	Insurance companies make money by betting on risk - the risk that insured won't die before his/her time and make the insurer payout, or the risk insured house won't burn down or insured SUV wo be totalled in a crash. The first task of any insurer is to price risk and charge a premium for assuming it. When a customer files a claim, the company must process it, check it for accuracy, and sub payment. This adjusting process is necessary to filter out fraudulent claims and minimize the risk of loss to the company. Insurance fraud costs the auto insurers and consumers billions of dollars each year, which forces insurance premium prices to go up annually. Fraud may be committed at different points in the transaction by applicants, policy holders, third-party claimants, or profession who provide services to claimants. Whether the matter involves a consumer misrepresenting a claim or a retailer or service provider misreporting products or labour, insurance fraud is considered crime. Investigating fraudulent claims is costly and time consuming for insurers. It is physically impossible for insurance companies to do a thorough check of the thousands of claims that enter the systems daily. Common frauds include 'padding', (which means to add damages, injuries and fictitious passengers to insurance claims) or inflating claims; misrepresenting facts on an insurance application; submitting claims for injuries or damage that never occurred; and staging accidents. Some lines of insurance are more vulnerable to fraud than others. Healthcare, workers compensation, and auto insurance are generally considered to be the sectors most affected. Any effort aiming to debug the fraudulent transactions in the above-mentioned businesses and probably in many other ones is named as a fraud detection process. In terms of insurance fraud detection, machine learning applies aspects of AI to give systems the ability to improve from experience with no extra programming by analyzing large, labelled data sets. Machine learning can improve fraud d
	Although it borrows underlying principals found in statistical models, the main focus of machine learning is producing predictions. These predictions are based on the analysis of known outcomes, known as "ground truth." Problem Statement: I will be creating a predictive model that predicts if an insurance claim is fraudulent or not. The answere between YES/NO, is a Binary Classification task. A comparison study has been performed understand which ML algorithm suits best to the dataset. import pandas as pd import numpy as np import seaborn as sns from sklearn.preprocessing import LabelEncoder, StandardScaler from sklearn.model_selection import train_test_split import numpy as np import numpy as np import numpy as np import rcParams
n [3]: n [4]:	<pre>%matplotlib inline import warnings warnings.filterwarnings("ignore") from sklearn import model_selection from sklearn.linear_model import LogisticRegressionCV from sklearn.linear_model import LogisticRegressionCV from sklearn.tree import DecisionTreeClassifier from sklearn.neighbors import KNeighborsClassifier from sklearn.svm import SVC from sklearn.ensemble import AdaBoostClassifier insr_data = pd.read_csv("https://raw.githubusercontent.com/dsrscientist/Data-Science-ML-Capstone-Projects/master/Automobile_insurance_fraud.csv") insr_data.shape</pre>
n [4]: n []: n [5]: ut[5]: n [6]:	<pre>(1000, 40) pd.set_option('display.max_columns', None) insr_data.head(10) insr_data['_c39'].value_counts() Series([], Name: _c39, dtype: int64) insr_data['_c39'].unique() array([nan])</pre>
n [7]: n [8]: ut[8]:	<pre># It is clear from above that '_c39','policy_number', 'incident_location' are not going to affect our study; so dropping it: auto_data = insr_data.drop(['_c39','policy_number', 'incident_location'], axis=1) Index(['months_as_customer', 'age', 'policy_bind_date', 'policy_state',</pre>
	<pre>'police_report_available', 'total_claim_amount', 'injury_claim', 'property_claim', 'vehicle_claim', 'auto_make', 'auto_model', 'auto_year', 'fraud_reported'], dtype='object') auto_data['fraud_reported'].value_counts()</pre>
	data_1.head(15) data_1['policy_csl'].value_counts() 250/500
[12]: [13]: [14]: [14]: [15]: [15]:	<pre>data_2 = data_1.drop(['injury_claim', 'property_claim', 'vehicle_claim'], axis=1) data_3 = data_2.drop(['insured_zip'], axis=1) data_3['policy_state'].unique() array(['OH', 'IN', 'IL'], dtype=object) data_3['policy_deductable'].unique() array([1000, 2000, 500], dtype=int64) data_3['umbrella_limit'].unique()</pre>
