

homework2

April 27, 2016

1 Homework 2

Titipat Achakulvisut

1.1 Problem 1: Classifier – back propagation

Use the code below to create three data clusters which are vertically stacked and therefore not linearly separable, as we discussed in class. Your job is to create code implementing back propagation for a two layer neural network which can perform this classification. Use a network with 4 hidden units, as indicated in the shell code. Don't worry about cross validation and all that – feel free to just use all the data.

```
In [1]: %matplotlib inline
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

In [2]: sd = .85
X1 = np.vstack((np.random.normal(0, sd, size=(50,1)),
                np.random.normal(0, sd, size=(50,1)),
                np.random.normal(0, sd, size=(50,1))))
X2 = np.vstack((np.random.normal(0, sd, size=(50,1)),
                np.random.normal(5, sd, size=(50,1)),
                np.random.normal(10, sd, size=(50,1))))
# X3 = np.ones((150, 1))
X = np.concatenate((X1, X2), axis=1)
y = np.zeros((150, 3))
y[0:50, 0] = 1
y[50:100, 1] = 1
y[100:150, 2] = 1

In [3]: plt.scatter(X[:,0], X[:,1], c=y);
```



```
In [4]: def sigmoid(z):
        return 1/(1 + np.exp(-z))

        def sigmoid_grad(z):
            return sigmoid(z)*(1 - sigmoid(z))

In [5]: ninput = 2
        nhidden = 4
        noutput = 3
        W = np.random.uniform(-1, 1, size=(ninput+1, nhidden)) - 0.5
