# ml_project

March 2, 2020

```
In [38]: import pandas as pd
         import numpy as np
         import scipy.linalg as la
         import scipy.stats as stats
         import statsmodels.api as sm
         from sklearn import linear_model, model_selection, metrics
         import sklearn
         import plotly.graph_objs as go
         import matplotlib.pyplot as plt
         from sklearn.datasets import load_breast_cancer, load_digits, load_diabetes
         from sklearn.model_selection import train_test_split
         from sklearn.discriminant_analysis import LinearDiscriminantAnalysis as LDA
         from sklearn.neighbors import KNeighborsClassifier as KNC
         from sklearn.discriminant_analysis import QuadraticDiscriminantAnalysis as GDA
         from sklearn.metrics import roc_auc_score as RAS
         from sklearn.metrics import roc_curve as ROC
         import seaborn as sns
         from sklearn.manifold import TSNE
         import umap
```

   jupyter nbconvert --to latex "Race and Incarceration in America.ipynb"
--TagRemovePreprocessor.remove_input_tags='{"hide_input"}'; pdflatex "Race and
Incarceration in America"

## 1 Introduction

Last semester I took a look a data set containing the information of more than 7.5 million individuals that have been processed by the criminal justice system. I found that racial minorities were more likely to receive extreme sentences, agreeing with existing research around bias in the criminal justice system. In this project I will be exploring the data from a machine learning perspective. My goal is to determine if this data can be classified in such a way that is predictive of race.

## 2 Possible questions

can we use ml techniques to correctly classify this data? which ones fail and why? can we create a predictive model for sentence lengths? should sentencing be offloaded to a ml algorithm.

```
In [14]: chunksize = 100000
         rdf = pd.read_csv('regression_df.csv', chunksize=chunksize)

In [39]: # reducer = umap.UMAP()
         # i = 0
         # for item in rdf:
         #     if i > 3:
         #         break
         #     else:
         #         i += 1

         #     data = item.drop(['Unnamed: 0','IDENTIFIER','LATEST ADMISSION DATE'], axis=1)
         #     target = data[['AMER IND','ASIAN','BLACK','HISPANIC']]
         #     data = data.drop(['AMER IND','ASIAN','BLACK','HISPANIC'], axis=1)
         #     trans = reducer.fit_transform(data)

         # #     create maskings
         #     black = target.BLACK == 1
         #     asian = target.ASIAN == 1
         #     hisp = target.HISPANIC == 1
         #     ind = target['AMER IND'] == 1
         #     white = ~(black + asian + hisp)

         #     plt.scatter(trans[:,0][black],trans[:,1][black],label='black',marker='.')
         # #     plt.legend()
         # #     plt.show()
         #     plt.scatter(trans[:,0][white],trans[:,1][white],label='white',marker='.')
         # #     plt.legend()
         # #     plt.show()
         #     plt.scatter(trans[:,0][hisp],trans[:,1][hisp],label='hisp',marker='.')
         # #     plt.legend()
         # #     plt.show()
         #     plt.scatter(trans[:,0][ind],trans[:,1][ind],label='ind',marker='.')
         # #     plt.legend()
         # #     plt.show()
         #     plt.scatter(trans[:,0][asian],trans[:,1][asian],label='asian',marker='.')
         #     plt.legend()
         #     plt.show()

In [29]: cols = ['RACE','GENDER','AGE','OFFENSE','FACILITY','DETAINER','SENTENCE DAYS']
         df = pd.read_csv('individuals.csv',usecols=cols)

In [45]: samp = df.sample(20000)
         samp.RACE = pd.factorize(samp['RACE'])[0] + 1
         samp.GENDER = pd.factorize(samp['GENDER'])[0] + 1
         samp.OFFENSE = pd.factorize(samp['OFFENSE'])[0] + 1
         samp.DETAINER = pd.factorize(samp['DETAINER'])[0] + 1
         samp.FACILITY = pd.factorize(samp['FACILITY'])[0] + 1
```
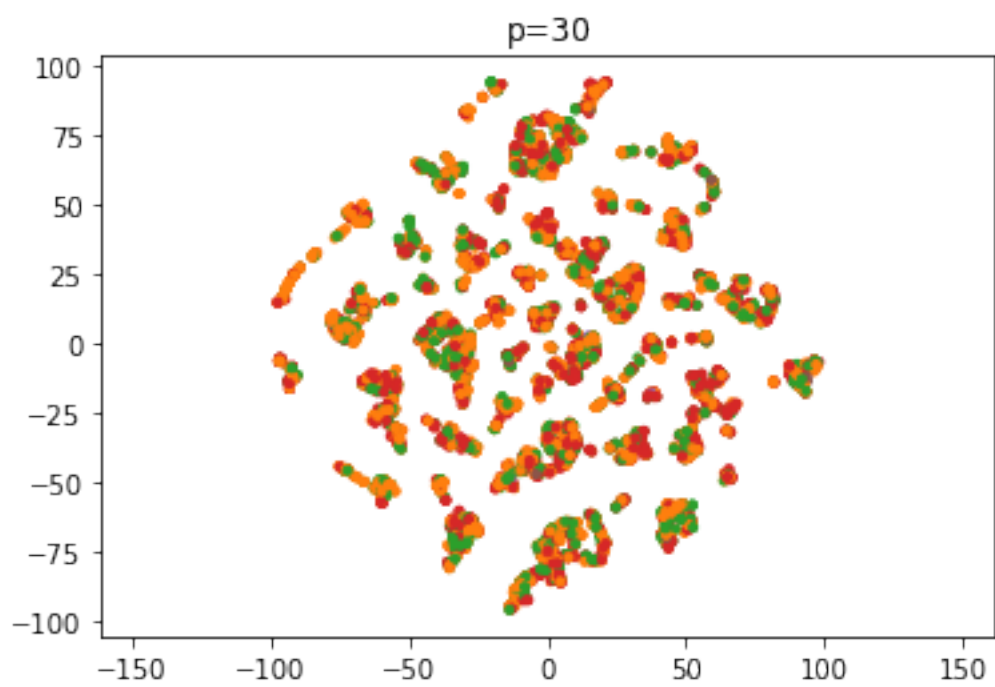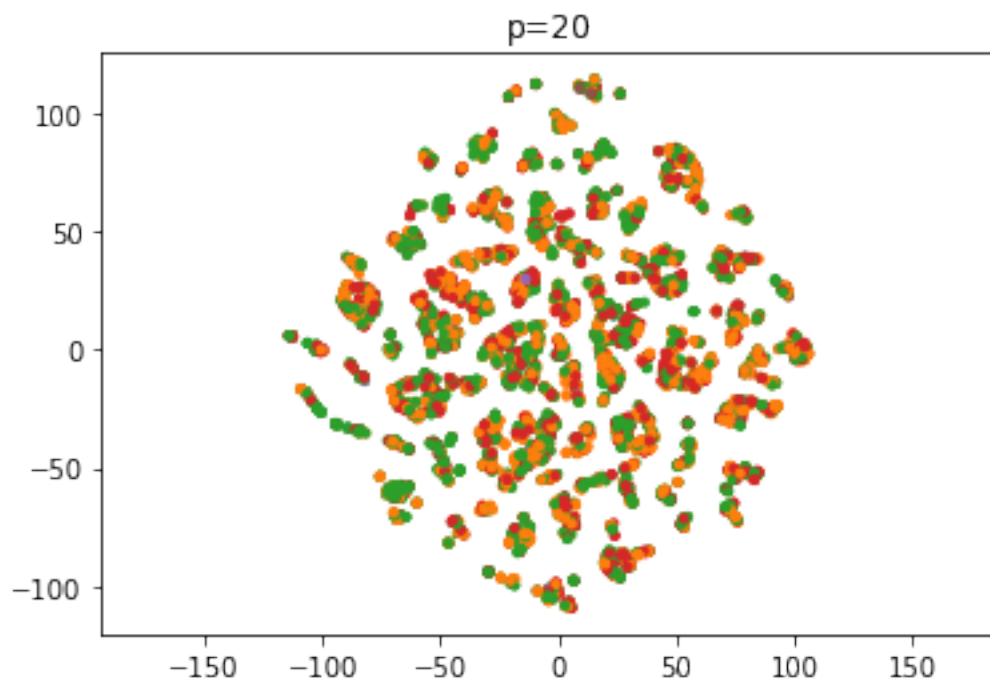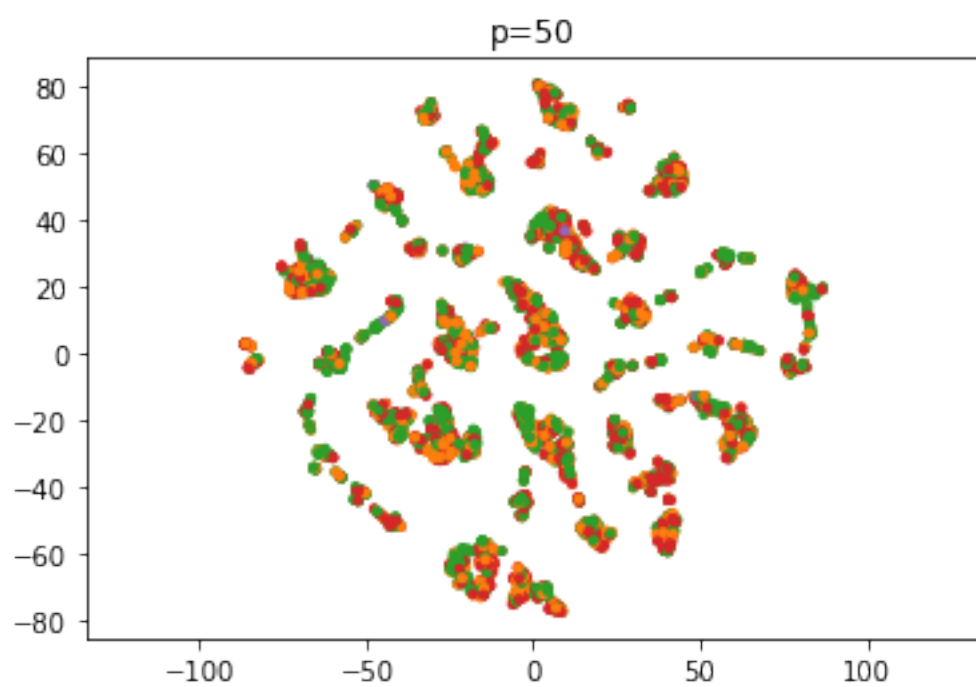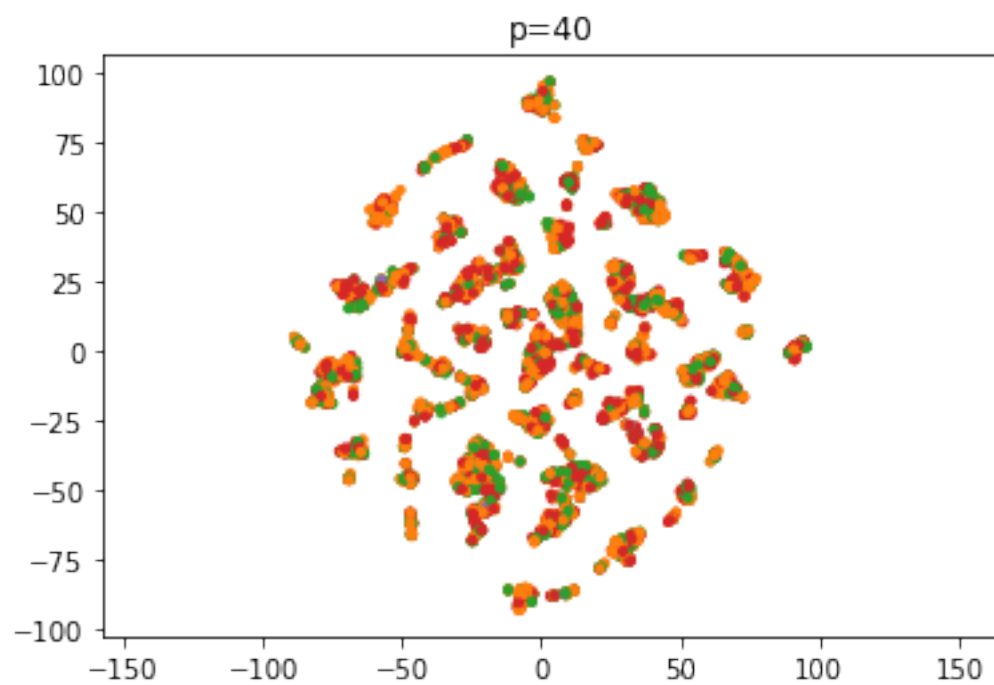
```python
        samp_y = samp.RACE
        samp_X = samp[['GENDER','AGE','OFFENSE','FACILITY','DETAINER',''SENTENCE DAYS']]

        trans = reducer.fit_transform(samp_X)

        plt.scatter(trans[:, 0], trans[:, 1], c=[sns.color_palette()[x] for x in samp_y], mar
        plt.gca().set_aspect('equal', 'datalim')
        plt.show()
```
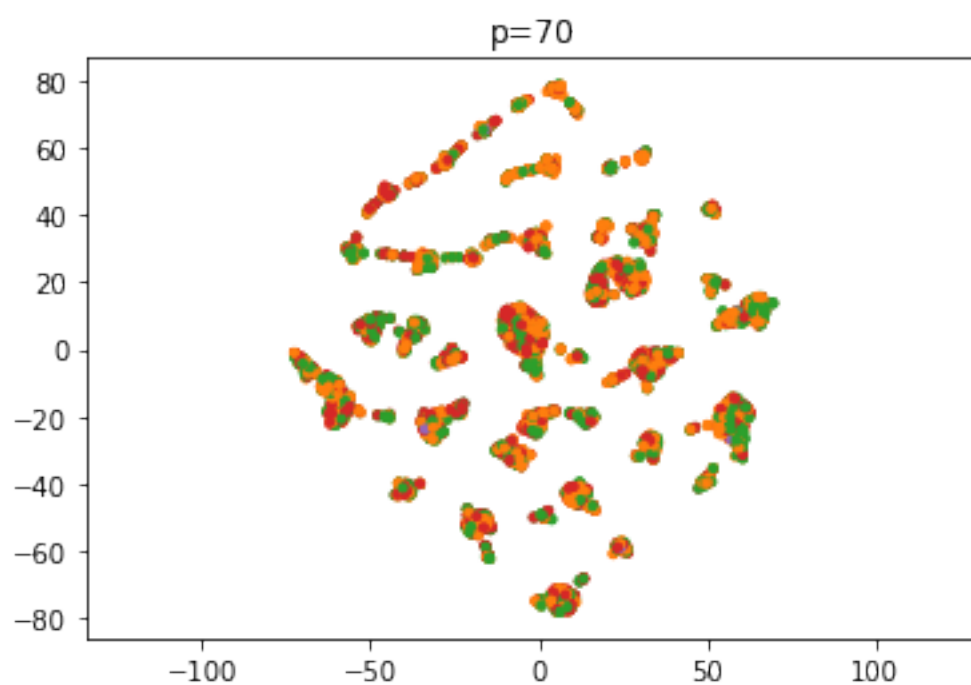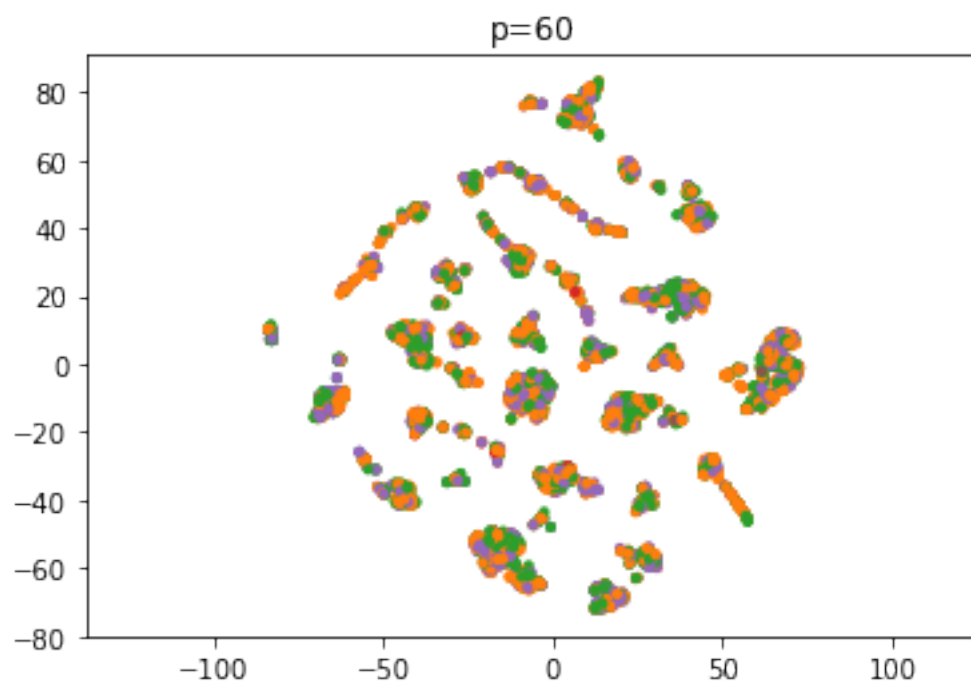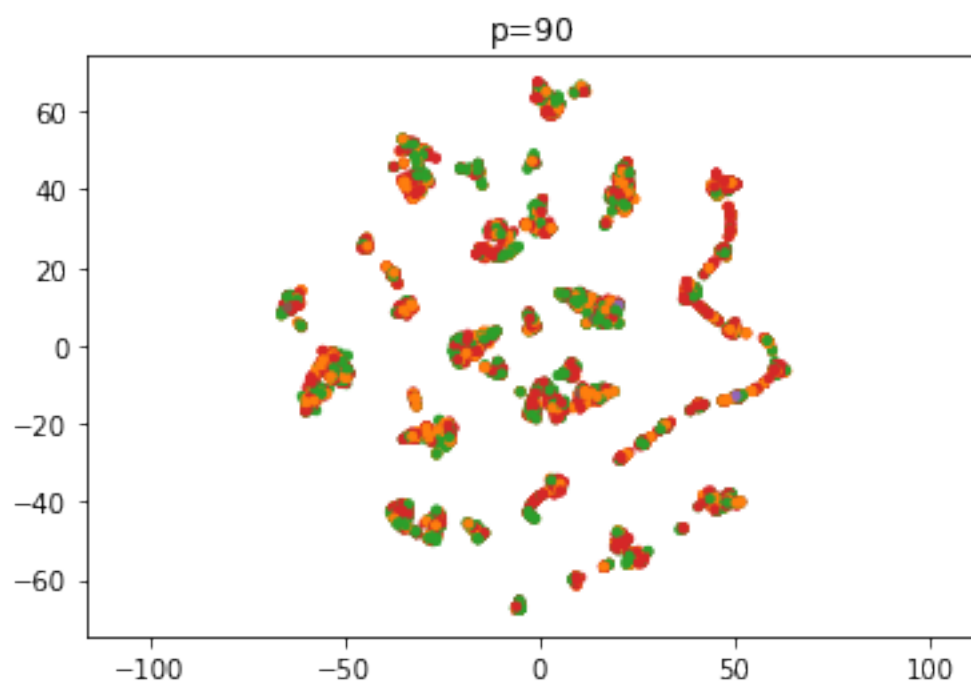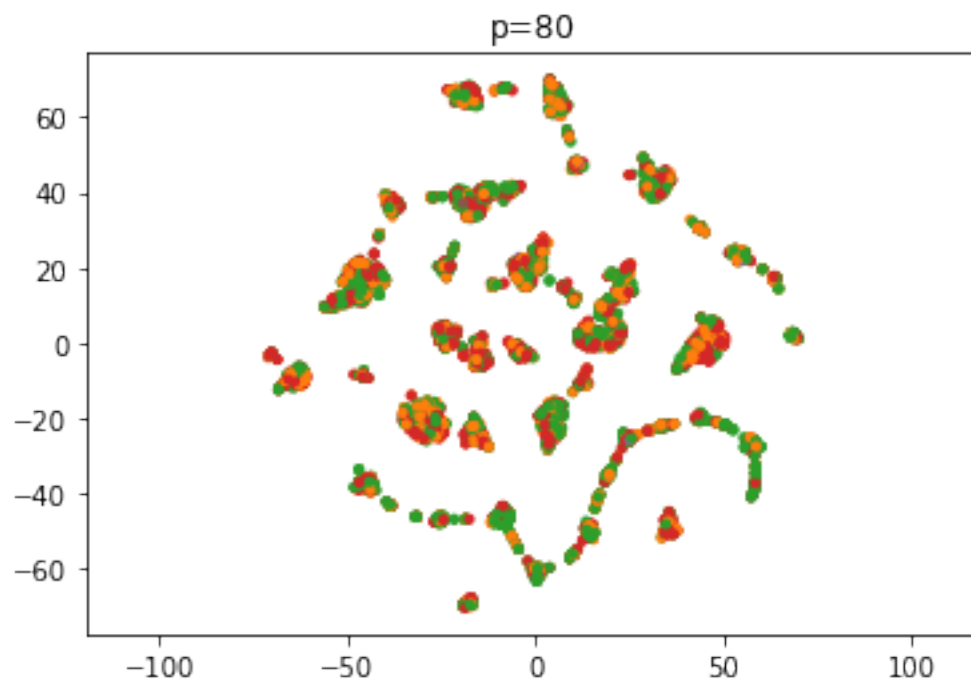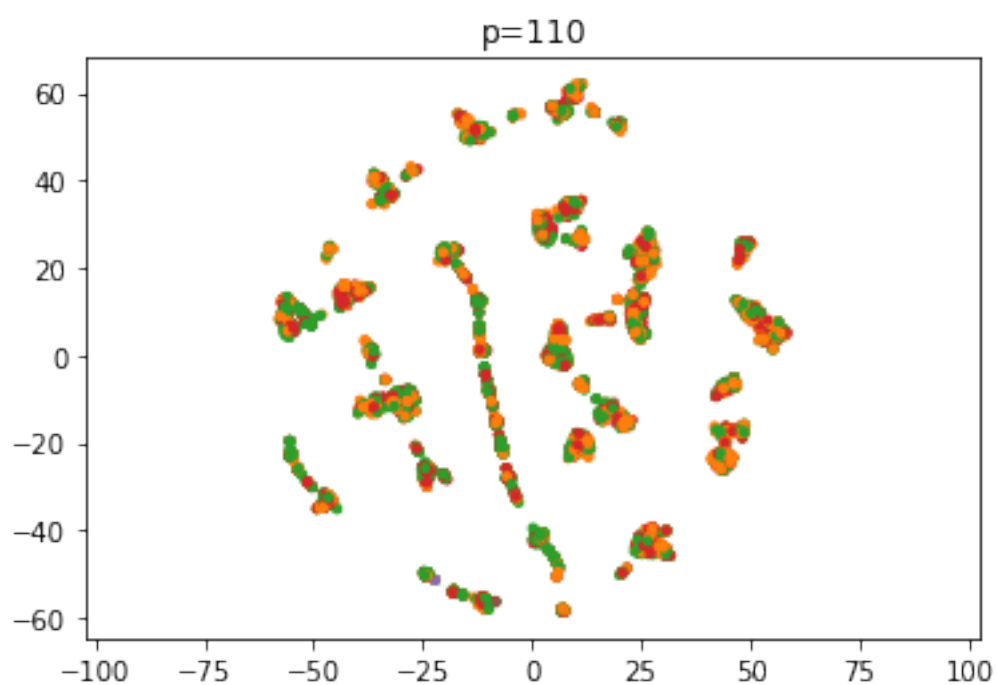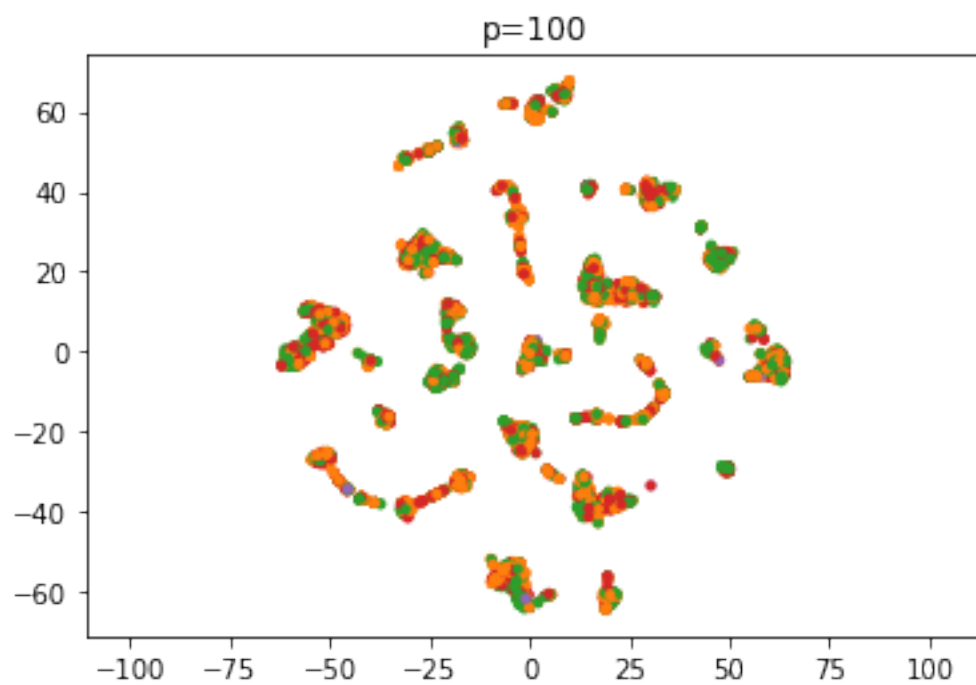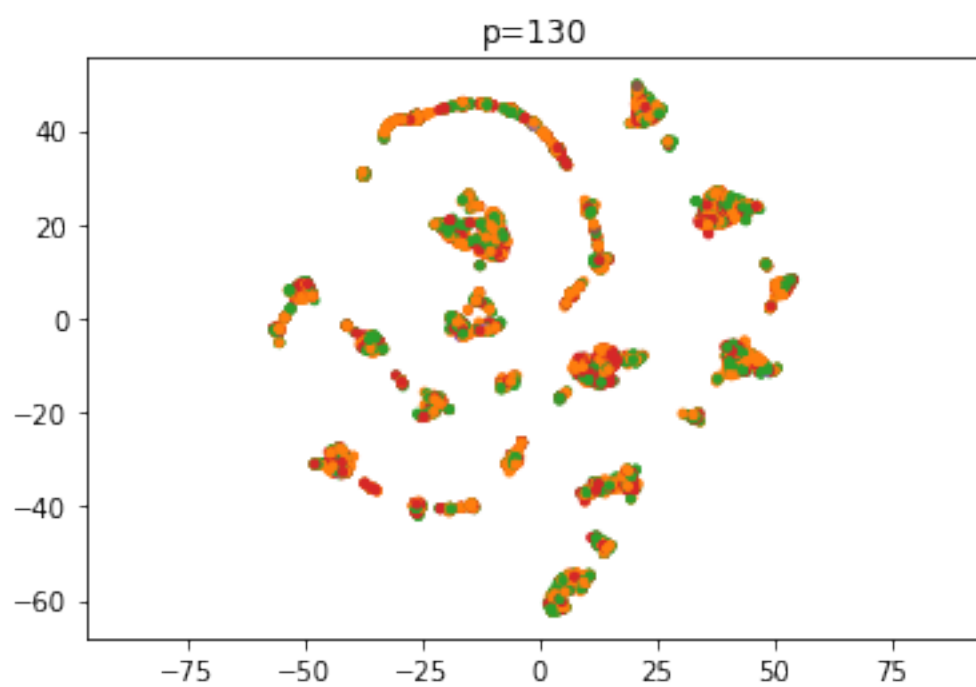
/home/ethan/.local/lib/python3.7/site-packages/numba/compiler.py:602: NumbaPerformanceWarning:

The keyword argument 'parallel=True' was specified but no transformation for parallel execution

To find out why, try turning on parallel diagnostics, see http://numba.pydata.org/numba-doc/la

File "../../../.local/lib/python3.7/site-packages/umap/nndescent.py", line 47:
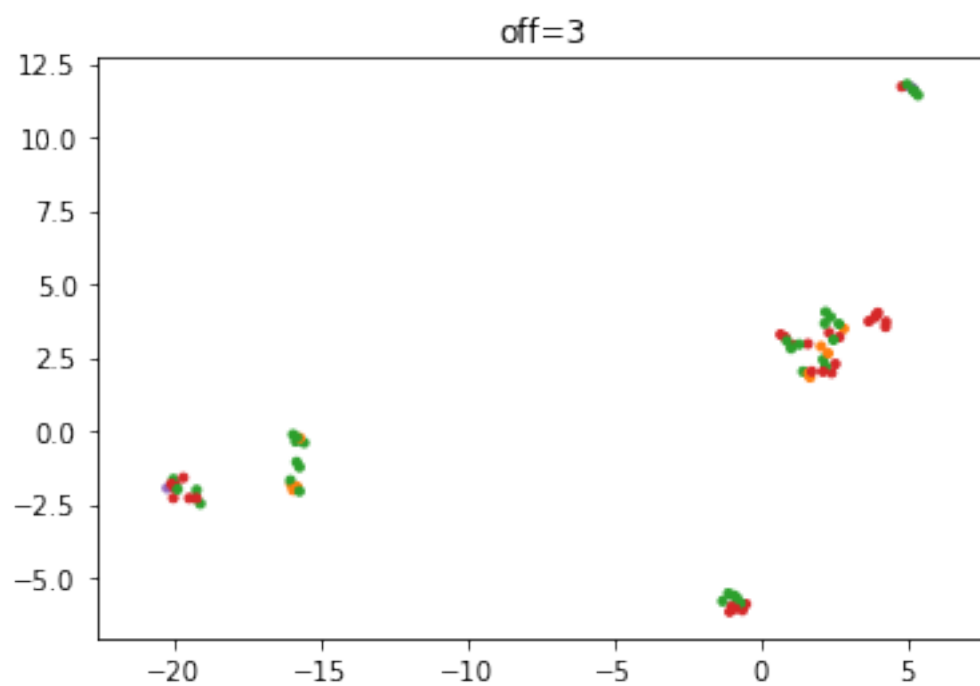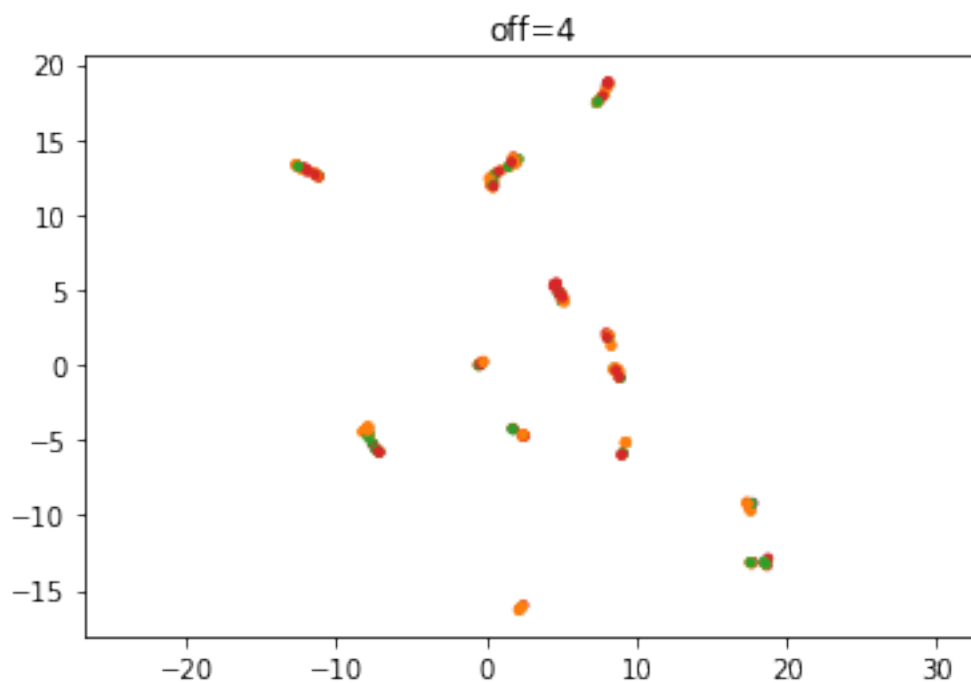    @numba.njit(parallel=True)
    def nn_descent(
    ^

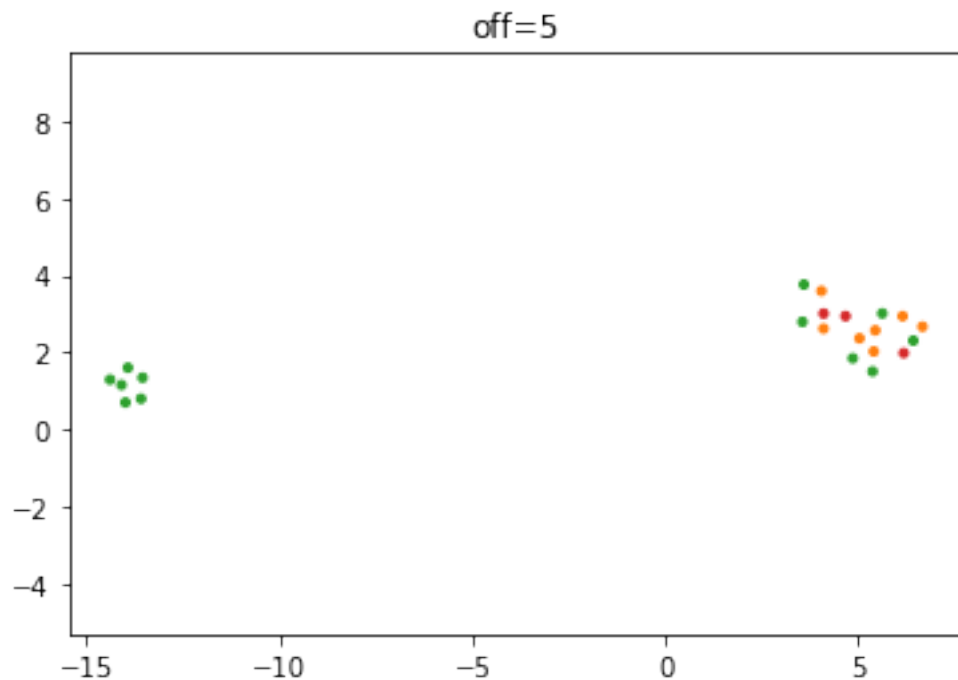/home/ethan/.local/lib/python3.7/site-packages/umap/spectral.py:229: UserWarning:

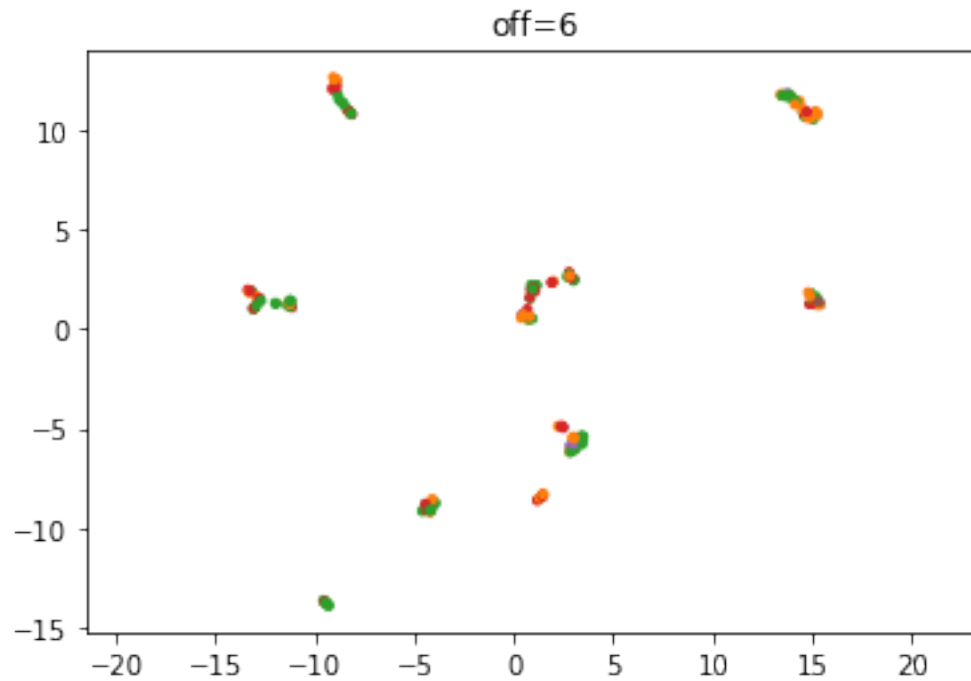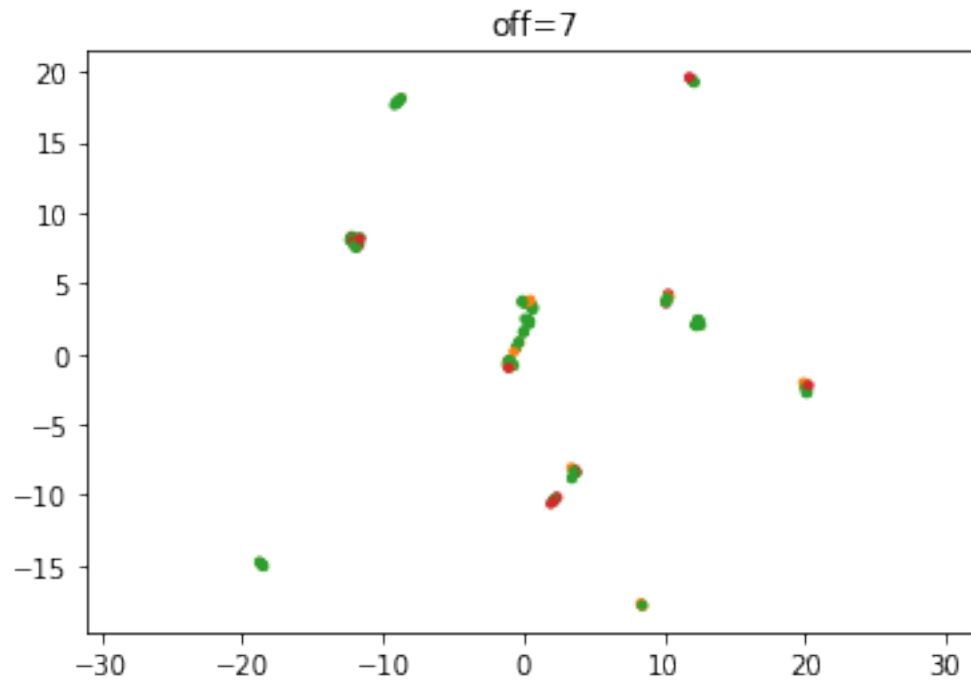Embedding a total of 50 separate connected components using meta-embedding (experimental)

/home/ethan/.local/lib/python3.7/site-packages/sklearn/manifold/spectral_embedding_.py:235: Use

