

In [12]: `#!/usr/bin/python`

```

"""
    Skeleton code for k-means clustering mini-project.
"""

import pickle
import numpy
import matplotlib.pyplot as plt
import sys
sys.path.append("../tools/")
from feature_format import featureFormat, targetFeatureSplit

def Draw(pred, features, poi, mark_poi=False, name="image.png", f1_name="feature_1", f2_name="feature_2"):
    """ some plotting code designed to help you visualize your clusters """

    ### plot each cluster with a different color--add more colors for
    ### drawing more than five clusters
    colors = ["b", "c", "k", "m", "g"]
    for ii, pp in enumerate(pred):
        plt.scatter(features[ii][0], features[ii][1], color = colors[pred[ii]])

    ### if you like, place red stars over points that are POIs (just for funsies)
    if mark_poi:
        for ii, pp in enumerate(pred):
            if poi[ii]:
                plt.scatter(features[ii][0], features[ii][1], color="r", marker="*")

    plt.xlabel(f1_name)
    plt.ylabel(f2_name)
    plt.savefig(name)
    plt.show()

    ### load in the dict of dicts containing all the data on each person in the dataset
    data_dict = pickle.load( open("../final_project/final_project_dataset.pkl", "r") )
    ### there's an outlier--remove it!
    data_dict.pop("TOTAL", 0)

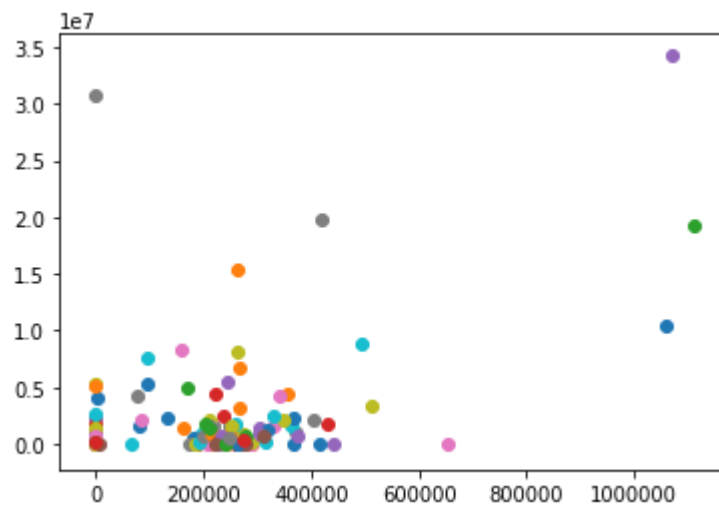
    ### the input features we want to use
    ### can be any key in the person-level dictionary (salary, director_fees, etc.)
    feature_1 = "salary"
    feature_2 = "exercised_stock_options"
    poi = "poi"
    features_list = [poi, feature_1, feature_2]
    data = featureFormat(data_dict, features_list )
    poi, finance_features = targetFeatureSplit( data )

```

The starter code can be found in `k_means/k_means_cluster.py`, which reads in the email + financial (E+F) dataset and gets us ready for clustering. You'll start with performing k-means based on just two financial features--take a look at the code, and determine which features the code uses for clustering.

Run the code, which will create a scatterplot of the data. Think a little bit about what clusters you would expect to arise if 2 clusters are created.

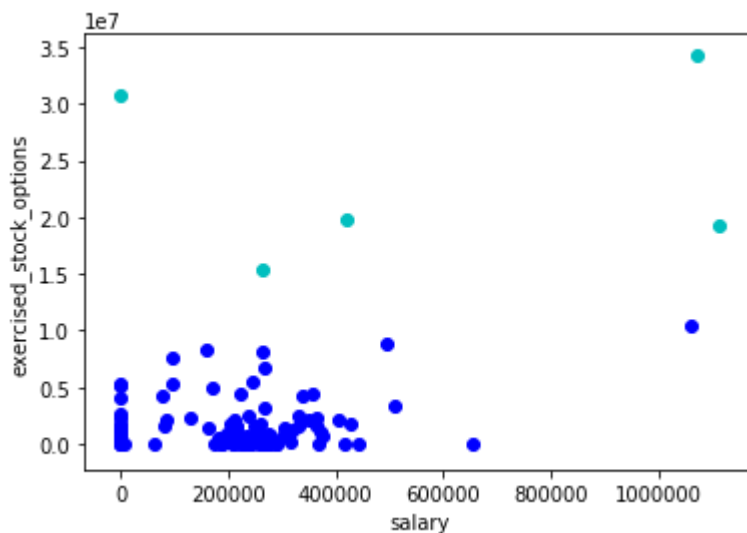
```
In [13]: ### in the "clustering with 3 features" part of the mini-project,  
### you'll want to change this line to  
### for f1, f2, _ in finance_features:  
### (as it's currently written, the line below assumes 2 features)  
for f1, f2 in finance_features:  
    plt.scatter( f1, f2 )  
plt.show()
```