        V = np.random.uniform(-1, 1, size=(nhidden+1, noutput)) - 0.5

In [6]: def predict(W, V, X):
        """Do forward propagation given first and second layers"""
        n, m = X.shape
        a1 = np.concatenate((np.ones((n,1)), X), axis=1)
        z2 = a1.dot(W)
        a2 = sigmoid(z2)
        a2 = np.concatenate((np.ones((n,1)), a2), axis=1)
        z3 = a2.dot(V)
        a3 = sigmoid(z3)
        h = a3
        return h

In [7]: def compute_grad_bp(W, V, X, Y):
        """
```

```
# initialize few parameters
n, m = X.shape

dW = np.zeros_like(W)
dV = np.zeros_like(V)

a1 = np.concatenate((np.ones((n,1)), X), axis=1)
z2 = a1.dot(W)
a2 = sigmoid(z2)
a2 = np.concatenate((np.ones((n,1)), a2), axis=1)
z3 = a2.dot(V)
a3 = sigmoid(z3)
h = a3

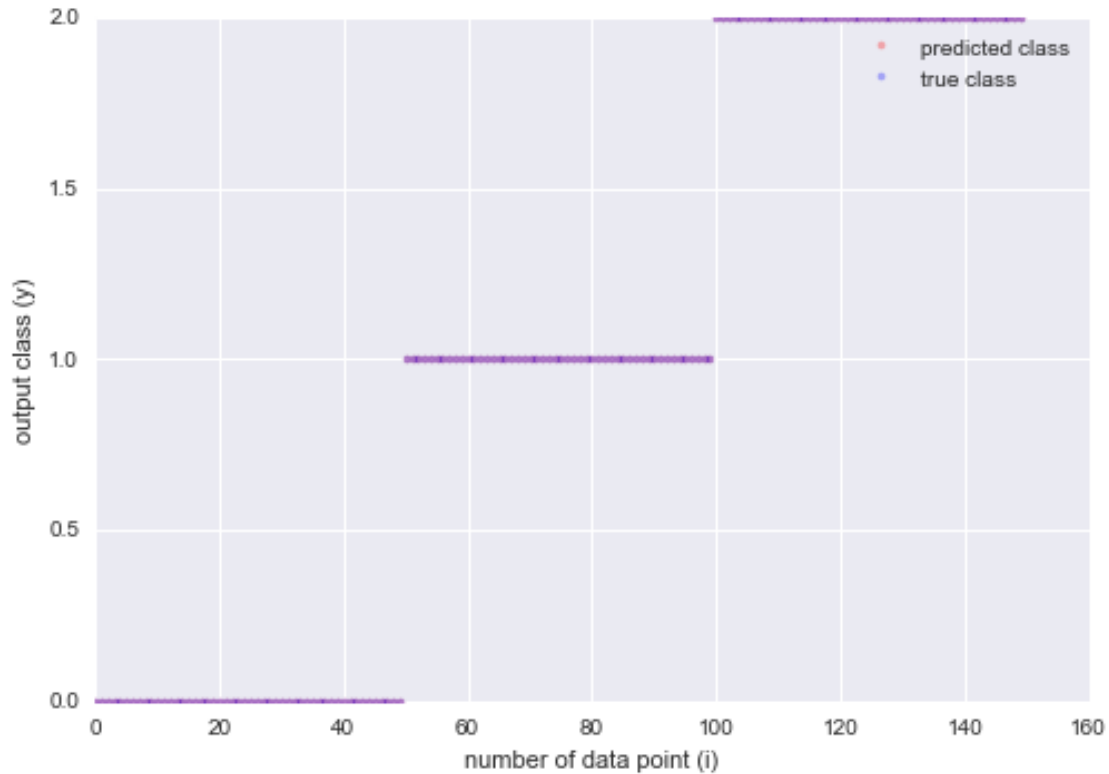
# back propagation
delta3 = h - Y
delta2 = delta3.dot(V[1:,:].T)*sigmoid_grad(z2)

dW = (1/m)*((delta2.T).dot(a1)).T
dV = (1/m)*((delta3.T).dot(a2)).T

return dW, dV
```

```
In [9]: y_hat = predict(W, V, X)
print('Predicted output class: \n', y_hat.argmax(axis=1))
print( )
```

```
In [10]: plt.plot(y_hat.argmax(axis=1), '.r', alpha=0.3)
plt.plot(y.argmax(axis=1), '.b', alpha=0.3) # output class for each datapoint
plt.xlabel('number of data point (i)')
plt.ylabel('output class (y)')
plt.legend(['predicted class', 'true class'])
plt.show()
```



