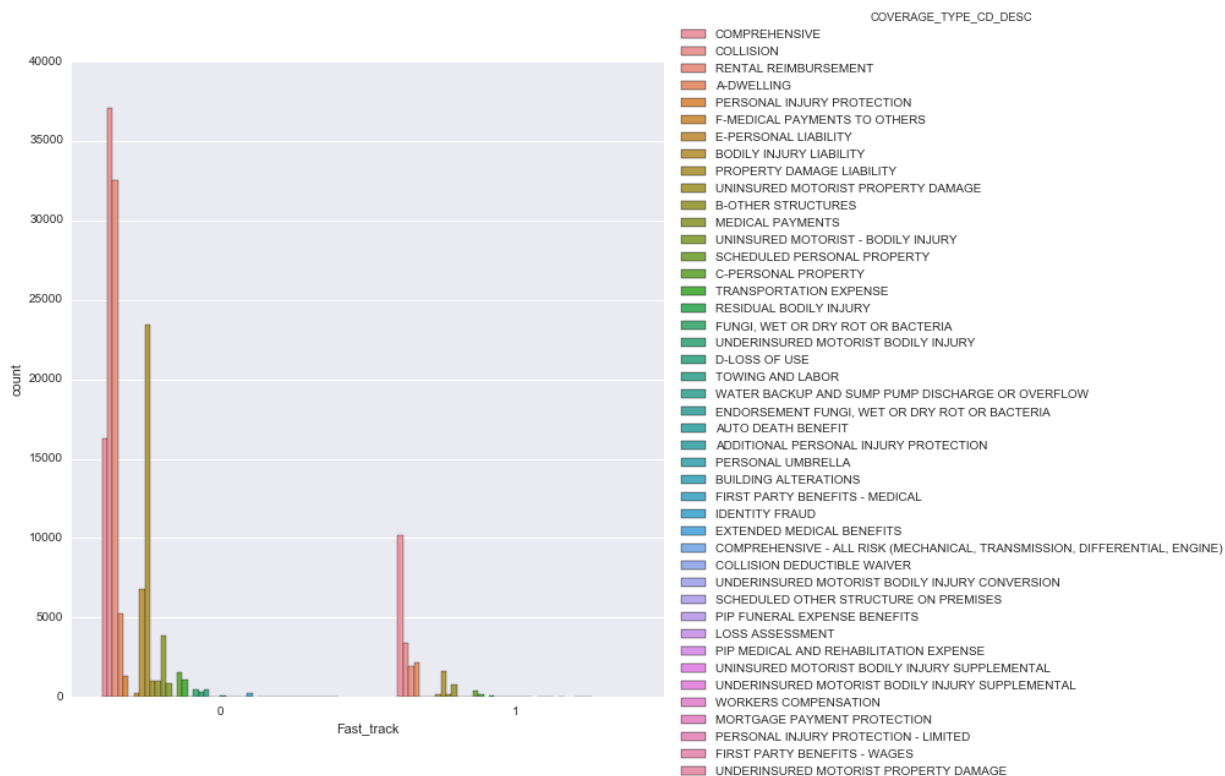
```
In [12]: # Plot the fault rating code and claim adjuster group.  
# df_report = df_claims[df_claims['CLAIM_ADJUSTER_GROUP_TYPE'] != 'FNOL']  
  
sns.factorplot('Fast_track',data=df_claims,hue='FAULT_RATING_CD_DESC',kind='count')  
  
Out[12]: <seaborn.axisgrid.FacetGrid at 0x115abcf8>
```



```
In [13]: # Plot the coverage code and claim adjuster group.
