

## PGM Assignment 2

### Group 19

Name	BITS ID
Vinayak Nayak	2021K04135
Shreyasi Kalra	2021K04586
Karan Khajuria	2021K04951

We have modelled the problem as a classification problem considering `Loan_Payment_Status` as the target variable and all other variables as independent ones. We have used the following techniques in modelling

- Logistic Regression with SGD Classifier from sklearn
- MultiLayer Perceptron (Deep Learning Model)
- Probabilistic Deep Learning Model

Due to limitation of compute resources and time, we have made certain modelling assumptions which we have stated as and when we make these modelling decisions. Also the auxiliary files, models and code which we used to perform this training could be found [here](#).

### System Configuration for model training

- For Training Logistic Regression with SGD Classifier & keras MLP DL Model (Local)
  - 12 GB RAM
    - 100 GB Disk Space
    - ~2 days of compute
    - 4 GB GPU RAM
      - 12 GB RAM
      - 128 GB Disk Space
      - 16 GB GPU RAM
      - ~6 hrs of compute
  - For training Probabilistic Deep Learning Model (Google Colab)
    - 12 GB RAM
    - 128 GB Disk Space
    - 16 GB GPU RAM
    - ~6 hrs of compute

```
In [ ]: # Import standard libraries needed for processing
from collections import defaultdict, Counter
from pathlib import Path
import pickle, re, csv
from datetime import datetime
from pathlib import Path
from tqdm import tqdm

# Import linear algebra libraries
import numpy as np
import pandas as pd
import random

# Import plotting libraries
import matplotlib.pyplot as plt
plt.style.use('ggplot')
import matplotlib
import seaborn as sns

# Import preprocessing libraries
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import OneHotEncoder, StandardScaler
from sklearn.linear_model import SGDClassifier

# Suppress unnecessary warnings
import warnings
warnings.filterwarnings("ignore")

In [ ]: # Get column names by reading the first line only
f = open("MultiFamily/FNMA_MF_Loan_Performance_Data_202203.csv", "r")
columns = f.readline().replace("\n", "").split(",")
print(columns)

# Loan Number, 'Acquisition Date', 'Note Date', 'Maturity Date at Acquisition', 'Loan Acquisition
# UPB', 'Amortization Type', 'Interest Type', 'Loan Product Type', 'Original UPB', 'Amortization Ter
# m', 'Original Interest Rate', 'Lien Position', 'Transaction ID', 'Issue Date', 'Loan Acquisitio
# n LTV', 'Underwritten DSCR Type', 'Underwritten DSCR Type', 'Original I/O Term', 'I/O
# End Date', 'Loan Ever 60+ Days Delinquent', 'Loss Sharing Type', 'Modified Loss Sharing Percenta
# ge', 'Number of Properties at Acquisition', 'Property Acquisition Total Unit Count', 'Specific Prop
# erty Type', 'Year Built', 'Property City', 'Property State', 'Property Zip Code', 'Metropolitan Sta
# tistical Area', 'Physical Occupancy %', 'Liquidation/Prepayment Code', 'Liquidation/Prepayment Da
# te', 'Foreclosure Date', 'Credit Event Date', 'Foreclosure Value', 'Lifetime Net Credit Loss Amou
# nt', 'Sale Price', 'Default Amount', 'Credit Event Type', 'Reporting Period Date', 'Loan Active Pro
# perty Count', 'Note Rate', 'Maturity Date - Current', 'UPB - Current', 'Delinquency UPB', 'Loan Pa
# yment Status', 'Prepayment Indicator', 'Most Recent Modification Date', 'Modification Indicator', 'Defea
# ant Date', 'Prepayment Provision', 'Prepayment Provision End Date', 'Affordable Housing Type', 'M
# CRT Deal ID', 'MCAS Deal ID', 'DUS Prepayment Outcomes', 'DUS Prepayment Segments', 'Loan Age'

In [ ]: # Read all the lines