1.2 Problem 2: Unsupervised learning – k-means clustering

In each of the next problems you are given data that appears to be produced by two clusters (use the code in `ps1 datasets.m`, 2a). Your job in this problem is to use k-means clustering to classify the data according to two clusters. Plot the trajectory of the means as they are updated over iterations, with this trajectory superimposed over the data sets.

```
In [11]: X1 = np.hstack((np.random.normal(6, 1, size=(100,1)),
                        np.random.normal(2, 1, size=(100, 1))))
X2 = np.hstack((np.random.normal(2, 1, size=(100,1)),
                np.random.normal(8, 1, size=(100, 1))))
X = np.vstack((X1, X2))

In [12]: def update(X, centers, alpha=0.05):
    """Update K-mean center and compute cost in that iteration"""
    # number of cluster
    K = len(centers)
    # compute distance
    D = np.vstack([np.linalg.norm(X - center, axis=1) for center in centers])
    clusters = np.argmin(D.T, axis=1) # assign clusters
    centers = np.vstack([X[clusters==k].mean(axis=0) for k in range(K)])

    # compute cost of K-mean
    J = np.sum([np.sum(np.linalg.norm(X[clusters==k] - centers[k], axis=1))
                for k in range(K)])
    return centers, J, clusters
```

```

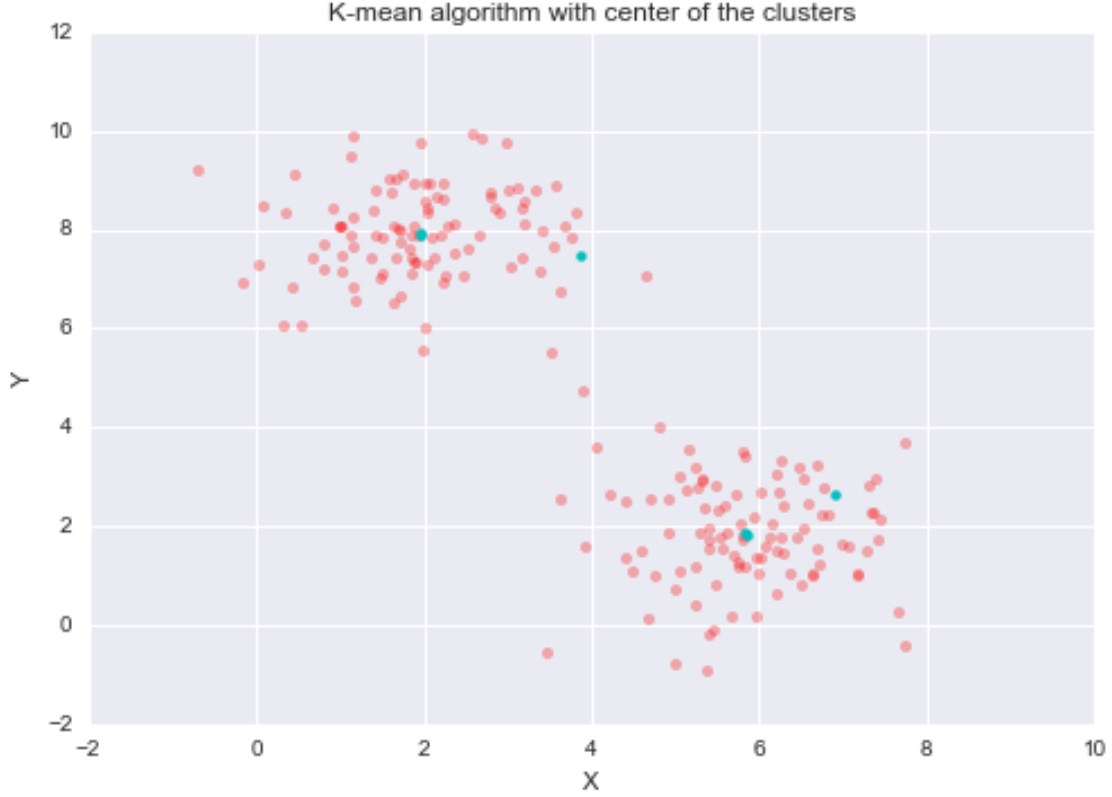
In [13]: def random_centers(X, K=2):
        """
        Randomly generate K centers from data X
        """
        X_min = X.min(axis=0)
        X_max = X.max(axis=0)
        centers = []
        for k in range(K):
            center = [np.random.uniform(X_min[i], X_max[i]) for i in range(len(X_min))]
            centers.append(center)
        return np.array(centers)

In [14]: n_iter = 100
        n_cluster = 2
        centers = random_centers(X, K=n_cluster) # random center
        print('Initial centers are \n', centers)
        clusters = np.zeros(X.shape[0])
        u = True
        J_list = []
        centers_list = []
        while u:
            clusters_prev = clusters
            centers_list.append(centers)
            centers, J, clusters = update(X, centers) # compute new center and cost
            J_list.append(J)
            if np.all(clusters_prev == clusters):
                u = False

Initial centers are
[[ 6.90985114  2.63489775]
 [ 3.88114647  7.45748333]]

In [15]: plt.scatter(X[:, 0], X[:, 1], color='r', alpha=0.3)
        for c in centers_list:
            for k in range(n_cluster):
                plt.scatter(c[k][0], c[k][1], color='c', lw = 0)
        plt.xlabel('X')
        plt.ylabel('Y')
        plt.title('K-mean algorithm with center of the clusters')
        plt.show()
        print('Center for final cluster are \n',
              list(centers_list[-1][0]), ' and \n',
              list(centers_list[-1][1]))

```



Center for final cluster are

[5.8350106130994046, 1.845847253250521] and
[1.9500944871562012, 7.9396316730666197]

1.3 Problem 3: Unsupervised learning – ML gradient descent

Your job is to classify the same dataset as used in problem 2 into two clusters using ML gradient descent to update the values of mean and standard deviation of two Gaussian clusters. Assume that the prior probabilities of each cluster are already given and fixed at 0.5 each.