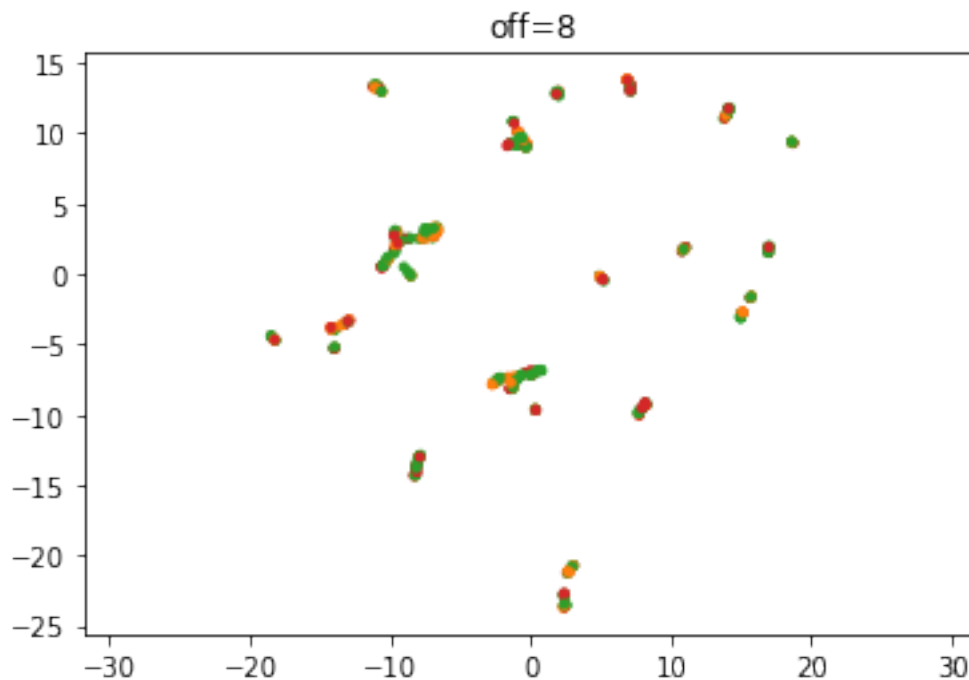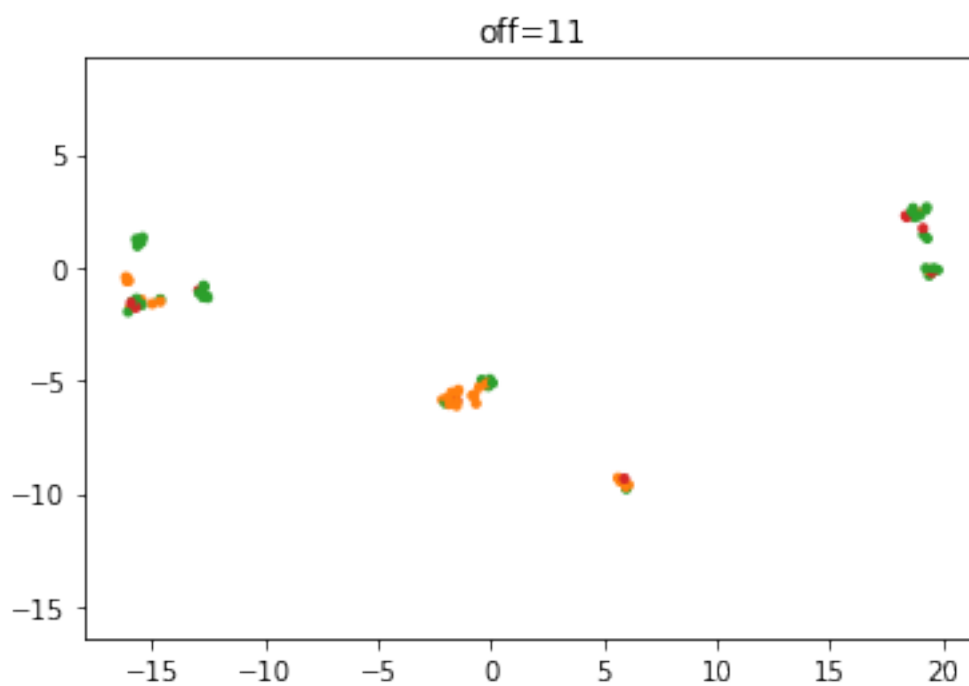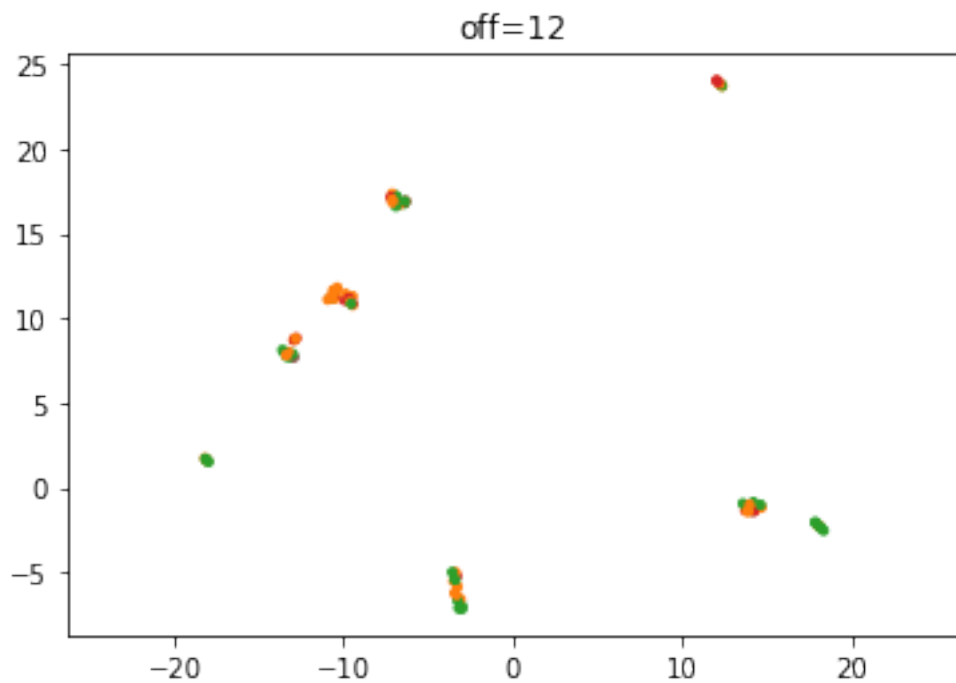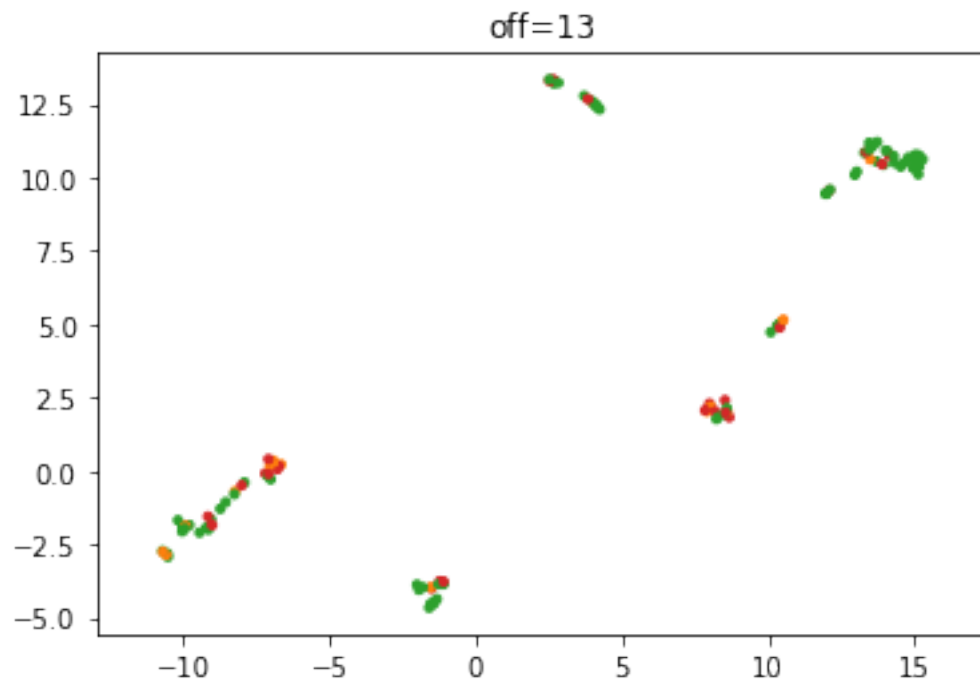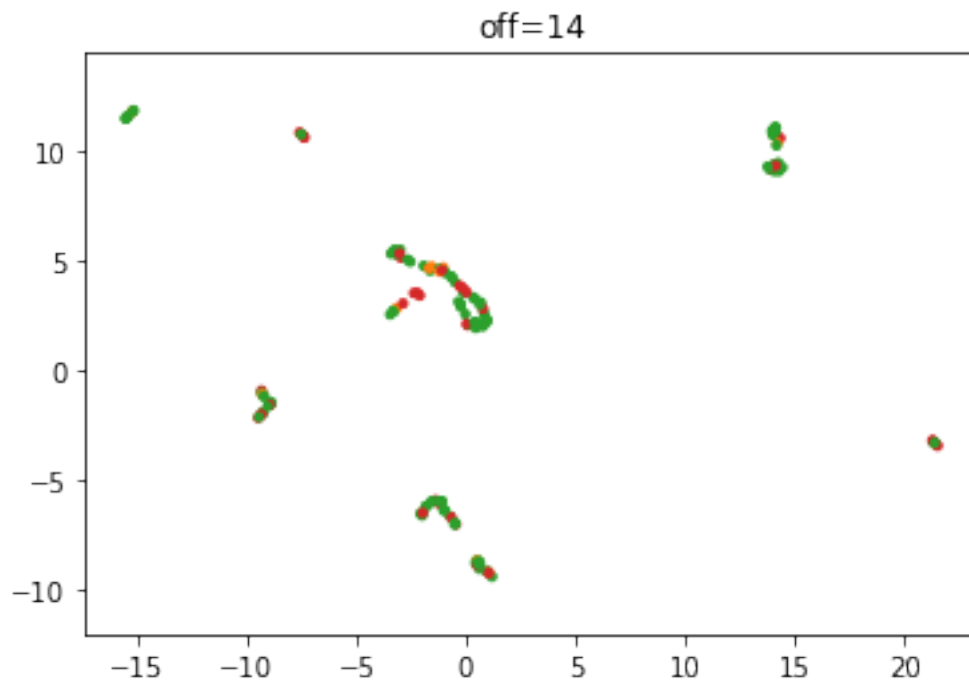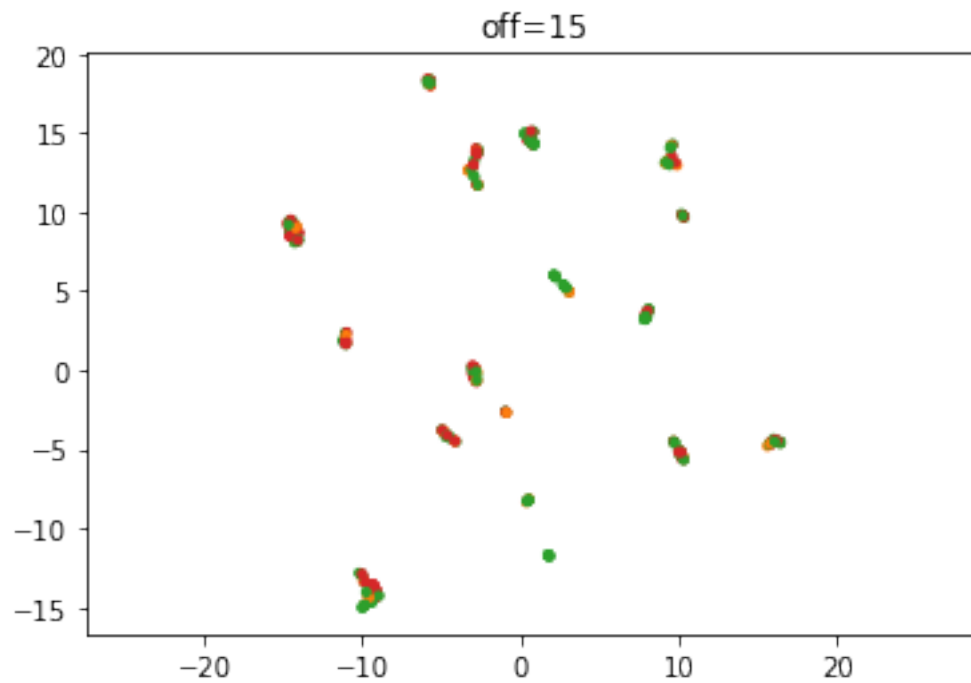Graph is not fully connected, spectral embedding may not work as expected.

It seems that UMAP doesnt work

```
In [47]: # params
         ps = [5,10,20,30,40,50,60,70,80,90,100,110,120,130,140,150,160,170,180,190,200]

         for p in ps:

             samp = df.sample(8000)
             samp.RACE = pd.factorize(samp['RACE'])[0] + 1
             samp.GENDER = pd.factorize(samp['GENDER'])[0] + 1
             samp.OFFENSE = pd.factorize(samp['OFFENSE'])[0] + 1
             samp.DETAINER = pd.factorize(samp['DETAINER'])[0] + 1
             samp.FACILITY = pd.factorize(samp['FACILITY'])[0] + 1
             samp_y = samp.RACE
             samp_X = samp[['GENDER','AGE','OFFENSE','FACILITY','DETAINER','SENTENCE DAYS']]


             trans = TSNE(perplexity=p).fit_transform(samp_X)

             plt.scatter(trans[:, 0], trans[:, 1], c=[sns.color_palette()[x] for x in samp_y],
             plt.gca().set_aspect('equal', 'datalim')
             plt.title(f'p={p}')
             plt.show()
```

p=5



p=10

## p=20



## p=30

p=40



p=50

## p=60



## p=70

p=80

p=90

## p=100



## p=110

p=120

p=130

## p=140



## p=150

p=160



p=170

p=180



p=190

p=200

```
In [70]: for off in set(samp.OFFENSE):

            mask = samp_X.OFFENSE == off

            data = samp_X[mask]

    #       trans = TSNE(perplexity=20).fit_transform(data)
            try:
                trans = reducer.fit_transform(data.values)

                plt.scatter(trans[:, 0], trans[:, 1], c=[sns.color_palette()[x] for x in samp_
                plt.gca().set_aspect('equal', 'datalim')
                plt.title(f'off={off}')
                plt.show()
            except:

                continue
```

/home/ethan/.local/lib/python3.7/site-packages/umap/spectral.py:229: UserWarning:

Embedding a total of 5 separate connected components using meta-embedding (experimental)

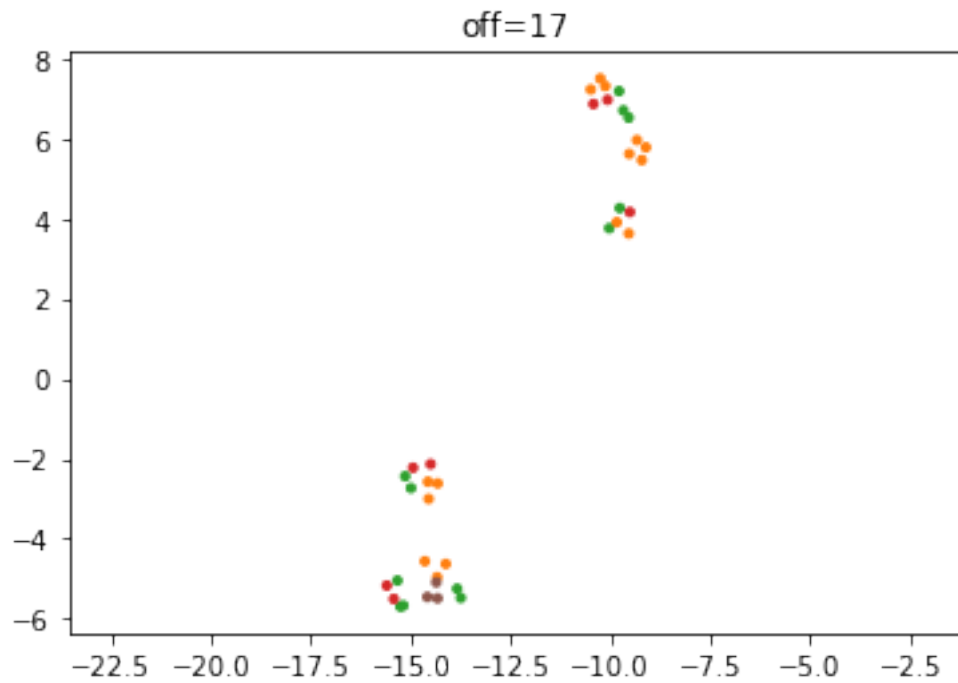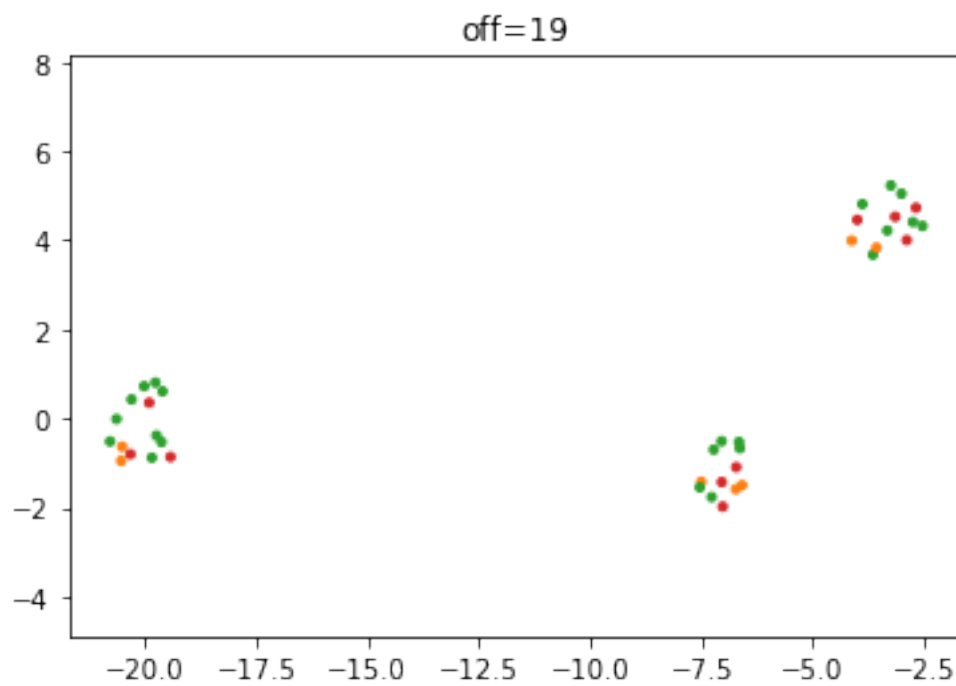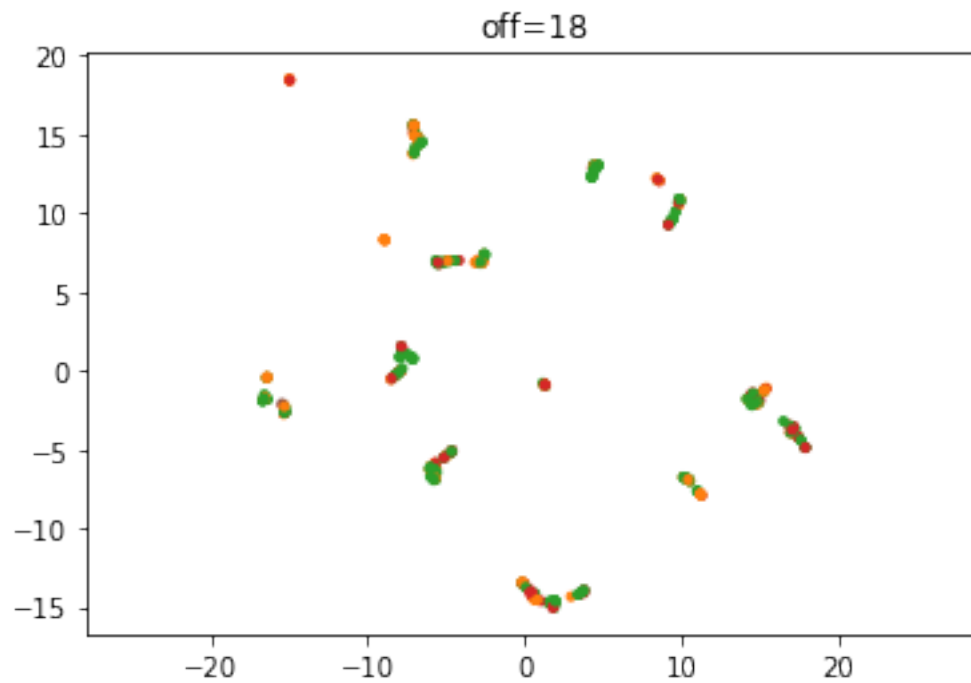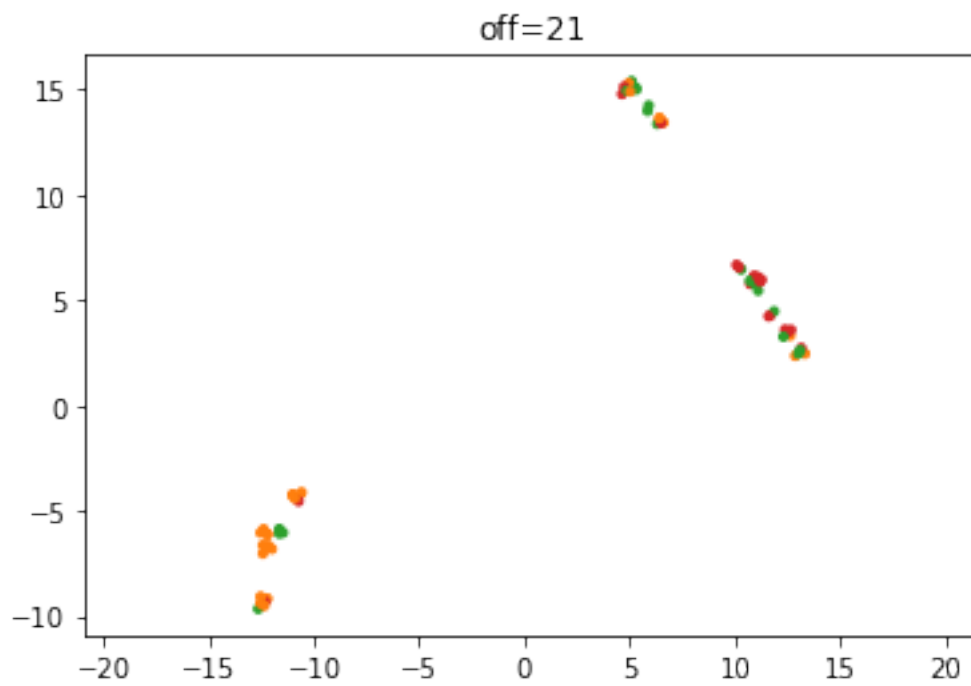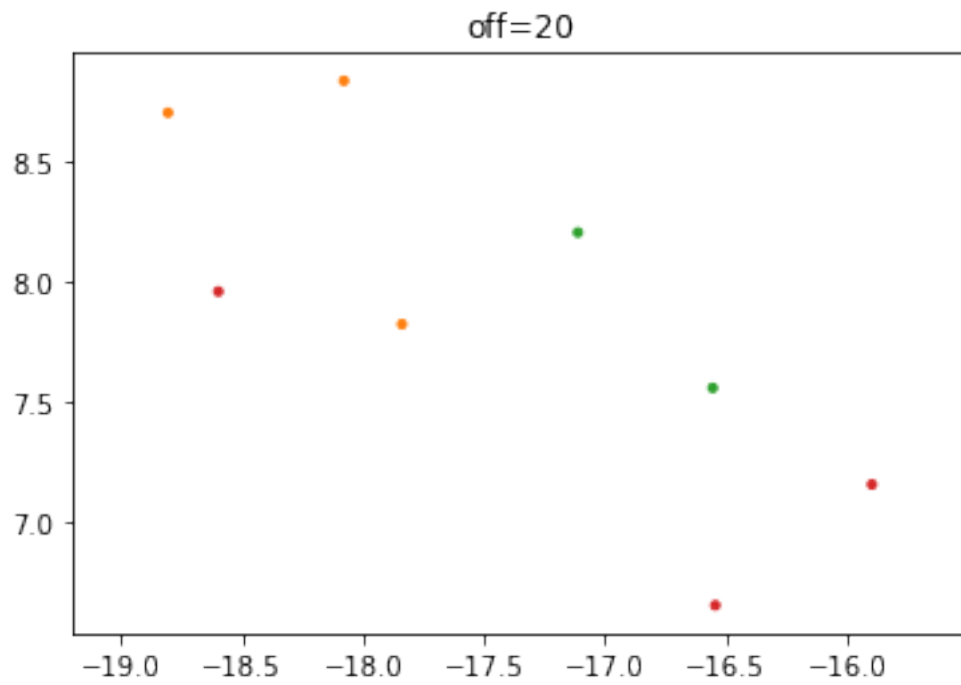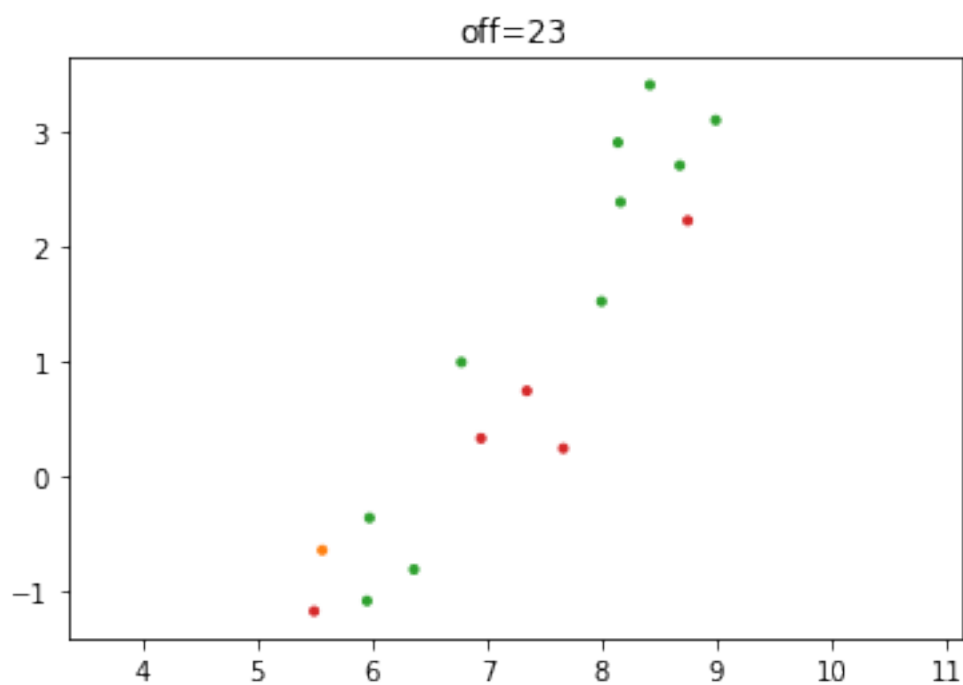/home/ethan/.local/lib/python3.7/site-packages/sklearn/manifold/spectral_embedding_.py:235: Use
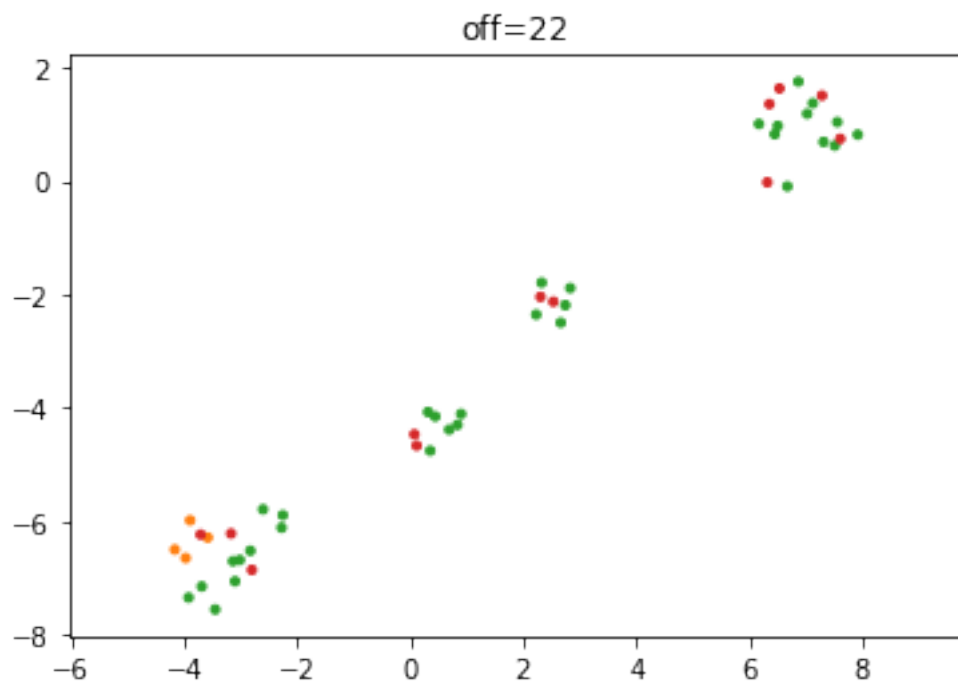
Graph is not fully connected, spectral embedding may not work as expected.



off=1

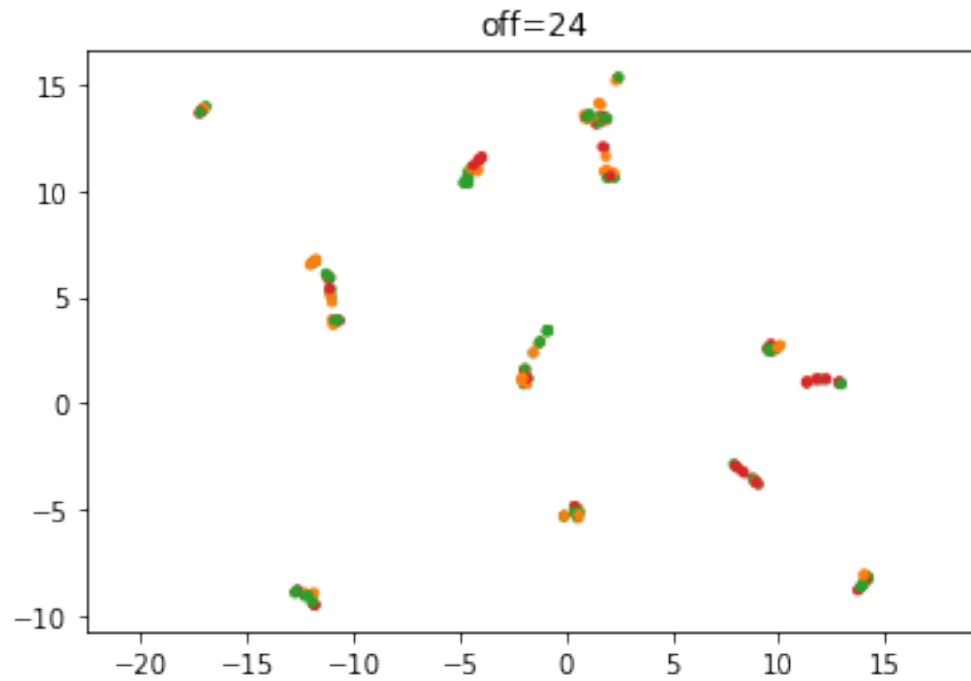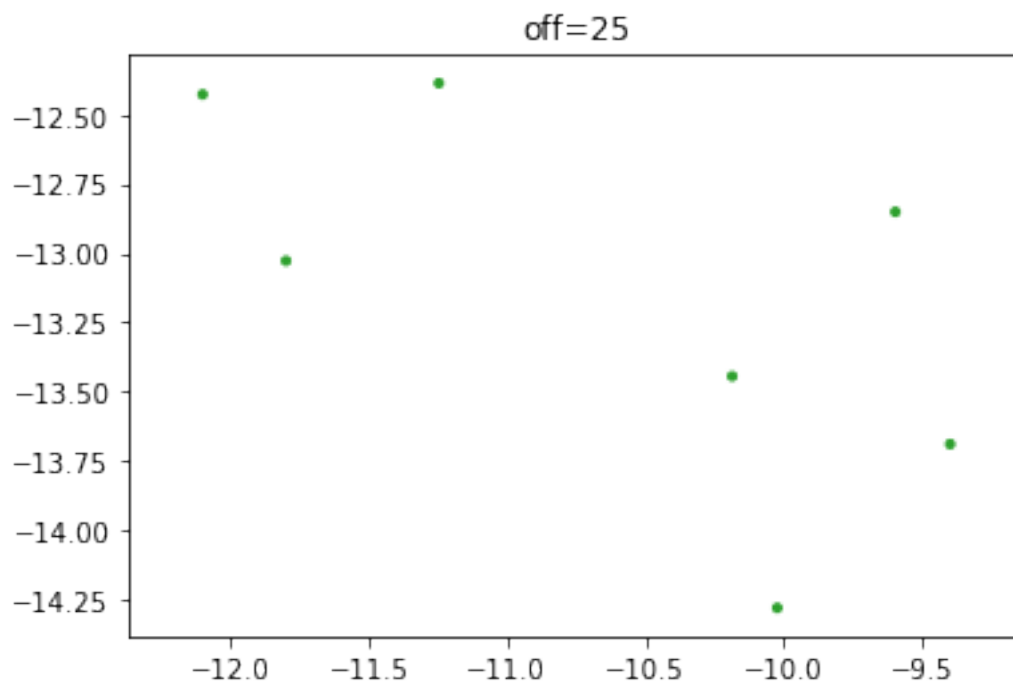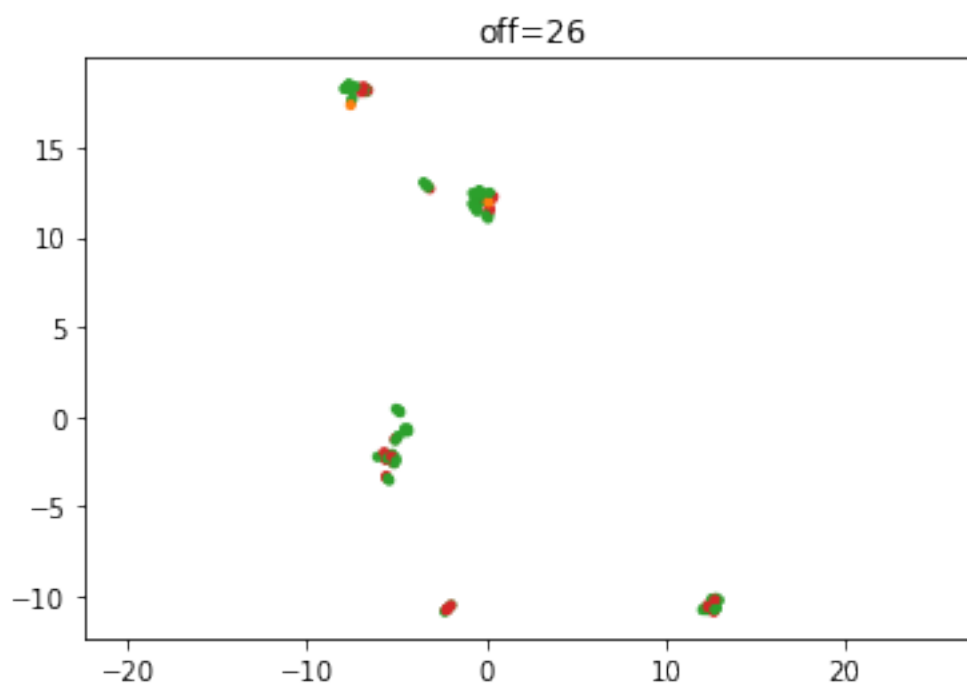/home/ethan/.local/lib/python3.7/site-packages/umap/spectral.py:229: UserWarning:

Embedding a total of 4 separate connected components using meta-embedding (experimental)

off=2

/home/ethan/.local/lib/python3.7/site-packages/umap/spectral.py:229: UserWarning:

Embedding a total of 2 separate connected components using meta-embedding (experimental)
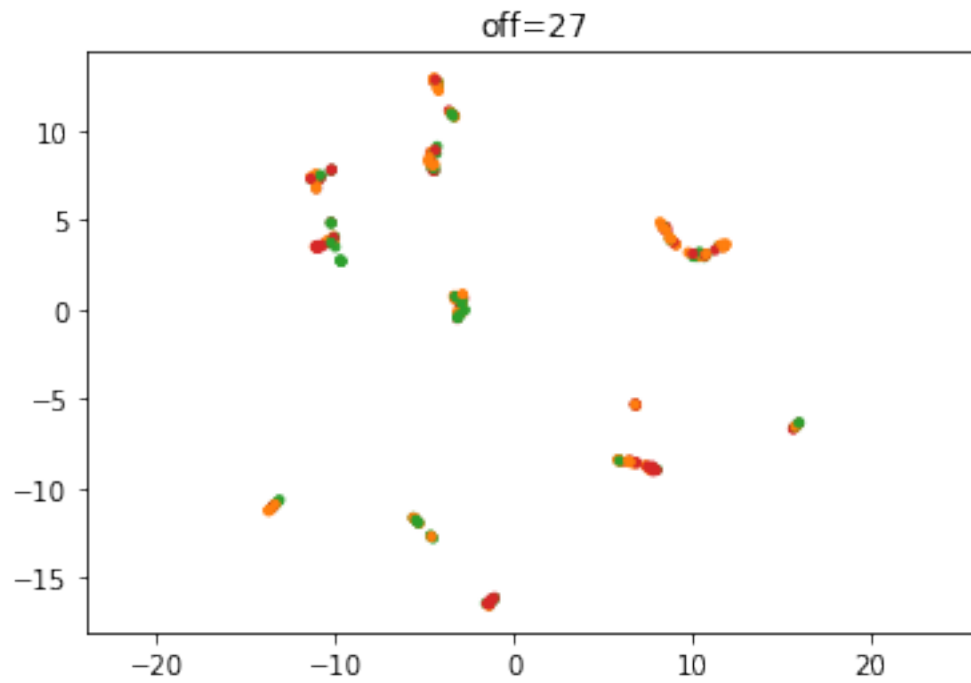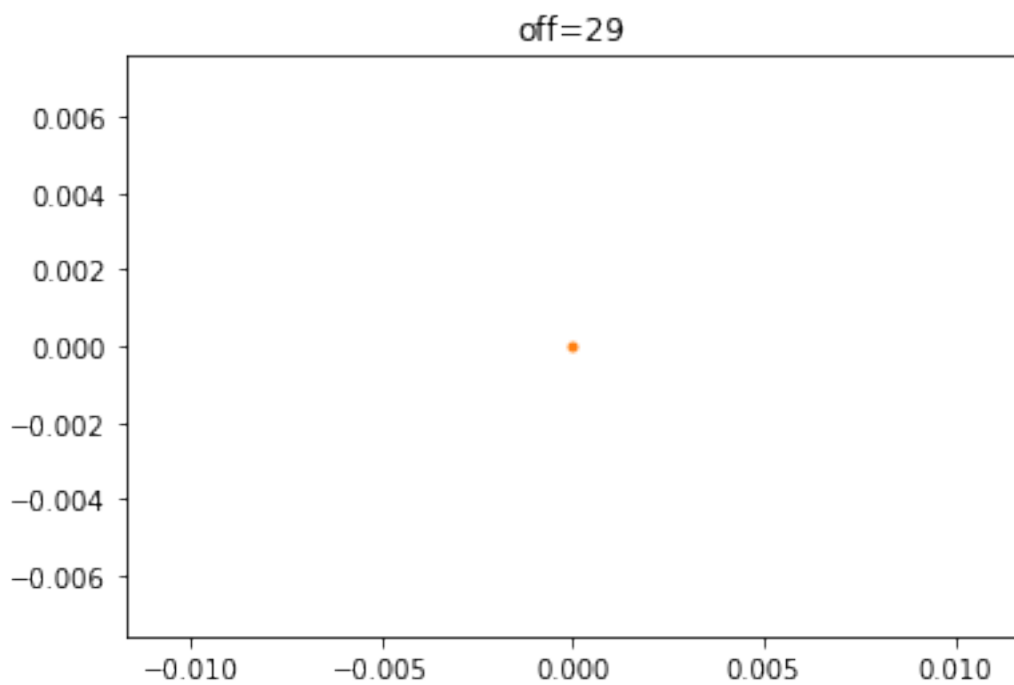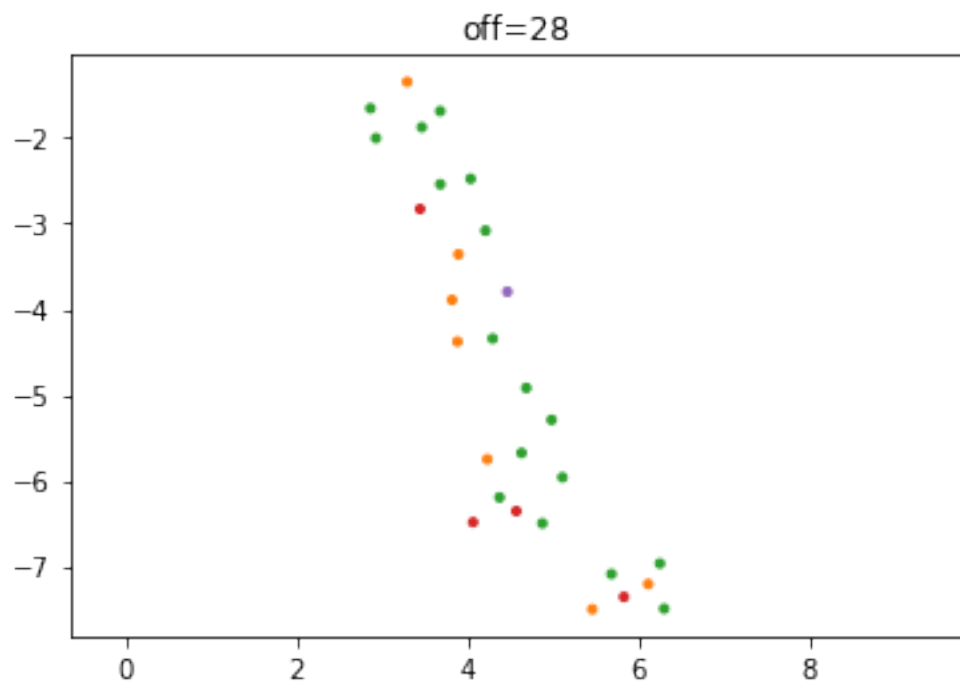


off=3

/home/ethan/.local/lib/python3.7/site-packages/umap/spectral.py:229: UserWarning:

Embedding a total of 4 separate connected components using meta-embedding (experimental)
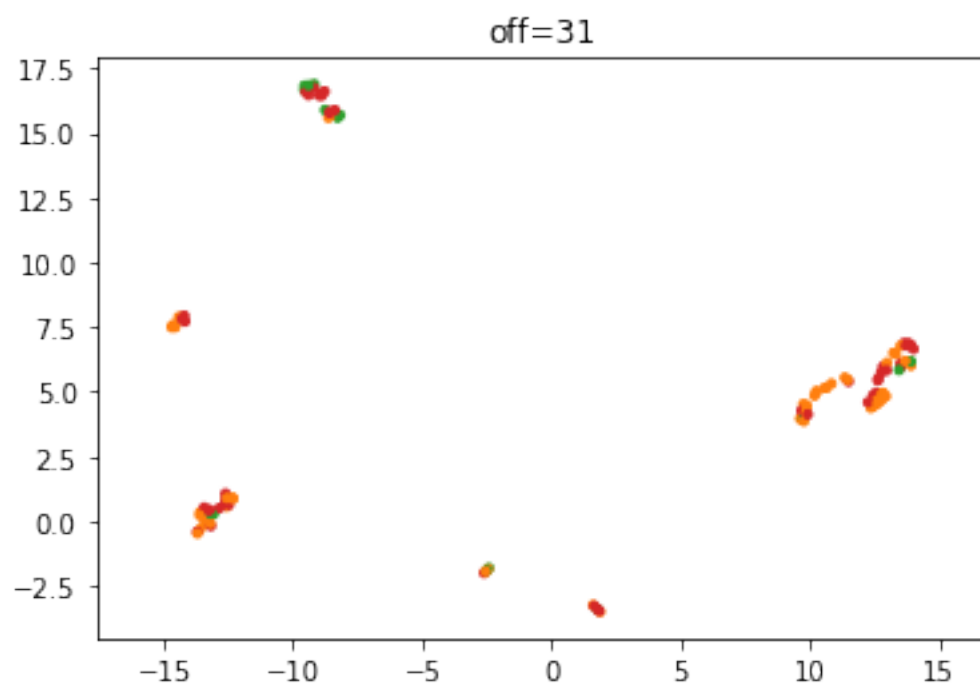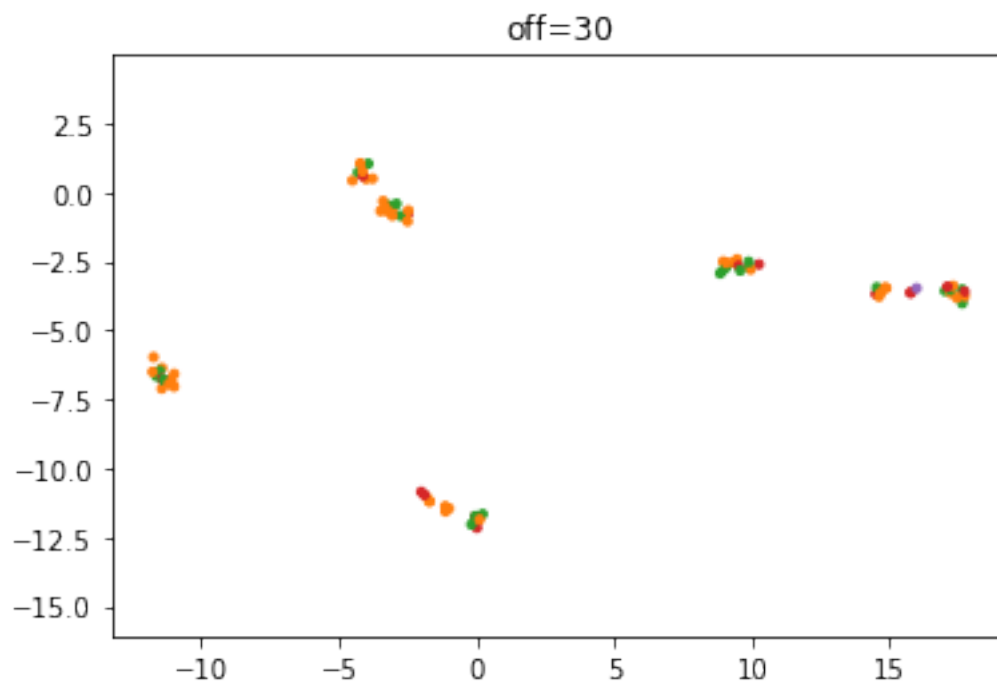
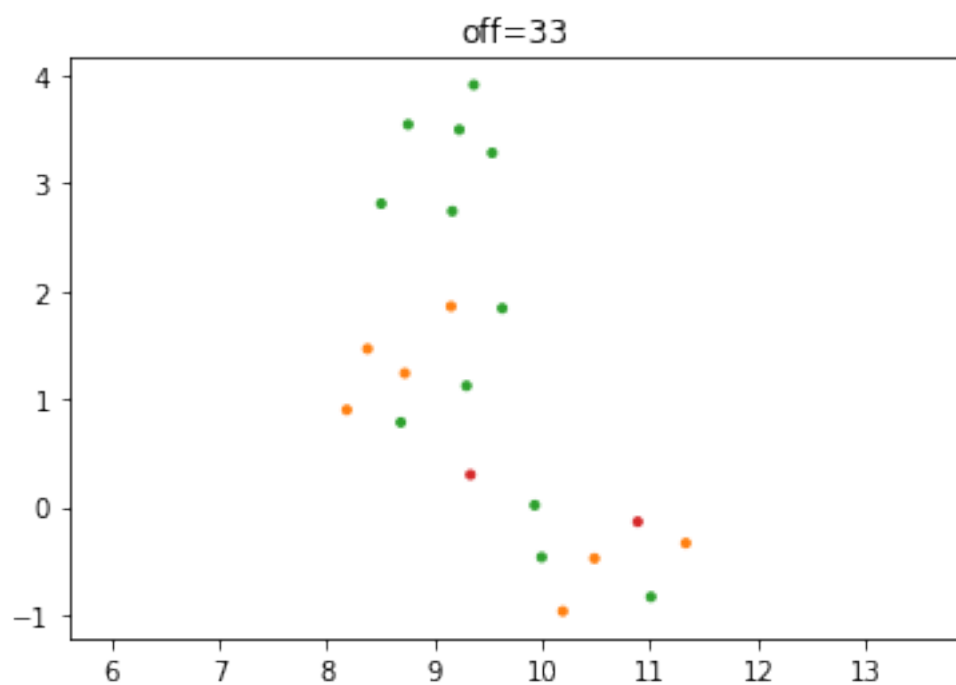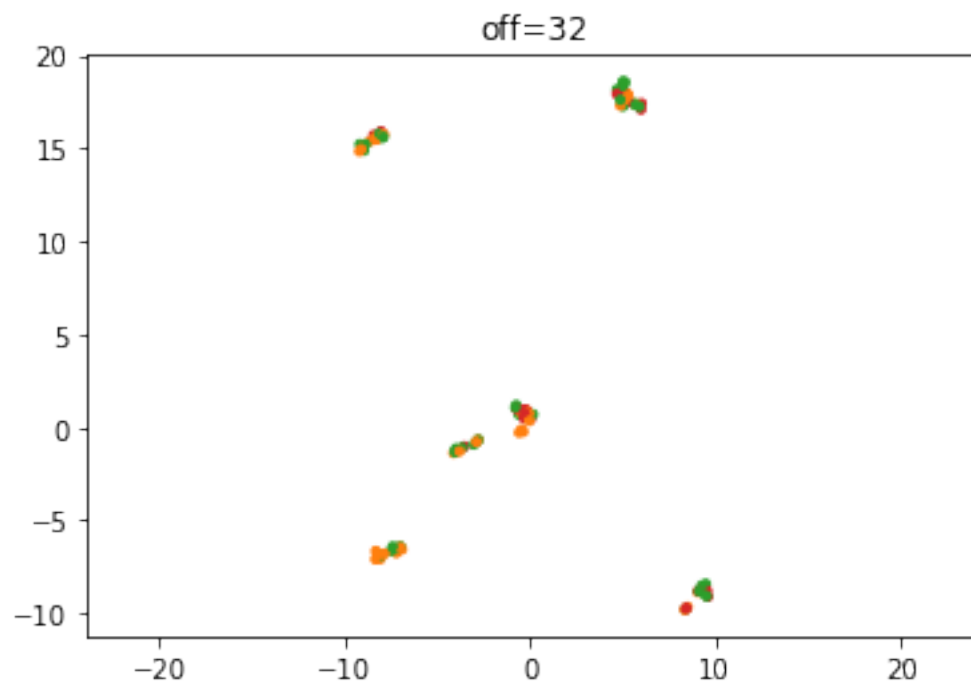

off=4

off=5

/home/ethan/.local/lib/python3.7/site-packages/umap/spectral.py:229: UserWarning:

Embedding a total of 6 separate connected components using meta-embedding (experimental)

/home/ethan/.local/lib/python3.7/site-packages/sklearn/manifold/spectral_embedding_.py:235: Use

Graph is not fully connected, spectral embedding may not work as expected.

off=6

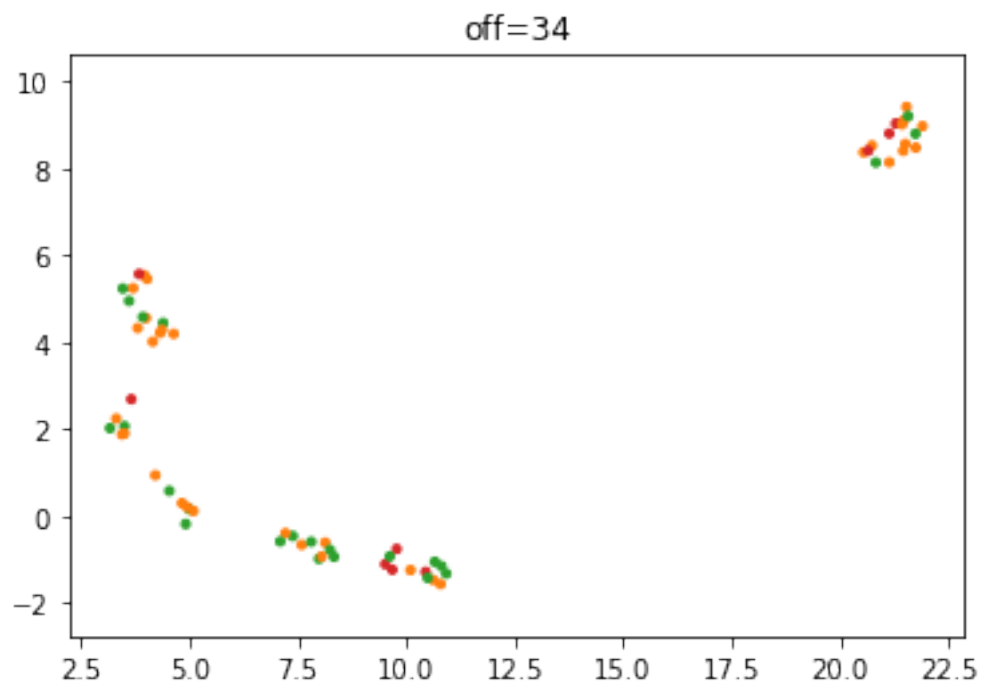/home/ethan/.local/lib/python3.7/site-packages/umap/spectral.py:229: UserWarning:

Embedding a total of 8 separate connected components using meta-embedding (experimental)
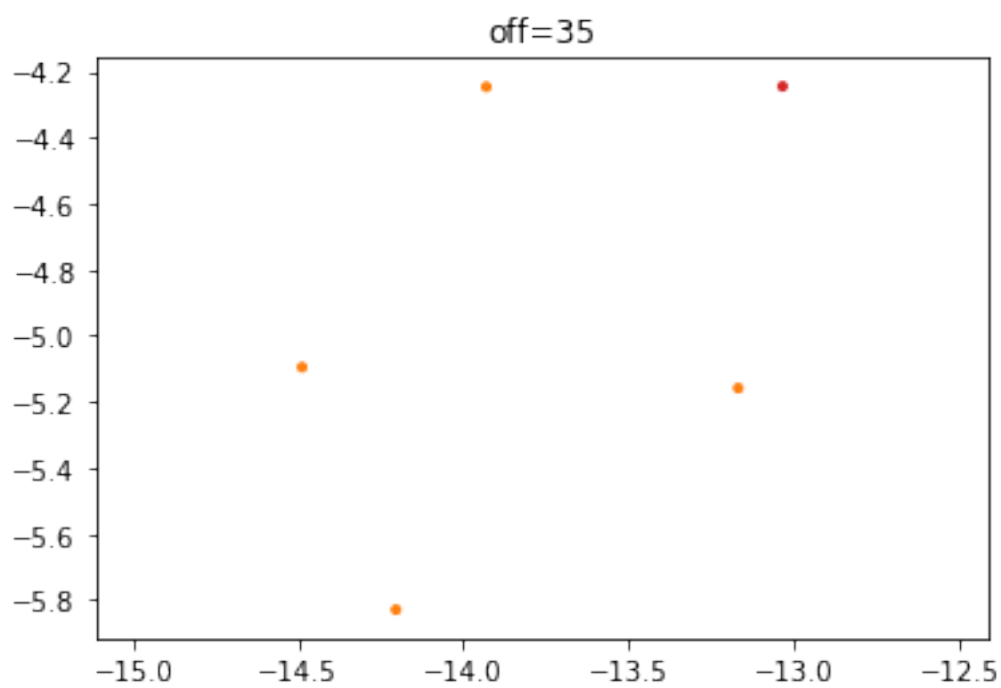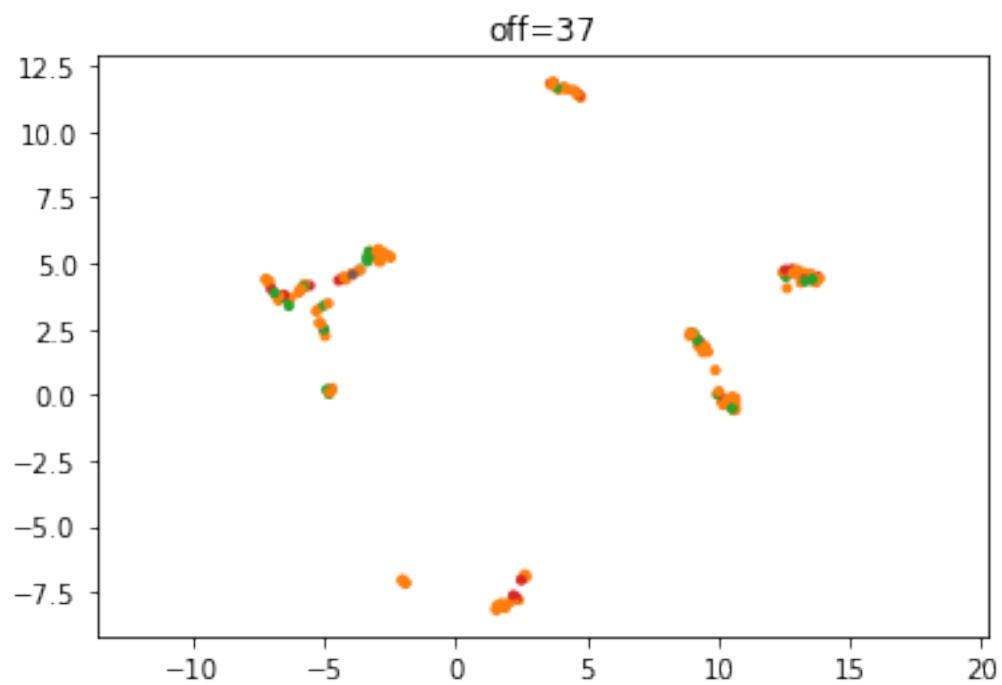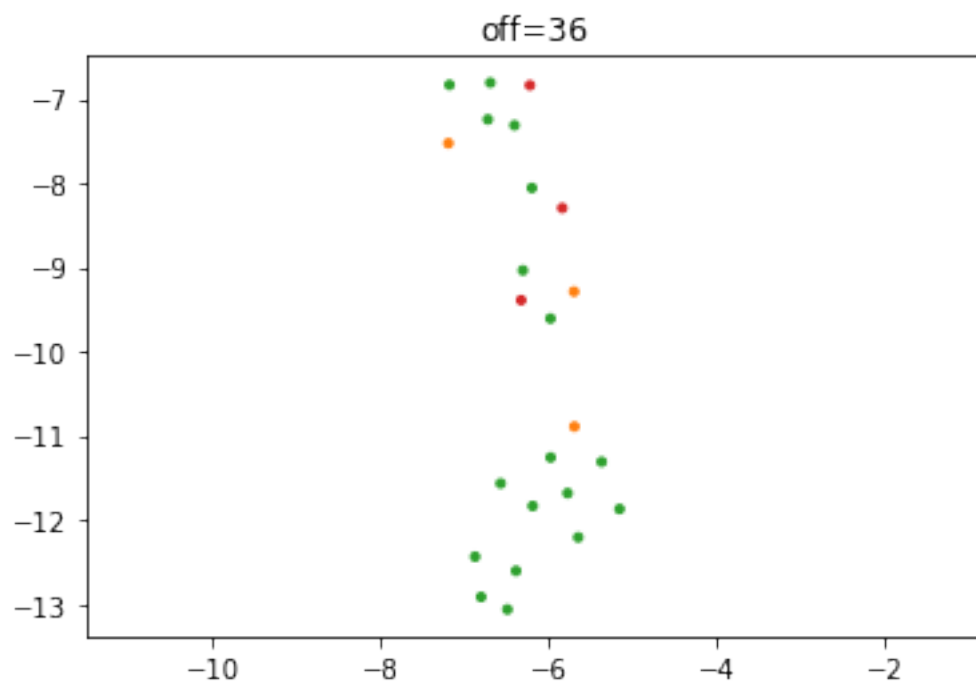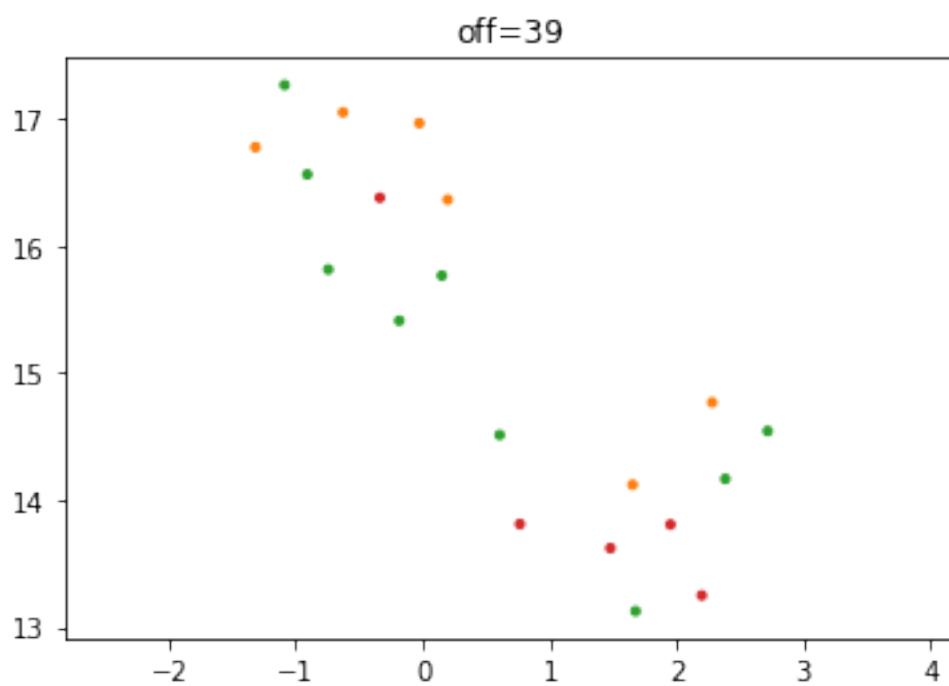
/home/ethan/.local/lib/python3.7/site-packages/sklearn/manifold/spectral_embedding_.py:235: Use

Graph is not fully connected, spectral embedding may not work as expected.

off=7

/home/ethan/.local/lib/python3.7/site-packages/umap/spectral.py:229: UserWarning:

Embedding a total of 9 separate connected components using meta-embedding (experimental)

/home/ethan/.local/lib/python3.7/site-packages/sklearn/manifold/spectral_embedding_.py:235: Us
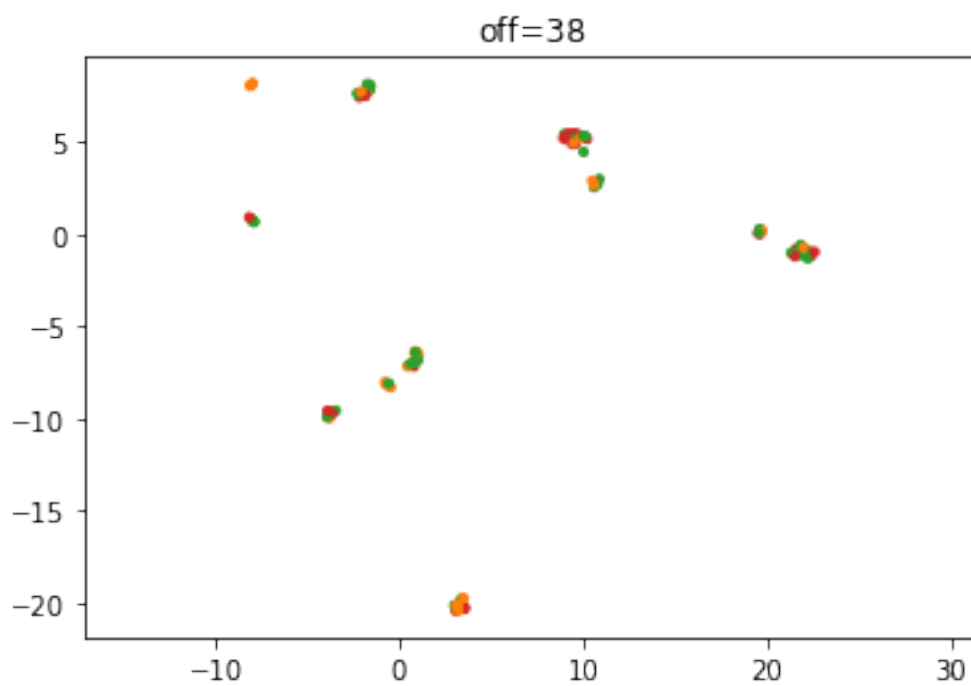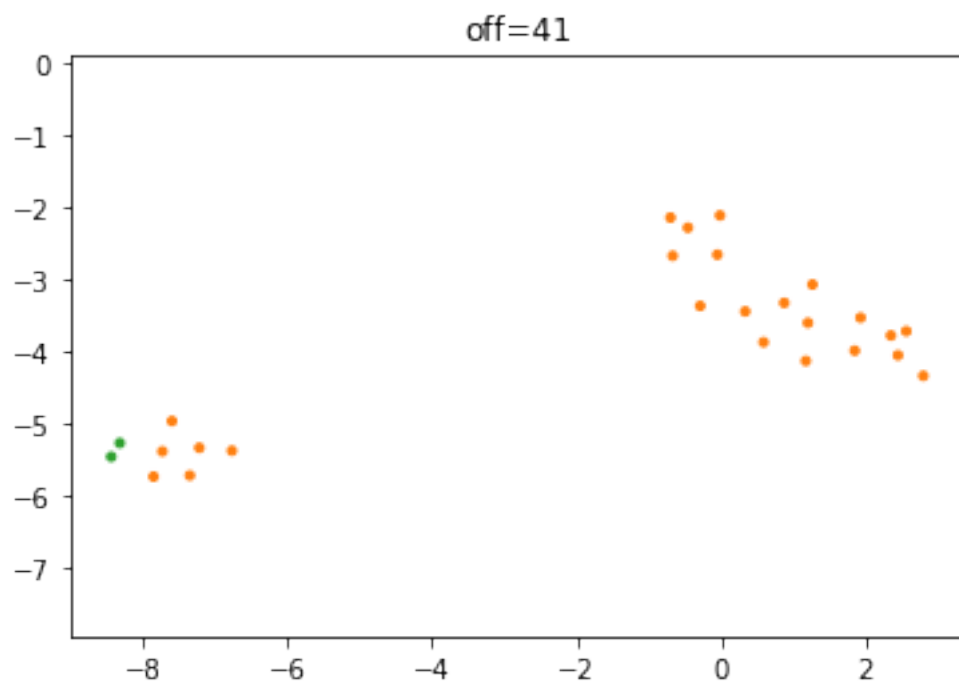
Graph is not fully connected, spectral embedding may not work as expected.

off=8

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

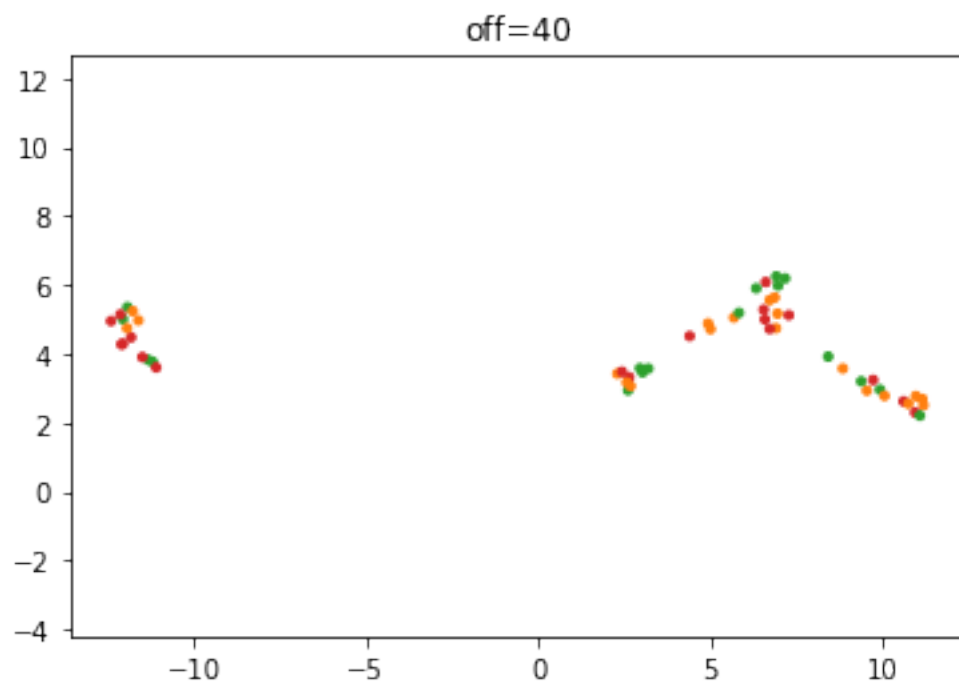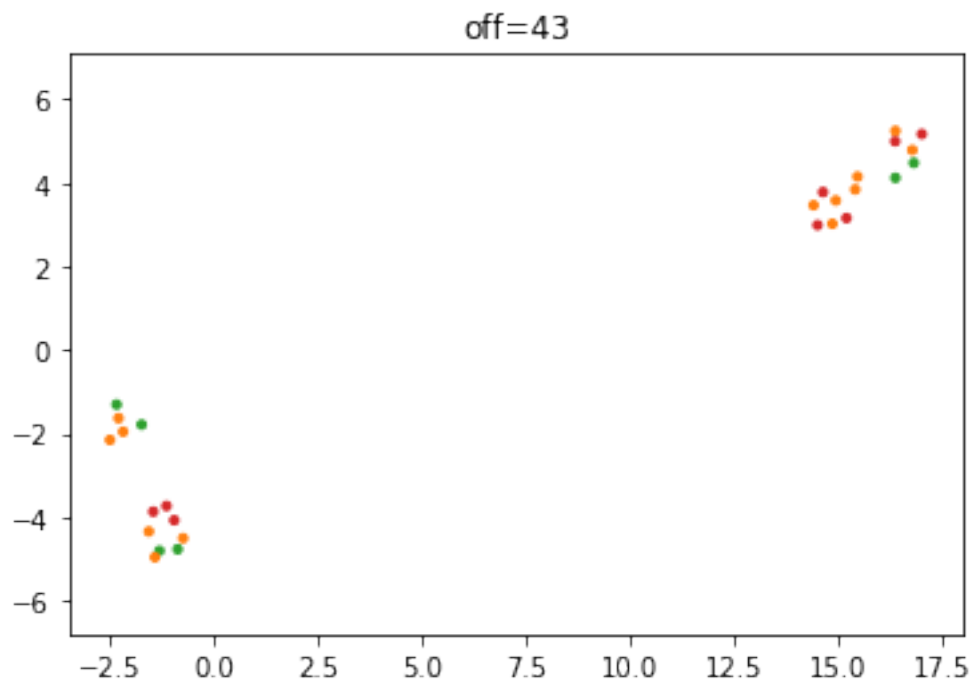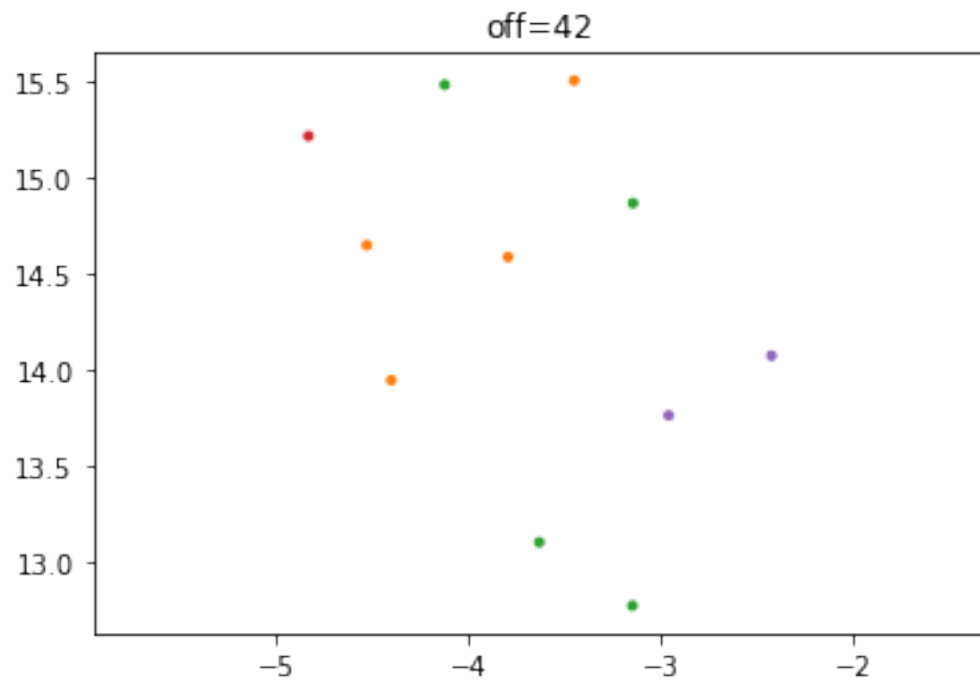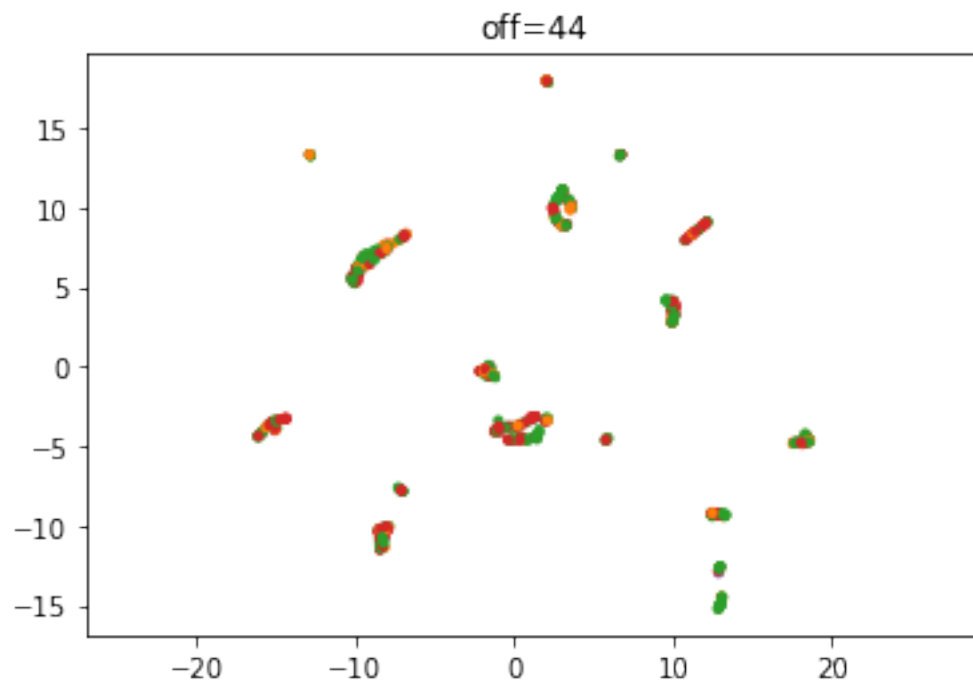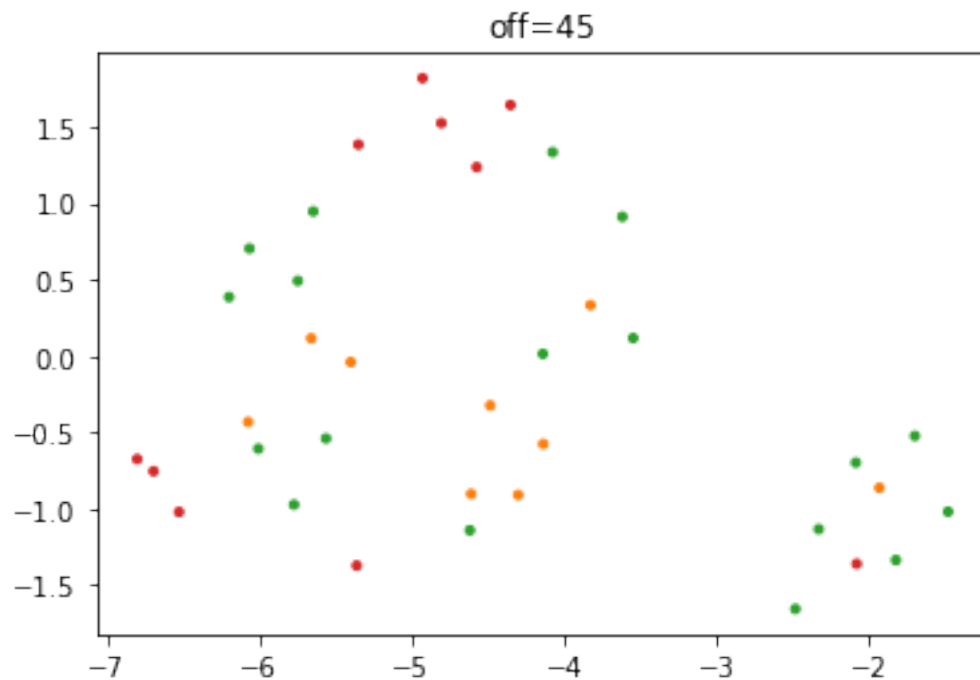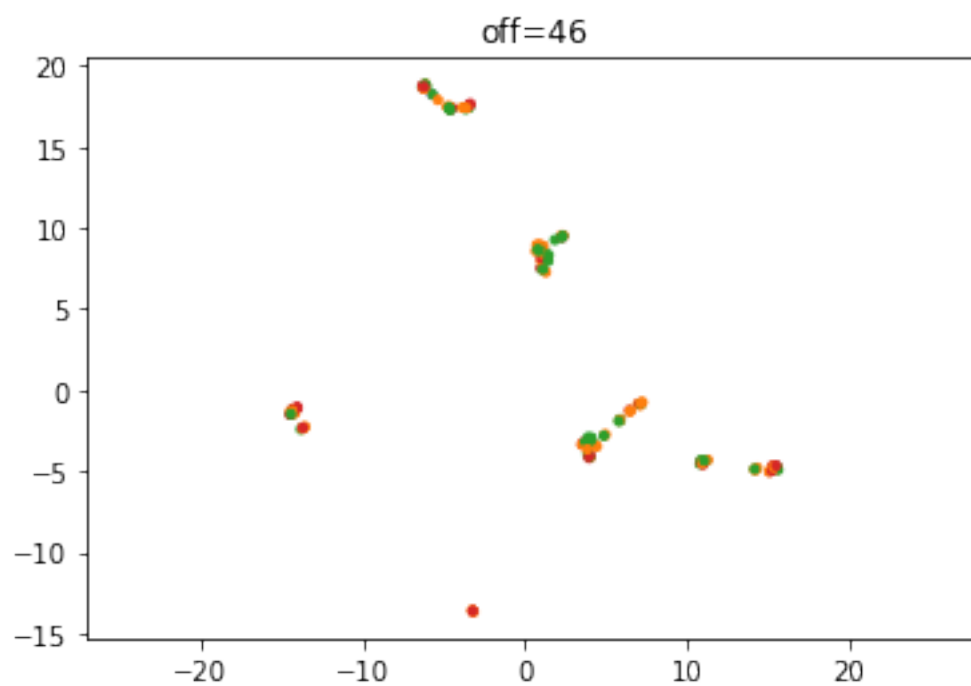n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

/home/ethan/.local/lib/python3.7/site-packages/scipy/sparse/linalg/eigen/arpack/arpack.py:1592

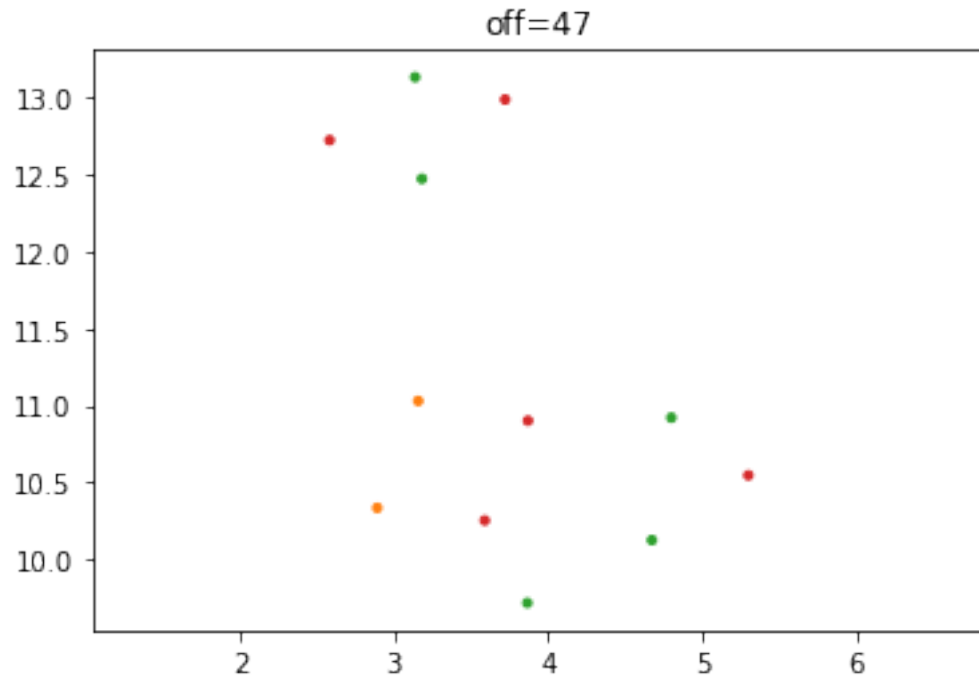k >= N for N * N square matrix. Attempting to use scipy.linalg.eigh instead.