# Read first one fourth of columns and iteratively read all slices to avoid running out of memory
# Do this iteratively with slice
# slice(0, 15), slice(15, 30), slice(30, 45), slice(45, 60)
lengths = []
d = defaultdict(lambda: [])
slc = slice(0, 15)

with open("MultiFamily/FNMA_MF_Loan_Performance_Data_202203.csv", newline="")
) as csvfile:
    filereader = csv.reader(csvfile, delimiter=",", quotechar='"')
    for row in filereader:
        feats = row[slc]
        for f, c in zip(feats, columns[slc]):
            if f != c:
                d[c].append(f)

In [ ]: # Write the individual columns to a single file for easier preprocessing
for k, v in items(d):
    with open(f"piecewise_columns_files/{k}.replace('/', '_').txt", "w") as f_out:
        for item in v:
            f_out.write(f"{item}\n")
        f_out.close()

In [ ]: # Free the memory occupied unnecessarily
import gc
del d, gc.collect()
```

## Read the individual column files and summarize them

```
In [ ]: # from IPython.display import IFrame
# IFrame(src='https://capitalmarkets.fanniemae.com/media/5986', width=700, height=600)

In [ ]: # If the percentage of missing values in a column is more than 90%, then it is not possible to model
# this variable; in such a case, Imputation might lead to highly unrealistic results
ACCEPTANCE_THRESHOLD_PCT = 80
cols_to_drop = []
reccs = []

def summarize_file(filepath):
    with open(filepath, "r") as f:
        coldata = [x.replace("\n", "") for x in f.readlines()]

        # Find out the number of unique values and proportion of missing data
        unique_vals = len(set([x for x in coldata if x != ""]))
        sample = [x.replace("\n", "") for x in list(set(coldata))[:7]]
        samples = [x for x in sample if x != ""]
        actual_data = [x for x in coldata if x != ""]

        # Compute missing data proportion as percentages
        missing_vals = len(coldata) - len(actual_data)
        missing_prop = round(missing_vals / len(coldata) * 100, 3)
        actual_prop = 100 - missing_prop

        # If prop of missing data is more than a particular threshold then drop those columns
        if missing_prop > ACCEPTANCE_THRESHOLD_PCT:
            cols_to_drop.append(filepath.stem)

        # Make a dataframe of summary_stats
        recs.append([filepath.stem, unique_vals, missing_vals, missing_prop, actual_prop])

    # Print each column summary
    print(f"Filepath: {filepath.stem}<20 | Unique: {str(unique_vals)<5} | Missing Count: {str(missing_vals)<5}
    | Missing Proportion: {str(missing_prop)<5} | Actual Proportion: {str(actual_prop)<5}")

In [ ]: # Get the names of all the files
files = [x for x in Path("piecewise_columns_files").glob("**/*.txt")
         if (x.is_file() and (not ".ipynb" in str(x)))]

# Read all the files
for f in files:
    summarize_file(f)

Loan Acquisition LTV | Unique: 2334 | Missing Count: 193 | Missing Proportion: 0.005% | Availab
le: 99.995%
Five examples of data in this column: ['24.72', '78.7', '14.0', '5.9', '57.15', '14.5']

SQD Indicator | Unique: 2 | Missing Count: 0 | Missing Proportion: 0.000% | Availab
le: 100.000%
Five examples of data in this column: ['Y', 'N']

Underwritten DSCR Type | Unique: 4 | Missing Count: 828 | Missing Proportion: 0.020% | Avail
able: 99.980%
Five examples of data in this column: ['Deal UW DSCR NCF', 'UW DSCR NCF', 'Lender UW DSCR', 'UW Ac
tual DSCR']

Property State | Unique: 54 | Missing Count: 0 | Missing Proportion: 0.000% | Availab
le: 100.000%
Five examples of data in this column: ['NH', 'GA', 'MS', 'PA', 'WI', 'PR', 'AK']

