





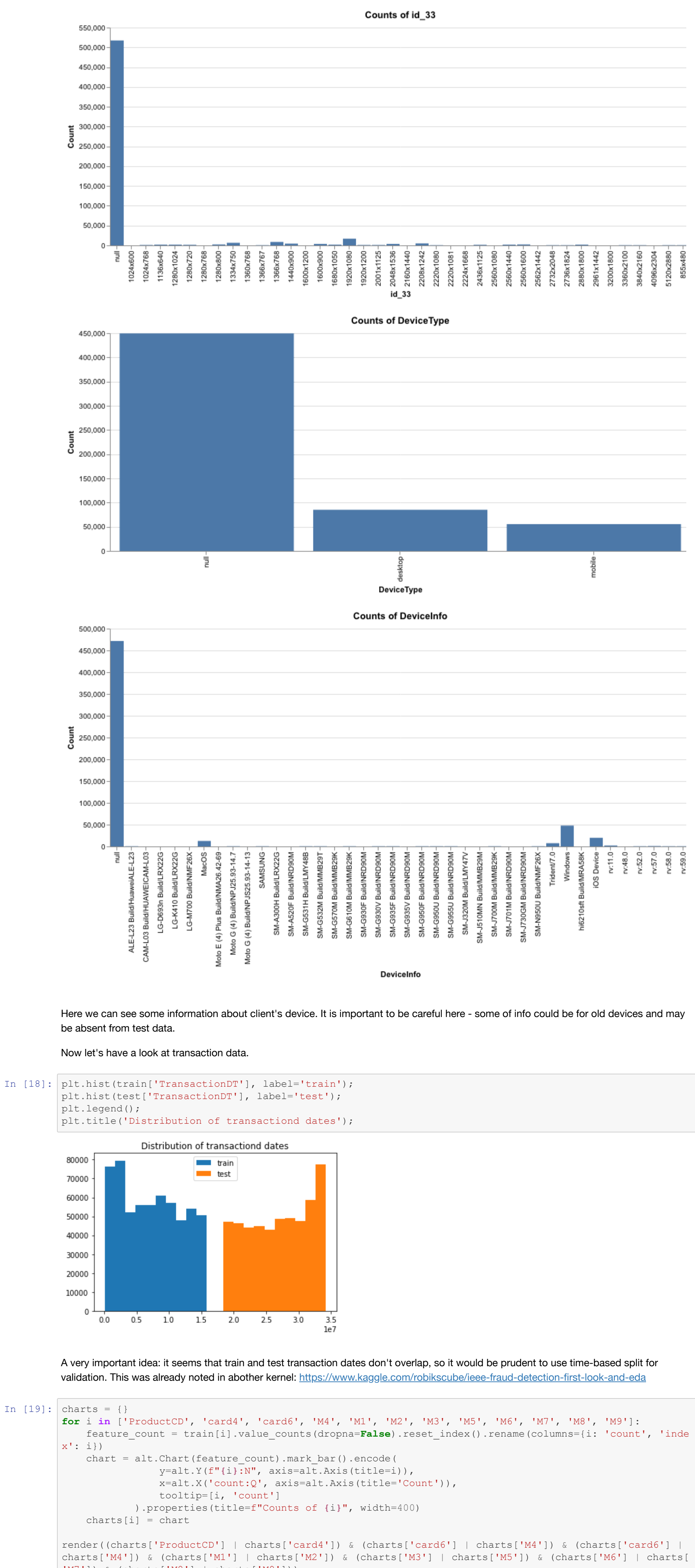


We have several features showing some kind of "found" status and several binary columns.

