

In [11]:

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import numpy as np
4 from sklearn.model_selection import train_test_split
```

In [12]:

```
1 df = pd.read_csv("fraudTest.csv")
2 df.head(5)
```

out[12]:

	ast	gender	street	...	lat	long	city_pop	job	dob		trans_num	unix_time	merch_lat	merch_long	is_fraud
	liott	M	351 Darlene Green	...	33.9659	-80.9355	333497	Mechanical engineer	1968-03-19	2da90c7d74bd46a0caf3777415b3ebd3	1371816865	33.986391	-81.200714		0
	ims	F	3638 Marsh Union	...	40.3207	-110.4360	302	Sales professional, IT	1990-01-17	324cc204407e99f51b0d6ca0055005e7	1371816873	39.450498	-109.960431		0
	pez	F	9333 Valentine Point	...	40.6729	-73.5365	34496	Librarian, public	1970-10-21	c81755dbbba9d5c77f094348a7579be	1371816893	40.495810	-74.196111		0
	ims	M	32941 Krystal Mill Apt. 552	...	28.5697	-80.8191	54767	Set designer	1987-07-25	2159175b9efe66dc301f149d3d5abf8c	1371816915	28.812398	-80.883061		0
	sey	M	5783 Evan Roads Apt. 465	...	44.2529	-85.0170	1126	Furniture designer	1955-07-06	57ff021bd3f328f8738bb535c302a31b	1371816917	44.959148	-85.884734		0

In [13]:

```
1 df.columns
```

Out[13]:

Index(['Unnamed: 0', 'trans_date_trans_time', 'cc_num', 'merchant', 'category', 'amt', 'first', 'last', 'gender', 'street', 'city', 'state', 'zip', 'lat', 'long', 'city_pop', 'job', 'dob', 'trans_num', 'unix_time', 'merch_lat', 'merch_long', 'is_fraud'], dtype='object')

In [47]:

```
1 dfcol = df[['amt', 'city_pop', 'unix_time', 'category']]
2 dfcol = pd.get_dummies(dfcol, columns=['category'])
3 dfcol.head()
```

Out[47]:

	category_grocery_net	category_grocery_pos	category_health_fitness	category_home	category_kids_pets	category_misc_net	category_misc_pos	category_pers
	0	0	0	0	0	0	0	
	0	0	0	0	0	0	0	
	0	0	1	0	0	0	0	
	0	0	0	0	0	0	1	
	0	0	0	0	0	0	0	

In [48]:

```
1 dfpred = df[['is_fraud']]
2 dfpred
```

Out[48]:

	is_fraud
0	0
1	0
2	0
3	0
4	0
...	...
555714	0
555715	0
555716	0
555717	0
555718	0
555719 rows × 1 columns	

```
In [49]: 1 X = dfcol
        2 y = dfpred
```

1. Use scikit learn to split data into 70/30

```
In [50]: 1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3)
```

2. Make 3 sets of training data (Oversample, Undersample and SMOTE)

```
In [51]: 1 from imblearn.over_sampling import RandomOverSampler, SMOTE
        2 from imblearn.under_sampling import RandomUnderSampler
        3
        4
        5 # Oversampled training data
        6 oversampler = RandomOverSampler()
        7 X_train_over, y_train_over = oversampler.fit_resample(X_train, y_train)
        8
        9 # Undersampled training data
       10 undersampler = RandomUnderSampler()
       11 X_train_under, y_train_under = undersampler.fit_resample(X_train, y_train)
       12
       13 # SMOTE training data
       14 smote = SMOTE()
       15 X_train_smote, y_train_smote = smote.fit_resample(X_train, y_train)
       16
```

3. Train 3 logistic regression models

```
In [55]: 1 from sklearn.linear_model import LogisticRegression
        2
        3 # Initialize logistic regression classifier
        4 log_reg = LogisticRegression()
        5
        6 # Fit the model using the oversampled training data
        7 log_reg_over = log_reg.fit(X_train_over, y_train_over)
        8
        9 # Fit the model using the undersampled training data
       10 log_reg_under = log_reg.fit(X_train_under, y_train_under)
       11
       12 # Fit the model using the SMOTE training data
       13 log_reg_smote = log_reg.fit(X_train_smote, y_train_smote)
       14
```

C:\Users\parzu\anaconda3\lib\site-packages\sklearn\utils\validation.py:993: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().

y = column_or_1d(y, warn=True)

C:\Users\parzu\anaconda3\lib\site-packages\sklearn\utils\validation.py:993: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().

y = column_or_1d(y, warn=True)

C:\Users\parzu\anaconda3\lib\site-packages\sklearn\utils\validation.py:993: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().

y = column_or_1d(y, warn=True)

4. Test the 3 models