Liquidation/Prepayment Date | Unique: 3205 | Missing Count: 404369 | Missing Proportion: 99.146%
| Available: 0.854%
Five examples of data in this column: ['2005-10-11', '2015-01-30', '2012-07-13', '2019-04-12', '20
22-01-24', '2015-10-08']

Specific Property Type | Unique: 8 | Missing Count: 0 | Missing Proportion: 0.000% | Avail
able: 100.000%
Five examples of data in this column: ['Other', 'Military', 'Seniors', 'Cooperative', 'Dedicated S
tudent', 'Multiple Properties', 'Manufactured Housing Community']

Modification Indicator | Unique: 2 | Missing Count: 0 | Missing Proportion: 0.000% | Avail
able: 100.000%
Five examples of data in this column: ['Y', 'N']

Foreclosure Value | Unique: 480 | Missing Count: 407717 | Missing Proportion: 99.982% | Avail
able: 0.018%
Five examples of data in this column: ['$5,206,000.00', '$2,740,000.00', '$700,000.00', '$1,240,00
0.00', '$2,675,000.00', '$2,480,000.00']

Delinquency UPB | Unique: 5688 | Missing Count: 4060399 | Missing Proportion: 99.558% | Avail
able: 0.442%
Five examples of data in this column: ['$964,268.99', '$3,823,637.93', '$472,943.74', '$20,687.91
5.19', '$2,338,705.99', '$10,771,206.03']

DUS Prepayment Outcomes | Unique: 8 | Missing Count: 716445 | Missing Proportion: 17.567% | Av
ailable: 82.433%
Five examples of data in this column: ['Paid Prior to Declining Premium End Date', 'Other', 'Invol
untary Prepayment', 'Paid Prior to Yield Maintenance End Date', 'Paid On or After Loan End Date',
'Paid On or After Yield Maintenance End Date']

Lien Position | Unique: 4 | Missing Count: 96 | Missing Proportion: 0.002% | Availab
le: 99.998%
Five examples of data in this column: ['Second', 'Third', 'Fourth or More Subordinate', 'First']

MCIRT Deal ID | Unique: 10 | Missing Count: 3659238 | Missing Proportion: 89.722% | Avail
able: 10.278%
Five examples of data in this column: ['MCIRT 2018-01', 'MCIRT 2017-01', 'MCIRT 2020-01', 'MCIRT 2
021-01', 'MCIRT 2021-02', 'MCIRT 2019-03']

Original Term | Unique: 281 | Missing Count: 0 | Missing Count: 0 | Missing Proportion: 0.000% | Availab
le: 100.000%
Five examples of data in this column: ['7.8', '100', '372', '238', '143', '139', '349']

Number of Properties at Acquisition | Unique: 19 | Missing Count: 57096 | Missing Proportion:
1.400% | Available: 98.600%
Five examples of data in this column: ['2', '4', '8', '18', '24', '6']

Note Date | Unique: 5565 | Missing Count: 0 | Missing Proportion: 0.000% | Availab
le: 100.000%
Five examples of data in this column: ['2018-05-02', '2013-06-10', '2009-06-03', '2020-01-04', '20
15-08-25', '2016-05-18', '2020-08-03']

Year Built | Unique: 179 | Missing Count: 8 | Missing Proportion: 0.000% | Availab
le: 100.000%
Five examples of data in this column: ['1994', '2019', '1961', '1860', '1881', '1899']

Metropolitan Statistical Area | Unique: 387 | Missing Count: 0 | Missing Proportion: 0.000%
| Available: 100.000%
Five examples of data in this column: ['KANNAKEE, IL METROPOLITAN STATISTICAL AREA', 'MAISON, WI
METROPOLITAN STATISTICAL AREA', 'SAN DIEGO-CHULA VISTA-CARLSBAD, CA METROPOLITAN STATISTICAL ARE
A', 'MODESTO, CA METROPOLITAN STATISTICAL AREA', 'MOUNT VERNON-ANACORTES, WA METROPOLITAN STATISTI
CAL AREA', 'FOCALTELO, ID METROPOLITAN STATISTICAL AREA', 'DAYTON-KETTERING, OH METROPOLITAN STATI
STICAL AREA']

Physical Occupancy % | Unique: 377 | Missing Count: 81425 | Missing Proportion: 1.996% | Availab
le: 98.004%
Five examples of data in this column: ['7.8', '86.2', '55.0', '70.6', '72.0', '96.2']

