

# AMATH 582: Final Project

Lee Burke  
Taylor Cunningham

March 16th, 2017

## Abstract

## Sec. I Introduction and Overview

Classifying bank transactions as fraudulent is a difficult task: (1) data is highly sensitive and so difficult to come by, (2) relevant features are not always clear, and (3) classification algorithms have trouble with unbalanced data where fraud is much less common than honest purchases.

To answer the first problem, we analyze real credit card data from a challenge on [link?](#). However, to preserve confidentiality, the features have been obscured by a singular value decomposition (SVD): we are given  $U\Sigma$  if the original data  $A = U\Sigma V^*$  in the usual SVD. That is, we have no say in the second problem, feature extraction. Indeed, the PCA modes have no discernable structure, and decay slowly. Instead, we focus on classifying skewed data.

That popular algorithms have trouble with skew is, in some sense, expected: minimizing the error rate inherently favors the major class. Guessing the larger class every time is already above 50% accuracy. In addition, overall error is not the metric we care most about: false positives (mistakenly flagging a transaction as fraud) have much lower cost to the bank and consumer than false negatives (someone gets away with fraud). Our dataset is a great example of imbalanced data. Only 0.17% of the data points are classified positive. It is difficult to improve on 99.83% accuracy!

We apply several methods learned in lecture, including linear discriminant analysis (LDA), a binary classification tree (BCT), and a support vector machine (SVM). A new method is used to reduce the skew in the class sizes in combination with LDA. We focus especially on how to compare the performance of these classifiers, because common performance metrics may not be as descriptive with high-cost, highly-skewed data.

## Sec. II Theoretical Background

A brief theoretical exploration of each classification scheme and its evaluation with respect to imbalanced data are given below.

### Sec. II.1 Classification

#### Sec. II.1.1 Linear Discriminant Analysis

The LDA algorithm builds normal probability distributions for each class with mean and covariance parameters  $(\vec{\mu}_0, \sigma_0)$   $(\vec{\mu}_1, \sigma_1)$  and  $(\vec{\mu}_1, \sigma_1)$ . These are referred to as the prior distributions. Classifying subsequent observations is then potentially just a matter of comparing their probability of being in class 0 or 1 and applying a threshold. MATLAB actually attempts to minimize the cost function:

$$\hat{y} = \arg \min_{y=1, \dots, K} \sum_{k=1}^K \hat{P}(k|x) C(y|k) \quad (1)$$

where  $\hat{y}$  is the predicted probability.  $K$  is the number of classes.  $\hat{P}$  is the posterior probability of class  $k$  for a given observation  $x$ , and  $C(y|k)$  is the cost of misclassifying  $k$  as  $y$ .  $\hat{P}$  is determined from Bayes

rule:

$$\hat{P}(k|x) = \frac{P(x|k)P(k)}{P(x)} \quad (2)$$

Here the posterior probability that  $x$  is in class  $k$  is what LDA estimates from the training data.  $P(k)$  and  $P(x)$  are readily calculated from the data. Note that for the case of an extremely small minority class such as ours  $P(n)$  will be extremely small (where  $k = n$ ); this means that  $\hat{P}(n|x)$  will be small as well. This is a source of bias for LDA toward the majority class.

### Sec. II.1.2 Binary Classification Trees

Binary decision trees are a common and effective choice for classification tasks. The MATLAB command to build a binary classification decision tree model, `fitctree`, uses an algorithm called CART, for Classification and Regression Tree. It works in a very intuitive way. You might think about how you would classify a transaction as fraud if I knew transaction locations and times. You might say well if a given transaction is more than 100 kilometers away from that person's home zip code then it's likely fraudulent and if it's less than that it is probably not. Similarly we might guess that if a transaction was made before 5am then it is likely fraud, otherwise it's likely not. Binary decision trees ask such questions of our data using the features we give them, figuring out good "decision boundaries" such as 100 and 5am to use. The result is a tree structure where you classify an unknown observation by starting at the root node and asking questions such as the above of that observation. At each node in the tree you make binary decision (meaning there are two branches) and you work your way up the tree until you hit a terminal leaf node. This node determines your prediction.

