### 4.1.3   MLE for an MVN

We now describe one way to estimate the parameters of an MVN, using MLE. In later sections, we will discuss Bayesian inference for the parameters, which can mitigate overfitting, and can provide a measure of confidence in our estimates.

**Theorem 4.1.1** (MLE for a Gaussian). *If we have $N$ iid samples $x_i \sim \mathcal{N}(\mu, \Sigma)$, then the MLE for the parameters is given by*

$$\hat{\mu}_{mle} = \frac{1}{N} \sum_{i=1}^{N} x_i \triangleq \bar{x} \tag{4.6}$$

$$\hat{\Sigma}_{mle} = \frac{1}{N} \sum_{i=1}^{N} (x_i - \bar{x})(x_i - \bar{x})^T = \frac{1}{N} \left( \sum_{i=1}^{N} x_i x_i^T \right) - \bar{x}\,\bar{x}^T \tag{4.7}$$

*That is, the MLE is just the empirical mean and empirical covariance. In the univariate case, we get the following familiar results:*

$$\hat{\mu} = \frac{1}{N} \sum_{i} x_i = \bar{x} \tag{4.8}$$

$$\hat{\sigma}^2 = \frac{1}{N} \sum_{i} (x_i - \bar{x})^2 = \left( \frac{1}{N} \sum_{i} x_i^2 \right) - (\bar{x})^2 \tag{4.9}$$

#### 4.1.3.1   Proof *

To prove this result, we will need several results from matrix algebra, which we summarize below. In the equations, a and b are vectors, and A and B are matrices. Also, the notation $tr(A)$ refers to the **trace** of a matrix, which is the sum of its diagonals: $tr(A) = \sum_i A_{ii}$.

$$\frac{\partial(b^T a)}{\partial a} = b$$

$$\frac{\partial(a^T A a)}{\partial a} = (A + A^T)a$$

$$\frac{\partial}{\partial A} tr(BA) = B^T \tag{4.10}$$

$$\frac{\partial}{\partial A} \log |A| = A^{-T} \triangleq (A^{-1})^T$$

$$tr(ABC) = tr(CAB) = tr(BCA)$$

The last equation is called the **cyclic permutation property** of the trace operator. Using this, we can derive the widely used **trace trick**, which reorders the scalar inner product $x^T A x$ as follows

$$x^T A x = tr(x^T A x) = tr(x x^T A) = tr(A x x^T) \tag{4.11}$$

*Proof.* We can now begin with the proof. The log-likelihood (dropping additive constants) is given by

$$\ell(\mu, \Sigma) = \log p(\mathcal{D}|\mu, \Sigma) = \frac{N}{2} \log |\Lambda| - \frac{1}{2} \sum_{i=1}^{N} (x_i - \mu)^T \Lambda (x_i - \mu) \tag{4.12}$$

where $\Lambda = \Sigma^{-1}$ is the precision matrix.

Using the substitution $y_i = x_i - \mu$ and the chain rule of calculus, we have

$$\frac{\partial}{\partial \mu} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) = \frac{\partial}{\partial y_i} y_i^T \Sigma^{-1} y_i \frac{\partial y_i}{\partial \mu} \tag{4.13}$$

$$= -1(\Sigma^{-1} + \Sigma^{-T}) y_i \tag{4.14}$$

Hence

$$\frac{\partial}{\partial \mu} \ell(\mu, \Sigma) = -\frac{1}{2} \sum_{i=1}^{N} -2\Sigma^{-1}(x_i - \mu) = \Sigma^{-1} \sum_{i=1}^{N} (x_i - \mu) = 0 \tag{4.15}$$

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^{N} x_i = \bar{x} \tag{4.16}$$

So the MLE of $\mu$ is just the empirical mean.

