

CS498AML_HW4_panz2

Group member: Pan Zhang, Xinyu Tian
9/29/2018

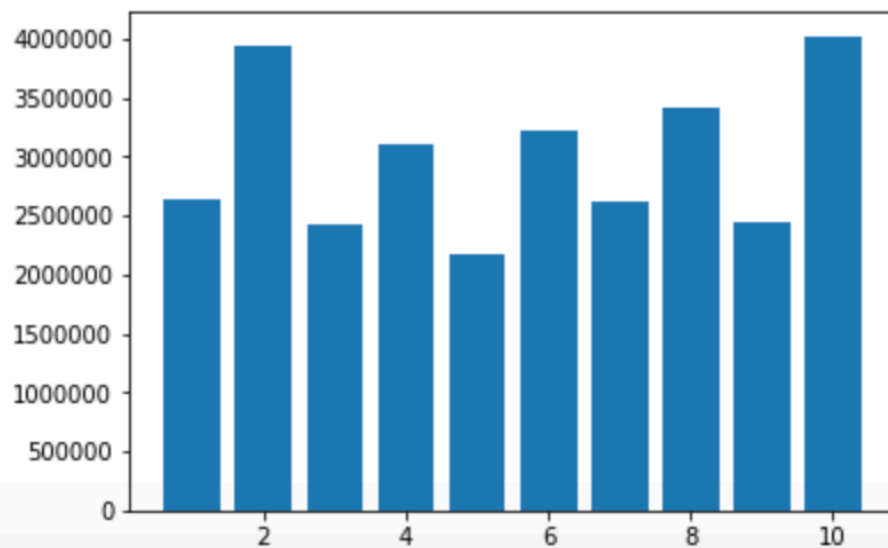
A single plot of error (use mean squared error) vs category (1-10):
Error is calculated based on training data (data_batch_1~5)

```
In [17]: li_MSE
```

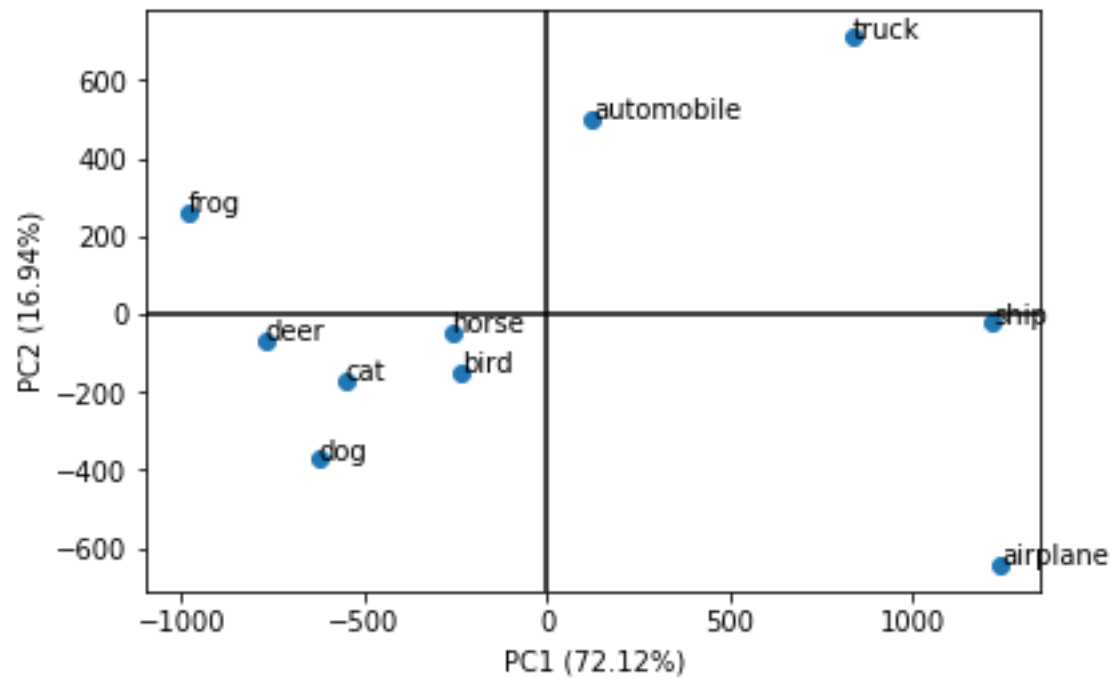
```
Out[17]: [2629467.1550547639,  
          3946267.9085392822,  
          2432107.4911725386,  
          3106124.7729624542,  
          2181905.7666432783,  
          3230245.0703085512,  
          2623410.9962643329,  
          3424435.2463813997,  
          2438641.3675135057,  
          4024349.8710604957]
```

```
In [18]: plt.bar(range(1,11,1), li_MSE)
```

```
Out[18]: <Container object of 10 artists>
```