Ans. Basically, we have mixture of Gaussians on variable \mathbf{x} as follows:

$$p(\mathbf{x}) = \sum_{k=1}^K \rho_k \frac{1}{\sqrt{\det(2\pi\Sigma_k)}} \exp\left[-\frac{1}{2}(\mathbf{x} - \mu_k)^T \Sigma_k^{-1}(\mathbf{x} - \mu_k)\right]$$

where ρ_k sum up to 1. Now the log-likelihood will be, as following

$$\ln p(\mathbf{x}) = \sum_{k=1}^K \log(\rho_k) - \log(C) - \log(\det(2\pi\Sigma_k)) - \frac{1}{2}((\mathbf{x} - \mu_k)^T \Sigma_k^{-1}(\mathbf{x} - \mu_k))$$

where C in this case is a constant term (in front of Gaussian distribution)

Thus the gradient of log-likelihood with respect to mean and variance (giving we assuming that we know ρ_k) is as follows

$$\frac{\partial \ln p(\mathbf{x})}{\partial \mu_k} = \Sigma_k^{-1}(\mathbf{x} - \mu_k)$$

and

$$\frac{\partial \ln p(\mathbf{x})}{\partial \Sigma_k} = \frac{1}{2}[-\Sigma_k^{-1} + \Sigma_k^{-1}(\mathbf{x} - \mu_k)(\mathbf{x} - \mu_k)^T \Sigma_k^{-1}]$$

note (some matrix tricks):

- use this trick in finding gradient with respect to mean $\frac{\partial}{\partial \mathbf{x}}(\mathbf{x} - \mathbf{s})^T \mathbf{W}(\mathbf{x} - \mathbf{s}) = 2\mathbf{W}(\mathbf{x} - \mathbf{s})$
- use this trick in finding gradient with respect to covariance matrix (in this case, we know that covariance matrix is symmetric so we can get rid of transpose) $\frac{\partial \ln |\det(\mathbf{X})|}{\partial (\mathbf{X})} = (\mathbf{X}^{-1})^T$

In this homework, we consider special case where Spherical covariance is as follows $\Sigma = \sigma_i^2 I$
The likelihood and log-likelihood can be written as follows:

$$p(\mathbf{x}) = \sum_{k=1}^K \rho_k \frac{1}{\sqrt{2\pi}\sigma_k} \exp\left[-\frac{1}{2} \frac{(\mathbf{x} - \mu_k)^T (\mathbf{x} - \mu_k)}{\sigma_k^2}\right]$$

$$\ln p(\mathbf{x}) = \sum_{k=1}^K \log(\rho_k) - \log(C) - d \log \sigma_k - \frac{1}{2} \frac{(\mathbf{x} - \mu_k)^T (\mathbf{x} - \mu_k)}{\sigma_k^2}$$

where d is dimension of the Gaussian (in our case, $d = 2$)