[16]: [17]: [17]: [18]:	<pre>#As umbrella limit cannot be negative so removing '-' : data_3['umbrella_limit'].replace(to_replace = -1000000, value = 1000000, inplace=True) data_3['umbrella_limit'].unique() data_3['collision_type'].unique() array(['Side Collision', '?', 'Rear Collision', 'Front Collision'],</pre>
[20]: [21]:	<pre>data_3['property_damage'].unique() array(['YES', '?', 'NO'], dtype=object) data_3['property_damage'].replace(to_replace='?', value = 'NA', inplace=True) data_3['property_damage'].value_counts() data_3['police_report_available'].unique() array(['YES', '?', 'NO'], dtype=object)</pre>
[22]: [23]: [24]:	<pre>data_3['police_report_available'].replace(to_replace='?', value = 'NA', inplace=True) data_3['police_report_available'].value_counts() Data Visualization import matplotlib.pyplot as plt fig = plt.figure(figsize=(7,5)) ax = (data_3['insured_sex'].value_counts()*100.0 /len(data_3))\ .plot.pie(autopct='%.1f%%', labels = ['Male', 'Female'], fontsize=12) ax.set_title('% Sex_Ratio') plt.show()</pre>
	Male 53.7% 46.3%
[25]: [26]:	<pre># trying to figure out that is there any affect of the month on incidents: # extracting month from incident_date data_3['month'] = pd.DatetimeIndex(data_3['incident_date']).month sizes = [12, 11, 10, 9, 8, 7, 6, 5, 4, 3, 2, 1] fig1, ax1 = plt.subplots(figsize=(12, 12)) fig1.subplots_adjust(0.3, 0, 1, 1) theme = plt.get_cmap('copper') ax1.set_prop_cycle("color", [theme(1. * i / len(sizes)))</pre>
	<pre>ax1.set_prop_cycle("color", [theme(1. * 1 / len(sizes))</pre>
	Jan Feb Mar Apr May June July Aug Sep Oct Nov Dec
	36.5% 1.9% Dec 2.3% Nov
	3.9% 3.5% Sep Aug July July Mar Apr May
[27]:	From this chart we can observe that the highest number of incidents are in the month of January followed by Feb. fig = plt.figure(figsize=(10,6)) ax = (data_3['insured_relationship'].value_counts()*100.0 /len(data_3))\ .plot.pie(autopct='%.1f%', labels = ['husband', 'wife', 'own-child', 'unmarried', 'other-relative', 'not-in-family'],
	OWB-child 17.4% 18.3% 14.1% 15.5% not-in-family other-relative
[28]:	<pre>fig = plt.figure(figsize=(10,6)) ax = (data_3['incident_severity'].value_counts()*100.0 /len(data_3))\ .plot.pie(autopct='%.1f%', labels = ['Major Damage', 'Total Loss', 'Minor Damage', 'Trivial Damage'],</pre>
[29]:	Total cost 9.0% Trivial Damage fig = plt.figure(figsize=(10,6)) ax = data_3.groupby('incident_state').fraud_reported.count().plot.bar(ylim=0)
	ax.set_ylabel('Fraud reported') plt.show() 250 - 200 - Property of the state of
[30]:	fig = plt.figure(figsize=(10,6)) ax = data_3.groupby('policy_state').fraud_reported.count().plot.bar(ylim=0) ax.set_ylabel('Fraud reported') plt.show()
	350 - 300 - 250 - Pe 150 -
[37]:	fig = plt.figure(figsize=(10,6)) ax = data_3.groupby('incident_type').fraud_reported.count().plot.bar(ylim=0) ax.set_xticklabels(ax.get_xticklabels(), rotation=20, ha="right") ax.set_ylabel('Fraud reported') plt.show()
	400 - 350 - 300 - Pe 250 - 150 -
[31]:	fig = plt.figure(figsize=(10,6)) ax = sns.countplot(x='incident_state', data=data_3)
	250 - 200 - 150 - 100 - 50 -
[32]:	fig = plt.figure(figsize=(10,6)) ax = sns.countplot(x = 'insured_education_level', data=data_3) ax.set_xlabel('policy_annual_premium') plt.show()
	120 - 100 -
[33]: [34]:	# extracting the age of Vehicle from auto_year, as all the incidents are of the year 2015: data_3['Vehicle_Age'] = 2015 - data_3['auto_year'] > Policy bind date and incident date can give the difference of time between policy taken and occurence of incident. Right now I am leaving it due to shartage of time. # Dealing with time it can tell us that at what time of the day major of the incidents occur: bins = [-1, 3, 6, 9, 12, 17, 20, 24] # Factorize according to the time period of the day. names = ["midnight", "early_morning", "late_morning", "afternoon", "evening", "night", "late_night"] data_3['incident_part_of_day'] = pd.cut(data_3.incident_hour_of_the_day, bins, labels=names).astype(object)
[34]:	data_3[['incident_hour_of_the_day', 'incident_part_of_day']].head(10) incident_hour_of_the_day incident_part_of_day 0
[35]:	8 21 late_night 9 14 evening
[36]: [37]: [38]: :[38]:	<pre>from sklearn.preprocessing import OneHotEncoder from sklearn.compose import ColumnTransformer data_4 = data_3.drop(columns = ['month', 'auto_year', 'incident_date', 'policy_bind_date']) cat_features = data_4.select_dtypes(include='0').keys() cat_features Index(['policy_state', 'policy_csl', 'insured_education_level',</pre>
[39]:	<pre>'police_report_available', 'auto_make', 'auto_model', 'incident_part_of_day'], dtype='object') cat_features = ['policy_state', 'policy_csl', 'insured_education_level',</pre>