# df_report = df_claims[df_claims['CLAIM_ADJUSTER_GROUP_TYPE'] != 'FNOL']

sns.factorplot('Fast_track', data=df_claims, hue='COVERAGE_TYPE_CD_DESC', kind='cou
```

```
Out[13]: <seaborn.axisgrid.FacetGrid at 0x11168e48>
```



```
In [14]: # drop fields which are not useful -
df_claims.drop(['DAMAGE_DESC', 'PROPERTY_DESC', 'CLAIM_TYPE_DESC', 'INCIDENT_TYPE_',
'NON_CHARGEABLE_CLAIM_IND', 'INCIDENT_ONLY_IND', 'LOSS_SUB_TYPE_DESC', 'CLAIM_FILE_
```

```
In [15]: # Encode the nominal features - String to integers.
for col in df_claims.columns:
    claim_labelEncode(df_claims,col)
```

```
In [16]: # Impute the missing values in the dataframe.
claim_imp = claim_Imputer(df_claims)
```

```
In [17]: # Convert to dataframe.
df_claim_imp = pd.DataFrame(data=claim_imp, columns=df_claims.columns)
```

```
In [18]: # Standardize the features so the features are on same - scale.
# array_std = claim_standardScaler(df_claim_imp[['LIABILITY_PERCENTAGE']])

# df_std = pd.DataFrame(data=array_std,columns=['LIABILITY_PERCENTAGE'])

# df_claim_imp.drop(['LIABILITY_PERCENTAGE'], inplace=True,axis=1)

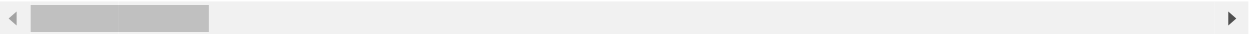
# df_claim_std = pd.concat([df_claim_imp,df_std],axis=1)

df_claim_std = df_claim_imp
df_claim_std.head()
```

Out[18]:

	codeCOVERAGE_NM	codeCOVERAGE_SUB_TYPE_CD_DESC	codeTOW_TYPE_CD_DE
0	11.0	10.0	3.0
1	11.0	10.0	3.0
2	11.0	10.0	3.0
3	11.0	10.0	3.0
4	11.0	10.0	3.0

5 rows × 23 columns



```
In [19]: # Create target dataframe. Set the fast track values to 1 and other groups to 0.
# df_target.loc[df_target['CLAIM_ADJUSTER_GROUP_TYPE'] != 'Fast_Track', 'CLAIM_AD
# df_target.loc[df_target['CLAIM_ADJUSTER_GROUP_TYPE'] == 'Fast_Track', 'CLAIM_AD
# df_target.groupby('CLAIM_ADJUSTER_GROUP_TYPE').count()
```

```

In [20]: # Use the Random forest classifier to identify the variable importance.
