# Credit EDA Case Study

**Submitted By :**

**ABHISHEK DAS**

Two types of risks are associated with the bank's decision:

1. TARGET(0)-If the applicant is likely to repay the loan, then not approving the loan results in a loss of business to the company
2. TARGET(1)-If the applicant is not likely to repay the loan, i.e. he/she is likely to default, then approving the loan may lead to a financial loss for the company.

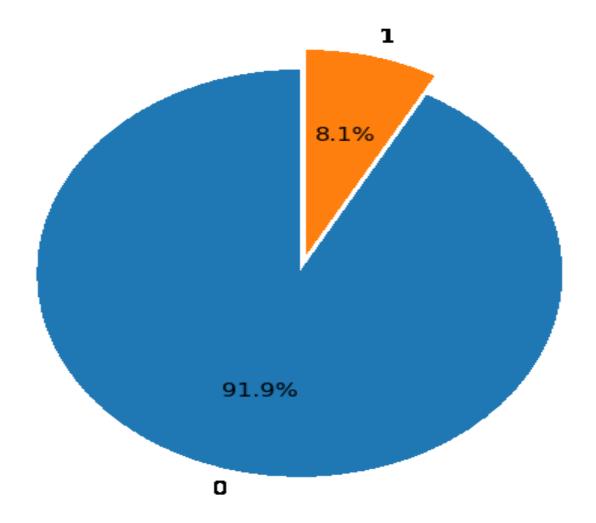ANALYSIS OF DATA IS DONE IN PYTHON ON A JUPITER NOTEBOOK

# DATA PRE PROCESSING

**STEPS:**

1. Check the structure of data(shape,size,info,describe)

2. Check missing value columns(which to drop,which to handle and how to handle and which column to keep as it is).

3. Drop any column which is irrelevant to TARGET variable

4. Check data type of each column,fix column with unappropriate data type and binary categorical column(into 0 and 1)

5. Check outliers and data imbalance .

6. Binning of continuous variables as required during analysis
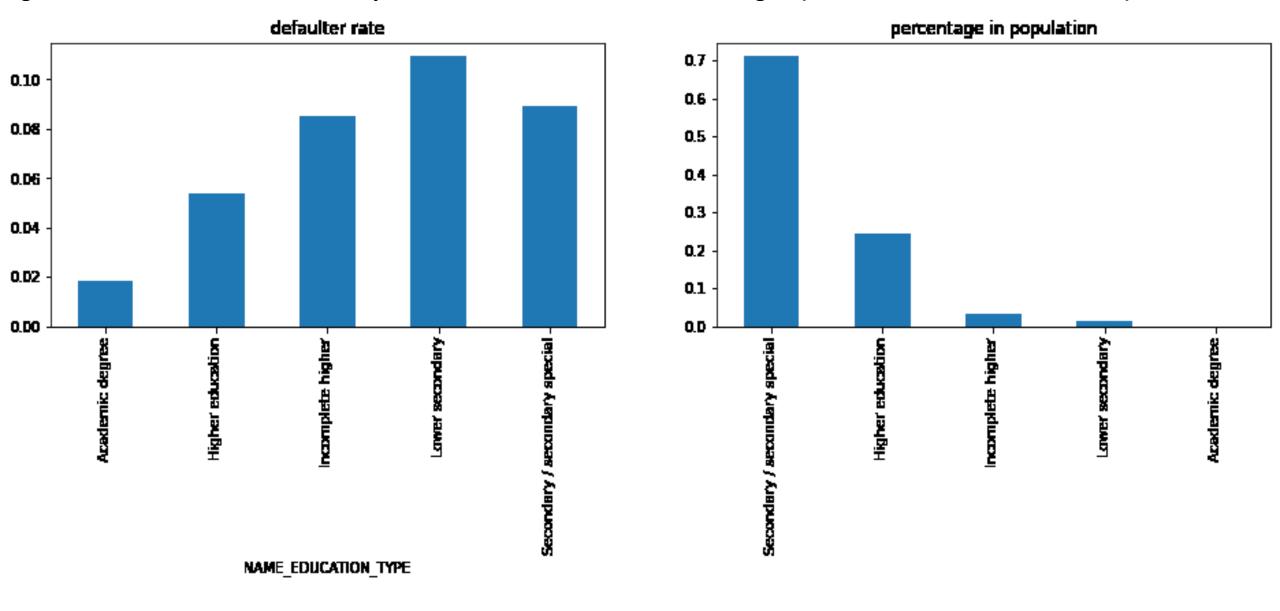
7. Top 10 correlation for defaulter and non defaulter

