

STA414: Statistical Methods for Machine Learning II

Homework 2

Tianyu (Shin) Ren

1002379188

University of Toronto

March 2019

In this assignment, I'll fit both generative and discriminative models to the MNIST dataset of handwritten numbers. Each datapoint in the MNIST [<http://yann.lecun.com/exdb/mnist/>] dataset is a 28x28 black-and-white image of a number in $\{0 \dots 9\}$, and a label indicating which number.

MNIST is the 'fruit fly' of machine learning - a simple standard problem useful for comparing the properties of different algorithms.

For this assignment, I'll binarize the dataset, converting the grey pixel values to either black or white (0 or 1) with > 0.5 being the cutoff. When comparing models, I'll need a training and test set. Use the first 60000 samples for training, and another 10000 for testing.

1. Fitting a Naïve Bayes Model

(a) Derive MLE for the class-conditional pixel means, $\hat{\theta}_{MLE}$

$$L(\theta) = P(X_i, c | \theta, \pi) = \pi_c \prod_{d=1}^{784} \theta_{cd}^{x_d^i} (1 - \theta_{cd})^{(1-x_d^i)}$$

$$\implies l(\theta) = \log(L(\theta)) = \sum_{i=1}^N (\log \pi_c^i + \sum_{d=1}^{784} (x_d^i \log \theta_{cd} + (1 - x_d^i) \log (1 - \theta_{cd})))$$

$$\operatorname{argmax}_{\theta} l(\theta) \quad \text{s.t.} \quad \sum_c \pi_c = 1$$

$$\implies \frac{\partial l}{\partial \theta_{cd}} = \sum_i \mathbb{1}(C^i = c) \left(\frac{x_d^i}{\theta_{cd}} - \frac{1 - x_d^i}{1 - \theta_{cd}} \right) = 0$$

$$\implies \sum_i \mathbb{1}(C^i = c) \theta_{cd} = \sum_i \mathbb{1}(C^i = c) x_d^i$$

$$\implies \hat{\theta}_{cd}^{MLE} = \frac{\sum_i \mathbb{1}(C^i = c) x_d^i}{\sum_i \mathbb{1}(C^i = c)} \quad \text{where } c \in \{0, \dots, 9\}, d \in \{1, \dots, 784\}$$

(b) Derive MAP for the class-conditional pixel means, $\hat{\theta}_{MAP}$

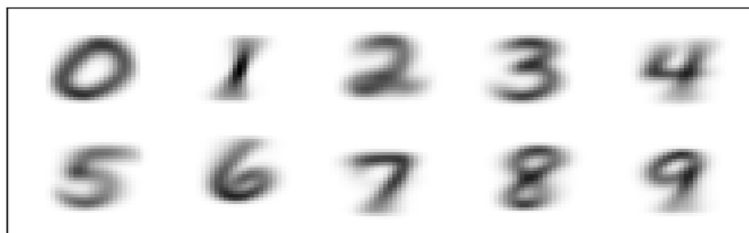
$$P(\theta | \mathbf{x}, c, \pi) \propto P(\theta) P(\mathbf{x}, c | \theta, \pi) \quad \text{where } \theta_{cd} \sim \text{Beta}(2, 2)$$

$$l(\theta) = \log P(\theta) + \log P(\mathbf{x}, c | \theta, \pi) + \text{constant}$$

$$= \log(\theta_{cd}(1 - \theta_{cd})) + \log \pi_c^i + \sum_{d=1}^{784} (x_d^i \log \theta_{cd} + (1 - x_d^i) \log(1 - \theta_{cd})) + \text{constant}$$

$$\begin{aligned} \implies \frac{\partial l}{\partial \theta_{cd}} &= \left(\frac{1}{\theta_{cd}} - \frac{1}{1 - \theta_{cd}} \right) + \sum_{i=1}^N \mathbb{1}(C^i = c) \left(\frac{x_d^i}{\theta_{cd}} - \frac{1 - x_d^i}{1 - \theta_{cd}} \right) = 0 \\ \implies (1 - \theta_{cd}) - \theta_{cd} + \sum_{i=1}^N \mathbb{1}(C^i = c) (x_d^i (1 - \theta_{cd}) - (1 - x_d^i) \theta_{cd}) &= 0 \\ \implies \hat{\theta}_{cd}^{MAP} &= \frac{\sum_{i=1}^N \mathbb{1}(C^i = c) x_d^i + 1}{\sum_{i=1}^N \mathbb{1}(C^i = c) + 2} \end{aligned}$$

- (c) Fit $\hat{\theta}_{MAP}$ to the training set. Plot $\hat{\theta}_{MAP}$ as 10 separate greyscale images, one for each class.