In [17]:

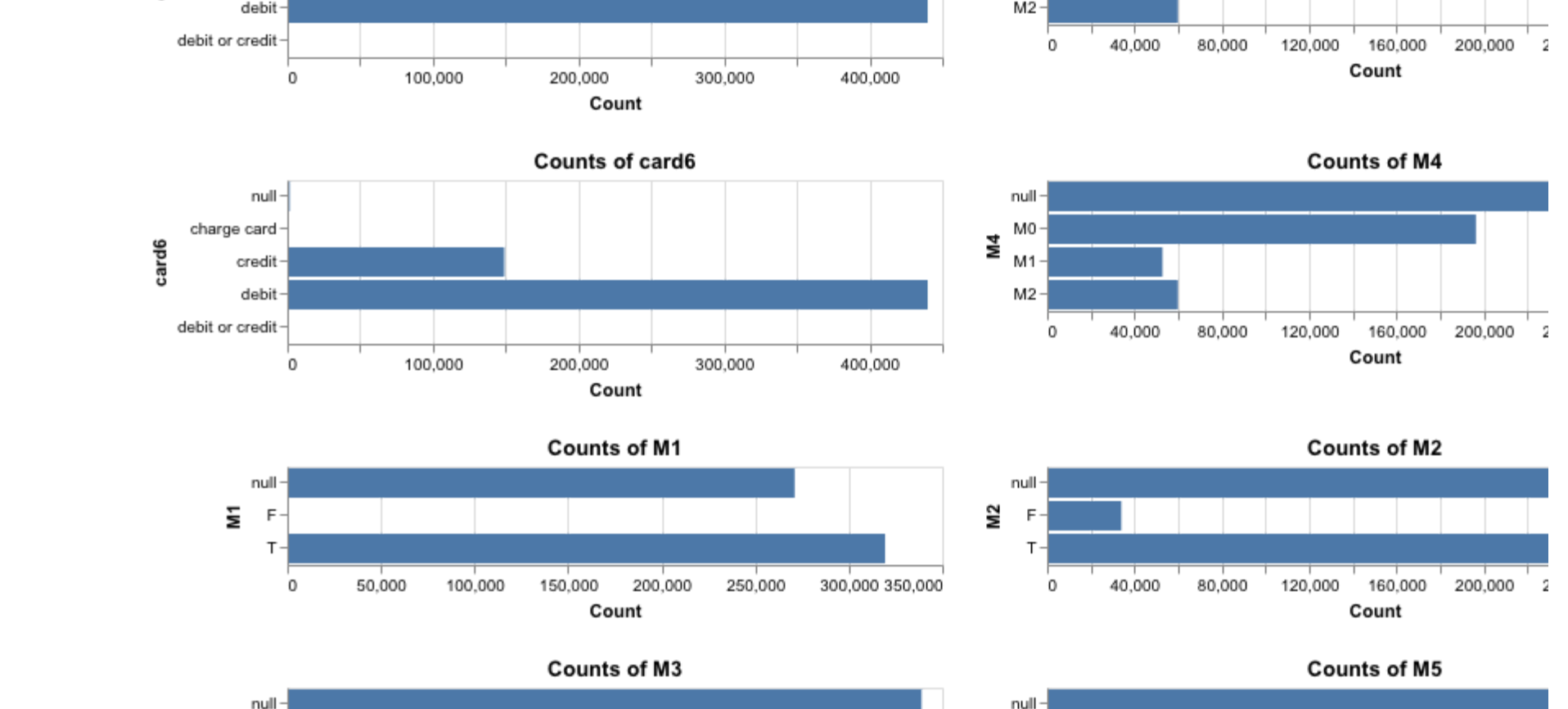
```
charts = {}
for i in ['id_30', 'id_31', 'id_33', 'DeviceType', 'DeviceInfo']:
    feature_count = train[i].value_counts(dropna=False).reset_index().rename(columns={i: 'count', 'index': i})
    chart = alt.Chart(feature_count).mark_bar().encode(
        y=alt.Y(f"{i}"), axis=alt.Axis(title=i),
        x=alt.X('count', title='Count'),
        tooltip=[i, 'count']
    ).properties(title=f'Counts of {i}', width=800)
    charts[i] = chart

render(charts['id_30'] | charts['id_31'] | charts['id_33'] | charts['DeviceType'] | charts['DeviceInfo'])
```

