In [1]:	import pandas as pd import numpy as np %matplotlib inline import matplotlib.pyplot as plt import matplotlib.lines as mlines from mpl_toolkits.mplot3d import Axes3D import seaborn as sns from sklearn.model_selection import train_test_split, learning_curve from sklearn.metrics import average_precision_score from xgboost.sklearn import XGBClassifier from xgboost import plot_importance, to_graphviz
	# Unzipping data on Colab #local_zip = './archive.zip' #zip_ref = zipfile.Zipfile(local_zip, 'r') #zip_ref.extractall('./') #zip_ref.close() import warnings warnings.filterwarnings("ignore", category=DeprecationWarning) Import data and correct spelling of original column headers for consistency
<pre>In [3]: In [4]: Out[4]:</pre>	df = pd.read_csv('PS_20174392719_1491204439457_log.csv') df = df.rename(cclumms={'oldbalanceOrg':'oldbalanceOrig', 'newbalanceOrig', 'newbalanceOrig', 'oldbalanceDest':'oldbalanceDest', 'newbalanceDest';'newBalanceDest')) step type amount nameOrig oldBalanceOrig newBalanceOrig (0 1 PAYMENT 9839.64 C1231006815 170136.0 160296.36 1 1 PAYMENT 1864.28 C1666544295 21249.0 19384.72 2 1 TRANSFER 181.00 C1305486145 181.0 0.00 3 1 CASH_OUT 181.00 C340083671 181.0 0.00 4 1 PAYMENT 11668.14 C2948537720 41554.0 29885.36 nameDest oldBalanceDest newBalanceDest isFraud isFlaggedFraud 0 M1979787155 0.0 0.0 0.0 0 0 1 M204428225 0.0 0.0 0.0 0 0 2 C553264965 0.0 0.0 0.0 1 0 3 C38997010 21182.0 0.0 1 0 4 M1239701703 0.0 0.0 0 0 0 df.isnull().values.any()
In [5]:	<pre>Exploratory Data Analysis print('\n The types of fraudulent transactions are {}'.format(\ list(df.loc[df.isFraud == 1].type.drop_duplicates().values))) # only 'CASH_OUT'</pre>
In [6]:	The types of fraudulent transactions are ['TRANSFER', 'CASH_OUT'] The number of fraudulent TRANSFERS = 4097 The number of fraudulent CASH_OUTS = 4116 print('\nThe type of transactions in which isFlaggedFraud is set: \ {}'.format(list(df.loc[df.isFlaggedFraud == 1].type.drop_duplicates()))) # only 'TRANSFER' dfTransfer = df.loc[df.type == 'TRANSFER'] dfFlagged = df.loc[df.isFlaggedFraud == 1] dfNotFlagged = df.loc[df.isFlaggedFraud == 0]
In [7]:	<pre>print('\nMin amount transacted when isFlaggedFraud is set= {}'\</pre>
	(dfTransfer.oldBalanceDest == 0) & (dfTransfer.newBalanceDest == 0)]))) # 4158 The number of TRANSFERs where isFlaggedFraud = 0, yet oldBalanceDest = 0 and newBalanceDest = 0: 4158 isFlaggedFraud being set cannot be thresholded on oldBalanceOrig since the corresponding range of values overlaps with that for TRANSFERs where isFlaggedFraud is not set (see below). Note that we do not need to consider newBalanceOrig since it is updated only after the transaction, whereas isFlaggedFraud would be set before the transaction takes place. print('\nMin, Max of oldBalanceOrig for isFlaggedFraud = 1 TRANSFERs: {}'.\ format([round(dfFlagged.oldBalanceOrig.min()), round(dfFlagged.oldBalanceOrig.max())])) print('\nMin, Max of oldBalanceOrig for isFlaggedFraud = 0 TRANSFERs where \ oldBalanceOrig = \ newBalanceOrig: {}'.format(\ [dfTransfer.loc[(dfTransfer.isFlaggedFraud == 0) & (dfTransfer.oldBalanceOrig \\ = dfTransfer.newBalanceOrig)].oldBalanceOrig.min(), \
	round(dfTransfer.loc[(dfTransfer.isFlaggedFraud == 0) & (dfTransfer.oldBalanceOrig \
	<pre>print('\nHave destinations for transactions flagged as fraud initiated\ other transactions? \ {}'.format((dfFlagged.nameDest.isin(dfNotFlagged.nameOrig)).any())) # False # Since only 2 destination accounts of 16 that have 'isFlaggedFraud' set have been # destination accounts more than once, # clearly 'isFlaggedFraud' being set is independent of whether a # destination account has been used before or not print('\nHow many destination accounts of transactions flagged as fraud have been \ destination accounts more than once?: {}'\ .format(sum(dfFlagged.nameDest.isin(dfNotFlagged.nameDest)))) # 2</pre> Have originators of transactions flagged as fraud transacted more than once? False
