

Problem 1.

$$x_n \sim \mathcal{N}(x_n | w z_n, \sigma^2 \mathbf{I}), \quad w \sim \left(\frac{\lambda}{2\pi}\right)^{\frac{dk}{2}} \exp\left\{-\frac{\lambda}{2} \text{trace}(w^T w)\right\}, \quad z_n \sim \mathcal{N}(z_n | 0, \mathbf{I})$$

Problem: Solve $w^* = \arg\max_w \ln p(x, w)$ We start from w_0 for $t=1, \dots, T$

$$\ln p(x, w) = \underbrace{\int q_t(z) \ln \frac{p(x, w, z)}{q_t(z)} dz}_{\mathcal{L}} + \underbrace{\int q_t(z) \ln \frac{q_t(z)}{p(z|w, x)} dz}_{KL(q_t(z) \| p(z|w, x))}$$

① In E-step

we assume $q_t(z) = p(z|w, x)$, which leads to $KL=0$

By Bayes rule, and the conjugate characteristic of Gaussian distribution

$$p(z|w, x) \propto \underbrace{p(z)}_{\mathcal{N}} \cdot \underbrace{p(x|z, w)}_{\mathcal{N}} \quad \begin{matrix} \nearrow \mathcal{N}(0, \mathbf{I}) \text{ and } \mathcal{N}(w z_n, \sigma^2 \mathbf{I}) \text{ respectively} \\ \searrow \end{matrix}$$

$$\Rightarrow p(z|w, x) = \mathcal{N}(z | \mu, \Sigma) \quad \text{where} \quad \begin{cases} \mu_n = \sum W^T (\sigma^2 \mathbf{I})^{-1} x_n = \sigma^{-2} \sum_n W^T x_n \\ \Sigma_n = (\mathbf{I} + W^T (\sigma^2 \mathbf{I})^{-1} W)^{-1} = (\mathbf{I} + \sigma^{-2} W^T W)^{-1} \end{cases}$$

② In M-step

we are solving $w_t = \arg\max_w \mathbb{E}_{q_t(z)} \left[\ln \frac{p(x, w, z)}{q_t(z)} \right]$, maximizing the \mathcal{L}

$$w_t = \arg\max_w \mathbb{E}_{q_t(z)} \left[\frac{\ln p(z, w, x)}{q_t(z)} \right] = \arg\max_w \left(\mathbb{E}_{q_t(z)} [\ln p(z, w, x)] - \mathbb{E}_{q_t(z)} [\ln q_t(z)] \right)$$

$$\ln p(z, w, x) = \ln (p(x|w, z) \cdot p(w) \cdot p(z))$$

$$= \sum_{n=1}^N \ln \mathcal{N}(x_n | w z_n, \sigma^2 \mathbf{I}) + \ln p(w) + \sum_{n=1}^N \ln \mathcal{N}(z_n | 0, \mathbf{I})$$

$$w_t = \arg\max_w \left(\ln p(w) + \sum_{n=1}^N \mathbb{E}_{q_t(z)} [\ln p(x_n | w, z_n)] \right)$$

$$\ln p(w) = \frac{dk}{2} \ln \left(\frac{\lambda}{2\pi} \right) - \frac{\lambda}{2} \text{trace}(w^T w)$$

$$\ln p(x_n | w, z_n) = C - \frac{1}{2\sigma^2} (x_n - w z_n)^T (x_n - w z_n) = C - \frac{1}{2\sigma^2} (x_n^T - z_n^T w^T) (x_n - w z_n)$$

nothing to do with w
we can ignore this term

we only consider the terms related to W

$$W_t = \underset{W}{\operatorname{argmax}} \left(-\frac{\lambda}{2} \operatorname{trace}(W^T W) + \sum_{n=1}^N \mathbb{E}_{q_t(z)} \left[-\frac{1}{2\sigma^2} (-z_n^T W^T X_n - X_n^T W z_n + z_n^T W^T W z_n) \right] \right)$$

To maximize \mathcal{L} , we set W_t s.t. $\frac{\partial \mathcal{L}}{\partial W} = 0$

$$\Rightarrow \frac{\partial \mathcal{L}}{\partial W} = -\lambda W - \frac{1}{2\sigma^2} \sum_{n=1}^N \mathbb{E}_{q_t(z)} \left[-\underbrace{z_n X_n^T - X_n^T z_n}_{\text{same}} + 2W z_n z_n^T \right] = 0$$

$$\Rightarrow \lambda W = \frac{1}{2\sigma^2} \sum_{n=1}^N \mathbb{E}_{q_t(z)} \left[2X_n^T z_n \right] - \frac{1}{2\sigma^2} \sum_{n=1}^N \mathbb{E}_{q_t(z)} \left[2W z_n z_n^T \right]$$

$$\Rightarrow W \left(\lambda + \sigma^{-2} \sum_{n=1}^N \mathbb{E}_{q_t(z)} [z_n z_n^T] \right) = \sigma^{-2} \sum_{n=1}^N \mathbb{E}_{q_t(z)} [X_n^T z_n]$$

$$\Rightarrow W = \left(\sigma^{-2} \sum_{n=1}^N X_n \mathbb{E}_{q_t(z)} [z_n^T] \right) \left(\lambda + \sigma^{-2} \sum_{n=1}^N \mathbb{E}_{q_t(z)} [z_n z_n^T] \right)^{-1}$$

