

Zillow's Price Feature Engineering

Chan Yu Yankai Liu Yifei Bi

Procedures Guide

01

02

- 1.Preparation
- 2.Impute Data 1
- 3.Creating New Features
- 4.Impute Data 2

03

04

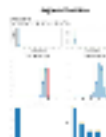
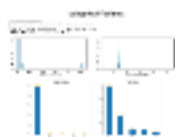
Reporter : Chan Yu

1.Convert Float64 to Float32

	parcelid	airconditioningtypeid	architecturalstyletypeid	basementsqft	bathroomcnt	bedroomcnt	buildingclasstypeid
0	10754147	1.0	7.0	535.0	0.0	0.0	4.0
1	10759547	1.0	7.0	535.0	0.0	0.0	4.0
2	10843547	1.0	7.0	535.0	0.0	0.0	5.0
3	10859147	1.0	7.0	535.0	0.0	0.0	3.0
4	10870947	1.0	7.0	535.0	0.0	0.0	4.0

5 rows x 8 columns

2.Group Features into Different Groups



3.Prepare Functions For Corresponding Group

1.Preparation

2.Impute Data 1

3.Creating New Features

4.Impute Data 2

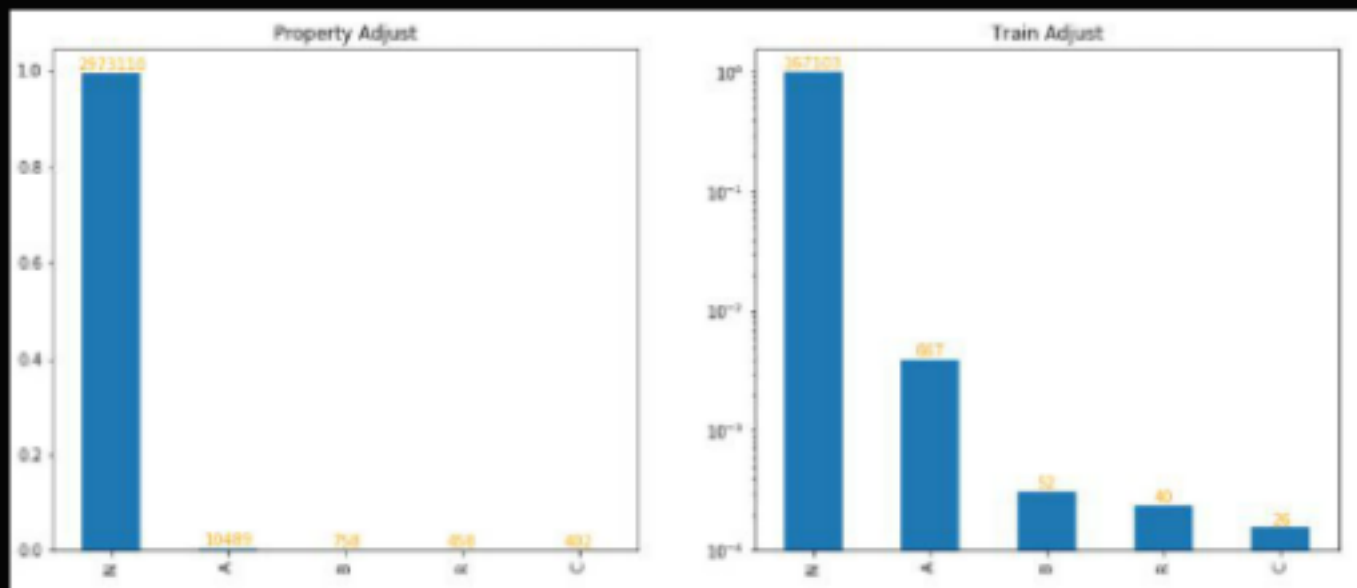
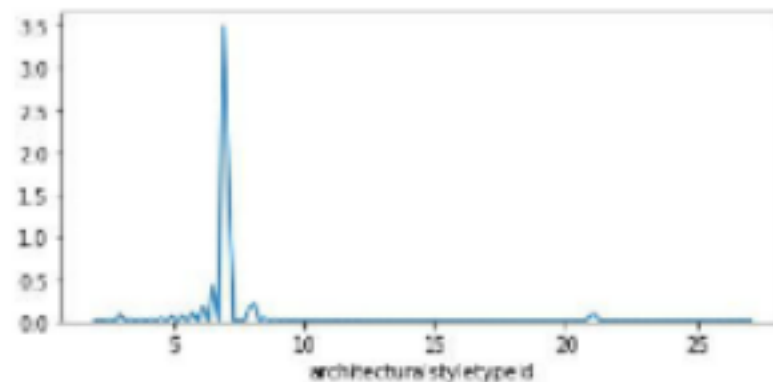
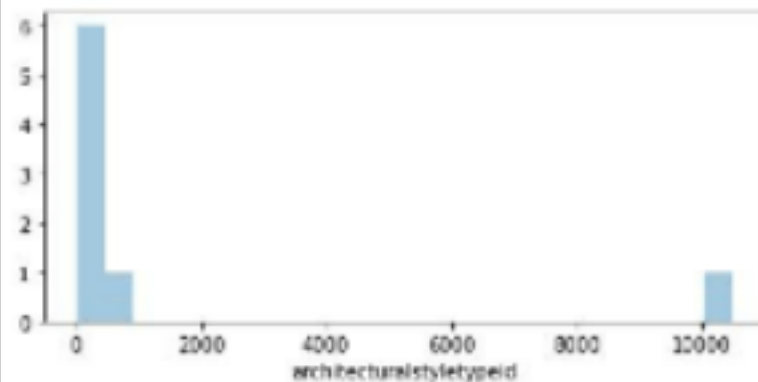
Categorical Features