- (d) Derive the log-likelihood $\log p(c|\mathbf{x}, \boldsymbol{\theta}, \pi)$ for a single training image

$$\begin{aligned} \log P(c|\mathbf{x}, \boldsymbol{\theta}, \pi) &= \log P(\mathbf{x}|c, \boldsymbol{\theta}) + \log P(c|\pi) - \log P(\mathbf{x}) \\ &= \log \prod_{d=1}^{784} P(x_d|c, \theta_{cd}) + \log P(c|\pi) - \log \sum_{c=0}^9 P(\mathbf{x}|C=c) P(C=c) \\ &= \sum_{d=1}^{784} (x_d \log \theta_{cd} + (1 - x_d) \log(1 - \theta_{cd})) + \log \pi_c - \log \sum_{c=0}^9 \pi_c \prod_{d=1}^{784} \theta_{cd}^{x_d} (1 - \theta_{cd})^{(1-x_d)} \end{aligned}$$

- (e) Given parameters fit to the training set, and $\pi_c = \frac{1}{10}$:
- the average log-likelihood per datapoint on the training set: -3.15
 - the average log-likelihood per datapoint on the test set: -2.98
 - the accuracy on the training set: 0.834
 - the accuracy on the test set: 0.845

2. Generating from a Naïve Bayes Model

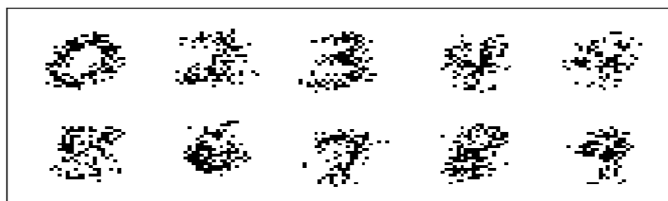
- (a) **True:** Given this model's assumptions, any two pixels x_i and x_j where $i \neq j$ are independent given c .

This is the assumption of Naïve Bayes model.

- (b) **False:** Given this model's assumptions, any two pixels x_i and x_j where $i \neq j$ are independent when marginalizing over c .

x_i, x_j are not independent since $p(x_i, x_j) = \sum_c p(x_i, x_j | c) = \sum_c p(x_i | c) p(x_j | c)$ and $p(x_i) p(x_j) = \sum_c p(x_i | c) \sum_c p(x_j | c)$. So, $p(x_i, x_j) \neq p(x_i) p(x_j)$

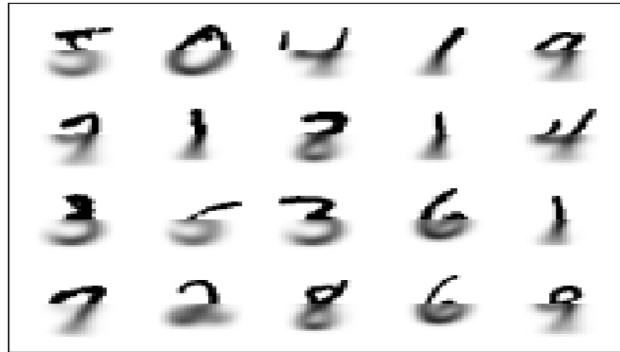
- (c) **Using the parameters fit in question 1, randomly sample and plot 10 binary images from the marginal distribution $p(\mathbf{x} | \theta, \pi)$.**



- (d) **Derive $p(\mathbf{x}_{i \in \text{bottom}} | \mathbf{x}_{\text{top}}, \theta, \pi)$**

$$\begin{aligned}
 p(\mathbf{x}_{i \in \text{bottom}} | \mathbf{x}_{\text{top}}) &= \sum_c p(x_i | x_t, c) p(c | x_t) \\
 &= \sum_c p(x_i | c) \frac{p(x_t | c) p(c)}{p(x_t)} \\
 &= \sum_c p(x_i | c) \frac{\prod_{d=1}^{392} p(x_t^d | c) p(c)}{\sum_{c'} \prod_{d=1}^{392} p(x_t^d | c') p(c')} \\
 &= \frac{\sum_{c=0}^9 \theta_{ci}^{x_i} (1 - \theta_{ci})^{(1-x_i)} \prod_{d=1}^{392} \theta_{cd}^{x_t^d} (1 - \theta_{cd})^{(1-x_t^d)} \pi_c}{\sum_{c'=0}^9 \prod_{d=1}^{392} \theta_{c'd}^{x_t^d} (1 - \theta_{c'd})^{(1-x_t^d)} \pi_{c'}}
 \end{aligned}$$

- (e) **For 20 images from the training set, plot the top half the image concatenated with the marginal distribution over each pixel in the bottom half.**