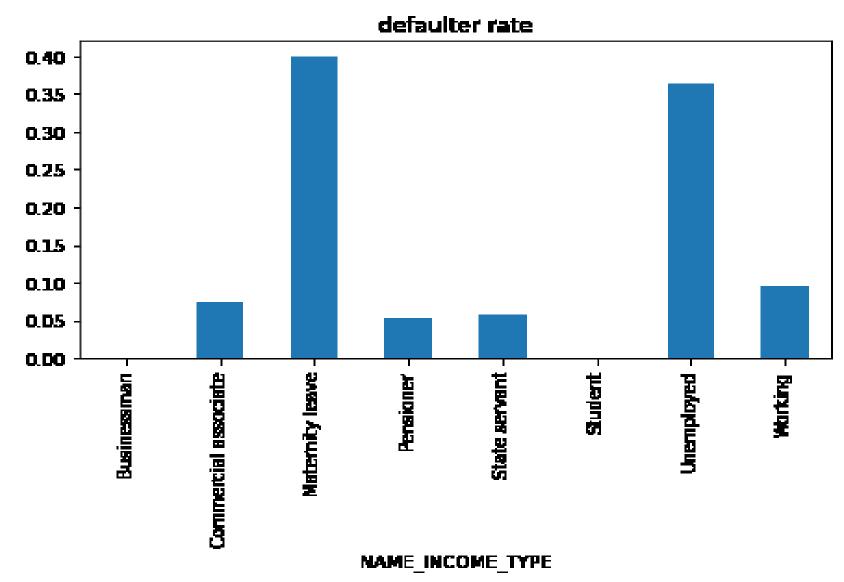imbalance in TARGET variable in percentage

Defaulter-8.1%

Non-defaulter-91.9%

# Higher default rate in lower secondary and least default rate in academic degree(default rate falls with education)
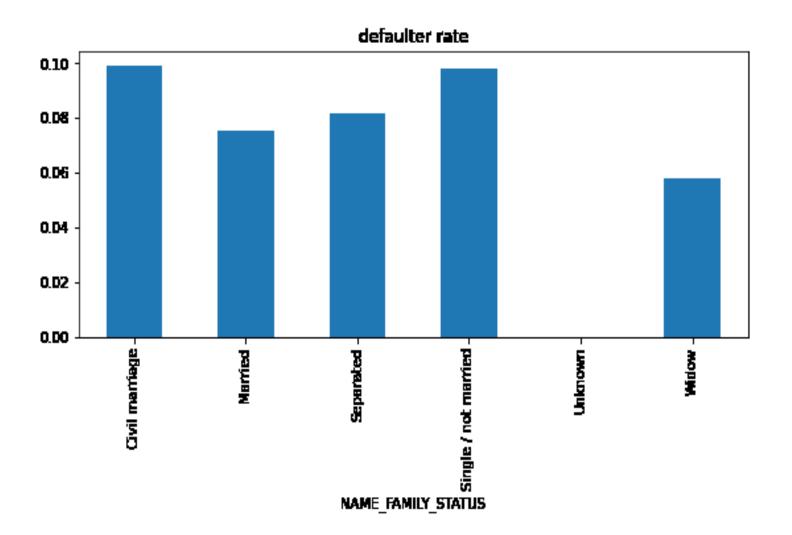
1)Student and buissnessman are safest.
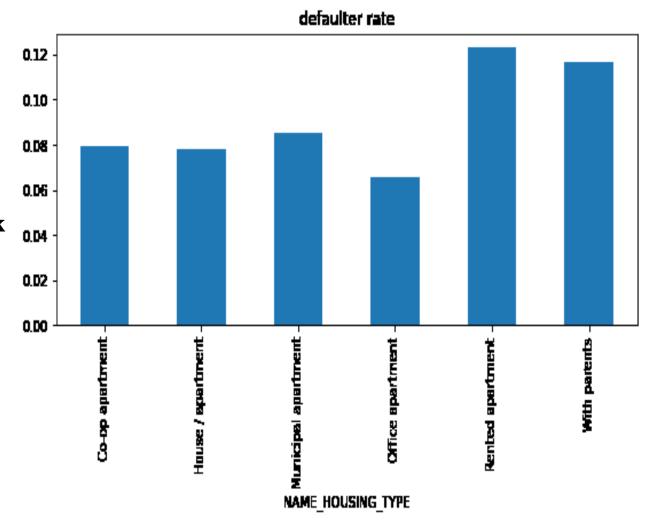2)Women in maternity leave and Unemployeds had highest default rate