/home/ethan/.local/lib/python3.7/site-packages/umap/spectral.py:229: UserWarning:
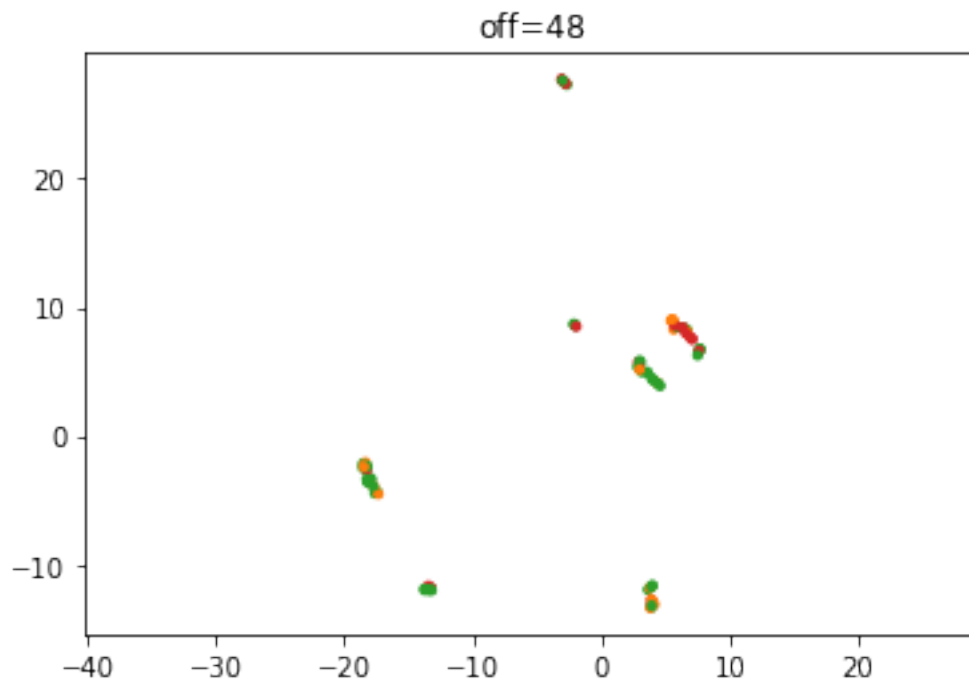
Embedding a total of 3 separate connected components using meta-embedding (experimental)

off=10

/home/ethan/.local/lib/python3.7/site-packages/umap/spectral.py:229: UserWarning:

Embedding a total of 2 separate connected components using meta-embedding (experimental)



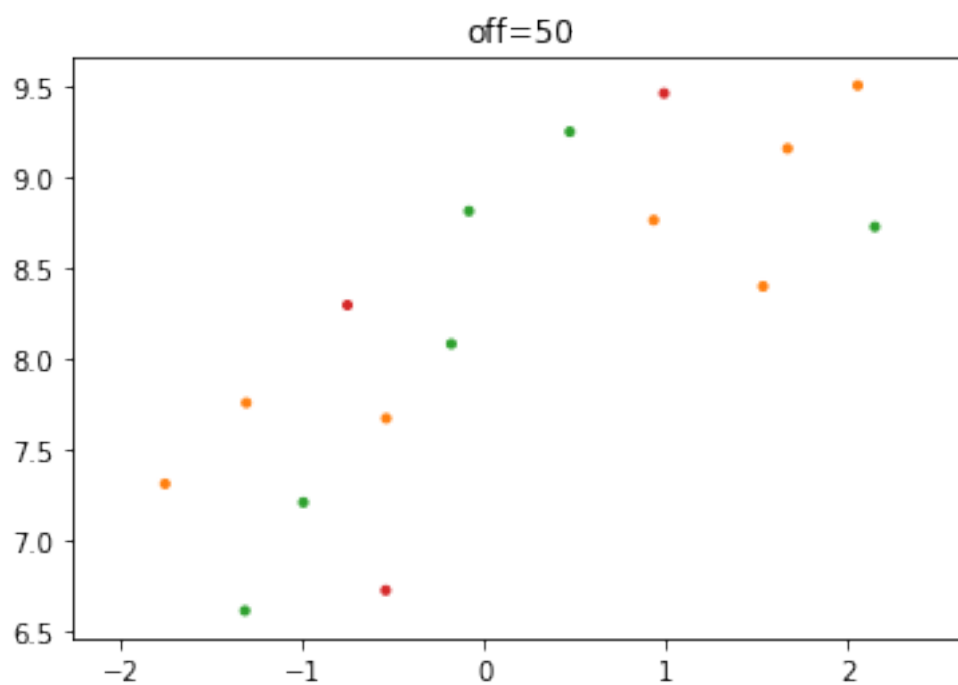off=11

/home/ethan/.local/lib/python3.7/site-packages/umap/spectral.py:229: UserWarning:

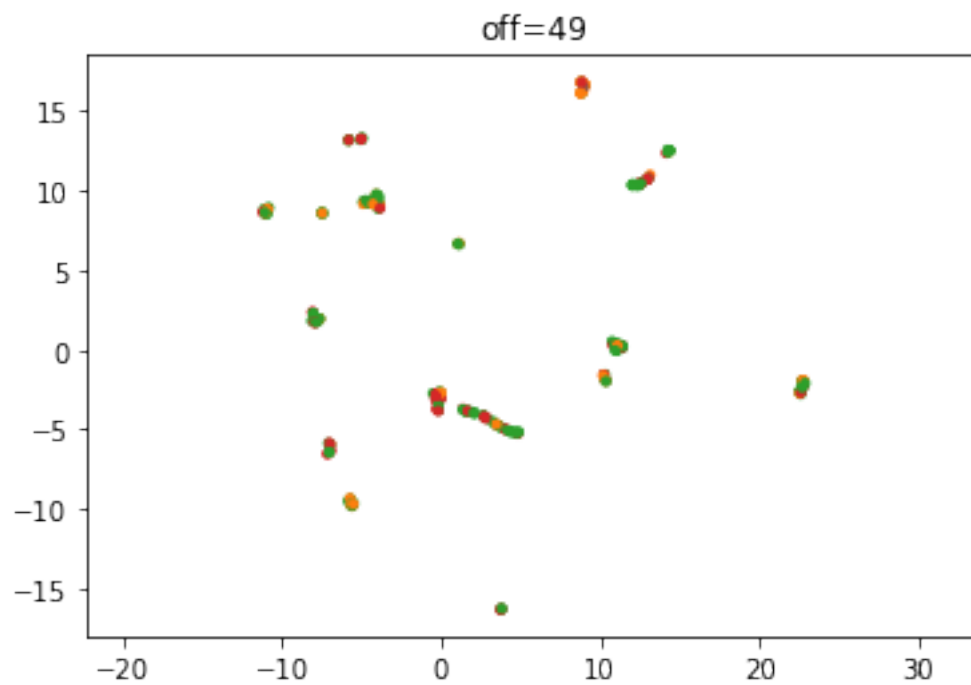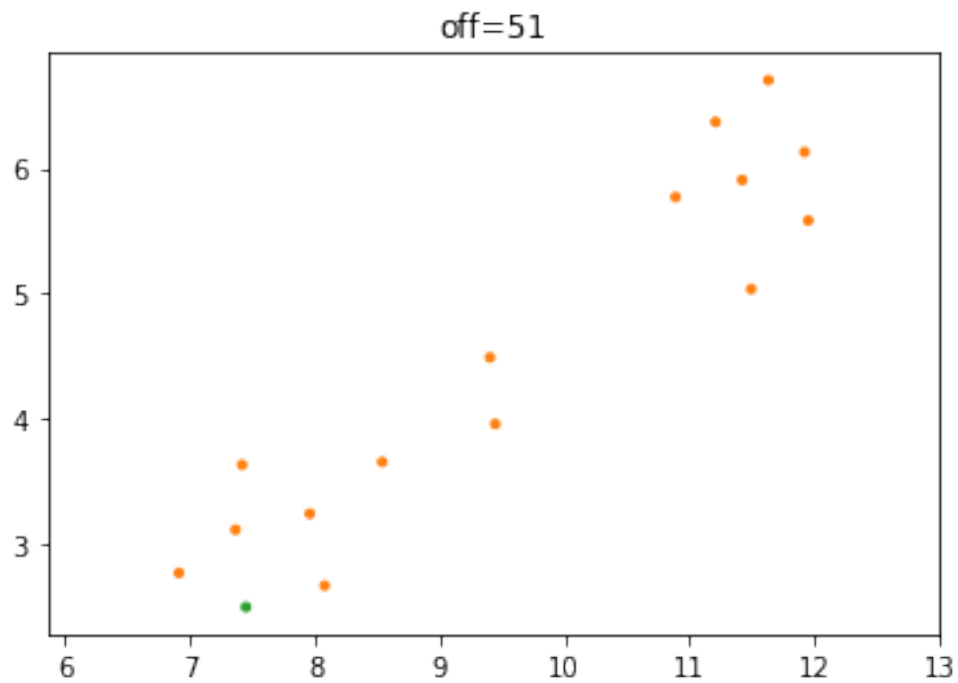Embedding a total of 5 separate connected components using meta-embedding (experimental)

/home/ethan/.local/lib/python3.7/site-packages/sklearn/manifold/spectral_embedding_.py:235: Use

Graph is not fully connected, spectral embedding may not work as expected.



off=12

/home/ethan/.local/lib/python3.7/site-packages/umap/spectral.py:229: UserWarning:

Embedding a total of 2 separate connected components using meta-embedding (experimental)

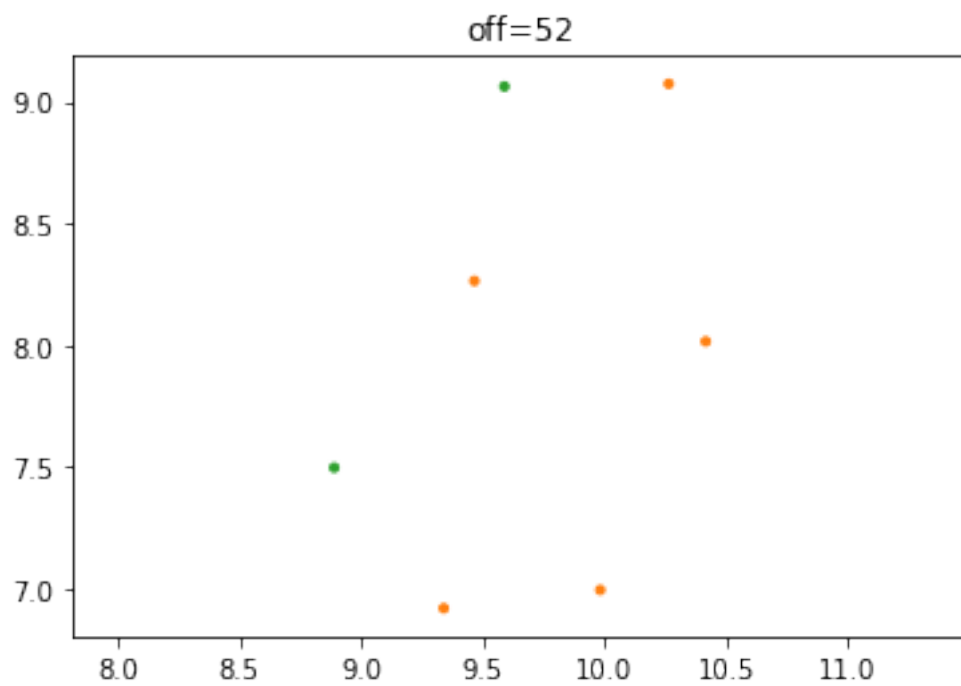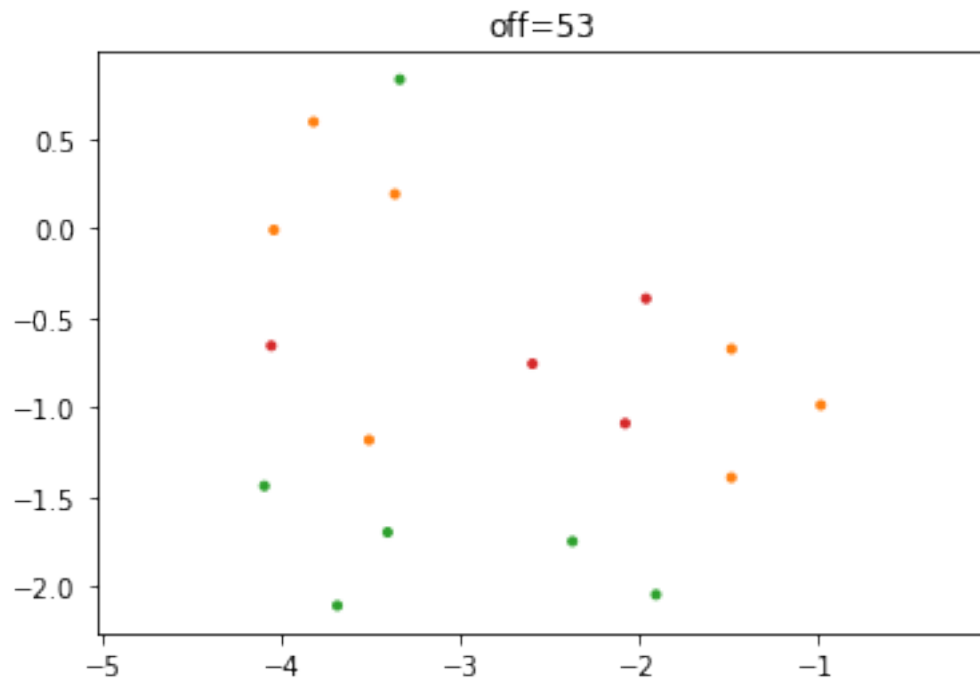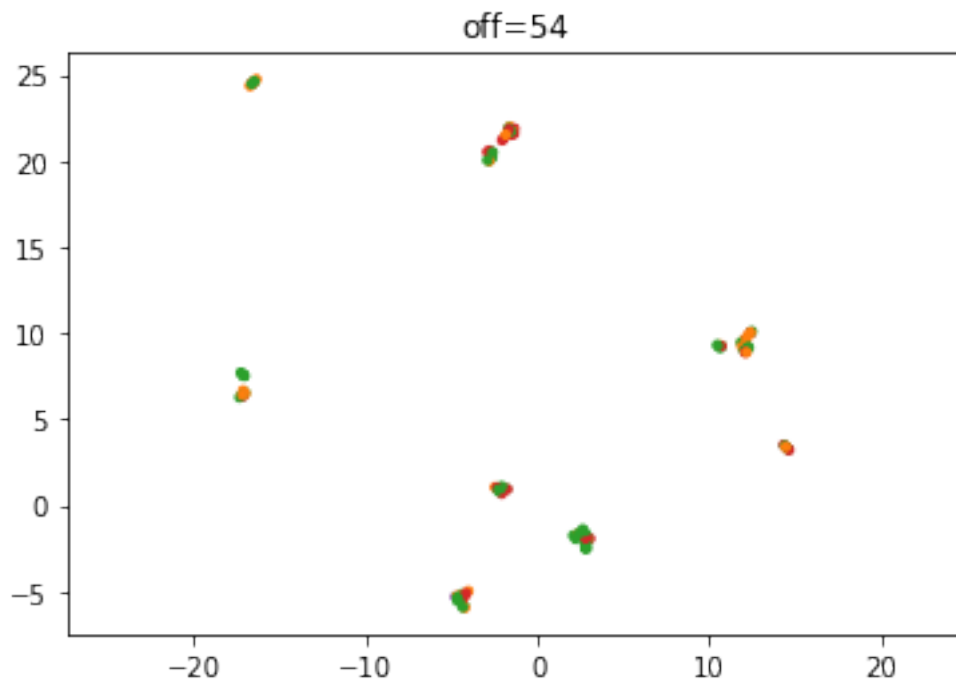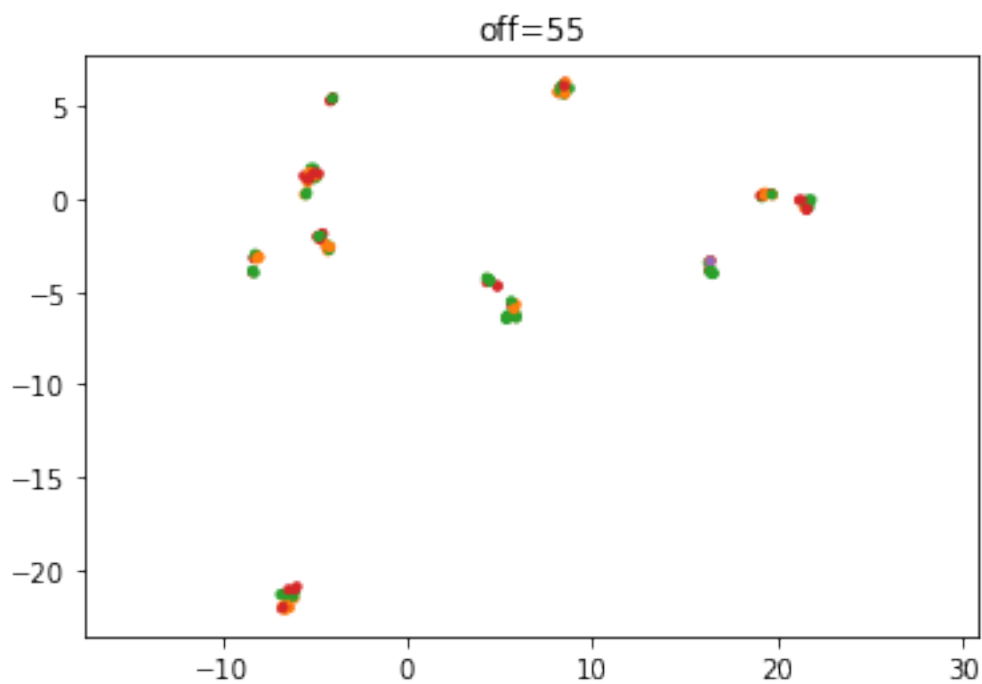off=13

/home/ethan/.local/lib/python3.7/site-packages/umap/spectral.py:229: UserWarning:
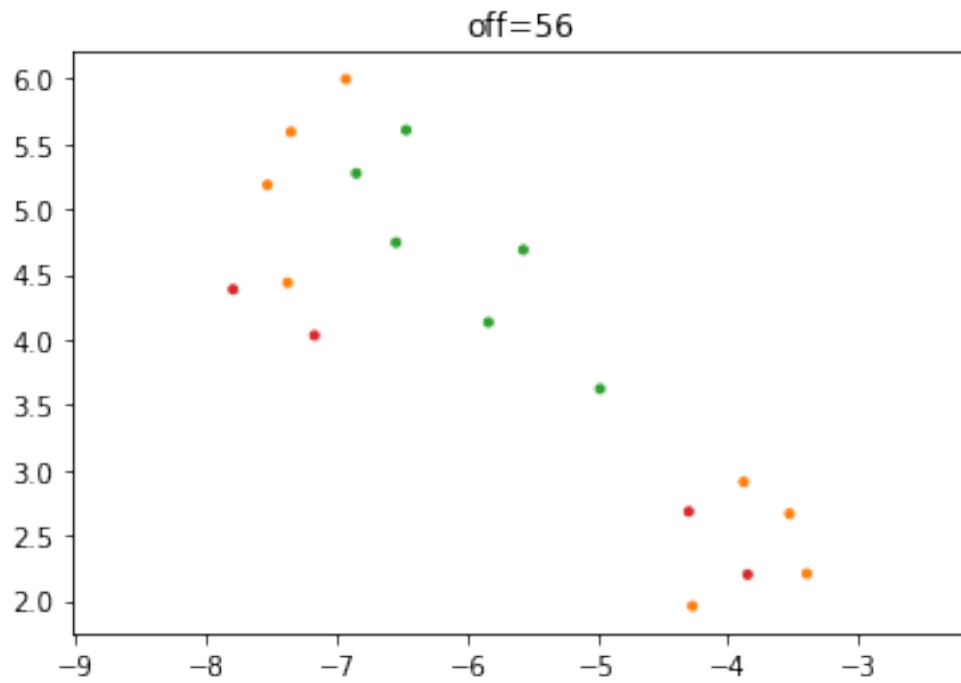
Embedding a total of 5 separate connected components using meta-embedding (experimental)

/home/ethan/.local/lib/python3.7/site-packages/sklearn/manifold/spectral_embedding_.py:235: Use

Graph is not fully connected, spectral embedding may not work as expected.
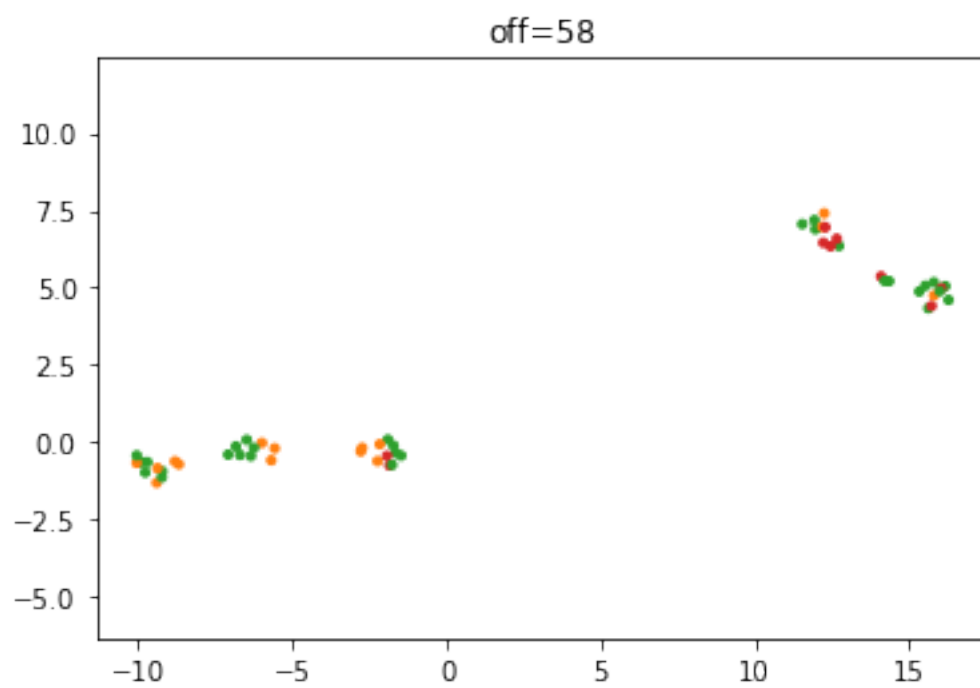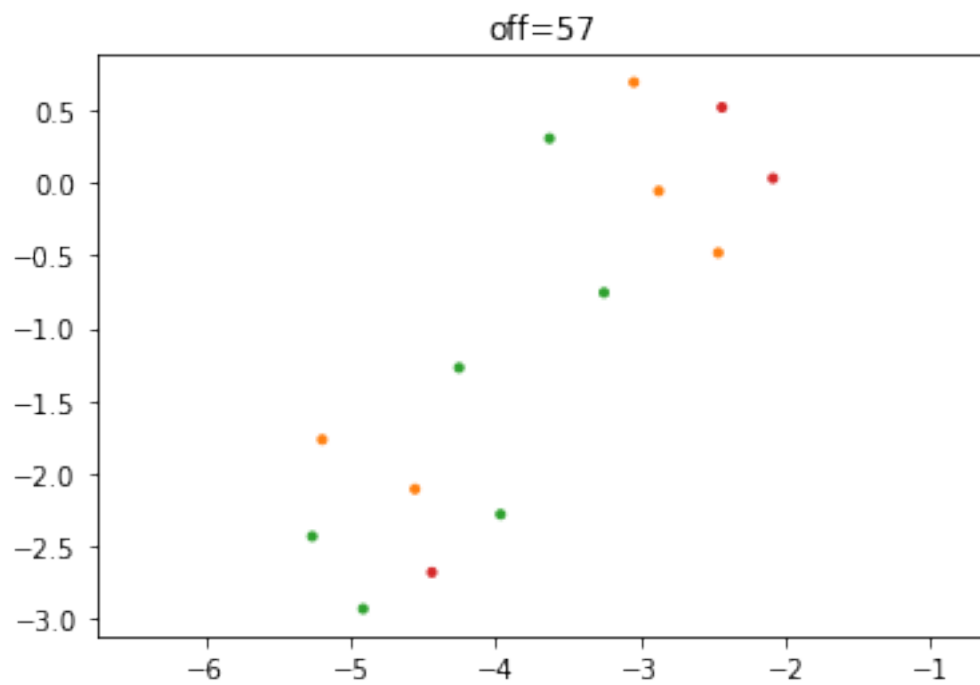
off=14

/home/ethan/.local/lib/python3.7/site-packages/umap/spectral.py:229: UserWarning:

Embedding a total of 7 separate connected components using meta-embedding (experimental)

/home/ethan/.local/lib/python3.7/site-packages/sklearn/manifold/spectral_embedding_.py:235: Use

Graph is not fully connected, spectral embedding may not work as expected.

off=15

/home/ethan/.local/lib/python3.7/site-packages/umap/spectral.py:229: UserWarning:

Embedding a total of 3 separate connected components using meta-embedding (experimental)



off=16

off=17

/home/ethan/.local/lib/python3.7/site-packages/umap/spectral.py:229: UserWarning:

Embedding a total of 8 separate connected components using meta-embedding (experimental)

/home/ethan/.local/lib/python3.7/site-packages/sklearn/manifold/spectral_embedding_.py:235: Use

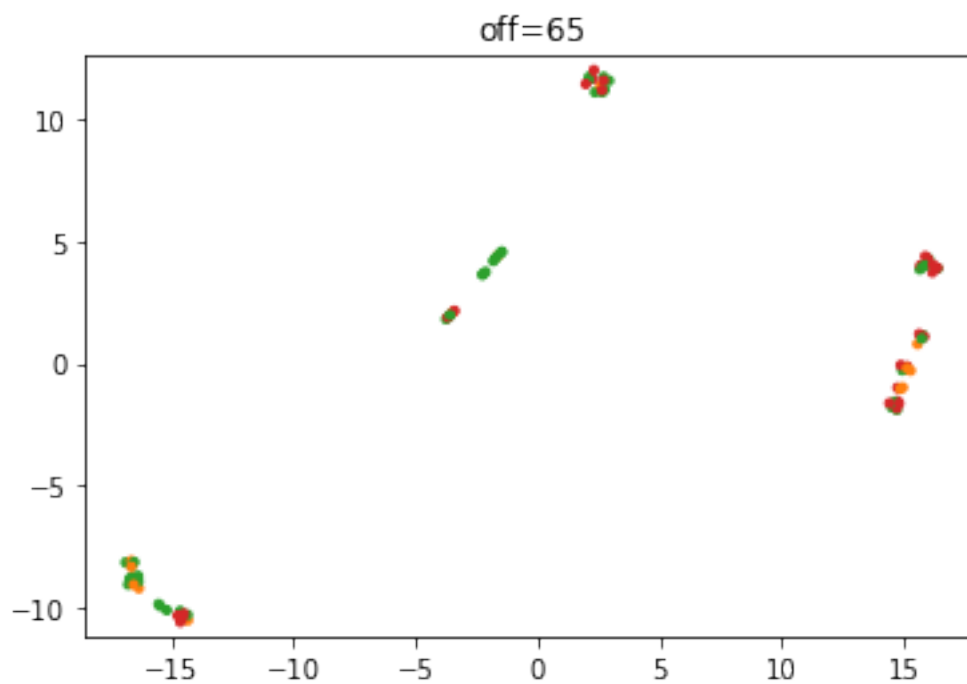Graph is not fully connected, spectral embedding may not work as expected.

off=18

off=19

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

off=20

off=21

off=22



off=23

/home/ethan/.local/lib/python3.7/site-packages/umap/spectral.py:229: UserWarning:

Embedding a total of 4 separate connected components using meta-embedding (experimental)



off=24

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

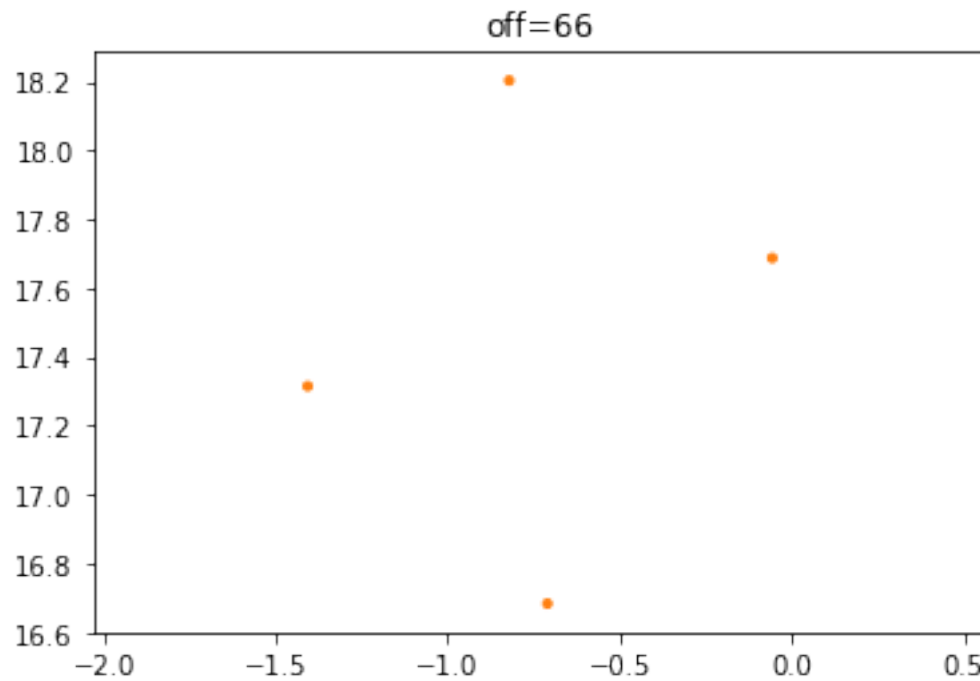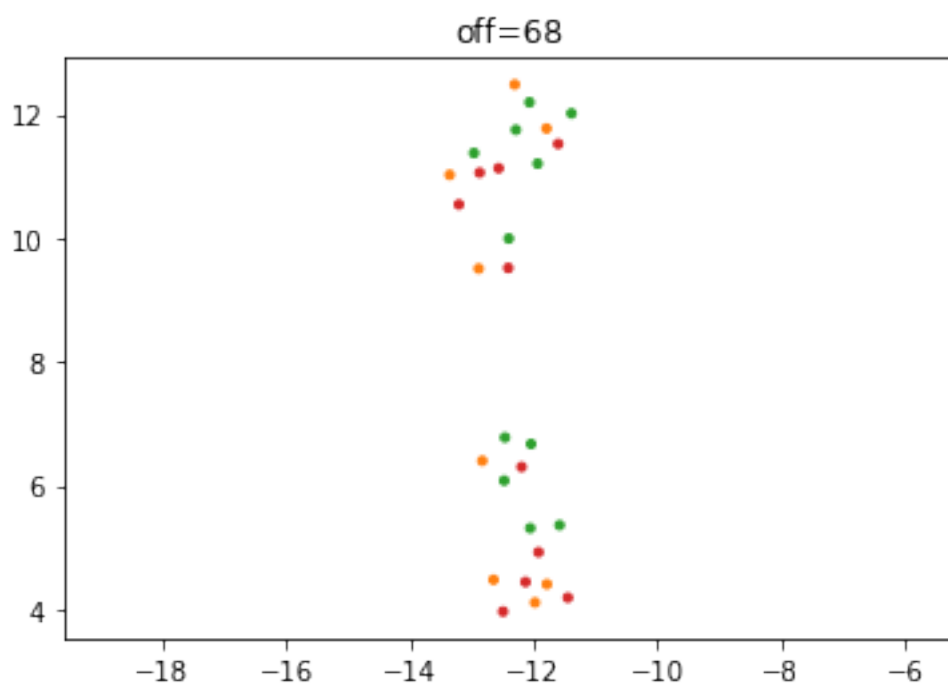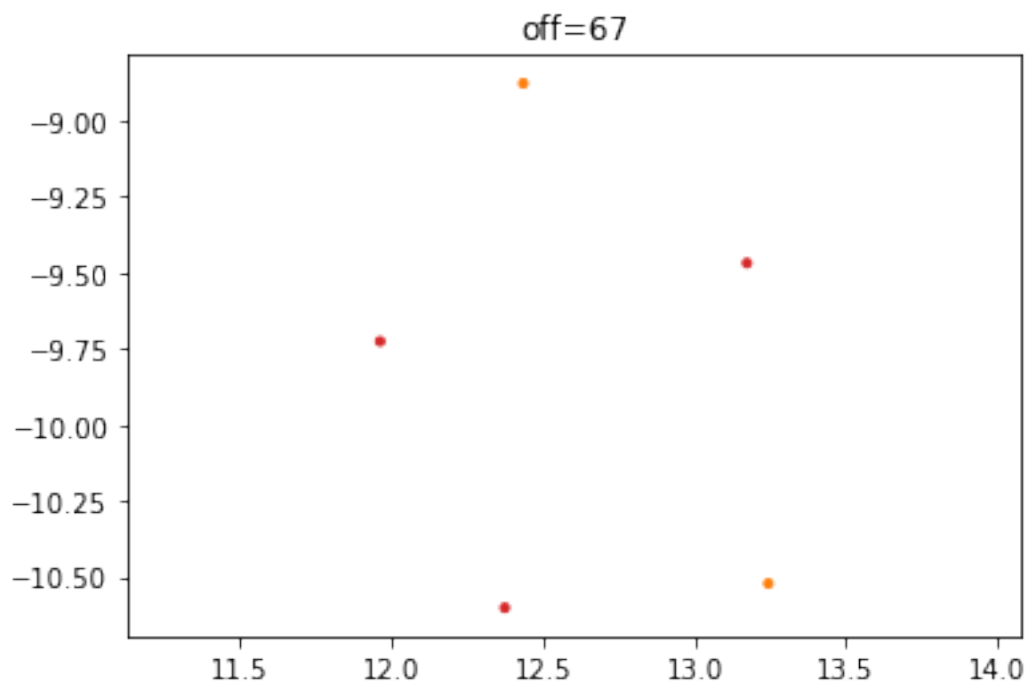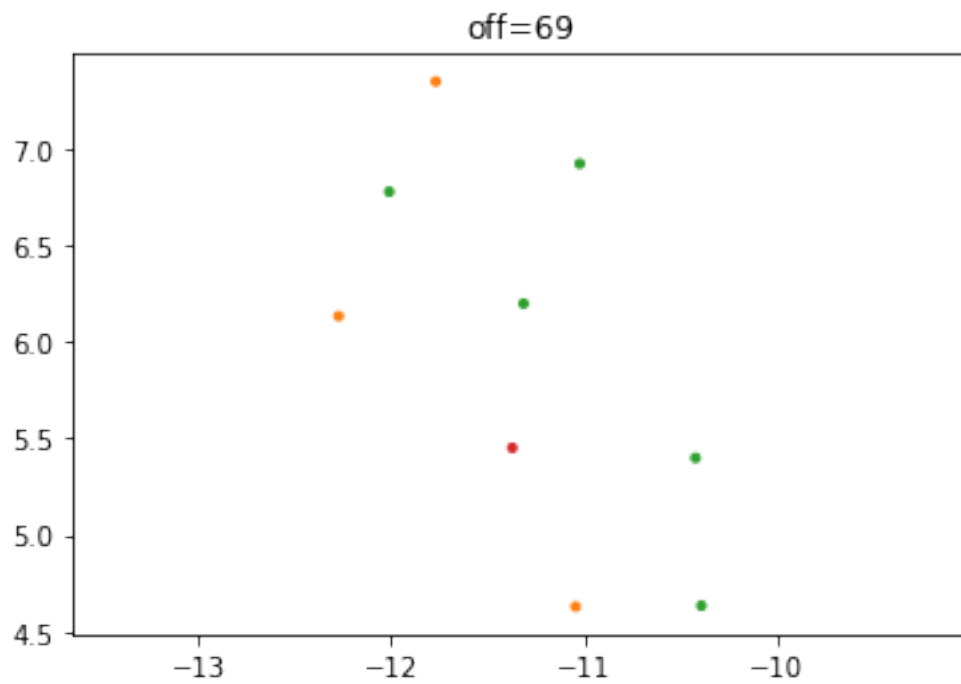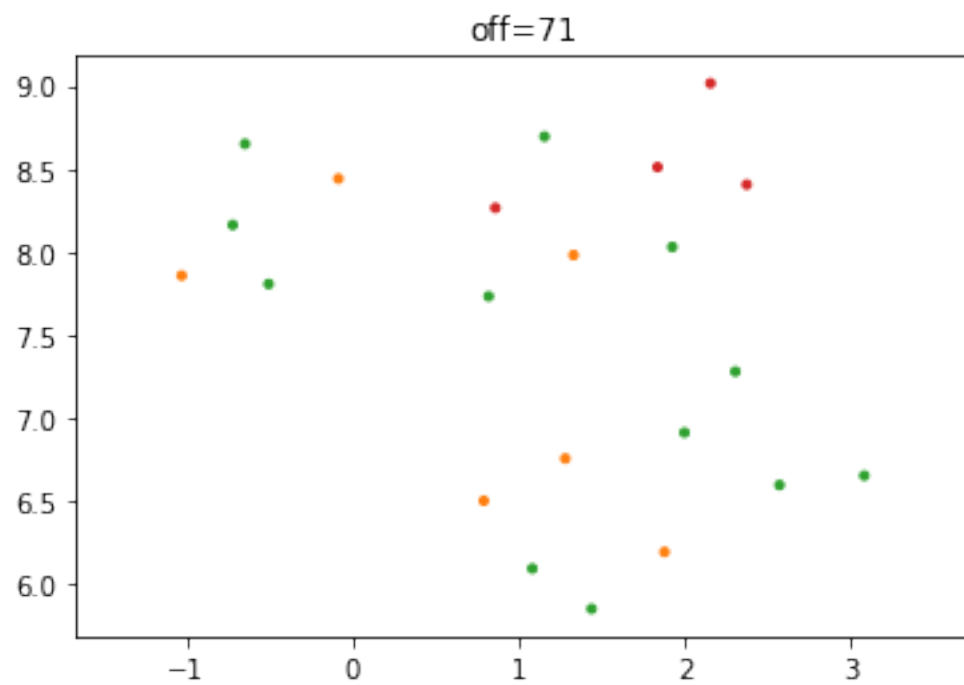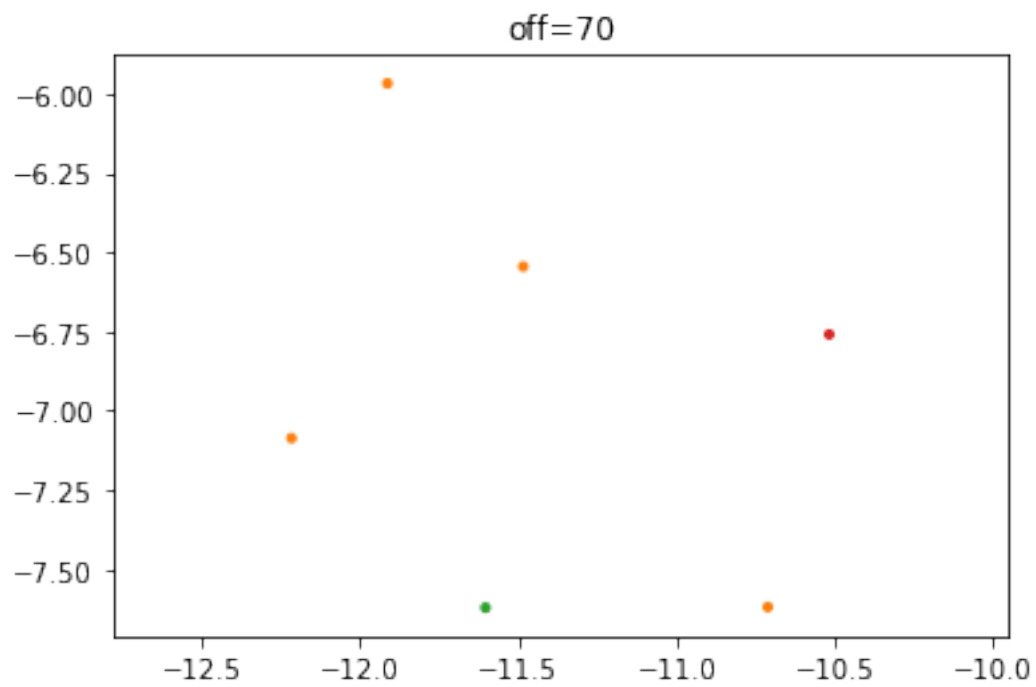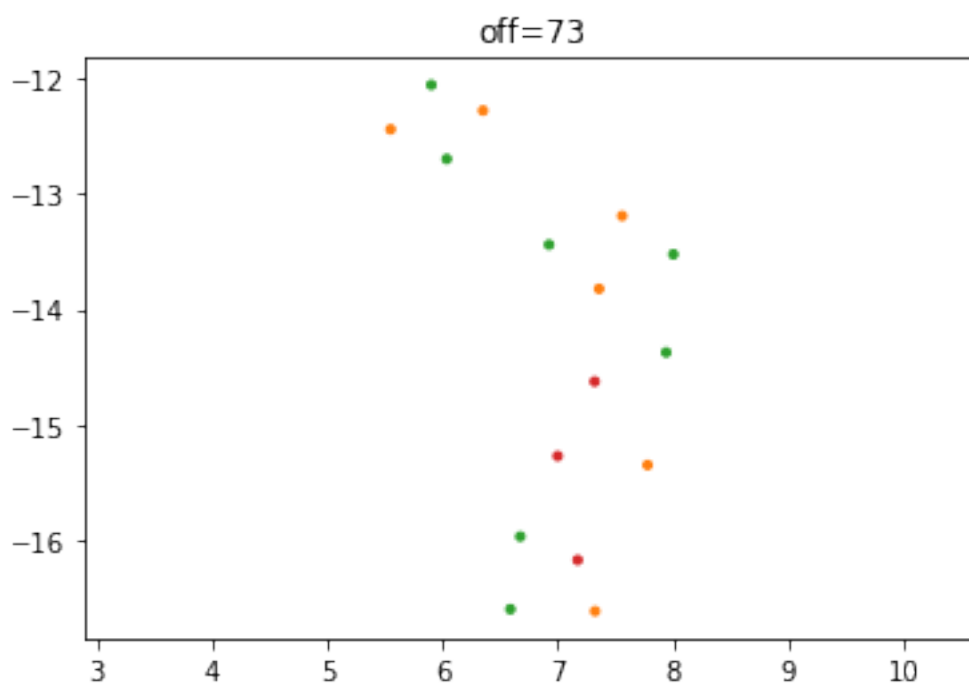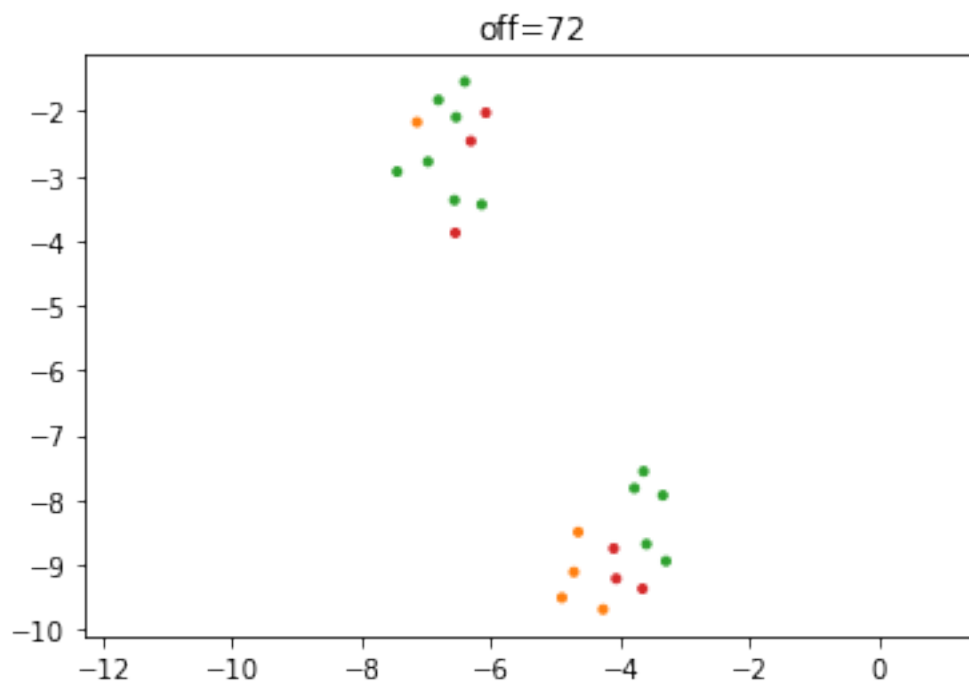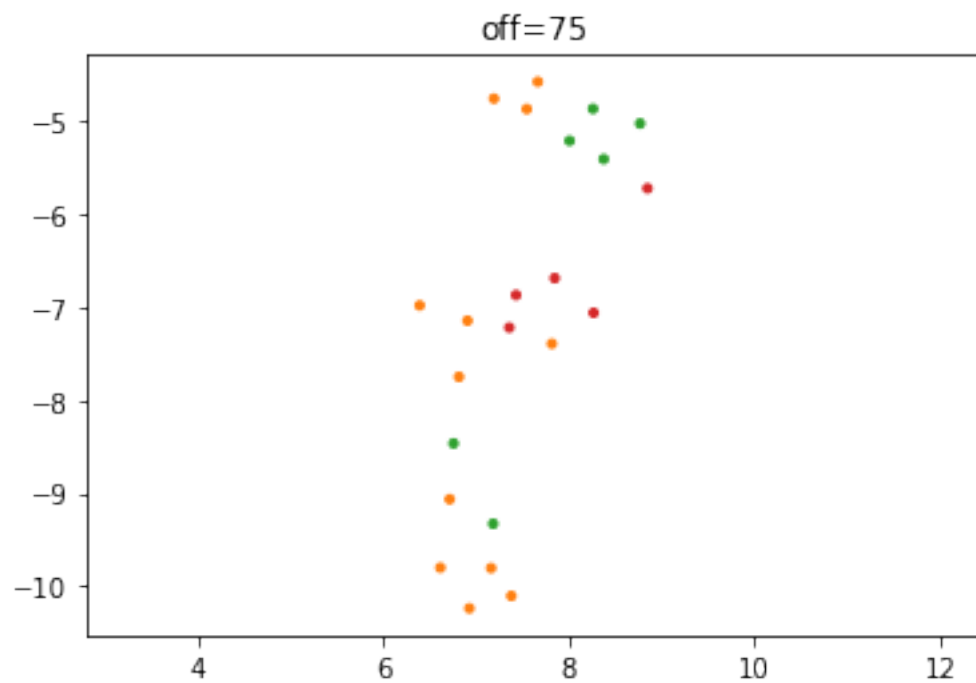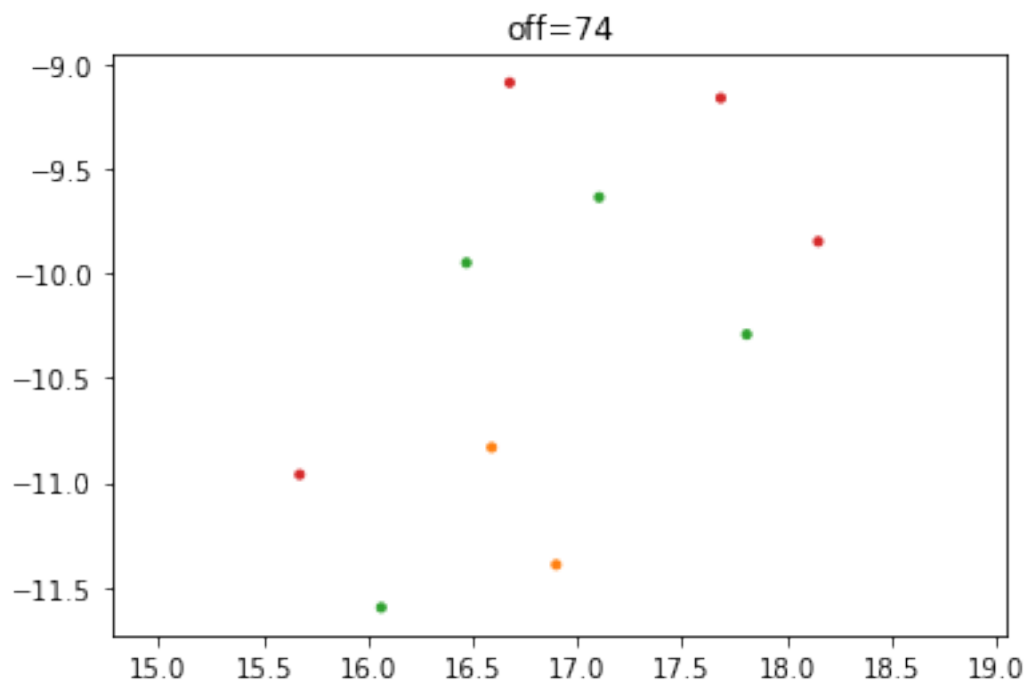n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

