


```
[35]: columns_to_unskew = []
columns_boxcox = [ ]
for col in df_copy.columns:
    skew = df_copy[col].skew()
    print("skewness for ",col," is ", skew)
    # if skew > 1 or skew < -1:
    #     if col == 'BALANCE_FREQUENCY' or col == 'TENURE':
    #         columns_boxcox.append(col)
    #     else:
    #         columns_to_unskew.append(col)
```

```
skewness for BALANCE is 2.393386842571886
skewness for BALANCE_FREQUENCY is 2.8523255185149678
skewness for PURCHASES is 8.144269864564051
skewness for ONEOFF_PURCHASES is 10.845982884786870
skewness for INSTALLMENTS_PURCHASES is 7.299119988745641
skewness for CASH_ADVANCE is 5.166689874897423
skewness for PURCHASES_FREQUENCY is 0.86016423558983591
skewness for ONEOFF_PURCHASES_FREQUENCY is 1.5356127835248519
skewness for PURCHASES_INSTALLMENTS_FREQUENCY is 0.589201649999882
skewness for CASH_ADVANCE_TRX is 1.52386856477852
skewness for CASH_ADVANCE_TRX is 5.721298293192298
skewness for PURCHASES_TRX is 4.638955269323469
skewness for CREDIT_LIMIT is 1.5225359591884323
skewness for PAYMENTS is 5.967619794397562
skewness for MINIMUM_PAYMENTS is 12.852446498665346
skewness for PRC_FULL_PAYMENT is 1.942619948971858
skewness for TENURE is -2.943817287199134
```

alot of the columns are skewed

```
In [36]: print(columns_to_unskew)

['BALANCE', 'BALANCE_FREQUENCY', 'PURCHASES', 'ONEOFF_PURCHASES', 'INSTALLMENTS_PURCHASES', 'CASH_ADVANCE', 'PAYMENTS', 'MINIMUM_PAYMENTS', 'PRC_FULL_PAYMENT', 'TENURE']
```

