


```
[64]: #https://www.kaggle.com/code/rageshcv/customer-loyalty-based-on-rfm-analysis/notebook

def rfm_feature(df, quantiles):
    """This function performs the RFM featurization on dataset by generating the RFM score
    and RFM index. It takes dataset as dataframe, and quantile values for scoring as list and
    returns the RFM features as dataframe."""
    agg_dict = {
        'card_id'      : ['count'],
        'purchase_date' : ['max'],
        'purchase_amount' : ['sum']
    }

    rfm_feature = agg_featurisation(historical_transactions, groupby, agg_dict)
    rfm_feature['recency'] = (datetime.date(2018, 3, 1) - rfm_feature['purchase_date_max']).dt.days
    rfm_feature.rename(columns = {'card_id count' : 'frequency', 'purchase_amount_sum' : 'monetary_value'}, inplace=True)
    rfm_feature = rfm_feature.drop(columns = ['purchase_date_max'])

    rfm_quantiles = rfm_feature.quantile(q = quantiles)

    rfm_feature['R_score'] = rfm_feature['frequency'].apply(RFM_Score, args = ('recency', rfm_quantiles))
    rfm_feature['F_score'] = rfm_feature['frequency'].apply(RFM_Score, args = ('frequency', rfm_quantiles))
    rfm_feature['M_score'] = rfm_feature['monetary_value'].apply(RFM_Score, args = ('monetary_value', rfm_quantiles))
    rfm_feature['RFM_Score'] = rfm_feature['R_score'] + rfm_feature['F_score'] + rfm_feature['M_score']
    rfm_feature['RFM_index'] = rfm_feature['RFM_Score'].map(atoi) + rfm_feature['F_score'].map(atoi) + rfm_feature['M_score'].map(atoi)
    rfm_feature['RFM_index'] = rfm_feature['RFM_index'].astype(int)

    rfm_feature = rfm_feature.drop(columns = ['recency', 'frequency', 'monetary_value'])

    return rfm_feature

In [65]: quantiles = [0.012, 0.02, 0.05, 0.2, 0.5, 0.8, 0.96, 0.992, 1.0]
hist_rfm_feature = rfm_feature(historical_transactions, quantiles)

In [66]: hist_rfm_feature = reduce_mem_usage(hist_rfm_feature)
hist_rfm_feature.head()
```

Mem. usage decreased to 13.65 Mb (40.6% reduction)

	card_id	R_score	F_score	M_score	RFM_Score	RFM_index
0	C_ID_0046db98a	9	7	7	23	977
1	C_ID_011b0d9794	4	5	6	15	456
2	C_ID_01904d743d	4	6	5	15	465
3	C_ID_01b098f01	4	7	6	17	476
4	C_ID_0382b662f4	3	7	6	16	376

Merging all transactions features.

```
In [67]: all_transaction_features = pd.merge(hist_transaction_features, new_transactions_features, on = 'card_id', how = 'left')
all_transaction_features = pd.merge(all_transaction_features, hist_category_features, on = 'card_id', how = 'left')
all_transaction_features = pd.merge(all_transaction_features, new_category_features, on = 'card_id', how = 'left')
all_transaction_features = pd.merge(all_transaction_features, hist_month_lag_features, on = 'card_id', how = 'left')
all_transaction_features = pd.merge(all_transaction_features, new_month_lag_features, on = 'card_id', how = 'left')
all_transaction_features = pd.merge(all_transaction_features, hist_installments_features, on = 'card_id', how = 'left')
all_transaction_features = pd.merge(all_transaction_features, new_installments_features, on = 'card_id', how = 'left')
all_transaction_features = pd.merge(all_transaction_features, hist_rfm_feature, on = 'card_id', how = 'left')

In [68]: all_transaction_features = reduce_mem_usage(all_transaction_features)
all_transaction_features.head()
```

Mem. usage decreased to 301.46 Mb (19.4% reduction)

	card_id	R_score	F_score	M_score	RFM_Score	RFM_index	hist_transc_count	hist_category_1.0_sum	hist_category_1.0_mean	hist_category_1.0_std
0	C_ID_0046db98a			150		142.0	0.946777	147.0		0.979980
1	C_ID_011b0d9794			27		27.0	1.000000	0.0		0.000000
2	C_ID_01904d743d			71		50.0	0.704102	71.0		1.000000
3	C_ID_01b098f01			182		128.0	0.703125	156.0		0.856934
4	C_ID_0382b662f4			170		162.0	0.953125	163.0		0.958984

5 rows x 341 columns

```
In [69]: all_transaction_features.to_csv('data/all_transaction_features.csv')
```

```
In [70]: del hist_transaction_features
del new_transactions_features
del hist_category_features
del new_category_features
del hist_month_lag_features
del new_month_lag_features
del hist_installments_features
del new_installments_features
del hist_rfm_feature
```

Merging transactions features to train and test set on card_id.

Finally we will merge all features engineered from transactions to train and test set based on card ids and features having null values are imputed with 0.

```
In [71]: train = pd.merge(train, all_transaction_features, on = 'card_id', how = 'left')
test = pd.merge(test, all_transaction_features, on = 'card_id', how = 'left')
```

```
In [72]: train.fillna(value = 0, inplace = True)
train.head()
```

Out[72]:

	first_active_month	card_id	feature_1	feature_2	feature_3	target	elapsed_time	first_active_year	outlier	hist_transc_count	...
0	6	C_ID_92a2005557	5	2	1	-0.820312	245	2017	0	260	...
1	1	C_ID_3d004492b4f	4	1	0	0.392822	396	2017	0	350	...
2	8	C_ID_0d399d96d	2	2	0	0.687988	549	2016	0	43	...
3	9	C_ID_186c6a6901	4	3	0	0.142456	153	2017	0	77	...
4	11	C_ID_cdb2c2db2	1	3	0	-0.159790	92	2017	0	133	...

5 rows x 349 columns

```
In [73]: test.fillna(value = 0, inplace = True)
test.head()
```

Out[73]:

	first_active_month	card_id	feature_1	feature_2	feature_3	elapsed_time	first_active_year	hist_transc_count	hist_authozied_flag_sum
0	4	C_ID_0ab67a2a2ab	3	3	1	306	2017	68	44
1	1	C_ID_130f60cbdd	2	3	0	396	2017	78	77
2	8	C_ID_b709037bc5	5	1	1	184	2017	13	9
3	12	C_ID_d27d835a9f	2	1	0	62	2017	26	26
4	12	C_ID_2b5e3d5f52	5	1	1	793	2015	110	87

5 rows x 347 columns

```
In [74]: train.to_csv('data/featurized_train.csv', index = False)
test.to_csv('data/featurized_test.csv', index = False)
```

Summary

1. Null values in merchant id columns of historical and new trasactions were imputed with mode while null values in category 2 and category 3 columns were imputed by training classifiers.
2. The purchase amount in historical and new transactions was de-anonymized and the categorical columns were one hot encoded.
3. Elapsed time feature for each card id was engineered from first active month and first active month was segregated into year and month.
4. Historical and New transactions were aggregated on card id and different features were engineered from transactions columns.
5. Historical and New transactions were aggregated on card id with different category 1, category 2 and category 3 values, and different features were engineered from purchase amount columns.
6. Historical and New transactions were aggregated on card id with different month lag values, and different features were engineered from purchase amount columns.
7. Historical and New transactions were aggregated successively aggregated on card id with installments columns, and different features were engineered from purchase amount and authorized flag columns.
8. RFM analysis was used to engineer RFM score and RFM index features for each card id.
9. All the features were merged with train and test set on card ids and the null values were filled with 0.