[40]: [40]: [41]: [41]: [42]:	<pre># apply le on categorical feature columns data_4[cat_features] = data_4[cat_features].apply(lambda col: le.fit_transform(col)) data_4 data_f = pd.get_dummies(data_4, columns = cat_features) data_f.head() data_f.shape (1000, 165) # Finding correlation between different variables:</pre>
[43]: [44]: [45]:	<pre>corr_matrix = data_4.corr() X = data_f.drop(columns = ['collision_type']) X.head(2) # seperating the target variable and predictor data: X = data_f.drop(columns = ['fraud_reported']) y = data_f['fraud_reported'] # splitting dataset into training and testing part:</pre>
[46]:	<pre>X_train, X_test, y_train, y_test = train_test_split(X, y, train_size=0.8, random_state=1234) print('length of X_train and X_test: ', len(X_train), len(X_test)) print('length of y_train and y_test: ', len(y_train), len(y_test)) length of X_train and X_test: 800 200 length of y_train and y_test: 800 200 # checking the quality of data: from sklearn.discriminant_analysis import LinearDiscriminantAnalysis from sklearn.model_selection import KFold from sklearn.model_selection import cross_val_score # evaluate an LDA model on the dataset using k-fold cross validation model = LinearDiscriminantAnalysis()</pre>
	<pre>model = LinearDiscriminantAnalysis() kfold = KFold(n_splits=5, random_state=7) result = cross_val_score(model, X, y, cv=kfold, scoring='accuracy') print(result.mean()) 0.841 It means our data is now ready for further analysis. MODEL BUILDING from sklearn.ensemble import RandomForestClassifier from sklearn.metrics import accuracy_score, recall_score, classification_report, cohen_kappa_score from sklearn import metrics # Baseline Random forest based Model rfc = RandomForestClassifier(criterion = 'gini', n_estimators=1000, verbose=1, n_jobs = -1,</pre>
	[Parallel(n_jobs=-1)]: Done 442 tasks
[48]:	0 0.74 0.94 0.83 139 1 0.67 0.26 0.38 61 accuracy 0.73 200 macro avg 0.71 0.60 0.60 200 weighted avg 0.72 0.73 0.69 200 from sklearn.metrics import confusion_matrix import itertools #Evaluation of RFModel - Confusion Matrix Plot def plot_confusion_matrix(cm, classes, title ='Confusion matrix', normalize=False, cmap = plt.cm.Blues): """ This function prints and plots the confusion matrix.
	<pre>Normalization can be applied by setting `normalize=True`. """ print('Confusion matrix') print(cm) fig = plt.figure(figsize=(10,6)) plt.style.use('fivethirtyeight') plt.imshow(cm, interpolation='nearest', cmap=cmap) plt.title(title) plt.colorbar() tick_marks = np.arange(len(classes)) plt.xticks(tick_marks, classes, rotation=45) plt.yticks(tick_marks, classes) fmt = '.2f' if normalize else 'd'</pre>
	<pre>thresh = cm.max() / 2. for i, j in itertools.product(range(cm.shape[0]), range(cm.shape[1])): plt.text(j, i, format(cm[i, j], fmt),</pre>
	# Plot non-normalized confusion matrix plt.figure() plot_confusion_matrix(cnf_matrix, classes=['Fraud reported_Y', 'Fraud_reported_N'],
	Fraud_reported_N 45 16 40 20
[49]: [50]:	Predicted label from sklearn.preprocessing import StandardScaler scaler = StandardScaler(with_mean=False) X_train_scaled = scaler.fit_transform(X_train) X_test_scaled = scaler.transform(X_test) from sklearn import model_selection from sklearn.linear_model import LogisticRegression from sklearn.linear_model import LogisticRegressionCV from sklearn.tree import DecisionTreeClassifier from sklearn.neighbors import KNeighborsClassifier
	<pre>from sklearn.svm import SVC from sklearn.ensemble import AdaBoostClassifier logreg2= LogisticRegressionCV(solver='lbfgs', cv=10) knn = KNeighborsClassifier(5) svcl = SVC() adb = AdaBoostClassifier() dtclf = DecisionTreeClassifier(max_depth=5) rfclf = RandomForestClassifier() # prepare configuration for cross validation test harness seed = 7 # prepare models models = [] models.append(('LR', LogisticRegressionCV(solver='lbfgs', max_iter=5000, cv=10))) models.append(('KNN', KNeighborsClassifier()))</pre>
	models.append(('LR', LogisticRegressionCV(solver='lbfgs', max_iter=5000, cv=10)))
[51]:	
	Algorithm Comparison 0.85 0.80 0.75
	D.70 LR KNN DT SVM RF ADA Above is the list of each algorithm, the mean accuracy and the standard deviation accuracy and a box & whisker plot showing the spread of the accuracy scores across each cross validation fold to each algorithm. From above we can conclude that LogisticRegression and DecisionTreeClassifier are having the higest accuracy score. Conclusion: We can use LogisticRegression Model to detect the fraud on insurance claims.