off=25

/home/ethan/.local/lib/python3.7/site-packages/umap/spectral.py:229: UserWarning:

Embedding a total of 4 separate connected components using meta-embedding (experimental)



off=26

/home/ethan/.local/lib/python3.7/site-packages/umap/spectral.py:229: UserWarning:

Embedding a total of 4 separate connected components using meta-embedding (experimental)



off=27

off=28
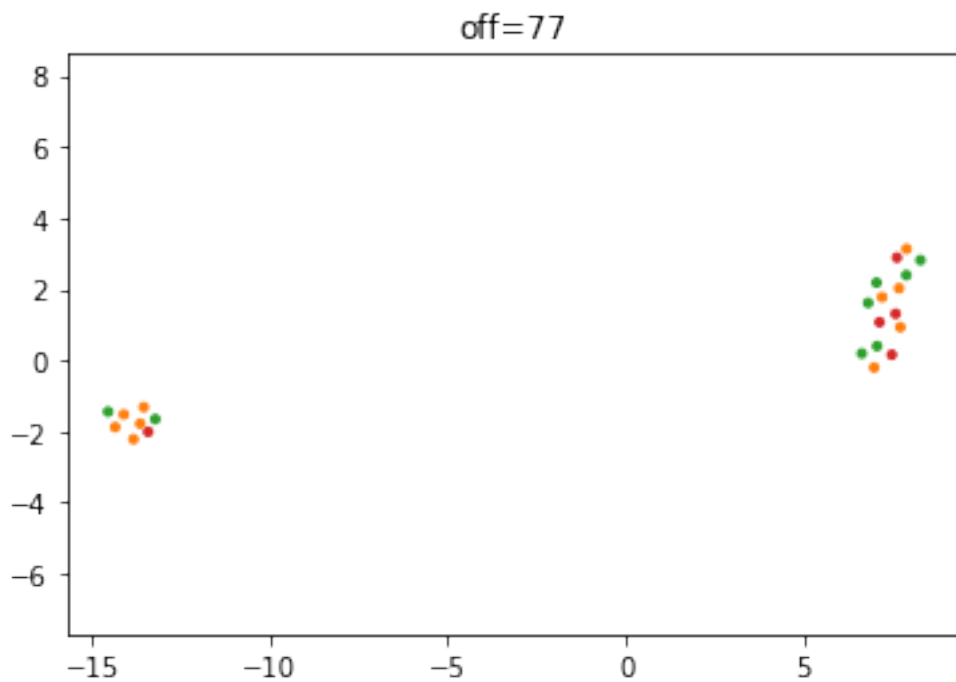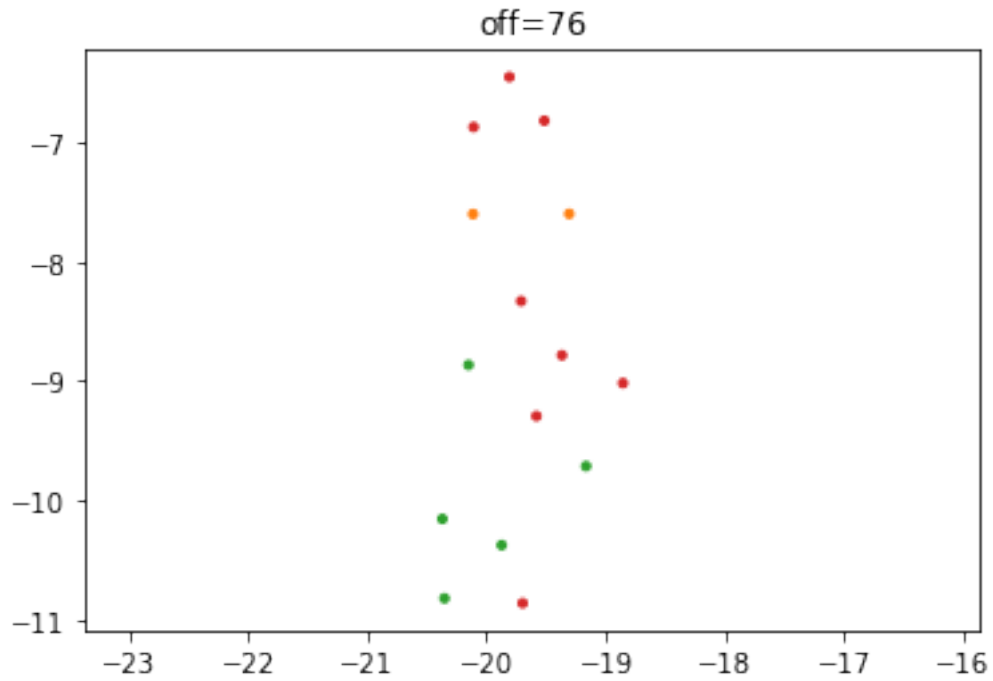


off=29

off=30



off=31

off=32



off=33

**off=34**



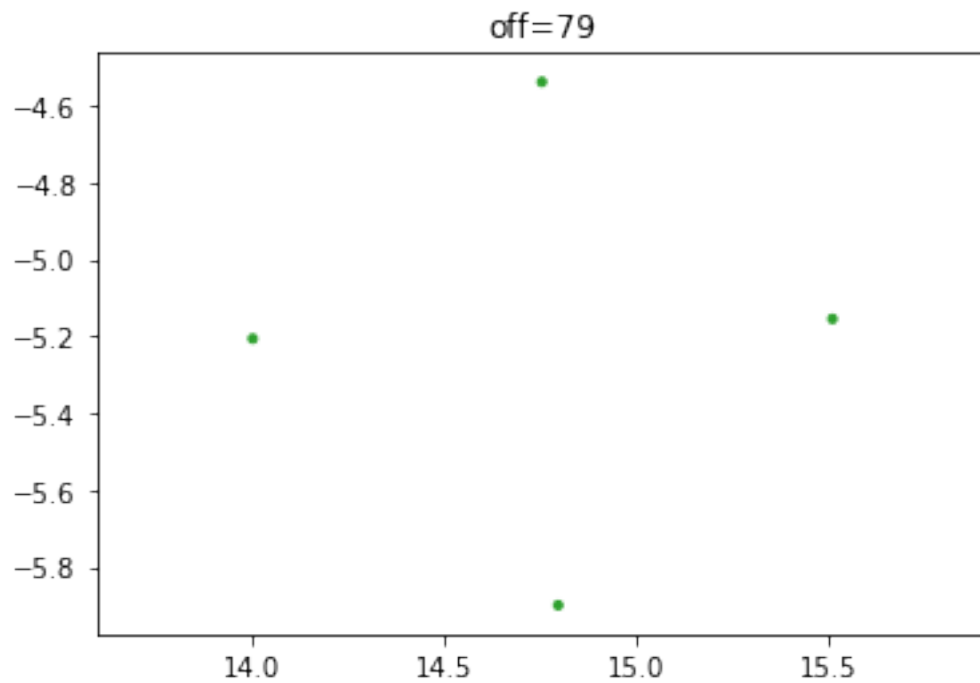/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

**off=35**

off=36

off=37

```
/home/ethan/.local/lib/python3.7/site-packages/umap/spectral.py:229: UserWarning:

Embedding a total of 2 separate connected components using meta-embedding (experimental)
```
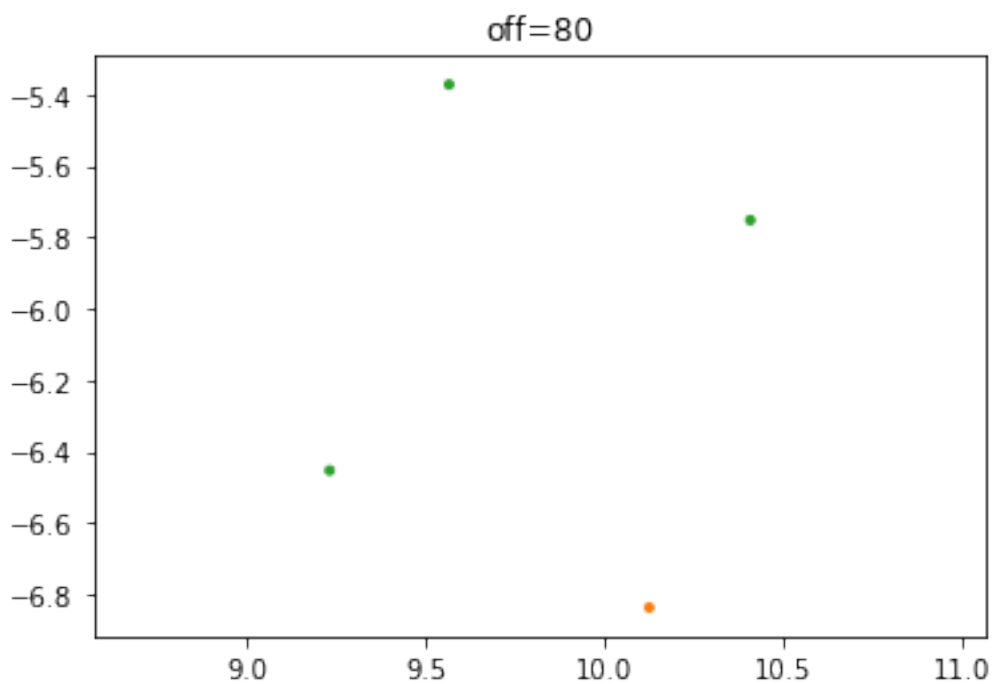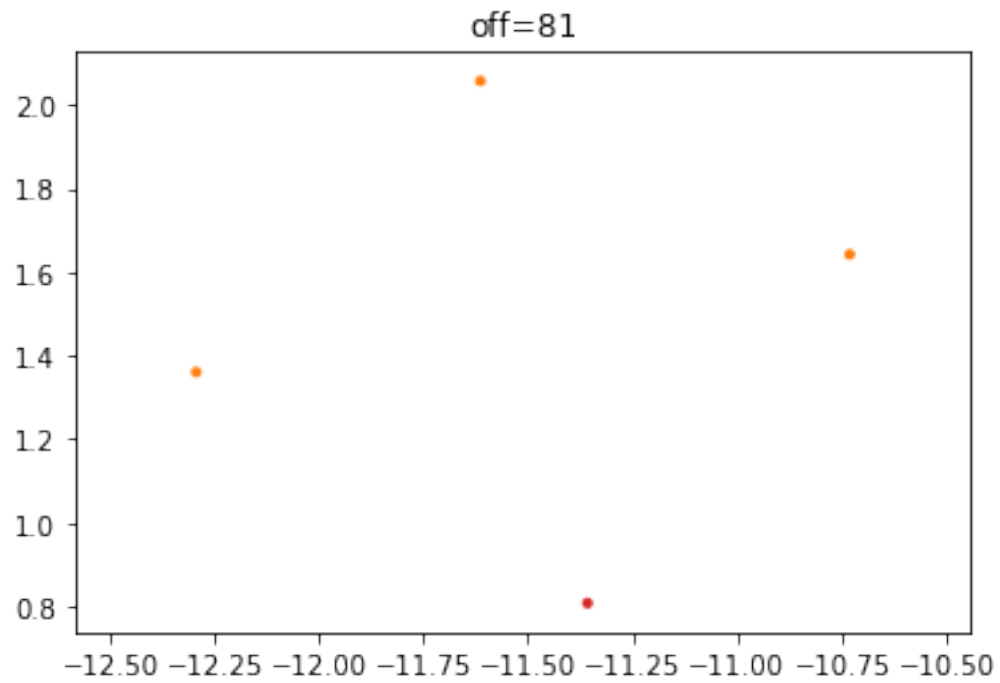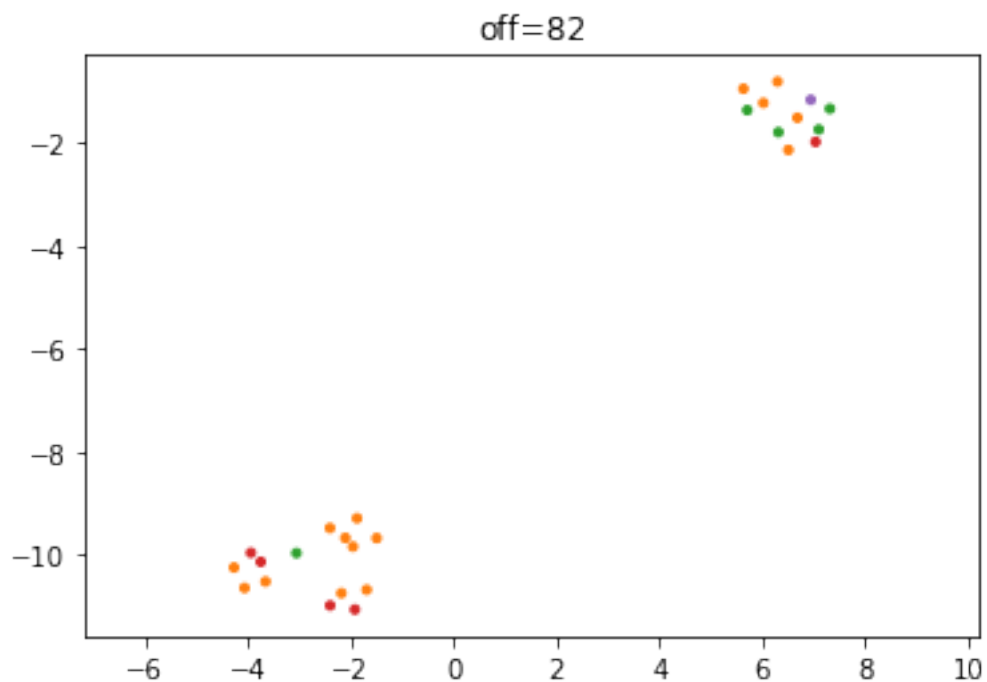


off=38



off=39

off=40

off=41

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

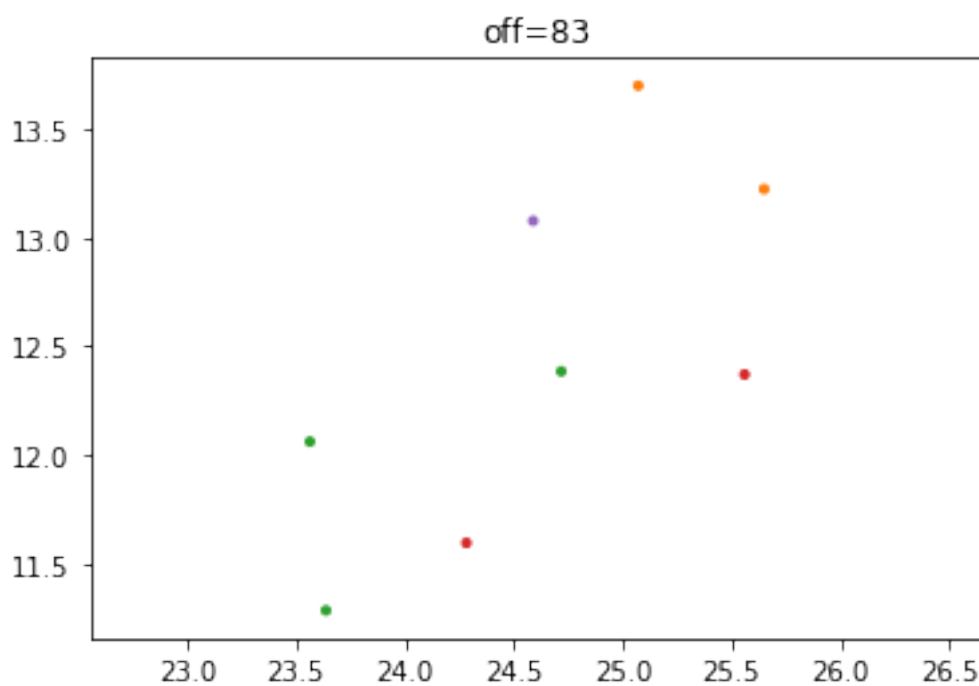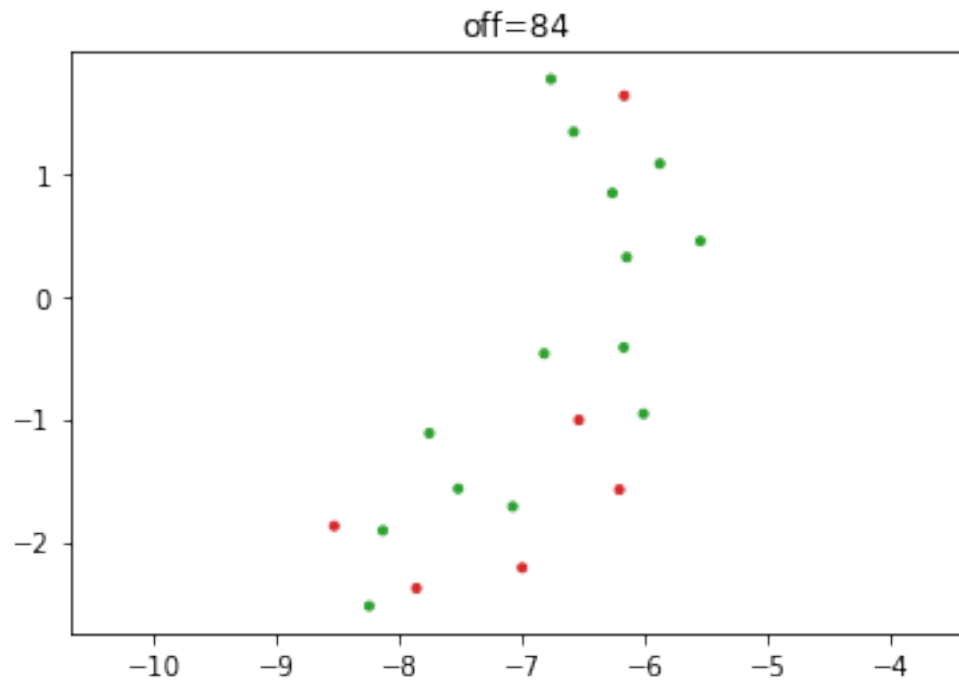n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1



off=42



off=43

/home/ethan/.local/lib/python3.7/site-packages/umap/spectral.py:229: UserWarning:

Embedding a total of 8 separate connected components using meta-embedding (experimental)

/home/ethan/.local/lib/python3.7/site-packages/sklearn/manifold/spectral_embedding_.py:235: Use

Graph is not fully connected, spectral embedding may not work as expected.

off=44

off=45

/home/ethan/.local/lib/python3.7/site-packages/umap/spectral.py:229: UserWarning:

Embedding a total of 3 separate connected components using meta-embedding (experimental)
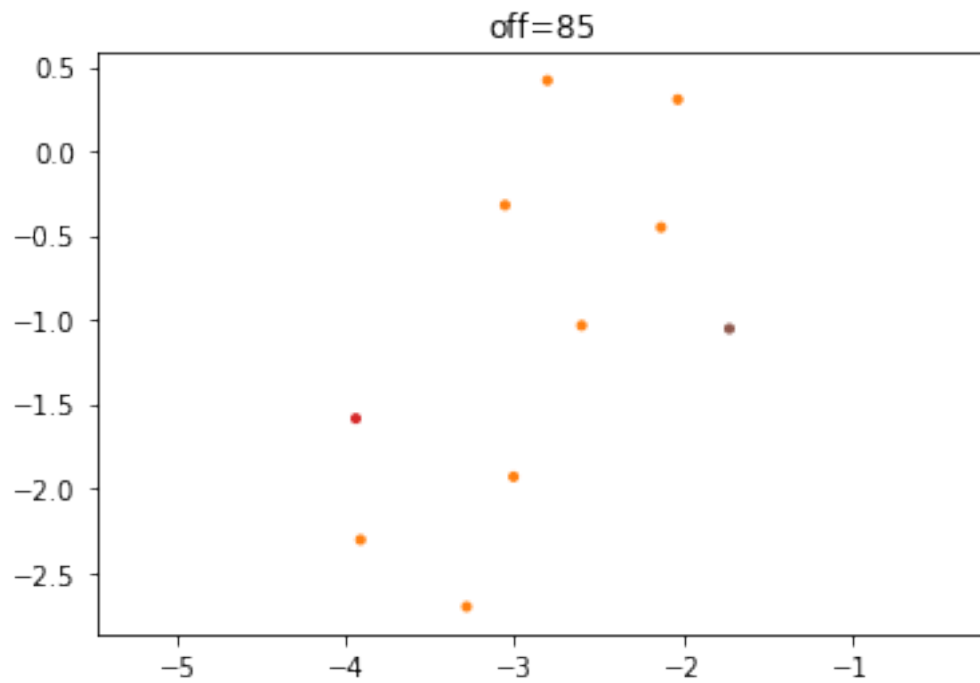


off=46

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

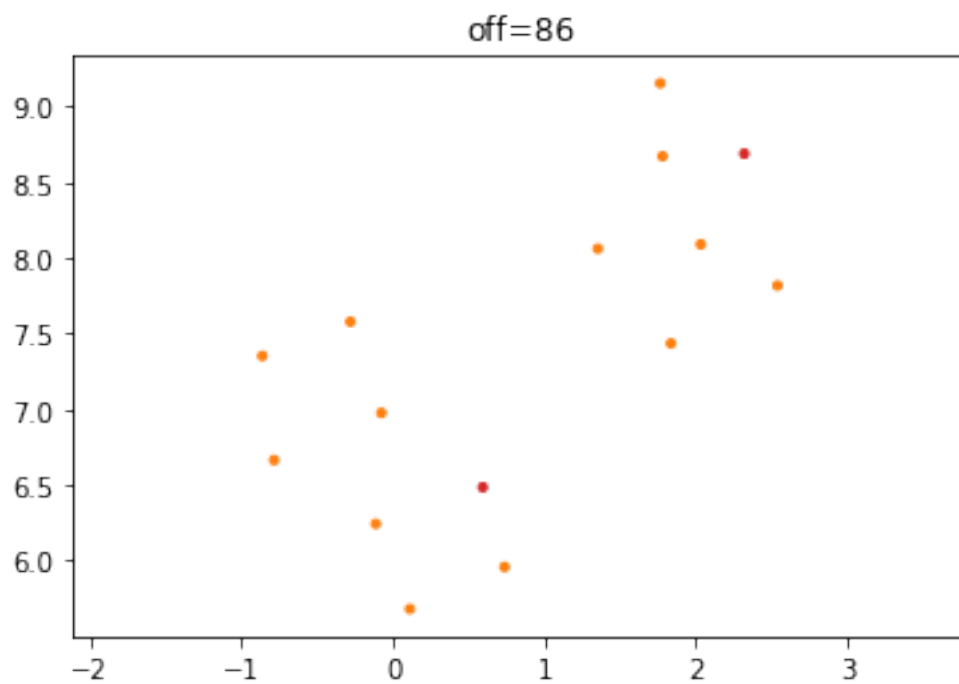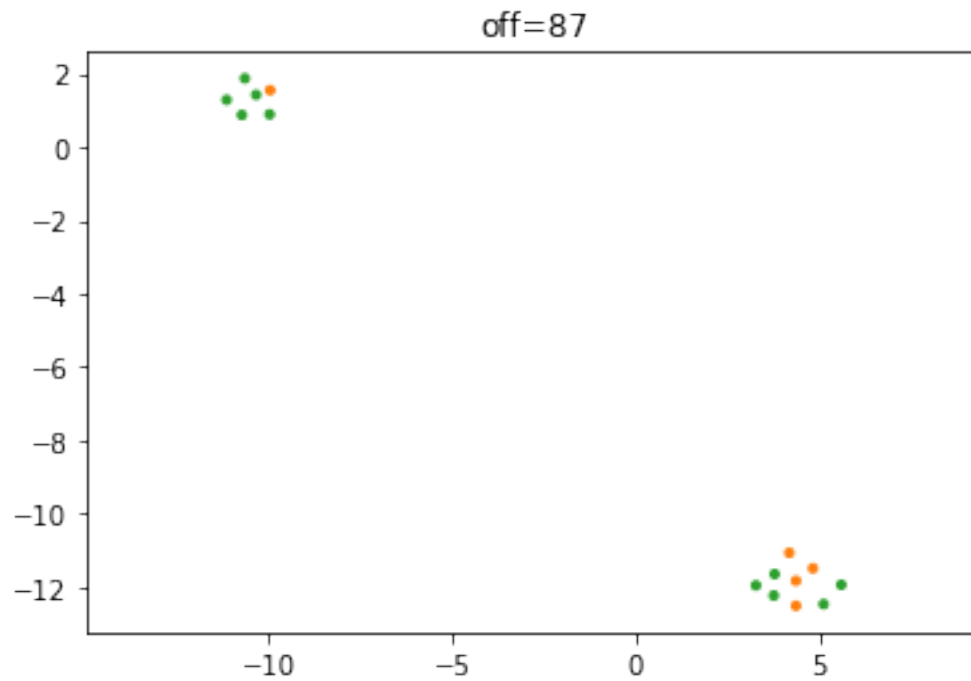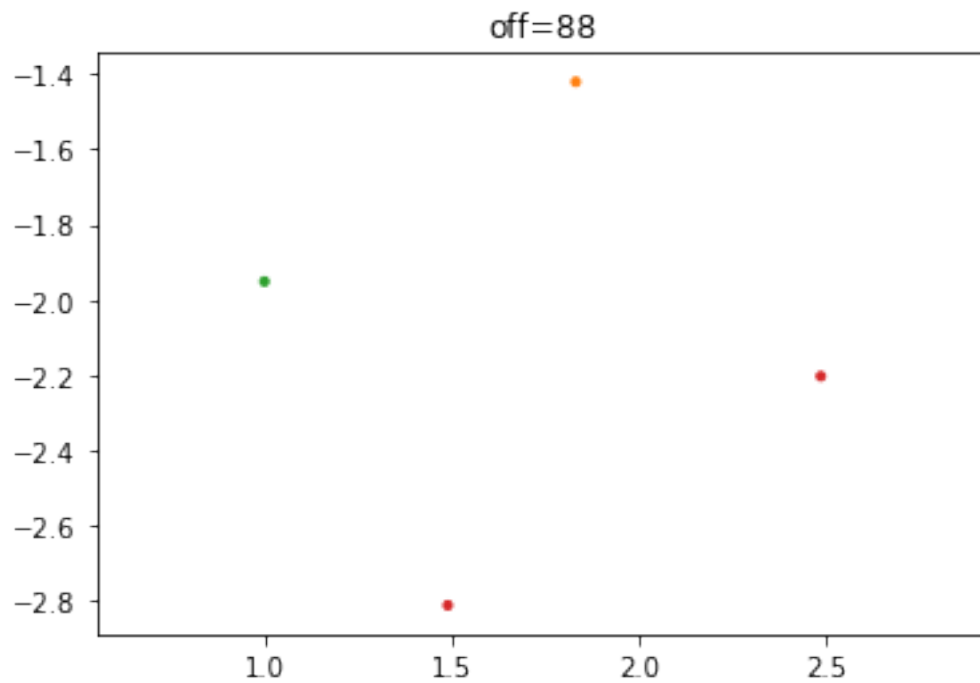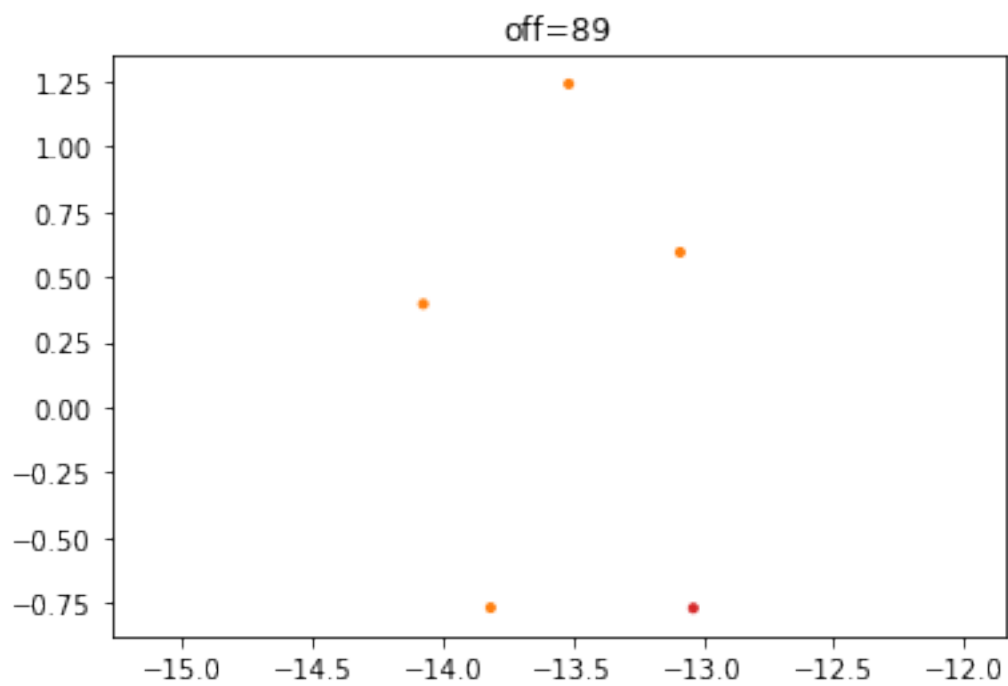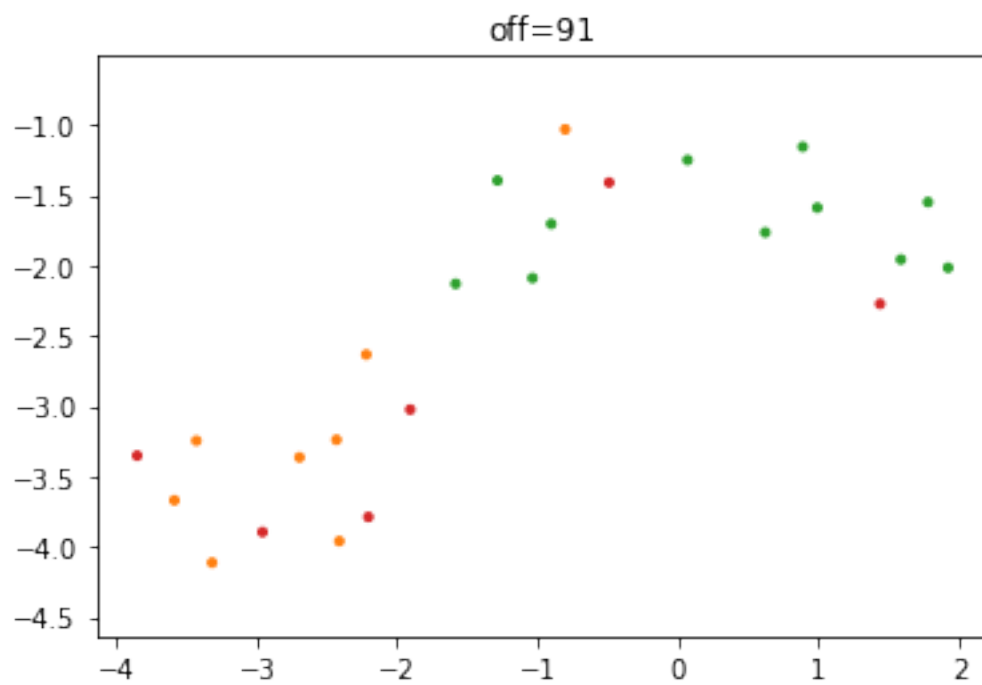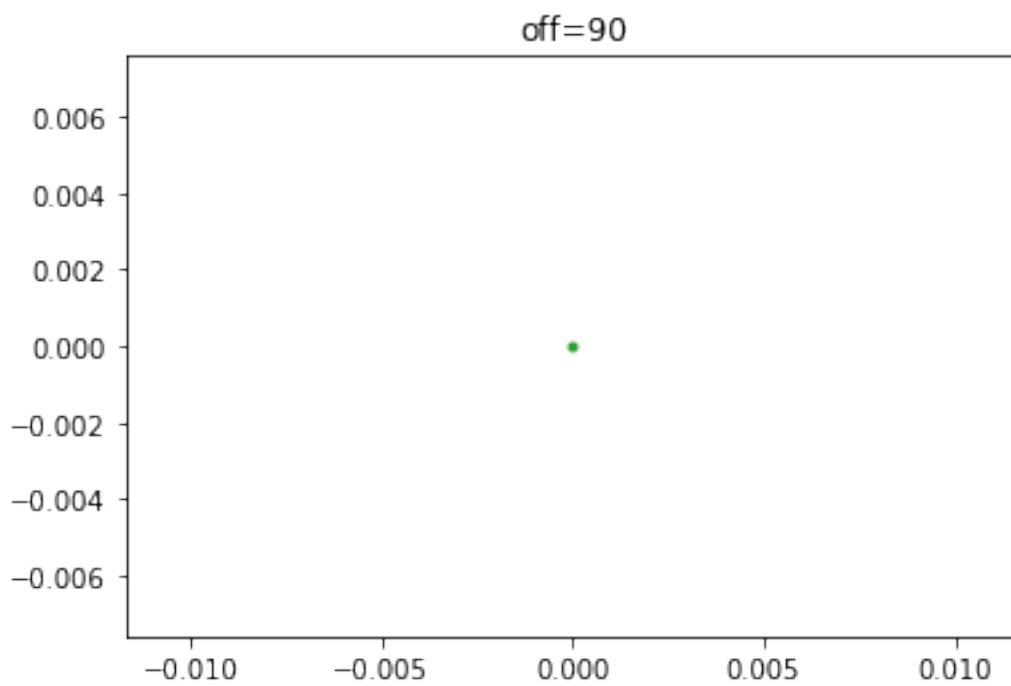n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1


off=47

/home/ethan/.local/lib/python3.7/site-packages/umap/spectral.py:229: UserWarning:

Embedding a total of 3 separate connected components using meta-embedding (experimental)

off=48

/home/ethan/.local/lib/python3.7/site-packages/umap/spectral.py:229: UserWarning:

Embedding a total of 6 separate connected components using meta-embedding (experimental)

/home/ethan/.local/lib/python3.7/site-packages/sklearn/manifold/spectral_embedding_.py:235: Use

Graph is not fully connected, spectral embedding may not work as expected.

off=49


off=50

off=51
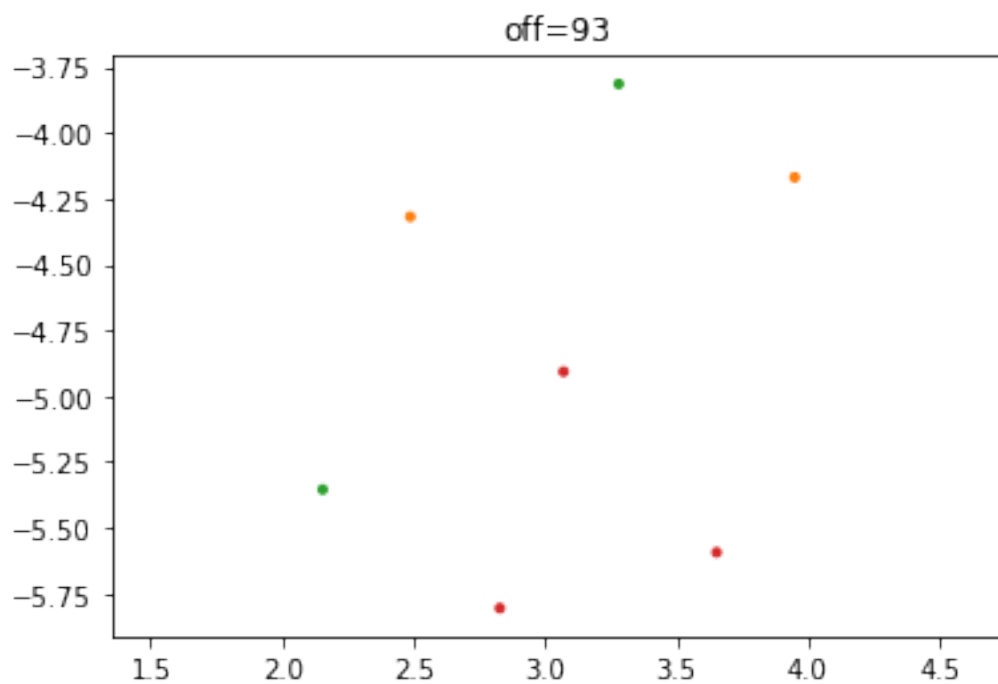
/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

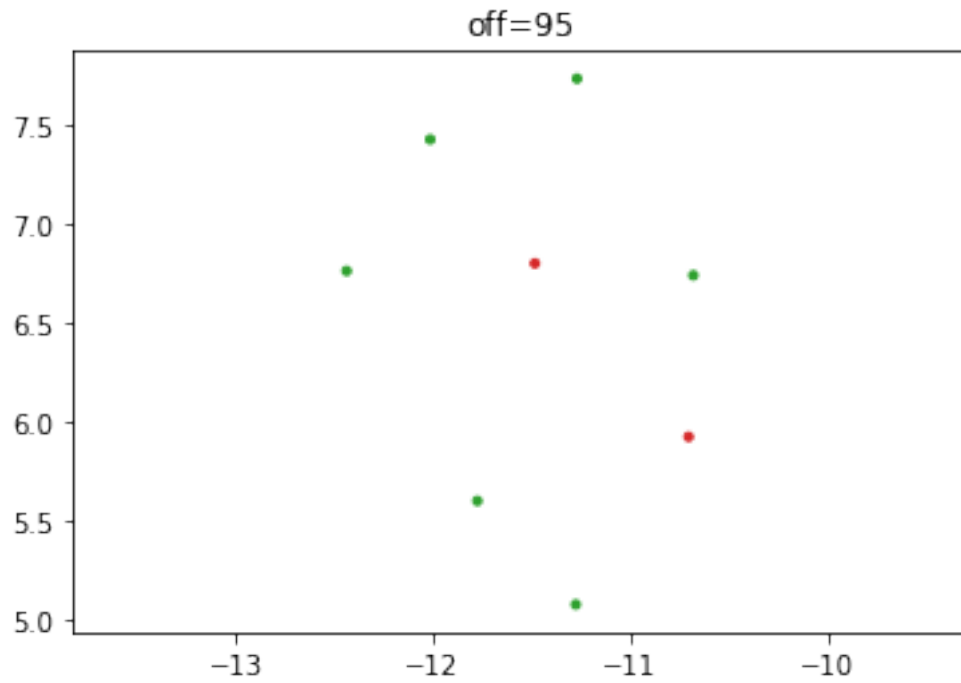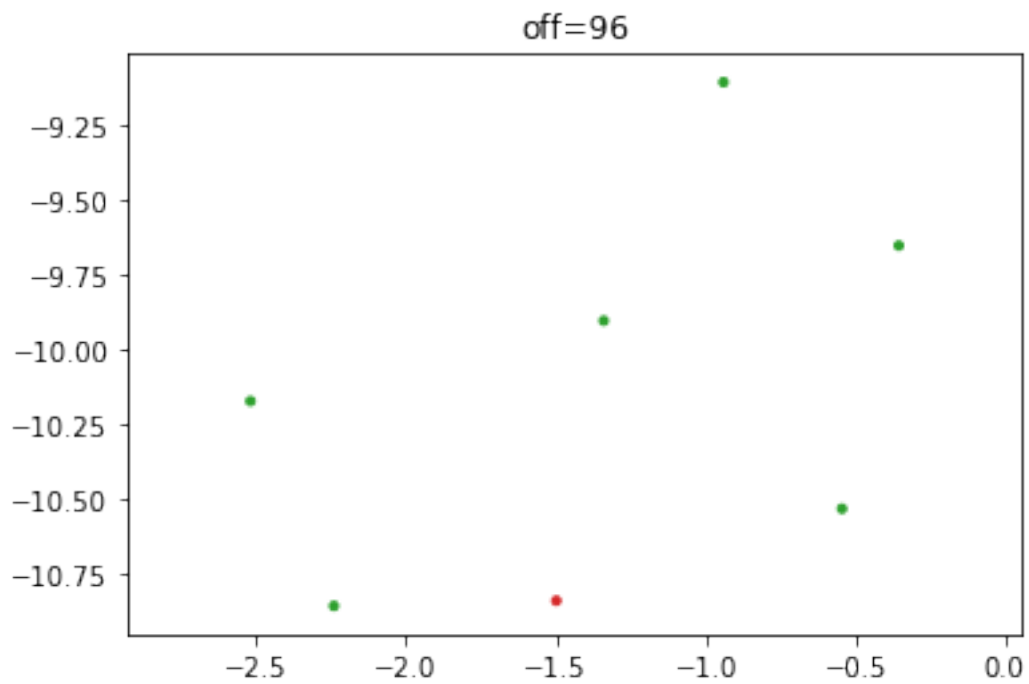n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1



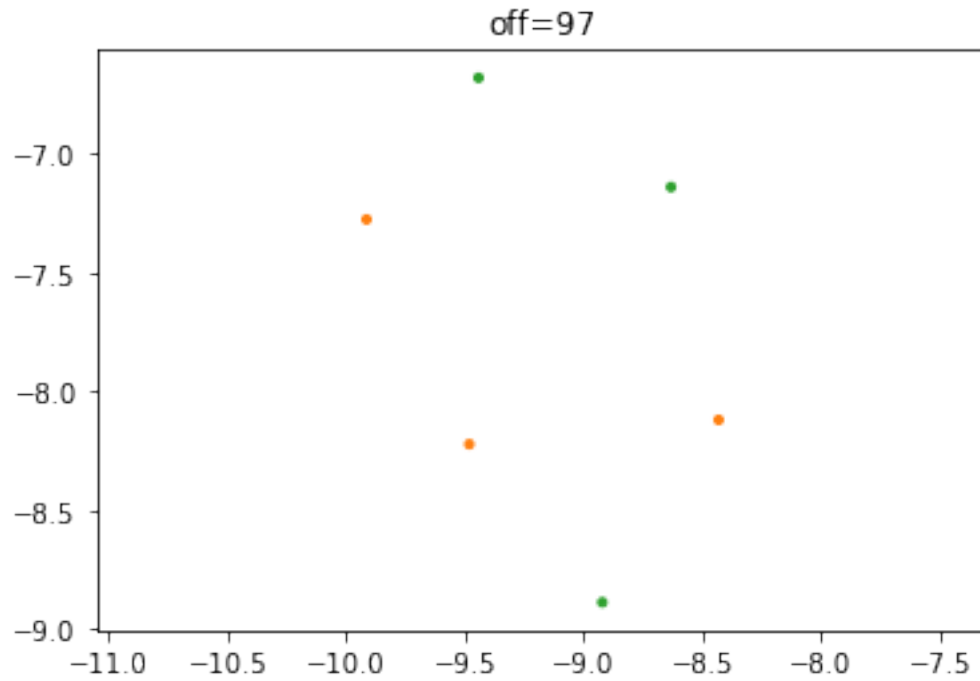off=52

off=53

/home/ethan/.local/lib/python3.7/site-packages/umap/spectral.py:229: UserWarning:

Embedding a total of 3 separate connected components using meta-embedding (experimental)

off=54

/home/ethan/.local/lib/python3.7/site-packages/umap/spectral.py:229: UserWarning:

Embedding a total of 2 separate connected components using meta-embedding (experimental)
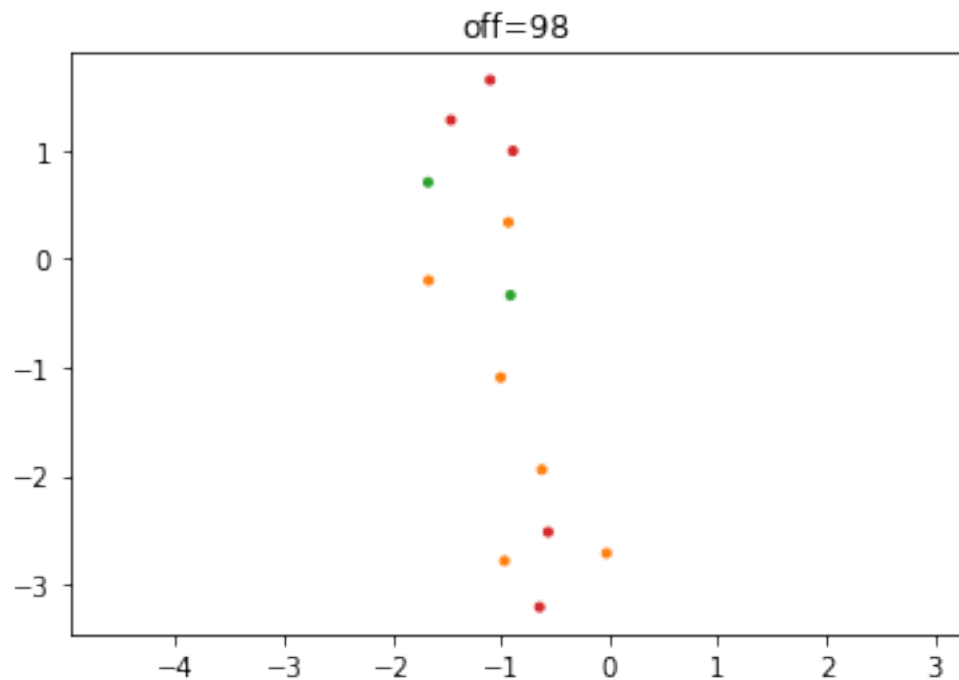


off=55

off=56

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

off=57

off=58

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

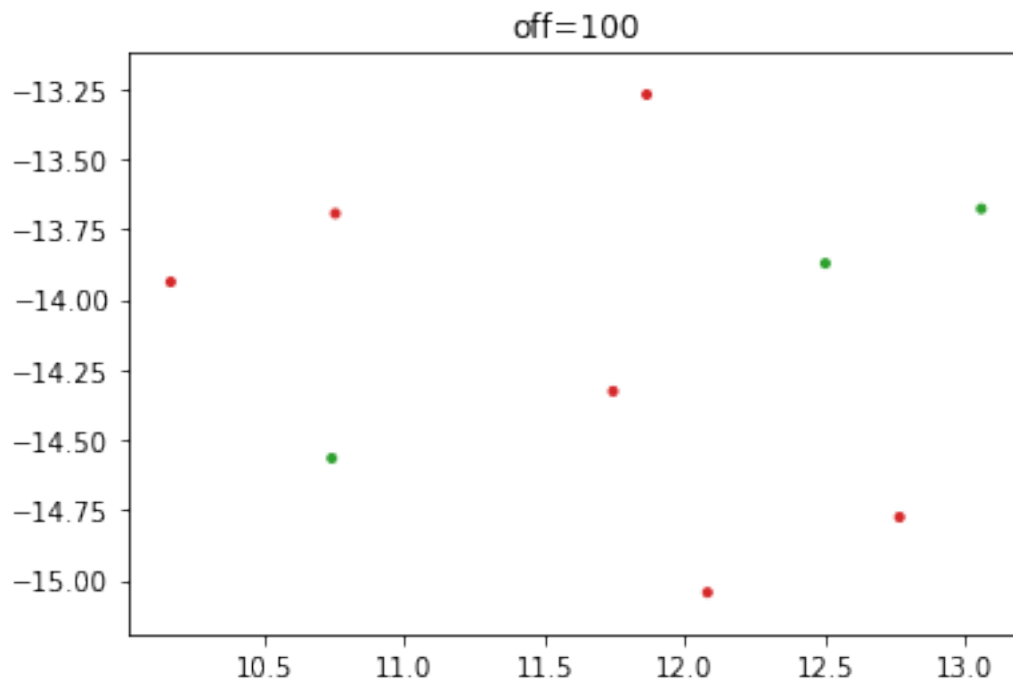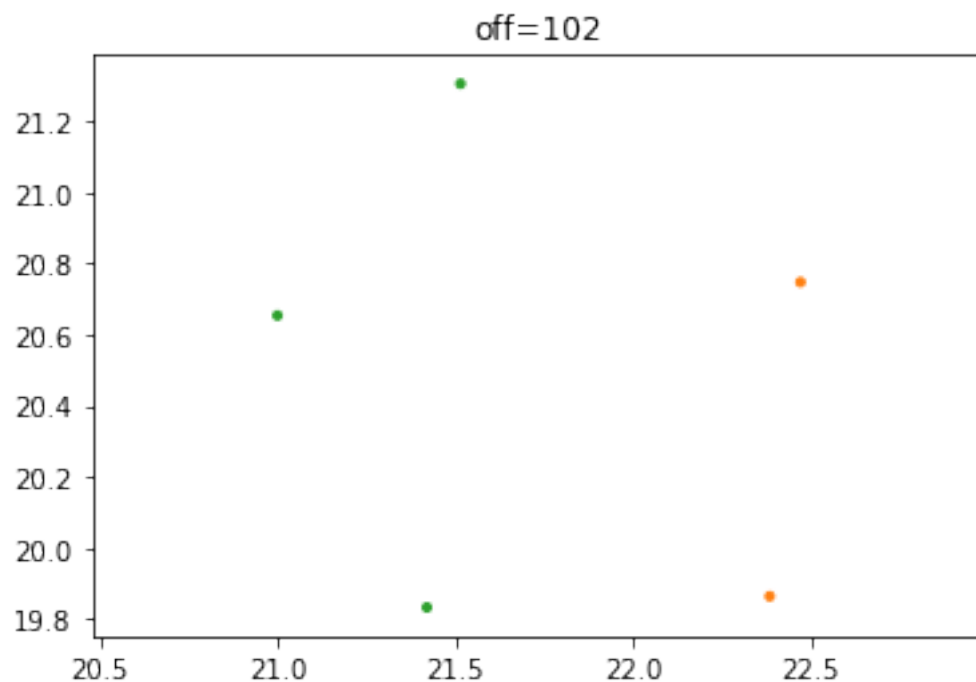/home/ethan/.local/lib/python3.7/site-packages/scipy/sparse/linalg/eigen/arpack/arpack.py:1592

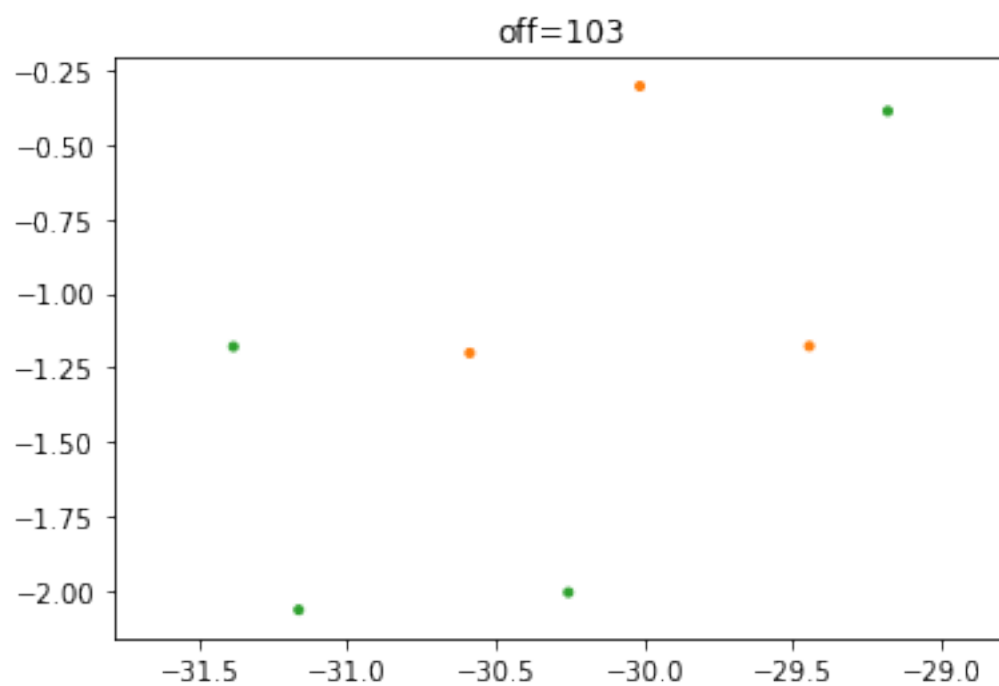k >= N for N * N square matrix. Attempting to use scipy.linalg.eigh instead.

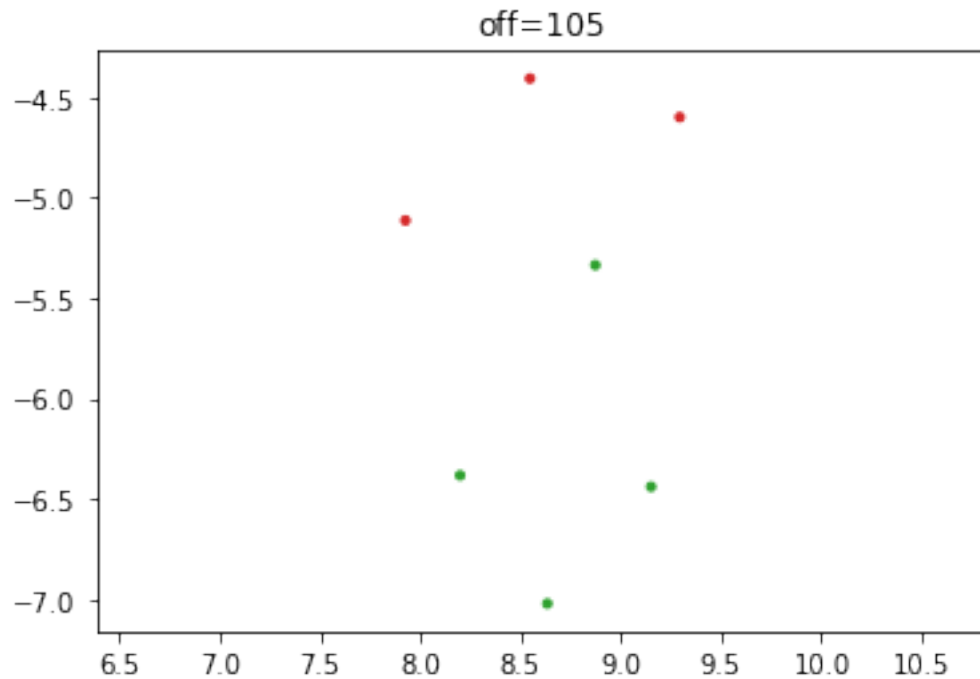/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

/home/ethan/.local/lib/python3.7/site-packages/scipy/sparse/linalg/eigen/arpack/arpack.py:1592

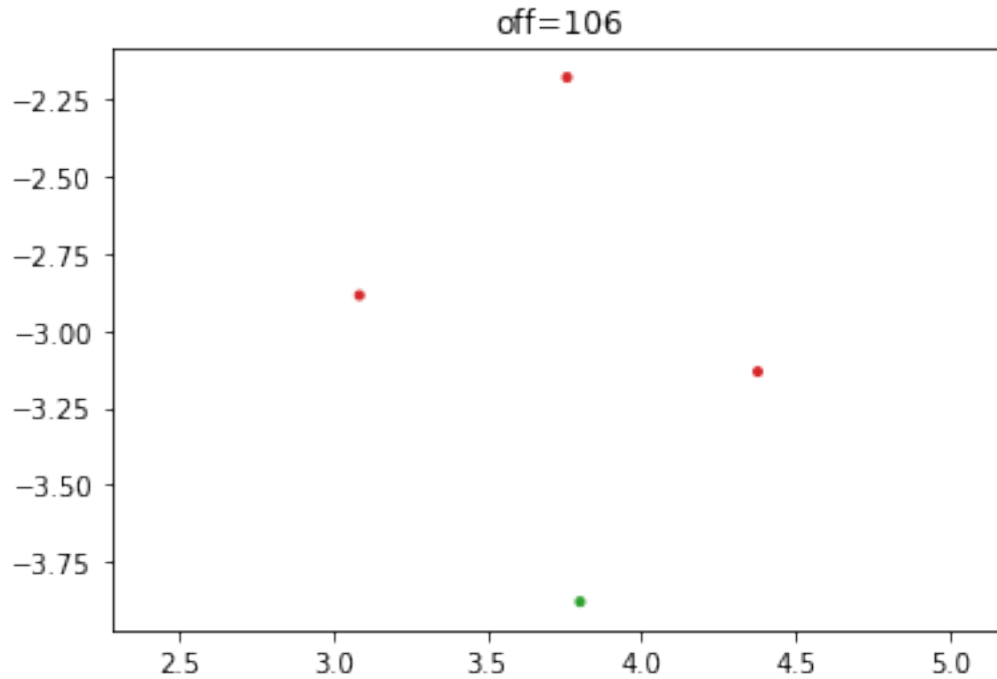k >= N for N * N square matrix. Attempting to use scipy.linalg.eigh instead.

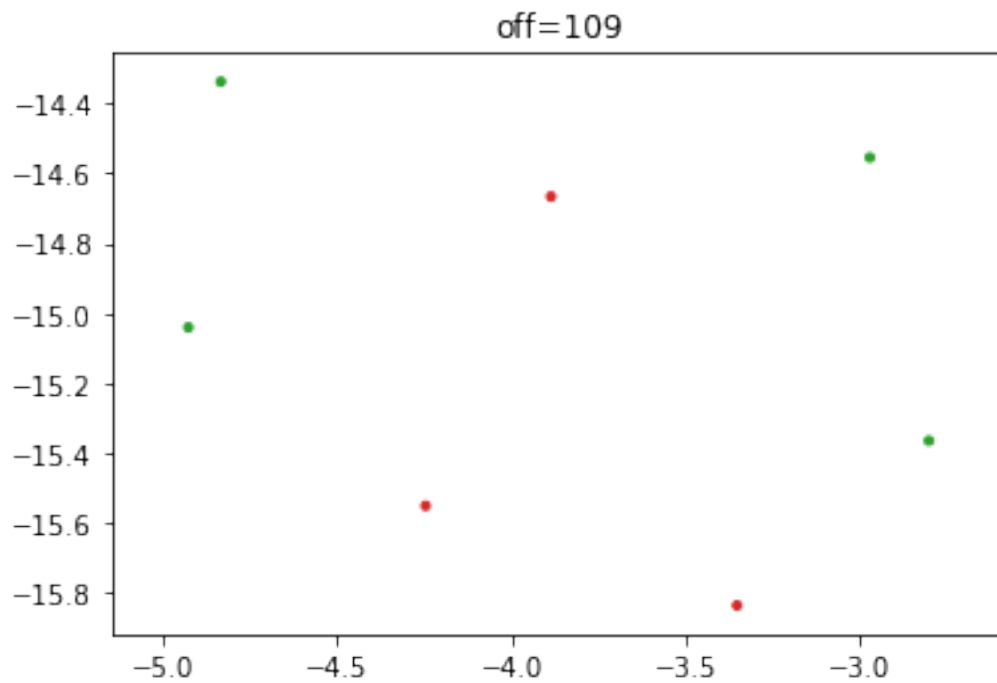/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

off=62



off=63

off=64

/home/ethan/.local/lib/python3.7/site-packages/umap/spectral.py:229: UserWarning:

Embedding a total of 3 separate connected components using meta-embedding (experimental)



off=65

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:
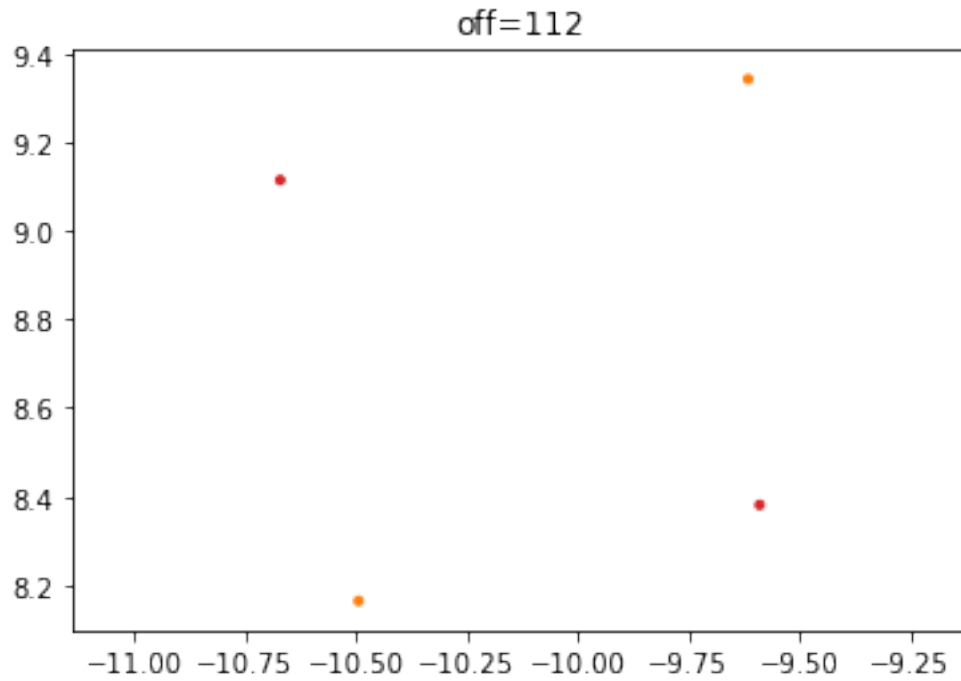
n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

off=66

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1
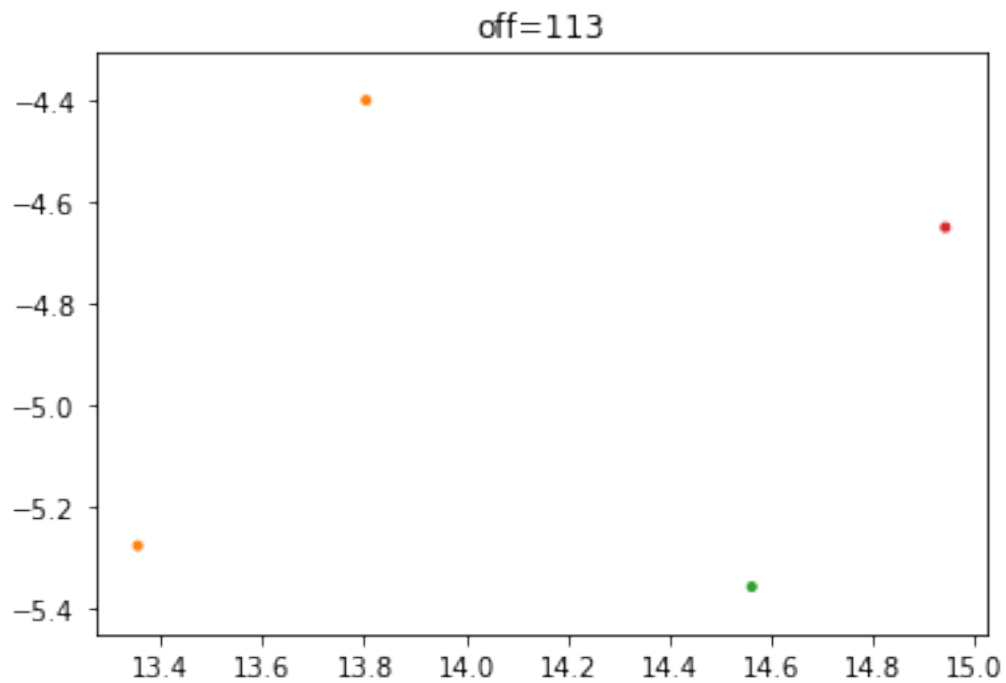
off=67



off=68

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1
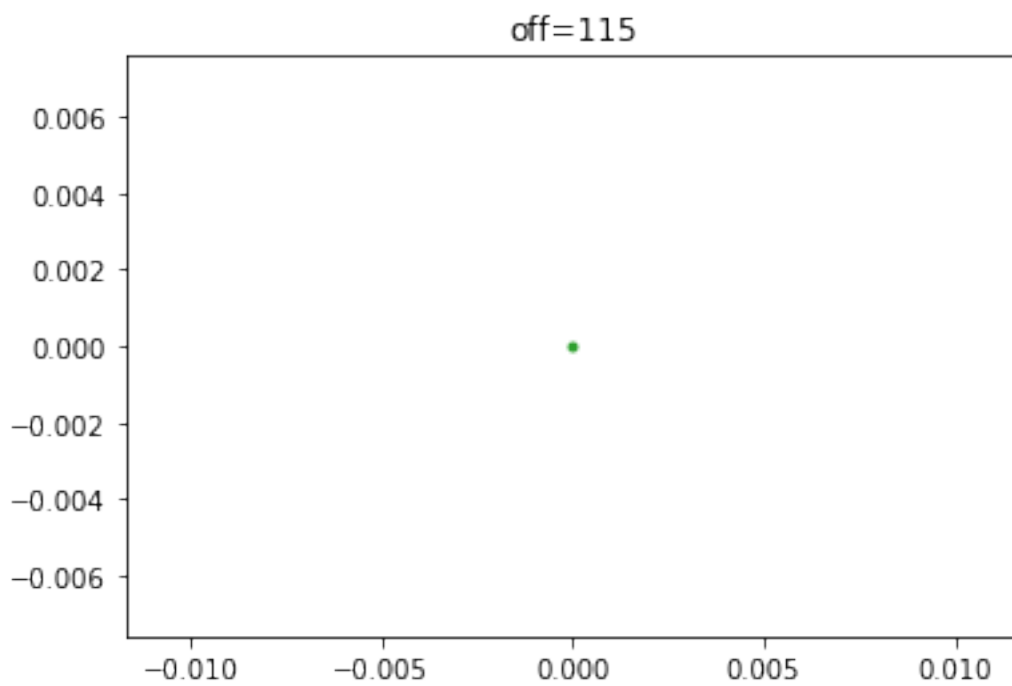


off=69

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

off=70



off=71

off=72



off=73

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

off=74

off=75

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

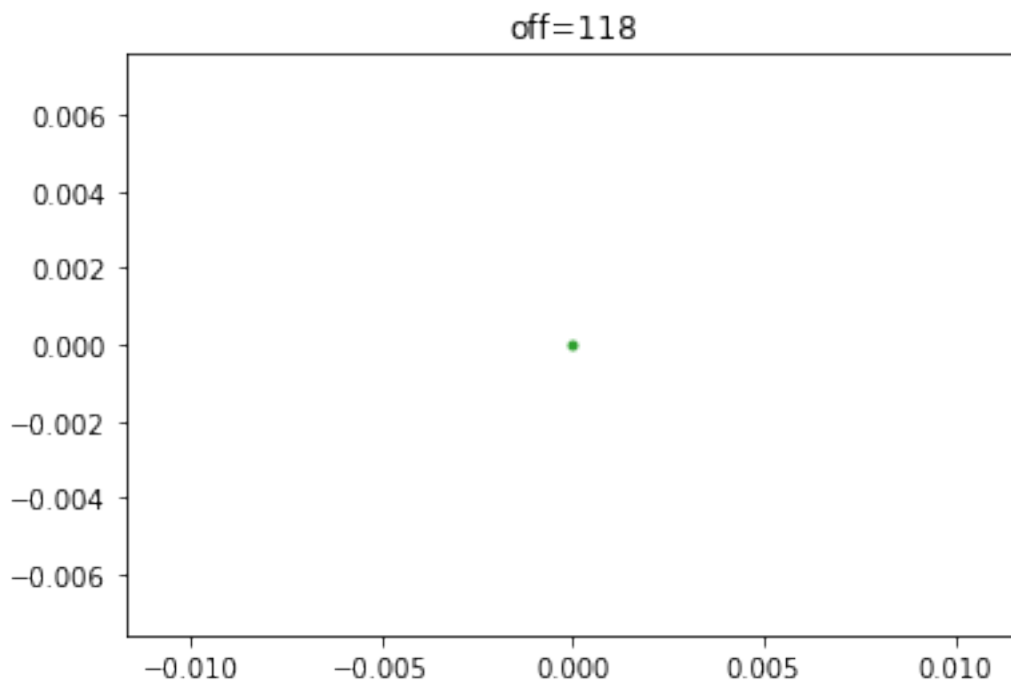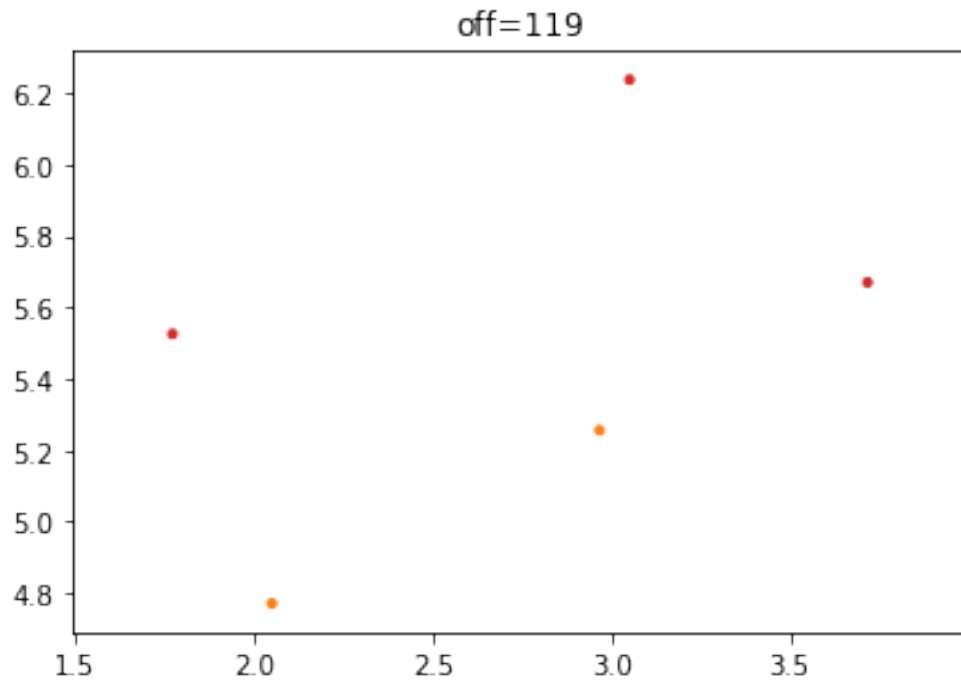n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1



off=76



off=77

off=78

## off=79



/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

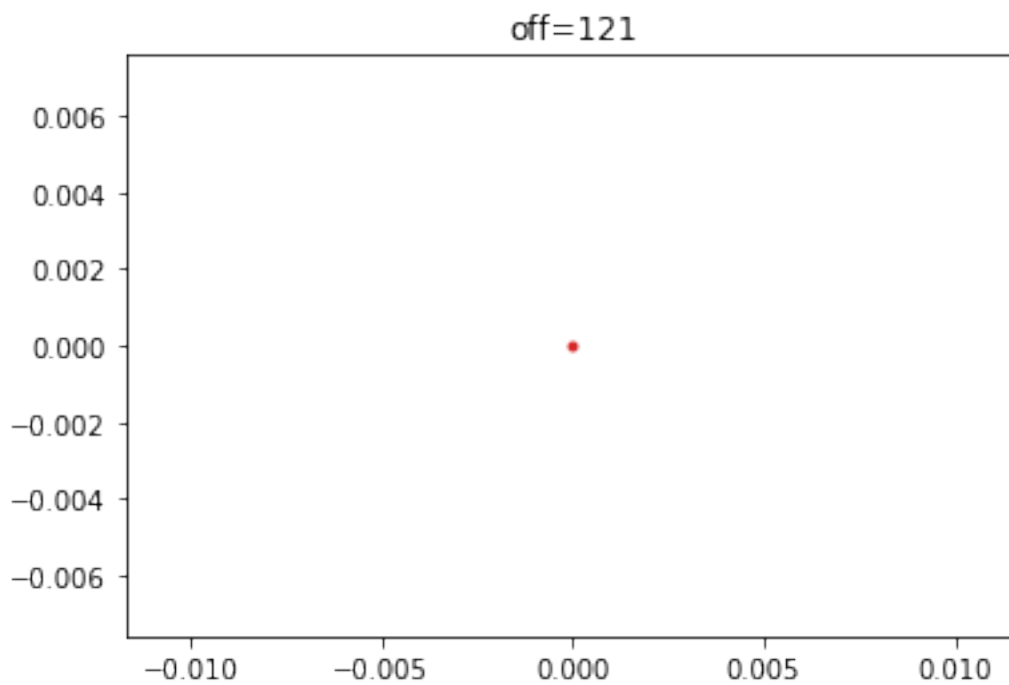n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

