

## Exercise 1: Orthogonal

- i) Show that the matrix in Equation 1 takes any vector  $v$  and projects it onto the space spanned by the columns of  $\Phi$ .

$$\Phi (\Phi^\top \Phi)^{-1} \Phi^\top \quad (1)$$

**Solution:** To prove: (1)  $\Phi (\Phi^\top \Phi)^{-1} \Phi^\top$  is a projection matrix; (2) the projection  $\Phi (\Phi^\top \Phi)^{-1} \Phi^\top v$  is in the column space of  $\Phi$ .

(1).

Need to prove that  $P = \Phi (\Phi^\top \Phi)^{-1} \Phi^\top$  is idempotent, i.e.,  $P^2 = P$ .

Simply applying the associativity of matrix multiplication:

$$P^2 = \Phi (\Phi^\top \Phi)^{-1} \Phi^\top \Phi (\Phi^\top \Phi)^{-1} \Phi^\top = \Phi (\Phi^\top \Phi)^{-1} \Phi^\top = P$$

(2).

W.l.o.g. let  $\Phi \in \mathbb{R}^{N \times M}$  and  $v^* = (\Phi^\top \Phi)^{-1} \Phi^\top v$ , then we know that  $v^*$  should be a  $M$ -dimensional vector (i.e.,  $v^* \in \mathbb{R}^{M \times 1}$ ).

By definition of matrix multiplication,

$$\begin{aligned} \Phi v^* &= (\varphi_0, \dots, \varphi_{M-1}) \begin{pmatrix} v_0^* \\ \vdots \\ v_{M-1}^* \end{pmatrix} \\ &= v_0^* \varphi_0 + \dots + v_{M-1}^* \varphi_{M-1} \end{aligned}$$

where  $\varphi_j \in \mathbb{R}^N$  denotes the  $j^{\text{th}}$  column of  $\Phi$ , and  $v_j^* \in \mathbb{R}$  is the  $j^{\text{th}}$  element of  $v^*$ .

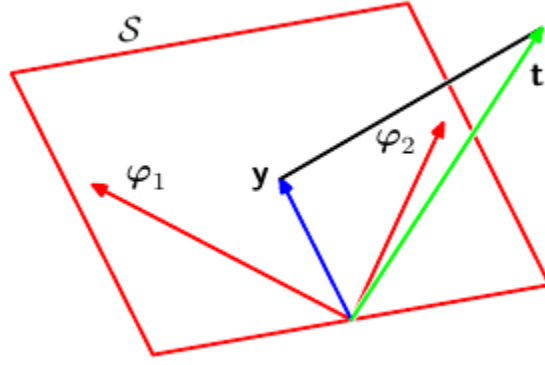
Thus, the projection lies in the space spanned by the columns of  $\Phi$ .

- ii) Use this result to show that the least-squares solution given in Equation 2 corresponds to an orthogonal projection of the vector  $t$  onto the manifold  $\mathcal{S}$  as shown in Figure 1.

$$w_{\text{ML}} = (\Phi^\top \Phi)^{-1} \Phi^\top t \quad (2)$$

**Solution:** To prove: the residual of the projection is orthogonal to the column space of  $\Phi$ .

By definition,  $y = \Phi w_{\text{ML}} = \Phi (\Phi^\top \Phi)^{-1} \Phi^\top t$ .



**Figure 1:** Geometrical interpretation of the least-squares solution: the least-squares regression function is obtained by finding the orthogonal projection of the data vector  $t$  onto the column space of  $\Phi$  ( $\varphi_j$  denotes the  $j^{\text{th}}$  column of  $\Phi$ ).

Thus, we have

$$\begin{aligned}
 (y - t)^\top \Phi &= \left( \Phi (\Phi^\top \Phi)^{-1} \Phi^\top t - t \right)^\top \Phi \\
 &= \left( \left( \Phi (\Phi^\top \Phi)^{-1} \Phi^\top t \right)^\top - t^\top \right) \Phi \\
 &= \left( t^\top \Phi (\Phi^\top \Phi)^{-1} \Phi^\top - t^\top \right) \Phi \\
 &= t^\top \left( \Phi (\Phi^\top \Phi)^{-1} \Phi^\top - I \right) \Phi \\
 &= t^\top \left( \Phi (\Phi^\top \Phi)^{-1} \Phi^\top \Phi - \Phi \right) \\
 &= 0
 \end{aligned}$$

Hence, we prove the orthogonality.