defaulter rate

NAME_INCOME_TYPE

**1)Civil married and singles are risky**

**2)Widows are safest**
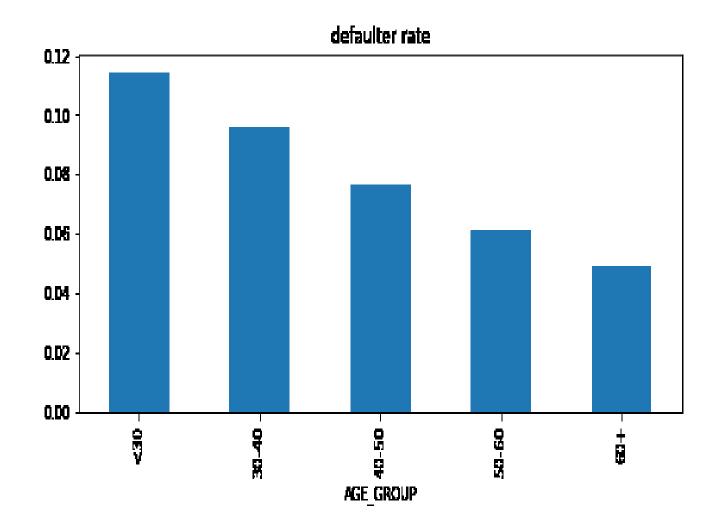


defaulter rate

**1)People living at rented apartment and with Parents are highly risky**

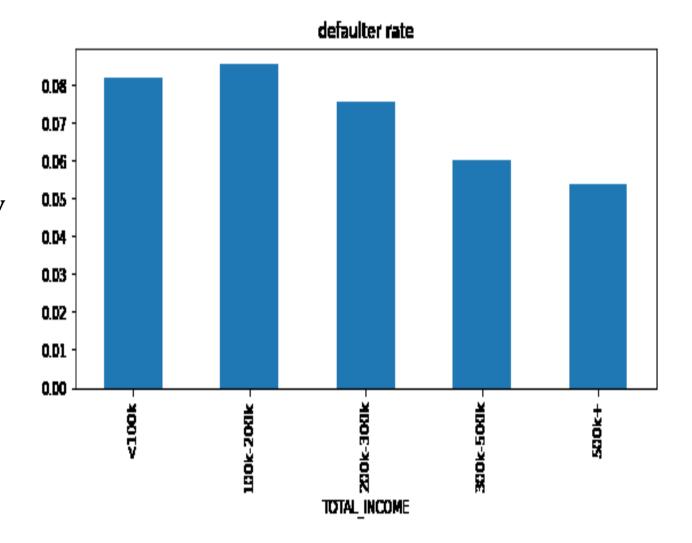**2)People living at office apartment are least at risk**

1)Default rate fall with age.