You can think about this process as splitting up  $m$ -dimensional space into "blocks" (if you imagine the 3-dimensional version), where  $m$  is the number of features. It starts with the full space and splits it into two, then one of those subspaces into two, and then one of those three subspaces into two, and onwards until it finishes. At that point Each of the block represents a terminal node in the tree, and thus a prediction. The algorithm chooses where to put the boundaries based on minimizing the Gini index:

$$G = \sum_k p_k(1 - p_k) \quad (3)$$

Here  $p_k$  is the proportion of points in a subsection of class  $k$ , summed over all classes. Note that the minimum you can achieve is 0 for a pure class and 0.50 for half of each class. It calculates the Gini index for all the potential split points and chooses the lowest one in a greedy fashion—it does not consider global error. The criterion for no longer continuing to split the subspaces is often that you have too few training instances at a node.

In the case of imbalanced data  $p_k$  will always be very small for the minority class in the general case where the classes are highly overlapping. Even if the algorithm has narrowed down to a box that just barely contains all the minority class observations and just a small proportion of the majority class observations, it is still very possible that the number of majority class points in that space will be of a similar number or greater than the minority class. The Gini index naturally gives high weight to the more represented class, which is a bias for the majority class.

### Sec. II.1.3 Support Vector Machines

## Sec. II.2 Handling Class Skew

There are some known methods for improving performance when classifying imbalanced data. The four common solution types are given here:

1. Under-sampling – Removing data from the larger class to balance the class sizes. May lead to a poorer choice of decision line due to losing data at the border of the classes.
2. Oversampling – Adding extra observations on top of existing minority class observations to balance out the class sizes. May lead to overfitting with some classification models.
3. Synthetic data generation – Generating artificial data from your existing data to balance classes. Generated data generally stays within the  $n$ -dimensional volume that minimally encloses the existing data.

4. Cost functions– Classification algorithms use cost functions (decision functions) to define their decision boundaries. With imbalanced data you would set the misclassification of the minority class to be much more costly, to encourage the algorithm to classify them correctly more often than the majority class.

The first three attempt to reduce the imbalance by reducing the number of majority data point or increasing the number of minority data points. These methods deal with the data only, not the classifier. The final option does alter the classifier however, to try to make it classify the minority class more reliably.

## Sec. II.3 Performance Evaluation

### Sec. II.3.1 Receiver Operating Characteristics (ROC) Curve

### Sec. II.3.2 Precision/Recall (PR) Curve

## Sec. III Algorithm Implementation and Development

As stated above our data was prepared and features provided as U columns of the SVD decomposition, so we did not alter or clean the data. We created a “full” data matrix  $X$  from the 28 columns of U and the amount of the transaction. We used standard Z-score of the  $X$  matrix,  $Z$ , for all implementations of the LDA classifier.

Then we split the data for cross validation... or did we do k-fold validation??

Using LDA as our classification algorithm we compared the synthetic data generation and alternate cost function approaches to the default algorithm. For synthetic data generation we chose to use ADASYN (Adaptive Synthetic Sampling Approach for Imbalanced Learning) a generation algorithm similar to the well known SMOTE (Synthetic Minority Over-Sampling Technique) algorithm, which is ADASYN’s precursor. The goal of both is to improve the class balance by increasing the number of minority class members. SMOTE places synthetic data between existing data points randomly (linear interpolation), with no preference shown to any specific points. ADASYN does the same thing but places more synthetic data points close to the boundary between classes because those are the original data points that are more difficult to learn. (Does this favor decision trees or SVMs or something? I imagine that all this would do for an LDA is to move the mean of the minority gaussian close to the boundary... SMOTE might be better here...)

We also tried changing the cost function of our LDA model to discourage FNs. Since LDA works by creating probability distributions the cost comes into play only when making predictions. Each observation that it is trying to predict has a calculated posterior probability. It tries to minimize the “classification cost”—this is a decision cost function. Unfortunately if an observation has an extremely small posterior probability then even with a very large cost for miscalculating that observation it may not change the classification cost by much, meaning there will be little change to the resulting prediction.

## Sec. IV Computational Results

### GIVE BASELINE LDA RESULTS

Altering the prediction cost function does not make a difference in our case. Weighting false negatives as very costly makes very little difference, even when the cost of a false negative is  $10^8$  as costly as a false positive. This seems odd, since we only expect class  $p$  to be about 500 times as likely class  $n$  ( $P(p) \approx 500 \times P(n)$ ). There must be something else going on... PERFORMANCE MATRICS

Synthetic data generation did improve recall. As expected this came at the cost of precision and accuracy. PERFORMANCE METRICS