3. Logistic Regression

Our model will be multiclass logistic regression:

$$p(c|\mathbf{x}, \mathbf{w}) = \frac{\exp(\mathbf{w}_c^T \mathbf{x})}{\sum_{c'=0}^9 \exp(\mathbf{w}_{c'}^T \mathbf{x})}$$

- (a) This model has **7840** parameters. Since each data point \mathbf{x}_i is a vector with dimension $784 * 1$, and there are 10 classes in total. So, it's $784 \times 10 = 7840$
- (b) **Derive the gradient of the predictive log-likelihood w.r.t. \mathbf{w} :** $\nabla_{\mathbf{w}} \log p(c|\mathbf{x}, \mathbf{w})$

$$\log P(C|\mathbf{x}, \mathbf{w}) = \mathbf{w}_c^T \mathbf{x} - \log \sum_{c'=0}^9 \exp(\mathbf{w}_{c'}^T \mathbf{x})$$

To compute $\nabla_{\mathbf{w}_j} \log P(C|\mathbf{x}, \mathbf{w})$,

$$\begin{aligned} \text{if } j = c &\implies \nabla_{\mathbf{w}_c} \log P(C|\mathbf{x}, \mathbf{w}) = \mathbf{x} - \frac{\exp(\mathbf{w}_c^T \mathbf{x})}{\sum_{c'=0}^9 \exp(\mathbf{w}_{c'}^T \mathbf{x})} \mathbf{x} \\ \text{if } j \neq c &\implies \nabla_{\mathbf{w}_j} \log P(C|\mathbf{x}, \mathbf{w}) = -\frac{\exp(\mathbf{w}_j^T \mathbf{x})}{\sum_{c'=0}^9 \exp(\mathbf{w}_{c'}^T \mathbf{x})} \mathbf{x} \end{aligned}$$

$$\begin{aligned} \text{Also, } \nabla_{\mathbf{w}_c} \log \prod_{i=1}^N P(C^i|\mathbf{x}_i, \mathbf{w}) &= \sum_{i=1}^N \nabla_{\mathbf{w}_c} \log P(C^i|\mathbf{x}_i, \mathbf{w}) \\ &= \sum_{i=1}^N \mathbb{1}(C^i = c) \left(\mathbf{x}_i - \frac{\exp(\mathbf{w}_c^T \mathbf{x}_i)}{\sum_{c'=0}^9 \exp(\mathbf{w}_{c'}^T \mathbf{x}_i)} \mathbf{x}_i \right) \end{aligned}$$

- (c) The plot of weights, W , one image per class:



(d) Given parameters fit to the training set:

the average log-likelihood per datapoint on the training set: -0.61

the average log-likelihood per datapoint on the test set: -0.58

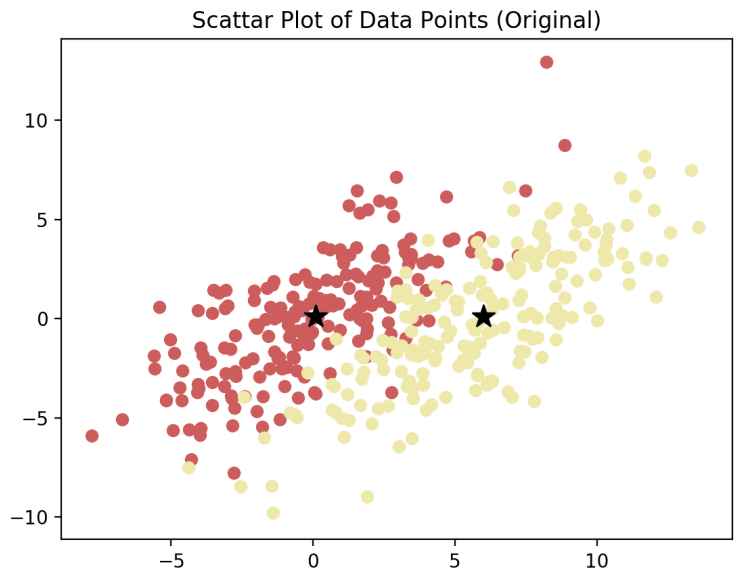
the accuracy on the training set: 0.85

the accuracy on the test set: 0.86

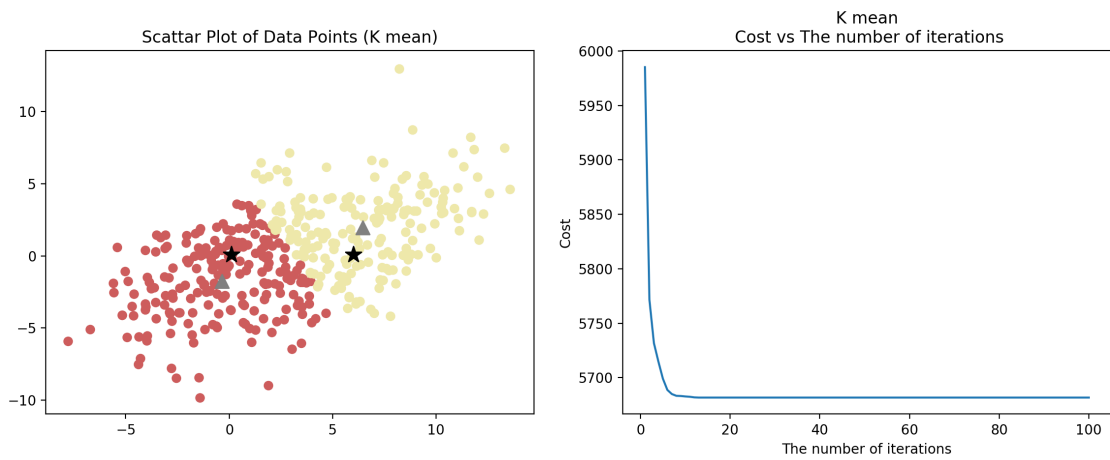
The logistic regression has a better performance than Naïve Bayes both in training and test sets.

4. Unsupervised Learning

(a) First, we will generate some data for this problem. Set the number of points $N = 400$, their dimension $D = 2$, and the number of clusters $K = 2$, and generate data from the distribution $p(x|\mathcal{C}_k) = \mathcal{N}(\mu_k, \Sigma_k)$. Sample 200 data points for \mathcal{C}_1 and 200 for \mathcal{C}_2

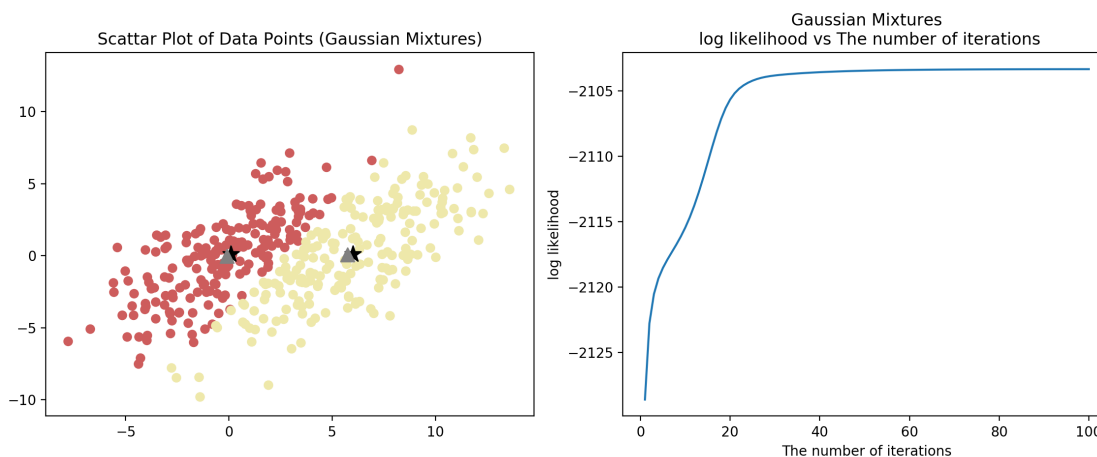


(b) Now, we assume that the true class labels are not known and implement the k-means algorithm.



The misclassification error for k mean: 0.255

(c) Next, implement the EM algorithm for Gaussian mixtures.



The misclassification error for Gaussian mixtures: 0.08

(d) **Comments:**

As we can see the plots from (b) and (c), resulting centers are very different from true means in the K mean algorithm while the centers are almost the same in the EM algorithm. Thus, we had a poor clustering performance in K mean and a good performance in EM. Regarding to the number of iterations, we can see that K mean converges in only a few steps, within 5 iterations while EM takes more than 30 iterations. If I change the covariace matrix to:

$$\Sigma_1 = \Sigma_2 = \begin{bmatrix} 1 & 7 \\ 7 & 1 \end{bmatrix}$$

then both K mean and EM have good performances and K mean converges faster than EM. So by hint, in most cases, k-means should fail to identify the correct cluster labels due to the covariance structure. There may be realizations that EM also fails to find the clusters correctly but in general it should work better than k-means.