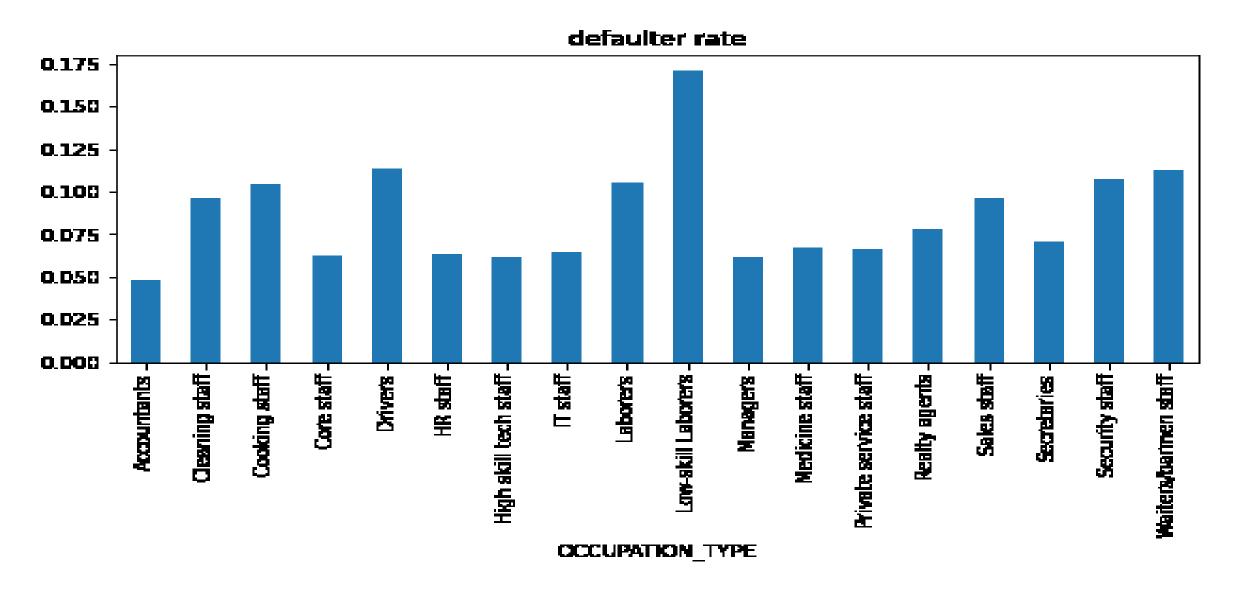2)People(<30) had highest default rate

3)people(>60) had least default rate



defaulter rate

1)Default rate fall with income.

2)People having income <200k are highest risky
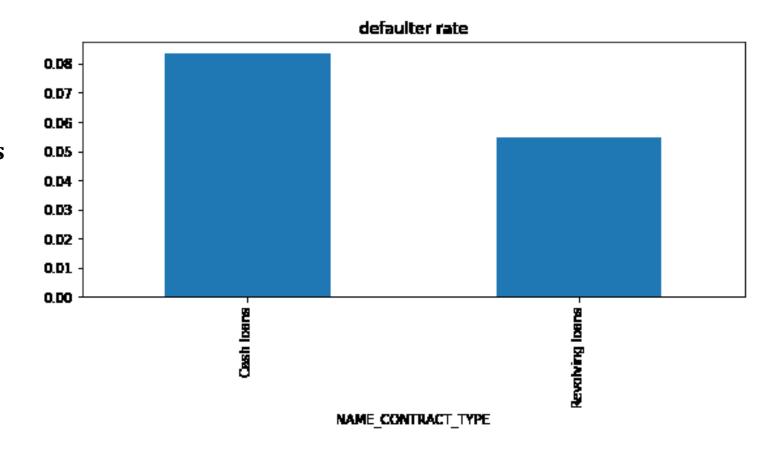
3)People with income >500k are safest

**1)** low skill labours(highest),waiters,drivers,labourers defaulter rate is high
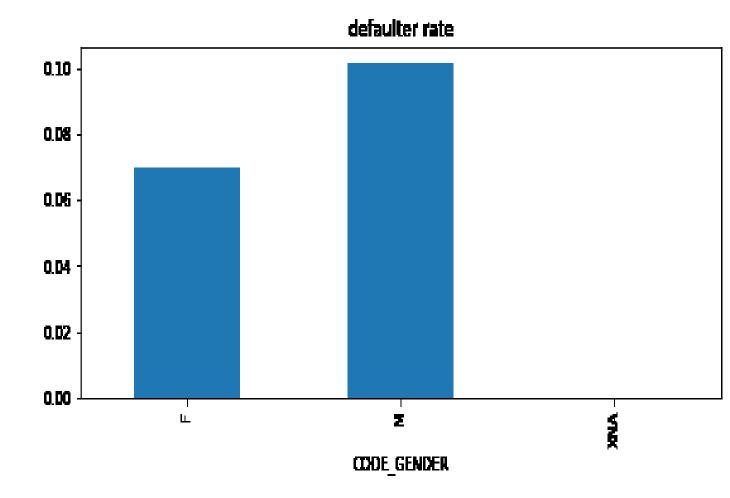
2) accountant,core staff,hr staff,IT staffs are safest