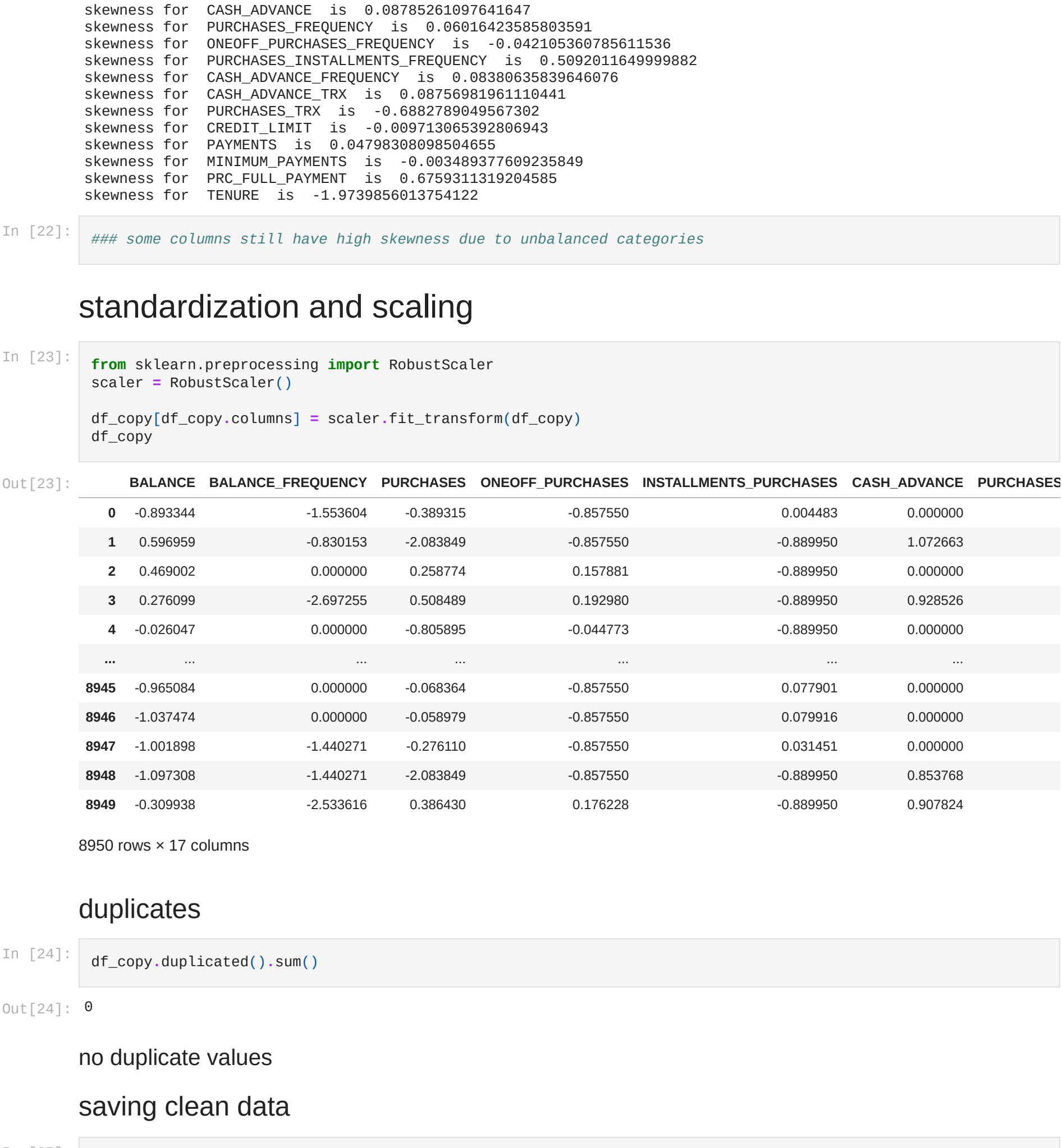
```
In [37]: from scipy import stats
```

```
In [38]: df_copy[columns_to_unskew] = np.log(df_copy[columns_to_unskew]+0.1)
df_copy[columns_to_unskew] = np.sort(df_copy[columns_to_unskew])
df_copy[columns_to_unskew] = np.log10(df_copy[columns_to_unskew]+0.1)
df_copy[columns_to_unskew] = np.arcsinh(df_copy[columns_to_unskew])
for col in columns_to_unskew:
    df_copy[col] = stats.boxcox(df_copy[col]*1e-6) # 1e-6 to deal with zero values
```

```
In [39]: df_copy.head()

Out[39]:
```

| | BALANCE | BALANCE_FREQUENCY | PURCHASES | ONEOFF_PURCHASES | INSTALLMENTS_PURCHASES | CASH_ADVANCE | PURCHASES_FREQUENCY |
|---|-----------|-------------------|-----------|------------------|------------------------|--------------|---------------------|
| 0 | 5.677545 | -0.159521 | 6.468611 | -13.227810 | | 4.912375 | -15.117524 |
| 1 | 21.718338 | -0.085202 | -1.549881 | -13.227810 | | -11.134216 | 8.294623 |
| 2 | 20.341074 | 0.000001 | 11.281173 | 6.792680 | | -11.134216 | -15.117524 |
| 3 | 18.264780 | -0.277124 | 13.048225 | 7.484694 | | -11.134216 | 5.148653 |
| 4 | 15.012643 | 0.000001 | 3.415678 | 2.797098 | | -11.134216 | -15.117524 |



```
In [21]: for col in df_copy.columns:
    skew = df_copy[col].skew()
    print("skewness for ",col," is ", skew)

skewness for BALANCE is -0.1571654364426791
skewness for BALANCE_FREQUENCY is -1.4774423861376662
skewness for PURCHASES is -0.7519311592951326
skewness for ONEOFF_PURCHASES is -0.84152181729153477
skewness for INSTALLMENTS_PURCHASES is -0.20899595968367456
skewness for CASH_ADVANCE is 0.88785261897641647
skewness for PURCHASES_FREQUENCY is 0.86016423558983591
skewness for ONEOFF_PURCHASES_FREQUENCY is -0.842305360785611536
skewness for PURCHASES_INSTALLMENTS_FREQUENCY is 1.5356127835248519
skewness for CASH_ADVANCE_FREQUENCY is 0.86388635839646876
skewness for CASH_ADVANCE_TRX is 0.88756981961118461
skewness for PURCHASES_TRX is -0.6882788949587382
skewness for CREDIT_LIMIT is -0.809713865392896943
skewness for PAYMENTS is 0.8479838898594655
skewness for MINIMUM_PAYMENTS is -0.803489377689235849
skewness for PRC_FULL_PAYMENT is 0.6759311319204589
skewness for TENURE is -1.9739856013754122
```

```
In [22]: ## some columns still have high skewness due to unbalanced categories
sns.boxplot(df_copy[['CASH_ADVANCE']])
```