## off=80

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:
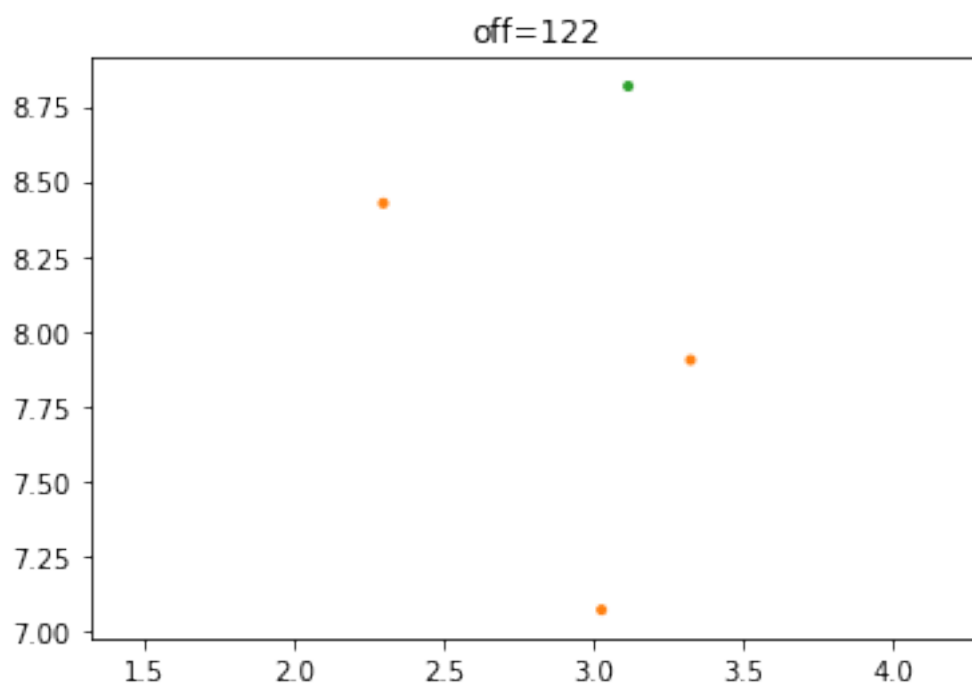
n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1



off=81

off=82

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1



off=83

off=84

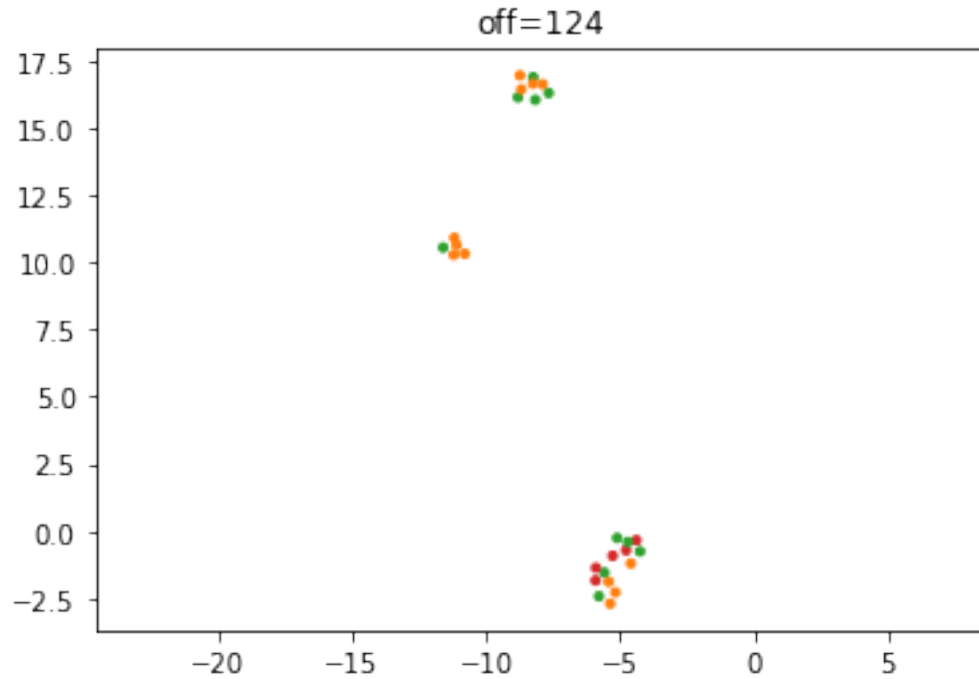/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

off=85

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1



off=86

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1



off=87

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:
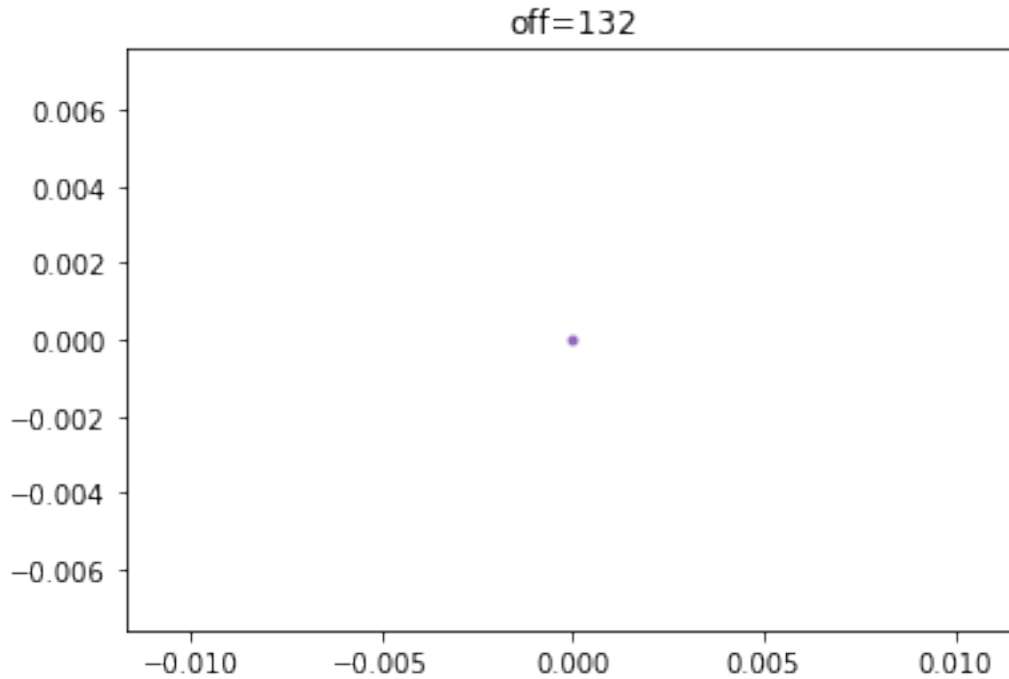
n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

off=88

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1



off=89

off=90



off=91

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:
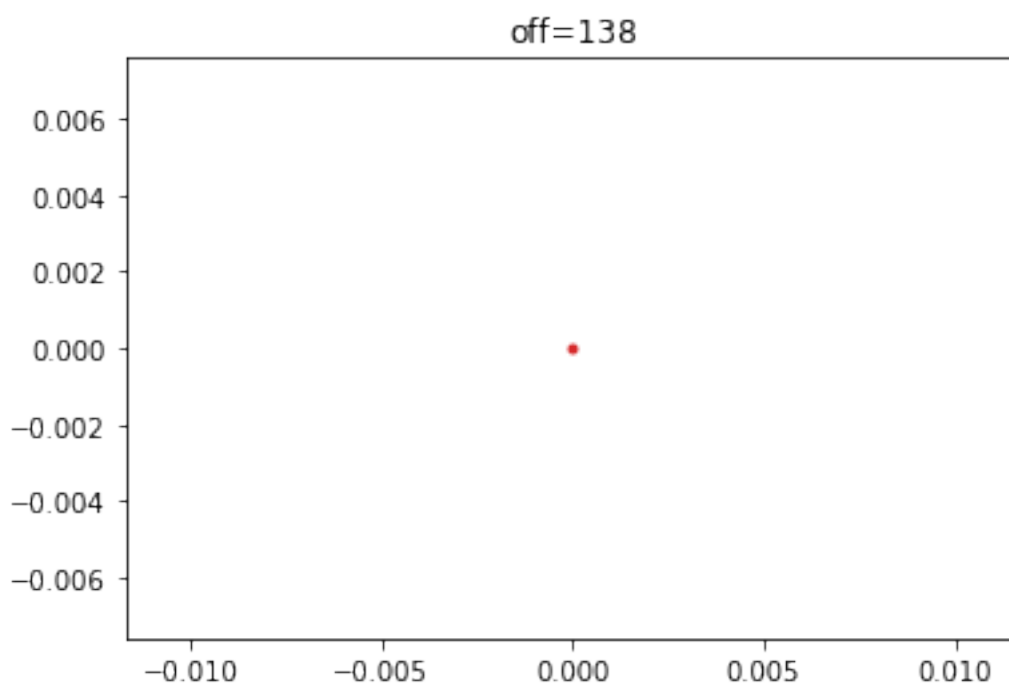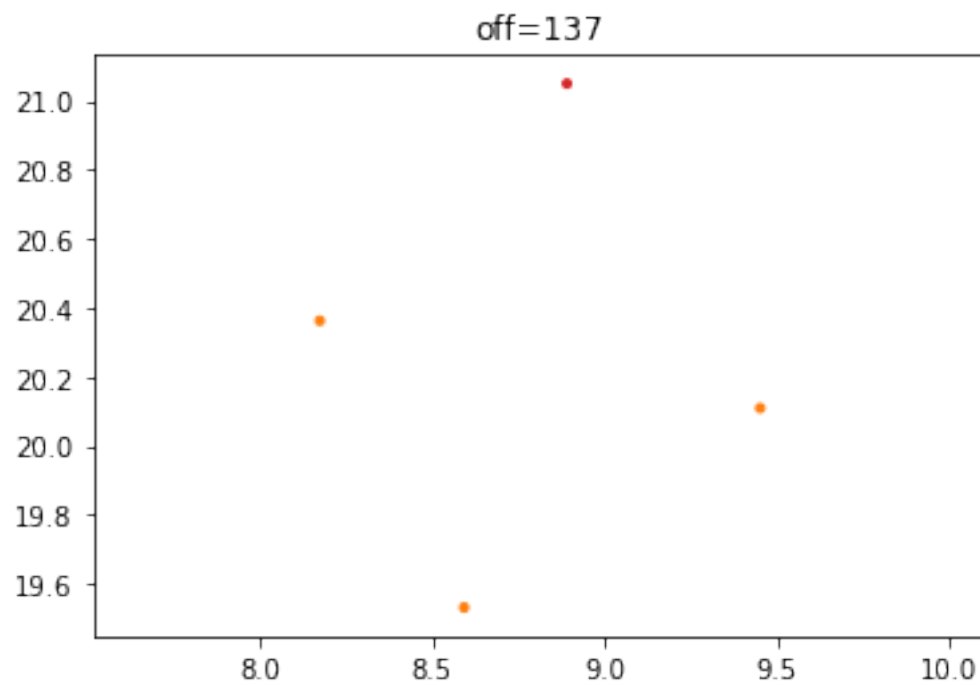
n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

/home/ethan/.local/lib/python3.7/site-packages/scipy/sparse/linalg/eigen/arpack/arpack.py:1592

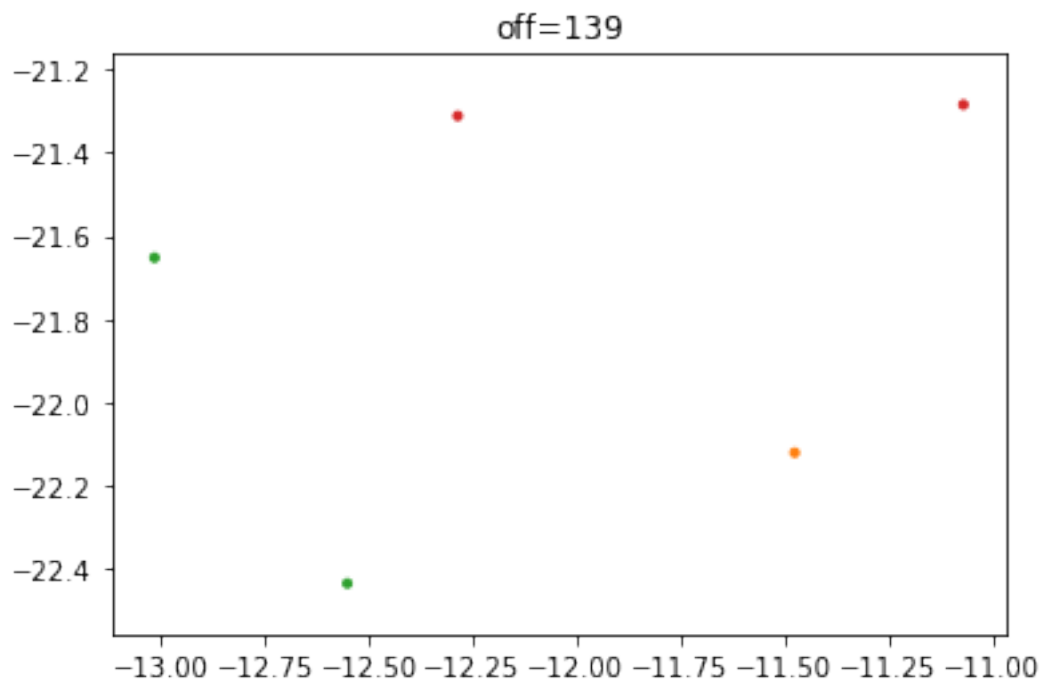k >= N for N * N square matrix. Attempting to use scipy.linalg.eigh instead.

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1



/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

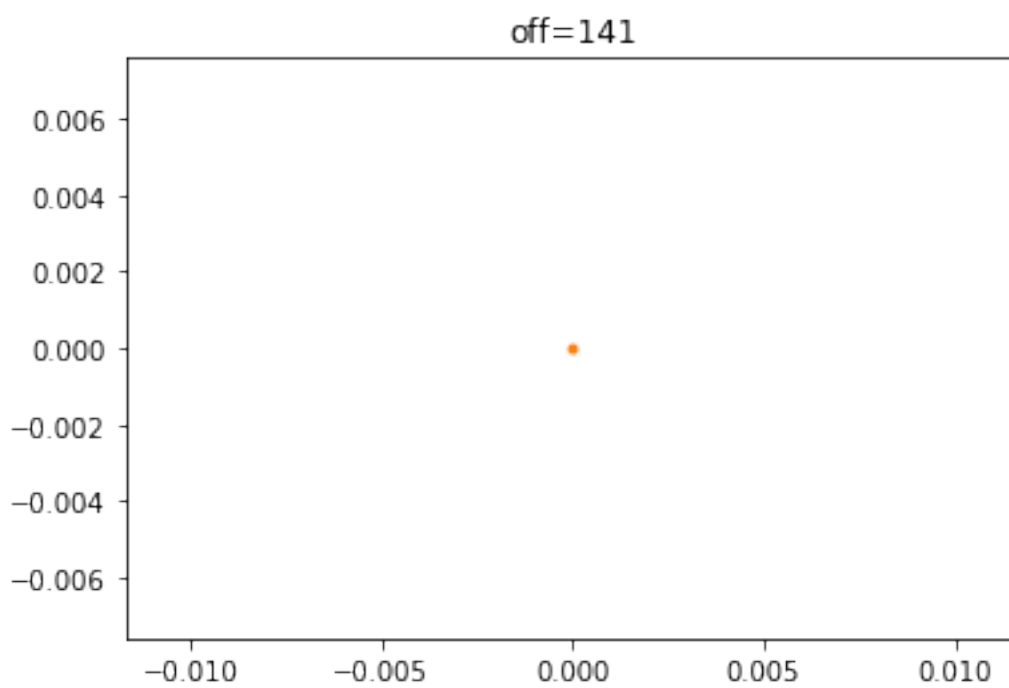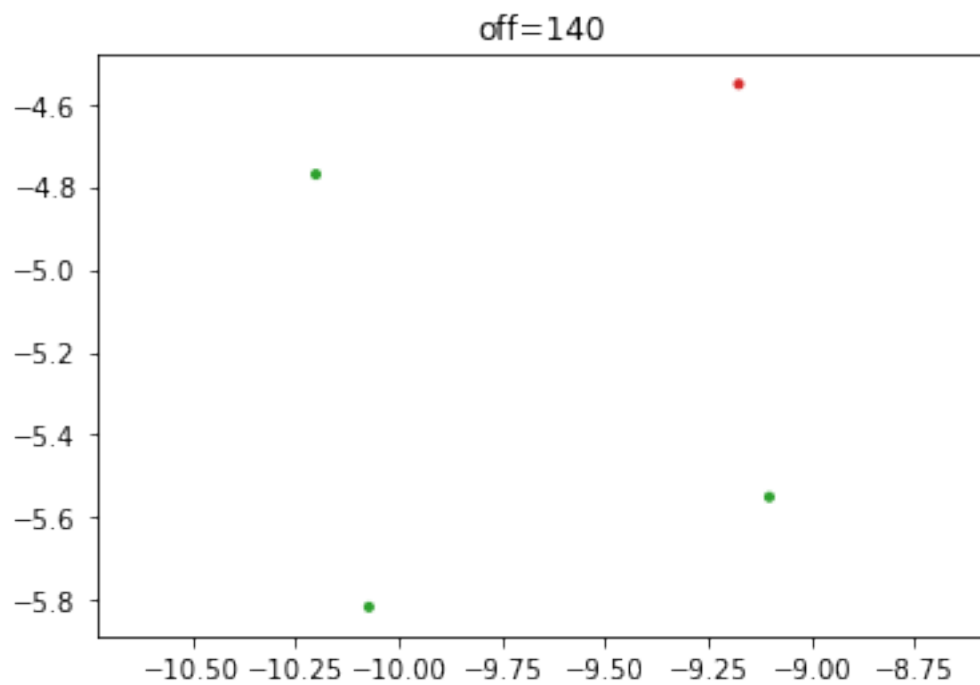n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

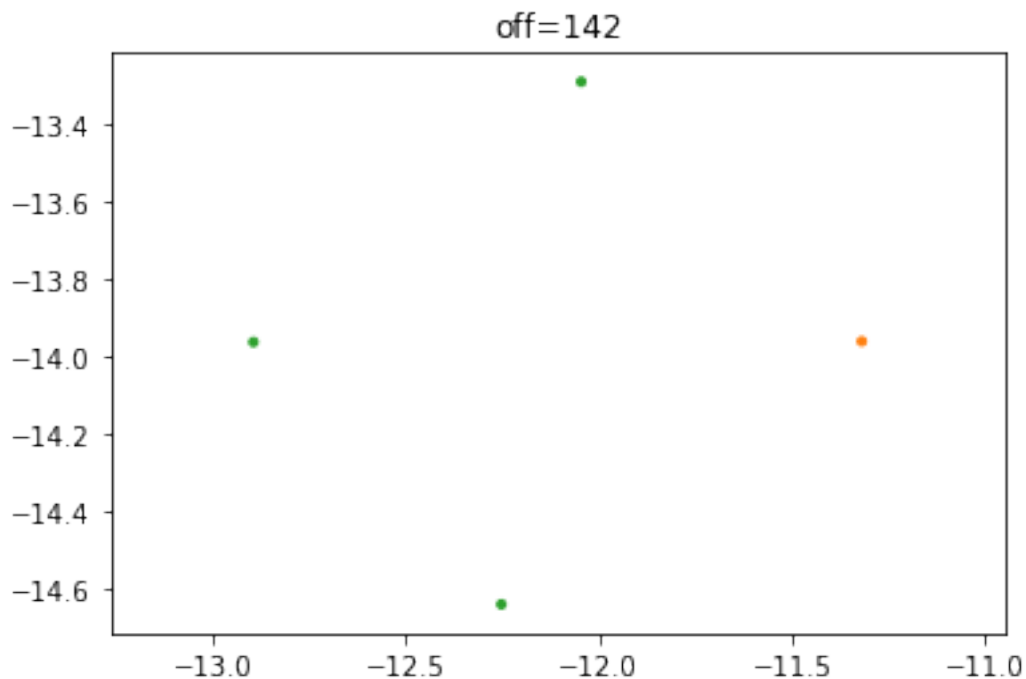n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

off=95

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1



off=96

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1



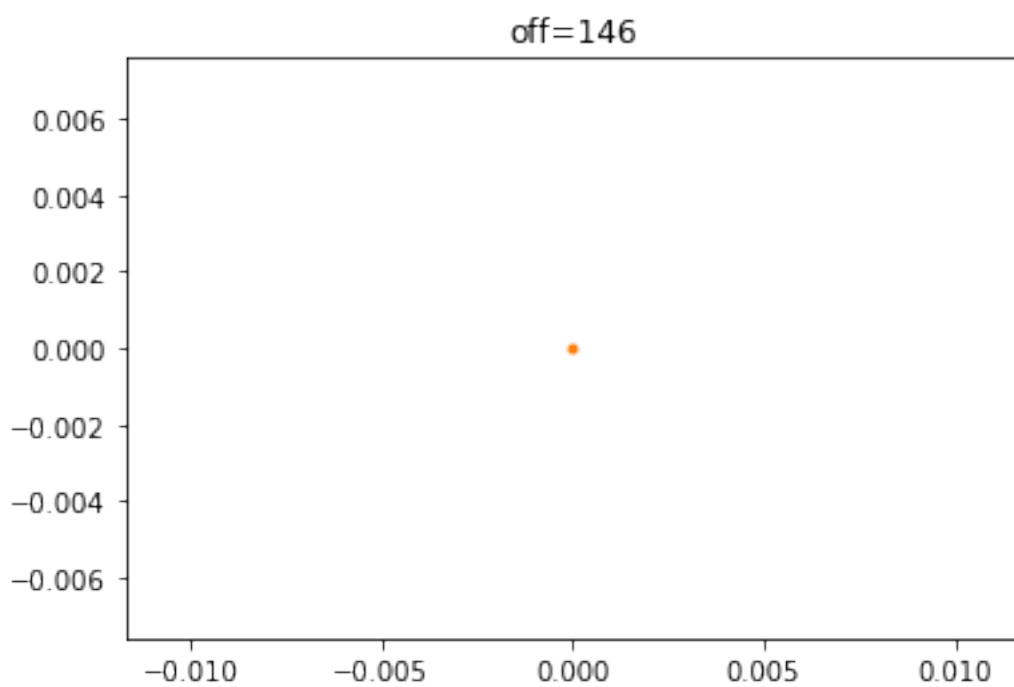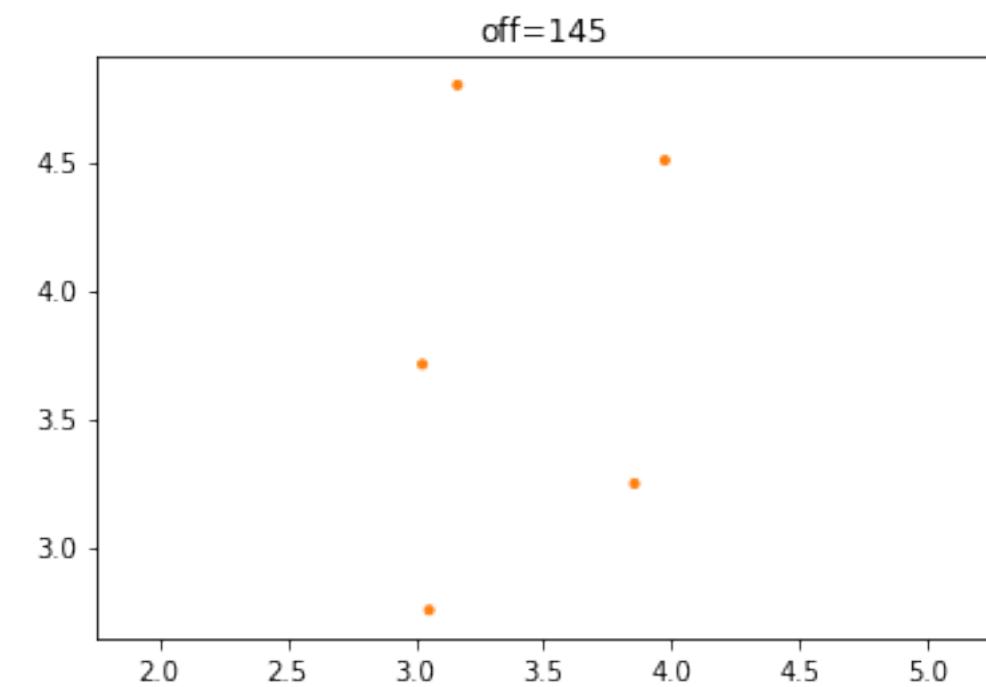/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

off=98

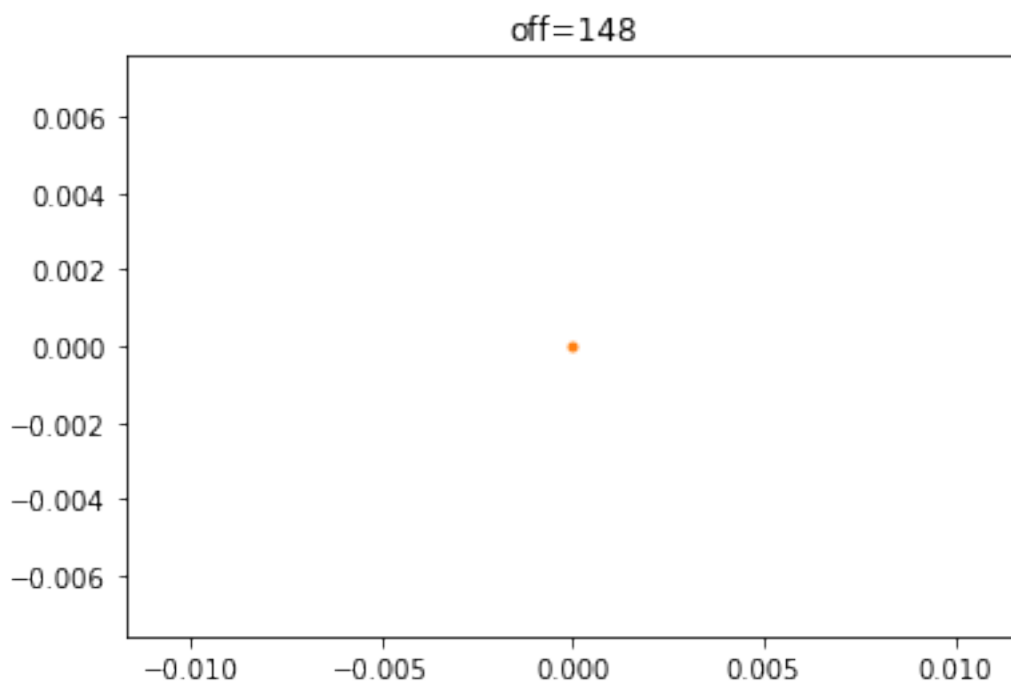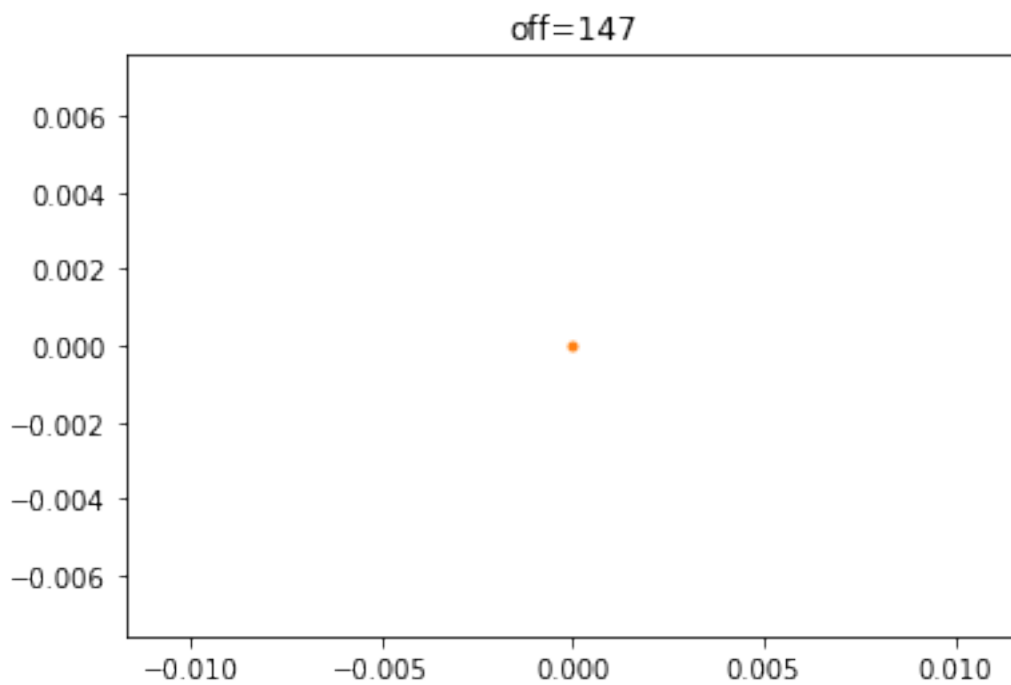/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

off=100

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

off=102

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1



off=103

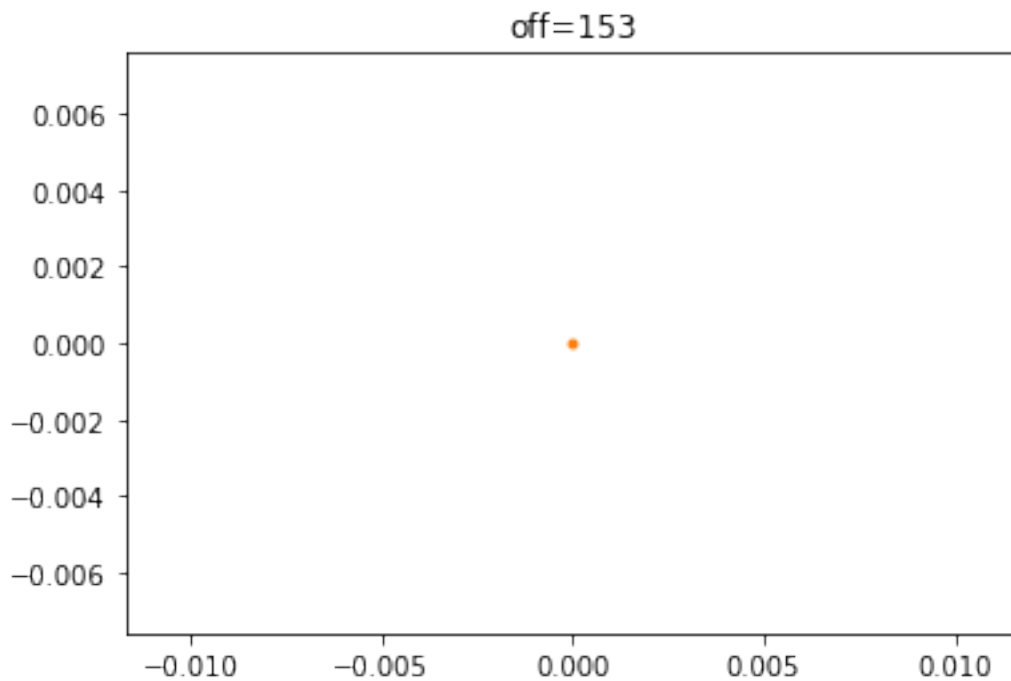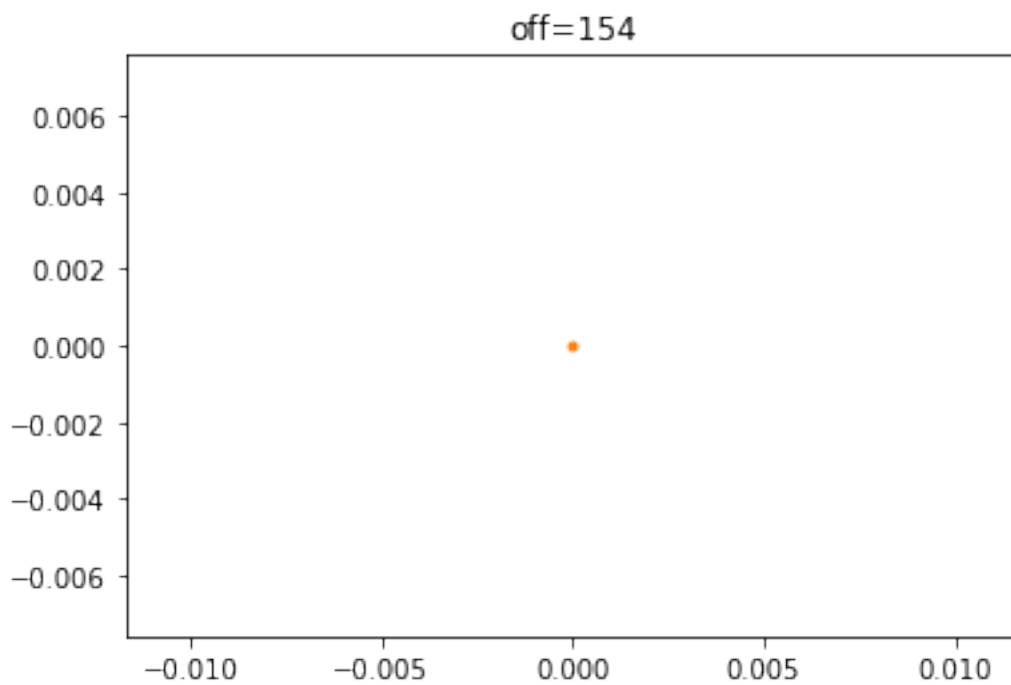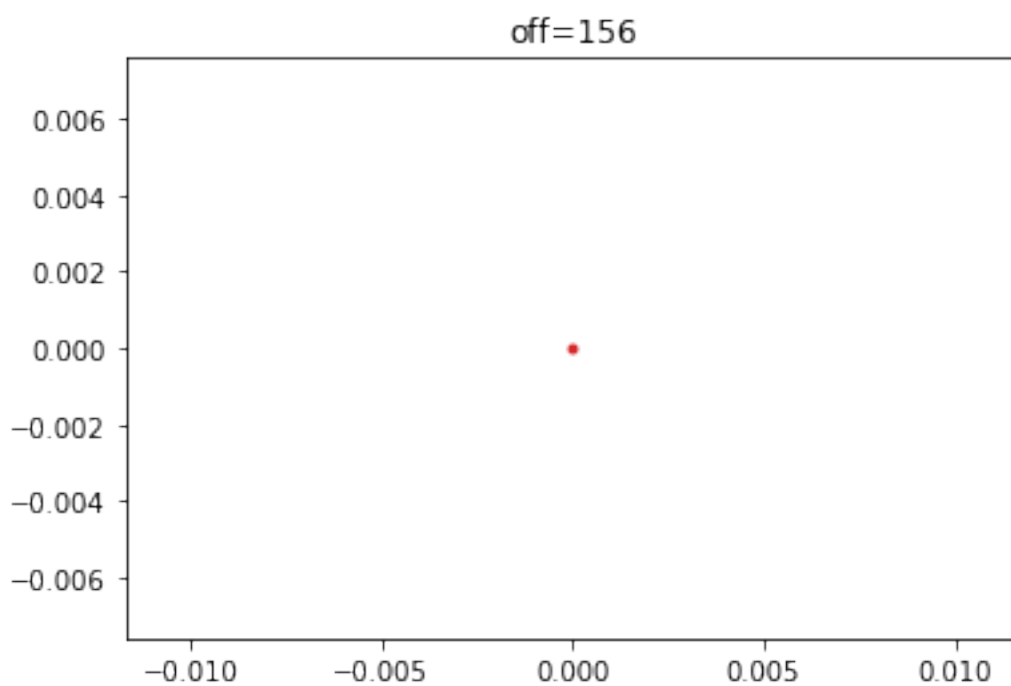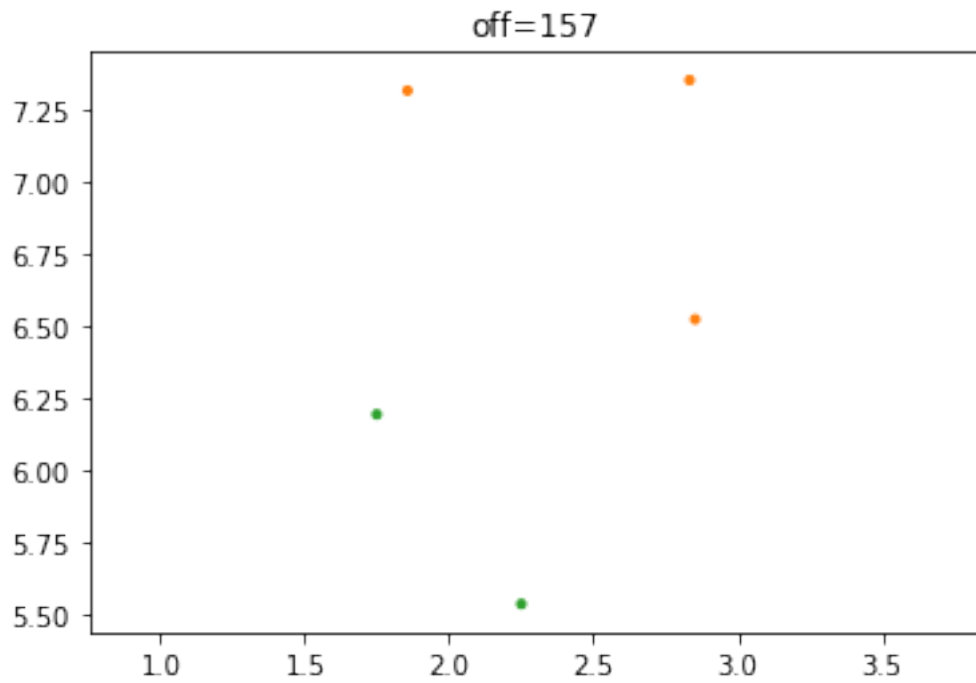/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1
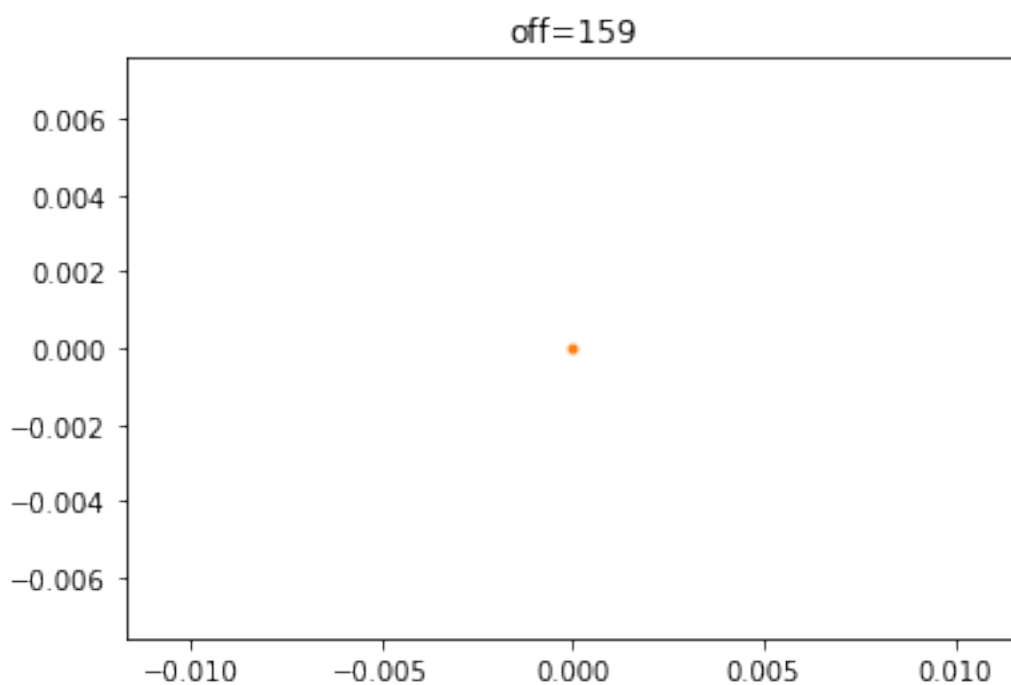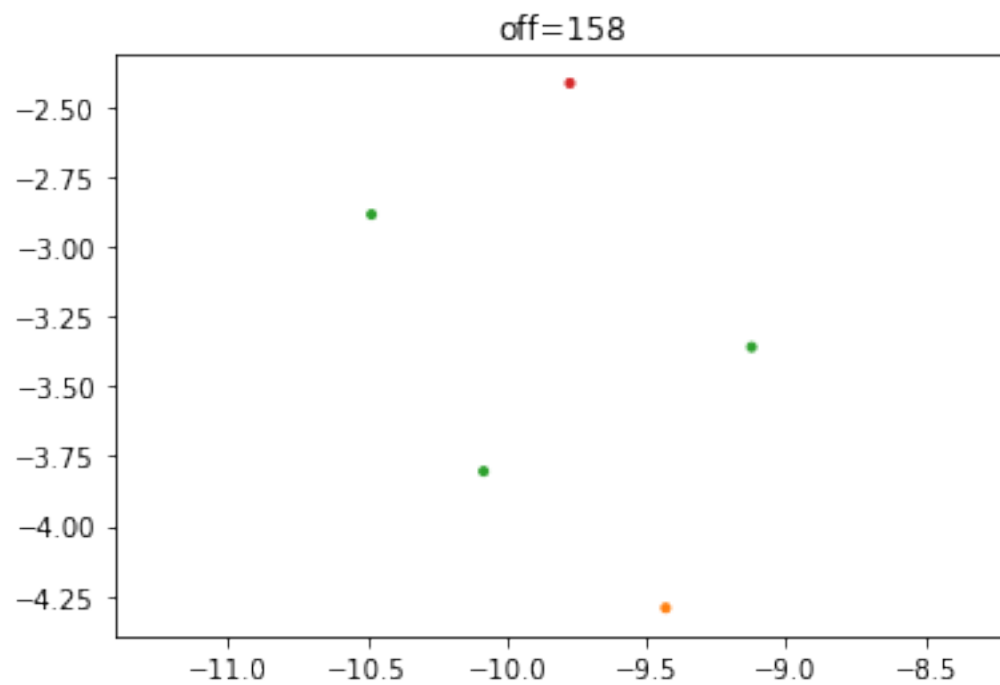
/home/ethan/.local/lib/python3.7/site-packages/scipy/sparse/linalg/eigen/arpack/arpack.py:1592

k >= N for N * N square matrix. Attempting to use scipy.linalg.eigh instead.