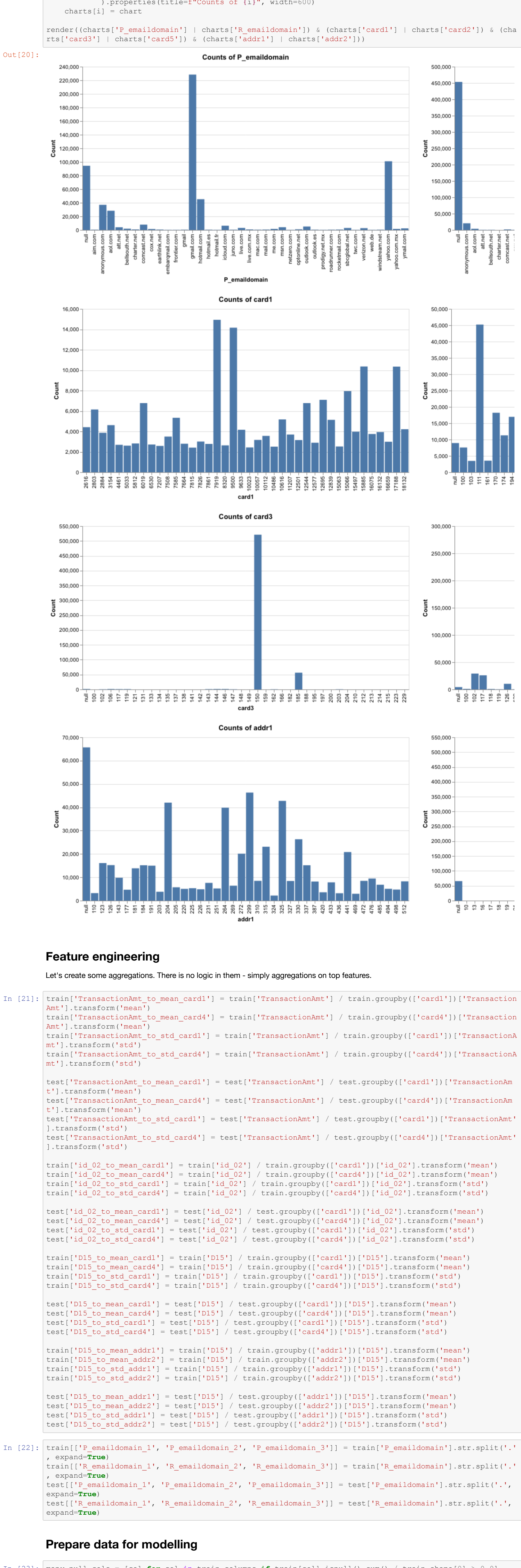
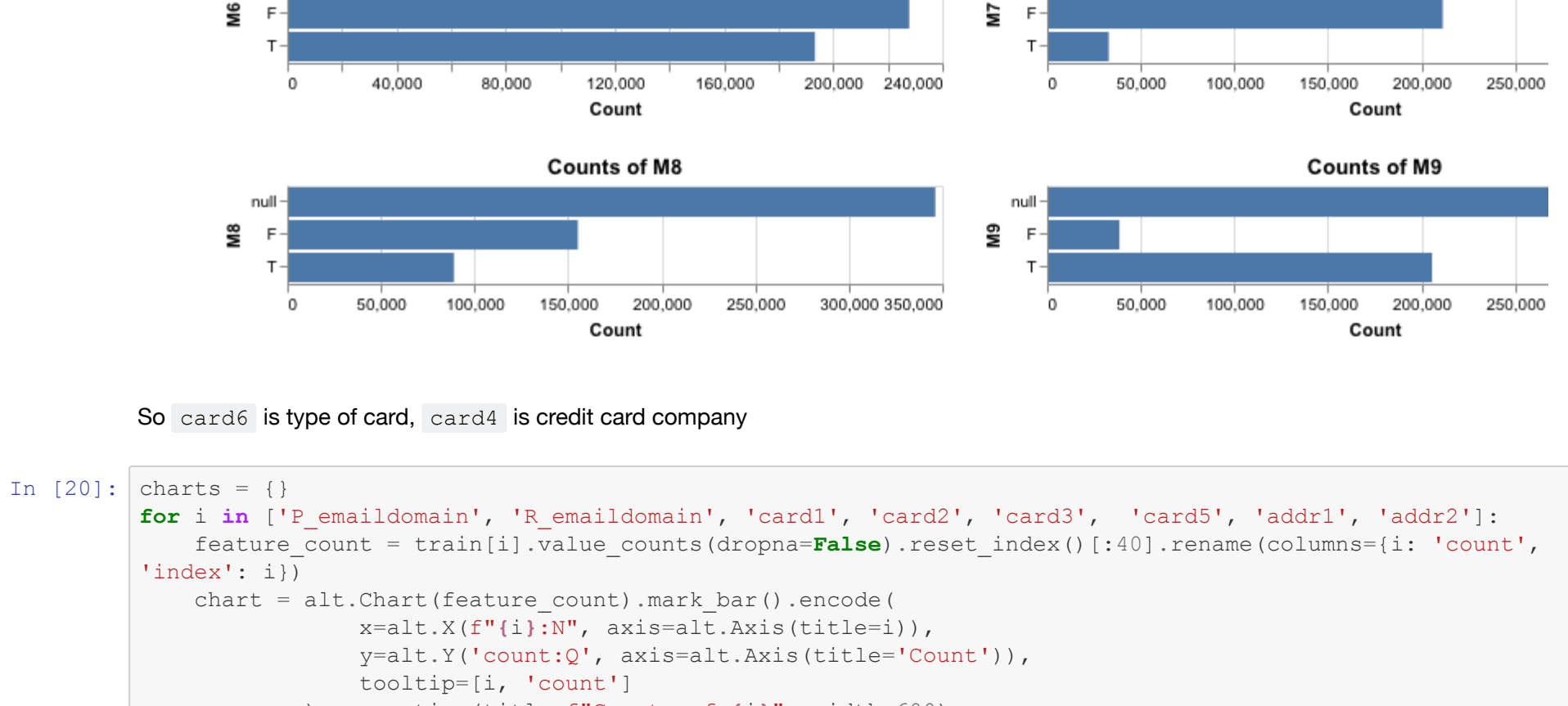