Appendix

Python Code

loadMNIST.py

```
1 from __future__ import absolute_import
2 from __future__ import print_function
3 from future.standard_library import install_aliases
4 install_aliases()
5
6 import numpy as np
7 import os
8 import gzip
9 import struct
10 import array
11
12 import matplotlib.pyplot as plt
13 import matplotlib.image
14 from urllib.request import urlretrieve
15
16 def download(url, filename):
17     if not os.path.exists('data'):
18         os.makedirs('data')
19     out_file = os.path.join('data', filename)
20     if not os.path.isfile(out_file):
21         urlretrieve(url, out_file)
22
23 def mnist():
24     base_url = 'http://yann.lecun.com/exdb/mnist/'
25
26     def parse_labels(filename):
27         with gzip.open(filename, 'rb') as fh:
28             magic, num_data = struct.unpack(">II", fh.read(8))
29             return np.array(array.array("B", fh.read()), dtype=np.uint8)
30
31     def parse_images(filename):
32         with gzip.open(filename, 'rb') as fh:
33             magic, num_data, rows, cols = struct.unpack(">IIII", fh.read(16))
34             return np.array(array.array("B", fh.read()), dtype=np.uint8).reshape(
35                 num_data, rows, cols)
36
37     for filename in ['train-images-idx3-ubyte.gz',
38                     'train-labels-idx1-ubyte.gz',
39                     't10k-images-idx3-ubyte.gz',
40                     't10k-labels-idx1-ubyte.gz']:
41         download(base_url + filename, filename)
42
43     train_images = parse_images('data/train-images-idx3-ubyte.gz')
44     train_labels = parse_labels('data/train-labels-idx1-ubyte.gz')
45     test_images = parse_images('data/t10k-images-idx3-ubyte.gz')
46     test_labels = parse_labels('data/t10k-labels-idx1-ubyte.gz')
47
48     return train_images, train_labels, test_images, test_labels
49
```

```

50 def load_mnist():
51     partial_flatten = lambda x : np.reshape(x, (x.shape[0], np.prod(x.shape[1:])))
52     one_hot = lambda x, k: np.array(x[:,None] == np.arange(k)[None, :], dtype=int)
53     train_images, train_labels, test_images, test_labels = mnist()
54     train_images = partial_flatten(train_images) / 255.0
55     test_images = partial_flatten(test_images) / 255.0
56     train_labels = one_hot(train_labels, 10)
57     test_labels = one_hot(test_labels, 10)
58     N_data = train_images.shape[0]
59
60     return N_data, train_images, train_labels, test_images, test_labels
61
62
63 def plot_images(images, ax, ims_per_row=5, padding=5, digit_dimensions=(28, 28),
64               cmap=matplotlib.cm.binary, vmin=None, vmax=None):
65     """Images should be a (N_images x pixels) matrix."""
66     N_images = images.shape[0]
67     N_rows = np.int32(np.ceil(float(N_images) / ims_per_row))
68     pad_value = np.min(images.ravel())
69     concat_images = np.full(((digit_dimensions[0] + padding) * N_rows + padding,
70                             (digit_dimensions[1] + padding) * ims_per_row + padding
71                             ), pad_value)
72     for i in range(N_images):
73         cur_image = np.reshape(images[i, :], digit_dimensions)
74         row_ix = i // ims_per_row
75         col_ix = i % ims_per_row
76         row_start = padding + (padding + digit_dimensions[0]) * row_ix
77         col_start = padding + (padding + digit_dimensions[1]) * col_ix
78         concat_images[row_start: row_start + digit_dimensions[0],
79                       col_start: col_start + digit_dimensions[1]] = cur_image
80     cax = ax.matshow(concat_images, cmap=cmap, vmin=vmin, vmax=vmax)
81     plt.xticks(np.array([]))
82     plt.yticks(np.array([]))
83     return cax
84
85 def save_images(images, filename, **kwargs):
86     fig = plt.figure(1)
87     fig.clf()
88     ax = fig.add_subplot(111)
89     plot_images(images, ax, **kwargs)
90     fig.patch.set_visible(False)
91     ax.patch.set_visible(False)
92     plt.savefig(filename)

```

model.py

```

1 from loadMNIST import *
2 from scipy.special import logsumexp
3 import time
4 np.random.seed(1)
5
6
7 COLORS = ["indianred", "palegoldenrod", "black", "gray"]
8
9
10 def get_images_by_label(images, labels, query_label):
11     """
12     Helper function to return all images in the provided array which match the query
13     label.