## Exercise 2: Sample weights

Consider a data set in which each data point  $(x_n, y_n)$  is associated with a weighting factor  $r_n > 0$ , so that the sum-of-squares loss function is as in Equation 3.

- Find an expression for the solution  $w^*$  that minimizes this error function.
- Give an interpretation of the weighted sum-of-squares error function in terms of data dependent noise variance.
- Give an interpretation of the weighted sum-of-squares error function in terms of replicated data points.

$$L(X, y, w) = \frac{1}{2} \sum_{n=1}^N r_n (y_n - w^\top \varphi(x_n))^2 \quad (3)$$

**Solution:**

i) To find the minimum, we need to take the derivative w.r.t.  $\mathbf{w}$ :

$$\begin{aligned}\nabla_{\mathbf{w}} L(\mathbf{X}, \mathbf{y}, \mathbf{w}) &= \nabla_{\mathbf{w}} \frac{1}{2} \sum_{n=1}^N r_n (y_n - \mathbf{w}^\top \varphi(\mathbf{x}_n))^2 \\ &= \frac{1}{2} \sum_{n=1}^N r_n \nabla_{\mathbf{w}} (-2y_n \mathbf{w}^\top \varphi(\mathbf{x}_n) + \mathbf{w}^\top \varphi(\mathbf{x}_n) \varphi(\mathbf{x}_n)^\top \mathbf{w}) \\ &= \sum_{n=1}^N r_n (-y_n \varphi(\mathbf{x}_n) + \varphi(\mathbf{x}_n) \varphi(\mathbf{x}_n)^\top \mathbf{w})\end{aligned}$$

The solution  $\mathbf{w}^*$  makes the derivative = 0 (Checking the second-order derivative we see that the Hessian  $\varphi(\mathbf{x}_n) \varphi(\mathbf{x}_n)^\top$  is indeed positive semi-definite):

$$\mathbf{w}^* = \left( \sum_{n=1}^N r_n \varphi(\mathbf{x}_n) \varphi(\mathbf{x}_n)^\top \right)^{-1} \left( \sum_{n=1}^N r_n y_n \varphi(\mathbf{x}_n) \right)$$

Let  $y'_n = \sqrt{r_n} y_n$ ,  $\mathbf{y}' = \begin{pmatrix} y'_1 \\ \vdots \\ y'_N \end{pmatrix}$ ,  $\Phi = \begin{pmatrix} \varphi(\mathbf{x}_1) \\ \vdots \\ \varphi(\mathbf{x}_N) \end{pmatrix}$ ,  $\Phi' = \begin{pmatrix} \sqrt{r_1} \varphi(\mathbf{x}_1) \\ \vdots \\ \sqrt{r_N} \varphi(\mathbf{x}_N) \end{pmatrix}$ , we can rewrite  $\mathbf{w}^*$  as:

$$\mathbf{w}^* = (\Phi'^\top \Phi')^{-1} \Phi'^\top \mathbf{y}'$$

ii) Recall the likelihood, log-likelihood, and the error functions of the vanilla least-squares regression:

$$\begin{aligned}p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta) &= \prod_{n=1}^N \mathcal{N}(y_n | \mathbf{w}^\top \varphi(\mathbf{x}_n), \beta^{-1}) = \prod_{n=1}^N \frac{1}{\sqrt{2\pi\beta^{-1}}} \exp\left(-\frac{(y_n - \mathbf{w}^\top \varphi(\mathbf{x}_n))^2}{2\beta^{-1}}\right) \\ \ln p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta) &= \sum_{n=1}^N \ln \mathcal{N}(y_n | \mathbf{w}^\top \varphi(\mathbf{x}_n), \beta^{-1}) \\ &= \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \frac{\beta}{2} \sum_{n=1}^N (y_n - \mathbf{w}^\top \varphi(\mathbf{x}_n))^2 \\ L(\mathbf{X}, \mathbf{y}, \mathbf{w}) &= \frac{1}{2} \sum_{n=1}^N (y_n - \mathbf{w}^\top \varphi(\mathbf{x}_n))^2\end{aligned}$$

