

Bayesian Mod for Machine Learning

Haoyang Chen hc2812 HW2

Prob 1. Denote $X = \{x_1, \dots, x_N\}$, $Z = \{z_1, \dots, z_N\}$

a) For E-step $q_t(z) = p(z|X, W_{t-1}) = \prod_{i=1}^N p(z_i|x_i, W_{t-1})$

$$p(z_i|x_i, W_{t-1}) \propto p(x_i|z_i, W_{t-1}) \cdot p(z_i) = \mathcal{N}(W_{t-1}z_i, \sigma^2 I) \cdot \mathcal{N}(0, I) = \mathcal{N}(z_i|u_i, \Sigma)$$

where $\Sigma = (I + \frac{1}{\sigma^2} W_{t-1}^T W_{t-1})^{-1}$, $u_i = \frac{\Sigma \cdot W_{t-1}^T x_i}{\sigma^2}$

$$L_t(W) = \int q_t(z) \ln p(X, W, z) dz - \int q_t(z) \ln q_t(z) dz$$

$$= E_{q_t(z)} [\ln p(X, W, z)] + \text{const.} \quad (\text{since } \int q_t(z) \ln q_t(z) dz \text{ does not contain } W, \text{ it's denoted as const.})$$

$$= E_{q_t(z)} [\ln p(W) + \sum_{i=1}^N \ln p(z_i) + \sum_{i=1}^N \ln p(x_i|z_i, W)] + \text{const.}$$

$$= \ln p(W) + \sum_{i=1}^N E_{q_t(z)} [\ln p(x_i|z_i, W)] + \text{const.} \quad (\text{Since } E_{q_t(z)} [\sum_{i=1}^N \ln p(z_i)] \text{ free from } W, \text{ it's}$$

$$= -\frac{\lambda}{2} \text{tr}(W^T W) + \sum_{i=1}^N E_{q_t(z)} \left[-\frac{1}{2\sigma^2} (x_i - Wz_i)^T (x_i - Wz_i) \right] + \text{const.} \quad \text{added into const. term)$$

$$= -\frac{\lambda}{2} \text{tr}(W^T W) + \sum_{i=1}^N \left[-\frac{1}{2\sigma^2} \{ -2\text{tr}(u_i x_i^T W) + \text{tr}(W^T W (u_i u_i^T + \Sigma)) \} \right] + \text{const.} \quad (\text{since}$$

$$\left(\begin{aligned} &\text{Since } E_{q_t(z)} [z_i] = u_i \Rightarrow E_{q_t(z)} [z_i z_i^T] = u_i u_i^T + \Sigma \\ &\Rightarrow E_{q_t(z)} [z_i^T W^T W z_i] = E_{q_t(z)} [\text{tr}(z_i^T W^T W z_i)] = \text{tr}(W^T W \cdot E_{q_t(z)} [z_i z_i^T]) \end{aligned} \right)$$

b) For M-step $W_t = \arg \max_W L_t(W)$

$$\frac{\partial L_t(W)}{\partial W} = \frac{1}{\sigma^2} \sum_{i=1}^N x_i u_i^T - W \left\{ \frac{1}{\sigma^2} \sum_{i=1}^N (u_i u_i^T + \Sigma) + \lambda I \right\} = 0$$

$$\therefore W_t = \left(\sum_{i=1}^N x_i u_i^T \right) \left(\sum_{i=1}^N u_i u_i^T + N \Sigma + \lambda \sigma^2 I \right)^{-1}$$

c) $\ln p(X, W_t) = \int q_t(z) \ln \frac{p(X, W_t, z)}{q_t(z)} dz + \int q_t(z) \ln \frac{q_t(z)}{p(z|X, W_t)} dz = E_{q_t(z)} [\ln p(X, z, W_t) - \ln p(z|X, W_t)]$

$$= \ln p(W_t) + \sum_{i=1}^N E_{q_t(z)} [\ln p(x_i|z_i, W_t) + \ln p(z_i) - \ln q_t(z_i)]$$

$$= -\frac{\lambda}{2} \text{tr}(W_t^T W_t) + \frac{dk}{2} \ln \frac{\lambda}{2\pi} + \sum_{i=1}^N \left[-\frac{d}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} x_i^T x_i + \frac{1}{\sigma^2} \text{tr}(u_i x_i^T W_t) - \frac{1}{2\sigma^2} \text{tr}(W_t^T W_t (u_i u_i^T + \Sigma)) \right]$$

$$+ \sum_{i=1}^N \left[-\frac{1}{2} \ln(2\pi) - \frac{1}{2} \text{tr}(u_i u_i^T + \Sigma) \right] + \sum_{i=1}^N \left[\frac{1}{2} \ln \det(2\pi\Sigma) \right]$$