architecturalstyletypeid

Full NaN : 90.89% Full NumOfCat : 9 Full CatRatio : 0.07%

Train NaN : 96.53% Train NumOfCat : 8 Train CatRatio : 1.11%

[10489L, 758L, 402L, 300L, 116L, 38L, 2L, 2L]



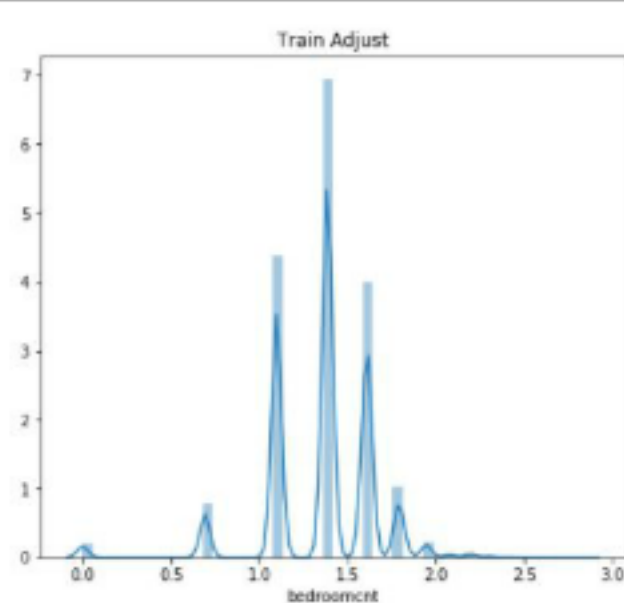
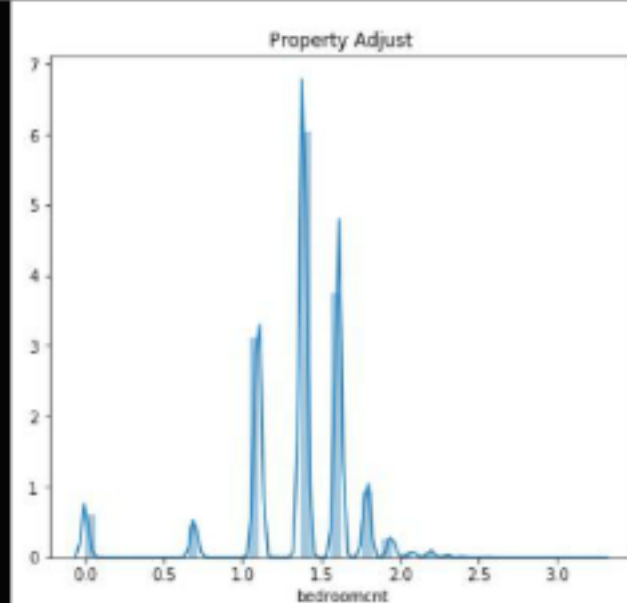
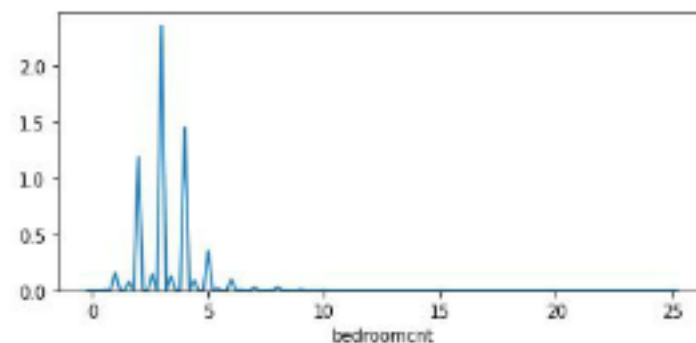
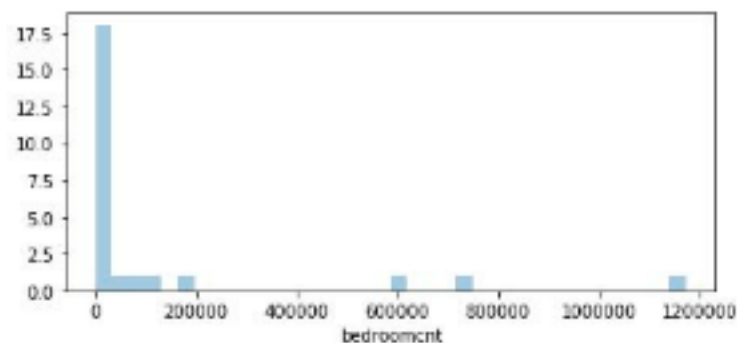
Numerical Features

