



# B2B Customer Segmentation and Churn Prediction

**Alon Benhaim**

# Business Challenge

- We are given random sample of merchant transaction activity (date, time, and amount), the goal is to try and answer the following questions:
  - Infer the types of merchants using the platform
  - Define, identify and then predict future churning merchants
- In data science language our goals are:
  - Merchant (B2B) segmentation – Use unsupervised learning
  - Merchant churn prediction – Build supervised learning model

# Dataset

- ~1.5M transactions of 14,351 different merchants over 2 years (from 01/01/2033 to 12/01/2034)
- Number of transactions per business ranges from 1 up to 25K
- There are no refunds, no missing data values and all transactions are within the 2 years range

Unnamed: 0		merchant	time	amount_usd_in_cents
0	1	faa029c6b0	2034-06-17 23:34:14	6349
1	2	ed7a7d91aa	2034-12-27 00:40:38	3854
2	3	5608f200cf	2034-04-30 01:29:42	789
3	4	15b1a0d61e	2034-09-16 01:06:23	4452
4	5	4770051790	2034-07-22 16:21:42	20203
...	...	...	...	...
1513714	1513715	72d37bedbf	2034-06-21 13:47:51	5274
1513715	1513716	5608f200cf	2034-04-20 02:23:59	754
1513716	1513717	fcbd1dae68	2033-09-19 14:02:33	13203
1513717	1513718	9843e52410	2034-12-28 20:07:59	4845
1513718	1513719	32acddd6cc	2034-08-23 09:07:07	3862

