

Data Science Practicum Report

Mar 8th 2017

Preprocessing

As agreed in the meeting on Thursday 23rd, in addition to remove columns with >80% missing values, we kept columns suggested by Joel in the excel file sent on Monday Feb 27th. Our final list of 79 features was as follows:

Sport	ECCNote	Eparent_welcome
EventsAttended	ECCNote_camp	Epost_event_email
Hprofileview	Ecoach_list_known_updated	Esms_update
Hcoachimport	ECoachEmailOpen	CollegeProspects
Hmessage	ECoachEval	MessagesReceived
Hsearchhit	ECoachImport	MessagesSent
Hcoacheval	ECoachSearchHit	CaptainU_CHURN
Hemailopen	ECoachVisit	NumYear
EAthlete newsletter	Ecolleges_going_to_the_event	NumMonth
Eathlete_new	Efailed_subscription	monthly_price
Eathlete_new_info_request	EEmailsDigest	Eparent_new
Gender	Hprofileview_Freq	Hcoachimport_Freq
Hmessage_Freq	Hsearchhit_Freq	Hcoacheval_Freq
Ecolleges_going_to_the_event_Freq	Efailed_subscription_Freq	Esms_update_Freq
Hits_Frequency	College_Prospects_Frequency	
Hemailopen_Freq	Ecoach_list_known_updated_Freq	ECoachVisit_Freq

We then took columns with continuous values and normalized the values by the mean and the standard deviation:

$$z = \frac{x - \mu}{\sigma}$$

Hence each feature had a mean of 0 and a standard deviation of 1. Below is a snapshot of the data:

```
In [379]: std_pd.head(5)
```

Out[379]:

Hprofileview	Hcoachimport	Hmessage	Hsearchhit	Hcoacheval	Hemailopen	Eathlete newsletter	Eathlete_new	Eathlete_new_info_request	...
-0.334910	-0.082896	-0.118345	-0.863017	-0.034355	-0.132922	-1.570523	-0.050501	-0.012864	...
-0.334910	-0.082896	-0.118345	-0.063949	-0.034355	-0.132922	1.506564	-0.050501	-0.012864	...
0.663747	0.843619	0.773367	1.001475	-0.034355	-0.132922	0.480869	-0.050501	-0.012864	...
-0.334910	-0.082896	2.556791	-0.197127	-0.034355	4.445944	0.480869	-0.050501	-0.012864	...
-0.334910	-0.082896	-0.118345	-0.330305	-0.034355	-0.132922	0.480869	-0.050501	-0.012864	...

We then created dummy variables from features with categorical features: gender and sport.

Final table shape: **16117 rows 79 columns**

Predicting Lifetime Churn

This week, we focused on predicting lifetime churn. We implemented three machine-learning models recorded and visualized Precision and Recall Values.

Below are the three models we chose:

- Decision Trees
- Logistic Regression
- Support Vector Machines (SVM)

When it came to building models, we used 80% of the data (12893 rows) for training and the remaining 20% for testing.

