# COMP3206: Exercises on Matrix calculus for optimisation

Srinandan Dasmahapatra

2017/8

## 1 Learning the parameters of a Gaussian by Maximum Likelihood Estimation (MLE)

**Estimation**: Data $\mathcal{X} = \{x^1, \ldots, x^n, \ldots x^N\}$ Find mean vector $\mu$ and variance-covariance matrix $\Sigma$.

Assuming data are drawn i.i.d. from Gaussian $\mathcal{N}(x|\mu, \Lambda)$, $\Lambda \triangleq \Sigma^{-1}$, the log likelihood $\mathcal{L}(\mu, \Lambda)$ is

$$\mathcal{L}(\mu, \Lambda) = -\sum_{n=1}^{N} \frac{1}{2}(x^n - \mu)^{\mathsf{T}}\Lambda(x^n - \mu) + \frac{N}{2}\log\det(\Lambda) + \text{constant}.$$

The exercises below will enable you to find the optimal mean and covariance matrix using maximum likelihood estimation.

- Optimal $\mu$: $\frac{\partial}{\partial\mu}\mathcal{L}(\mu, \Lambda) = 0 = \sum_n \Lambda(x^n - \mu)$

$$\sum_n \Lambda x^n = \sum_n \Lambda\mu = N\Lambda\mu \Rightarrow \mu = \frac{1}{N}\sum_n x^n.$$

- Optimal $\Lambda$: $\frac{\partial}{\partial\Lambda}\mathcal{L}(\mu, \Lambda) = 0$.

$$\Sigma = \Lambda^{-1} = \frac{1}{N}\sum_{n=1}^{N}(x^n - \mu)(x^n - \mu)^{\mathsf{T}}.$$

### 1.1 Exercises

1. Verify that the log likelihood function for

$$p(x^n|\mu, \Sigma) = (\sqrt{(2\pi)^p|\Sigma|})^{-1}\exp\left(-\frac{1}{2}(x^n - \mu)^{\mathsf{T}}\Sigma^{-1}(x^n - \mu)\right)$$

is as shown above.

2. For $p \times p$ matrices $\mathbf{A}, \mathbf{B}$ with matrix elements $(\mathbf{A})_{ij} = a_{ij}$ and $(\mathbf{B})_{ij} = b_{ij}$, show that tr($\mathbf{AB}$)=tr($\mathbf{BA}$) by writing tr($\mathbf{A}$) $= \sum_i a_{ii}$) and the product of matrices as
$$(\mathbf{AB})_{ij} = \sum_k a_{ik}b_{kj}.$$

3. For $n \times p$ matrix $\mathbf{A}$ with matrix elements $(\mathbf{A})_{ij} = a_{ij}$, show that the sum of the squares of the matrix elements
$$\sum_{ij} a_{ij}^2 = \text{tr}(\mathbf{AA}^\top), \text{ where tr is the matrix trace.}$$

4. For the Kronecker delta $\delta_{ij}$ defined as
$$\delta_{ij} = \begin{cases} 1 & \text{for } i = j \\ 0 & \text{otherwise.} \end{cases}$$
show that

   (i) $\sum_j a_{ij}\delta_{kj} = a_{ik}$,
   (ii) The diagonal elements of the product of matrices $\mathbf{A}, \mathbf{B}$ is $\sum_{jk} a_{ij}b_{jk}\delta_{ki}$,
   (iii) Trace tr($\mathbf{A}$) $= \sum_{ij} a_{ij}\delta_{ij}$,
   (iv) Make sure you grok $\dfrac{\partial}{\partial x_i}x_j = \delta_{ij}$.

5. For $p \times p$ matrix $\mathbf{A}$ with matrix elements $(\mathbf{A})_{ij} = a_{ij}$ $1 \leqslant i, j \leqslant p$ and vector $\mathbf{x} = (x_1, \ldots, x_p)^\top$ the $i$-th element of vector $(\mathbf{Ax})$ is $(\mathbf{Ax})_i = \sum_{j=1}^p a_{ij}x_j$. Show that $\frac{\partial}{\partial \mathbf{x}}(\mathbf{Ax}) = \mathbf{A}^\top$ by writing out the indices explicitly:
$$\left(\frac{\partial}{\partial \mathbf{x}}(\mathbf{Ax})\right)_{ij} = \frac{\partial}{\partial x_i}(\mathbf{Ax})_j = \frac{\partial}{\partial x_i}\sum_{k=1}^p a_{jk}x_k.$$

6. Show, by writing out the matrix elements as above, that the gradient of the scalar quadratic form $\mathbf{x}\mathbf{Ax}$ is $\frac{\partial}{\partial \mathbf{x}}\mathbf{x}^\top\mathbf{Ax} = (\mathbf{A} + \mathbf{A}^\top)\mathbf{x}$. *Hint:* the $i$-th matrix element of the gradient is
$$\frac{\partial}{\partial x_i}\left(\sum_{r,s=1}^p x_r a_{rs} x_s\right).$$

This should lead to the expression for the MLE of the mean.

7. The partial derivative of the quadratic form $\mathbf{x}\mathbf{A}\mathbf{x}$ with respect to $\mathbf{A}$ can be evaluated for each matrix element $a_{ij}$, $1 \leqslant i, j \leqslant p$:

$$\frac{\partial}{\partial a_{ij}} \left( \sum_{r,s=1}^{p} x_r a_{rs} x_s \right).$$

Show that the result is $\mathbf{x}\mathbf{x}^\top$ (a $p \times p$ matrix).

8. Remember that the determinant of a matrix can be written as a sum

$$\det(\mathbf{A}) = \sum_{j=1}^{p} a_{ij} \text{cof}(a_{ij})$$

where $\text{cof}(a_{ij})$ is $(-1)^{i+j}$ times the determinant of the submatrix of $\mathbf{A}$ obtained by deleting the $i$-th row and $j$-th column. In particular, $\text{cof}(a_{ij})$ does not contain $a_{ij}$. Show that

$$\left( \frac{\partial}{\partial \mathbf{A}} \ln \det \mathbf{A} \right)_{ij} = \frac{\partial}{\partial a_{ij}} \ln \left( \sum_{s=1}^{p} a_{rs} \text{cof}(a_{rs}) \right) = \left( \mathbf{A}^{-1} \right)_{ij}.$$

This and the previous problem should help you derive the MLE of the covariance matrix.

## 2  Regularised linear regression

In regression problems, we have been minimising the residual sum of errors (RSS) with respect to the parameters $\theta$ that are the weight vectors $\mathbf{w}$ in

$$f(\mathbf{x}; \mathbf{w}) = w_0 + \sum_{i=1}^{p} w_i \phi_i(\mathbf{x}).$$

If we introduce $x_0 = 1$, we can write, for each data point $(\mathbf{x}^n, y^n) = (1, x_1^n, \ldots, x_p^n, y^n)$, and the RSS is

$$\text{RSS} = \sum_{n=1}^{N} (r^n)^2 = \sum_{n=1}^{N} (y^n - f(\mathbf{x}^n; \mathbf{w}))^2.$$

When we introduce a $L_2$ regularisation term $\|\mathbf{w}\|_2 = \mathbf{w}^\top \mathbf{w}$ for the weights $\mathbf{w}$, the minimisation is then over a loss function $\ell(\mathbf{w})$:

$$\ell(\mathbf{w}) = -\mathcal{L}(\mathbf{w}) = \sum_{n=1}^{N} (y^n - f(\mathbf{x}^n; \mathbf{w}))^2 + \lambda \mathbf{w}^\top \mathbf{w},$$

where $\lambda$ controls the trade-off between where the data wants the learnt functions to go and how small the modeller wants to keep $\|w\|_2$. Minimising $\ell(w)$ is equivalent to maximising $\mathcal{L}(w)$. This $L_2$ regularised version of linear regression is called **ridge regression**.

## 2.1 Exercises

1. Take the gradient of the loss function $\ell(w)$ with respect to the weight vector $w$ and set it equal to zero. For each vector component $w_i$ compute

$$\frac{\partial}{\partial w_i} \sum_{n=1}^{N} \left(y^n - \sum_{j=0}^{p} w_j \phi_j(x^n)\right)\left(y^n - \sum_{k=0}^{p} w_k \phi_k(x^n)\right) + \lambda \sum_{j,k=0}^{p} w_j w_k \delta_{jk}.$$

Show that the derivative reduces to $-2$ multiplied by

$$\sum_{n=1}^{N} \left\{ \phi_i(x^n) y^n - \sum_{j=0}^{p} (w_j + \lambda \delta_{ij}) \phi_j(x^n) \phi_i(x^n) \right\},$$

a quantity that we will set to zero for max/minimisation.

2. Keep in mind that the data index $n$ in the superscript is a row index while the $i, j, k$ indices for weights stand for columns. Introduce the matrix $\boldsymbol{\Phi}$ with matrix elements $(\boldsymbol{\Phi})_{nj} = \phi_j(x^n)$. Also, the column vector of $y$ values is $\mathbf{y}$. Use this to rewrite the above as a matrix equation

$$\boldsymbol{\Phi}^\top \mathbf{y} = (\boldsymbol{\Phi}^\top \boldsymbol{\Phi} + \lambda \mathbb{I}) w \Rightarrow w = (\boldsymbol{\Phi}^\top \boldsymbol{\Phi} + \lambda \mathbb{I})^{-1} \boldsymbol{\Phi}^\top \mathbf{y}.$$

3. Identify the negative of the loss function $\ell(w)$ with the quantity you worked with for the maximum likelihood estimation problem for the Gaussian. Think about the correspondence. The $\lambda \|w\|_2$ term becomes a prior distribution on weights in a Bayesian interpretation.