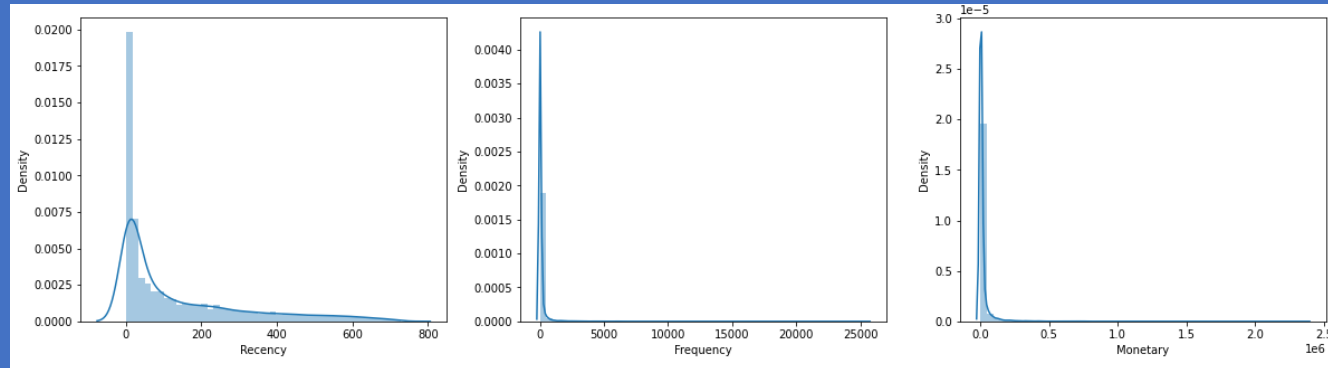


Kmeans Clustering

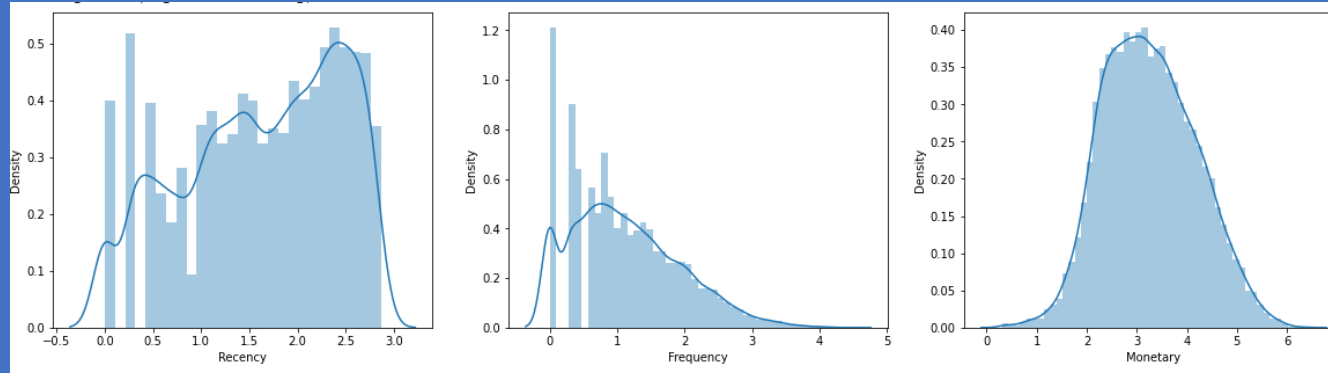
# RFM Analysis and data preprocessing

- Recency – # days to recent purchase
- Frequency – Total # of transactions
- Monetary (USD) – Total amount of transactions

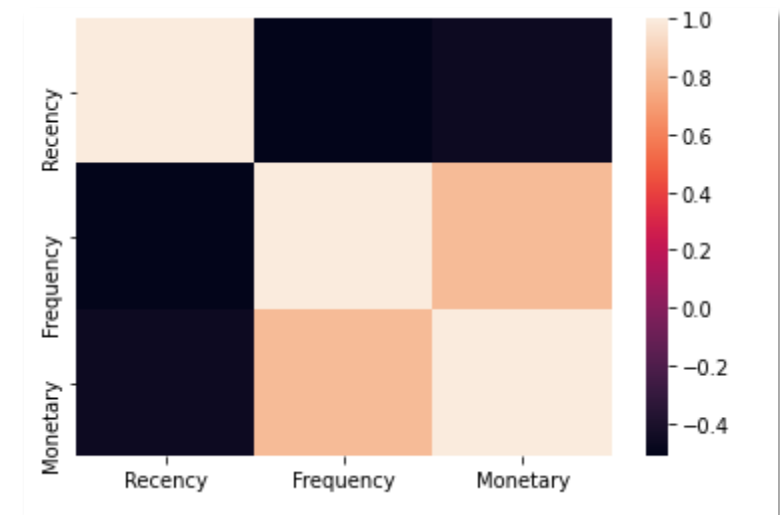
	Recency	Frequency	Monetary
merchant			
0002b63b92	595	1	33.79
0002d07bba	17	4	892.78
00057d4302	515	28	295.21
000bcff341	510	1	78.26
000ddb0ca	578	1	102.99