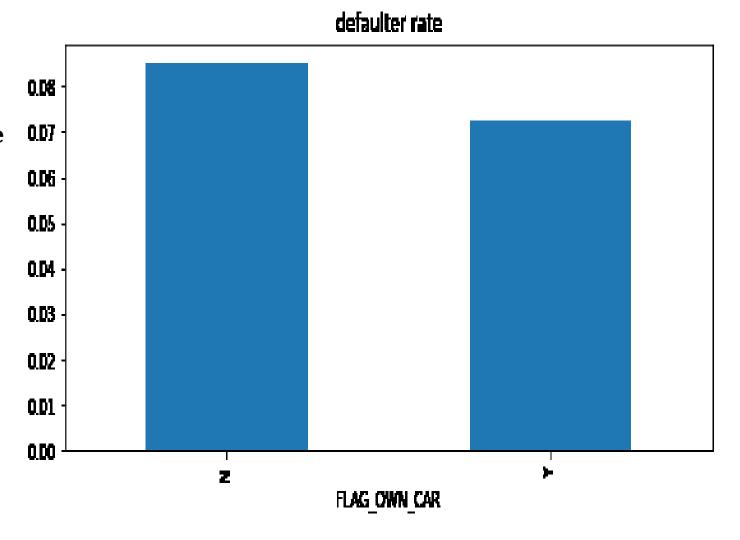

defaulter rate

**1)Defaulter seem to apply for cash loan**

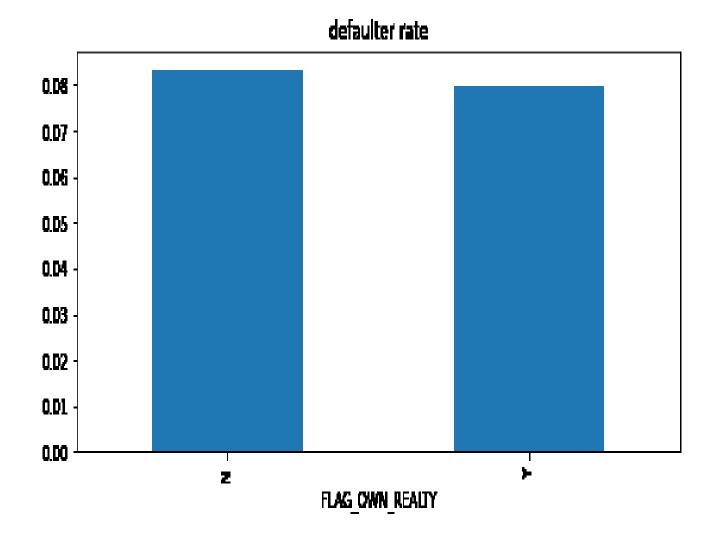**2)non-defaulter applies for revolving loans**
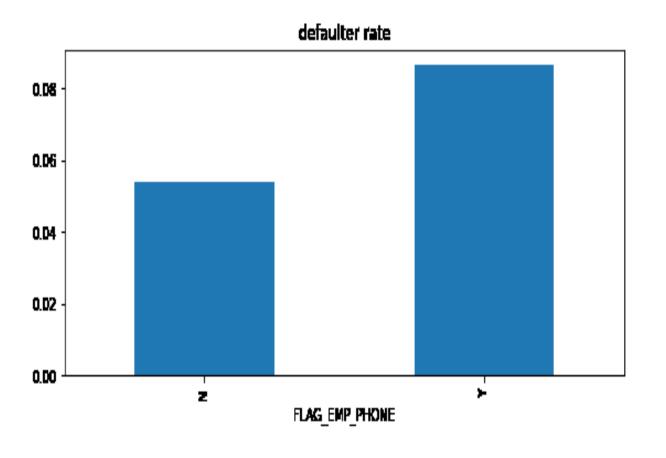
**Males seem to be more prone to default**