standardization and scaling

```
In [23]: from sklearn.preprocessing import RobustScaler
scaler = RobustScaler()
```

```
df_copy[df_copy.columns] = scaler.fit_transform(df_copy)
df_copy
```

```
Out[23]:
```

| | BALANCE | BALANCE_FREQUENCY | PURCHASES | ONEOFF_PURCHASES | INSTALLMENTS_PURCHASES | CASH_ADVANCE | PURCHASES_FREQUENCY |
|------|-----------|-------------------|-----------|------------------|------------------------|--------------|---------------------|
| 0 | -0.893344 | -1.553604 | -0.389315 | -0.857550 | | 0.004483 | 0.000000 |
| 1 | 0.596959 | -0.830153 | -2.083849 | -0.857550 | | -0.899950 | 1.072683 |
| 2 | 0.468002 | 0.000000 | 0.258774 | 0.157881 | | -0.899950 | 0.000000 |
| 3 | 0.276099 | -2.697255 | 0.508489 | 0.192980 | | -0.899950 | 0.928526 |
| 4 | -0.26047 | 0.000000 | -0.805995 | -0.044773 | | -0.899950 | 0.000000 |
| ... | ... | ... | ... | ... | ... | ... | ... |
| 8945 | -0.965984 | 0.000000 | -0.068364 | -0.857550 | | 0.077901 | 0.000000 |
| 8946 | -1.037474 | 0.000000 | -0.058979 | -0.857550 | | 0.079916 | 0.000000 |
| 8947 | -1.001898 | -1.440271 | -0.276110 | -0.857550 | | 0.031451 | 0.000000 |
| 8948 | -1.097908 | -1.440271 | -0.283849 | -0.857550 | | -0.899950 | 0.853768 |
| 8949 | -0.309938 | -2.533616 | 0.386430 | 0.176228 | | -0.899950 | 0.907824 |

8950 rows × 7 columns

duplicates

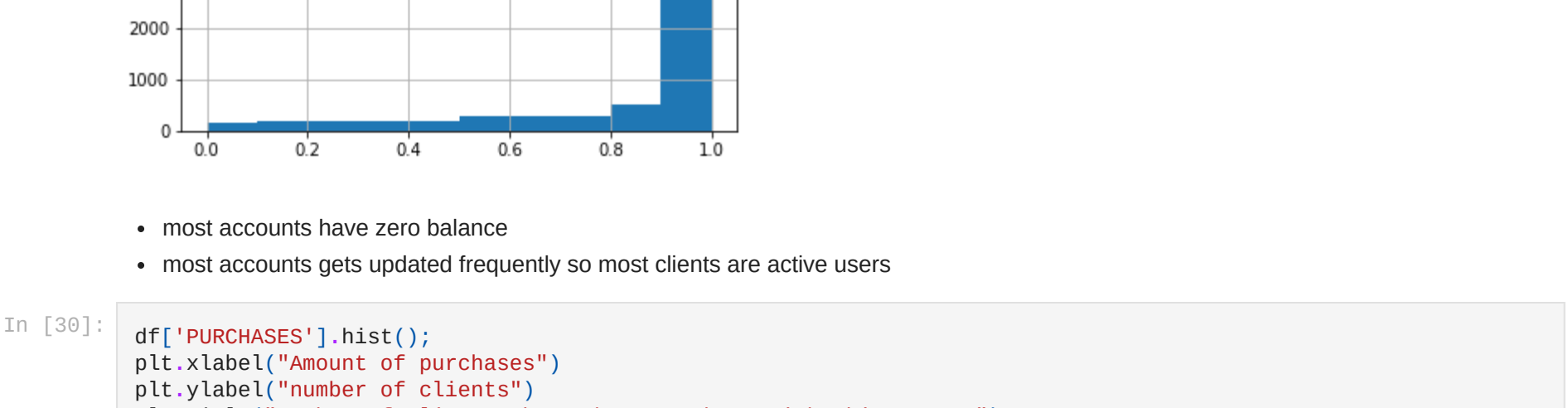
```
In [24]: df_copy.duplicated().sum()
```

```
Out[24]: 0
```

saving clean data

```
In [25]: df_clean = df_copy.copy()
df_clean.to_csv("credit_card_cleanData")
```