## Log transformation

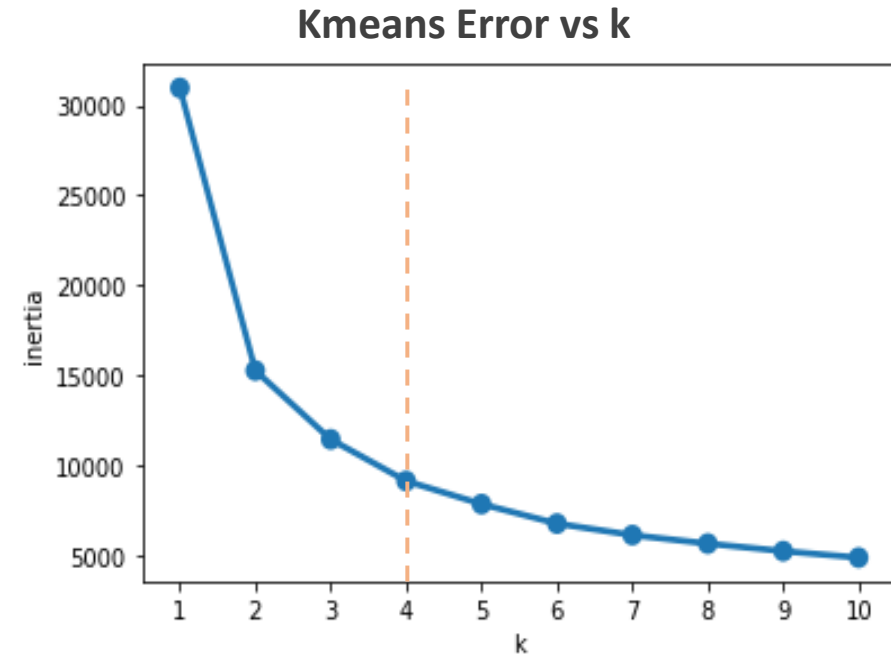
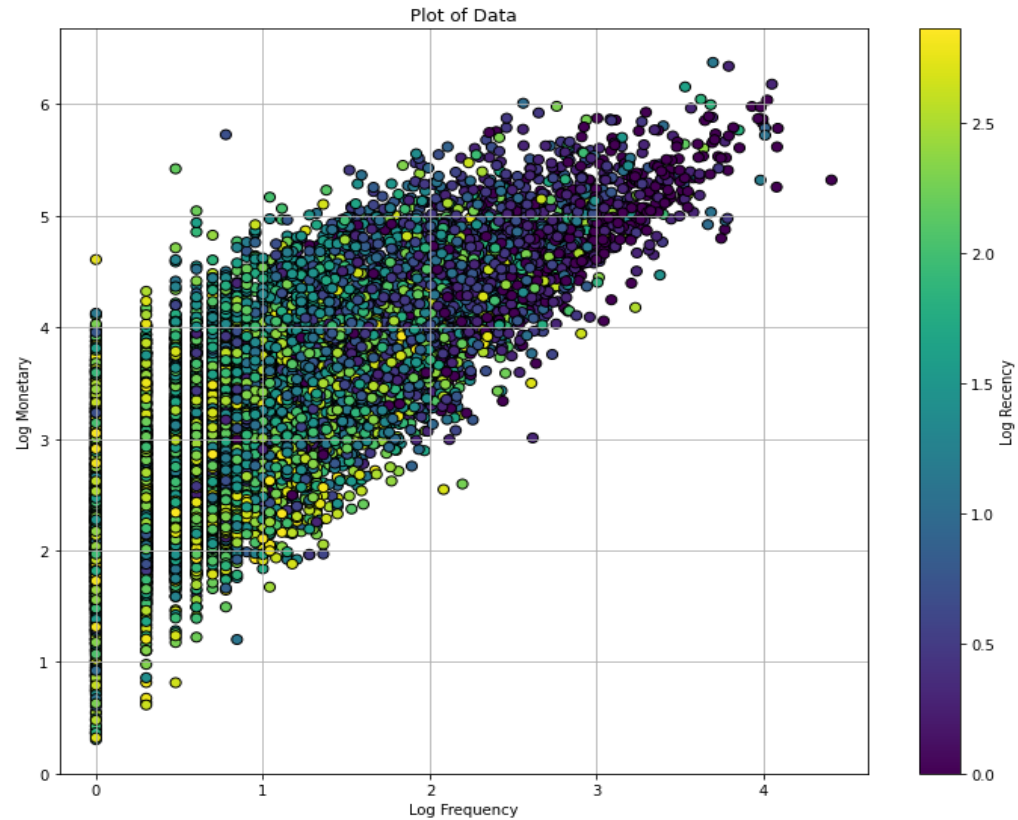