## Sec. V Summary and Conclusions

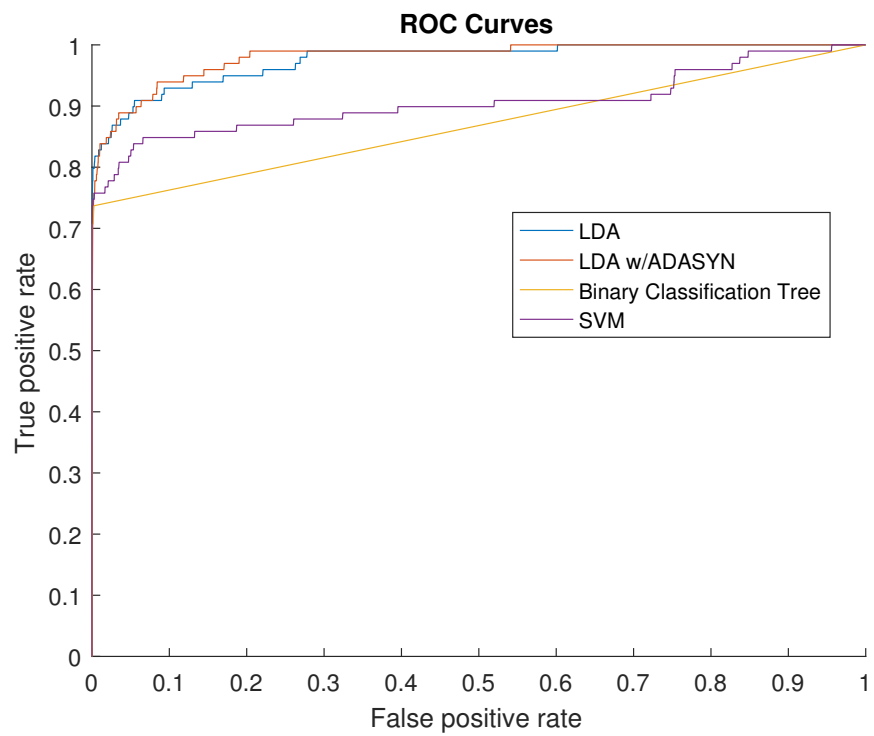


Figure 1:

Table 1:

	AUCs	AUCPR
LDA	0.976	0.124
LDA w/ADASYN	0.980	0.538
Binary Classification Tree	0.868	0.416
SVM	0.905	0.660

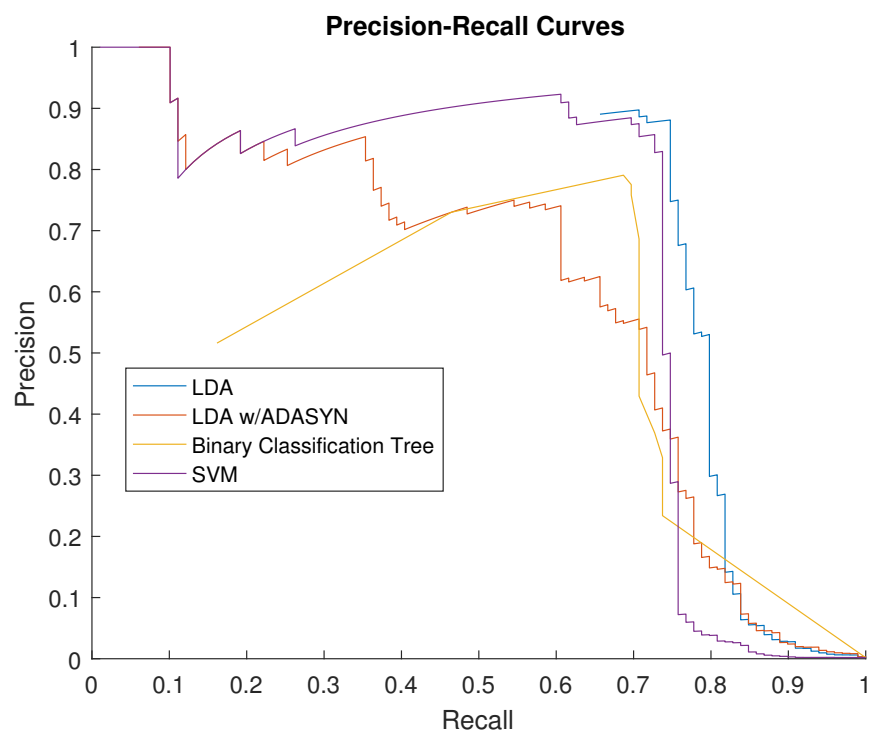


Figure 2:

## References

- [1] <https://www3.nd.edu/~dial/publications/dalpozzolo2015calibrating.pdf>
- [2] wiki: ROC, AUC
- [3] <http://machinelearningmastery.com/classification-and-regression-trees-for-machine-learning/>

