

which is the mean of the observed set of data points. The maximization of (3.102) with respect to Σ is rather more involved. The simplest approach is to ignore the symmetry constraint and show that the resulting solution is symmetric as required. Alternative derivations of this result, which impose the symmetry and positive definiteness constraints explicitly, can be found in Magnus and Neudecker (1999). The result is as expected and takes the form

$$\Sigma_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu}_{\text{ML}})(\mathbf{x}_n - \boldsymbol{\mu}_{\text{ML}})^T, \quad (3.106)$$

which involves $\boldsymbol{\mu}_{\text{ML}}$ because this is the result of a joint maximization with respect to $\boldsymbol{\mu}$ and Σ . Note that the solution (3.105) for $\boldsymbol{\mu}_{\text{ML}}$ does not depend on Σ_{ML} , and so we can first evaluate $\boldsymbol{\mu}_{\text{ML}}$ and then use this to evaluate Σ_{ML} .

If we evaluate the expectations of the maximum likelihood solutions under the true distribution, we obtain the following results

$$\mathbb{E}[\boldsymbol{\mu}_{\text{ML}}] = \boldsymbol{\mu} \quad (3.107)$$

$$\mathbb{E}[\Sigma_{\text{ML}}] = \frac{N-1}{N} \Sigma. \quad (3.108)$$

We see that the expectation of the maximum likelihood estimate for the mean is equal to the true mean. However, the maximum likelihood estimate for the covariance has an expectation that is less than the true value, and hence, it is biased. We can correct this bias by defining a different estimator $\tilde{\Sigma}$ given by

$$\tilde{\Sigma} = \frac{1}{N-1} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu}_{\text{ML}})(\mathbf{x}_n - \boldsymbol{\mu}_{\text{ML}})^T. \quad (3.109)$$

Clearly from (3.106) and (3.108), the expectation of $\tilde{\Sigma}$ is equal to Σ .

3.2.8 Sequential estimation

Our discussion of the maximum likelihood solution represents a *batch* method in which the entire training data set is considered at once. An alternative is to use *sequential* methods, which allow data points to be processed one at a time and then discarded. These are important for online applications and for large data when the batch processing of all data points at once is infeasible.

Consider the result (3.105) for the maximum likelihood estimator of the mean $\boldsymbol{\mu}_{\text{ML}}$, which we will denote by $\boldsymbol{\mu}_{\text{ML}}^{(N)}$ when it is based on N observations. If we

Exercise 3.28

Exercise 3.29