bedroomcnt

Full NaN : 0.10% Full NumOfCat : 26 Full CatRatio : 0.00%

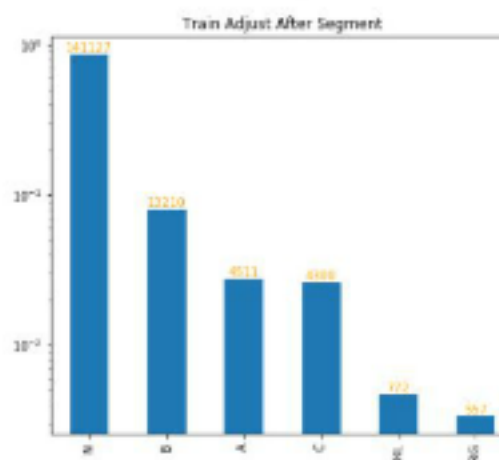
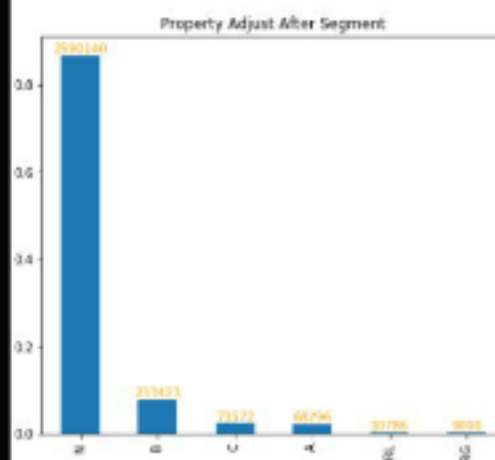
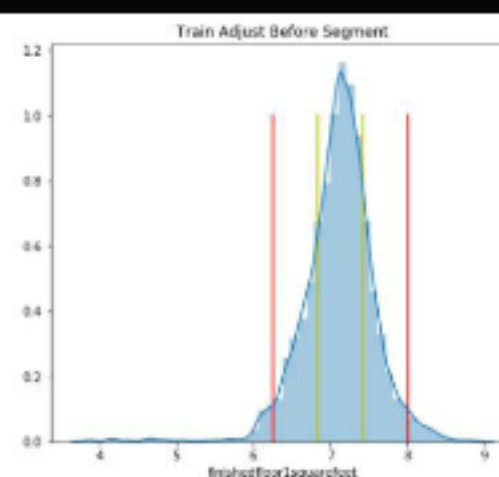
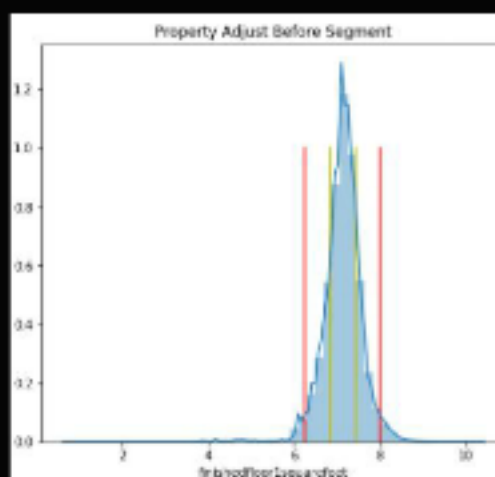
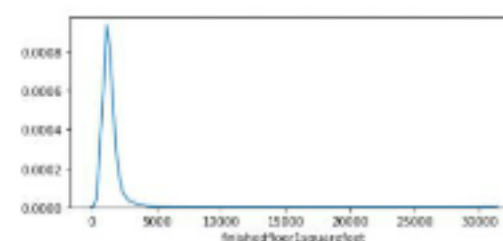
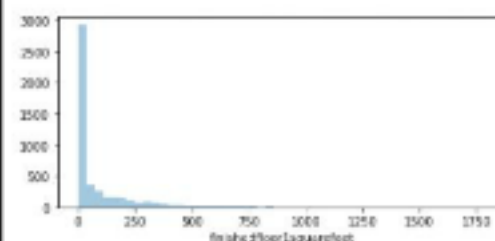
Train NaN : 0.03% Train NumOfCat : 18 Train CatRatio : 0.01%

