

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:

28.6865797126

Average log likelihood for test:

27.4861999423

I also report the log-likelihood averaged over each class for each class. This is a more useful metric because it shows that the classifier is assigning a higher probability to the correct class relative to the others. As one might expect, the diagonals have the highest number, i.e. the highest probability, so the classifier is working correctly.

Average log-likelihood for training set (ith row is average over all points with ith label):

```
[[ 33.92402221 -88.45111834 -16.89326222 -25.83617441 -95.99459913 -28.0469699 -33.43930375 -128.00676019 -40.12341711 -120.7505451 ]
 [-68.22979333 43.8330528 -8.62313018 -41.79892253 -6.16921495 -42.54976231 -22.89991813 -33.06720503 4.20662396 -19.78080586]
 [-48.25937897 -55.92309386 16.88110064 -28.9959108 -65.50319452 -47.8188899 -54.71585182 -84.2803214 -35.97571408 -79.85123569]
 [-37.76535226 -60.50596233 -11.52881147 25.34749836 -80.23781057 -5.942891 -61.89235394 -74.20703712 -7.77750941 -63.81771045]
 [-63.57874438 -38.1128814 -27.37434988 -43.54955558 24.20650818 -37.62118033 -38.11452022 -49.30644691 -19.02665868 -3.93941094]
 [-31.92250461 -45.89294551 -23.0214228 0.2051562 -56.8635482 24.12660812 -36.96891002 -84.4180825 -2.6016538 -53.5127435 ]
 [-16.08499354 -65.52399532 -17.22797111 -38.51566436 -59.26386047 -19.90746009 29.24659472 -165.40709904 -29.97466814 -126.96745943]
 [-47.34810481 -22.74497808 0.17304175 -13.15842334 -22.22716495 -39.09215712 -80.81237908 38.53051509 -7.68591751 9.77505843]
 [-51.89004965 -33.31890348 -15.94135661 -21.54837236 -45.81703868 -18.8929108 -47.39728619 -66.1336154 20.02956548 -37.4736449 ]
 [-65.1618034 -46.64166345 -24.55400204 -32.3243134 -3.74088864 -39.22988767 -85.57166051 -7.67267147 -5.99624188 30.74033154]]
```

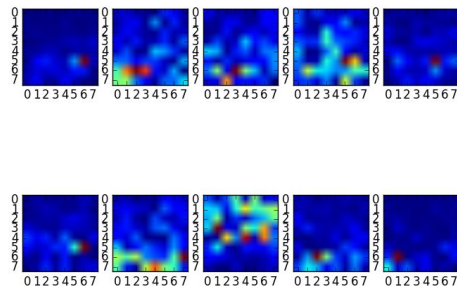
Average log likelihood for test:

```
[[ 33.21271943 -85.10474932 -14.92668303 -25.03919668 -92.24055178 -28.41301758 -31.98166934 -121.11671514 -38.9435702 -115.68602049]
 [-65.98273597 43.75667838 -7.89988668 -40.78305048 -4.21465107 -41.85439817 -19.49559598 -27.26307905 5.51581709 -17.43312634]
 [-48.40011136 -56.79648782 15.76987273 -32.02961238 -64.05056948 -48.58909001 -52.90835733 -87.94971095 -36.4708241 -81.13238684]
 [-36.61811894 -59.0998393 -10.81252291 23.18843659 -79.32152079 -5.49772638 -60.31175001 -75.53668363 -8.21771127 -64.68566461]
 [-63.6224916 -38.31251625 -27.68532145 -42.79967361 23.16501442 -36.98680559 -38.8213785 -48.45629329 -19.54241646 -2.7608024 ]
 [-28.61602964 -40.69825848 -20.36105698 3.672266 -55.92112179 24.03785909 -32.49333123 -81.85722903 -1.63104965 -53.09672336]
 [-17.94182886 -70.13173546 -19.32409691 -40.17621616 -61.330969 -20.45038953 26.35234161 -169.74575806 -30.4058043 -128.24679529]
 [-48.50416378 -22.52592109 -0.48298009 -13.97483349 -20.82250406 -38.86921001 -81.1159767 37.69449875 -8.12088509 9.78420061]
 [-53.05155379 -34.54053281 -17.60860408 -22.40460407 -46.75337517 -19.61290705 -49.00099401 -68.0970773 18.6596498 -38.65401494]
 [-64.02644638 -46.67289562 -22.92041583 -30.60657673 -4.95359756 -38.41795983 -86.48303501 -6.37184554 -5.88093645 29.02492862]]
```

b) Accuracy on training set: 0.981429

Accuracy on test set: 0.972750

c)



Question 2

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

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

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

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

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

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

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

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

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

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

How to do $\int d\theta$?

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

Question 3

3. a) $Z \sim \mathcal{N}(0, 1)$, $x|Z \sim \mathcal{N}(Z\mu, \Sigma)$, $\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_p^2)$ scalar

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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