	Have destinations for transactions flagged as fraud initiated other transactions? False How many destination accounts of transactions flagged as fraud have been destination accounts more than once?: 2 It can be easily seen that transactions with isFlaggedFraud set occur at all values of step, similar to the complementary set of transactions. Thus isFlaggedFraud does not correlate with step either and is therefore seemingly unrelated to any explanatory variable or feature in the data **Conclusion:* Although isFraud* is always set when isFlaggedFraud* is set, since isFlaggedFraud* is set just 16 times in a seemingly meaningless way, we can treat this feature as insignificant and discard it in the dataset without loosing information. Are expected merchant accounts accordingly labelled? **print('\nare there any merchants among originator accounts for CASH_IN \ transactions? {}'.format(\ (df.loc[df.type == 'CASH_IN'].nameOrig.str.contains('M')).any())) # False
In [11]: In [12]:	Are there any merchants among originator accounts for CASH_IN transactions? False Similarly, it was stated that CASH_OUT involves paying a merchant. However, for CASH_OUT transactions there are no merchants among the destination accounts. print('\nAre there any merchants among destination accounts for CASH_OUT \ transactions? {}'.format(\ (df.loc[df.type == 'CASH_OUT'].nameDest.str.contains('M')).any())) # False Are there any merchants among destination accounts for CASH_OUT transactions? False In fact, there are no merchants among any originator accounts. Merchants are only present in destination accounts for all PAYMENTS. print('\nAre there merchants among any originator accounts? {}'.format(\ df.nameOrig.str.contains('M').any())) # False
	print('\nAre there any transactions having merchants among destination accounts\ other than the PAYMENT type? {}'.format(\ (df.loc[df.nameDest.str.contains('M')].type != 'PAYMENT').any())) # False Are there merchants among any originator accounts? False Are there any transactions having merchants among destination accounts other than the PAYMENT type? False Conclusion: Among the account labels nameOrig and nameDest, for all transactions, the merchant prefix of 'M' occurs in an unexpected way. Are there account labels common to fraudulent TRANSFERs and CASH_OUTS? print('\nWithin fraudulent transactions, are there destinations for TRANSFERS \ that are also originators for CASH_OUTS? {}'.format(\ (dfFraudTransfer.nameDest.isin(dfFraudCashout.nameOrig)).any())) # False dfNotFraud = df.loc[df.isFraud == 0]
In [14]:	Within fraudulent transactions, are there destinations for TRANSFERS that are also originators for CASH_OUTs? False Could destination accounts for fraudulent TRANSFERs originate CASHOUTs that are not detected and are labeled as genuine? It turns out there are 3 such accounts. print('\nFraudulent TRANSFERs whose destination accounts are originators of \ genuine CASH_OUTs: \n\n\{\}'.format(dfFraudTransfer.loc[dfFraudTransfer.nameDest.\) isin(dfNotFraud.loc[dfNotFraud.type == 'CASH_OUT'].nameOrig.drop_duplicates())])) Fraudulent TRANSFERs whose destination accounts are originators of genuine CASH_OUTs: step type amount nameOrig oldBalanceOrig \ 1030443 65 TRANSFER 1282971.57 C1175896731 1282971.57 6039814 486 TRANSFER 1282971.57 C1175896731 1282971.57 6039814 486 TRANSFER 214793.32 C2140495649 214793.32 6362556 738 TRANSFER 814689.88 C2029041842 814689.88
	newBalanceOrig nameDest oldBalanceDest isFraud \ 1030443
	dfNotFraud.loc[(dfNotFraud.type == 'CASH_OUT') & (dfNotFraud.nameOrig == \
	<pre>#X = X.loc[np.random.choice(X.index, 100000, replace = False)] Y = X['isFraud'] del X['isFraud'] # Eliminate columns shown to be irrelevant for analysis in the EDA X = X.drop(['nameOrig', 'nameDest', 'isFlaggedFraud'], axis = 1) # Binary-encoding of labelled data in 'type' X.loc[X.type == 'TRANSFER', 'type'] = 0 X.loc[X.type == 'CASH_OUT', 'type'] = 1 X.type = X.type.astype(int) # convert dtype('0') to dtype(int)</pre> Imputation of Latent Missing Values