Here we can see some information about client's device. It is important to be careful here - some of info could be for old devices and may be absent from test data.

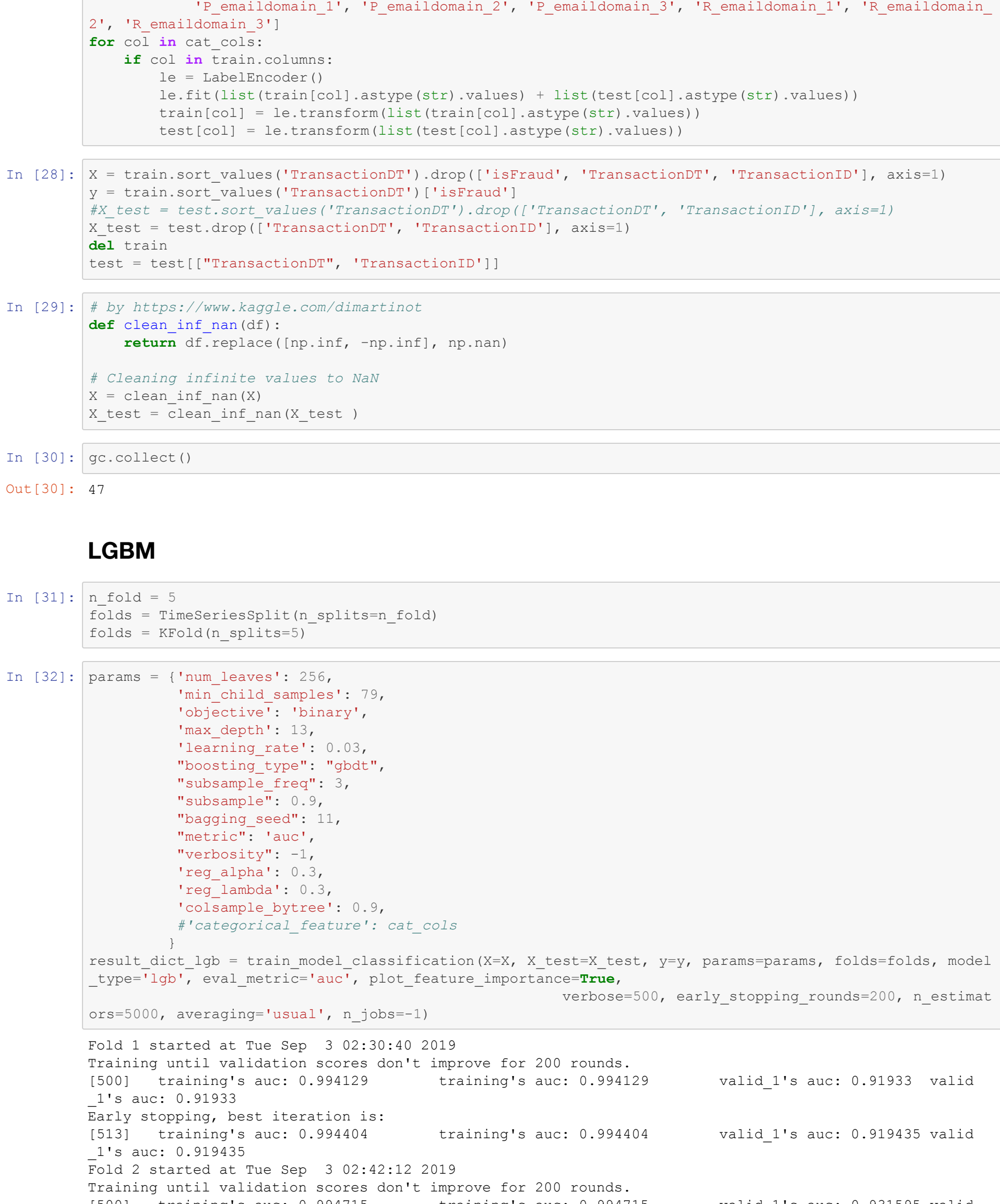
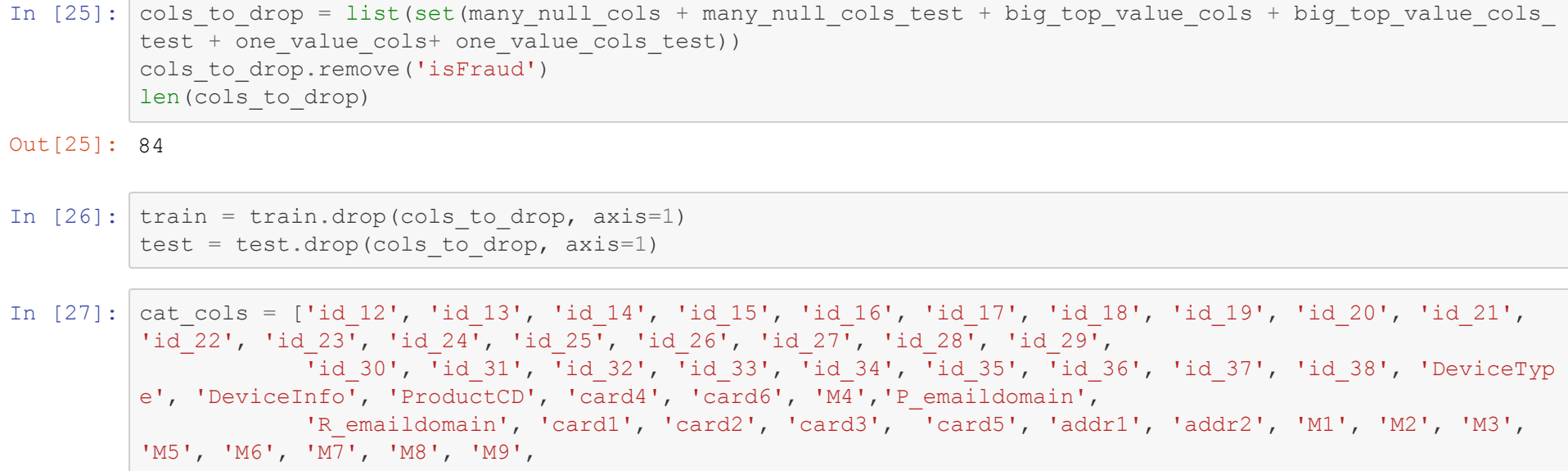
Now let's have a look at transaction data.