forest = RandomForestClassifier(n_estimators=1000, random_state=0,n_jobs=-1)

X_train, X_test, y_train, y_test = train_test_split(df_claim_std, list(df_target[

labels = df_claim_std.columns
forest.fit(X_train,y_train)
importances = forest.feature_importances_
indices = np.argsort(importances) [::-1]

for f in range(X_train.shape[1]):
    print(labels[f], importances[indices[f]])

feat_imp = pd.Series(importances, labels).sort_values(ascending=False)
feat_imp.plot(kind='bar', title='Feature Importances',figsize=(10,10))
plt.ylabel('Feature Importance Score')

('codeCOVERAGE_NM', 0.20923911046891788)
('codeCOVERAGE_SUB_TYPE_CD_DESC', 0.12500975677965381)
('codeTOW_TYPE_CD_DESC', 0.11884671739233212)
('codeCOVERAGE_TYPE_CD_DESC', 0.096151750235224973)
('codeVEHICLE_TOW_STORAGE_IND', 0.056325027359171838)
('codeATTORNEY_ON_CLAIM_IND', 0.047165505340229778)
('codePRODUCT_TYPE_DESC', 0.045992582021588381)
('codeGLASS_ONLY_IND', 0.045766589627370881)
('codeLIABILITY_CD_DESC', 0.045716173797628372)
('codeCATASTROPHE_IND', 0.03951539276309491)
('codeFAULT_RATING_CD_DESC', 0.039077431635860749)
('codeCLAIM_TIER_CD_DESC', 0.029695507439374654)
('codeCLAIM_COMPLEXITY_CD_DESC', 0.029674353595676734)
('codeLOSS_TYPE_CD_DESC', 0.022212382841006562)
('codeCLAIM_JURISDICTION_ST_CD_DESC', 0.01692482449312464)
('codeWEATHER_CD_DESC', 0.01112378868903534)
('codeLEGAL_ENTITY_NM', 0.010790759563695329)
('codeACCIDENT_FAULT_IND', 0.0057055044526657713)
('codeHIGH_VALUE_CLAIM_IND', 0.0050012914157365005)
('codeSUBROGATION_IND', 2.2502583056957403e-05)
('codeFATALITY_IND', 2.1974548443165455e-05)
('codeCASUALTY_UNIT_IND', 2.1072957110490578e-05)
('codeCLAIM_FRAUD_CD_DESC', 0.0)

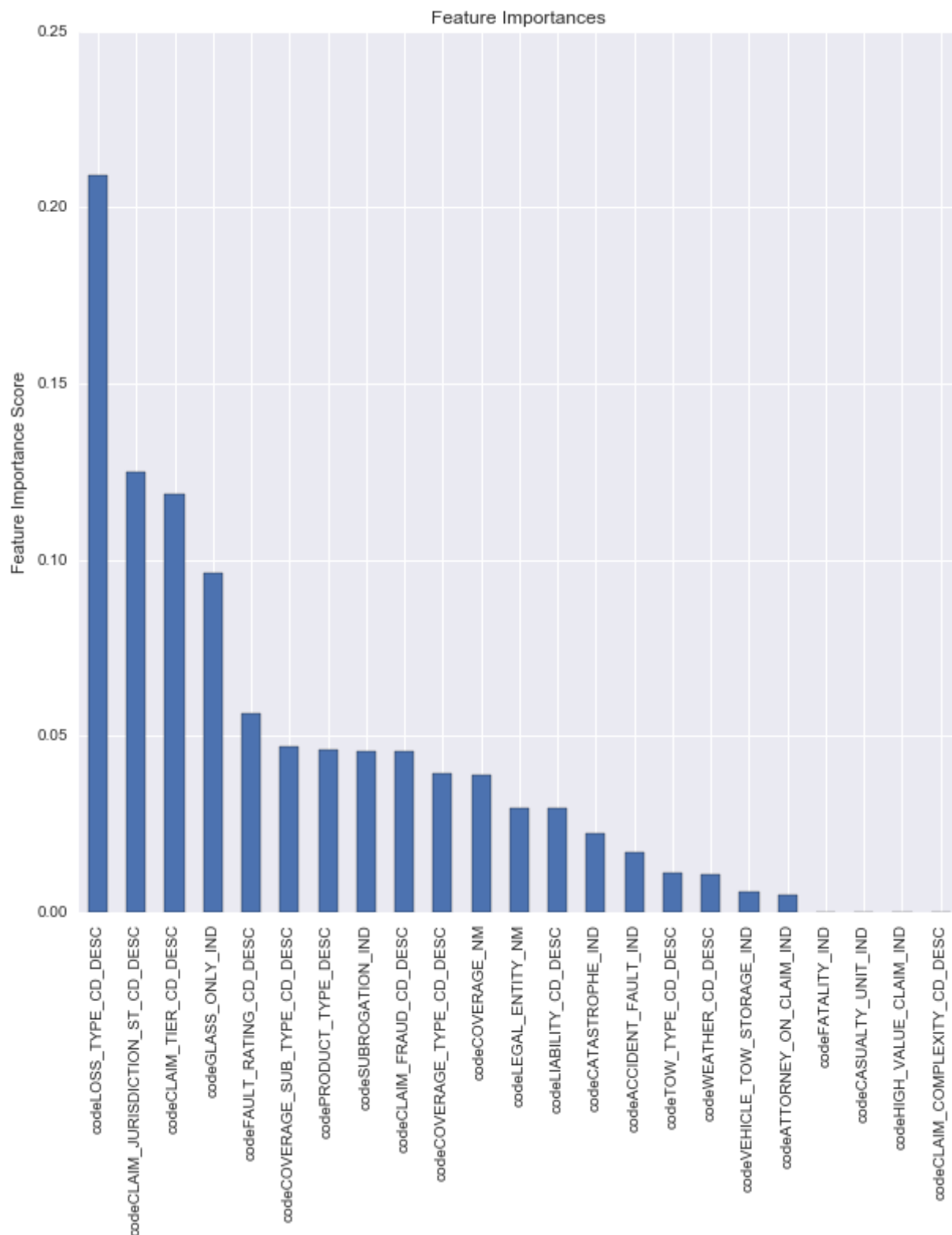
```

```

Out[20]: <matplotlib.text.Text at 0x44c5af98>

```





```
In [21]: # first 23 fields needs to be hot-encoded.
