== Goals ==

Create a model to predict churn and help churn reduction.
Try different tools provided by scikit-learn package, including PCA and PLS.


== Work flow ==

Data was extracted from the spreadsheet and samples were shuffled.
Data forms predictor and response matrices (X,Y) (test,train) with float values.
Three unique area codes were converted to numeric variable.
All these codes were from California, and state column was decided to be ignored.

Churn prediction problem can be solved as binary classification.
The classifier output can be a real value (continuous).
For example, it can be probability of churn.
In this case the boundary between classes is determined by a threshold value.
Result of the prediction can contain errors of two types:
1) false positive (churn predicted, actually no churn),
2) false negative (no churn predicted, actually was churn).
Taking different threshold values, we can reduce one, enlarging other.
So this value is determined by these error comparative costs.
Quality of classifier can be estimated
    by ROC curve (receiver operating characteristic) and AUC value.

Linear regression model was successfully tried.
PLS model behaved a slightly better than simple linear, but not notably.
AUC > 0.8 (blue ROC curves).

PCA was tried and score and loadings plots reviewed.
PCA is a way to reduce dimensionality of data --
    compose principal components from redundant variables (linear combination).
Score plots show samples in the space of some pair of principal components.
Loadings plots show how original variables affect pair of principal components.
One of score plots showed two clusters of samples.
The corresponding loadings plot showed connection with vmail* and night* stuff.
However clusters were not stable and can disappear with another training subset.

PCA mostly works better with scaled data and can be used in a PCR pipeline, like:
    scaling --> PCA --> linear regression.
Dimensionality reduces to about 13 components.
But the real power of PCR and PLS can be seen
    when number of variables is more than number of samples.
In our case this worked as good as simple linear regression.

Logistic regression and support vector machines were briefly tried,
but not successfully due to bad tuning or overall incompatibility with the task.

1

The best result was given by using decision tree classifier (green ROC curve).
Tree with maximal depth 5 seems to work well.
Making it bigger had not much improved the model or even made it overfitted.
Pipelining with PCA was not successful.

Decision tree is well observable (white box).
However it greatly changes appearance with training set reselection.
Variables, that was mostly mentioned in the tree:
    CustServ Calls, Day Mins, Int'l Plan, VMail Message, Eve Mins.
Suggestions are that customers are likely to churn, having:
-- int'l plan,
-- talking much w/o vmail messages,
-- calling support (m. b. angry with problems) and not talking much.
These can happen due to customers personal unstable behaviour.
These can happen due to corresponding competitor's proposal.

== More things to try ==

Clusters discovered by PCA can be then handled by separate models,
then accuracy may rise.

Variables have different nature.
It may be good to observe their distribution and preprocess them separately,
    for example, applying logarithmic scale.

States and codes are not comparable as numbers.
It may be better to map each of them to
    as many variables as there are different options, getting sparse matrix.
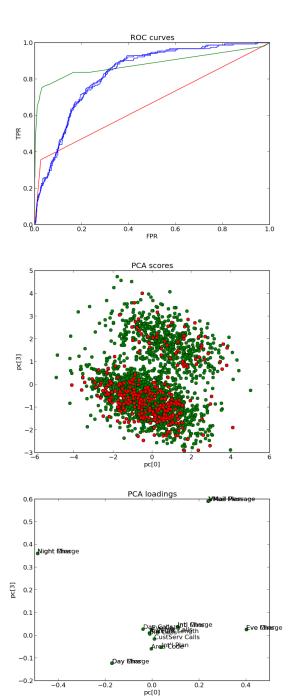For example, we'll have zeroes and ones in column California.
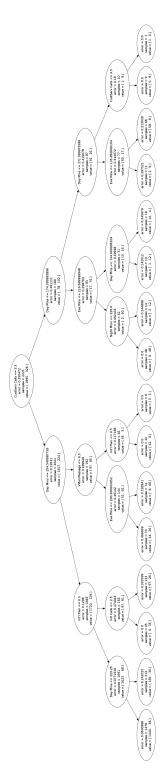
Model may work in pair.
Samples may be handled by decision tree,
    but unclear ones forwarded to another model.

Bayesian networks and neural networks can be tried.

Decision tree options (like min_samples_leaf) can be tuned.

