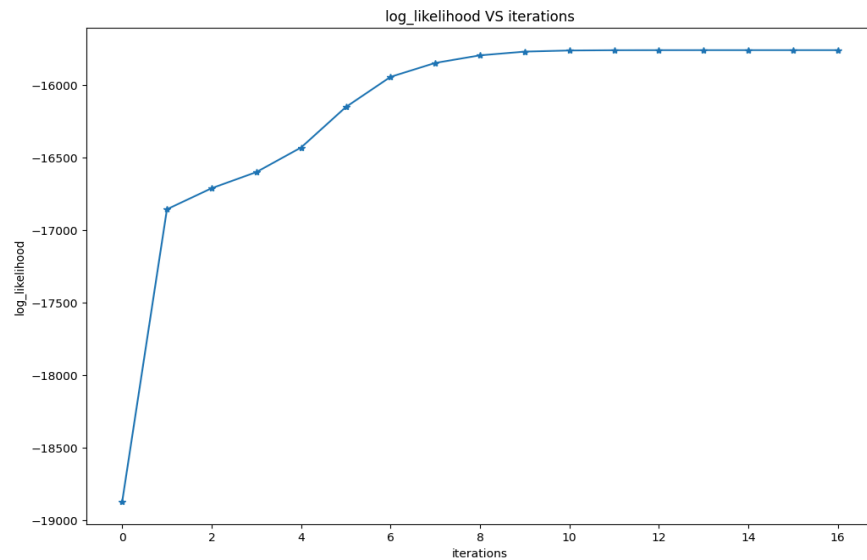


Homework 3

Xiaofan Jiao

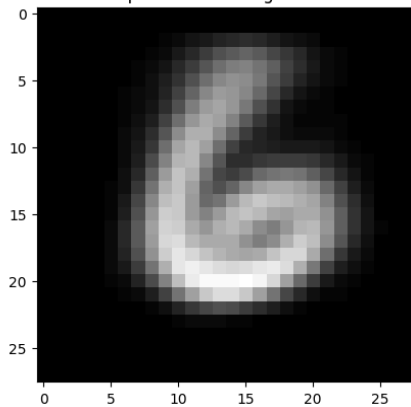
Question 1. Implementing EM for MNIST dataset.

1). The log-likelihood increases with each iteration, indicating that the EM algorithm is successfully improving the model parameters to better fit the data.

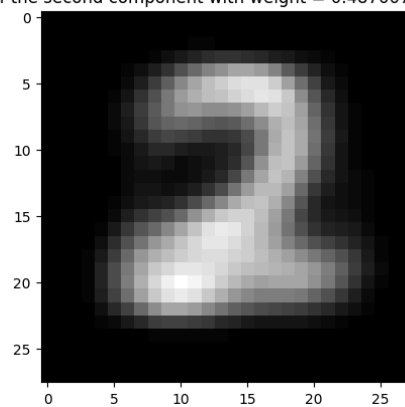


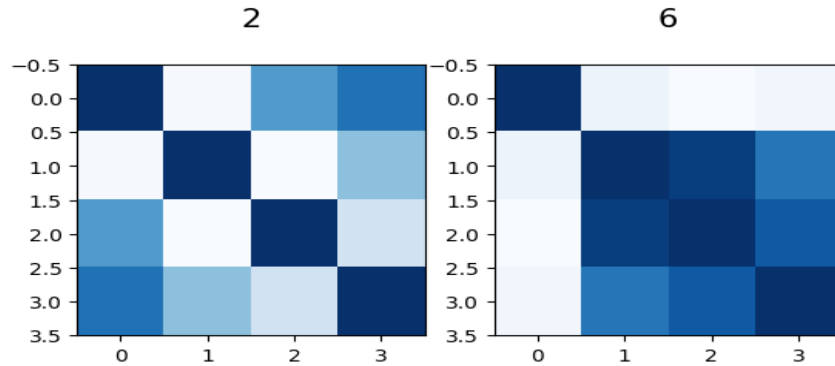
2). After performing the GMM algorithm, the weights for μ_1 and μ_2 were 0.52 and 0.49 respectively.

Mean of the first component with weight = 0.5129921989362725



Mean of the second component with weight = 0.4870078010637276





3). GMM achieves better performance overall.