The EM algorithm pseudo-code:

1. Initialize W_0 as a $d \times k$ matrix with zeros

2. For $t = 1, \dots, T$, do:

(i) E-Step: Calculate $\Sigma = (-I + \frac{1}{\sigma^2} W_{t-1}^T W_{t-1})^{-1}$, $\mu_i = \frac{\Sigma \cdot W_{t-1}^T x_i}{\sigma^2}$, $i = 1, \dots, N$

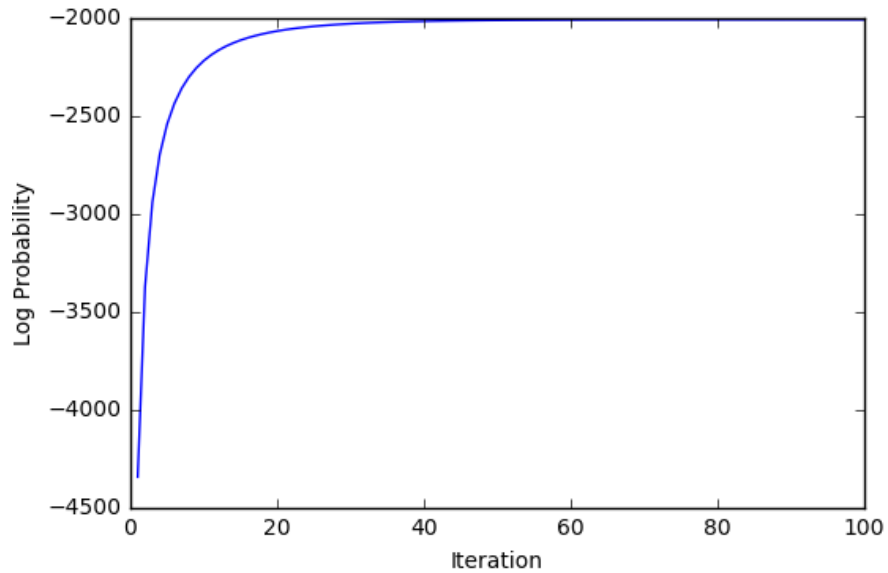
(ii) M-Step: Calculate $W_t = (\sum_{i=1}^N x_i \mu_i^T) (\sum_{i=1}^N \mu_i \mu_i^T + N \cdot \Sigma + \lambda \sigma^2 I)^{-1}$

(iii) Calculate $\ln \mathcal{P}(X, W_t)$, if $\ln \mathcal{P}(X, W_t) < \ln \mathcal{P}(X, W_{t-1})$, stop.

2.

a). Shown in the code

b). Draw the plot:



c). The confusion matrix is shown below:

	Actual = 0	Actual = 1
Predicted = 0	930	52
Predicted = 1	77	932

Accuracy = 0.9352

d). 3 misclassified images

Image Index: 46

True label: 0

Predicted label: 1

Prob: 0.7819897

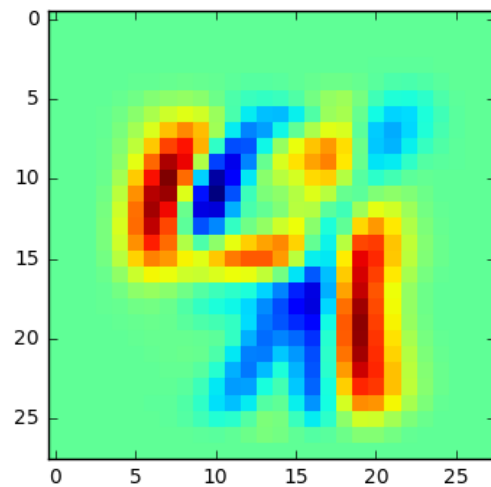


Image Index: 156

True label: 0

Predicted label: 1

Prob: 0.9310753

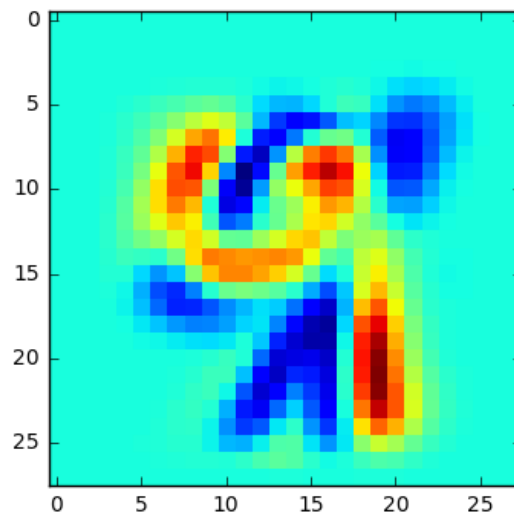
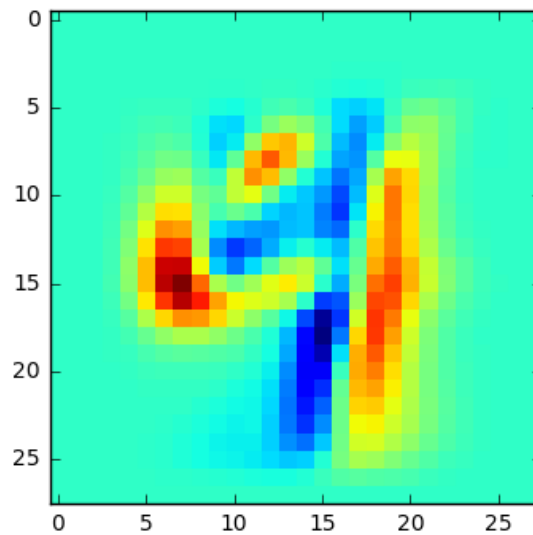


Image Index: 259

True label: 0

Predicted label: 1

Prob: 0.8018295



e). 3 most ambiguous predictions

Image Index: 586

True label: 0

Predicted label: 1

Prob: 0.5002605

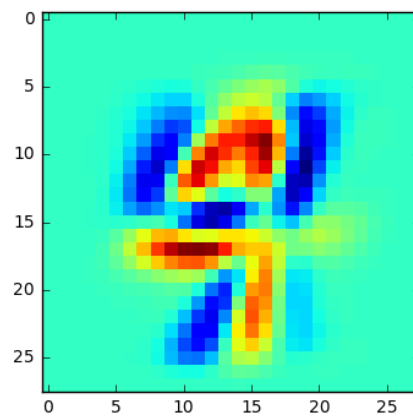


Image Index: 340

True label: 0

Predicted label: 1

Prob: 0.5046497

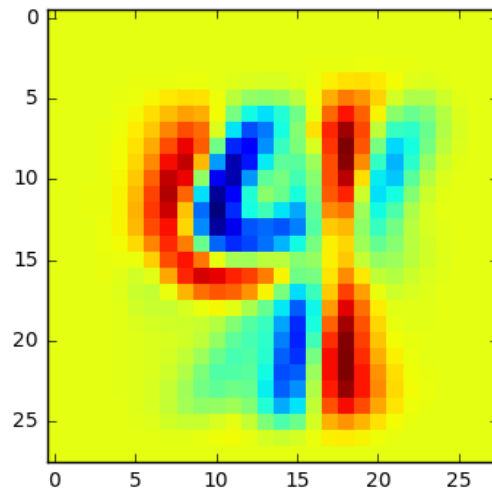
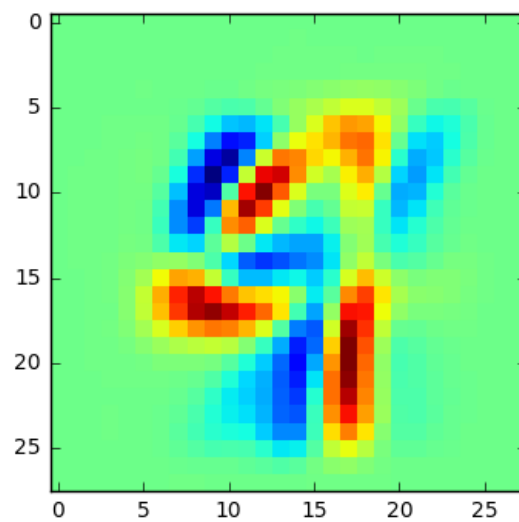


Image Index: 210

True label: 0

Predicted label: 1

Prob: 0.5061759



f). Reconstruct W

The weights images look like the number images. The difference between images becomes smaller when t gets larger. It means when t become larger, the convergence rate of weights would be slower. The plot in problem b also shows this fact.