```

```

13     """
14     assert images.shape[0] == labels.shape[0]
15     matching_indices = labels == query_label
16     return images[matching_indices]
17
18
19 class NaiveBayes:
20     """
21     Q1, Naive Bayes model.
22     """
23     def __init__(self, train_images, train_labels):
24         self.train_images = train_images
25         self.train_labels = train_labels
26
27     def map_naive_bayes(self, plot=False):
28         """
29         return a matrix 10 * 784 where each row represents a class.
30         """
31         theta = np.zeros((10, 784))
32         for c in range(10):
33             images = get_images_by_label(self.train_images, self.train_labels, c)
34             theta[c] = np.divide(np.sum(images, axis=0) + 1., images.shape[0] + 2.)
35         if plot:
36             save_images(theta, "theta_map.png")
37         return theta
38
39     def log_likelihood(self, X, y, theta):
40         """
41         return a matrix N * 10 where each row represents log likelihood of a data point,
42         and each column represents log lilihood of a class.
43         """
44         ll = np.zeros((X.shape[0], 10))
45         log_p_x = logsumexp(np.log(0.1) + np.dot(X, np.log(theta.T)) + np.dot((1. - X),
46         np.log(1. - theta.T)), axis=1)
47         for c in range(10):
48             ll[:, c] = np.dot(X, np.log(theta[c])) + np.dot((1. - X), np.log(1. - theta[c]
49             )) + np.log(0.1) - log_p_x
50         return ll
51
52     def avg_log_likelihood(self, X, y, theta):
53         ll = 0
54         for c in range(10):
55             X_c = get_images_by_label(X, y, c)
56             log_p_x = logsumexp(np.log(0.1) + np.dot(X_c, np.log(theta.T)) + np.dot((1. -
57             X_c), np.log(1. - theta.T)), axis=1)
58             ll += np.sum(np.dot(X_c, np.log(theta[c])) + np.dot((1. - X_c), np.log(1. -
59             theta[c])) + np.log(0.1) - log_p_x)
60         return ll / X.shape[0]
61
62     def predict(self, X, y, theta, train=False, test=False):
63         ll = self.log_likelihood(X, y, theta)
64         pred = np.argmax(ll, axis=1)
65         avg_ll = self.avg_log_likelihood(X, y, theta)
66         accuracy = np.mean(pred == y)
67         name = "test" if test else "train"
68         print("average log-likelihood of naive bayes model on the {} set: ".format(name)
69         + str(avg_ll))
70         print("accuracy of naive bayes model on the {} set: ".format(name) + str(

```

```

accuracy))
66
67
68 class GenerativeNaiveBayes:
69     """
70     Q2, Generating from a Naive Bayes Model
71     """
72     def __init__(self, theta):
73         self.theta = theta
74
75     def sample_plot(self):
76         """
77         randomly sample and plot 10 binary images from the marginal distribution, p(x|
78         theta, pi)
79         """
80         c = np.random.multinomial(10, [0.1]*10)
81         images = np.zeros((10, 784))
82         count = 0
83         for i in range(10):
84             for j in range((c[i])):
85                 images[count] = np.random.binomial(1, self.theta[i]).reshape((1, 784))
86                 count += 1
87         save_images(images, "samples.png")
88
89     def predict_half(self, X_top):
90         """
91         plot the top half the image concatenated with the marginal distribution over
92         each pixel in the bottom half.
93         """
94         X_bot = np.zeros((X_top.shape[0], X_top.shape[1]))
95         theta_top, theta_bot = self.theta[:, :392].T, self.theta[:, 392:].T
96         for i in range(392):
97             constant = np.dot(X_top, np.log(theta_top)) + np.dot(1 - X_top, np.log(1 -
98             theta_top))
99             X_bot[:, i] = logsumexp(np.add(constant, np.log(theta_bot[i])), axis=1) -
100             logsumexp(constant, axis=1)
101             save_images(np.concatenate((X_top, np.exp(X_bot)), axis=1), "predict_half.png")
102
103 class LogisticRegression:
104     """
105     Q3, Fitting a simple predictive model using gradient descent.
106     Our model will be multiclass logistic regression.
107     """
108     def __init__(self, train_images, train_labels):
109         self.train_images = train_images
110         self.train_labels = train_labels
111         self.W = np.zeros((10, 784))
112
113     def softmax(self, X, W):
114         """
115         return a N * 10 vector where each row is a data point
116         and each column is the probability of that class.
117         """
118         return (np.exp(np.dot(X, W.T)).T / np.exp(logsumexp(np.dot(X, W.T), axis=1))).T
119
120     def grad_pred_ll(self, X, W, c):
121         """