Exploratory Data Analysis

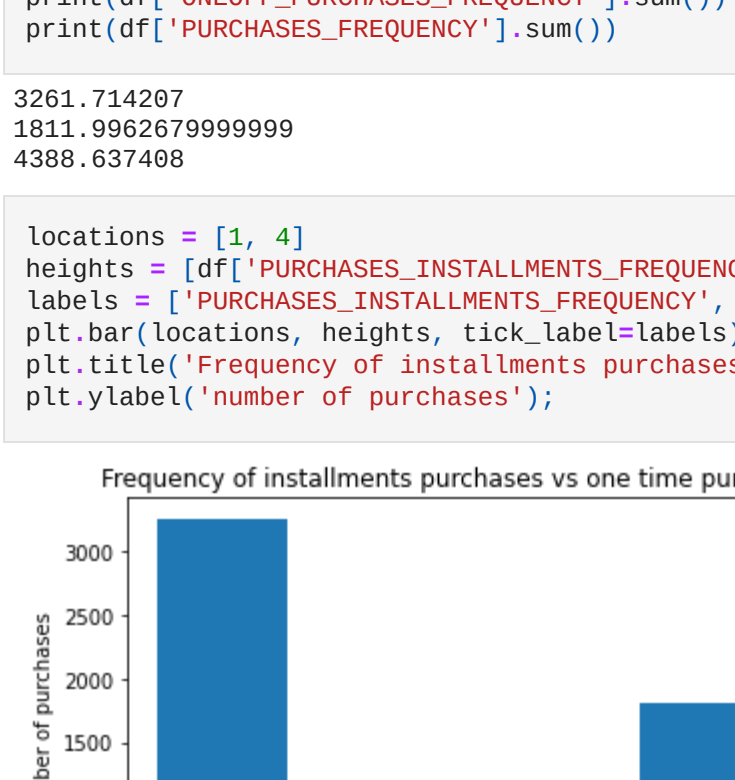


```
In [27]: #sns.pairplot(df_clean)
#plt.show()
```