# defaulter rate

**1)Those who don't have car are more prone
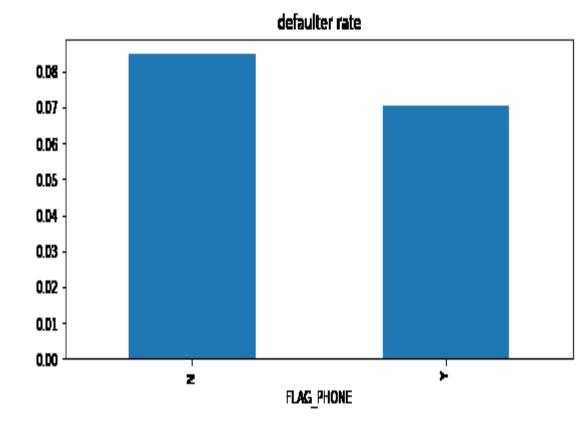To default**

**1)Those who don't own flat/house had high Chance of default**
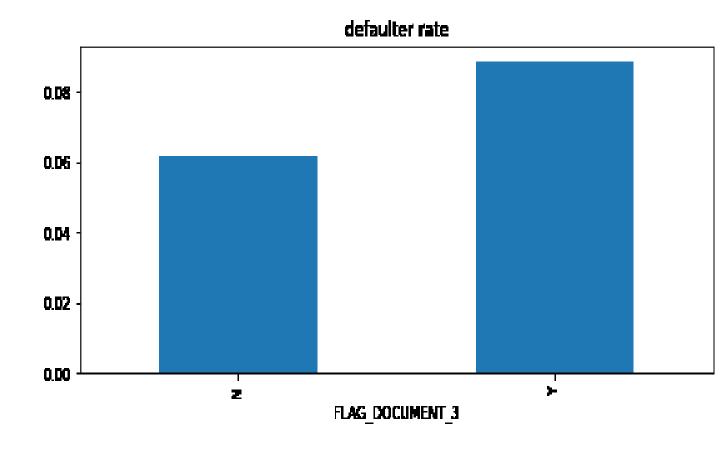


defaulter rate

**Defaulter had a tendency of giving work phone but not home phone in application**
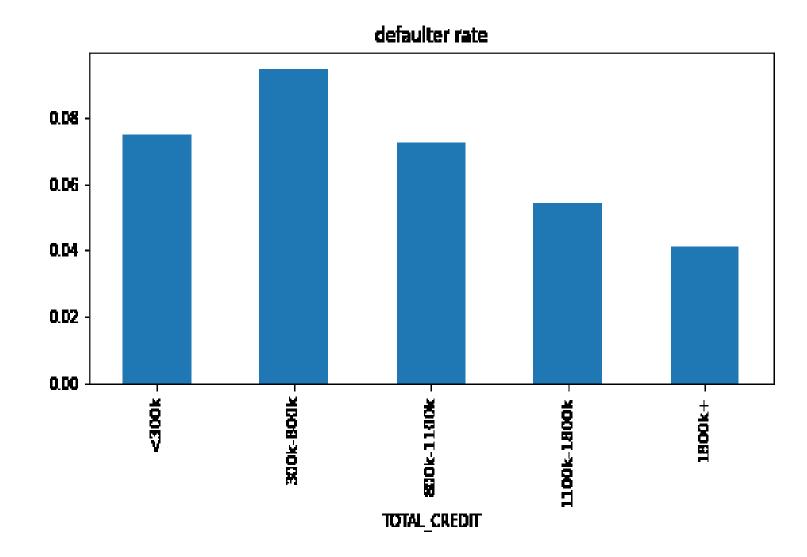
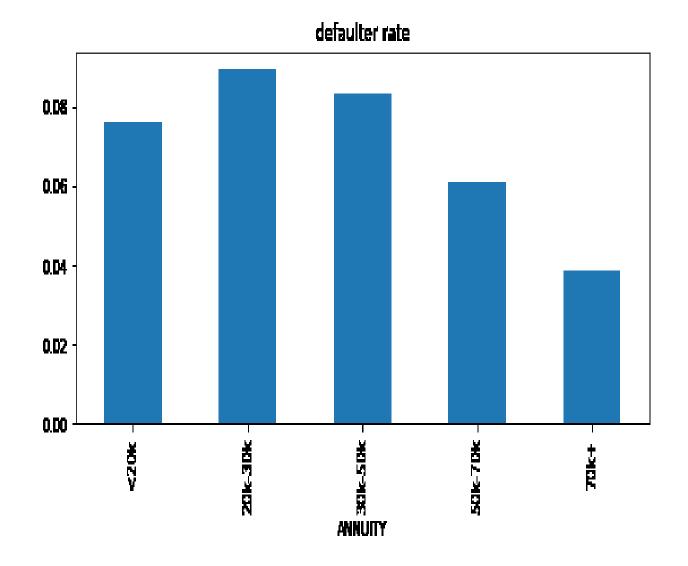**1)People who submitted document 3 had high Default rate**