```

```

119     This function calculate the gradient of the predictive log-likelihood.
120     return a 1 * 784 vector
121     """
122     constant = np.exp(logsumexp(np.dot(X, W.T), axis=1))
123     return np.sum(X - (X.T * np.divide(np.exp(np.dot(X, W[c])), constant)).T, axis
124     =0)
125
126 def gradient_ascent(self, lr=0.00001, iters=100):
127     for _ in range(iters):
128         prob = self.softmax(self.train_images, self.W)
129         pred = np.argmax(prob, axis=1)
130         accuracy = np.mean(pred == self.train_labels)
131         print("training accuracy: {}, iterations: {}/{}".format(round(accuracy, 2), _,
132         iters))
133         for c in range(10):
134             X_c = get_images_by_label(self.train_images, self.train_labels, c)
135             self.W[c] = self.W[c] + lr * self.grad_pred_ll(X_c, self.W, c)
136
137 def log_likelihood(self, X, y, W):
138     ll = 0
139     for c in range(10):
140         X_c = get_images_by_label(X, y, c)
141         ll += np.sum(np.dot(X_c, W[c]) - logsumexp(np.dot(X_c, W.T), axis=1))
142     return ll / X.shape[0]
143
144 def predict(self, X, y, train=False, test=False):
145     if train:
146         self.gradient_ascent()
147         save_images(self.W, "weights.png")
148         avg_ll = self.log_likelihood(X, y, self.W)
149         pred = np.argmax(self.softmax(X, self.W), axis=1)
150         accuracy = np.mean(pred == y)
151         name = "test" if test else "train"
152         print("average log-likelihood of softmax model on the {} set: ".format(name) +
153         str(avg_ll))
154         print("accuracy of softmax model on the {} set: ".format(name) + str(accuracy))
155
156
157 class EM:
158     """
159     Q4, EM algorithm for K means and Gaussian mixtures.
160     """
161     def __init__(self, initials, c1, c2):
162         self.initials = initials
163         self.data = np.concatenate((c1, c2), axis=0)
164         self.N, self.D = self.data.shape # Data is a N * D matrix, and here N=400,
165         D=2
166
167         # Initial values for K mean and GMM
168         self.miu_hat = np.concatenate((self.initials['MIU1_HAT'], self.initials['
169         MIU2_HAT']), axis=0).reshape((2,2))
170         self.clusters = np.concatenate((np.zeros(int(self.N/2)), np.ones(int(self.N/2)))
171         , axis=0)
172         self.costs_iter = [[], []]
173
174     def plot_clusters(self, km=False, gmm=False):
175         """
176         a scatter plot of the data points showing the true cluster assignment of each

```

```

171     point.
172     """
173     Also plot a scatter plot of K mean or gaussian mixtures.
174     """
175     f2 = plt.figure()
176     ax2 = f2.add_subplot(111)
177     for i in range(self.D):
178         plt.scatter(self.data[self.clusters == i][:, 0], self.data[self.clusters == i]
179                    [:, 1], c=COLORS[i])
180         plt.scatter(self.initials['MIU'+str(i+1)][0], self.initials['MIU'+str(i+1)]
181                    [1], marker='*', c=COLORS[2], s=150)
182         plt.title("Scattar Plot of Data Points (Original)")
183         if km or gmm:
184             name = "K mean" if km else "Gaussian Mixtures"
185             plt.scatter(self.miu_hat[:,0], self.miu_hat[:,1], marker='^', c=COLORS[3], s
186                       =100)
187             plt.title("Scattar Plot of Data Points ({}).format(name)")
188
189     def misclassification_error(self):
190         return (np.sum(self.clusters[:int(self.N/2)] == 1) + np.sum(self.clusters[int(
191         self.N/2):] == 0)) / self.N
192
193     class KMean(EM):
194         def __init__(self, initials, c1, c2):
195             super().__init__(initials, c1, c2)
196
197         def cost(self):
198             cost = 0
199             for i in range(self.D):
200                 cost += np.sum(np.linalg.norm(self.data[self.clusters == i] - self.miu_hat[i],
201                 axis=1) ** 2)
202             return cost
203
204         def km_e_step(self):
205             distances = np.zeros((self.N, 2))
206             for i in range(self.D):
207                 distances[:, i] = np.linalg.norm(self.data - self.miu_hat[i], axis=1)
208             self.clusters = np.argmin(distances, axis=1)
209
210         def km_m_step(self):
211             for i in range(self.D):
212                 self.miu_hat[i] = np.mean(self.data[self.clusters == i], axis=0)
213
214         def train(self, max_iter=100):
215             i = 1
216             while i <= max_iter:
217                 self.km_e_step()
218                 self.km_m_step()
219                 self.costs_iter[0].append(self.cost())
220                 self.costs_iter[1].append(i)
221                 i += 1
222             f3 = plt.figure()
223             ax3 = f3.add_subplot(111)
224             plt.plot(self.costs_iter[1], self.costs_iter[0])
225             plt.title("K mean\n Cost vs The number of iterations")
226             plt.xlabel("The number of iterations")
227             plt.ylabel("Cost")
228             self.plot_clusters(km=True)

