

Question 1

a) It can be shown that the gaussian parameters μ_k, Σ_k that maximize the likelihood function for a *given* class k is the mean and sample variance of the observed points of the class. As per the specifications of the question statement, we assume a prior for each class that is uniform. Ordinarily, we might instead assume a multinomial distribution and take the prior as the number of occurrences of a class divided by total number of occurrences of all classes.

Average log-likelihood for training set:

-0.124624436669

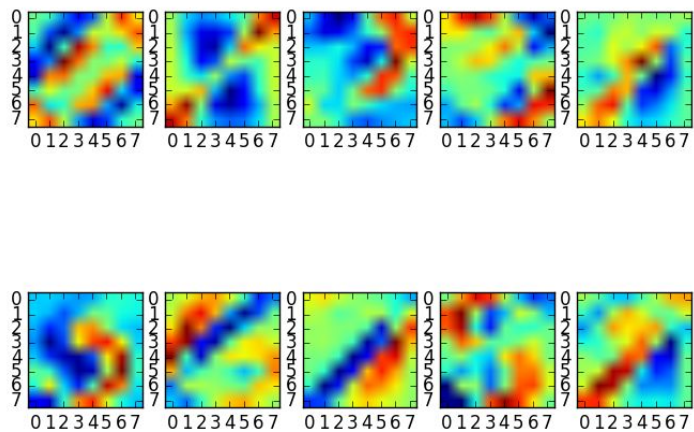
Average log likelihood for test:

-0.196673203255

b) Accuracy on training set: 0.981429

Accuracy on test set: 0.972750

c)



Question 2

$$2a) P(\theta|D) = \frac{P(D|\theta)P(\theta)}{P(D)} \sim P(D|\theta)P(\theta)$$

$$\text{As given, } P(\theta) = \theta_1^{a_1-1} \dots \theta_K^{a_K-1} = \prod_{i=1}^K \theta_i^{a_i-1}$$

$$P(x_1, \dots, x_N | \theta) = P(x_1^N | \theta) \dots P(x_K^N | \theta) \text{ assuming independence}$$

$$= \left(\prod_{n=1}^N \theta_1^{x_{1n}} \right) \dots \left(\prod_{n=1}^N \theta_K^{x_{Kn}} \right)$$

$$= \prod_{n=1}^N \prod_{k=1}^K \theta_k^{x_{kn}} = \theta_1^{N_1} \dots \theta_K^{N_K}$$

$$\therefore P(\theta|D) = \left(\prod_{n=1}^N \prod_{k=1}^K \theta_k^{x_{kn}} \right) \left(\prod_{i=1}^K \theta_i^{a_i-1} \right) = P(\theta_1, \dots, \theta_K | D)$$

$$P(D'|D) = \int P(\theta|D) P(D'|\theta) d\theta$$

$$P(D'|\theta) = P(x_{b+1}^{N+1} = 1 | \theta) = \theta_b$$

$\rightarrow \theta_i \in (0, 1)$ But

$$\sum_{i=1}^K \theta_i = 1$$

How to do $\int d\theta$?

$$\therefore P(D'|D) = \int P(\theta|D) \theta_b d\theta =$$

$$\int \dots \int \left(\prod_{n=1}^N \prod_{k=1}^K \theta_k^{x_{kn}} \right) \left(\prod_{i=1}^K \theta_i^{a_i-1} \right) \theta_b d\theta_1 \dots d\theta_K$$

$$= \int \dots \int \theta_1^{N_1} \dots \theta_K^{N_K} \cdot \theta_1^{a_1-1} \dots \theta_K^{a_K-1} \cdot \theta_b d\theta_1 \dots d\theta_K$$

$$= \int \dots \int \theta_1^{N_1+a_1-1} \dots \theta_K^{N_K+a_K-1} \theta_b d\theta_1 \dots d\theta_K$$

$$\underbrace{\text{Dirichlet} \sim (N_1+a_1, \dots, N_K+a_K)}_{\text{Dirichlet}}$$

$$E[\theta_b] \text{ given } (\theta_1, \dots, \theta_K) \sim \text{Dirichlet}(N_1+a_1, \dots, N_K+a_K)$$

$$P(x_{b+1}^{N+1} = 1 | \theta) = P(D'|D) = E[\theta_b] = \frac{N_b + a_b}{\sum_{i=1}^K (N_i + a_i)}$$

$$\text{where } N_i = \sum_{n=1}^N x_i^{(n)}$$

$$2b) \quad P(\theta|D) \sim P(D|\theta) P(\theta) \sim \log P(D|\theta) + \log P(\theta)$$

$$\text{Recall: } P(D|\theta) = \theta_1^{N_1} \dots \theta_k^{N_k}$$

$$P(\theta) = \theta_1^{a_1-1} \dots \theta_k^{a_k-1}$$

$$\therefore P(\theta|D) \sim N_1 \log \theta_1 + \dots + N_k \log \theta_k + (a_1-1) \log \theta_1 + \dots + (a_k-1) \log \theta_k$$