- For digit '2', both GMM and K-means have relatively similar misclassification rates, with GMM being slightly worse (6.4922%) compared to K-means (6.1047%).
- For digit '6', GMM performs significantly better with a misclassification rate of 0.9395%, whereas K-means has a higher misclassification rate of 7.9332%.
- The misclassification rates are as follows:

	'2'	'6'
GMM	6.4922%	0.9395%
K-means	6.1047%	7.9332%

Question 2. Optimization

1).

HW 3. Question 2. Optimization:

1.
$$l(\theta) = \sum_{i=1}^m [y^i \theta^T x^i - \log(1 + \exp(\theta^T x^i))]$$

① differentiate $l(\theta)$ with respect to θ

$$\frac{\partial l(\theta)}{\partial \theta} = \sum_{i=1}^m \left[y^i x^i - \frac{\exp(\theta^T x^i)}{1 + \exp(\theta^T x^i)} x^i \right]$$

② Simplify term using $\sigma(z) = \frac{1}{1 + \exp(-z)}$

$$\frac{\exp(\theta^T x^i)}{1 + \exp(\theta^T x^i)} = \sigma(\theta^T x^i)$$

③ Thus, the gradient becomes:

$$\frac{\partial l(\theta)}{\partial \theta} = \sum_{i=1}^m [(y^i - \sigma(\theta^T x^i)) x^i]$$

2).

2. ① Initialize:

Set the initial value of θ , define learning rate $\theta = \theta_0$

② Iterate:

Until convergence (the change in θ is smaller than a threshold or a max number of iterations)

Compute the gradient of the cost function with respect to θ

Update θ by taking a step in the direction of the gradient.

1 Initialize $\theta = \theta_0$

1 Set learning rate α

1 Repeat until convergence

1 { gradient = 0

1 for $i = 1$ to m { gradient = gradient + $(y^i - \delta(\theta^T x^i)) x^i$

1 $\theta = \theta + \alpha \times \text{gradient}$ }

3).

3. The Stochastic Gradient Descent (SGD) updates the parameters for each training example instead of the whole dataset.

① Initialize:

Initialize $\theta = \theta_0$

Set learning rate α

② Repeat until convergence

• For each training example $i = 1$ to m

• Compute the gradient

$$\text{gradient} = (y^i - \delta(\theta^T x^i)) x^i$$

• Update the parameters:

$$\theta = \theta + \alpha \cdot \text{gradient}$$

• Gradient Descent uses the entire dataset to compute the gradient, leading to more stable convergence but can be computationally expensive. For larger datasets, it will require more computing power.

• Stochastic Gradient Descent updates the gradient for each training example, making it faster and able to handle large datasets, but the path to convergence may be noisier. When we have large dataset, SGD is often more computationally efficient when compared with GD.

4).

4. From previous questions we know that

Log-likelihood Function:

$$l(\theta) = \sum_{i=1}^m [y^i \theta^T x^i - \log(1 + \exp(\theta^T x^i))]$$

Gradient of $l(\theta)$:

$$\frac{dl(\theta)}{d\theta} = \sum_{i=1}^m [y^i - \sigma(\theta^T x^i)] x^i$$

where $\sigma(z) = \frac{1}{1 + \exp(-z)}$

Hessian Matrix: Second-order derivative of $l(\theta)$

$$H(\theta) = \frac{d^2 l(\theta)}{d\theta^2}$$

First Layer derivative

$$\frac{dl(\theta)}{d\theta} = \sum_{i=1}^m \left(y^i - \frac{1}{1 + \exp(\theta^T x^i)} \right) x^i$$

$$H(\theta) = \frac{d^2 l(\theta)}{d\theta^2} = \sum_{i=1}^m \left(\frac{x^i x^{iT} \exp(\theta^T x^i)}{(1 + \exp(\theta^T x^i))^2} \right) x^i$$

$$\frac{d^2 l(\theta)}{d\theta^2} = \sum_{i=1}^m \left(\frac{\exp(\theta^T x^i)}{(1 + \exp(\theta^T x^i))^2} \right)$$

$$\frac{d^2 l(\theta)}{d\theta^2} = \sum_{i=1}^m \frac{1}{(1 + \exp(\theta^T x^i))} \times \frac{\exp(\theta^T x^i)}{(1 + \exp(\theta^T x^i))}$$

$$\frac{d^2 l(\theta)}{d\theta^2} = \sum_{i=1}^m \left(\frac{1}{(1 + \exp(\theta^T x^i))} \left(\frac{1}{1 + \exp(\theta^T x^i)} - 1 \right) \right)$$

Since $0 < \frac{1}{(1 + \exp(\theta^T x^i))} < 1$

Thus $\sum_{i=1}^m \left(\frac{1}{(1 + \exp(\theta^T x^i))} \left(\frac{1}{1 + \exp(\theta^T x^i)} - 1 \right) \right) < 0$

This ensures that $H(\theta)$ is negative semi-definite and $l(\theta)$ is concave. The concavity of $l(\theta)$ makes sure that there is a unique global maximum. This property is crucial because it means that optimization algorithms like gradient descent will converge to this unique global optimizer. Gradient descent is a reliable algorithm for finding the maximum of $l(\theta)$ because the concavity ensures no local maxima exist apart from the global maximum. The negative semi-definite Hessian means the function curves downwards, facilitating efficient optimization.