```

```

223     print(" misclassification error for k mean: " + str(self.misclassification_error
224           ()))
225
226 class GaussianMixtures(EM):
227     def __init__(self, initials, c1, c2):
228         super().__init__(initials, c1, c2)
229         # Initial values for Gaussian mixtures
230         self.simga_hat = [np.eye(self.D)] * 2
231         self.pi_hat = [0.5, 0.5] # Mixing proportions
232         self.R = np.zeros((self.N, 2)) # This is the posterior/responsibilities
233         self.Nk = [] # The number of data in class K
234
235     def normal_density(self, X, miu, sigma):
236         """
237         This is a vectorized normal_density, where X is N*D, miu is 1*D, sigma is D*D
238         Output is N*1, where each element is a pdf value.
239         """
240         constant = 1 / np.sqrt((2 * np.pi) ** self.D * np.linalg.det(sigma))
241         return constant * np.diag(np.exp(-0.5 * np.dot(np.dot((X - miu), np.linalg.inv(
242           sigma)), (X - miu).T)))
243
244     def log_likelihood(self):
245         normal_sum = np.zeros(self.N)
246         for i in range(self.D):
247             normal_sum += self.pi_hat[i] * self.normal_density(self.data, self.miu_hat[i],
248               self.simga_hat[i])
249         return np.sum(np.log(normal_sum))
250
251     def em_e_step(self):
252         for i in range(self.D):
253             self.R[:, i] = self.pi_hat[i] * self.normal_density(self.data, self.miu_hat[i],
254               self.simga_hat[i])
255         # Normalize R
256         self.R = (self.R.T / np.sum(self.R, axis=1)).T
257         # assign datapoints to each gaussian
258         self.Nk = np.sum(self.R, axis = 0)
259
260     def em_m_step(self):
261         for i in range(self.D):
262             self.miu_hat[i] = 1. / self.Nk[i] * np.sum(self.R[:, i] * self.data.T, axis
263               =1).T
264             diff = self.data - self.miu_hat[i]
265             self.simga_hat[i] = 1. / self.Nk[i] * np.dot(np.multiply(diff.T, self.R[:, i]
266               ), diff)
267             self.pi_hat[i] = self.Nk[i] / self.N
268
269     def train(self, max_iter=100):
270         i = 1
271         while i <= max_iter:
272             self.em_e_step()
273             self.em_m_step()
274             self.costs_iter[0].append(self.log_likelihood())
275             self.costs_iter[1].append(i)
276             i += 1
277         f4 = plt.figure()
278         ax4 = f4.add_subplot(111)
279         plt.plot(self.costs_iter[1], self.costs_iter[0])

```

```

275     plt.title("Gaussian Mixtures\n log likelihood vs The number of iterations")
276     plt.xlabel("The number of iterations")
277     plt.ylabel("log likelihood")
278     self.clusters = np.argmax(self.R, axis=1)
279     self.plot_clusters(gmm=True)
280     print("misclassification error for gmm: " + str(self.misclassification_error()))
281
282
283 if __name__ == '__main__':
284     start = time.time()
285     print("loading data...")
286     N_data, train_images, train_labels, test_images, test_labels = load_mnist()
287     train_labels = np.argmax(train_labels, axis=1)
288     test_labels = np.argmax(test_labels, axis=1)
289
290     print("training a Naive Bayes model...")
291     nb_model = NaiveBayes(train_images, train_labels)
292     theta_map = nb_model.map_naive_bayes(plot=True)
293     nb_model.predict(train_images, train_labels, theta_map, train=True)
294     nb_model.predict(test_images, test_labels, theta_map, test=True)
295
296     print("training a generative Naive Bayes model...")
297     gnb = GenerativeNaiveBayes(theta_map)
298     gnb.sample_plot()
299     gnb.predict_half(train_images[:20, :392])
300
301     print("training a softmax model...")
302     lr_model = LogisticRegression(train_images, train_labels)
303     lr_model.predict(train_images, train_labels, train=True)
304     lr_model.predict(test_images, test_labels, test=True)
305
306     print("training K mean and GMM algorithms...")
307     initials = {'Nk': 200,
308                'MIU1': np.array([0.1, 0.1]),
309                'MIU2': np.array([6., 0.1]),
310                'COV': np.array([[10., 7.], [7., 10.]]),
311                'MIU1.HAT': np.array([0., 0.]),
312                'MIU2.HAT': np.array([1., 1.])
313                }
314     # Sampling data from a multivariate gaussian distribution
315     c1 = np.random.multivariate_normal(initials['MIU1'], initials['COV'], initials['Nk'])
316     c2 = np.random.multivariate_normal(initials['MIU2'], initials['COV'], initials['Nk'])
317     kmean = KMean(initials, c1, c2)
318     kmean.plot_clusters()
319     kmean.train()
320     gmm = GaussianMixtures(initials, c1, c2)
321     gmm.train()
322     end = time.time()
323     print("running time: {}s".format(round(end - start, 2)))
324     plt.show()

```