ohe = OneHotEncoder(categorical_features=np.arange(0,21))
df_claim_hot = ohe.fit_transform(df_claim_std)
```

```
In [22]: # Create data frame of hot-encoded array.
df_claim_mod = pd.DataFrame(df_claim_hot.toarray())
```

```
In [23]: # Split the dataset to train and test dataset. 70% of data is trained and 30% of
X_train, X_test, y_train, y_test = train_test_split(df_claim_mod, list(df_target[
```

```
In [67]: # Train the Logistic Regression model.
param_range = [0.01,0.02,0.03,0.04,0.05,0.06,0.07,0.08,0.09]
```

```
for param in param_range:
    print('C : ',param)
    lr = LogisticRegression(penalty='l2',C=param)
    lr.fit(X_train,y_train)
    print('training accuracy:', lr.score(X_train,y_train))
    print('test accuracy: ', lr.score(X_test, y_test))
```

```
('C : ', 0.01)
('training accuracy:', 0.86705969589058651)
('test accuracy: ', 0.86775888585099115)
('C : ', 0.02)
('training accuracy:', 0.8671146223349232)
('test accuracy: ', 0.86790840738209163)
('C : ', 0.03)
('training accuracy:', 0.86716039437187031)
('test accuracy: ', 0.86820745044429259)
('C : ', 0.04)
('training accuracy:', 0.86717870318664925)
('test accuracy: ', 0.86797248803827753)
('C : ', 0.05)
('training accuracy:', 0.86716954877925978)
('test accuracy: ', 0.86792976760082019)
('C : ', 0.06)
('training accuracy:', 0.8671878575940386)
('test accuracy: ', 0.86792976760082019)
('C : ', 0.07)
('training accuracy:', 0.86720616640881754)
('test accuracy: ', 0.86792976760082019)
('C : ', 0.08)
('training accuracy:', 0.867215320816207)
('test accuracy: ', 0.86797248803827753)
('C : ', 0.09)
('training accuracy:', 0.86646465941027306)
('test accuracy: ', 0.86739576213260428)
```

```
In [ ]: # Train the SVM model.
param_range = [0.0001,0.001,0.01,0.1,1.0,10.0,100.0,1000.0]
```

```
for param in param_range:
    svm = SVC(kernel='sigmoid',C=param,random_state= 0)
    svm.fit(X_train,y_train)

    print('training accuracy:', svm.score(X_train,y_train))
    print('test accuracy: ', svm.score(X_test, y_test))
```

```

In [ ]: # Plot the training Curve
pipe_lr = Pipeline([('clf', LogisticRegression(penalty='l2', random_state=0, C=0.01))

train_sizes, train_scores, test_scores = learning_curve(estimator=pipe_lr,
                                                         X=X_train,
                                                         y=y_train,
                                                         train_sizes=np.linspace(0.1, 1.0, 10),
                                                         cv=10,
                                                         n_jobs=1
                                                         )

# for x,y,z in zip(train_sizes, train_scores, test_scores):
#     print('train size ', x)
#     print('train scores ', y)
#     print('test scores ', z)

train_mean = np.mean(train_scores, axis=1)
train_std = np.std(train_scores, axis=1)
test_mean = np.mean(test_scores, axis=1)
test_std = np.std(test_scores, axis=1)

print('train_mean ', train_mean)
print('train_std ', train_std)
print('test_mean ', test_mean)
print('test_std ', test_std)

plt.plot(train_sizes, train_mean, color='blue', marker='o', markersize=5, label='train')