A very important idea: it seems that train and test transaction dates don't overlap, so it would be prudent to use time-based split for validation. This was already noted in another kernel: <https://www.kaggle.com/robikublee/ieee-fraud-detection-first-look-and-sql>



So card6 is type of card, card4 is credit card company



## Feature engineering

Let's create some aggregations. There is no logic in them - simply aggregations on top features.

In [21]:

```
train['TransactionAmt_to_mean_card1'] = train['TransactionAmt'] / train.groupby(['card1'])['TransactionAmt'].transform('mean')
train['TransactionAmt_to_std_card1'] = train['TransactionAmt'] / train.groupby(['card1'])['TransactionAmt'].transform('std')
train['TransactionAmt_to_std_card4'] = train['TransactionAmt'] / train.groupby(['card4'])['TransactionAmt'].transform('std')
train['TransactionAmt_to_std_card4'] = train['TransactionAmt'] / train.groupby(['card4'])['TransactionAmt'].transform('std')

test['TransactionAmt_to_mean_card1'] = test['TransactionAmt'] / test.groupby(['card1'])['TransactionAmt'].transform('mean')
test['TransactionAmt_to_std_card1'] = test['TransactionAmt'] / test.groupby(['card1'])['TransactionAmt'].transform('std')
test['TransactionAmt_to_std_card4'] = test['TransactionAmt'] / test.groupby(['card4'])['TransactionAmt'].transform('std')
test['TransactionAmt_to_std_card4'] = test['TransactionAmt'] / test.groupby(['card4'])['TransactionAmt'].transform('std')

train['id_02_to_mean_card1'] = train['id_02'] / train.groupby(['card1'])['id_02'].transform('mean')
train['id_02_to_std_card1'] = train['id_02'] / train.groupby(['card1'])['id_02'].transform('std')
train['id_02_to_std_card4'] = train['id_02'] / train.groupby(['card4'])['id_02'].transform('std')
train['id_02_to_std_card4'] = train['id_02'] / train.groupby(['card4'])['id_02'].transform('std')

test['id_02_to_mean_card1'] = test['id_02'] / test.groupby(['card1'])['id_02'].transform('mean')
test['id_02_to_std_card1'] = test['id_02'] / test.groupby(['card1'])['id_02'].transform('std')
test['id_02_to_std_card4'] = test['id_02'] / test.groupby(['card4'])['id_02'].transform('std')
test['id_02_to_std_card4'] = test['id_02'] / test.groupby(['card4'])['id_02'].transform('std')

train['D15_to_mean_card1'] = train['D15'] / train.groupby(['card1'])['D15'].transform('mean')
train['D15_to_std_card1'] = train['D15'] / train.groupby(['card1'])['D15'].transform('std')
train['D15_to_std_card4'] = train['D15'] / train.groupby(['card4'])['D15'].transform('std')
train['D15_to_std_card4'] = train['D15'] / train.groupby(['card4'])['D15'].transform('std')

test['D15_to_mean_card1'] = test['D15'] / test.groupby(['card1'])['D15'].transform('mean')
test['D15_to_std_card1'] = test['D15'] / test.groupby(['card1'])['D15'].transform('std')
test['D15_to_std_card4'] = test['D15'] / test.groupby(['card4'])['D15'].transform('std')
test['D15_to_std_card4'] = test['D15'] / test.groupby(['card4'])['D15'].transform('std')

train['D15_to_mean_addr1'] = train['D15'] / train.groupby(['addr1'])['D15'].transform('mean')
train['D15_to_std_addr1'] = train['D15'] / train.groupby(['addr1'])['D15'].transform('std')
train['D15_to_std_addr2'] = train['D15'] / train.groupby(['addr2'])['D15'].transform('std')
train['D15_to_std_addr2'] = train['D15'] / train.groupby(['addr2'])['D15'].transform('std')

test['D15_to_mean_addr1'] = test['D15'] / test.groupby(['addr1'])['D15'].transform('mean')
test['D15_to_std_addr1'] = test['D15'] / test.groupby(['addr1'])['D15'].transform('std')
test['D15_to_std_addr2'] = test['D15'] / test.groupby(['addr2'])['D15'].transform('std')
test['D15_to_std_addr2'] = test['D15'] / test.groupby(['addr2'])['D15'].transform('std')
```