[1172760L, 731476L, 606787L, 182768L, 118705L, 86942L, 48915L, 13542L, 12763L, 4279L]



Segment Features

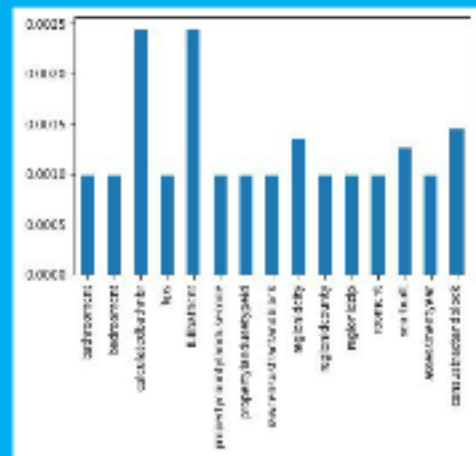
FinishedFloorSquareFeet
 Full Count: 16,796 Full Mean: 4941 Full StdDev: 1,176
 Train Count: 85,896 Train Mean: 2881 Train StdDev: 15,848
 [178EL, 1706L, 1366L, 1191L, 1160L, 1140L, 1113L, 1101L, 1107L, 1126L]



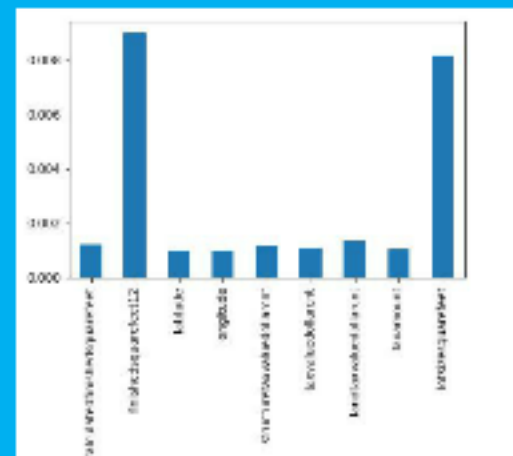
2

- 1.Preparation
- 2.Impute Data 1
- 3.Creating New Features
- 4.Impute Data 2

Impute By Median



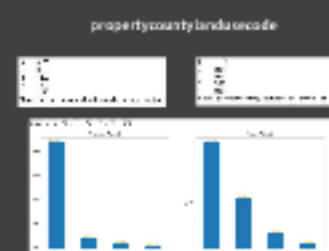
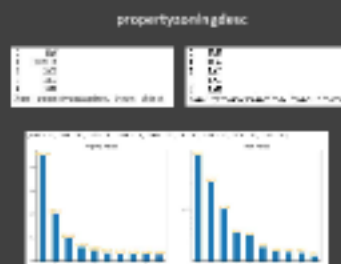
Impute By Mode