## ROC curves



## PCA scores



## PCA loadings



3

Decision tree (read top-to-bottom, left-to-right):

- **CustServ Calls <= 3.5** — error = 0.20472, samples = 2221, value = [1895., 326.]
  - **Day Mins <= 254.549897793** — error = 0.15041, samples = 2041, value = [1817., 224.]
    - **Intl Plan <= 0.5** — error = 0.1288, samples = 1849, value = [1720., 129.]
      - **Day Mins <= 223.25** — error = 0.077316, samples = 1691, value = [1622., 69.]
        - error = 0.090988, samples = 1478, value = [1442., 36.]
        - error = 0.242015, samples = 213, value = [180., 33.]
      - **Intl Calls <= 2.5** — error = 0.474043, samples = 158, value = [97., 61.]
        - error = 0.0, samples = 35, value = [0., 35.]
        - error = 0.333399, samples = 123, value = [97., 26.]
    - **VMail Message <= 6.5** — error = 0.499946, samples = 192, value = [97., 95.]
      - **Eve Mins <= 200.300002052** — error = 0.49542, samples = 144, value = [52., 92.]
        - error = 0.466959, samples = 70, value = [44., 26.]
        - error = 0.192841, samples = 74, value = [8., 66.]
      - **Intl Plan <= 0.5** — error = 0.111288, samples = 48, value = [45., 3.]
        - error = 0.0, samples = 42, value = [42., 0.]
        - error = 0.5, samples = 6, value = [3., 3.]
  - **Day Mins <= 174.399993896** — error = 0.481111, samples = 180, value = [78., 102.]
    - **Eve Mins <= 216.699996948** — error = 0.298753, samples = 93, value = [17., 76.]
      - **Night Mins <= 226.0** — error = 0.0526435, samples = 62, value = [2., 60.]
        - error = 0.0, samples = 48, value = [0., 48.]
        - error = 0.244898, samples = 14, value = [2., 12.]
      - **Day Mins <= 134.600000104** — error = 0.49948, samples = 31, value = [15., 16.]
        - error = 0.142012, samples = 13, value = [1., 12.]
        - error = 0.340579, samples = 18, value = [14., 4.]
    - **Day Mins <= 271.099975586** — error = 0.110078, samples = 87, value = [61., 26.]
      - **Eve Mins <= 135.800000104** — error = 0.344672, samples = 77, value = [60., 17.]
        - error = 0.297021, samples = 11, value = [2., 9.]
        - error = 0.213039, samples = 66, value = [58., 8.]
      - **CustServ Calls <= 4.5** — error = 0.18, samples = 10, value = [1., 9.]
        - error = 0.0, samples = 9, value = [0., 9.]
        - error = 0.0, samples = 1, value = [1., 0.]