In [22]:

```
train['P_emaildomain_1', 'P_emaildomain_2', 'P_emaildomain_3'] = train['P_emaildomain'].str.split('.', expand=True)
train['R_emaildomain_1', 'R_emaildomain_2', 'R_emaildomain_3'] = train['R_emaildomain'].str.split('.', expand=True)
test['P_emaildomain_1', 'P_emaildomain_2', 'P_emaildomain_3'] = test['P_emaildomain'].str.split('.', expand=True)
test['R_emaildomain_1', 'R_emaildomain_2', 'R_emaildomain_3'] = test['R_emaildomain'].str.split('.', expand=True)
```

## Prepare data for modelling

In [23]:

```
many_null_cols = [col for col in train.columns if train[col].isnull().sum() / train.shape[0] > 0.9]
many_null_cols_test = [col for col in test.columns if test[col].isnull().sum() / test.shape[0] > 0.9]
```

In [24]:

```
big_top_value_cols = [col for col in train.columns if train[col].value_counts(dropna=False, sort=True).values[0] > 0.9]
big_top_value_cols_test = [col for col in test.columns if test[col].value_counts(dropna=False, sort=True).values[0] > 0.9]
```

In [25]:

```
cols_to_drop = list(set(many_null_cols + many_null_cols_test + big_top_value_cols + big_top_value_cols_test - one_value_cols - one_value_cols_test))
cols_to_drop.remove('isFraud')
```

Out [25]:

```
84
```

In [26]:

```
train = train.drop(cols_to_drop, axis=1)
test = test.drop(cols_to_drop, axis=1)
```

In [27]:

```
cat_cols = ['id_12', 'id_13', 'id_14', 'id_15', 'id_16', 'id_17', 'id_18', 'id_19', 'id_20', 'id_21', 'id_22', 'id_23', 'id_24', 'id_25', 'id_26', 'id_27', 'id_28', 'id_29', 'id_30', 'id_31', 'id_32', 'id_33', 'id_34', 'id_35', 'id_36', 'id_37', 'id_38', 'DeviceInfo', 'DeviceType', 'ProductCD', 'card4', 'card6', 'M4', 'P_emaildomain', 'R_emaildomain', 'card1', 'card2', 'card3', 'card5', 'addr1', 'addr2', 'M1', 'M2', 'M3', 'M5', 'M6', 'M7', 'M8', 'M9', 'P_emaildomain_1', 'P_emaildomain_2', 'P_emaildomain_3', 'R_emaildomain_1', 'R_emaildomain_2', 'R_emaildomain_3']
for col in cat_cols:
    le = LabelEncoder()
    train[col] = le.fit_transform(train[col].astype(str).values)
    test[col] = le.transform(test[col].astype(str).values)
```

In [28]:

```
X = train.sort_values('TransactionDT').drop(['isFraud', 'TransactionDT', 'TransactionID'], axis=1)
y = train['isFraud']
X_test = test.sort_values('TransactionDT').drop(['TransactionDT', 'TransactionID'], axis=1)
X_test = test.drop(['TransactionDT', 'TransactionID'], axis=1)
del train
del test
```

In [29]:

```
# by https://www.kaggle.com/dmarinot
def clean_inf_nan(df):
    return df.replace([np.inf, -np.inf], np.nan)

# Cleaning infinite values to NaN
X = clean_inf_nan(X)
X_test = clean_inf_nan(X_test)
```

In [30]:

```
gc.collect()
```

Out [30]:

```
47
```

## LGBM

In [31]:

```
n_fold = 5
folds = TimeSeriesSplit(n_splits=n_fold)
folds = KFold(n_splits=5)
```

In [32]:

```
params = {'num_leaves': 256,
          'min_child_samples': 79,
          'max_depth': 5,
          'objective': 'binary:logistic',
          'subsample': 0.85,
          'boosting_type': 'gbdt',
          'subsample_freq': 3,
          'subsample': 0.5,
          'bagging_seed': 11,
          'metric': 'auc',
          'verbosity': -1,
          'reg_alpha': 0.3,
          'reg_lambda': 0.3,
          'colsample_bytree': 0.9,
          'categorical_feature': cat_cols}

result_dict_lgb = train_model_classification(X=X, X_test=X_test, y=y, params=params, folds=folds, model_type='lgb', eval_metric='auc', plot_feature_importance=True,
                                             n_jobs=500, early_stopping_rounds=200, n_estimators=5000, averaging='rank')
```

