```python
1   import matplotlib.pyplot as plt
2   import numpy as np
3   import os
4
5   PLOT_COLORS = ['red', 'green', 'blue', 'orange']  # Colors for your plots
6   K = 4               # Number of Gaussians in the mixture model
7   NUM_TRIALS = 3   # Number of trials to run (can be adjusted for debugging)
8   UNLABELED = -1   # Cluster label for unlabeled data points (do not change)
9
10
11  def main(is_semi_supervised, trial_num):
12      """Problem 3: EM for Gaussian Mixture Models (unsupervised and semi-supervised)"""
13      print('Running {} EM algorithm...'
14            .format('semi-supervised' if is_semi_supervised else 'unsupervised'))
15
16      # Load dataset
17      train_path = os.path.join('.', 'train.csv')
18      x_all, z_all = load_gmm_dataset(train_path)
19
20      # Split into labeled and unlabeled examples
21      labeled_idxs = (z_all != UNLABELED).squeeze()
22      x_tilde = x_all[labeled_idxs, :]    # Labeled examples
23      z_tilde = z_all[labeled_idxs, :]    # Corresponding labels
24      x = x_all[~labeled_idxs, :]         # Unlabeled examples
25
26      # *** START CODE HERE ***
27      # (1) Initialize mu and sigma by splitting the n_examples data points uniformly at random
28      # into K groups, then calculating the sample mean and covariance for each group
29      n, d = x.shape
30      group = np.random.choice(K, n)
31      mu = [np.mean(x[group == g, :], axis=0) for g in range(K)]
32      sigma = [np.cov(x[group == g, :].T) for g in range(K)]
33
34      # (2) Initialize phi to place equal probability on each Gaussian
35      # phi should be a numpy array of shape (K,)
36      phi = np.full((K,), fill_value=(1. / K), dtype=np.float32)
37
38      # (3) Initialize the w values to place equal probability on each Gaussian
39      # w should be a numpy array of shape (m, K)
40      w = np.full((n, K), fill_value=(1. / K), dtype=np.float32)
41      # *** END CODE HERE ***
42
43      if is_semi_supervised:
44          w = run_semi_supervised_em(x, x_tilde, z_tilde, w, phi, mu, sigma)
45      else:
46          w = run_em(x, w, phi, mu, sigma)
47
48      # Plot your predictions
49      z_pred = np.zeros(n)
50      if w is not None:  # Just a placeholder for the starter code
51          for i in range(n):
52              z_pred[i] = np.argmax(w[i])
53
54      plot_gmm_preds(x, z_pred, is_semi_supervised, plot_id=trial_num)
55
56
57  def run_em(x, w, phi, mu, sigma):
58      """Problem 3(d): EM Algorithm (unsupervised).
59
60      See inline comments for instructions.
61
62      Args:
63          x: Design matrix of shape (n_examples, dim).
64          w: Initial weight matrix of shape (n_examples, k).
65          phi: Initial mixture prior, of shape (k,).
66          mu: Initial cluster means, list of k arrays of shape (dim,).
67          sigma: Initial cluster covariances, list of k arrays of shape (dim, dim).
68
69      Returns:
70          Updated weight matrix of shape (n_examples, k) resulting from EM algorithm.
71          More specifically, w[i, j] should contain the probability of
72          example x^(i) belonging to the j-th Gaussian in the mixture.
```

```python
 73          """
 74          # No need to change any of these parameters
 75          eps = 1e-3   # Convergence threshold
 76          max_iter = 1000
 77
 78          # Stop when the absolute change in log-likelihood is < eps
 79          # See below for explanation of the convergence criterion
 80          it = 0
 81          ll = prev_ll = None
 82          while it < max_iter and (prev_ll is None or np.abs(ll - prev_ll) >= eps):
 83              pass  # Just a placeholder for the starter code
 84              # *** START CODE HERE
 85              # (1) E-step: Update your estimates in w
 86              w = e_step(x, w, phi, mu, sigma)
 87
 88              # (2) M-step: Update the model parameters phi, mu, and sigma
 89              phi, mu, sigma = m_step(x, w, mu, sigma)
 90
 91              # (3) Compute the log-likelihood of the data to check for convergence.
 92              # By log-likelihood, we mean `ll = sum_x[log(sum_z[p(x|z) * p(z)])]`.
 93              # We define convergence by the first iteration where abs(ll - prev_ll) < eps.
 94              # Hint: For debugging, recall part (a). We showed that ll should be monotonically increasing.