	The data has several transactions with zero balances in the destination account both before and after a non-zero amount is transacted. The fraction of such transactions, where zero likely denotes a missing value, is much larger in fraudulent (50%) compared to genuine transactions (0.06%). Xfraud = X.loc[Y == 1] XnonFraud = X.loc[Y == 0] print('\nThe fraction of fraudulent transactions with \'oldBalanceDest\' = \\ \'newBalanceDest\' = 0 although the transacted \'amount\' is non-zero is: {}'.\\ format(len(Xfraud.loc[(Xfraud.oldBalanceDest == 0) & \((Xfraud.newBalanceDest == 0) & (Xfraud.amount)]) / (1.0 * len(Xfraud)))) print('\nThe fraction of genuine transactions with \'oldBalanceDest\' = \\ newBalanceDest\' = 0 although the transacted \'amount\' is non-zero is: {}'.\\ format(len(XnonFraud.loc[(XnonFraud.oldBalanceDest == 0) & \((XnonFraud.newBalanceDest == 0) & (XnonFraud.amount)]) / (1.0 * len(XnonFraud))))
In [18]:	The fraction of fraudulent transactions with 'oldBalanceDest' = 'newBalanceDest' = 0 although the transacted 'amount' is non-zero is: 0.4955558261293072 The fraction of genuine transactions with 'oldBalanceDest' = newBalanceDest' = 0 although the transacted 'amount' is non-zero is: 0.0006176245277308345 Since the destination account balances being zero is a strong indicator of fraud, we do not impute the account balance (before the transaction is made) with a statistic or from a distribution with a subsequent adjustment for the amount transacted. Doing so would mask this indicator of fraud and make fraudulent transactions appear genuine. Instead, below we replace the value of 0 with -1 which will be more useful to a suitable machine-learning (ML) algorithm detecting fraud. X.loc[(X.oldBalanceDest == 0) & (X.newBalanceDest == 0) & (X.amount != 0), \ ['oldBalanceDest', 'newBalanceDest']] = -1 The data also has several transactions with zero balances in the originating account both before and after a non-zero amount is transacted. In this case, the fraction of such transactions is much smaller in fraudulent (0.3%) compared to genuine transactions (47%). Once again, from similar reasoning as above, instead of imputing a numerical value we replace the value of 0 with a null value.
In [20]:	X.loc[(X.oldBalanceOrig == 0) & (X.newBalanceOrig == 0) & (X.amount != 0), \ ['oldBalanceOrig', 'newBalanceOrig']] = np.nan back to top Feature-engineering Motivated by the possibility of zero-balances serving to differentiate between fraudulent and genuine transactions, we take the data-imputation of section 3.1 a step further and create 2 new features (columns) recording errors in the originating and destination accounts for each transaction. These new features turn out to be important in obtaining the best performance from the ML algorithm that we will finally use. X['errorBalanceOrig'] = X.newBalanceOrig + X.amount - X.oldBalanceOrig X['errorBalanceDest'] = X.oldBalanceDest + X.amount - X.newBalanceDest
	Data visualization The best way of confirming that the data contains enough information so that a ML algorithm can make strong predictions, is to try and directly visualize the differences between fraudulent and genuine transactions. Motivated by this principle, I visualize these differences in several ways in the plots below. limit = len(X) def plotStrip(x, y, hue, figsize = (14, 9)): fig = plt.figure(figsize = figsize) colours = plt.cm.tab10(np.linspace(0, 1, 9)) with sns.axes_style('ticks'): ax = sns.stripplot(x, y, \ hue = hue, jitter = 0.4, marker = '.', \ size = 4, palette = colours) ax.set_xlabel('')