See python sources at https://github.com/rewlad/ttu_andm

```
1
2  import glob
3  import csv
4  import numpy as np
5  import matplotlib.pyplot as plt
6  from sklearn import pipeline, preprocessing, decomposition,
       linear_model, pls, tree, metrics, utils
7
8  def load_rows():
9      with open('Churn.csv','rb') as csvfile:
10         reader = csv.reader(csvfile,delimiter=',',quotechar='"')
11         rows = []
12         for rn,row in enumerate(reader):
13             if rn==0: head = row[0:22]
14             else: rows.append(row[0:22])
15         return rows, head
16
17 def rows_to_predictor_response(rows, head):
18     codes = [row[19] for row in rows]
19     codes_h = { v:j for j,v in enumerate(set(codes)) }
20     for row in rows:
21         for cn, v in enumerate(row):
22             row[cn] = float(v) if cn<18 else codes_h[v] if cn==19
                   else 0
23     arows = utils.shuffle(np.array(rows), random_state=0)
24     x_row_indexes = [0,1,2,3,4,5,6, 8,9,10,11,12,13,14,15,16,17,
           19]
25     x_head = np.array(head)[x_row_indexes]
26     X = arows[:, x_row_indexes]
27     Y = arows[:, [7]]
28     rng, l = np.arange(len(X)), len(X)/3
29     train, test = rng > l, rng <= l
30     return X, Y, x_head, train, test
31
32 def mode_check_rows():
33     rows, head = load_rows()
34     print x_head
35     phones = [row[20] for row in rows]
36     print len(phones), len(set(phones)) #checking phones are uniq
37     codes = [row[19] for row in rows]
38     states = [row[18] for row in rows]
39     print set(codes), set(states) #codes '415', '510', '408' all
           from California
40
41     X, Y, x_head, train, test = rows_to_predictor_response(rows,
           head)
42     print head, x_head, head[7]
43     print X.shape, Y.shape
44     print len(Y[train]), len(Y[test])
45     print np.count_nonzero(Y[train]), np.count_nonzero(Y[test])
46     print np.count_nonzero(Y)/float(len(Y))
47     print np.count_nonzero(Y[train])/float(len(Y[train]))
48     print np.count_nonzero(Y[test])/float(len(Y[test]))
49
50 def plot_rocs(nm,color,y_true,y_pred):
51     fpr, tpr, thresholds = metrics.roc_curve(y_true=y_true, y_score
```

```
                     =y_pred)
52          print metrics.auc(fpr, tpr)
53          if len(fpr)<4: print fpr, tpr
54          path = "out/roc."+nm+".npz"
55          np.savez(path,fpr=fpr, tpr=tpr, thresholds=thresholds, color=np
                 .array([color]))
56          plt.title('ROC␣curves')
57          plt.xlabel('FPR')
58          plt.ylabel('TPR')
59          for fn in glob.glob("out/roc.*.npz"):
60              v = np.load(fn)
61              plt.plot(v['fpr'], v['tpr'], v['color'][0])
62          plt.savefig("out/roc.png")
63          plt.cla()
64
65
66  def mode_pca():
67      rows, head = load_rows()
68      X, Y, x_head, train, test = rows_to_predictor_response(rows,
             head)
69
70      pca = decomposition.PCA(n_components=13)
71      re_pipeline = pipeline.Pipeline([ ('scaler',preprocessing.
             Scaler()), ('pca',pca) ])
72      pc = re_pipeline.fit_transform(X[train])
73
74      churn = Y[train,0] > 0.5
75      for i in range(1,13):
76          plt.title('PCA␣scores')
77          plt.xlabel('pc[0]')
78          plt.ylabel('pc['+str(i)+']')
79          plt.plot(pc[:,0], pc[:,i], 'go')
80          plt.plot(pc[churn,0], pc[churn,i], 'ro')
81          plt.savefig("out/churn_scores_0_"+str(i)+".png")
82          plt.cla()
83
84      loadings = pca.components_
85      for i in range(1,13):
86          plt.title('PCA␣loadings')
87          plt.xlabel('pc[0]')
88          plt.ylabel('pc['+str(i)+']')
89          plt.plot(loadings[0], loadings[i], 'go')
90          for j,l in enumerate(loadings[[0,i]].T):
91              plt.annotate(x_head[j],l); #-(j%3)*0.02
92              #print x_head[i],l
93          plt.savefig("out/churn_loadings_0_"+str(i)+".png")
94          plt.cla()
95
96      y_scaler = preprocessing.Scaler(with_std=False)
97      linre = linear_model.LinearRegression()
98      linre.fit(X=pc, y=y_scaler.fit_transform(Y[train,0]))
99
100     y_pred = y_scaler.inverse_transform( linre.predict( re_pipeline
             .transform(X[test]) ) )
101
102     plot_rocs('pca','b-',Y[test,0],y_pred)
103
```

```
104  def mode_pls():
105      rows, head = load_rows()
106      X, Y, x_head, train, test = rows_to_predictor_response(rows,
             head)
107
108      re = pls.PLSRegression(n_components=11)
109      re.fit(X=X[train], Y=Y[train])
110      y_pred = re.predict( X[test] )[:,0]
111
112      plot_rocs('pls','b-',Y[test,0],y_pred)
113
114  def mode_linre():
115      rows, head = load_rows()
116      X, Y, x_head, train, test = rows_to_predictor_response(rows,
             head)
117
118      linre = linear_model.LinearRegression()
119      linre.fit(X=X[train], y=Y[train,0])
120      y_pred = linre.predict( X[test] )
121
122      plot_rocs('linre','b-',Y[test,0],y_pred)
123
124  def mode_tree():
125      rows, head = load_rows()
126      X, Y, x_head, train, test = rows_to_predictor_response(rows,
             head)
127
128      clf = tree.DecisionTreeClassifier(max_depth=5)#,
             min_samples_leaf?
129      clf.fit(X=X[train], y=Y[train,0]) #*2-1
130      #y_pred = clf.predict( X[test] )
131      y_pred = clf.predict_proba( X[test] )[:,1]
132
133      plot_rocs('tree','g-',Y[test,0],y_pred)
134
135      tree.export_graphviz(clf, out_file='dtree.graphviz',
             feature_names=x_head)
136      #convert using: dot dtree.graphviz -Tpng > dtree2.png
137
138  mode_linre()
139  mode_pca()
140  mode_pls()
141  mode_tree()
```