 95              prev_ll = ll
 96              ll = log_likelihood(x, phi, mu, sigma)
 97              it += 1
 98              print('[iter: {:03d}, log-likelihood: {:.4f}]'.format(it, ll))
 99              # *** END CODE HERE ***
100
101      return w
102
103
104  def run_semi_supervised_em(x, x_tilde, z_tilde, w, phi, mu, sigma):
105      """Problem 3(e): Semi-Supervised EM Algorithm.
106
107      See inline comments for instructions.
108
109      Args:
110          x: Design matrix of unlabeled examples of shape (n_examples_unobs, dim).
111          x_tilde: Design matrix of labeled examples of shape (n_examples_obs, dim).
112          z_tilde: Array of labels of shape (n_examples_obs, 1).
113          w: Initial weight matrix of shape (n_examples, k).
114          phi: Initial mixture prior, of shape (k,).
115          mu: Initial cluster means, list of k arrays of shape (dim,).
116          sigma: Initial cluster covariances, list of k arrays of shape (dim, dim).
117
118      Returns:
119          Updated weight matrix of shape (n_examples, k) resulting from semi-supervised EM algorithm.
120          More specifically, w[i, j] should contain the probability of
121          example x^(i) belonging to the j-th Gaussian in the mixture.
122      """
123      # No need to change any of these parameters
124      alpha = 20.  # Weight for the labeled examples
125      eps = 1e-3    # Convergence threshold
126      max_iter = 1000
127
128      # Stop when the absolute change in log-likelihood is < eps
129      # See below for explanation of the convergence criterion
130      it = 0
131      ll = prev_ll = None
132      while it < max_iter and (prev_ll is None or np.abs(ll - prev_ll) >= eps):
133          pass  # Just a placeholder for the starter code
134          # *** START CODE HERE ***
135          # (1) E-step: Update your estimates in w
136          w = e_step(x, w, phi, mu, sigma)
137
138          # (2) M-step: Update the model parameters phi, mu, and sigma
139          phi, mu, sigma = m_step_ss(x, x_tilde, z_tilde, w, phi, mu, sigma, alpha)
140
141          # (3) Compute the log-likelihood of the data to check for convergence.
142          # Hint: Make sure to include alpha in your calculation of ll.
143          # Hint: For debugging, recall part (a). We showed that ll should be monotonically increasing.
144          prev_ll = ll
```

```python
145             ll = log_likelihood(x, phi, mu, sigma)
146             ll += alpha * log_likelihood(x_tilde, phi, mu, sigma, z_tilde)
147             it += 1
148             print('[iter: {:03d}, log-likelihood: {:.4f}]'.format(it, ll))
149             # *** END CODE HERE ***
150
151     return w
152
153
154 # *** START CODE HERE ***
155 # Helper functions
156
157 def e_step(x, w, phi, mu, sigma):
158     """E-step for both unsupervised and semi-supervised EM."""
159     n, d = x.shape
160     k = len(mu)
161
162     for i in range(n):
163         for j in range(k):
164             w[i, j] = p_x_given_z(x[i], mu[j], sigma[j]) * phi[j]
165
166     w /= np.sum(w, axis=1, keepdims=True)
167
168     return w
169
170
171 def m_step(x, w, mu, sigma):
172     """M-step for unsupervised EM."""
173     n, d = x.shape
174     k = len(mu)
175
176     phi = np.mean(w, axis=0)
177
178     for j in range(k):
179         w_j = w[:, j:j + 1]
180         mu[j] = np.sum(w_j * x, axis=0) / np.sum(w_j)
181
182         sigma[j] = np.zeros_like(sigma[j])
183         for i in range(n):
184             x_minus_mu = x[i] - mu[j]
185             sigma[j] += w[i, j] * np.outer(x_minus_mu, x_minus_mu)
186         sigma[j] /= np.sum(w_j)
187
188     return phi, mu, sigma
189
190
191 def m_step_ss(x, x_tilde, z_tilde, w, phi, mu, sigma, alpha):
192     """M-step for semi-supervised EM."""
193     n, _ = x.shape
194     n_tilde, _ = x_tilde.shape
195     k = len(mu)
196
197     w_colsums = np.sum(w, axis=0)
198     k_counts = [np.sum(z_tilde == j) for j in range(k)]
199     for j in range(k):
200         phi[j] = (w_colsums[j] + alpha * k_counts[j]) / (n + alpha * n_tilde)
201
202         w_j = w[:, j:j + 1]
203         mu[j] = ((np.sum(w_j * x, axis=0)
204                   + alpha * np.sum(x_tilde[(z_tilde == j).squeeze(), :], axis=0))
205                  / (np.sum(w_j) + alpha * k_counts[j]))
206
207         sigma[j] = np.zeros_like(sigma[j])
208         for i in range(n):
209             x_minus_mu = x[i] - mu[j]
210             sigma[j] += w[i, j] * np.outer(x_minus_mu, x_minus_mu)
211         for i in range(n_tilde):
212             if z_tilde[i] == j:
213                 x_minus_mu = x_tilde[i] - mu[j]
214                 sigma[j] += alpha * np.outer(x_minus_mu, x_minus_mu)
215         sigma[j] /= (np.sum(w_j) + alpha * k_counts[j])
216
```

```python
217         return phi, mu, sigma
218
219
220 def log_likelihood(x, phi, mu, sigma, z=None):
221     """Get log-likelihood of the data `x` given model parameters
222     `phi`, `mu`, and `sigma`.