Deploy k-means clustering on the `financial_features` data, with 2 clusters specified as a parameter. Store your cluster predictions to a list called `pred`, so that the `Draw()` command at the bottom of the script works properly. In the scatterplot that pops up, are the clusters what you expected?

```
In [14]: ### cluster here; create predictions of the cluster labels
### for the data and store them to a list called pred
from sklearn.cluster import KMeans as km
clf = km(n_clusters = 2)
clf.fit(finance_features, poi)
pred = clf.predict(finance_features)

### rename the "name" parameter when you change the number of features
### so that the figure gets saved to a different file
try:
    Draw(pred, finance_features, poi, mark_poi=False, name="clusters.pdf", f1_name="clusters.pdf")
except NameError:
    print "no predictions object named pred found, no clusters to plot"
```



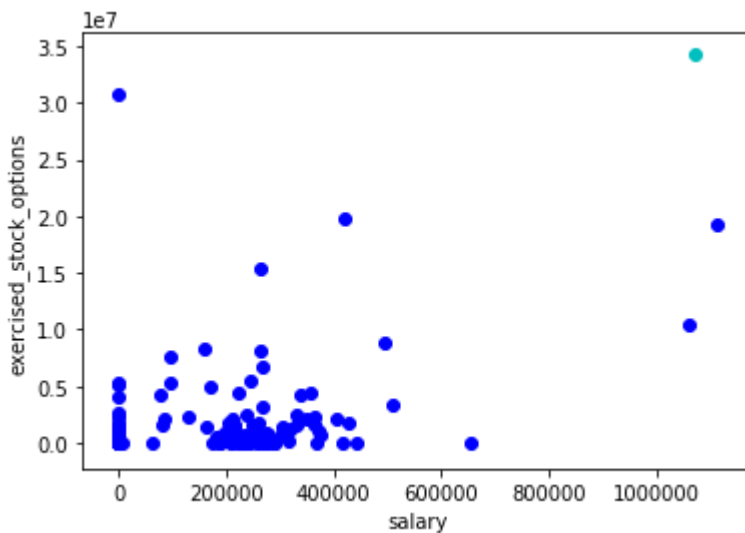
Add a third feature to features\_list, "total\_payments". Now rerun clustering, using 3 input features instead of 2 (obviously we can still only visualize the original 2 dimensions). Compare the plot with the clusterings to the one you obtained with 2 input features. Do any points switch clusters? How many? This new clustering, using 3 features, couldn't have been guessed by eye--it was the k-means algorithm that identified it.

(You'll need to change the code that makes the scatterplot to accommodate 3 features instead of 2, see the comments in the starter code for instructions on how to do this.)

```
In [16]: feature_3 = "total_payments"
feature_1 = "salary"
feature_2 = "exercised_stock_options"
poi = "poi"
features_list = [poi, feature_1, feature_2, feature_3]
data = featureFormat(data_dict, features_list)
poi, finance_features = targetFeatureSplit(data)

from sklearn.cluster import KMeans as km
clf = km(n_clusters = 2)
clf.fit(finance_features, poi)
pred = clf.predict(finance_features)

try:
    Draw(pred, finance_features, poi, mark_poi=False, name="clusters.pdf", f1_name="poi")
except NameError:
    print "no predictions object named pred found, no clusters to plot"
```



In the next lesson, we'll talk about feature scaling. It's a type of feature preprocessing that you should perform before some classification and regression tasks. Here's a sneak preview that should call your attention to the general outline of what feature scaling does.

What are the maximum and minimum values taken by the "exercised\_stock\_options" feature used in this example?

(NB: if you look at finance\_features, there are some "NaN" values that have been cleaned away and replaced with zeroes--so while those might look like the minima, it's a bit deceptive because they're more like points for which we don't have information, and just have to put in a number. So for this question, go back to data\_dict and look for the maximum and minimum numbers that show up there, ignoring all the "NaN" entries.)

```
In [20]: stock = featureFormat(data_dict, ["exercised_stock_options"])
print max(stock), min(stock)
```

```
[34348384.] [3285.]
```

What are the maximum and minimum values taken by “salary”?

(NB: same caveat as in the last quiz. If you look at finance\_features, there are some "NaN" values that have been cleaned away and replaced with zeroes--so while those might look like the minima, it's a bit deceptive because they're more like points for which we don't have information, and just have to put in a number. So for this question, go back to data\_dict and look for the maximum and minimum numbers that show up there, ignoring all the "NaN" entries.)

```
In [22]: sal = featureFormat(data_dict, ["salary"])  
print max(sal),min(sal)  
  
[1111258.] [477.]
```

```
In [ ]:
```