Amortization Term | Unique: 219 | Missing Count: 87752 | Missing Proportion: 2.152% | Availab
le: 97.848%
Five examples of data in this column: ['100', '74', '238', '143', '139', '81']

Credit Event Type | Unique: 2 | Missing Count: 4034976 | Missing Proportion: 98.934% | Avail
able: 1.066%
Five examples of data in this column: ['Non-REO', 'REO']

Loan Payment Status | Unique: 4 | Missing Count: 0 | Missing Proportion: 0.000% | Availab
le: 100.000%
Five examples of data in this column: ['60-89 Days Delinquent', 'Current', '90+ Days Delinquent',
'30-59 Days Delinquent']

Transaction ID | Unique: 52255 | Missing Count: 472276 | Missing Proportion: 11.580% | Avail
able: 88.420%
Five examples of data in this column: ['467722', '465857', 'BL3138', 'B50835', 'BL7169', 'BL3391']

Issue Date | Unique: 267 | Missing Count: 475272 | Missing Proportion: 11.653% | Avail
able: 88.347%
Five examples of data in this column: ['2014-01-01', '2010-06-01', '2018-06-01', '2007-10-01', '20
18-07-01', '2003-12-01']

UPB - Current | Unique: 3068727 | Missing Count: 0 | Missing Proportion: 0.000% | Avail
able: 100.000%
Five examples of data in this column: ['$1,562,444.34', '$8,202,852.00', '$1,895,899.28', '$1,324,
298.62', '$2,324,111.55', '$923,119.63', '$27,584,690.03']

Loan Ever 60+ Days Delinquent | Unique: 2 | Missing Count: 0 | Missing Proportion: 0.000%
| Available: 100.000%
Five examples of data in this column: ['Y', 'N']

Underwritten DSCR | Unique: 1805 | Missing Count: 820 | Missing Proportion: 0.020% | Availab
le: 99.980%
Five examples of data in this column: ['14.1', '5.9', '8.11', '20.68', '4.68']

Maturity Date at Acquisition | Unique: 572 | Missing Count: 0 | Missing Proportion: 0.000% | Avail
able: 100.000%
Five examples of data in this column: ['2040-10-01', '2044-07-01', '2007-10-01', '2018-07-01', '20
36-03-01', '2003-08-01', '2050-01-01']

Affordable Housing Type | Unique: 8 | Missing Count: 3903565 | Missing Proportion: 95.712% | A
vailable: 4.288%
Five examples of data in this column: ['Other - Special Public Sponsor', 'Other', 'LIHTC or Project
Based HAP/Sec 8', 'Non HAP', 'Multiple Properties', 'Other - Sponsor Initiated Affordability']

Property Zip Code | Unique: 8368 | Missing Count: 0 | Missing Proportion: 0.000% | Availab
le: 100.000%
Five examples of data in this column: ['41001', '44234', '90210', '77018', '71112', '40504', '2390
1']

Modified Loss Sharing Percentage | Unique: 63 | Missing Count: 3898377 | Missing Proportion: 98.900
% | Available: 1.100%
Five examples of data in this column: ['62.5000000000', '25.0000000000', '19.2000000000', '85.1000
000000', '58.0000000000', '17.1000000000']

Original Interest Rate | Unique: 2084 | Missing Count: 0 | Missing Proportion: 0.000% | Avail
able: 100.000%
Five examples of data in this column: ['7.385', '5.9', '7.72', '2.0825', '4.68', '7.68', '2.513']

Credit Event Date | Unique: 476 | Missing Count: 4077744 | Missing Proportion: 99.983% | Avail
able: 0.017%
Five examples of data in this column: ['2005-10-11', '2022-01-24', '2012-07-13', '2012-02-15', '20
12-01-24', '2015-04-15']

Original I/O Term | Unique: 150 | Missing Count: 2736500 | Missing Proportion: 67.097% | Avail
able: 32.903%
Five examples of data in this column: ['74', '100', '23', '139', '81', '103']

Acquisition Date | Unique: 5029 | Missing Count: 0 | Missing Proportion: 0.000% | Availab
le: 100.000%
Five examples of data in this column: ['2018-05-02', '2013-06-10', '2009-06-03', '2015-08-25', '20
16-05-18', '2016-04-19', '2007-11-15']

Loan Number | Unique: 61048 | Missing Count: 0 | Missing Proportion: 0.000% | Availab