plt.fill_between(train_sizes, train_mean + train_std, train_mean - train_std, alpha=0.1)
plt.plot(train_sizes, test_mean, color='green', marker='x', markersize=5, label='valid')
plt.fill_between(train_sizes, test_mean + test_std, test_mean - test_std, alpha=0.1)

plt.grid()
plt.xlabel('Number of training samples')
plt.ylabel('Accuracy')
plt.legend(loc='lower right')
# plt.ylim([0.8, 1.0])
plt.show()

```

```

In [25]: # confusion matrix

#y_predict = lr.predict(X_test)

#confusion_matrix(y_test, y_predict)

#print ("Accuracy : %.4g" % metrics.accuracy_score(y_test, y_predict))
#print ("AUC Score (Train): %f" % metrics.roc_auc_score(y_test, y_predict))

#cv_score = cross_validation.cross_val_score(alg, dtrain[predictors], dtrain['Dis'])
# print "CV Score : Mean - %.7g | Std - %.7g | Min - %.7g | Max - %.7g" % (np.mean(cv_score),

```

```
In [26]: # Gradient boosting classifier.
param_test = {'n_estimators':range(20,81,10)}
gsearch1 = GridSearchCV(estimator = GradientBoostingClassifier(learning_rate=0.1,
param_grid = param_test, scoring='roc_auc',n_jobs=4,iid=False, cv=5)
gsearch1.fit(X_train,y_train)
```

```
Out[26]: GridSearchCV(cv=5, error_score='raise',
    estimator=GradientBoostingClassifier(criterion='friedman_mse', init=None,
    learning_rate=0.1, loss='deviance', max_depth=8,
    max_features='sqrt', max_leaf_nodes=None,
    min_impurity_split=1e-07, min_samples_leaf=50,
    min_samples_split=500, min_weight_fraction_leaf=0.0,
    n_estimators=100, presort='auto', random_state=10,
    subsample=0.8, verbose=0, warm_start=False),
    fit_params={}, iid=False, n_jobs=4,
    param_grid={'n_estimators': [20, 30, 40, 50, 60, 70, 80]},
    pre_dispatch='2*n_jobs', refit=True, return_train_score=True,
    scoring='roc_auc', verbose=0)
```

```
In [27]: gsearch1.grid_scores_, gsearch1.best_params_, gsearch1.best_score_

C:\Users\gbu4moh\AppData\Local\Enthought\Canopy\User\lib\site-packages\sklearn
\model_selection\_search.py:662: DeprecationWarning: The grid_scores_ attribute
was deprecated in version 0.18 in favor of the more elaborate cv_results_ attr
ibute. The grid_scores_ attribute will not be available from 0.20
DeprecationWarning)
```

```
Out[27]: ([mean: 0.83821, std: 0.00146, params: {'n_estimators': 20},
    mean: 0.84184, std: 0.00150, params: {'n_estimators': 30},
    mean: 0.84368, std: 0.00168, params: {'n_estimators': 40},
    mean: 0.84508, std: 0.00192, params: {'n_estimators': 50},
    mean: 0.84590, std: 0.00195, params: {'n_estimators': 60},
    mean: 0.84648, std: 0.00194, params: {'n_estimators': 70},
    mean: 0.84702, std: 0.00198, params: {'n_estimators': 80}],
    {'n_estimators': 80},
    0.84701706237241292)
```

```
In [28]: y_predict = gsearch1.predict(X_test)

confusion_matrix(y_test, y_predict)
```

```
Out[28]: array([[39421, 1012],
    [ 5109, 1274]])
```

```
In [29]: print('training accuracy:', gsearch1.score(X_train,y_train))
print('test accuracy: ', gsearch1.score(X_test, y_test))

('training accuracy:', 0.85074250258692097)
('test accuracy: ', 0.84597644256213966)
```

```
In [ ]:
```