Training data: **12893 rows 79 columns**

Test data: : **3224 rows 79 columns**

Class distribution of Churners and Non-Churners

Training set

Non-Churners	8083
Churners	4810

Testing Set

Non-Churners	1991
Churners	1233

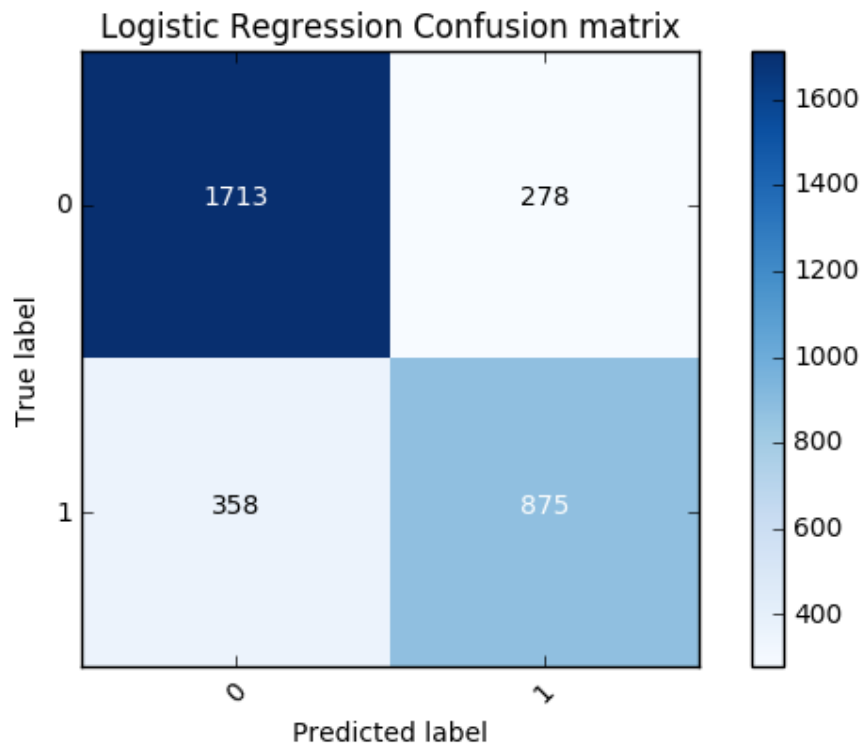
Summary of Precision Recall Values of Churning

Model	Precision	Recall	F1 Score
Decision Trees	0.73	0.72	0.73
Logistic Regression	0.76	0.71	0.73
SVM	0.77	0.68	0.72

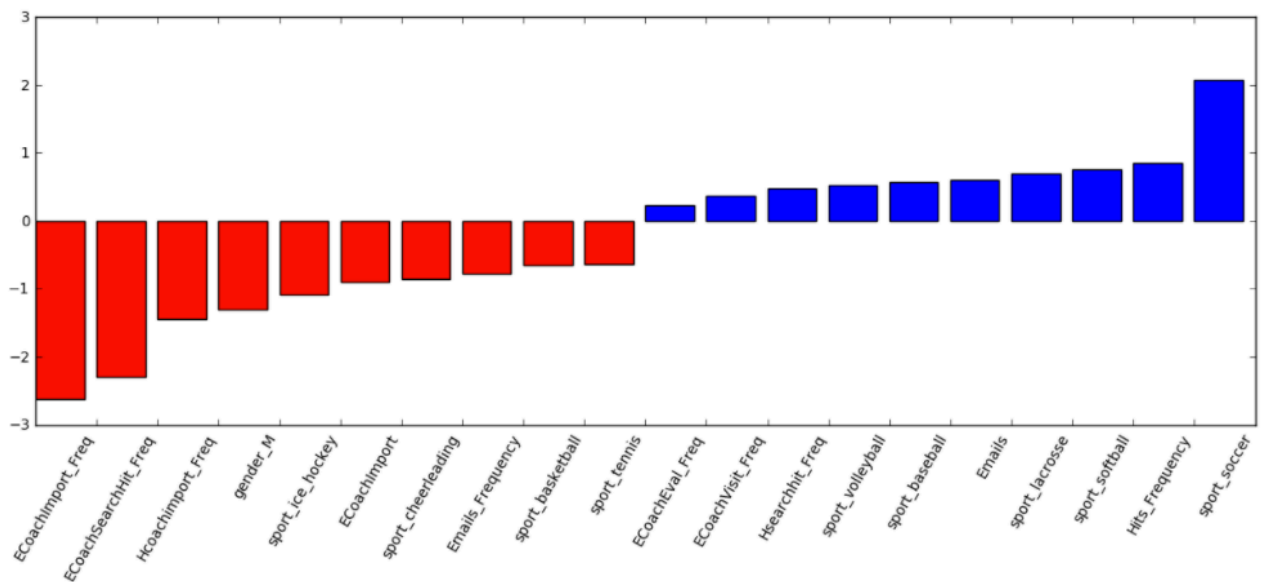
Summary of Precision Recall Values of Retention

Model	Precision	Recall	F1 Score
Decision Trees	0.83	0.83	0.83
Logistic Regression	0.83	0.86	0.84
SVM	0.82	0.82	0.82