le: 100.000%
Five examples of data in this column: ['8300002588', '1720005003', '1717464188', '1717460470', '17
1178769', '1717479774', '1717484703']

Defeasance Date | Unique: 226 | Missing Count: 4065726 | Missing Proportion: 99.640% | Avail
able: 0.360%
Five examples of data in this column: ['2015-01-30', '2006-05-26', '2013-10-31', '2007-11-02', '20
10-12-30', '2013-03-21']

Property City | Unique: 3640 | Missing Count: 0 | Missing Proportion: 0.000% | Availab
le: 100.000%
Five examples of data in this column: ['REGO PARK', 'OOLTEWAH', 'KIHIE', 'FLEMINGTON', 'STUTTGAR
T', 'DAMA POINT', 'ROMULUS']

Amortization Type | Unique: 5 | Missing Count: 8 | Missing Proportion: 0.000% | Availab
le: 100.000%
Five examples of data in this column: ['Interest Only/Balloon', 'Amortizing Balloon', 'Interest On
ly/Fully Amortizing', 'Fully Amortizing', 'Interest Only/Amortizing/Balloon']

I/O End Date | Unique: 411 | Missing Count: 2733418 | Missing Proportion: 67.021% | Avail
able: 32.979%
Five examples of data in this column: ['2007-10-01', '2018-07-01', '2003-12-01', '2036-03-01', '20
03-08-01', '2015-08-01']

Liquidation/Prepayment Code | Unique: 11 | Missing Count: 4043609 | Missing Proportion: 99.146%
| Available: 0.854%
Five examples of data in this column: ['LI Third Party Sale', 'Liquidation', 'Dissolution', 'Fully P
aid, Refinanced', 'Other Liquidation', 'Fully Paid, Matured']

Foreclosure Date | Unique: 419 | Missing Count: 4077862 | Missing Proportion: 99.984% | Avail
able: 0.016%
Five examples of data in this column: ['2010-06-04', '2012-02-15', '2010-08-03', '2016-03-25', '20
06-01-02', '2012-03-12']

MCAS Deal ID | Unique: 2 | Missing Count: 4056178 | Missing Proportion: 99.454% | Avail
able: 0.546%
Five examples of data in this column: ['MCAS 2020-01', 'MCAS 2019-01']

Loan Active Property Count | Unique: 64 | Missing Count: 77943 | Missing Proportion: 1.889% | A
vailable: 98.111%
Five examples of data in this column: ['74', '22', '23', '4', '8', '45']

Interest Type | Unique: 2 | Missing Count: 0 | Missing Proportion: 0.000% | Availab
le: 100.000%
Five examples of data in this column: ['Fixed', 'ARM']

Default Amount | Unique: 656 | Missing Count: 4077774 | Missing Proportion: 99.984% | Avail
able: 0.016%
Five examples of data in this column: ['$6,348,242.94', '$12,476,382.12', '$8,997,903.99', '$933,0
10.06', '$5,859,497.29', '$5,768,640.00']

Prepayment Provision | Unique: 1826 | Missing Count: 22210 | Missing Proportion: 0.545% | Avail
able: 99.455%
Five examples of data in this column: ['YM(80), 0(6)', 'YM(168), See Issuance Documents(36)', 'YM
(122), 0*(5)', 'YM(36), 0*(24)', 'YM(72), See Issuance Documents(24)', 'YM(132), 0(12)']

Prepayment Provision End Date | Unique: 3366 | Missing Count: 1112384 | Missing Proportion: 27.27
5% | Available: 72.725%
Five examples of data in this column: ['YM(09/30/2023), 1%(06/30/2027), 0(10/01/2027)', 'YM(05/31/
2027), 0(12/31/2027)', 'YM(09/30/2024), See Issuance Documents(10/01/2025)', 'YM(04/30/2025), See
Issuance Documents(05/01/2040)', 'YM(01/31/2028), See Issuance Documents(02/01/2031)', 'YM(11/30/2
020), 1%(02/28/2021), See Issuance Documents(06/01/2021)']

Maturity Date - Current | Unique: 583 | Missing Count: 0 | Missing Proportion: 0.000% | Avail
able: 100.000%
Five examples of data in this column: ['2040-10-01', '2044-07-01', '2007-10-01', '2018-07-01', '20
36-03-01', '2003-08-01', '2050-01-01']

Lifetime Net Credit Loss Amount | Unique: 666 | Missing Count: 4077744 | Missing Proportion: 99.