## Appendix A MATLAB functions used and brief implementation explanation

- fitctree
- fitcsvm
- fitcdiscr
- ADASYN
- predict
- perfcurve

## Appendix B MATLAB codes

Listing 1: Main

```
1 clear,close all
2 tic
3 disp('Loading data and initializing')
4 %% Parameters
5 % Ratio of training to test subsets
6 ratio = 0.8;
7 % Desired ratio of class 0 to class 1 for ADASYN
8 classBalance = 0.1;
9 % Cost of false negative for cost models
10 costs = [10 100 10^5 10^10];
11
12 %% Initialize variables
13 % amount | amount of transaction in cents
14 % time | time of transaction in seconds
15 % class | boolean fraudulent/legitimate-->1/0
16 % data | features (US in svd USV')
17
18 if exist('data.mat', 'file') == 2
19     load('data.mat')
20 else
21     readData
22 end
23
24 data = [data amount]; % time does not seem to improve results
25
26 %% Test and train subsets
27 [trainData, testData, trainClasses, testClasses] = ...
28     splitBinaryClassData(ratio, data, class);
29
30 %% Undersampling / Oversampling / ADASYN
31 disp(['Running ADASYN (',num2str(toc),')'])
32 % Scaling is required for ADASYN's internal KNN search
33 Ztrain = zscore(trainData);
34 Ztest = zscore(testData);
35 [synZ, synClassOut] = ADASYN(Ztrain, trainClasses, classBalance);
36 synFull = [Ztrain; synZ];
```

```

37 synClass = [trainClasses; synClassOut];
38
39 %% Classify
40 disp(['Building/ running classification models (' ,num2str(toc),')'])
41 % Binary classification decision tree
42 disp(['Binary classification decision tree (' ,num2str(toc),')']);
43 tree = evaluate(fitctree(Ztrain,trainClasses), Ztest, testClasses);
44
45 % Support vector machine
46 disp(['Support vector machine (' ,num2str(toc),')']);
47 svm = evaluate(fitcsvm(Ztrain, trainClasses), Ztest, testClasses);
48
49 % Linear discriminant analysis
50 disp(['Linear discriminant analyses (' ,num2str(toc),')']);
51 lda = evaluate(fitcdiscr(Ztrain, trainClasses), Ztest, testClasses);
52
53 % LDA with synthetic data -- train with syn data but test with original
54 ldaSyn = evaluate(fitcdiscr(synFull, synClass), Ztest, testClasses);
55
56 % LDA with custom cost function
57 costModel = fitcdiscr(Ztrain, trainClasses,'ClassNames',[0,1]);
58 for i=1:length(costs) % TODO is something wrong here?? all perf. values are ...
    exactly the same as vanilla LDA.....
59     costModel.Cost = [0 1; costs(i) 0];
60     ldaCosts(i) = evaluate(costModel, testData, testClasses);
61 end
62
63 disp(['Done classifying (' ,num2str(toc),')'])
64
65
66 %% Plot
67 % lda
68 plotModels = {lda};
69 modelNames = {'LDA'};
70 AUCs = lda.AUC;
71 AUCPRs = lda.AUCPR;
72 % ldaSyn
73 plotModels = {plotModels{:},ldaSyn};
74 modelNames = {modelNames{:}, 'LDA w/ADASYN'};
75 AUCs = [AUCs;ldaSyn.AUC];
76 AUCPRs = [AUCPRs;ldaSyn.AUCPR];
77 % tree
78 plotModels = {plotModels{:},tree};
79 modelNames = {modelNames{:}, 'Tree'};
80 AUCs = [AUCs;tree.AUC];
81 AUCPRs = [AUCPRs;tree.AUCPR];
82 % svm
83 plotModels = {plotModels{:},svm};
84 modelNames = {modelNames{:}, 'SVM'};
85 AUCs = [AUCs;svm.AUC];
86 AUCPRs = [AUCPRs;svm.AUCPR];
87 % ldaCosts
88 for i=1:length(costs)
89     plotModels = {plotModels{:},ldaCosts(i)};
90     modelNames = {modelNames{:}, ['LDA w/Costs ' ,num2str(i)]};
91     AUCs = [AUCs;ldaCosts(i).AUC];
92     AUCPRs = [AUCPRs;ldaCosts(i).AUCPR];
93 end
94
95 T = table(AUCs,AUCPRs,'RowNames',modelNames)
96
97 % ROC Curve
98 figure
99 hold on
100 for i=1:length(plotModels)