## Feature Correlation



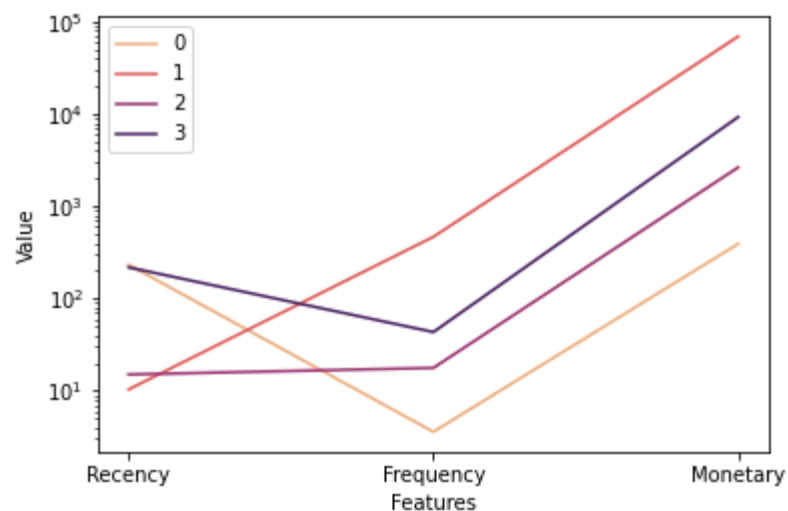
# Modelling

- Elbow method confirms  $k=4$  clusters segmentation



# Results

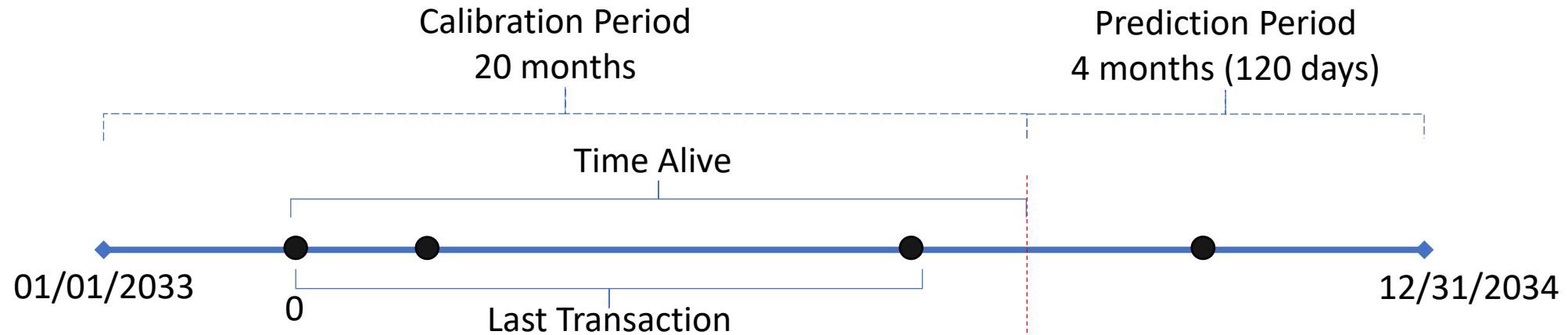
- Types of merchants for each cluster type



	Recency	Frequency	Monetary	Cluster
Cluster				
0	232.848101	3.587223	392.237116	5056
1	10.348753	467.886881	69718.396979	2767
2	15.030303	17.700406	2651.202390	3201
3	218.172528	43.366396	9321.851178	3327

Cluster	Merchant Type	RFM attributes
0	Churner merchant (lost)	Old last transaction, with low frequency and monetary.
1	Best merchant	Recent last transaction, with high frequency and monetary
2	New merchant	Recent last transaction, with low frequency and monetary.
3	At risk/churner merchant	Old last transaction, with moderate frequency and monetary.