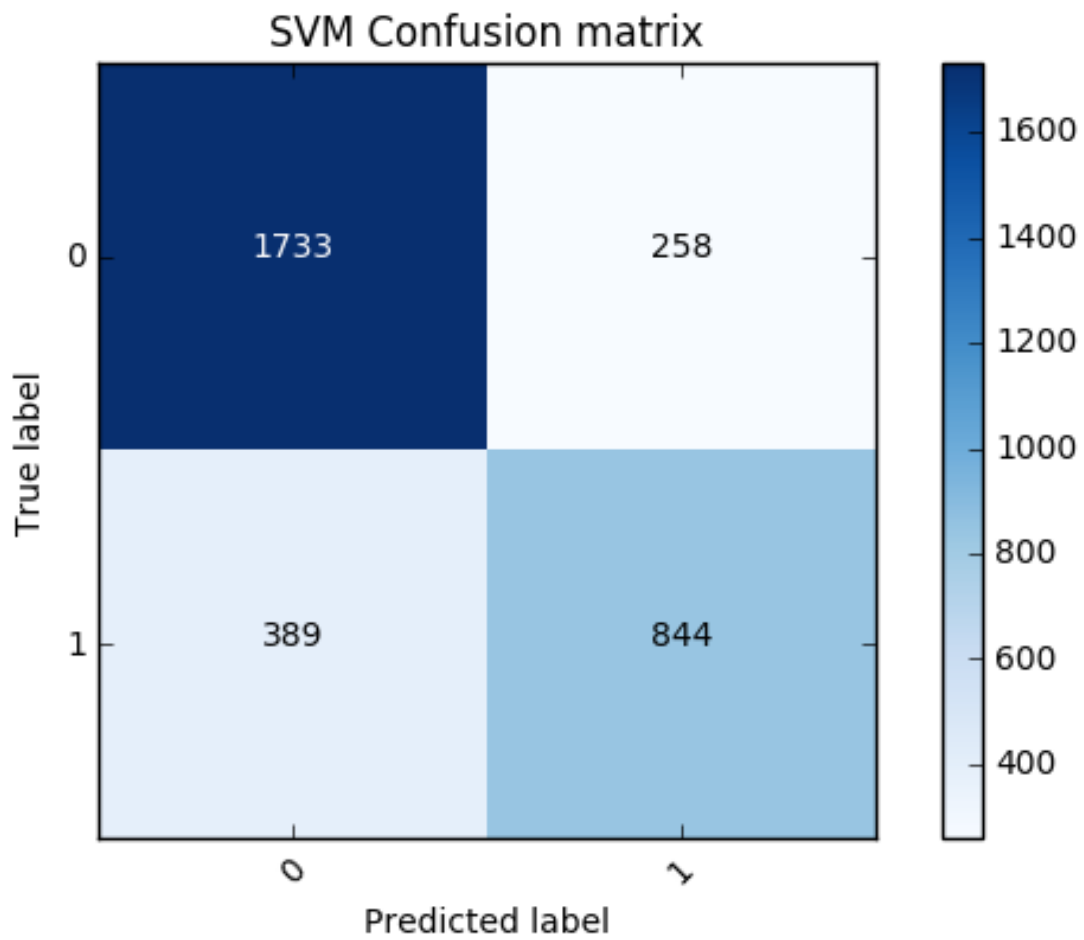
Logistic Regression Visualization



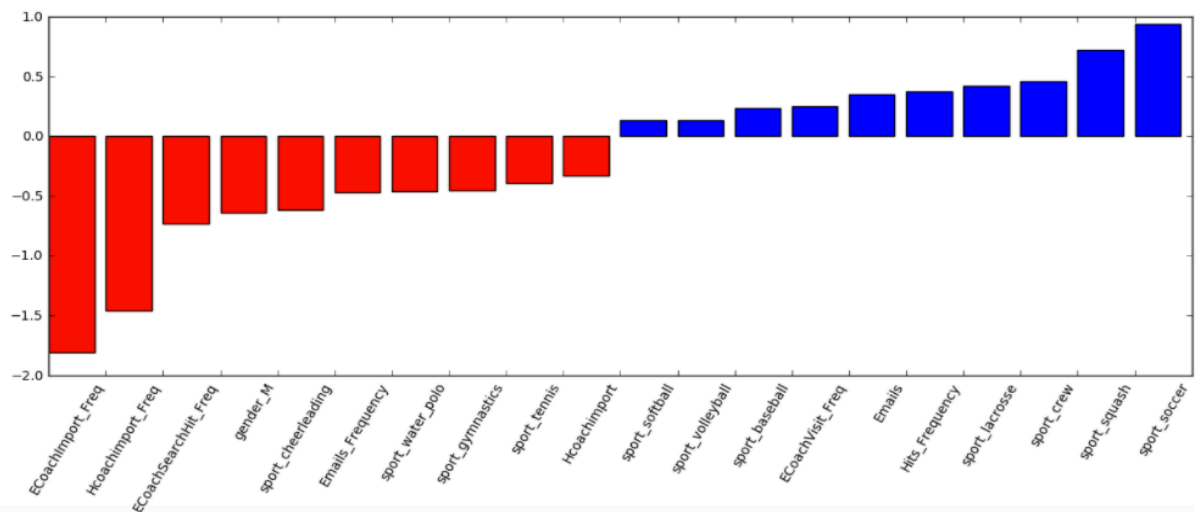
Important Features According to Logistic Regression