983% | Available: 0.017%
Five examples of data in this column: ['$618,805.14', '$933,010.06', '$2,768,773.21', '$93,706.
96', '$2,641,120.12', '$233,578.40']

Sale Price | Unique: 414 | Missing Count: 4077866 | Missing Proportion: 99.986% | Avail
able: 0.014%
Five examples of data in this column: ['$825,000.00', '$5,206,000.00', '$21,250,000.00', '$700,00
0.00', '$405,000.00', '$14,600,000.00']

Loan Age | Unique: 228 | Missing Count: 1808265 | Missing Proportion: 44.337% | Avail
able: 55.663%
Five examples of data in this column: ['74', '100', '23', '143', '139', '81']

Loan Product Type | Unique: 4 | Missing Count: 0 | Missing Proportion: 0.000% | Availab
le: 100.000%
Five examples of data in this column: ['Credit Facility', 'Bulk Delivery', 'DUS', 'Non-DUS']

Loan Acquisition UPB | Unique: 23804 | Missing Count: 0 | Missing Proportion: 0.000% | Availab
le: 100.000%
Five examples of data in this column: ['$29,239,000.00', '$6,233,500.00', '$1,333,405.31', '$2,04
4,500.00', '$11,123,000.00', '$25,094,000.00', '$23,230,000.00']

Reporting Period Date | Unique: 267 | Missing Count: 0 | Missing Proportion: 0.000% | Availa
ble: 100.000%
Five examples of data in this column: ['2014-01-01', '2010-06-01', '2018-06-01', '2007-10-01', '20
18-07-01', '2003-12-01', '2002-09-01']

Note Rate | Unique: 6764 | Missing Count: 4361 | Missing Proportion: 0.107% | Availab
le: 99.893%
Five examples of data in this column: ['2.2812', '3.6116', '4.124', '9.78', '8.555', '6.697']

DUS Prepayment Segments | Unique: 5 | Missing Count: 716409 | Missing Proportion: 17.566% | Av
ailable: 82.434%
Five examples of data in this column: ['ARM 7-6', 'Fixed Rate, Declining Premium', 'Nonstandard',
'Standard', 'SARM']

Most Recent Modification Date | Unique: 228 | Missing Count: 4040183 | Missing Proportion: 99.06
2% | Available: 0.938%
Five examples of data in this column: ['2014-01-01', '2010-06-01', '2014-02-04', '2013-10-31', '20
18-07-01', '2013-04-10']

Loss Sharing Type | Unique: 3 | Missing Count: 167134 | Missing Proportion: 4.098% | Availa
ble: 95.902%
Five examples of data in this column: ['Pari Passu', 'Standard DUS', 'No Lender Loss Sharing']

Original UPB | Unique: 19577 | Missing Count: 0 | Missing Proportion: 0.000% | Availab
le: 100.000%
Five examples of data in this column: ['$5,208,000.00', '$29,239,000.00', '$1,801,000.00', '$15,66
0,000.00', '$1,337,000.00', '$5,571,640.00', '$6,530,000.00']

Property Acquisition Total Unit Count | Unique: 902 | Missing Count: 81425 | Missing Proportion:
1.996% | Available: 98.004%
Five examples of data in this column: ['100', '603', '426', '143', '139', '784']
```

```
In [ ]: # Make a dataframe of records and display the columns which would be dropped based on the threshold
df = pd.DataFrame(reccs.values())
print(len(cols_to_drop))
print([".".join(cols_to_drop)])
```

```
16
Liquidation/Prepayment Date | Foreclosure Value | Delinquency UPB | MCIRT Deal ID | Credit Event Ty
pe | Affordable Housing Type | Modified Loss Sharing Percentage | Credit Event Date | Defeasance T