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

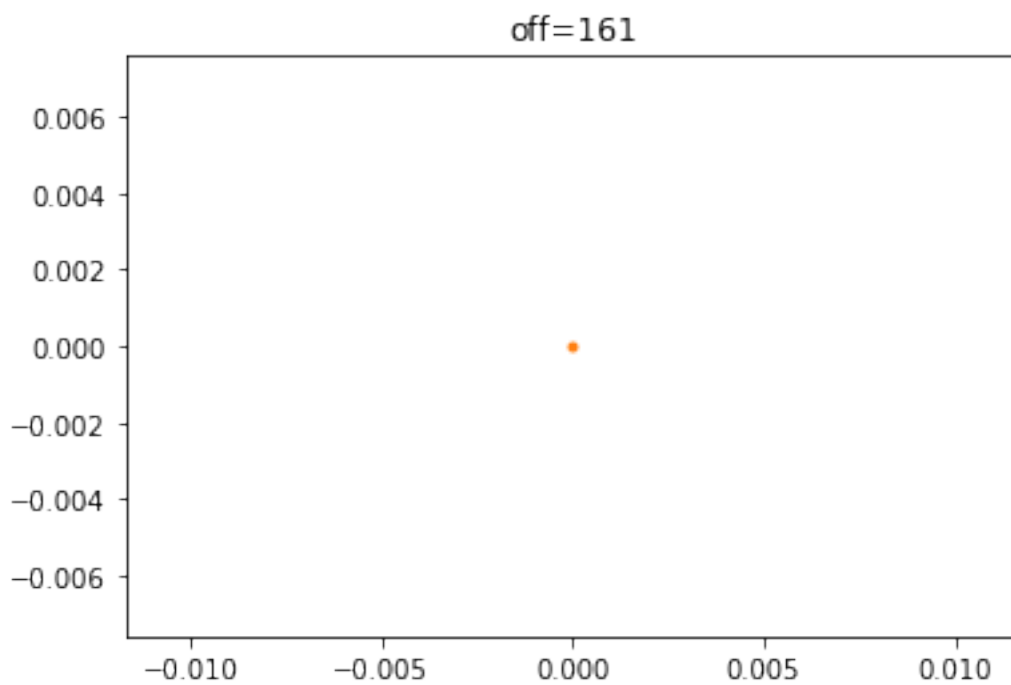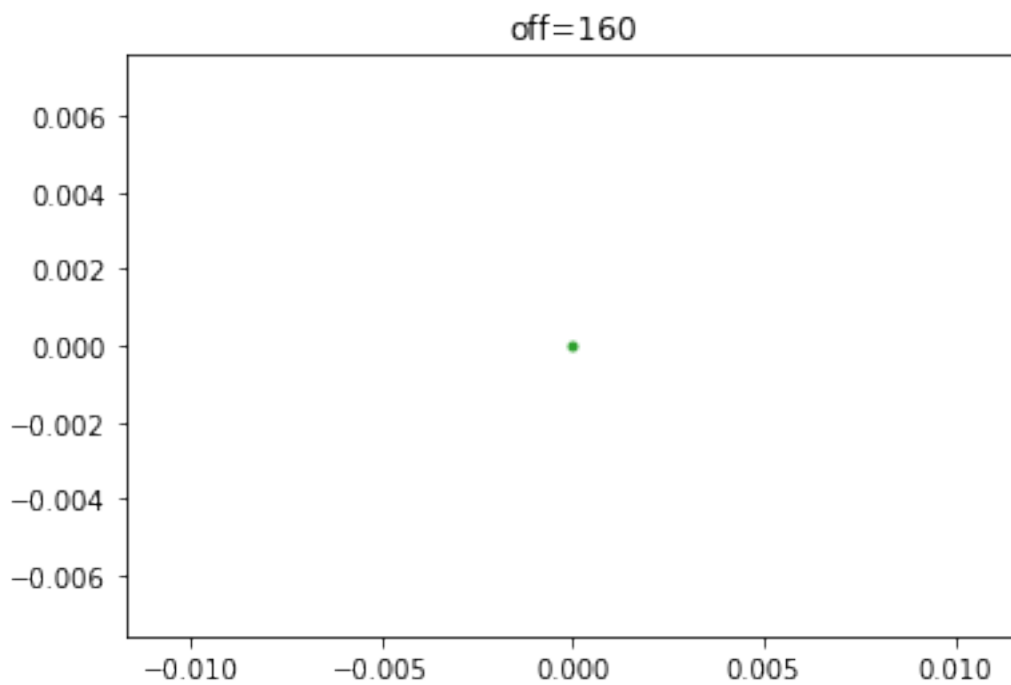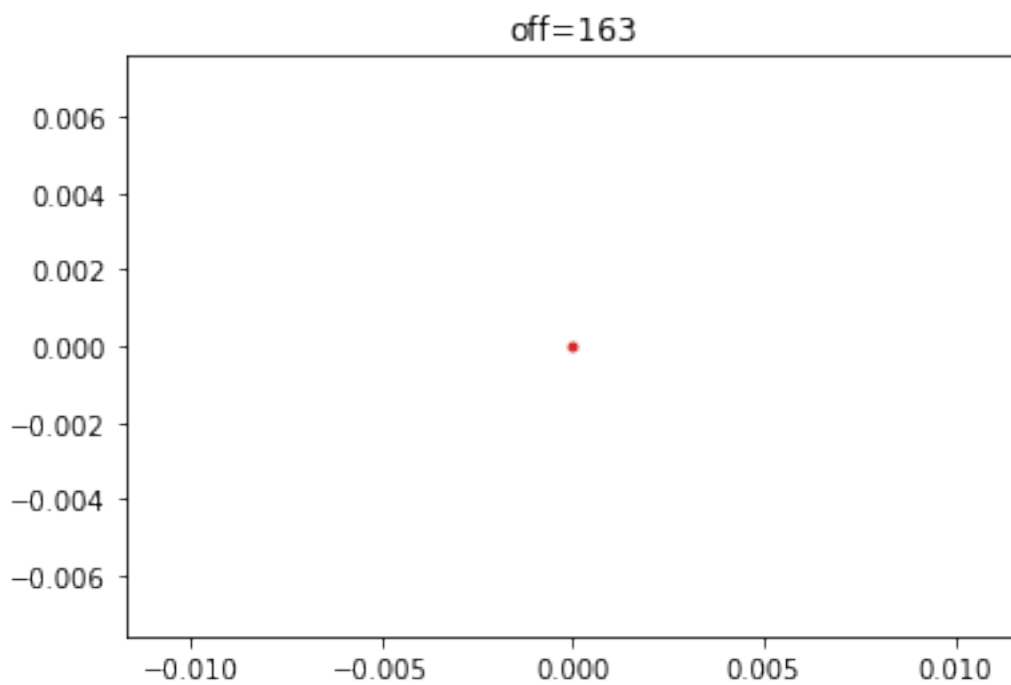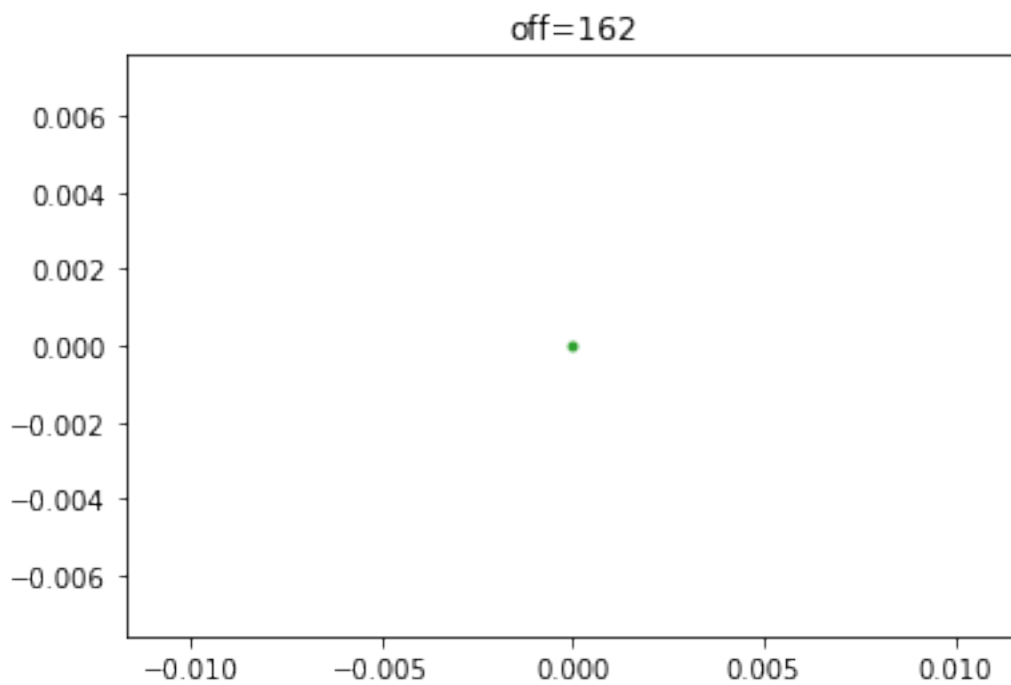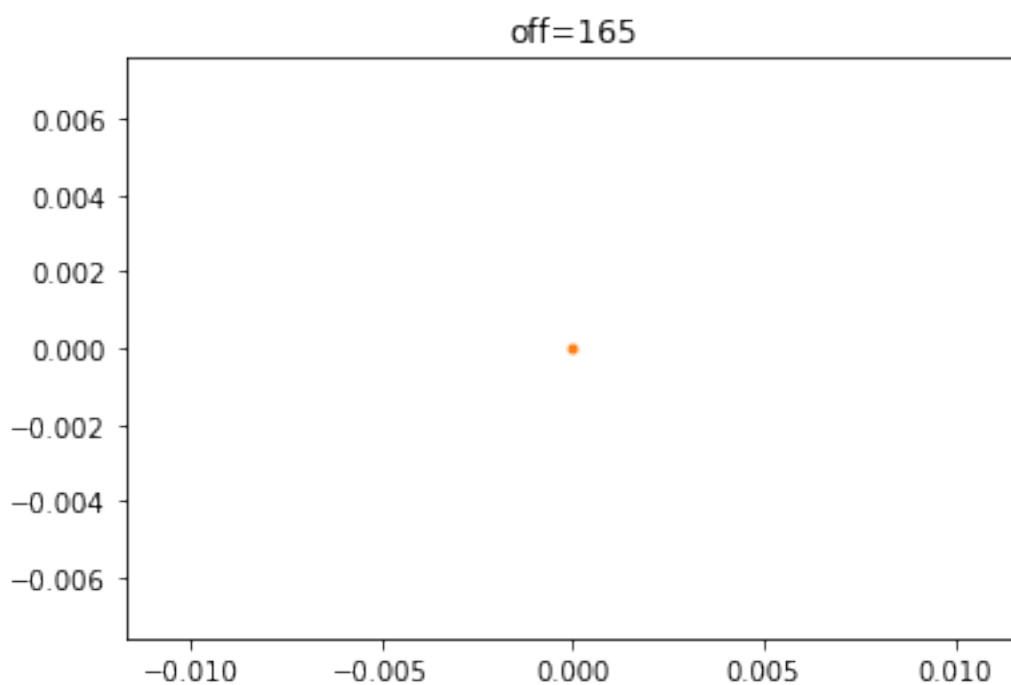n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1



/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

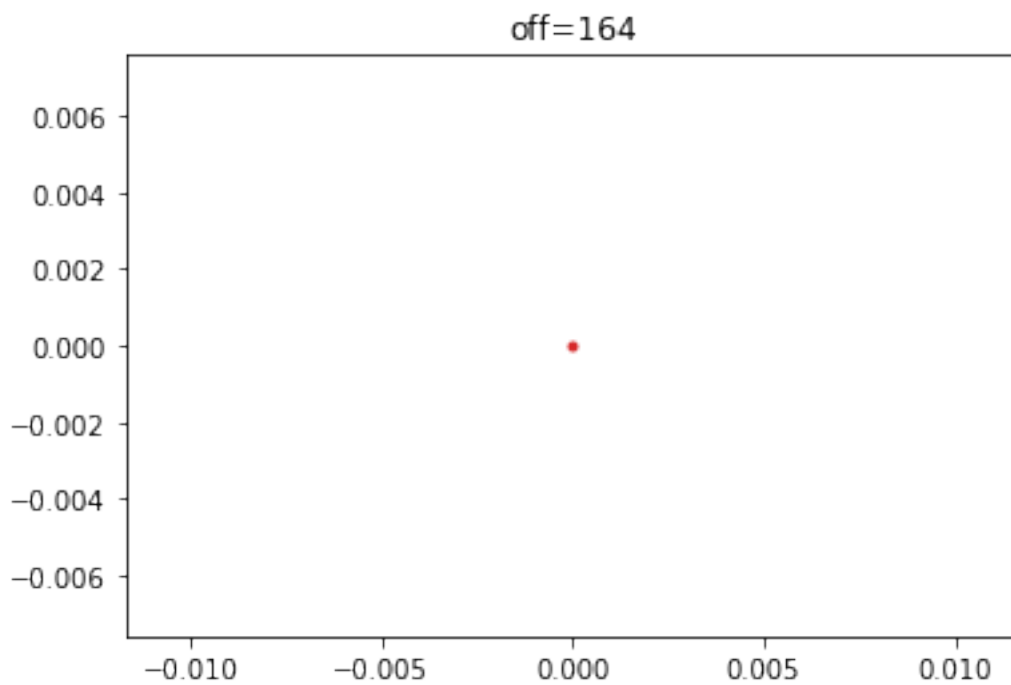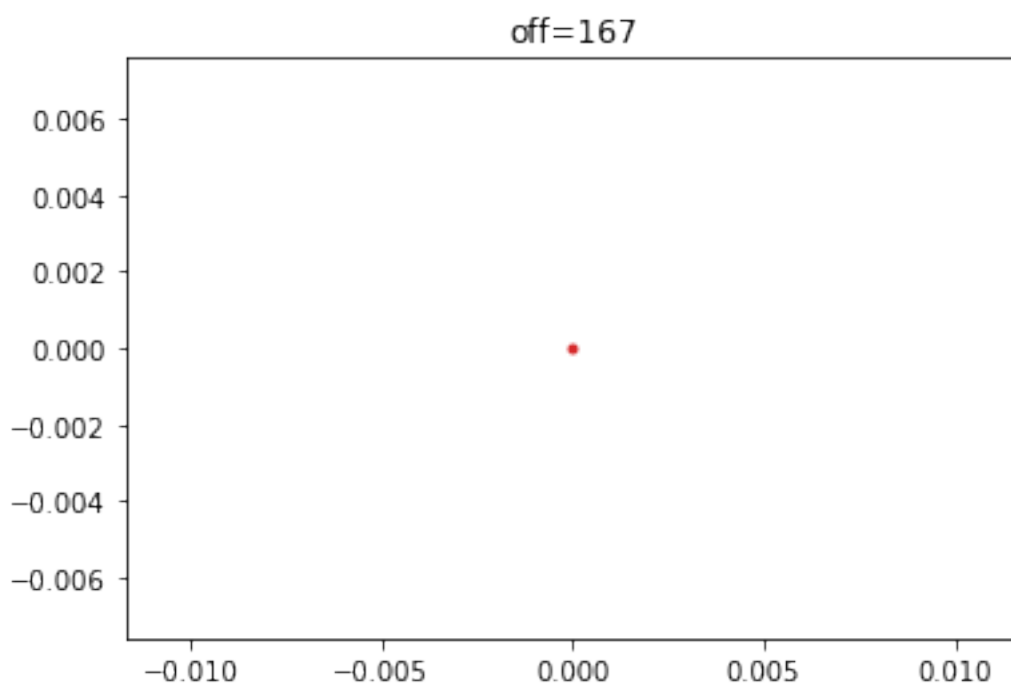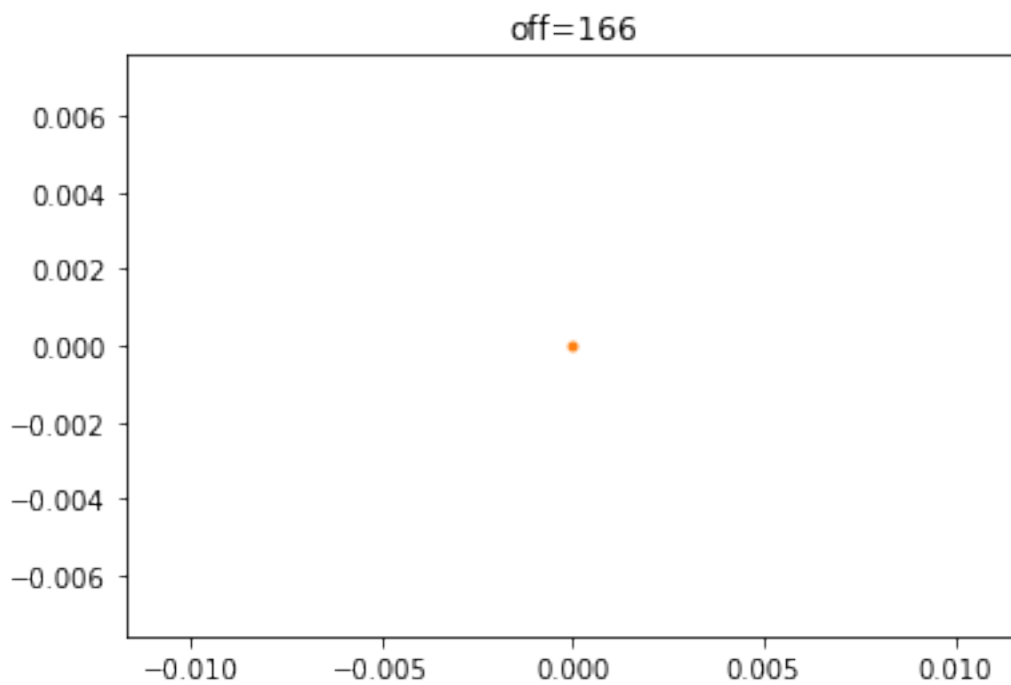n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

off=106

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

/home/ethan/.local/lib/python3.7/site-packages/scipy/sparse/linalg/eigen/arpack/arpack.py:1592

k >= N for N * N square matrix. Attempting to use scipy.linalg.eigh instead.

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

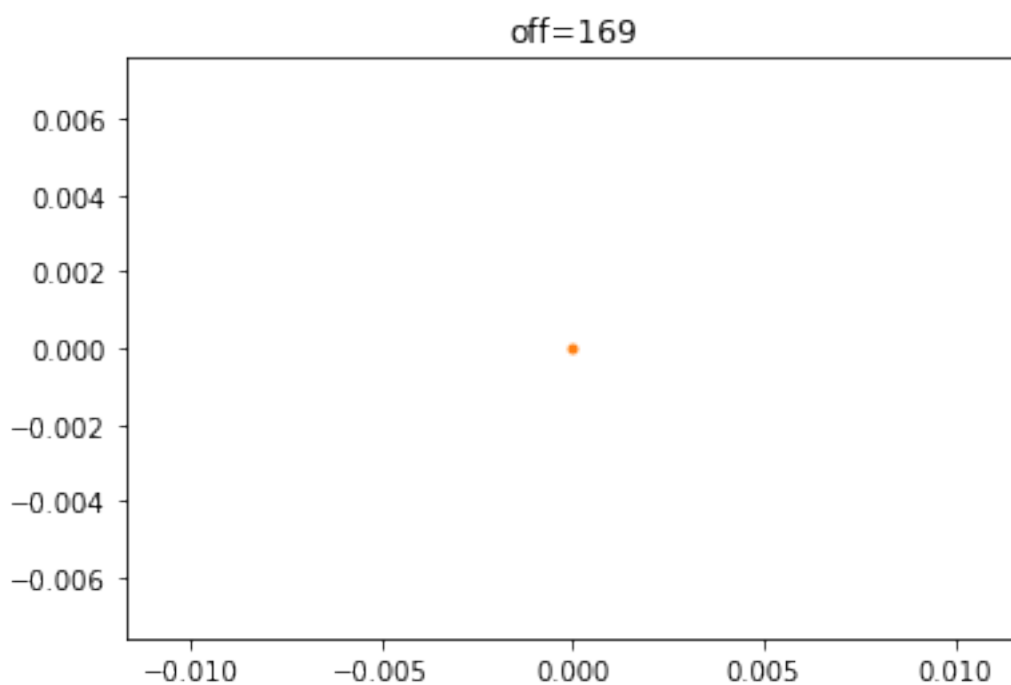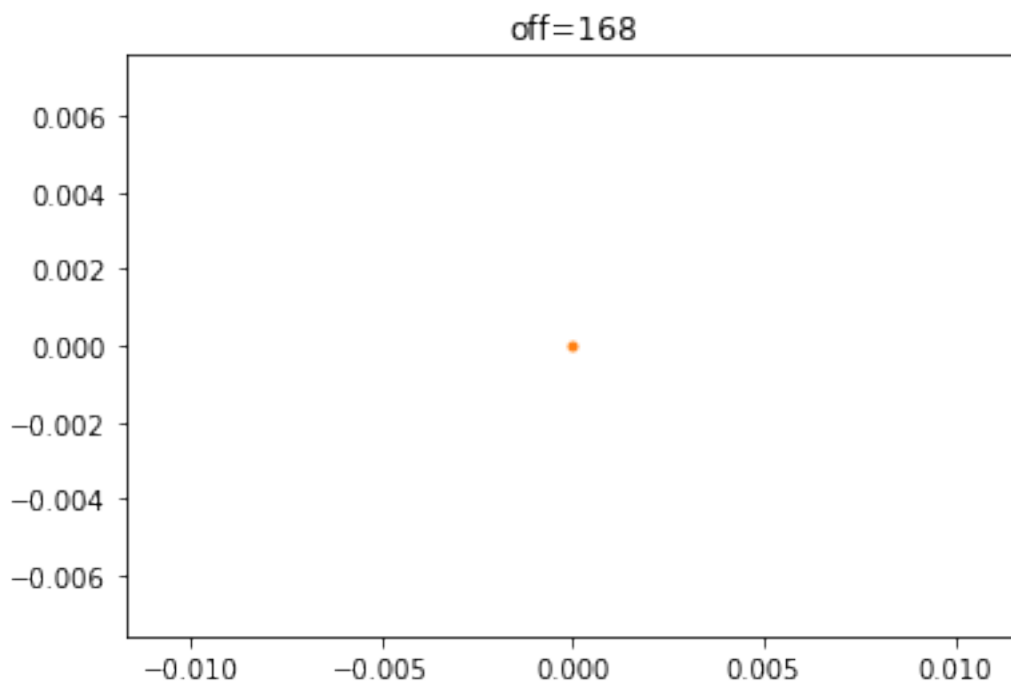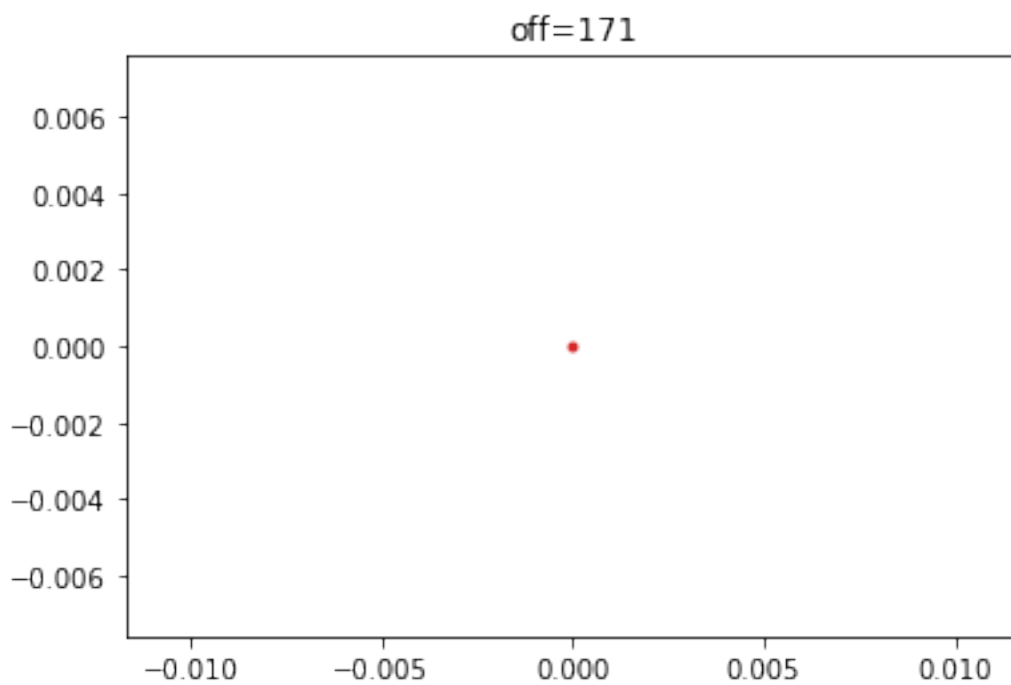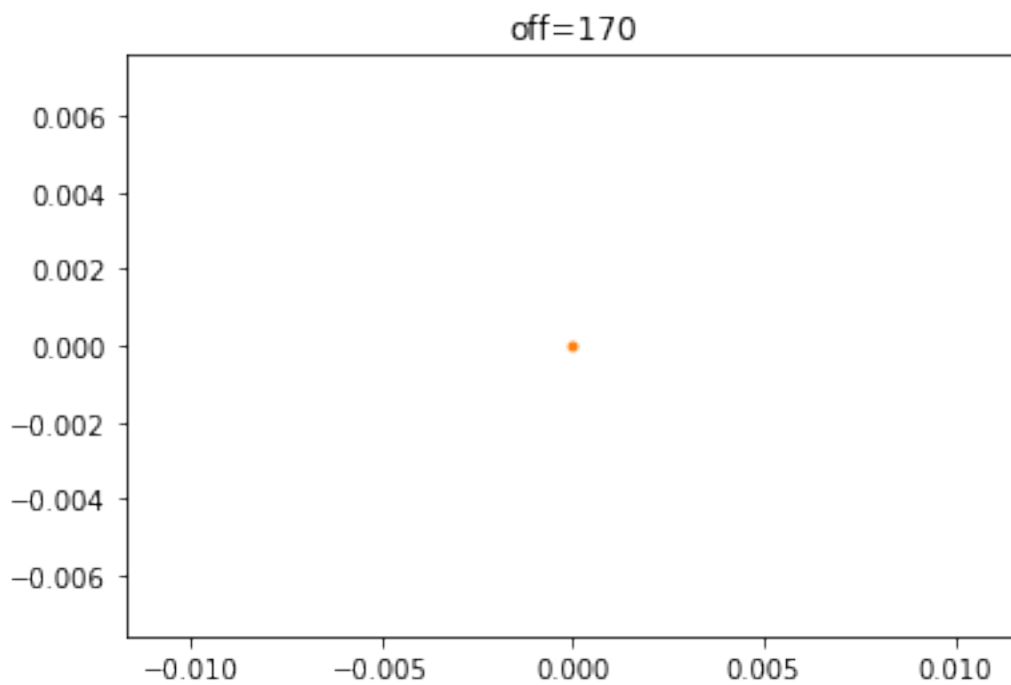n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

/home/ethan/.local/lib/python3.7/site-packages/scipy/sparse/linalg/eigen/arpack/arpack.py:1592

k >= N for N * N square matrix. Attempting to use scipy.linalg.eigh instead.

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

off=109

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1



off=111

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1



off=112

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

off=113

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1



off=115

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

off=118

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

off=119

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

/home/ethan/.local/lib/python3.7/site-packages/scipy/sparse/linalg/eigen/arpack/arpack.py:1592

k >= N for N * N square matrix. Attempting to use scipy.linalg.eigh instead.

off=121

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1



off=122

```
/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1
```



off=123

off=124

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

/home/ethan/.local/lib/python3.7/site-packages/scipy/sparse/linalg/eigen/arpack/arpack.py:1592

k >= N for N * N square matrix. Attempting to use scipy.linalg.eigh instead.

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

off=127

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

/home/ethan/.local/lib/python3.7/site-packages/scipy/sparse/linalg/eigen/arpack/arpack.py:1592

k >= N for N * N square matrix. Attempting to use scipy.linalg.eigh instead.

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

/home/ethan/.local/lib/python3.7/site-packages/scipy/sparse/linalg/eigen/arpack/arpack.py:1592

k >= N for N * N square matrix. Attempting to use scipy.linalg.eigh instead.

off=132

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

/home/ethan/.local/lib/python3.7/site-packages/scipy/sparse/linalg/eigen/arpack/arpack.py:1592

k >= N for N * N square matrix. Attempting to use scipy.linalg.eigh instead.

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

## off=137



## off=138



/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1



off=139

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

off=140



off=141

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1



off=142

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

/home/ethan/.local/lib/python3.7/site-packages/scipy/sparse/linalg/eigen/arpack/arpack.py:1592

k >= N for N * N square matrix. Attempting to use scipy.linalg.eigh instead.

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

/home/ethan/.local/lib/python3.7/site-packages/scipy/sparse/linalg/eigen/arpack/arpack.py:1592

k >= N for N * N square matrix. Attempting to use scipy.linalg.eigh instead.

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

off=145



off=146

off=147



off=148

off=149

off=150

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

/home/ethan/.local/lib/python3.7/site-packages/scipy/sparse/linalg/eigen/arpack/arpack.py:1592

k >= N for N * N square matrix. Attempting to use scipy.linalg.eigh instead.

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

off=153

off=154

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1



off=156

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1



off=157

/home/ethan/.local/lib/python3.7/site-packages/umap/umap_.py:1383: UserWarning:

n_neighbors is larger than the dataset size; truncating to X.shape[0] - 1

off=158



off=159

off=160



off=161

off=162



off=163

off=164



off=165

off=166



off=167

off=168



off=169

off=170



off=171

off=172



off=173

off=174

In [73]: **for** off **in** set(samp.OFFENSE):

        mask = samp_X.OFFENSE == off

        data = samp_X[mask]

        **try**:
            trans = TSNE(perplexity=12).fit_transform(data)

            plt.scatter(trans[:, 0], trans[:, 1], c=[sns.color_palette()[x] **for** x **in** samp_
            plt.gca().set_aspect('equal', 'datalim')
            plt.title(f'off={off}')
            plt.show()
        **except**:

            **continue**

108

off=1



off=2

off=3



off=4

off=5



off=6

## off=7



## off=8

off=9

off=10

off=11

off=12

off=13



off=14

off=15



off=16

off=17

off=18

off=19



off=20

off=21


off=22

off=23

off=24

off=25


off=26

off=27



off=28

off=30

off=31

off=32



off=33

off=34

off=35

off=36

off=37

off=38



off=39

off=40



off=41

off=42



off=43

off=44



off=45

off=46



off=47

off=48

off=49

off=50



off=51

off=52



off=53

off=54



off=55

off=56



off=57

off=58



off=59

137

off=60



off=61

off=62



off=63

off=64



off=65

off=66



off=67

off=68



off=69

off=70



off=71

off=72

off=73

off=74



off=75

off=76



off=77

off=78



off=79

off=80

off=81

off=82



off=83

off=84



off=85

off=86



off=87

off=88



off=89

off=91



off=92

off=93



off=94

off=95



off=96

off=97



off=98

off=99



off=100

off=101



off=102

off=103



off=104

off=105



off=106

off=107



off=108

off=109



off=110

off=111



off=112

off=113



off=114

off=116



off=117

off=119



off=120

off=122



off=123

off=124

off=125

off=126



off=127

off=128

off=129

off=130



off=131

off=133



off=134

off=135



off=136

off=137



off=139

off=140



off=142

## off=143



## off=144

off=145



off=151
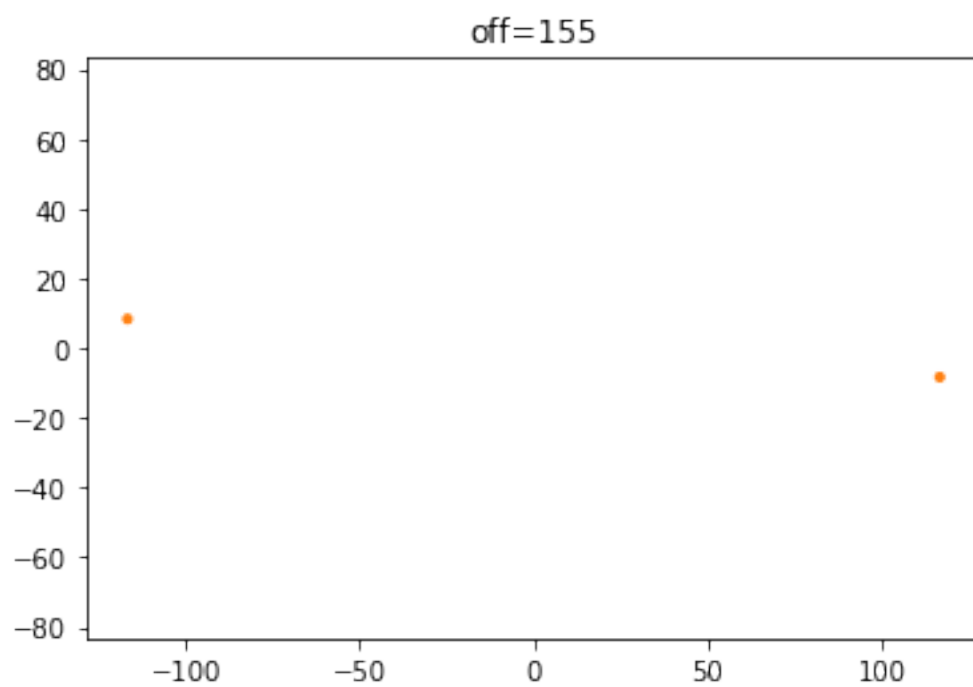
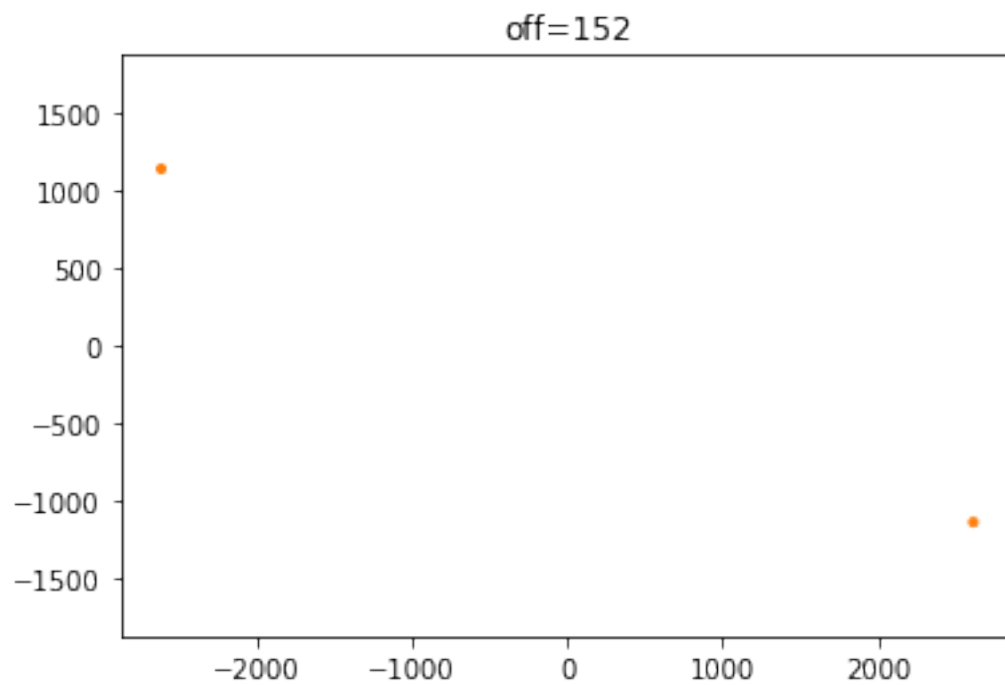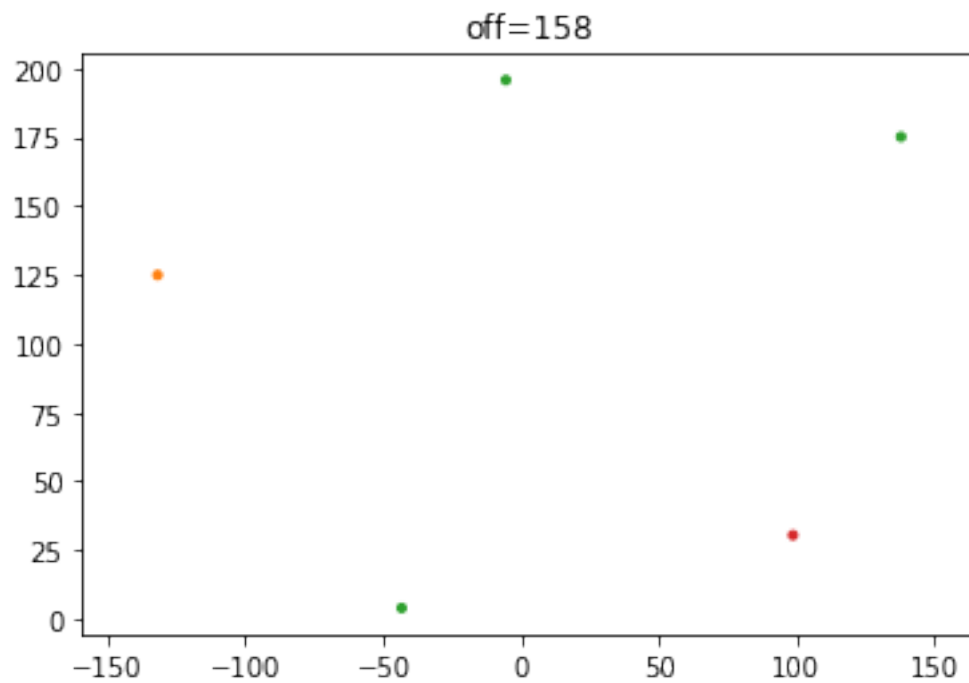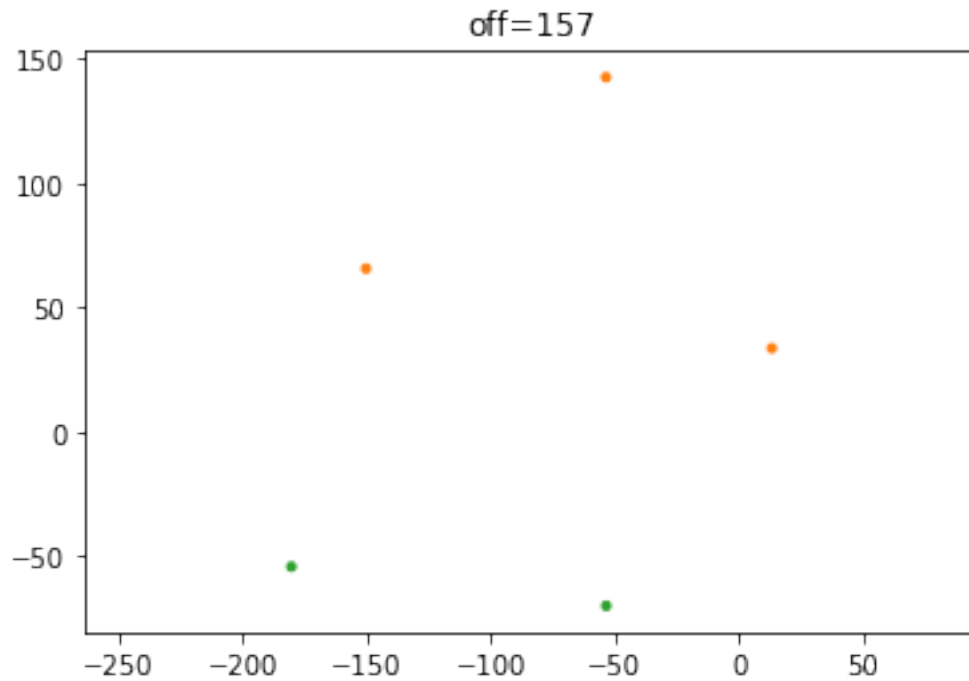off=152



off=155

off=157



off=158

TSNE is having mized results while UMAP doesn't seem to capture any good groupings. Maybe I should go straight to kmeans or kdtrees\ I should block by crime and examine a few "big" crimes the try several different perplexities.

In [ ]: