

1 a) Model 2: $S := S_1 = S_2$

$$g_i(x) = \frac{-d}{2} \log(2\pi) - \frac{1}{2} \log(|\Sigma|) - \frac{1}{2} (x - \mu_i)^T \Sigma^{-1} (x - \mu_i) + \log(P(c_i))$$

$$\begin{aligned} 2(\Sigma^{-1} | x) &\equiv \sum_{t=1}^N \log(g_i(x^t)) \equiv \sum_{t=1}^N g_i(x^t) \\ &= \sum_{t=1}^N \left[\frac{-d}{2} \log(2\pi) - \frac{1}{2} \log(|\Sigma|) - \frac{1}{2} (x^t - \mu_i)^T \Sigma^{-1} (x^t - \mu_i) + \log(P(c_i)) \right] \\ &\equiv \sum_{t=1}^N \left[-\frac{1}{2} \log(|\Sigma|) - \frac{1}{2} (x^t - \mu_i)^T \Sigma^{-1} (x^t - \mu_i) \right] \\ &= -\frac{N}{2} \log(|\Sigma|) - \frac{1}{2} \sum_{t=1}^N (x^t - \mu_i)^T \Sigma^{-1} (x^t - \mu_i) \end{aligned}$$

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \Sigma^{-1}} &= -\frac{N}{2} \frac{\partial (\log(|\Sigma|))}{\partial \Sigma^{-1}} - \frac{1}{2} \sum_{t=1}^N \frac{\partial [(x^t - \mu_i)^T \Sigma^{-1} (x^t - \mu_i)]}{\partial \Sigma^{-1}} \\ &= -\frac{N}{2} (-\Sigma^T) - \frac{1}{2} \sum_{t=1}^N (x^t - \mu_i)^T (x^t - \mu_i) = 0 \end{aligned}$$

$$\frac{N}{2} \Sigma = \frac{1}{2} \sum_{t=1}^N (x^t - \mu_i)^T (x^t - \mu_i)$$

$$S = \Sigma = \frac{\sum_{t=1}^N (x^t - \mu_i)^T (x^t - \mu_i)}{N}$$

1 a) Model 3: $S_1 = \alpha_1 I$ $S_2 = \alpha_2 I$

$$g_i(x) = \frac{-d}{2} \log(2\pi) - \frac{1}{2} \log(|\Sigma_i|) - \frac{1}{2} (x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i) + \log(P(c_i)).$$

$$-\frac{1}{2} \log(|\Sigma_i|) = -\frac{1}{2} \log(|\alpha_i I|) = -\frac{1}{2} \log(\alpha_i^d) = -\frac{d}{2} \log(\alpha_i)$$

$$\Sigma_i^{-1} = (\alpha_i I)^{-1} = (\alpha_i^{-1}) I$$

$$g_i(x) = \frac{-d}{2} \log(2\pi) - \frac{d}{2} \log(\alpha_i) - \frac{1}{2} \alpha_i^{-1} (x - \mu_i)^T (x - \mu_i) + \log(P(c_i))$$

$$2(\alpha_i | x) \equiv \sum_{t=1}^N \log(g_i(x^t))$$

$$\equiv -\frac{1}{2} \sum_{t=1}^N \log(d \log(\alpha_i) + \alpha_i^{-1} (x^t - \mu_i)^T (x^t - \mu_i))$$

$$\equiv -\frac{1}{2} \sum_{t=1}^N \left[d \log(\alpha_i) + \alpha_i^{-1} (x^t - \mu_i)^T (x^t - \mu_i) \right]$$

$$= -\frac{1}{2} \cdot N d \log(\alpha_i) + \alpha_i^{-1} \sum_{t=1}^N (x^t - \mu_i)^T (x^t - \mu_i)$$

$$\frac{\partial 2}{\partial \alpha_i} = -\frac{1}{2} \cdot N d \frac{\partial (\log(\alpha_i))}{\partial \alpha_i} + \frac{\partial (\alpha_i^{-1})}{\partial \alpha_i} \sum_{t=1}^N (x^t - \mu_i)^T (x - \mu_i)$$

$$= -\frac{Nd}{2\alpha_i} + \frac{-1}{\alpha_i^2} \cdot \sum_{t=1}^N (x^t - \mu_i)^T (x - \mu_i) = 0.$$

$$\alpha_i = \frac{2}{Nd} \sum_{t=1}^N (x^t - \mu_i)^T (x - \mu_i)$$

1c) Error Rates:

	#	error
<u>Model 1:</u>	test-data 1	— 0.22
	test-data 2	— 0.23
	test-data 3	— 0.11

<u>Model 2:</u>	test-data 1	— 0.17
	test-data 2	— 0.55
	test-data 3	— 0.45

<u>Model 3:</u>	test-data 1	— 0.34
	test-data 2	— 0.38
	test-data 3	— 0.07

data 3
model 3

Based on the results, it seems most likely that test-data 3 was distributed with S_1 , very different from S_2 , which is why models 1 and 3 were much more effective than model 2. It was also likely to have $S_1 \approx 2I$, which meant that Model 3, a less complex model, had less error than model 1.

data 1
model 2

It is also likely that test-data 1 had $S_1 \approx S_2$, which meant that Model 2 (less complex than model 1) had less error. It is also likely that there were many non-zero covariances, which would explain model 3's high error rate.

data 2
model 1

None of the models were particularly good at learning test-data 2, which means that $S_1 \neq S_2$ and $S_1 \neq 2I$, but it is possible that a different simplification would allow better fit than model 1.