Fold 1 started at Tue Sep 3 02:30:40 2019  
Training until validation scores don't improve for 200 rounds.  
[500] training's auc: 0.994129 training's auc: 0.994129 valid\_1's auc: 0.91933 valid\_1's auc: 0.91933  
Early stopping, best iteration is:  
[513] training's auc: 0.994402 training's auc: 0.994402 valid\_1's auc: 0.919435 valid\_1's auc: 0.919435  
Fold 2 started at Tue Sep 3 02:42:12 2019  
Training until validation scores don't improve for 200 rounds.  
[500] training's auc: 0.994715 training's auc: 0.994715 valid\_1's auc: 0.931595 valid\_1's auc: 0.931595  
Early stopping, best iteration is:  
[546] training's auc: 0.995664 training's auc: 0.995664 valid\_1's auc: 0.931982 valid\_1's auc: 0.931982  
Fold 3 started at Tue Sep 3 02:54:10 2019  
Training until validation scores don't improve for 200 rounds.  
[500] training's auc: 0.994741 training's auc: 0.994741 valid\_1's auc: 0.926174 valid\_1's auc: 0.926174  
Early stopping, best iteration is:  
[372] training's auc: 0.990799 training's auc: 0.990799 valid\_1's auc: 0.926607 valid\_1's auc: 0.926607  
Fold 4 started at Tue Sep 3 02:03:28 2019  
Training until validation scores don't improve for 200 rounds.  
[500] training's auc: 0.994673 training's auc: 0.994673 valid\_1's auc: 0.946452 valid\_1's auc: 0.946452  
Early stopping, best iteration is:  
[464] training's auc: 0.993779 training's auc: 0.993779 valid\_1's auc: 0.946655 valid\_1's auc: 0.946655  
Fold 5 started at Tue Sep 3 03:13:55 2019  
Training until validation scores don't improve for 200 rounds.  
[500] training's auc: 0.995027 training's auc: 0.995027 valid\_1's auc: 0.922803 valid\_1's auc: 0.922803  
Early stopping, best iteration is:  
[405] training's auc: 0.992298 training's auc: 0.992298 valid\_1's auc: 0.923309 valid\_1's auc: 0.923309  
CV mean score: 0.9296, std: 0.0095.



In [33]:

```
sub['isFraud'] = result_dict_lgb['prediction']
sub.to_csv('submission_csv', index=False)
```

Out [34]:

```
sub.head()
```

Out [34]:

```
TransactionID isFraud
0 3663549 0.001571974543870
1 3663550 0.00182337471988
2 3663551 0.00348749177351
3 3663552 0.001593881699954
4 3663553 0.002441298688730
```

In [35]:

```
pd.DataFrame(result_dict_lgb[['oof']], to_csv('lgb_oof.csv', index=False)
```

## Blending

In [36]:

```
# xgb_params = {'eta': 0.04,
#               'max_depth': 5,
#               'subsample': 0.85,
#               'objective': 'binary:logistic',
#               'eval_metric': 'auc',
#               'silent': True,
#               'nthread': -1,
#               'tree_method': 'gpu_hist'})
# result_dict_xgb = train_model_classification(X=X, X_test=X_test, y=y, params=xgb_params, folds=folds, model_type='xgb', eval_metric='auc', plot_feature_importance=False,
#                                             n_jobs=500, early_stopping_rounds=200, n_estimators=5000, averaging='rank')
```

In [37]:

```
# test = test.sort_values('TransactionDT')
# result_dict_xgb = result_dict_lgb['prediction']
# sub.to_csv('blend_csv', index=False)
```

In [38]:

```
# test = test.sort_values('TransactionID')
# result_dict_xgb = result_dict_lgb['prediction'] + result_dict_xgb['prediction']
# sub['isFraud'] = pd.merge(sub, test, on='TransactionID')['prediction']
# sub.to_csv('blend_csv', index=False)
```