defaulter rate

1)Default rate fall with credit amount.
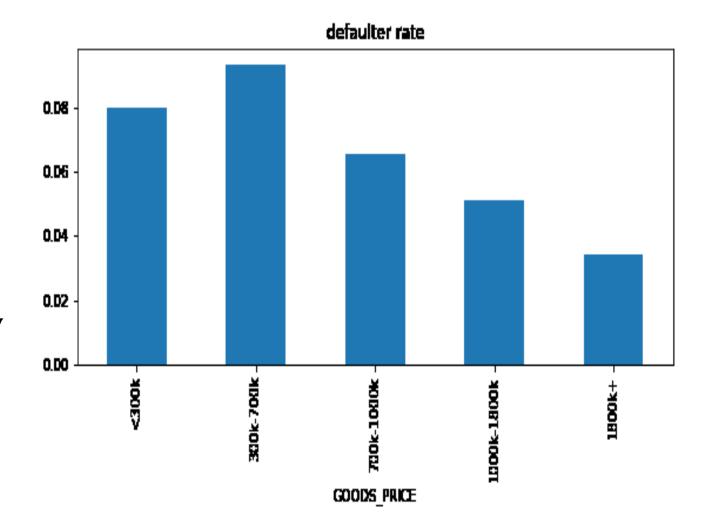2)Highest default rate(300k-800k)
3)Lowest default rate(>1800k)

**1)Default rate fall with annuity amount.**
**2)Highest default rate(20k-30k)**
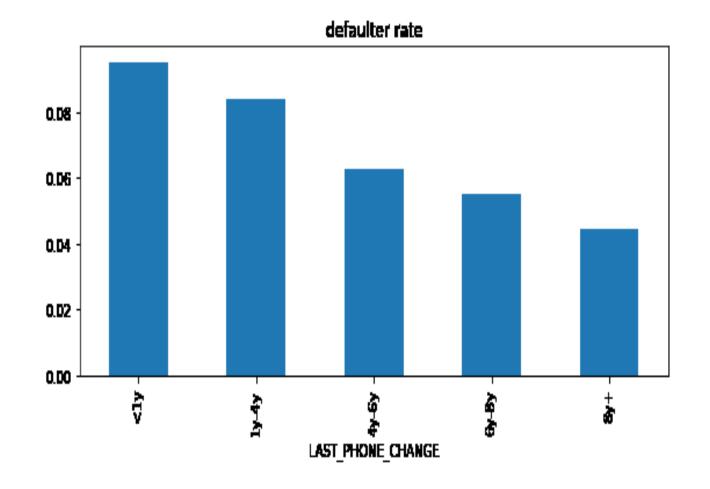**3)Lowest default rate(>70k)**



defaulter rate

**1)Default rate fall with goods_price.**
**2)Highest default rate(300k-700k)**
**3)Lowest default rate(>1800k)**


**So credit,annuity and good_price showing very**
**High correlation**



defaulter rate

GOODS_PRICE

**1)Defaulter seem to change their phone before Applying**

**2)People who changed(<1y)**

**3)People who changed(>8y) are least defaulter**



defaulter rate

correlation in defaulters

1)amt_credit vs amt_goods_price(corr=0.98)
2)cnt_family_member vs cnt_children(corr=0.89)
3)DEF_30_CNT_SOCIAL_CIRCLE vs DEF_60_CNT_SOCIAL_CIRCLE(corr=0.87)
4)amt_credit vs amt_annuity(corr=0.75)
5)amt_annuity vs amt_goods_price(corr=0.75)
6)days_birth vs days_employed(corr=0.58)
 7)DEF_30_CNT_SOCIAL_CIRCLE vs OBS_60_CNT_SOCIAL_CIRCLE(corr=0.34)
8)DEF_30_CNT_SOCIAL_CIRCLE vs OBS_30_CNT_SOCIAL_CIRCLE(corr=0.33)
9)DEF_60_CNT_SOCIAL_CIRCLE vs OBS_30_CNT_SOCIAL_CIRCLE(corr=0.26)
10)days_birth vs days_id_publish(corr=0.25)

1)amt_credit vs amt_goods_price(corr=0.99)
2)cnt_family_member vs cnt_children(corr=0.88)
3)DEF_30_CNT_SOCIAL_CIRCLE vs DEF_60_CNT_SOCIAL_CIRCLE(corr=0.86)
4)amt_annuity vs amt_goods_price(corr=0.78)
5)amt_credit vs amt_annuity(corr=0.77)
6)days_birth vs days_employed(corr=0.63)
7)DEF_30_CNT_SOCIAL_CIRCLE vs OBS_60_CNT_SOCIAL_CIRCLE(corr=0.33)
8)DEF_30_CNT_SOCIAL_CIRCLE vs OBS_30_CNT_SOCIAL_CIRCLE(corr=0.33)
9)days_employed vs days_id_publish(corr=0.28)
10)DEF_60_CNT_SOCIAL_CIRCLE vs OBS_30_CNT_SOCIAL_CIRCLE(corr=0.25)

observation-sequence of top 10 correlation in defaulter and non-defaulter is almost same