Missing ratio < 10%

3

- 1.Preparation
- 2.Impute Data 1
- 3.Creating New Features
- 4.Impute Data 2

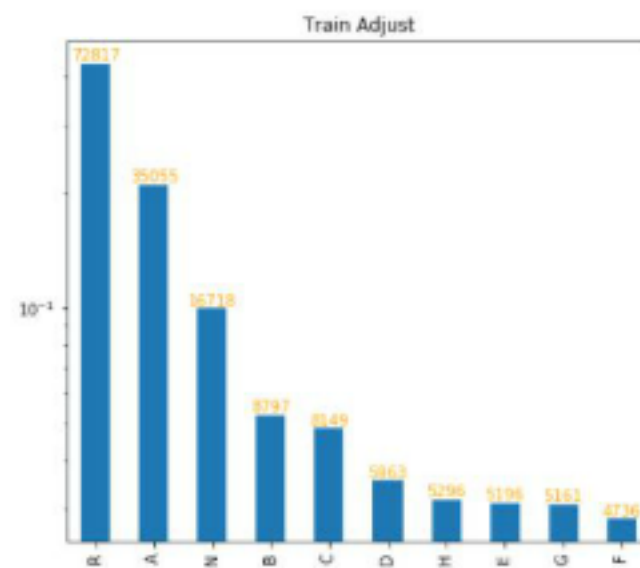
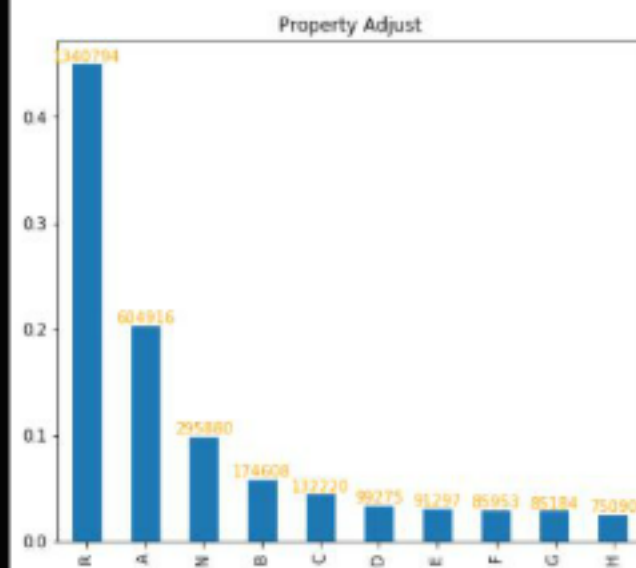


propertyzoningdesc

```
0      NaN
1  LCA11*
2      LAC2
3      LAC2
4      LAM1
Name: propertyzoningdesc, dtype: object
```

```
0      NaN
1      LCA
2      LAC
3      LAC
4      LAM
Name: propertyzoningdesc, dtype: object
```

('LARS': 'F', 'LER': 'E', 'LER': 'E', 'SCUR': 'H', 'LARD': 'C', 'N': 'N', 'LARE': 'G', 'LCA': 'D', 'LAR': 'A')