223     """
224     n, d = x.shape
225     k = len(phi)
226     ll = 0.
227     for i in range(n):
228         if z is None:  # Unsupervised case
229             p_x = 0.
230             for j in range(k):
231                 p_x += p_x_given_z(x[i], mu[j], sigma[j]) * phi[j]
232         else:  # Supervised case
233             j = int(z[i])
234             p_x = p_x_given_z(x[i], mu[j], sigma[j]) * phi[j]
235         ll += np.log(p_x)
236
237     return ll
238
239
240 def p_x_given_z(x, mu, sigma):
241     """Get probability of a single example `x` given model parameters
242     `mu` and `sigma` (corresponding to cluster z = j).
243     """
244     d = len(x)
245     assert d == len(mu) and sigma.shape == (d, d), 'Shape mismatch.'
246
247     c = 1. / ((2. * np.pi) ** (d / 2) * np.sqrt(np.linalg.det(sigma)))
248     x_minus_mu = x - mu
249     sigma_inv = np.linalg.inv(sigma)
250     p_val = c * np.exp(-.5 * x_minus_mu.dot(sigma_inv).dot(x_minus_mu.T))
251
252     return p_val
253 # *** END CODE HERE ***
254
255
256 def plot_gmm_preds(x, z, with_supervision, plot_id):
257     """Plot GMM predictions on a 2D dataset `x` with labels `z`.
258
259     Write to the output directory, including `plot_id`
260     in the name, and appending 'ss' if the GMM had supervision.
261
262     NOTE: You do not need to edit this function.
263     """
264     plt.figure(figsize=(12, 8))
265     plt.title('{} GMM Predictions'.format('Semi-supervised' if with_supervision else 'Unsupervised'))
266     plt.xlabel('x_1')
267     plt.ylabel('x_2')
268
269     for x_1, x_2, z_ in zip(x[:, 0], x[:, 1], z):
270         color = 'gray' if z_ < 0 else PLOT_COLORS[int(z_)]
271         alpha = 0.25 if z_ < 0 else 0.75
272         plt.scatter(x_1, x_2, marker='.', c=color, alpha=alpha)
273
274     file_name = 'pred{}_{}.pdf'.format('_ss' if with_supervision else '', plot_id)
275     save_path = os.path.join('.', file_name)
276     plt.savefig(save_path)
277
278
279 def load_gmm_dataset(csv_path):
280     """Load dataset for Gaussian Mixture Model.
281
282     Args:
283         csv_path: Path to CSV file containing dataset.
284
285     Returns:
286         x: NumPy array shape (n_examples, dim)
287         z: NumPy array shape (n_exampls, 1)
288
```

```python
289        NOTE: You do not need to edit this function.
290        """
291
292        # Load headers
293        with open(csv_path, 'r') as csv_fh:
294            headers = csv_fh.readline().strip().split(',')
295
296        # Load features and labels
297        x_cols = [i for i in range(len(headers)) if headers[i].startswith('x')]
298        z_cols = [i for i in range(len(headers)) if headers[i] == 'z']
299
300        x = np.loadtxt(csv_path, delimiter=',', skiprows=1, usecols=x_cols, dtype=float)
301        z = np.loadtxt(csv_path, delimiter=',', skiprows=1, usecols=z_cols, dtype=float)
302
303        if z.ndim == 1:
304            z = np.expand_dims(z, axis=-1)
305
306        return x, z
307
308
309    if __name__ == '__main__':
310        np.random.seed(229)
311        # Run NUM_TRIALS trials to see how different initializations
312        # affect the final predictions with and without supervision
313        for t in range(NUM_TRIALS):
314            main(is_semi_supervised=False, trial_num=t)
315
316            # *** START CODE HERE ***
317            # Once you've implemented the semi-supervised version,
318            # uncomment the following line.
319            # You do not need to add any other lines in this code block.
320
321            # main(is_semi_supervised=True, trial_num=t)
322
323            # *** END CODE HERE ***
```