# observations

1. women in maternity leave who applied for high credit had very high chance of default

2. to give small credit loan(below 800k),students are good choice

3. unemployed people are comparatively safer to give big credit(above 8ook)

4. married person are comparative safe to give loan,only when you are giving above 800k credit(just increase credit amount to ensure further more safety)

5. below 500k credit,person with academic degree is very safe

6. for giving loan to person with secondary education go beyond 1000k(toward high credit)

7. person with lower secondary education is highestrisk,choosing those who ask for credit beyond 1500k ensure safety

as you move toward higher days of employment,defaulting chance reduces

**For previous applicants,8.7% are defaults**



imbalance in TARGET variable in percentage

defaulter rate

## CONCLUSIONS ABOUT DEFAULTERS

1)firstly,around 8% of population is found to be defaulter and mens are more likely to be defaulter then females
2)education=defaulter rate increases with decline in education so people with lower secondary education and then secondary education had highest default rate
3)income type=unemployed(second highest defaulter rate ) and women in maternity(highest defaulter rate) had highest default rate
4)family status=in defaulter rate,civil married is highest followed by single(second)
5)housing status=people living in rented apartment had highest share of defaulters,followed by those who live with parents
6)age=defaulter rate falls with age so people under 30 and then in 30's had highest default rate
7)income=default rate increase with fall of income. people with total income below 200000 are highest in risk
8)occupation type=low skill labourers had highest default rate.other serving class like waiter,driver,labourers are also at high risk
9)contract type=people applying for cash loan had high default rate
10)own_car=people who dont own car had high chance of defaulting
11)house/flat=people who dont own real state(house/flat) had higher risk of default
12)in defaulter there is a pattern that they provide work phone but not home phone
13)doc_3=people who submitted this document had higher default rate
14)total_credit=people who applied for credit of 3lac to 8lac had highest default rate then default rate fall with credit

15)annuity and goods price-following same trend as credit, people with annuity(20k-30k) and good price(300k-700k) had highest default rate,falling with increase in value

16)changing phone=defaulter had tendency of changing phone right before application.people who change phone(<1yr) are highest in risk.default rate decreases with time

17)women in maternity leave who applied for high credit had very high chance of default

18)employment_days=people recently employed had highest chance of default

19)persons who was rejected in prev app due to reject code SCOFR,turned defaulter in most number

20)defaulter had a tendency to refuse to name the goal for which they applied for cash loan,as observed

21)defaulters seem to apply for next loan application quicker than non-defaulter

1)firstly,around 92% of population is found to be non-defaulter and females are more likely to be non-defaulter then males

2)education=non-defaulter tendency increases with education level so people with academic degree had least chance of default

3)income type=students and buissnessman are most likely to be non-defaulter(minimum risk)

4)family status=widow is least risky

5)housing status=people living at office apartment are least risky

6)age=defaulter rate falls with age so people with age 60+ are least risky

7)income=default rate falls with increase in income. people with total income above 500k are least in risk

8)occupation type=accountant,core staff,hr staff,IT staffs are safest

9)contract type=people applying for revolving loan are safe

10)own_car=people who own car had non-defaulting tendency

11)house/flat=people who own real state(house/flat) had non-default tendency

12)non-defaulter provides home phone

13)total_credit= default rate fall with credit amount,people applied for 1800k credit are safest

14)annuity and goods price-following same trend as credit, people with annuity(>70k) and good price(>1800k) had least default rate

15)changing phone=as high the duration with which person didnt change his phone from app date,that much low is risk.person who didnt change phone from 8y+ had leat risk

16)for small credit(which is risky),students are safest.

17)for unemployed people(who are risky),its safer to give big credit(above 8ook)

18)married person are comparative safe to give loan,only when you are giving above 800k credit(just increase credit amount to ensure further more safety)

19)for giving loan to person with secondary eucation(second high in risk) go beyond 1000k(toward high credit)

20)for person with lower secondary education(highest risk),choosing those who ask for credit beyond 1500k ensure safety

21)as you move toward higher days of employment,defaulting chance reduces

22)person whose early application rejected due to system(code reject reason) is least risky