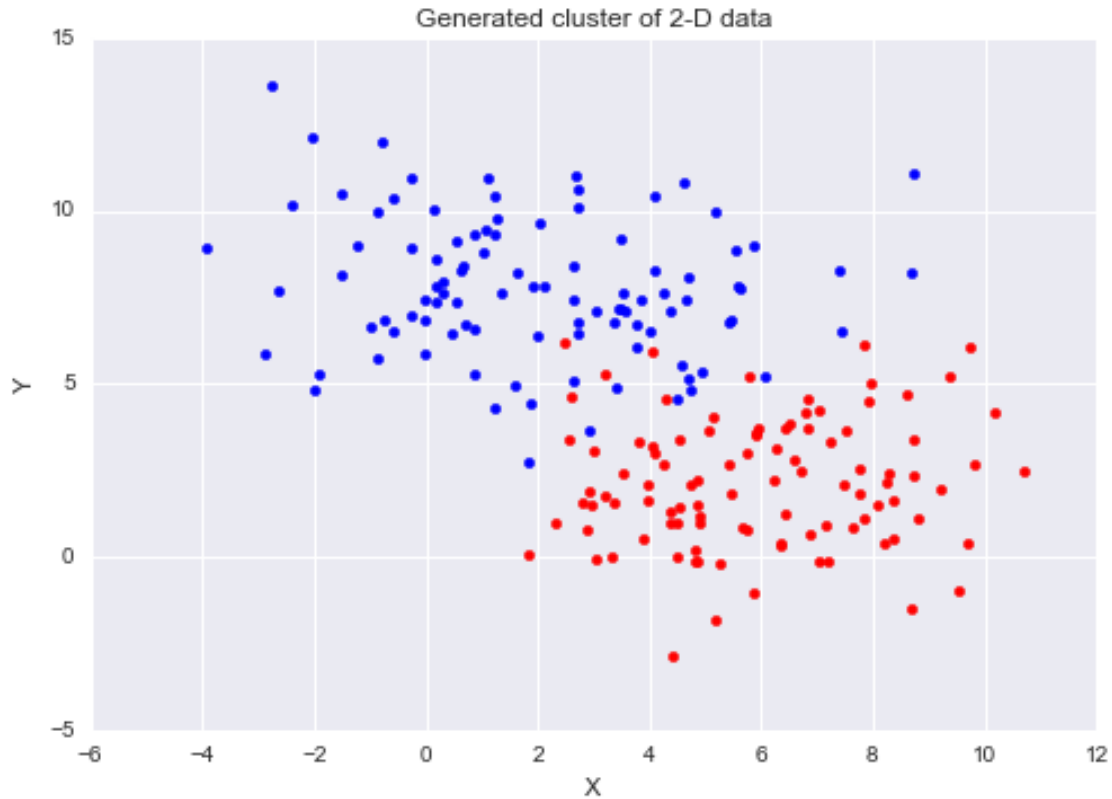
Therefore, the partial derivative with respect to mean and standard deviation is as follows:

$$\frac{\partial \ln p(\mathbf{x})}{\partial \mu_k} = \frac{(\mathbf{x} - \mu_k)}{\sigma_k^2}$$

$$\frac{\partial \ln p(\mathbf{x})}{\partial \sigma_k} = -\frac{d}{\sigma_k} + \frac{(\mathbf{x} - \mu_k)^T (\mathbf{x} - \mu_k)}{\sigma_k^3}$$

```
In [16]: import numpy.linalg as la
def inv(A):
    """return inverse of input matrix A"""
    return la.inv(A)

In [17]: # generate data for problem
X1 = np.hstack((np.random.normal(6, 2, size=(100,1)),
                np.random.normal(2, 2, size=(100, 1))))
X2 = np.hstack((np.random.normal(2, 3, size=(100,1)),
                np.random.normal(8, 2, size=(100, 1))))
X = np.vstack((X1, X2))
y = np.vstack([np.ones((100, 1)), np.zeros((100, 1))]).flatten()
plt.scatter(X[(y==1).flatten(), 0], X[(y==1).flatten(), 1], color='r')
plt.scatter(X[(y==0).flatten(), 0], X[(y==0).flatten(), 1], color='b')
plt.xlabel('X')
plt.ylabel('Y')
plt.title('Generated cluster of 2-D data')
plt.show()
print('Mean for each clusters is \n',
      np.mean(X1, axis=0),
      np.mean(X2, axis=0))
```



Mean for each clusters is
 [5.92435433 2.15143611] [1.99348394 7.75918159]

This is the case when we want to Gaussian mixture of mean and arbitrary covariance matrix. Here we print out final center of clusters and covariance matrix.

```
In [18]: def data_cluster_matrix(x, centers, sigma):
         """compute class of data point x for all centers and sigma"""
         p = [float(rbf_kernel_matrix(x, center, s))
               for (center, s) in zip(centers, sigma)]
         c = np.argmax(p)
         return c
```

```
In [19]: def rbf_kernel_matrix(x, y, S):
         """
         Compute radial basis kernel of two vector
         k(x, y) = exp(-||x-y||^2/s^2)
         """
         nom = np.exp(-(x-y).T.dot(np.linalg.inv(S)).dot(x-y))
         denom = np.sqrt(np.linalg.det(2.*np.pi*S))
         return nom/denom
```

```
In [23]: mu = 0.005
         mu_sigma = 0.001
         n_iter = 150
         # centers = random_centers(X, K=2)
```