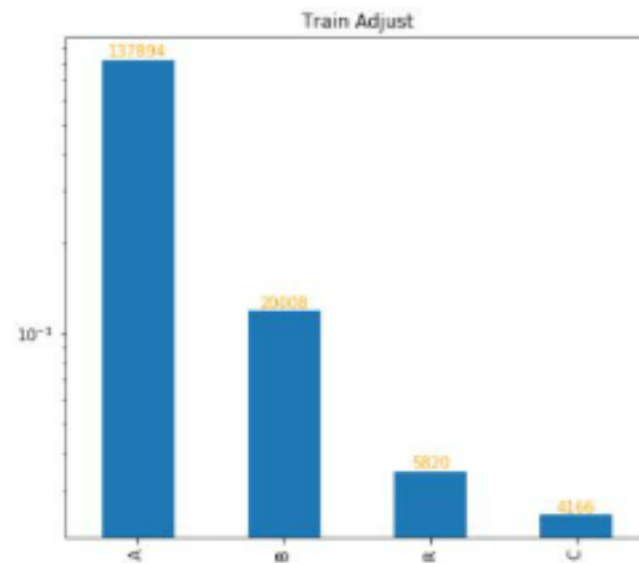
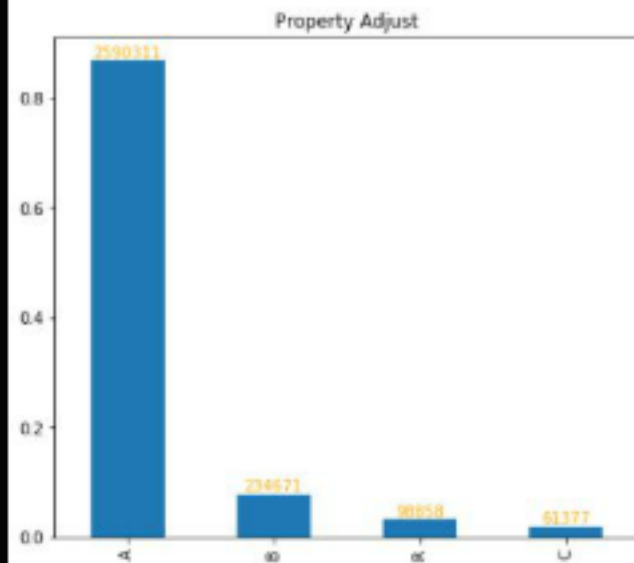


propertycountylandusecode

```
0    0100  
1    0109  
2    1200  
3    1200  
4    1210  
Name: propertycountylandusecode, dtype: object
```

```
0      D  
1    NaN  
2    NaN  
3    NaN  
4    NaN  
Name: propertycountylandusecode, dtype: object
```

(['X&R1': 'A', 'C': 'B', 'D': 'C', 'X': 'N'])



Cluster Geo Info

	regionidcity	regionidzip	rawcensustractandblock	censustractandblock	latitude	longitude	fips
0	37688.0	98337.0	60378004.0	6.059032e+13	34144440.0	-118654080.0	A
1	37688.0	98337.0	60378000.0	6.059032e+13	34140432.0	-118625360.0	A
2	51617.0	96095.0	60377032.0	6.059032e+13	33989360.0	-118384632.0	A
3	12447.0	98424.0	60371412.0	6.059032e+13	34148884.0	-118437208.0	A
4	12447.0	96450.0	60371232.0	6.059032e+13	34194168.0	-118385916.0	A

7 algorithms and each for 8 classes

```
cluster_geo_info = pd.DataFrame({'regionidcity': regionidcity, 'regionidzip': regionidzip, 'rawcensustractandblock': rawcensustractandblock, 'censustractandblock': censustractandblock, 'latitude': latitude, 'longitude': longitude, 'fips': fips})
cluster_geo_info['cluster_location'] = cluster_geo_info['rawcensustractandblock'].astype(int)
model = LinearSVC()
model.fit(cluster_geo_info[['rawcensustractandblock', 'latitude', 'longitude']])
projected = cluster_geo_info[['rawcensustractandblock', 'latitude', 'longitude']]
projected['cluster_location'] = model.predict(projected[['rawcensustractandblock', 'latitude', 'longitude']])
print(...)
```

```
model = LinearSVC()
model.fit(cluster_geo_info[['rawcensustractandblock', 'latitude', 'longitude']])
projected = cluster_geo_info[['rawcensustractandblock', 'latitude', 'longitude']]
projected['cluster_location'] = model.predict(projected[['rawcensustractandblock', 'latitude', 'longitude']])
print(...)
```

```
model = LinearSVC()
model.fit(cluster_geo_info[['rawcensustractandblock', 'latitude', 'longitude']])
projected = cluster_geo_info[['rawcensustractandblock', 'latitude', 'longitude']]
projected['cluster_location'] = model.predict(projected[['rawcensustractandblock', 'latitude', 'longitude']])
print(...)
```

```
model = LinearSVC()
model.fit(cluster_geo_info[['rawcensustractandblock', 'latitude', 'longitude']])
projected = cluster_geo_info[['rawcensustractandblock', 'latitude', 'longitude']]
projected['cluster_location'] = model.predict(projected[['rawcensustractandblock', 'latitude', 'longitude']])
print(...)
```

```
model = LinearSVC()
model.fit(cluster_geo_info[['rawcensustractandblock', 'latitude', 'longitude']])
projected = cluster_geo_info[['rawcensustractandblock', 'latitude', 'longitude']]
projected['cluster_location'] = model.predict(projected[['rawcensustractandblock', 'latitude', 'longitude']])
print(...)
```

```
model = LinearSVC()
model.fit(cluster_geo_info[['rawcensustractandblock', 'latitude', 'longitude']])
projected = cluster_geo_info[['rawcensustractandblock', 'latitude', 'longitude']]
projected['cluster_location'] = model.predict(projected[['rawcensustractandblock', 'latitude', 'longitude']])
print(...)
```

```
model = LinearSVC()
model.fit(cluster_geo_info[['rawcensustractandblock', 'latitude', 'longitude']])
projected = cluster_geo_info[['rawcensustractandblock', 'latitude', 'longitude']]
projected['cluster_location'] = model.predict(projected[['rawcensustractandblock', 'latitude', 'longitude']])
print(...)
```

```
projected['cluster_location'] = cluster_geo_info['cluster_location']
print(...)
```

