

UML2015
Solutions and other comments
to some exercises from sets 4 and 6

Set 4 Exercise 1:

1.1

$$\mathbb{E}(\mathcal{Y}_1 \mathcal{Y}_2) \approx \frac{1}{K} \sum_{k=1}^K (\mathbf{y}_1)_k (\mathbf{y}_2)_k = \frac{1}{K} \mathbf{y}_1 \cdot \mathbf{y}_2 \quad (1)$$

Assume that A is some real random variable, and then assume that (A_1, A_2, \dots, A_K) is iid each A_k having the same distribution as A . If (a_1, a_2, \dots, a_K) is a sample draw, an approximation

$$\mathbb{E}(A) \approx \frac{1}{K} \sum_{k=1}^K a_k \quad (2)$$

holds. This is true for all random variables for which the mean exist, and in particular for $A = \mathcal{Y}_1 \mathcal{Y}_2$ too.

If \mathcal{Y}_1 and \mathcal{Y}_2 are uncorrelated, then $\mathbb{E}(\mathcal{Y}_1 \mathcal{Y}_2) = 0$. This will not imply that the vectors \mathbf{y}_1 and \mathbf{y}_2 would be orthogonal, but since it does imply $\frac{1}{K} \mathbf{y}_1 \cdot \mathbf{y}_2 \approx 0$, we can say that the vectors are approximately orthogonal in this sense. The point of the exercise is to understand that theoretical uncorrelatedness of \mathcal{Y}_1 and \mathcal{Y}_2 implies approximate orthogonality for \mathbf{y}_1 and \mathbf{y}_2 .

1.2

$$\frac{1}{K} \overline{X} \overline{X}^T = \frac{1}{K} (U^T X) (X^T U) = U^T \underbrace{\left(\frac{1}{K} X X^T \right)}_{=U\Lambda} U = \underbrace{U^T U}_{=\text{id}} \Lambda = \Lambda \quad (3)$$

where Λ is an $M \times M$ diagonal matrix containing the most significant eigenvalues of $\frac{1}{K} X X^T$ on its diagonal.

$$\implies \quad \overline{X}_{m*} \cdot \overline{X}_{m'*} = (\overline{X} \overline{X}^T)_{mm'} = 0 \quad \forall m \neq m' \quad (4)$$

Here we used the knowledge that eigenvectors (here U_{*m}) of a symmetric matrix (here $\frac{1}{K} X X^T$) can always be assumed to be orthogonal. Notice that we did not assume that U would have fully diagonalized the $\frac{1}{K} X X^T$ (which is $N \times N$). The calculation works under the assumption that some of the eigenvectors were placed into the rows of U .

A little mistake: Eigenvectors of a symmetric matrix are always orthogonal if the corresponding eigenvalues are distinct. The eigenvectors are not necessarily orthogonal, if some of the eigenvalues appear multiple times (some eigenvalues have degree greater than 1). The eigenvectors are never unique, and even if some

eigenvalues appear multiple times, the eigenvectors can be chosen so that they are orthogonal. In this exercise, it was not mentioned that the eigenvalues of $\frac{1}{K}XX^T$ would be distinct, and also it was not mentioned that the columns of U would be assumed orthogonal. This can be seen as a mistake in the exercise, and it should have been clarified that the columns are assumed to be orthogonal so that $U^TU = \text{id}$.

If X was a real data sample, it would be extremely unlikely that some of the eigenvalues of $\frac{1}{K}XX^T$ would be the same, so in this sense the mistake was very small.

Further comments to the explanation request in 1.1. If \mathbf{X} was a zero mean $N \times 1$ random vector, whose elements X_1, \dots, X_N were not uncorrelated, it would be nice, if we could find a matrix U with orthogonal columns such that a new random vector $\bar{\mathbf{X}} = U^T\mathbf{X}$ did have uncorrelated elements so that $\mathbb{E}(\bar{X}_m\bar{X}_{m'}) = 0$ for all $m \neq m'$. It is not possible to find this kind of U based on any $N \times K$ sample X due to the randomness of the sample points. However, we learned that it is possible to find U such that the rows of U^TX were precisely orthogonal. In doing this, we found U such that the elements of $\bar{\mathbf{X}} = U^T\mathbf{X}$ are approximately uncorrelated so that $\mathbb{E}(\bar{X}_m\bar{X}_{m'}) \approx 0$ for all $m \neq m'$. So the relationship of uncorrelatedness and orthogonality works in this direction too.

Set 4 Exercise 2:

2.1 This exercise deals with basics of linear algebra. We assume that A is a real $N \times M$ matrix and $M < N$, and that the columns of A are linearly independent. The given task is to prove *that AA^T has M positive eigenvalues, and zero appears as the eigenvalues $N - M$ times.* The trick to this is that we don't attempt to solve the non-zero eigenvalues, and also don't solve any eigenvectors either, but we seek to prove that $\dim(\ker(AA^T)) = N - M$, which will be sufficient to accomplish the given task. Notice the obvious: $\mathbf{x} \in \ker(AA^T)$ is equivalent with $AA^T\mathbf{x} = 0$, which is equivalent with $AA^T\mathbf{x} = \lambda\mathbf{x}$ for $\lambda = 0$. Hence the kernel contains the eigenvectors corresponding to the eigenvalue zero and nothing else (ignoring the origin). By considering the diagonalization of AA^T we see that the order of the eigenvalue zero (how many times it appears on the diagonal after diagonalization) is the same as the dimension of the kernel.

First we seek to prove that

$$\ker(AA^T) = \ker(A^T) \quad (5)$$

One possibility to (5): First we prove that $\ker(A) = \{\mathbf{0}\}$. So we want to prove

$$A\mathbf{y} = \mathbf{0} \iff \mathbf{y} = \mathbf{0}. \quad (6)$$

Notice that $A\mathbf{y}$ is the same thing as $\sum_{m=1}^M A_{*m}y_m$, which is a linear combination of the columns of A with real coefficients y_m . By definition of the linear independence, this kind of linear combination cannot be zero with non-zero coefficients, so $\ker(A) = \{\mathbf{0}\}$ is clear. Now we know that

$$(AA^T)\mathbf{x} = \mathbf{0} \iff A(A^T\mathbf{x}) = \mathbf{0} \iff A^T\mathbf{x} = \mathbf{0}. \quad (7)$$

which is the same thing as (5).

Second possibility to (5): The direction

$$AA^T\mathbf{x} = \mathbf{0} \iff A^T\mathbf{x} = \mathbf{0} \quad (8)$$

is obvious, so it suffices to prove the other direction. We assume $AA^T\mathbf{x} = \mathbf{0}$. Then also $0 = \mathbf{x}^T(AA^T\mathbf{x}) = (\mathbf{x}^TA)(A^T\mathbf{x}) = \|A^T\mathbf{x}\|^2$, which implies $A^T\mathbf{x} = \mathbf{0}$.

Once (5) has been proven, we seek to prove

$$\dim \ker(A^T) = N - M \quad (9)$$

next. By the rank-nullity theorem equation

$$N = \dim \operatorname{im}(A^T) + \dim \ker(A^T) \quad (10)$$

holds. Therefore it will be sufficient to prove

$$\dim \operatorname{im}(A^T) = M. \quad (11)$$

We make an anti-thesis that $\dim \operatorname{im}(A^T) < M$ would hold. Another basic result is that if $V \subset \mathbb{R}^M$ is some vector space, $V + V^\perp = \mathbb{R}^M$ holds where V^\perp is the orthogonal complement of V , and the sum of the dimensions of V and V^\perp equals the M . The anti-thesis $\dim \operatorname{im}(A^T) < M$ will then imply that the orthogonal complement of the image must have a dimension greater than or equal to 1: $\dim(\operatorname{im}(A^T))^\perp \geq 1$. This implies that in particular $(\operatorname{im}(A^T))^\perp \neq \{\mathbf{0}\}$, and at least one non-zero vector $\mathbf{x} \in (\operatorname{im}(A^T))^\perp$ can be found. In other words, the anti-thesis implies that there exists a vector $\mathbf{x} \in \mathbb{R}^M \setminus \{\mathbf{0}\}$ with the property, that it is orthogonal to all vectors in the $\operatorname{im}(A^T)$. Then, in particular it will be orthogonal to the columns of A^T (which belong to the image), and

$$0 = \mathbf{x} \cdot (A^T)_{*n} \quad \forall 1 \leq n \leq N \implies \mathbf{0} = \mathbf{x}^T A^T \implies \mathbf{0} = A\mathbf{x} \quad (12)$$

and this contradicts the assumption that the columns of A are linearly independent.

We have now proven both equations (5) and (9), and together they imply

$$\dim \ker(AA^T) = N - M. \quad (13)$$

This in turn implies that the matrix AA^T has M non-zero eigenvalues. The exercise requested a proof for a claim that AA^T has M positive eigenvalues. To accomplish this, we must prove that negative eigenvalues cannot exist. Assume \mathbf{x} is a normalized eigenvector of AA^T with an eigenvalue λ . Then

$$\lambda = \lambda \|\mathbf{x}\|^2 = \mathbf{x}^T(\lambda \mathbf{x}) = \mathbf{x}^T AA^T \mathbf{x} = \|A^T \mathbf{x}\|^2 \geq 0, \quad (14)$$

so we see that the eigenvalues of AA^T can never be negative.

2.2 The answer is 2^N .

2.3

$$\bar{A} = A(A^T A)^{-\frac{1}{2}} \quad (15)$$

$$\bar{A}^T \bar{A} = ((A^T A)^{-\frac{1}{2}} A^T) (A(A^T A)^{-\frac{1}{2}}) = (A^T A)^{-\frac{1}{2}} \underbrace{A^T A (A^T A)^{-\frac{1}{2}}}_{=(A^T A)^{\frac{1}{2}}} = \text{id} \quad (16)$$

Set 6 Exercise 1:

1.1: We first check the case $M = 2$.

$$\begin{aligned} \left(\prod_{m=1}^2 \sum_{a_m=1}^D \right) \left(\prod_{m'=1}^2 f_{m'}(a_{m'}) \right) &= \sum_{a_2=1}^D \sum_{a_1=1}^D (f_1(a_1) f_2(a_2)) \\ &= \sum_{a_2=1}^D \left(\sum_{a_1=1}^D (f_1(a_1) f_2(a_2)) \right) \\ &= \sum_{a_2=1}^D \left(f_2(a_2) \sum_{a_1=1}^D f_1(a_1) \right) \\ &= \left(\sum_{a_1=1}^D f_1(a_1) \right) \left(\sum_{a_2=1}^D f_2(a_2) \right) \\ &= \prod_{m=1}^2 \left(\sum_{a_m=1}^D f_m(a_m) \right) \end{aligned} \quad (17)$$

We write the sum as an iterated sum so that the sum over a_1 is inside, and sum over a_2 is outside. In the inner sum we see that $f_2(a_2)$ is a constant with respect to the sum over a_1 , and hence it can be taken in front of the sum. Then the entire sum $\sum_{a_1=1}^D f_1(a_1)$ is a constant with respect to the outer sum over a_2 , and hence this can be taken in front as a constant too.

For arbitrary M the result can be proven with an induction step.

$$\begin{aligned}
\left(\prod_{m=1}^M \sum_{a_m=1}^D \right) \left(\prod_{m'=1}^M f_{m'}(a_{m'}) \right) &= \sum_{a_M=1}^D \cdots \sum_{a_1=1}^D (f_1(a_1) \cdots f_M(a_M)) \\
&= \sum_{a_M=1}^D \cdots \sum_{a_2=1}^D \left(\sum_{a_1=1}^D (f_1(a_1) f_2(a_2) \cdots f_M(a_M)) \right) \\
&= \sum_{a_M=1}^D \cdots \sum_{a_2=1}^D \left(f_2(a_2) \cdots f_M(a_M) \sum_{a_1=1}^D f_1(a_1) \right) \\
&= \left(\sum_{a_1=1}^D f_1(a_1) \right) \left(\sum_{a_M=1}^D \cdots \sum_{a_2=1}^D f_2(a_2) \cdots f_M(a_M) \right) \quad (18) \\
&= \left(\sum_{a_1=1}^D f_1(a_1) \right) \left(\left(\prod_{m=2}^M \sum_{a_m=1}^D \right) \left(\prod_{m'=2}^M f_{m'}(a_{m'}) \right) \right) \\
&= \left(\sum_{a_1=1}^D f_1(a_1) \right) \left(\prod_{m=2}^M \left(\sum_{a_m=1}^D f_m(a_m) \right) \right) \\
&= \prod_{m=1}^M \left(\sum_{a_m=1}^D f_m(a_m) \right)
\end{aligned}$$

We write the sum as an iterated sum, placing the sum with respect to a_1 inside. Then $f_2(a_2) \cdots f_M(a_M)$ is a constant with respect to the sum, and can be taken in front. Then the sum $\sum_{a_1=1}^D f_1(a_1)$ is a constant with respect to all other sums, so it can be taken in front too. Then we use the induction assumption, that the result is known for $M - 1$ dimensional sum, and finally compose the product into one product over m .