Question 3. Bayes Classifier for spam filtering

1).

HW3 Question 3

1. • Total messages $m = 7$.
 • No. of spam $N_{\text{spam}} = 3$.
 • No. of non-spam $N_{\text{non-spam}} = 4$.
 Thus $P(y=0) = \frac{N_{\text{spam}}}{m} = \frac{3}{7} \approx 0.429$.
 $P(y=1) = \frac{N_{\text{non-spam}}}{m} = \frac{4}{7} \approx 0.571$.

2).

2. Given $m=7$ log-likelihood function:

$$l(\theta) = \sum_{i=1}^7 \sum_{k=1}^{15} x_k^{(i)} \log \theta_{y(i),k}$$
 We introduce Lagrangian multipliers.

$$\sum_{k=1}^d \theta_{c,k} = 1$$

$$\mathcal{L} = l(\theta) + \lambda_0 \left(1 - \sum_{k=1}^{15} \theta_{0,k}\right) + \lambda_1 \left(1 - \sum_{k=1}^{15} \theta_{1,k}\right)$$

① Compute Partial Derivatives & Set to 0

$$\theta_{0,k}: \frac{\partial \mathcal{L}}{\partial \theta_{0,k}} = \sum_{i: y(i)=0} \frac{x_k^{(i)}}{\theta_{0,k}} - \lambda_0 = 0 \quad \sum_{i: y(i)=0} x_k^{(i)} = \lambda_0 \theta_{0,k}$$

$$\theta_{1,k}: \frac{\partial \mathcal{L}}{\partial \theta_{1,k}} = \sum_{i: y(i)=1} \frac{x_k^{(i)}}{\theta_{1,k}} - \lambda_1 = 0 \quad \sum_{i: y(i)=1} x_k^{(i)} = \lambda_1 \theta_{1,k}$$

$$\theta_{0,k} = \frac{\sum_{i: y(i)=0} x_k^{(i)}}{\lambda_0} \quad \theta_{1,k} = \frac{\sum_{i: y(i)=1} x_k^{(i)}}{\lambda_1}$$

② Use constraints $= 1$

$$\lambda_0 = \frac{\sum_{k=1}^{15} \sum_{i: y(i)=0} x_k^{(i)}}{\lambda_0} = 1 \quad \lambda_0 = \sum_{k=1}^{15} \sum_{i: y(i)=0} x_k^{(i)}$$

$$\lambda_1 = \sum_{k=1}^{15} \sum_{i: y(i)=1} x_k^{(i)}$$

$$\theta_{0,1} = 0.3 \quad \theta_{1,1} \approx 0.059$$

$$\theta_{0,7} = 0.2 \quad \theta_{1,15} \approx 0.059$$

3).

3. From prior we know that

$$P(y=0) = \frac{3}{7} \quad P(y=1) = \frac{4}{7}$$

Spam Non-Spam

today is secret

$$X_{\text{test}} = \{1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0\}$$

$$\begin{bmatrix} \theta_{0,k} = 0.3 & \theta_{0,7} = 0.2 & \theta_{0,11} = 0.1 \\ \theta_{1,1} \approx 0.059 & \theta_{1,7} \approx 0.176 & \theta_{1,11} \approx 0.118 \end{bmatrix}$$

For $y=0$

$$P(X_{\text{test}}|y=0) \propto \theta_{0,1}^{x_1} \cdot \theta_{0,7}^{x_7} \cdot \theta_{0,11}^{x_{11}} = 0.3 \cdot 0.2 \cdot 0.1 = 0.006$$

For $y=1$

$$P(X_{\text{test}}|y=1) \propto \theta_{1,1}^{x_1} \cdot \theta_{1,7}^{x_7} \cdot \theta_{1,11}^{x_{11}} = 0.059 \cdot 0.176 \cdot 0.118 = 0.001226$$

$$P(X_{\text{test}}) = 0.006 \cdot \frac{3}{7} + 0.001226 \cdot \frac{4}{7} \approx 0.00257454$$