Pseudo code:

1. Initialize W_0 to a vector of all zeros
2. For iteration $t=1, \dots, T$

Since $q_t(z) = \mathcal{N}(z_n | \mu_n', \Sigma_n')$

E-step:

we update $q_t(z)$ by
$$\begin{cases} \mu_{n(t)}' = \sigma^{-2} \sum_{n(t)} W_{t-1}^T X_n & \dots \dots \dots (E_{q_{t(n)}}[z_n^T]) \\ \Sigma_{n(t)}' = (I + \sigma^{-2} W_{t-1}^T W_{t-1})^{-1} & \dots \dots \dots (E_{q_{t(n)}}[z_n z_n^T]) \end{cases}$$

M-step:

$$W_t = \left(\sigma^{-2} \sum_{n=1}^N X_n \mathbb{E}_{q_t(z)} [z_n^T] \right) \left(\lambda + \sigma^{-2} \sum_{n=1}^N \mathbb{E}_{q_t(z)} [z_n z_n^T] \right)^{-1} \quad \text{for next iteration}$$

Calculate $\ln p(X, W) = \ln p(W) + \sum_{i=1}^N \ln p(x_i)$

$$= \frac{dK}{2} \ln \left(\frac{\lambda}{2\pi} \right) - \frac{\lambda}{2} \operatorname{trace}(W^T W) + \sum_{n=1}^N \left[-\frac{n}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} (X_n - W z_n)^T (X_n - W z_n) \right]$$

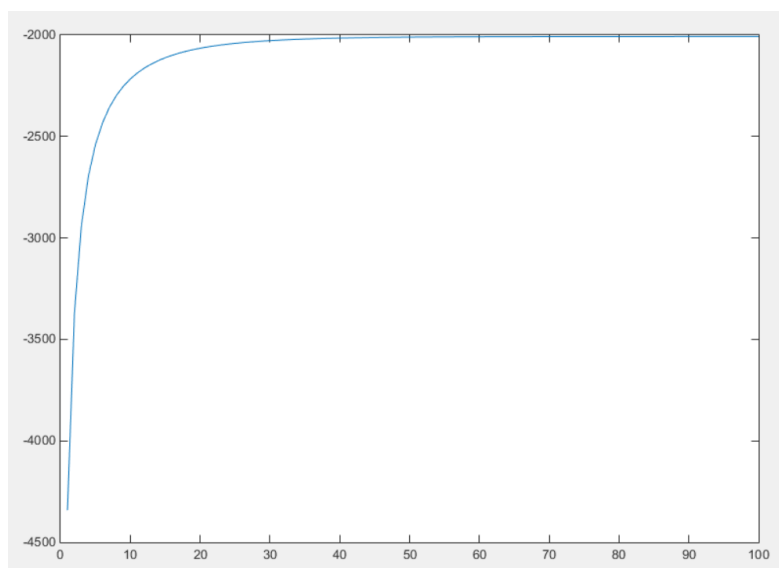
until the criterion is satisfied.

Problem 2

a) Attached please find the code of implementation.

Accuracy after 100 iterations: 93.57% (1863/1991)

b) Figure:



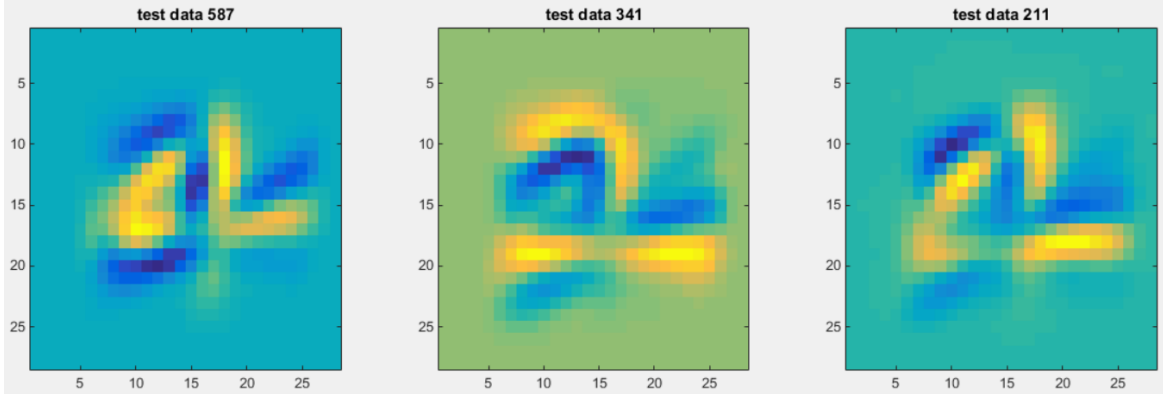
c) Confusion matrix

	Classified as 0 (4)	Classified as 1 (9)
Label 0 (4)	931	51
Label 1 (9)	77	932

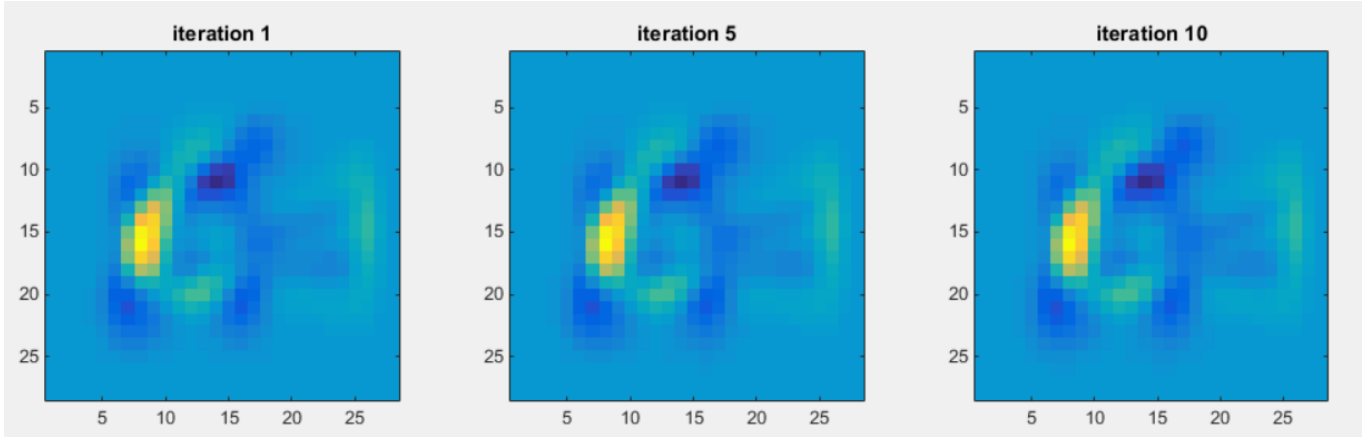
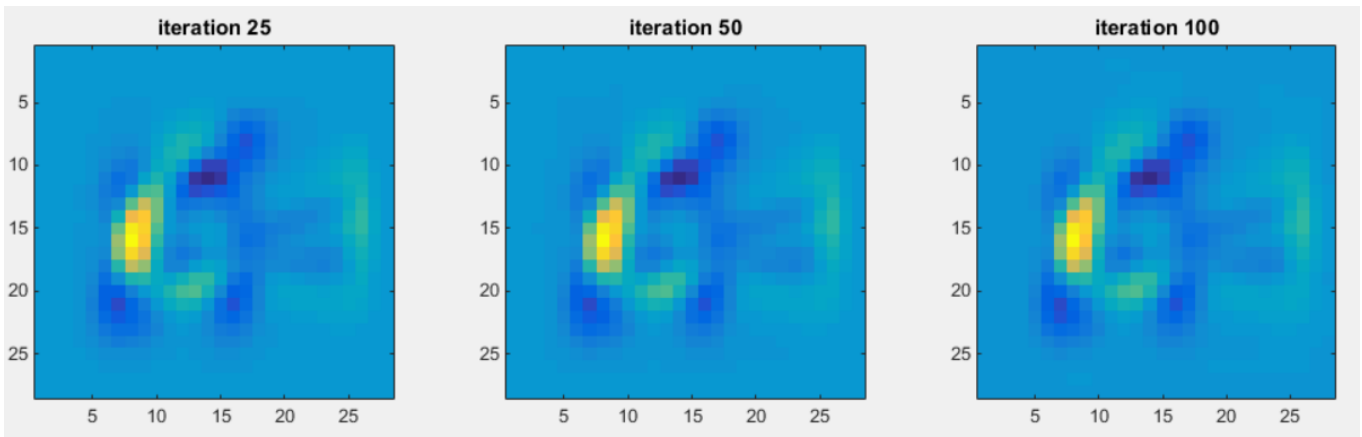
d) Three misclassified digits: (labeled as 4 but misclassified as 9)

Test data index	41	47	65
Images			
Predictive prob.	0.6775	0.6994	0.9062

e) Three most ambiguous digits:

Test data index	587	341	211
Images			
Predictive prob.	0.4995	0.5038	0.5038

f) Vector w_t for different t

Iteration 1	Iteration 5	Iteration 10
		
Iteration 25	Iteration 50	Iteration 100
		

Through the figure in b), it can be shown that the vector w_t is converging after 20~30 iterations, thus the last three figures above are almost the same (both visually and numerically). Compare w_1 to w_{100} , although they are visually the same, there are slight difference on vector value, and the color depth are also distinct. In conclusion, using EM algorithm for optimizing gives an acceptable approximate solution in a few iterations.