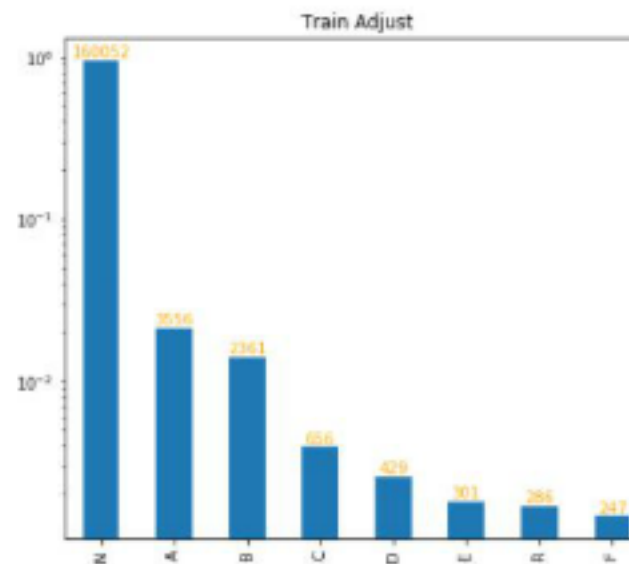
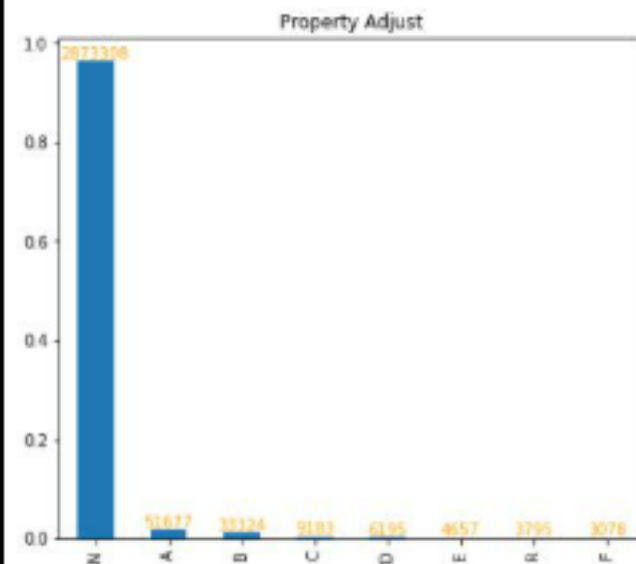
Date Features

	assessmentyear	yearbuilt
0	2016.0	1963.0
1	2015.0	1963.0
2	2016.0	1959.0
3	2016.0	1948.0
4	2016.0	1947.0

	assessmentyear	yearbuilt
0	1.0	62.0
1	2.0	62.0
2	1.0	58.0
3	1.0	69.0
4	1.0	70.0

taxdelinquencyyear

{10.0: 'F', 11.0: 'E', 12.0: 'D', 13.0: 'C', 14.0: 'B', 15.0: 'A', 'N': 'N'}



10% < Missing ratio < 30%

- 1.Preparation
- 2.Impute Data 1
- 3.Creating New Features
- 4.Impute Data 2