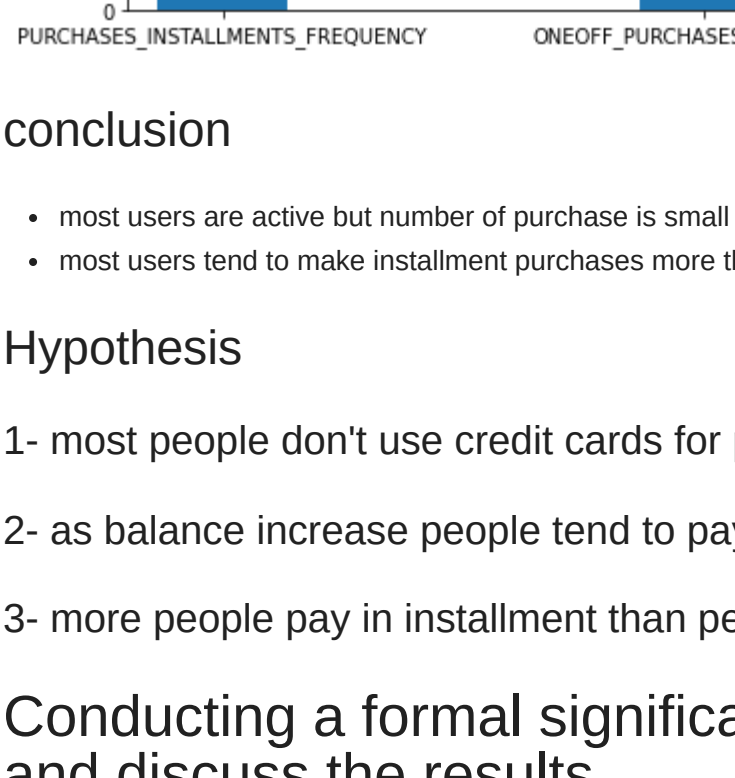
multiple columns are highly correlated

- Balance is positively correlated with Minimum payment, cash advance, balance frequency
- Balance is negatively correlated with prc_full_payment which is percentage of full payment paid by the user
- so when balance increase people tend to pay in installment
- purchases are negatively correlated with cash advance, cash advance freq and cash advance tx
- so as purchases increase cash paid in advance by user decreases and the rate of down payments decrease

```
In [28]: df['BALANCE'].hist();
```

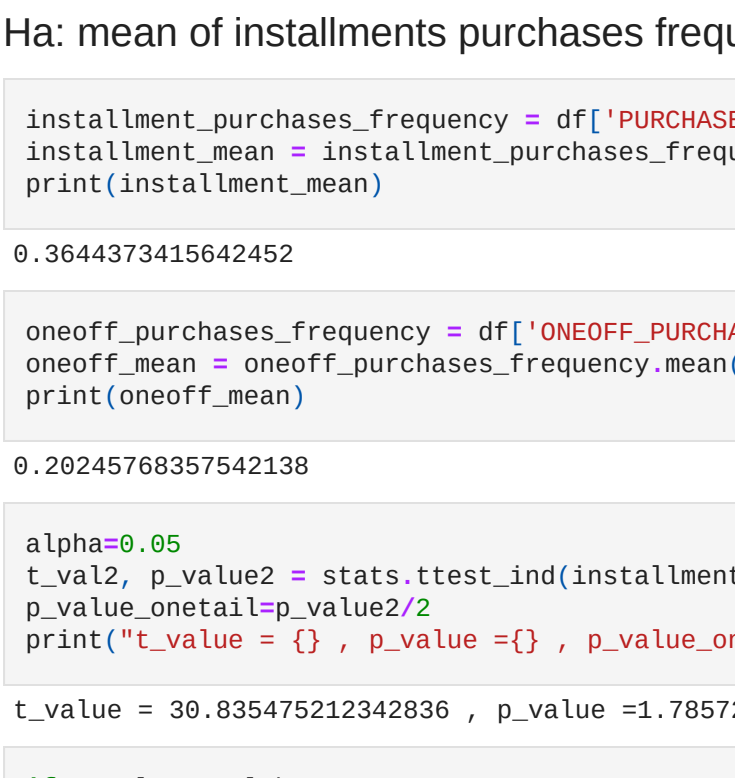


```
In [29]: df['BALANCE_FREQUENCY'].hist();
```



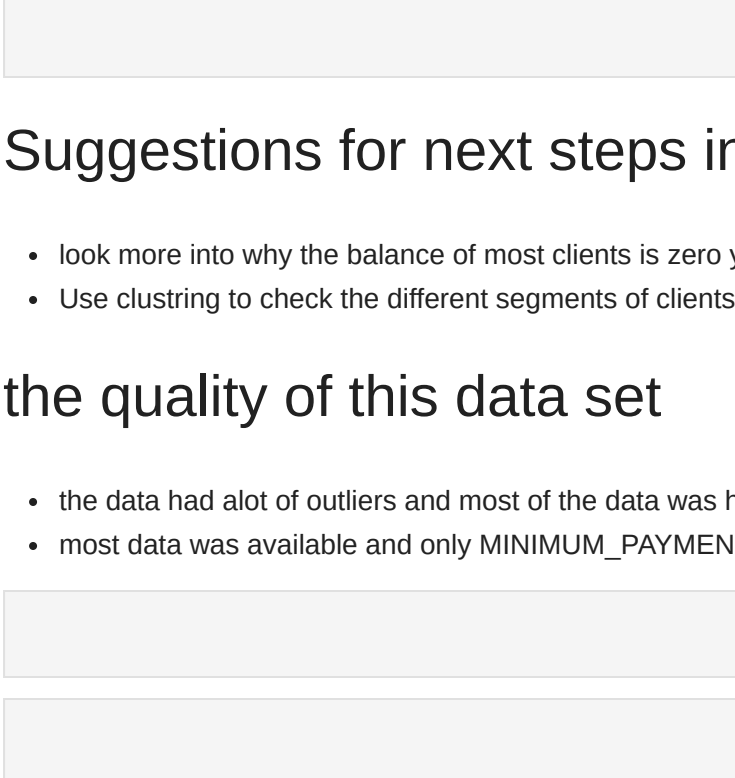
- most accounts have zero balance
- most accounts get updated frequently so most clients are active users

```
In [30]: df['PURCHASES'].hist();
plt.xlabel('Amount of purchases')
plt.ylabel('number of clients who made a purchase with this amount')
plt.show()
```