Posterior Probabilities:

For $y=0$

$$P(y=0|X_{\text{test}}) = \frac{0.006 \times 0.429}{0.00257454} \approx 0.999 > 0.272$$

For $y=1$

$$P(y=1|X_{\text{test}}) = \frac{0.001226 \times 0.571}{0.00257454} \approx 0.272$$

\therefore 'today is secret' is classified as spam because $y=0 > y=1$

theta_0_1 (secret|spam) = 0.273
 theta_0_7 (today|spam) = 0.182
 theta_1_1 (secret|non-spam) = 0.059
 theta_1_15 (pizza|non-spam) = 0.059

Question 4. Comparing classifiers: Divorce classification/prediction

1).

	Model Name	Training Accuracy	Testing Accuracy
0	Naive Bayes	98.529412	94.117647
1	Logistic Regressor	100.000000	91.176471
2	KNN	98.529412	94.117647

The best classifiers are both KNN and Naive Bayes. Because of its ease of use and reliance on the concept of conditional independence, the Naive Bayes classifier is beneficial and appears to perform well on this dataset. By utilizing the local data structure, KNN also exhibits remarkably high performance. Even though logistic regression achieves flawless training accuracy, overfitting causes testing accuracy to

decline. Consequently, Naive Bayes and KNN are the best models; their fair treatment of training and testing accuracy makes them appropriate for this specific classification job.

2).

Naive Bayes (PCA) training accuracy: 98.53%

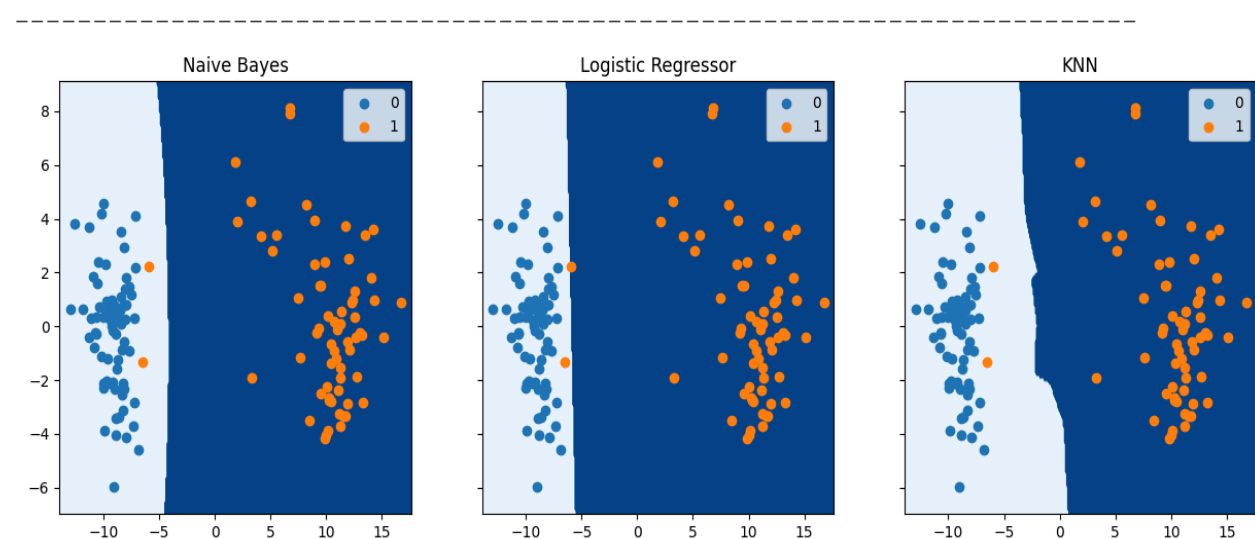
Naive Bayes (PCA) testing accuracy: 94.12%

Logistic Regressor (PCA) training accuracy: 99.26%

Logistic Regressor (PCA) testing accuracy: 100.00%

KNN (PCA) training accuracy: 98.53%

KNN (PCA) testing accuracy: 94.12%



For this dataset, the Logistic Regressor performs best; this is probably because the PCA modification made the space where logistic regression works best linearly separable. Naive Bayes and KNN also exhibit strong performance, when reduced to two dimensions using PCA. Logistic Regression shows perfect testing accuracy, making it the best model for this dataset.