Set 6 Exercise 2:

2.3: The obvious answer is

$$\ell(\theta; a, b) = \log(p(A = a, B = b; \theta)) \quad (19)$$

where the $p(A = a, B = b; \theta)$ is probability density with respect to a , and probability with respect to b . To the second request the intended answer is

$$\ell(\theta; a, b) = \log(p_A(a|B = b; \theta)) + \log(P(B = b; \theta)) \quad (20)$$

Here $p_A(a|B = b; \theta)$ is probability density, and $P(B = b; \theta)$ is probability, so there are no mixed probability quantities present.

Set 6 Exercise 4:

4.1 The answer is:

$$\begin{aligned}\ell_{\text{full}}(\theta; \mathbf{x}, r) = & -\frac{1}{2} \sum_{t=1}^T \left(\log \det C_{r(t)} + (\mathbf{x}(t) - \boldsymbol{\mu}_{r(t)})^T C_{r(t)}^{-1} (\mathbf{x}(t) - \boldsymbol{\mu}_{r(t)}) \right) \\ & + \sum_{t=1}^T \log \pi_{r(t)} + \text{const.}\end{aligned}\tag{21}$$

4.3

$$\begin{aligned}\ell_{\text{estim.}}(\theta; \mathbf{x}, \theta_i) = & \sum_{r \in \{1, 2, \dots, K\}^T} \ell_{\text{full}}(\theta; \mathbf{x}, r) p(R = r | \mathbf{X} = \mathbf{x}; \theta_i) \\ = & \sum_{t=1}^T \sum_{r \in \{1, 2, \dots, K\}^T} \log(p(\mathbf{x}(t), r(t); \theta)) p(R = r | \mathbf{X} = \mathbf{x}; \theta_i) \\ = & \sum_{t=1}^T \left(\prod_{t'=1}^T \sum_{r(t')=1}^K \right) \log(p(\mathbf{x}(t), r(t); \theta)) \left(\prod_{t''=1}^T p(r(t'') | \mathbf{x}(t''); \theta_i) \right) \\ = & \sum_{t=1}^T \left(\sum_{r(t)=1}^K \log(p(\mathbf{x}(t), r(t); \theta)) p(r(t) | \mathbf{x}(t); \theta_i) \right) \\ & \underbrace{\left(\prod_{\substack{t''=1 \\ t'' \neq t}}^T \left(\sum_{r(t'')=1}^K p(r(t'') | \mathbf{x}(t''); \theta_i) \right) \right)}_{=1} \\ = & \sum_{t=1}^T \sum_{r(t)=1}^K \log(p(\mathbf{x}(t), r(t); \theta)) p(r(t) | \mathbf{x}(t); \theta_i)\end{aligned}\tag{22}$$

First we write the definition of the estimated log-likelihood, which is otherwise the same as definition in Equation (19) in the original exercise sheet, except that the integral over \mathbf{s} has been replaced with a sum over r . In the sum r goes over all possible vectors $r = (r(1), r(2), \dots, r(T))$, where each $r(t)$ can obtain values from the set $\{1, 2, \dots, K\}$. The total amount of terms in this sum is K^T . We don't immediately substitute the full formula of ℓ_{full} , because it would make the equations unnecessarily complicated. Instead we only substitute the iid assumption, and take the product out of the logarithm, making it a sum. Then the order of sums can be changed. Then we use the iid assumption to $p(R = r | \mathbf{X} = \mathbf{x}; \theta_i)$ and denote

the sum over r in the same way as in the first exercise of set 6. Then with fixed t the product is arranged like

$$\left(\log (p(\mathbf{x}(t), r(t); \theta)) p(r(t) | \mathbf{x}(t); \theta_i) \right) \left(\prod_{\substack{t''=1 \\ t'' \neq t}}^T p(r(t'') | \mathbf{x}(t''); \theta_i) \right) \quad (23)$$

and the result from first exercise is used to change the order of sums and products. In the end we can substitute the known formula for $\log(p(\mathbf{x}(t), r(t); \theta))$ and get

$$\begin{aligned} \ell_{\text{estim.}}(\theta; \mathbf{x}, \theta_i) = \sum_{t=1}^T \sum_{r(t)=1}^K & \left(-\frac{1}{2} \log \det C_{r(t)} - \frac{1}{2} (\mathbf{x}(t) - \boldsymbol{\mu}_{r(t)})^T C_{r(t)}^{-1} (\mathbf{x}(t) - \boldsymbol{\mu}_{r(t)}) \right. \\ & \left. + \log \pi_{r(t)} \right) p(r(t) | \mathbf{x}(t); \theta_i) + \text{const.} \end{aligned} \quad (24)$$