Compare with the weighted sum-of-squares error function:

$$L(\mathbf{X}, \mathbf{y}, \mathbf{w}) = \frac{1}{2} \sum_{n=1}^N r_n (y_n - \mathbf{w}^\top \varphi(\mathbf{x}_n))^2$$

We see that the weighted sum-of-squares error function corresponds to log-likelihood of the form:

$$\ln p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta) \sim \frac{\beta}{2} \sum_{n=1}^N r_n (y_n - \mathbf{w}^\top \varphi(\mathbf{x}_n))^2 \sim \sum_{n=1}^N \ln \mathcal{N}(y_n | \mathbf{w}^\top \varphi(\mathbf{x}_n), r_n^{-1} \beta^{-1})$$

which is Gaussian with data-dependent noise variance.

iii) It can be seen from the log-likelihood function that  $r_n$  can be viewed as the effective number of observation of  $(\mathbf{x}_n, y_n)$ , i.e., you can treat  $(\mathbf{x}_n, y_n)$  as repeatedly occurring  $r_n$  times in the data set.

### Exercise 3: Independent noise

Consider the linear model given in Equation 4 together with the sum-of-squares loss function given in 5. Now suppose that Gaussian noise  $\varepsilon_i$  with zero mean and variance  $\sigma^2$  is added independently to each of the input variables  $x_i$ . By making use of  $\mathbb{E}[\varepsilon_i] = 0$  and  $\mathbb{E}[\varepsilon_i \varepsilon_j] = \delta_{ij} \sigma^2$  show that minimizing  $L(\mathbf{X}, \mathbf{y}, \mathbf{w})$  averaged over the noise distribution is equivalent to minimizing the sum of squares loss for noise-free input variables with the addition of a weight-decay regularization term, in which the bias parameter  $w_0$  is omitted from the regularizer.

$$f(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{i=1}^D w_i x_i \quad (4)$$

$$L(\mathbf{X}, \mathbf{y}, \mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (f(\mathbf{x}_n, \mathbf{w}) - y_n)^2 \quad (5)$$

**Solution:**

Let

$$\begin{aligned} \tilde{y}_n &= w_0 + \sum_{i=1}^D w_i (x_{ni} + \varepsilon_{ni}) \\ &= w_0 + \sum_{i=1}^D w_i x_{ni} + \sum_{i=1}^D w_i \varepsilon_{ni} \\ &= f(\mathbf{x}_n, \mathbf{w}) + \sum_{i=1}^D w_i \varepsilon_{ni} \end{aligned}$$

From Equation 5, we then have:

$$\begin{aligned} \tilde{L}(\mathbf{X}, \mathbf{y}, \mathbf{w}) &= \frac{1}{2} \sum_{n=1}^N (\tilde{y}_n - y_n)^2 \\ &= \frac{1}{2} \sum_{n=1}^N (\tilde{y}_n^2 - 2\tilde{y}_n y_n + y_n^2) \\ &= \frac{1}{2} \sum_{n=1}^N \left( \left( f(\mathbf{x}_n, \mathbf{w}) + \sum_{i=1}^D w_i \varepsilon_{ni} \right)^2 - 2 \left( f(\mathbf{x}_n, \mathbf{w}) + \sum_{i=1}^D w_i \varepsilon_{ni} \right) y_n + y_n^2 \right) \\ &= \frac{1}{2} \sum_{n=1}^N \left( f(\mathbf{x}_n, \mathbf{w})^2 + \left( \sum_{i=1}^D w_i \varepsilon_{ni} \right)^2 + 2f(\mathbf{x}_n, \mathbf{w}) \sum_{i=1}^D w_i \varepsilon_{ni} - 2f(\mathbf{x}_n, \mathbf{w}) y_n - 2 \sum_{i=1}^D w_i \varepsilon_{ni} y_n + y_n^2 \right) \end{aligned}$$