- most people have purchase of 0

```
In [31]: df['PURCHASES_FREQUENCY'].hist();
```



- number of people who make frequent purchase are aproximatly equal to number of people who make no purchases

```
In [32]: print(df['PURCHASES_INSTALLMENTS_FREQUENCY'].sum())
print(df['ONEOFF_PURCHASES_FREQUENCY'].sum())
print(df['PURCHASES_FREQUENCY'].sum())
```

```
3261.714287
3831.9962679999999
4388.637488
```

```
In [33]: locations = [1, 4]
heights = [df['PURCHASES_INSTALLMENTS_FREQUENCY'].sum(), df['ONEOFF_PURCHASES_FREQUENCY'].sum()]
labels = ['PURCHASES_INSTALLMENTS_FREQUENCY', 'ONEOFF_PURCHASES_FREQUENCY']
plt.bar(locations, heights, tick_label=labels)
plt.title('frequency of installments purchases vs one time purchases')
plt.ylabel('number of purchases');
```



conclusion

- most users are active but number of purchase is small for most users and most users keep balance at zero
- most users tend to make installment purchases more than one time purchases

Hypothesis

1- most people don't use credit cards for purchases

2- as balance increase people tend to pay in installment

3- more people pay in installment than people who make one time purchases

Conducting a formal significance test for one of the hypotheses and discuss the results

Hypothesis three

more people pay in installment than people who make one time purchases

Ho: mean of installments purchases frequency <= mean of one time purchases frequency

Ha: mean of installments purchases frequency > mean of one time purchases frequency

```
In [34]: installment_purchases_frequency = df['PURCHASES_INSTALLMENTS_FREQUENCY']
installment_mean = installment_purchases_frequency.mean()
print(installment_mean)
```

```
0.3644373415642452
```

```
In [35]: oneoff_purchases_frequency = df['ONEOFF_PURCHASES_FREQUENCY']
oneoff_mean = oneoff_purchases_frequency.mean()
print(oneoff_mean)
```

```
0.262457668357542138
```

```
In [36]: alpha=0.05
t_val2, p_value2 = stats.ttest_ind(installment_purchases_frequency, oneoff_purchases_frequency)
p_value_oneetail=p_value2/2
print("t_value = {}, p_value = {}, p_value_oneetail = {}".format(t_val2, p_value2, p_value_oneetail))
```

```
t_value = 38.835475212342836, p_value = 1.7857215613197782e-203, p_value_oneetail = 8.928687866598851e-204
```

```
In [37]: if p_value2 < alpha :
    print("Conclusion: Since p_value {} is less than alpha {}".format(p_value_oneetail,alpha) )
    print("Reject null hypothesis that mean of installments purchases frequency <= mean of one time purchases ")
else:
    print("Conclusion: Since p_value {} is greater than alpha {}".format(p_value_oneetail,alpha) )
    print("Failed to reject null hypothesis that mean of installments purchases frequency > mean of one time purchases ")
```

```
Conclusion: Since p_value 8.928687866598851e-204 is less than alpha 0.05
Reject null hypothesis that mean of installments purchases frequency <= mean of one time purchases frequency.
```

```
In [ ]:
```

Suggestions for next steps in analyzing this data

- look more into why the balance of most segments is zero yet most clients are active
- Use clustering to check the different segments of clients and check their behaviours

the quality of this data set

- the data had alot of outliers and most of the data was heavily skewed
- most data was available and only MINIMUM_PAYMENTS had missing data

```
In [ ]:
```

```
In [ ]:
```