$$= (N_1 + a_1 - 1) \log \theta_1 + \dots + (N_k + a_k - 1) \log \theta_k$$

$$\frac{dP(\theta|D)}{d\theta_i} = \frac{N_i + a_i - 1}{\theta_i} = 0 \quad \times. \text{ Need Lagrange again.}$$

$$\text{Maximize } P(\theta|D) \quad := f$$

$$\text{subject to } 1 - (\theta_1 + \dots + \theta_k) = 0 \quad := g$$

$$\text{Lagrangian Method: } \nabla_{\theta} f = \lambda \nabla_{\theta} g \quad - K \text{ equations}$$

$$g = 0 \quad 1 \text{ equation}$$

$$\theta_1, \dots, \theta_k, \lambda \rightarrow K+1 \text{ unknowns}$$

$$\nabla_{\theta} f = \lambda \nabla_{\theta} g \Rightarrow \frac{N_i + a_i - 1}{\theta_i} = \lambda (-1)$$

$$\theta_i = \frac{1 - a_i - N_i}{\lambda}$$

$$\therefore 1 - \sum_{i=1}^K \theta_i = 0$$

$$1 - \sum_{i=1}^K \frac{1 - a_i - N_i}{\lambda} = 0$$

$$1 = \frac{1}{\lambda} (\sum 1 - \sum a_i - \sum N_i)$$

$$\lambda = K - \sum a_i - N$$

$$\text{Note } \sum 1 = K$$

$$\theta_i = \frac{1 - a_i - N_i}{K - (\sum a_i) - N} = \frac{N_i + a_i - 1}{N + \sum a_i - K}$$

Question 3

3. a) $Z \sim \mathcal{N}(0, I)$, $x|Z \sim \mathcal{N}(Z\mu, \Sigma)$, $\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_p^2)$

Find $P(z|x;\theta)$

From the appendix: $\mathbb{P}(Z) = \mathcal{N}(Z | \mu, \Sigma^{-1})$, $P(x|z) \sim \mathcal{N}(x | Az + b$

Then $p(z|x) \sim \mathcal{N}(z | C(A^T L(x-b) + \lambda \mu), C)$

$$C = (N + A^T L A)^{-1}$$

$$A = u \quad b = 0 \quad \Lambda = \mathbb{I} = \Lambda^{-1} \quad \underline{f} = \underline{z}, \mu = 0$$

$$C = (1 + \mu^T \Sigma^{-1} \mu)^{-1}$$

$$C(A^T L(x-b) + \lambda \mu) = (1 + \mu^T \Sigma' \mu)^{-1} (\mu^T \Sigma' (x) + 0) \\ = (1 + \mu^T \Sigma' \mu)^{-1} (\mu^T \Sigma' x)$$

$$\therefore P(z|x; \theta) \sim \mathcal{N}\left(z \mid (1 + \mu^T \Sigma^{-1} \mu)^{-1} \mu^T \Sigma^{-1} x, (1 + \mu^T \Sigma^{-1} \mu)^{-1}\right)$$

$$m = E[z|x] = (1 + \mu^T \Sigma^{-1} \mu)^{-1} \mu^T \Sigma^{-1} x$$

$$\text{scalar} - S = E[z^2|x] = \text{Var}(z|x) + (E[z|x])^2$$

$$= (1 + \mu^T \Sigma^{-1} \mu)^{-1} + [(1 + \mu^T \Sigma^{-1} \mu)^{-1} \mu^T \Sigma^{-1} x]^2$$

$$x \in \mathbb{R}^D, z \in \mathbb{R}^1, \mu \in \mathbb{R}^{D \times 1}$$

$$(1 + \mu^T \Sigma^{-1} \mu)^{-1} \mu^T \Sigma^{-1} x \quad \text{univariate}$$

$$q = P(z|x) \sim \mathcal{N}\left(\underbrace{\mu_1}_{\mu}, \underbrace{\Sigma_1}_{\Sigma}\right)$$

$$p(x|z) \sim \mathcal{N}(z\mu, \Sigma)$$

$$p(z) \sim \mathcal{N}(0, I)$$

b) $\log p(z^n, x^n | \theta) = \log p(x^n | z^n, \theta) + \log p(z^n | \theta) = \log p(x^n | z^n, \theta) + \log p(z^n | \theta)$
 $p(x|z; \theta) \sim \mathcal{N}(z\mu, \Sigma)$ $p(z) \sim \mathcal{N}(0, I)$

$$\textcircled{1} = \mathbb{E}_{p(z^n | x^n, \mu, \Sigma)} [\log p(x^n | z^n; \mu, \Sigma) + \log p(z^n)]$$