# Churn Prediction Model

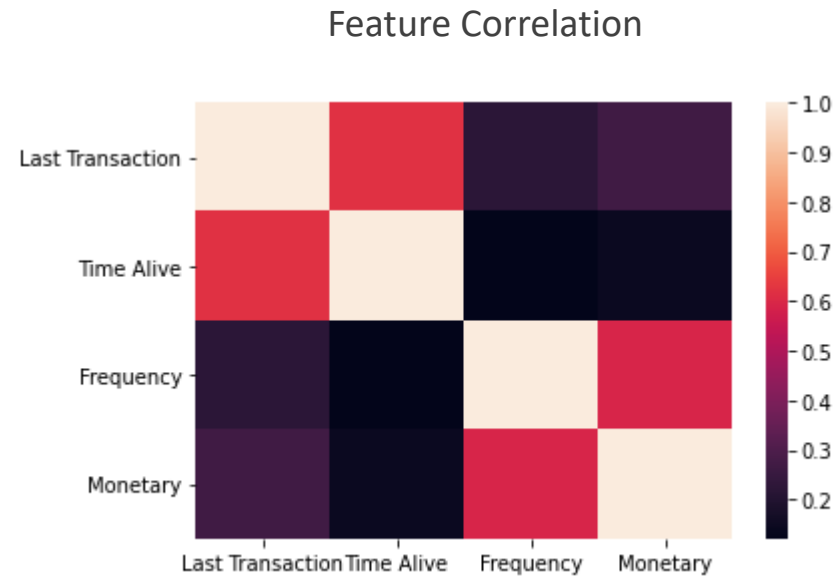


- Labels: Churner(**1**)/Non-Churner(**0**) definition (11,630/14,351 remaining in model)
- Features (LT, TA, F, M): using only historical data from calibration period

# Data Preprocessing

- 45% of labels are churners - data set is balanced.
- Similar RFM statistics to original data
- 80%-20% train-test random split and standardization (mean=0 and stdev=1).

	Last Transaction	Time Alive	Frequency	Monetary	Churn
merchant					
0002b63b92	0	475	1	33.79	1
00057d4302	66	461	28	295.21	1
000bcff341	0	390	1	78.26	1
000ddb0ca	0	458	1	102.99	1
000ed1585f	549	562	59	15754.72	0
...	...	...	...	...	...
ffd3e45675	23	607	5	726.26	1
ffe1f6b51a	260	456	53	2816.16	1
ffe26b900d	254	255	53	7164.12	0
ffec05edb9	20	221	3	159.34	1
fff1754102	382	391	41	4988.59	0