```

In [56]: 1 from sklearn.metrics import accuracy_score, f1_score, recall_score, precision_score
2
3 # Make predictions on the test data using the oversampled model
4 y_pred_over = log_reg_over.predict(X_test)
5
6 # Make predictions on the test data using the undersampled model
7 y_pred_under = log_reg_under.predict(X_test)
8
9 # Make predictions on the test data using the SMOTE model
10 y_pred_smote = log_reg_smote.predict(X_test)
11
12 # Evaluate the oversampled model
13 print("Oversampled model evaluation:")
14 print("Accuracy:", accuracy_score(y_test, y_pred_over))
15 print("F1-Score:", f1_score(y_test, y_pred_over))
16 print("Recall:", recall_score(y_test, y_pred_over))
17 print("Precision:", precision_score(y_test, y_pred_over))
18
19 # Evaluate the undersampled model
20 print("\nUndersampled model evaluation:")
21 print("Accuracy:", accuracy_score(y_test, y_pred_under))
22 print("F1-Score:", f1_score(y_test, y_pred_under))
23 print("Recall:", recall_score(y_test, y_pred_under))
24 print("Precision:", precision_score(y_test, y_pred_under))
25
26 # Evaluate the SMOTE model
27 print("\nSMOTE model evaluation:")
28 print("Accuracy:", accuracy_score(y_test, y_pred_smote))
29 print("F1-Score:", f1_score(y_test, y_pred_smote))
30 print("Recall:", recall_score(y_test, y_pred_smote))
31 print("Precision:", precision_score(y_test, y_pred_smote))
32

```

Oversampled model evaluation:
 Accuracy: 0.16996569015571392
 F1-Score: 0.008455144740613357
 Recall: 0.8575581395348837
 Precision: 0.004248516619620946

Undersampled model evaluation:
 Accuracy: 0.16996569015571392
 F1-Score: 0.008455144740613357
 Recall: 0.8575581395348837
 Precision: 0.004248516619620946

SMOTE model evaluation:
 Accuracy: 0.16996569015571392
 F1-Score: 0.008455144740613357
 Recall: 0.8575581395348837
 Precision: 0.004248516619620946

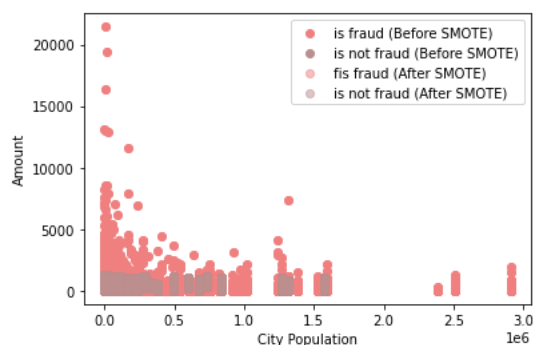
5. The accuracy is the same for all three models which is unusual, suggesting that the model may have issues.

6. Pick 2 features and plot the 2 classes before and after SMOTE

```

In [69]: 1 import matplotlib.pyplot as plt
2 from imblearn.over_sampling import SMOTE
3
4 # Get the data
5 X = df[['city_pop', 'amt']]
6 y = df['is_fraud']
7
8 # Split the data into training and test sets
9 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=0)
10
11 # Plot the original data
12 plt.scatter(X_train[y_train == 0]['city_pop'], X_train[y_train == 0]['amt'], color='lightcoral', label='is fraud (Before SMOTE)')
13 plt.scatter(X_train[y_train == 1]['city_pop'], X_train[y_train == 1]['amt'], color='rosybrown', label='is not fraud (Before SMOTE)')
14
15 # Apply SMOTE to the training data
16 smote = SMOTE(sampling_strategy='minority')
17 X_train_smote, y_train_smote = smote.fit_resample(X_train, y_train)
18
19 # Plot the SMOTE data
20 plt.scatter(X_train_smote[y_train_smote == 0]['city_pop'], X_train_smote[y_train_smote == 0]['amt'], color='lightcoral', label='is fraud (After SMOTE)')
21 plt.scatter(X_train_smote[y_train_smote == 1]['city_pop'], X_train_smote[y_train_smote == 1]['amt'], color='rosybrown', label='is not fraud (After SMOTE)')
22
23 plt.xlabel('City Population')
24 plt.ylabel('Amount')
25 plt.legend()
26 plt.show()
27

```



```

In [ ]: 1

```