```

```

101     mdl = plotModels{i};
102     % ...
103     errorbar(mdl.rocX,mdl.rocY(:,1),mdl.rocY(:,1)-mdl.rocY(:,2),mdl.rocY(:,3)-mdl.rocY(:,1));
104     plot(plotModels{i}.rocX,plotModels{i}.rocY);
105 end
106 hold off
107 legend(modelNames,'Location','Best');
108 xlabel('False positive rate')
109 ylabel('True positive rate')
110 title('ROC Curves')
111
112 % PR Curve
113 figure
114 hold on
115 for i=1:length(plotModels)
116     mdl = plotModels{i};
117     % ...
118     errorbar(mdl.prX,mdl.prY(:,1),mdl.prY(:,1)-mdl.prY(:,2),mdl.prY(:,3)-mdl.prY(:,1));
119     plot(plotModels{i}.prX,plotModels{i}.prY);
120 end
121 hold off
122 legend(modelNames,'Location','Best')
123 xlabel('Recall')
124 ylabel('Precision')
125 title('Precision-Recall Curves')
126
127
128
129
130
131
132
133
134
135
136 %% Explore feature space
137 % Singular values and POD modes
138 % S = zeros(size(data,2),1);
139 % U = zeros(size(data));
140 % for i=1:size(data,2)
141 %     S(i)=norm(data(:,i));
142 %     U(:,i)=data(:,i)/S(i);
143 % end
144 % plot(S)
145 % plotFeatureSpace(U(class==1,:), U(class==0,:), [1 2 3]);
146 % plotFeatureSpace(posData, negData, [3 10 29]);
147 %% Run model nIter times
148 % nIter = 1;
149 % disp(['Starting trials ',num2str(toc),''])
150 % for i=1:nIter
151 %     disp(['Running test ',num2str(i),' out of ',num2str(nIter)])
152 %     t = tic;
153 %     t = toc(t);
154 %     disp(['Test ',num2str(i),' took ',num2str(t),' seconds.'])
155 % end
156 %% NOTE: if we are going to do cross validation MATLAB suggests:
157 %     cvmodel = crossval(lda,'kfold',2);
158 %     cverror = kfoldLoss(cvmodel)

```

Listing 2: Split data into training and test sets



```

1 function [train, test, trainClass, testClass] = splitBinaryClassData(ratio, ...
    data, class)
2     n0 = sum(class==0);
3     n1 = sum(class==1);
4
5     data0 = data(class == 0,:);
6     data1 = data(class == 1,:);
7
8     [iTrain0, iTest0] = splitIndices(n0,ratio);
9     [iTrain1, iTest1] = splitIndices(n1,ratio);
10
11     train = [data1(iTrain1,:); data0(iTrain0,:)];
12     test = [data1(iTest1,:); data0(iTest0,:)];
13
14     trainClass = [ones(length(iTrain1),1); zeros(length(iTrain0),1)];
15     testClass = [ones(length(iTest1),1); zeros(length(iTest0),1)];
16 end

```

Listing 3: Random permutation for indices

```

1 function [i_train, i_test, numTrain, numTest] = splitIndices(n, ratio)
2     numTrain = floor(ratio*n);
3     numTest = n - numTrain;
4
5     % Mix it up
6     ri = randperm(n);
7     i_train = ri(1:numTrain);
8     i_test = ri(numTrain+1:end);
9 end

```

Listing 4: Performance characteristics of a classification model

```

1 function s = evaluate(model, testData, testClasses)
2     s = struct;
3     s.model = model;
4     [s.labels,s.scores,s.costs] = predict(model, testData);
5     s.confusion = confusionmat(testClasses,s.labels);
6     [s.rocX,s.rocY,s.rocT,s.AUC] = perfcurve(testClasses,s.scores(:,2),'1');
7     [s.prX,s.prY,s.prT,s.AUCPR] = ...
        perfcurve(testClasses,s.scores(:,2),'1','XCrit','reca','YCrit','prec');
8     %     disp('Bootstrapping ROC')
9     %     [s.rocX,s.rocY,s.rocT,s.AUC] = ...
        perfcurve(testClasses,s.scores(:,2),'1','NBoot',1,'TVals',0:0.05:1);
10    %     disp('Bootstrapping PR')
11    %     [s.prX,s.prY,s.prT,s.AUCPR] = ...
        perfcurve(testClasses,s.scores(:,2),'1','XCrit','reca','YCrit','prec','NBoot',1,'TVals',0:0.05:1);
12 end

```