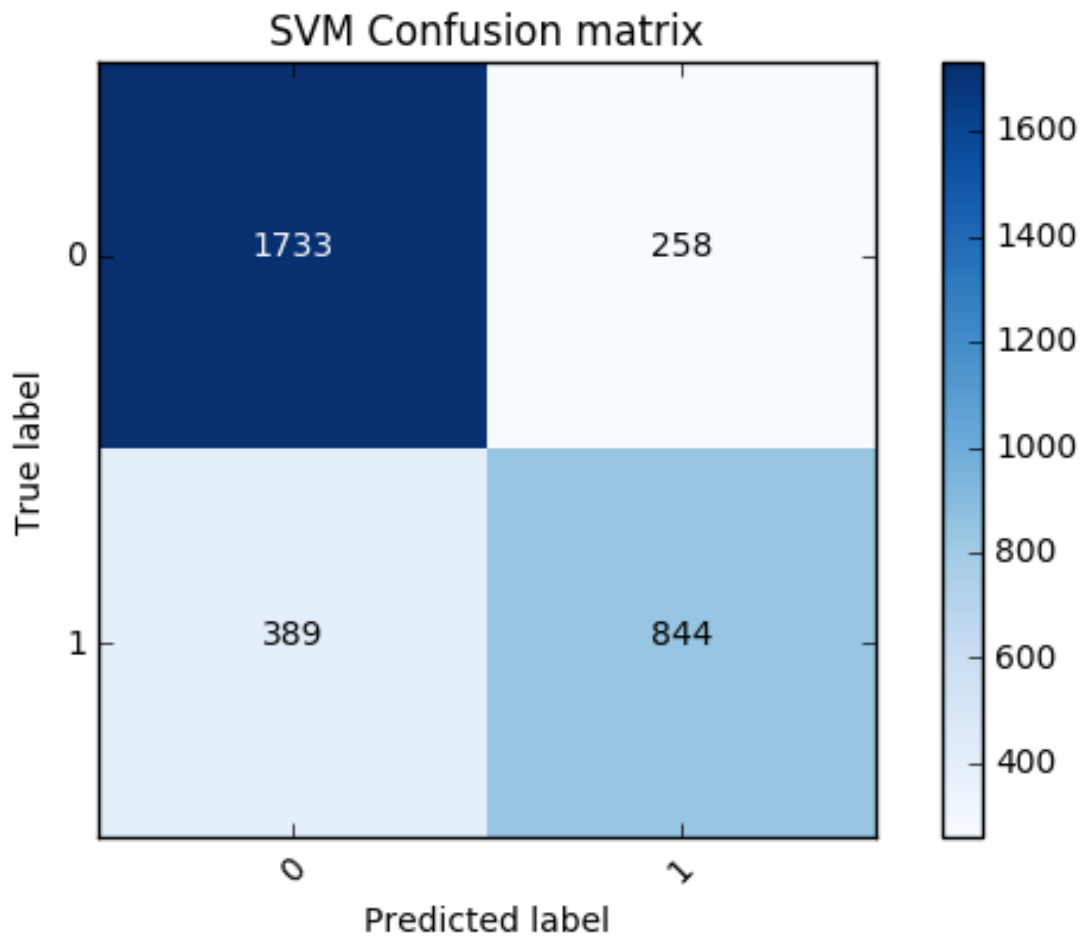
Support Vector Machine Visualization



Important Features According to SVM



Decision Trees



Tree Visualization