This is precisely the same thing as

$$\begin{aligned} \ell_{\text{estim.}}(\theta; \mathbf{x}, \theta_i) = \sum_{t=1}^T \sum_{k=1}^K & \left(-\frac{1}{2} \log \det C_k - \frac{1}{2} (\mathbf{x}(t) - \boldsymbol{\mu}_k)^T C_k^{-1} (\mathbf{x}(t) - \boldsymbol{\mu}_k) \right. \\ & \left. + \log \pi_k \right) p(R(t) = k | \mathbf{x}(t); \theta_i) + \text{const.} \end{aligned} \quad (25)$$

which is perhaps clearer.

4.4 The most difficult part of exercise 4 of set 6 was probably getting the sums right in 4.3. Once the correct formula has been found for $\ell_{\text{estim.}}$, its gradient can be computed with the same formulas that were used earlier in the exercise set 3. The final place where to avoid mistakes is to remember, that if the expression of $\ell_{\text{estim.}}$ contains a sum over k , then you must not use the same index for partial derivatives. So if you want to apply $\nabla_{\boldsymbol{\mu}_k}$ and ∇_{C_k} on $\ell_{\text{estim.}}$, then the sum must be replaced into a sum over k' for example (or keep $r(t)$ which was the original index). The Kronecker delta $\delta_{kk'}$ (or $\delta_{k,r(t)}$) will appear in very trivial manner in 4.4, but in 4.5 it is not so trivial.

4.5

$$\pi_k = \frac{e^{\gamma_k}}{\sum_{k'=1}^K e^{\gamma_{k'}}} \quad (26)$$

$$\ell_{\text{estim.}}(\theta; \mathbf{x}, \theta_i) = \sum_{t=1}^T \sum_{k=1}^K \log \left(\frac{e^{\gamma_k}}{\sum_{k'=1}^K e^{\gamma_{k'}}} \right) p(R(t) = k | \mathbf{x}(t); \theta_i) + \text{const.} \quad (27)$$

Here we have denoted as constants all terms not depending on π_1, \dots, π_K .

$$\frac{\partial \ell_{\text{estim.}}(\theta; \mathbf{x}, \theta_i)}{\partial \gamma_m} = \sum_{t=1}^T \sum_{k=1}^K \left(\frac{\partial}{\partial \gamma_m} \log \left(\frac{e^{\gamma_k}}{\sum_{k'=1}^K e^{\gamma_{k'}}} \right) \right) p(R(t) = k \mid \mathbf{x}(t); \theta_i) \quad (28)$$

We take a closer look at the derivative inside:

$$\begin{aligned} \frac{\partial}{\partial \gamma_m} \log \left(\frac{e^{\gamma_k}}{\sum_{k'=1}^K e^{\gamma_{k'}}} \right) &= \frac{1}{\pi_k} \frac{\partial}{\partial \gamma_m} \frac{e^{\gamma_k}}{\sum_{k'=1}^K e^{\gamma_{k'}}} \\ &= \frac{1}{\pi_k} \left(\frac{e^{\gamma_k} \delta_{km}}{\sum_{k'=1}^K e^{\gamma_{k'}}} - \frac{e^{\gamma_k} e^{\gamma_m}}{\left(\sum_{k'=1}^K e^{\gamma_{k'}} \right)^2} \right) \\ &= \delta_{km} - \pi_m \end{aligned} \quad (29)$$

Here the relation (26) was used multiple times. We set the partial derivative as zero:

$$\begin{aligned} 0 &= \frac{\partial \ell_{\text{estim.}}(\theta; \mathbf{x}, \theta_i)}{\partial \gamma_m} = \sum_{t=1}^T \sum_{k=1}^K (\delta_{km} - \pi_m) p(R(t) = k \mid \mathbf{x}(t); \theta_i) \\ &= \sum_{t=1}^T p(R(t) = m \mid \mathbf{x}(t); \theta_i) - \pi_m \underbrace{\sum_{t=1}^T \sum_{k=1}^K p(R(t) = k \mid \mathbf{x}(t); \theta_i)}_{=T} \end{aligned} \quad (30)$$

The π_m can be solved out of this, and this is how the iteration formula is found.