Taking the expectation over  $\varepsilon_{ni}$ , we have:

$$\mathbb{E} \left[ 2f(\mathbf{x}_n, \mathbf{w}) \sum_{i=1}^D w_i \varepsilon_{ni} \right] = 0 \text{ and } \mathbb{E} \left[ 2 \sum_{i=1}^D w_i \varepsilon_{ni} y_n \right] = 0, \text{ since } \mathbb{E}[\varepsilon_{ni}] = 0$$

$$\mathbb{E} \left[ \left( \sum_{i=1}^D w_i \varepsilon_{ni} \right)^2 \right] = \mathbb{E} \left[ \sum_{i=1}^D \sum_{j=1}^D w_i w_j \varepsilon_{ni} \varepsilon_{nj} \right] = \sum_{i=1}^D w_i^2 \sigma^2, \text{ since } \mathbb{E}[\varepsilon_{ni} \varepsilon_{nj}] = \delta_{ij} \sigma^2$$

Thus,

$$\begin{aligned}
\mathbb{E} [\tilde{L}(\mathbf{X}, \mathbf{y}, \mathbf{w})] &= \mathbb{E} \left[ \frac{1}{2} \sum_{n=1}^N (f(\mathbf{x}_n, \mathbf{w})^2 - 2f(\mathbf{x}_n, \mathbf{w})y_n + y_n^2) \right] + \frac{1}{2} \sum_{n=1}^N \mathbb{E} \left[ \left( \sum_{i=1}^D w_i \varepsilon_{ni} \right)^2 \right] \\
&= \mathbb{E} [L(\mathbf{X}, \mathbf{y}, \mathbf{w})] + \frac{N\sigma^2}{2} \sum_{i=1}^D w_i^2 \\
&= L(\mathbf{X}, \mathbf{y}, \mathbf{w}) + \frac{N\sigma^2}{2} \sum_{i=1}^D w_i^2
\end{aligned}$$

where

- Left-hand side: the sum of squares loss (for *noisy* input) averaged over the noise distribution
- Right-hand side: the sum of squares loss (for *noise-free* input) + a weight-decay regularization term

## Exercise 4: Linear basis functions

Consider a linear basis function regression model for a multivariate target variable  $\mathbf{y}$  having a Gaussian distribution as given in Equation 6 where  $f(\mathbf{x}, \mathbf{W}) = \mathbf{W}^\top \varphi(\mathbf{x})$  together with a training data set comprising input basis vectors  $\varphi(\mathbf{x}_n)$  and corresponding target vectors  $\mathbf{y}_n$ , with  $n = 1, \dots, N$ .

$$p(\mathbf{y} | \mathbf{W}, \Sigma) = \mathcal{N}(\mathbf{y} | f(\mathbf{x}, \mathbf{W}), \Sigma) \quad (6)$$

- Show that the maximum likelihood solution  $\mathbf{W}_{\text{ML}}$  for the parameter matrix  $\mathbf{W}$  has the property that each column is given by an expression shown in Equation 2, which was the solution for an isotropic noise distribution. Note that this is independent of the covariance matrix  $\Sigma$ .
- Show that the maximum likelihood solution for  $\Sigma$  is given by

$$\Sigma_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N (\mathbf{y}_n - \mathbf{W}_{\text{ML}}^\top \varphi(\mathbf{x}_n)) (\mathbf{y}_n - \mathbf{W}_{\text{ML}}^\top \varphi(\mathbf{x}_n))^\top$$

**Solution:** The likelihood and log-likelihood functions are:

$$\begin{aligned}
p(\mathbf{y} | \mathbf{X}, \mathbf{W}, \Sigma) &= \prod_{n=1}^N \mathcal{N}(\mathbf{y}_n | \mathbf{W}^\top \varphi(\mathbf{x}_n), \Sigma) = \prod_{n=1}^N \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} \exp \left( -\frac{1}{2} (\mathbf{y}_n - \mathbf{W}^\top \varphi(\mathbf{x}_n))^\top \Sigma^{-1} (\mathbf{y}_n - \mathbf{W}^\top \varphi(\mathbf{x}_n)) \right) \\
\ln p(\mathbf{y} | \mathbf{X}, \mathbf{W}, \Sigma) &= \sum_{n=1}^N \ln \mathcal{N}(\mathbf{y}_n | \mathbf{W}^\top \varphi(\mathbf{x}_n), \Sigma) \\
&= -\frac{N}{2} \ln |\Sigma| - \frac{ND}{2} \ln (2\pi) - \frac{1}{2} \sum_{n=1}^N (\mathbf{y}_n - \mathbf{W}^\top \varphi(\mathbf{x}_n))^\top \Sigma^{-1} (\mathbf{y}_n - \mathbf{W}^\top \varphi(\mathbf{x}_n))
\end{aligned}$$

- Take the derivative of log-likelihood w.r.t.  $\mathbf{W}$ :

$$\begin{aligned}
\nabla_{\mathbf{W}} \ln p(\mathbf{y} | \mathbf{X}, \mathbf{W}, \Sigma) &= -\frac{1}{2} \nabla_{\mathbf{W}} \sum_{n=1}^N (\mathbf{y}_n - \mathbf{W}^\top \varphi(\mathbf{x}_n))^\top \Sigma^{-1} (\mathbf{y}_n - \mathbf{W}^\top \varphi(\mathbf{x}_n)) \\
&= -\frac{1}{2} \nabla_{\mathbf{W}} \sum_{n=1}^N (-\mathbf{y}_n^\top \Sigma^{-1} \mathbf{W}^\top \varphi(\mathbf{x}_n) - \varphi(\mathbf{x}_n)^\top \mathbf{W} \Sigma^{-1} \mathbf{y}_n + \varphi(\mathbf{x}_n)^\top \mathbf{W} \Sigma^{-1} \mathbf{W}^\top \varphi(\mathbf{x}_n))
\end{aligned}$$

Let's investigate each term separately:

$$\begin{aligned}\nabla_{\mathbf{W}} (-\mathbf{y}_n^\top \boldsymbol{\Sigma}^{-1} \mathbf{W}^\top \varphi(\mathbf{x}_n)) &= \nabla_{\mathbf{W}} \text{tr}(-\mathbf{y}_n^\top \boldsymbol{\Sigma}^{-1} \mathbf{W}^\top \varphi(\mathbf{x}_n)) \\ &= -\nabla_{\mathbf{W}} \text{tr}(\mathbf{W}^\top \varphi(\mathbf{x}_n) \mathbf{y}_n^\top \boldsymbol{\Sigma}^{-1}) \\ &= -\varphi(\mathbf{x}_n) \mathbf{y}_n^\top \boldsymbol{\Sigma}^{-1}\end{aligned}$$

$$\begin{aligned}\nabla_{\mathbf{W}} (-\varphi(\mathbf{x}_n)^\top \mathbf{W} \boldsymbol{\Sigma}^{-1} \mathbf{y}_n) &= \nabla_{\mathbf{W}} \text{tr}(-\varphi(\mathbf{x}_n)^\top \mathbf{W} \boldsymbol{\Sigma}^{-1} \mathbf{y}_n) \\ &= -\text{tr}(\mathbf{W} \boldsymbol{\Sigma}^{-1} \mathbf{y}_n \varphi(\mathbf{x}_n)^\top) \\ &= -(\boldsymbol{\Sigma}^{-1} \mathbf{y}_n \varphi(\mathbf{x}_n)^\top)^\top \\ &= -\varphi(\mathbf{x}_n) \mathbf{y}_n^\top \boldsymbol{\Sigma}^{-1}\end{aligned}$$

$$\begin{aligned}\nabla_{\mathbf{W}} (\varphi(\mathbf{x}_n)^\top \mathbf{W} \boldsymbol{\Sigma}^{-1} \mathbf{W}^\top \varphi(\mathbf{x}_n)) &= \nabla_{\mathbf{W}} \text{tr}(\varphi(\mathbf{x}_n)^\top \mathbf{W} \boldsymbol{\Sigma}^{-1} \mathbf{W}^\top \varphi(\mathbf{x}_n)) \\ &= \nabla_{\mathbf{W}} \text{tr}(\mathbf{W} \boldsymbol{\Sigma}^{-1} \mathbf{W}^\top \varphi(\mathbf{x}_n) \varphi(\mathbf{x}_n)^\top) \\ &= \varphi(\mathbf{x}_n) \varphi(\mathbf{x}_n)^\top \mathbf{W} \boldsymbol{\Sigma}^{-1} + \varphi(\mathbf{x}_n) \varphi(\mathbf{x}_n)^\top \mathbf{W} (\boldsymbol{\Sigma}^{-1})^\top \\ &= 2\varphi(\mathbf{x}_n) \varphi(\mathbf{x}_n)^\top \mathbf{W} \boldsymbol{\Sigma}^{-1}\end{aligned}$$