ype | Liquidation/Prepayment Code | Foreclosure Date | MCAS Deal ID | Default Amount | Lifetime N
et Credit Loss Amount | Sale Price | Most Recent Modification Date
```

```
In [ ]: # Get the files which are needed and which are not needed
all_text_files = [x for x in Path("piecewise_columns_files").glob("**/*.txt") if x.is_file()]
val_text_files = [x for x in all_text_files if (not x.stem in cols_to_drop)]
unwanted_text_files = [x for x in all_text_files if x not in val_text_files]
len(all_text_files), len(val_text_files), len(unwanted_text_files)
```

```
Out[ ]: (56, 44, 16)
```

```
In [ ]: # View all the stats at once in increasing order of missing percentage
df.sort_values(by=["MissingProp(%)", "reset_index(drop=True)"]
```

23	Maturity Date - Current	583	0	0.000	100.000
24	Loan Product Type	4	0	0.000	100.000
25	Lien Position	4	96	0.002	99.998
26	Loan Acquisition LTV	2334	193	0.005	99.995
27	Underwritten DSCR	1005	820	0.020	99.980
28	Underwritten DSCR Type	4	820	0.020	99.980
29	Note Rate	6764	4361	0.107	99.893
30	Prepayment Provision	1826	22210	0.545	99.455
31	Number of Properties at Acquisition	19	57096	1.400	98.600
32	Loan Active Property Count	64	77043	1.889	98.111
33	Property Acquisition Total Unit Count	902	81425	1.996	98.004
34	Physical Occupancy %	377	81425	1.996	98.004
35	Amortization Term	219	87752	2.152	97.848
36	Loss Sharing Type	3	167134	4.098	95.902
37	"Transaction ID"	52255	472276	11.580	88.420
38	Issue Date	267	475272	11.653	88.347
39	DUS Prepayment Segments	5	716409	17.566	82.434
40	DUS Prepayment Outcomes	8	716445	17.567	82.433
41	Prepayment Provision End Date	3366	1112384	27.275	72.725
42	Loan Age	228	1808265	44.337	55.663
43	L/O End Date	411	2733418	67.021	32.979
44	Original L/O Term	150	2736500	67.097	32.903
45	MCIRT Deal ID	10	3659238	67.720	10.278
46	Modified Loss Sharing Percentage	63	3898377	95.585	4.415
47	Affordable Housing Type	8	3903565	95.712	4.288
48	Credit Event Type	2	4034976	98.934	1.066
49	Most Recent Modification Date	228	4040183	99.062	0.938
50	Liquidation_Payment Code	11	4043609	99.146	0.854
51	Liquidation_Payment Date	3205	4043609	99.146	0.854
52	MCAS Deal ID	2	4056178	99.454	0.546
53	Delinquency UPB	5688	4060399	99.558	0.442
54	Defeasance Date	226	4063770	99.640	0.360
55	Foreclosure Value	480	4077717	99.982	0.018
56	Lifetime Net Credit Loss Amount	666	4077744	99.983	0.017
57	Credit Event Date	476	4077744	99.983	0.017
58	Foreclosure Date	419	4077802	99.984	0.016
59	Default Amount	656	4077774	99.984	0.016
60	Sale Price	414	4077866	99.986	0.014

# Look at the dataframe by sorting cardinality and missing proportion of all the columns  
df = df[df["MissingProp(%)"] < ACCEPTANCE\_THRESHOLD\_PCT].reset\_index(drop=True)  
df.sort\_values(by=["Unique", "MissingProp(%)"], reset\_index(drop=True))

	Filepath	Unique	MissingCount	MissingProp(%)	ActualProp(%)
0	SDQ Indicator	2	0	0.000	100.000
1	Modification Indicator	2	0	0.000	100.000
2	Loan Ever 60+ Days Delinquent	2	0	0.000	100.000
3	Interest Type	2	0	0.000	100.000
4	Loss Sharing Type	3	167134	4.098	95.902