# Modelling

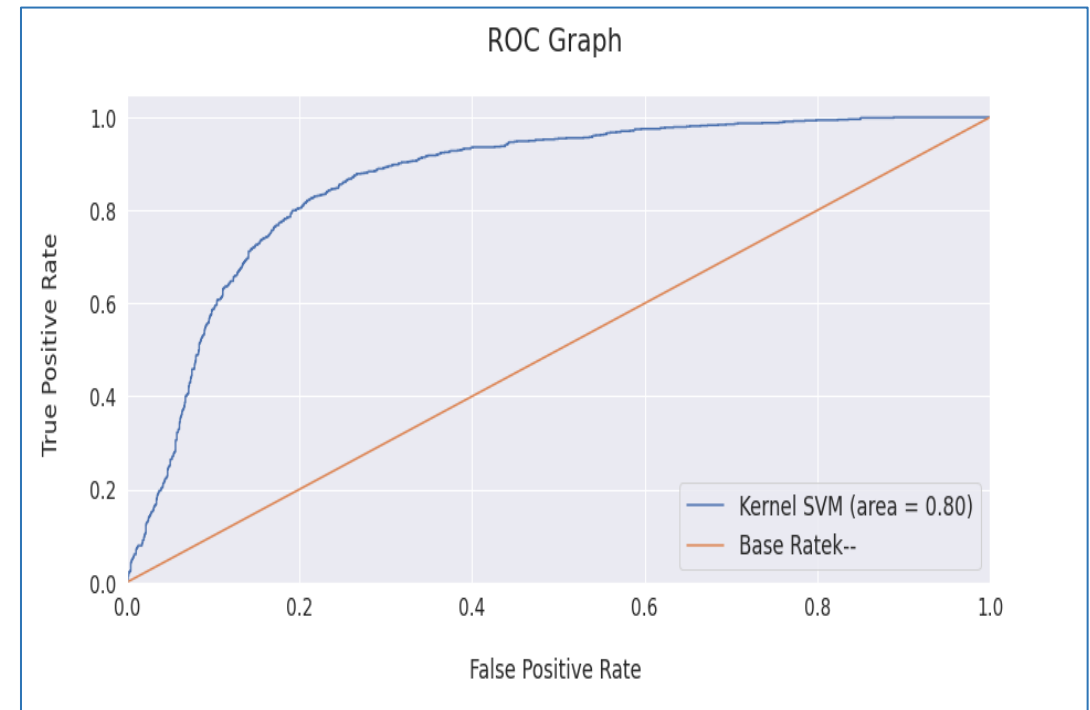
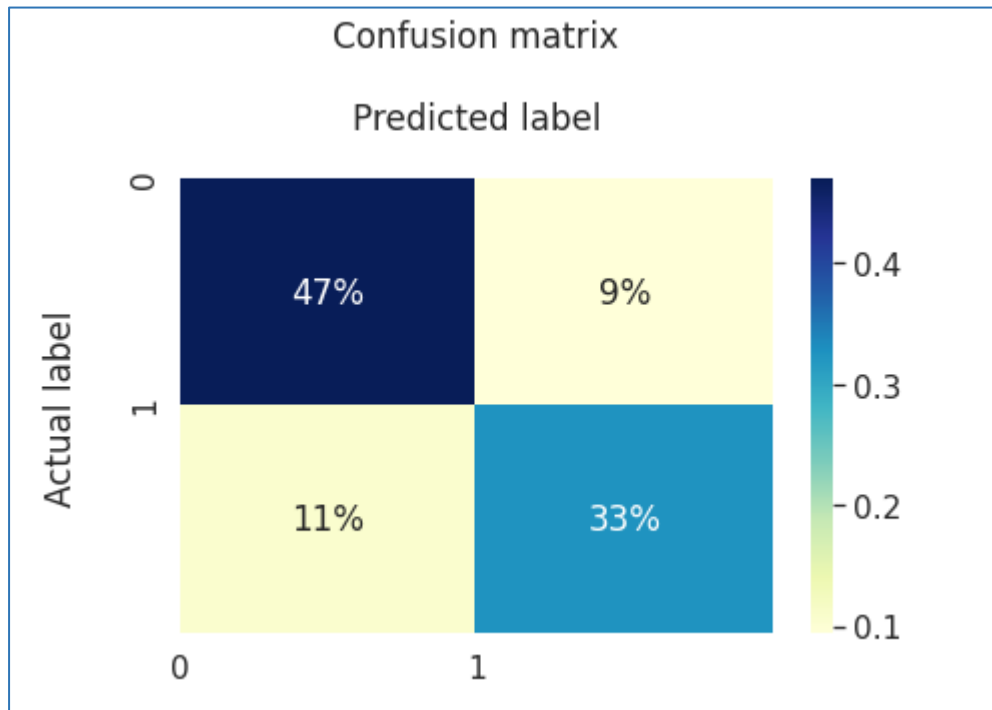
- FN cost is higher than cost of FP => Precision, **Recall** and F1 score
- Grid search over SVM parameters to fine tune model

	Algorithm	ROC	AUC	Mean	ROC	AUC	STD	Accuracy	Mean	Accuracy	STD
1	Kernel SVM			87.36			0.87		80.41		1.09
0	Logistic Regression			86.96			1.06		79.42		1.33
2	Random Forest			86.50			1.12		79.02		1.08
3	Gradient Boosting			86.39			1.02		78.60		1.20

	Model	Accuracy	Precision	Recall	F1 Score	F2 Score
0	Logistic Regression	0.79	0.79	0.70	0.74	0.72
1	Kernel SVM	0.80	0.78	0.75	0.76	0.76
2	Random Forest	0.79	0.76	0.76	0.76	0.76
3	Gradient Boosting	0.78	0.79	0.68	0.73	0.70

# Results

- 10-Fold cross validation accuracy:  $0.8 \pm 0.02$
- Recall: 0.75



# Future Work

- Consider removing merchants from cluster 0 from model as they churned very fast and may not be indicative of an active merchant in the calibration period
- Tweak the calibration and prediction periods to get better classifiers
- Fit a neural network classifier
- Analyze the results with a business/marketing team that can have real world insights
- Develop a survival-based model and compare to the supervised model
- Develop an online learning model that learns after each transaction is logged