A single 2D plot with the results of principal coordinate analysis applied to pairs of means:



```

In [1]: 1 import pickle
        2 import os
        3 import numpy as np
        4 import pandas as pd
        5 from sklearn.decomposition import PCA
        6 from scipy.spatial.distance import euclidean
        7 from skbio.stats.ordination import pcoa
        8 import matplotlib.pyplot as plt

In [2]: 1 def unpickle(file):
        2     with open(file, 'rb') as fo:
        3         dict = pickle.load(fo, encoding='bytes')
        4     return dict

In [3]: 1 batchName = [i for i in os.listdir() if 'data_batch' in i]
        2 for i in batchName:
        3     globals()[i] = unpickle(i)

In [4]: 1 classLevels = set()
        2 df_class = {}
        3 for eachBatch in batchName:
        4     # get all class Levels
        5     classLevels = classLevels.union(set(globals()[eachBatch][b'labels']))
        6     # initialize - create null dataframes
        7     for classValue in classLevels:
        8         df_class[classValue] = pd.DataFrame(columns = list(range(0, 3072, 1)))

In [5]: 1 def sepPickleClass(batch, classValueRange):
        2     for classValue in classValueRange:
        3         ind = [i for i, val in enumerate(batch[b'labels']) if val == classValue]
        4         df_class[classValue] = df_class[classValue].append(pd.DataFrame(batch[b'data'][ind]))

In [6]: 1 for eachBatch in batchName:
        2     sepPickleClass(globals()[eachBatch], classLevels)

In [7]: 1 def getMeanImage(classValue):
        2     return df_class[classValue].mean(axis = 0)

In [8]: 1 mean_class = {}
        2 for classValue in classLevels:
        3     mean_class[classValue] = getMeanImage(classValue)

In [9]: 1 def getPcaRepr(classValue, n_components):
        2     pca_model = PCA(n_components = n_components, svd_solver = 'full')
        3     pca_model.fit(df_class[classValue])
        4     pcaRepr = pca_model.inverse_transform(pca_model.transform(df_class[classValue]))
        5     return pcaRepr

In [10]: 1 pca_repr_class = {}
        2 for classValue in classLevels:
        3     pca_repr_class[classValue] = getPcaRepr(classValue, 20)

In [11]: 1 def getMSE(y_true, y_pred):
        2     assert y_true.shape[0] == y_pred.shape[0]
        3     return np.sum(np.sum(np.square(np.subtract(y_pred, y_true)))) / y_true.shape[0]

In [12]: 1 li_MSE = []
        2 for classValue in classLevels:
        3     li_MSE.append(getMSE(df_class[classValue], pca_repr_class[classValue]))

In [18]: 1 plt.bar(list(classLevels), li_MSE)
        2 plt.xticks(list(classLevels))

In [19]: 1 dTab = []
        2 for classValue1 in classLevels:
        3     dRow = []
        4     for classValue2 in classLevels:
        5         d = euclidean(mean_class[classValue1], mean_class[classValue2])
        6         dRow.append(d)
        7     dTab.append(dRow)
        8 dMat = np.matrix(dTab)

In [20]: 1 def create2DMap(distanceMat, axes, axis_labels = None):
        2     ordResult = pcoa(distanceMat)
        3     coord_matrix = ordResult.samples.values.T
        4     prop_expl = ordResult.proportion_explained
        5     fig = plt.figure()
        6     xs = coord_matrix[axes[0]]
        7     ys = coord_matrix[axes[1]]
        8     plt.scatter(xs, ys)
        9     names = ['airplane', 'automobile', 'bird', 'cat', 'deer', 'dog', 'frog', 'horse', 'ship', 'truck']
        10    for i in range(len(xs)):
        11        plt.annotate(names[i], (xs[i], ys[i]))
        12    if axis_labels is None:
        13        axis_labels = ['PC{0}'.format((axis+1), prop_expl[axis]) for axis in axes]
        14    plt.xlabel(axis_labels[0])
        15    plt.ylabel(axis_labels[1])
        16    plt.axvline(x=0, color = 'k')
        17    plt.axhline(y=0, color = 'k')
        18    return fig

In [21]: 1 pcoa_2dMap = create2DMap(dMat, [0,1])

```