	ax.set_xticklabels(['genuine', 'fraudulent'], size = 16) for axis in ['top', 'bottom', 'left', 'right']: ax.spines[axis].set_linewidth(2) handles, labels = ax.get_legend_handles_labels() plt.legend(handles, ['Transfer', 'Cash out'], bbox_to_anchor=(1, 1), \
In [22]:	ax = plotStrip(Y[:limit], X.step[:limit], X.type[:limit]) ax.set_ylabel('time [hour]', size = 16) ax.set_title('Striped vs. homogenous fingerprints of genuine and fraudulent \ transactions over time', size = 20); Striped vs. homogenous fingerprints of genuine and fraudulent transactions over time
	genuine fraudulent
	Dispersion over amount The two plots below shows that although the presence of fraud in a transaction can be discerned by the original amount feature, the new errorBalanceDest feature is more effective at making a distinction. limit = len(X) ax = plotStrip(Y[:limit], X.amount[:limit], X.type[:limit], figsize = (14, 9)) ax.set_ylabel('amount', size = 16) ax.set_jlabel('same-signed fingerprints of genuine \ and fraudulent transactions over amount', size = 18); Same-signed fingerprints of genuine and fraudulent transactions over amount Transfer Cash out
	0.6 - UDOU 0.4 -
In [24]:	Dispersion over error in balance in destination accounts
	<pre>limit = len(X) ax = plotStrip(Y[:limit], - X.errorBalanceDest[:limit], X.type[:limit], \</pre>
	e de la companya de l
	genuine fraudulent 5. 4. Separating out genuine from fraudulent transactions The 3D plot below distinguishes best between fraud and non-fraud data by using both of the engineered error-based features. Clearly, the original <i>step</i> feature is ineffective in seperating out fraud. Note the striped nature of the genuine data vs time which was aniticipated from the figure in section 5.1. # Long computation in this cell (-2.5 minutes) x = 'errorBalanceDest'
	<pre>y = 'step' z = 'errorBalanceOrig' zoffset = 0.02 limit = len(X) sns.reset_orig() # prevent seaborn from over-riding mplot3d defaults fig = plt.figure(figsize = (10, 12)) ax = fig.add_subplot(111, projection='3d') ax.scatter(X.loc[Y == 0, x][:limit], X.loc[Y == 0, y][:limit], \ -np.log10(X.loc[Y == 0, z][:limit] + 20ffset), c = 'g', marker = '.', \ s = 1, label = 'genuine') ax.scatter(X.loc[Y == 1, x][:limit], X.loc[Y == 1, y][:limit], \ -np.log10(X.loc[Y == 1, z][:limit] + zoffset), c = 'r', marker = '.', \ s = 1, label = 'frauddlent')</pre>
	<pre>ax.set_xlabel(x, size = 16); ax.set_ylabel(y + ' [hour]', size = 16); ax.set_zlabel('- log\$_{10}\$ (' + z + ')', size = 16) ax.set_title('Error-based features separate out genuine and fraudulent \ transactions', size = 20) plt.axis('tight') ax.grid(1) noFraudMarker = mlines.Line2D([], [], linewidth = 0, color='g', marker='.',</pre>
	Error-based features separate out genuine and fraudulent transactions
	John Standard Standar
In [26]:	Fingerprints of genuine and fraudulent transactions Smoking gun and comprehensive evidence embedded in the dataset of the difference between fraudulent and genuine transactions is obtained by examining their respective correlations in the heatmaps below. Xfraud = X.loc[Y == 1] # update Xfraud & XnonFraud with cleaned data XnonFraud = X.loc[Y == 0]
	<pre>correlationNonFraud = XnonFraud.loc[:, X.columns != 'step'].corr() mask = np.zeros_like(correlationNonFraud) indices = np.triu_indices_from(correlationNonFraud) mask[indices] = True grid_kws = {"width_ratios": (.9, .9, .05), "wspace": 0.2} f, (ax1, ax2, cbar_ax) = plt.subplots(1, 3, gridspec_kw=grid_kws, \</pre>
	<pre>correlationFraud = Xfraud.loc[:, X.columns != 'step'].corr() ax2 = sns.heatmap(correlationFraud, vmin = -1, vmax = 1, cmap = cmap, \ ax = ax2, square = False, linewidths = 0.5, mask = mask, yticklabels = False, \ cbar_ax = cbar_ax, cbar_kws={'orientation': 'vertical', \ 'ticks': [-1, -0.5, 0, 0.5, 1]}) ax2.set_xticklabels(ax2.get_xticklabels(), size = 16); ax2.set_title('Fraudulent \n transactions', size = 20); cbar_ax.set_yticklabels(cbar_ax.get_yticklabels(), size = 14); Genuine transactions Fraudulent transactions type -</pre> 1.0