Now we can use the trace-trick to rewrite the log-likelihood for $\Lambda$ as follows:

$$\ell(\Lambda) = \frac{N}{2} \log |\Lambda| - \frac{1}{2} \sum_{i} \text{tr}[(x_i - \mu)(x_i - \mu)^T \Lambda] \tag{4.17}$$

$$= \frac{N}{2} \log |\Lambda| - \frac{1}{2} \text{tr}[S_\mu \Lambda] \tag{4.18}$$

$$\tag{4.19}$$

where

$$S_\mu \triangleq \sum_{i=1}^{N} (x_i - \mu)(x_i - \mu)^T \tag{4.20}$$

is the scatter matrix centered on $\mu$. Taking derivatives of this expression with respect to $\Lambda$ yields

$$\frac{\partial \ell(\Lambda)}{\partial \Lambda} = \frac{N}{2} \Lambda^{-T} - \frac{1}{2} S_\mu^T = 0 \tag{4.21}$$

$$\Lambda^{-T} = \Lambda^{-1} = \Sigma = \frac{1}{N} S_\mu \tag{4.22}$$

so

$$\hat{\Sigma} = \frac{1}{N} \sum_{i=1}^{N} (x_i - \mu)(x_i - \mu)^T \tag{4.23}$$

which is just the empirical covariance matrix centered on $\mu$. If we plug-in the MLE $\mu = \bar{x}$ (since both parameters must be simultaneously optimized), we get the standard equation for the MLE of a covariance matrix.

□

### 4.1.4 Maximum entropy derivation of the Gaussian *

In this section, we show that the multivariate Gaussian is the distribution with maximum entropy subject to having a specified mean and covariance (see also Section 9.2.6). This is one reason the Gaussian is so widely used: the first two moments are usually all that we can reliably estimate from data, so we want a distribution that captures these properties, but otherwise makes as few additional assumptions as possible.

To simplify notation, we will assume the mean is zero. The pdf has the form

$$p(\mathbf{x}) = \frac{1}{Z} \exp(-\frac{1}{2}\mathbf{x}^T \boldsymbol{\Sigma}^{-1}\mathbf{x}) \tag{4.24}$$

If we define $f_{ij}(\mathbf{x}) = x_i x_j$ and $\lambda_{ij} = \frac{1}{2}(\boldsymbol{\Sigma}^{-1})_{ij}$, for $i, j \in \{1, \ldots, D\}$, we see that this is in the same form as Equation 9.74. The (differential) entropy of this distribution (using log base $e$) is given by

$$h(\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})) = \frac{1}{2} \ln \left[(2\pi e)^D |\boldsymbol{\Sigma}|\right] \tag{4.25}$$

We now show the MVN has maximum entropy amongst all distributions with a specified covariance $\boldsymbol{\Sigma}$.

**Theorem 4.1.2.** *Let $q(\mathbf{x})$ be any density satisfying $\int q(\mathbf{x})x_i x_j\, d\mathbf{x} = \Sigma_{ij}$. Let $p = \mathcal{N}(0, \boldsymbol{\Sigma})$. Then $h(q) \le h(p)$.*

*Proof.* (From (Cover and Thomas 1991, p234).) We have

$$0 \le \text{KL}(q\|p) = \int q(\mathbf{x}) \log \frac{q(\mathbf{x})}{p(\mathbf{x})} d\mathbf{x} \tag{4.26}$$

$$= -h(q) - \int q(\mathbf{x}) \log p(\mathbf{x}) d\mathbf{x} \tag{4.27}$$

$$=^* -h(q) - \int p(\mathbf{x}) \log p(\mathbf{x}) d\mathbf{x} \tag{4.28}$$

$$= -h(q) + h(p) \tag{4.29}$$

where the key step in Equation 4.28 (marked with a *) follows since $q$ and $p$ yield the same moments for the quadratic form encoded by $\log p(\mathbf{x})$.

□

## 4.2 Gaussian discriminant analysis

One important application of MVNs is to define the class conditional densities in a generative classifier, i.e.,

$$p(\mathbf{x}|y = c, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c) \tag{4.30}$$