$$= \int_{-\infty}^{\infty} dz^n \frac{\exp\left(-\frac{1}{2}(z^n - \mu)^T \Sigma^{-1} (z^n - \mu)\right)}{\sqrt{(2\pi)^D |\Sigma|}} \left(\underbrace{\log p(x^n | z^n; \mu, \Sigma)}_{\text{old}} + \underbrace{\log p(z^n)}_{\text{old}} \right)$$

$$= \int_{-\infty}^{\infty} dz^n \frac{\exp\left(-\frac{1}{2}(z^n - \mu)^T \Sigma^{-1} (z^n - \mu)\right)}{\sqrt{(2\pi)^D |\Sigma|}} \left(-\frac{1}{2} (x^n - z^n \mu)^T \Sigma^{-1} (x^n - z^n \mu) - \frac{1}{2} \log((2\pi)^D |\Sigma|) - \frac{1}{2} (z^n)^T z^n - \frac{1}{2} \log(2\pi) \right)$$

Recall:
 $\int_{-\infty}^{\infty} f(x) dx = 1$
 if $f(x)$ is PDF

$$\begin{aligned} & ((x^n)^T - \mu^T (z^n)^T) \Sigma^{-1} (x^n - z^n \mu) \\ &= (x^n)^T \Sigma^{-1} x^n - (x^n)^T \Sigma^{-1} z^n \mu - \mu^T (z^n)^T \Sigma^{-1} x^n + \mu^T (z^n)^T \Sigma^{-1} z^n \mu \\ &= (x^n)^T \Sigma^{-1} x^n - z^n (x^n)^T \Sigma^{-1} \mu - z^n \mu^T \Sigma^{-1} x^n + (z^n)^T \mu^T \Sigma^{-1} \mu \end{aligned}$$

$$\textcircled{1} = (x^n)^T \Sigma^{-1} x^n - (x^n)^T \Sigma^{-1} \mu s^n - \mu^T \Sigma^{-1} x^n s^n + \mu^T \Sigma^{-1} \mu s^n - \frac{1}{2} \log((2\pi)^D |\Sigma|) - \frac{1}{2} s^n - \frac{1}{2} \log(2\pi)$$

$$\mathcal{L}(q, \theta) = \frac{1}{N} \sum_{n=1}^N \textcircled{1}$$

$$\begin{aligned} \frac{d\mathcal{L}}{d\mu} &= \frac{1}{N} \sum_{n=1}^N - (x^n)^T \Sigma^{-1} s^n - (\Sigma^{-1} x^n)^T s^n + 2 s^n \mu^T \Sigma^{-1} \\ &= \frac{1}{N} \sum_{n=1}^N - (x^n)^T \Sigma^{-1} s^n - (x^n)^T \Sigma^{-1} s^n + 2 s^n \mu^T \Sigma^{-1} \quad \text{Note: } (\Sigma^{-1})^T = \Sigma^{-1} \\ 0 &= \frac{1}{N} \sum_{n=1}^N - 2 (x^n)^T \Sigma^{-1} s^n + 2 s^n \mu^T \Sigma^{-1} \quad \Sigma^T = \Sigma \end{aligned}$$

$$\begin{aligned} \Sigma (x^n)^T \Sigma^{-1} s^n &= \Sigma s^n \mu^T \Sigma^{-1} \\ \Sigma s^n \Sigma^{-1} x^n &= \Sigma s^n \Sigma^{-1} \mu \quad \text{mult. both sides by } \Sigma \\ \Sigma s^n x^n &= \Sigma s^n \mu \end{aligned}$$

$$\boxed{\mu = \frac{\Sigma s^n x^n}{\Sigma s^n}}$$

Note: it's a sum at the end, not the covariance. I dropped the bounds for convenience.

Proof of some derivatives used above:

What is $\frac{d}{d\mu} (\mu^T \overset{1 \times D}{\Sigma^{-1}} \overset{D \times D}{x^n} \overset{D \times 1}{m^n})$?
 $\frac{d}{d\mu} (\mu^T b)$, $b \in D \times 1$

$$\begin{bmatrix} \mu_1 & \dots & \mu_D \end{bmatrix} \begin{bmatrix} b_1 \\ \vdots \\ b_D \end{bmatrix} = \mu_1 b_1 + \dots + \mu_D b_D$$

$$\frac{d}{d\mu} (\mu_1 b_1 + \dots + \mu_D b_D) = \begin{bmatrix} b_1 & \dots & b_D \end{bmatrix} = b^T$$

Previously we showed: $\frac{d(B^T A B)}{dB} = 2B^T A$ if A is diag matrix

What is $\frac{d}{d\mu} ((x^n)^T \overset{D \times D}{\Sigma^{-1}} \mu)$?

$$\frac{d}{d\mu} (b^T \mu)$$

$$\rightarrow = \begin{bmatrix} b_1 & \dots & b_D \end{bmatrix} \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_D \end{bmatrix} = \begin{bmatrix} b_1 \mu_1 & \dots & b_D \mu_D \end{bmatrix}$$

$$\therefore \frac{d}{d\mu} (b^T \mu) = b^T$$