As the derivative vanishes at  $\mathbf{W}_{\text{ML}}$ , we have:

$$\sum_{n=1}^N \varphi(\mathbf{x}_n) \varphi(\mathbf{x}_n)^\top \mathbf{W}_{\text{ML}} \boldsymbol{\Sigma}^{-1} = \sum_{n=1}^N \varphi(\mathbf{x}_n) \mathbf{y}_n^\top \boldsymbol{\Sigma}^{-1}$$

Right multiply both side by  $\boldsymbol{\Sigma}$  and rewrite the equation using the design matrix  $\boldsymbol{\Phi}$  and target matrix  $\mathbf{Y}$ :

$$\boldsymbol{\Phi}^\top \boldsymbol{\Phi} \mathbf{W}_{\text{ML}} = \boldsymbol{\Phi}^\top \mathbf{Y}$$

i.e.,

$$\mathbf{W}_{\text{ML}} = (\boldsymbol{\Phi}^\top \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^\top \mathbf{Y}$$

ii) Take the derivative of log-likelihood function w.r.t.  $\boldsymbol{\Sigma}^{-1}$ :

$$\begin{aligned}\nabla_{\boldsymbol{\Sigma}^{-1}} \ln p(\mathbf{y}|\mathbf{X}, \mathbf{W}, \boldsymbol{\Sigma}) &= \nabla_{\boldsymbol{\Sigma}^{-1}} \frac{N}{2} \ln |\boldsymbol{\Sigma}|^{-1} - \frac{1}{2} \nabla_{\boldsymbol{\Sigma}^{-1}} \sum_{n=1}^N (\mathbf{y}_n - \mathbf{W}^\top \varphi(\mathbf{x}_n))^\top \boldsymbol{\Sigma}^{-1} (\mathbf{y}_n - \mathbf{W}^\top \varphi(\mathbf{x}_n)) \\ &= \nabla_{\boldsymbol{\Sigma}^{-1}} \frac{N}{2} \ln |\boldsymbol{\Sigma}^{-1}| - \frac{1}{2} \nabla_{\boldsymbol{\Sigma}^{-1}} \sum_{n=1}^N \text{tr}(\boldsymbol{\Sigma}^{-1} (\mathbf{y}_n - \mathbf{W}^\top \varphi(\mathbf{x}_n)) (\mathbf{y}_n - \mathbf{W}^\top \varphi(\mathbf{x}_n))^\top) \\ &= \frac{N}{2} \boldsymbol{\Sigma}^\top - \frac{1}{2} \sum_{n=1}^N (\mathbf{y}_n - \mathbf{W}^\top \varphi(\mathbf{x}_n)) (\mathbf{y}_n - \mathbf{W}^\top \varphi(\mathbf{x}_n))^\top \\ &= \frac{N}{2} \boldsymbol{\Sigma} - \frac{1}{2} \sum_{n=1}^N (\mathbf{y}_n - \mathbf{W}^\top \varphi(\mathbf{x}_n)) (\mathbf{y}_n - \mathbf{W}^\top \varphi(\mathbf{x}_n))^\top\end{aligned}$$

Hence, we have

$$\boldsymbol{\Sigma}_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N (\mathbf{y}_n - \mathbf{W}^\top \varphi(\mathbf{x}_n)) (\mathbf{y}_n - \mathbf{W}^\top \varphi(\mathbf{x}_n))^\top$$

Reference: Chapter 15 of [3]

## References

- [1] C. M. Bishop. *Pattern recognition and machine learning*. springer, 2006.
- [2] J. Friedman, T. Hastie, R. Tibshirani, et al. *The elements of statistical learning*, volume 1. Springer series in statistics New York, 2001.
- [3] J. R. Magnus and H. Neudecker. *Matrix differential calculus with applications in statistics and econometrics*. John Wiley & Sons, 2019.