	amount - oldBalanceOrig - newBalanceOrig - oldBalanceDest -
	newBalanceOrig - errorBalanceOrig - errorBalanceOri
	6. Machine Learning to Detect Fraud in Skewed Data Having obtained evidence from the plots above that the data now contains features that make fraudulent transactions clearly detectable, the remaining obstacle for training a robust ML model is the highly imbalanced nature of the data. print('skew = {}'.format(len(Xfraud) / float(len(X)))) skew = 0.002964544224336551 Split the data into training and test sets in a 80:20 ratio
In [28]:	<pre>trainX, testX, trainY, testY = train_test_split(X, Y, test_size = 0.2, \</pre>
	The figure below shows that the new feature errorBalanceOrig that we created is the most relevant feature for the model. The features are ordered based on the number of samples affected by splits on those features. fig = plt.figure(figsize = (14, 9)) ax = fig.add_subplot(111) colours = plt.cm.Set1(np.linspace(0, 1, 9)) ax = plot_importance(clf, height = 1, color = colours, grid = False, \ show_values = False, importance_type = 'cover', ax = ax); for axis in ['top', 'bottom', 'left', 'right']: ax.spines[axis].set_linewidth(2) ax.set_xlabel('importance score', size = 16); ax.set_ylabel('features', size = 16); ax.set_yticklabels(ax.get_yticklabels(), size = 12);
	ax.set_title('Ordering of features by importance to the model learnt', size = 20); Ordering of features by importance to the model learnt errorBalanceOrig step newBalanceOrig oldBalanceDest
	newBalanceDest amount errorBalanceDest oldBalanceOrig type
	6.2. Visualization of ML model The root node in the decision tree visualized below is indeed the feature <i>errorBalanceOrig</i> , as would be expected from its high significance to the model. to_graphviz(clf)
	errorBalanceOrig<1.13687e-13
	newBalanceOrig<0.105 newBalanceOrig<1.34055e+06
	yes no, missing yes, missing no
	leaf=0.2 amount<989592 step<39 oldBalanceDest<-0.5 yes, missing no yes, missing no leaf=-0.199344 leaf=-0.199149 leaf=-0.190305 leaf=-0.199465 leaf=-0.199851 leaf=-0.196429
	Bias-variance tradeoff The model we have learnt has a degree of bias and is slighly underfit. This is indicated by the levelling in AUPRC as the size of the training set is increased in the cross-validation curve below. The easiest way to improve the performance of the model still further is to increase the max_depth parameter of the XGBClassifier at the expense of the longer time spent learning the model. Other parameters of the classifier that can be adjusted to correct for the effect of the modest underfitting include decreasing min_child_weight and decreasing reg_lambda. # Long computation in this cell (~6 minutes) trainSizes, trainScores, crossValScores = learning_curve(\ XGBClassifier(max_depth = 3, scale_pos_weight = weights, n_jobs = 4), trainX,\ trainY, scoring = 'average_precision')
in [33]:	<pre>trainScoresMean = np.mean(trainScores, axis=1) trainScoresStd = np.std(trainScores, axis=1) crossValScoresMean = np.mean(crossValScores, axis=1) crossValScoresStd = np.std(crossValScores, axis=1) colours = plt.cm.tab10(np.linspace(0, 1, 9)) fig = plt.figure(figsize = (14, 9)) plt.fill_between(trainSizes, trainScoresMean - trainScoresStd,</pre>
	0.998 - 0.996
	0.994 0.992
	7. Conclusion We thoroughly interrogated the data at the outset to gain insight into which features could be discarded and those which could be valuably engineered. The plots provided visual confirmation that the data could be indeed be discriminated with the aid of the new features. To deal with the large skew in the data, we chose an appropriate metric and used an ML algorithm based on an ensemble of decision trees which works best with strongly imbalanced classes. The method used in this kernel should therefore be broadly applicable to a range of such problems.