

# Maximum Likelihood Estimation

In the last couple of lectures, we developed a pretty good understanding of how to predict the outcomes of random variables given observations of related random variables. Here we consider the different, but related problem, of **parameter estimation**. We observe one or more outcomes of a random variable whose distribution is controlled (parameterized) by one or more variables that we collect in the vector  $\boldsymbol{\theta}$ . More mathematically, the distribution of our data is assumed to be<sup>1</sup>

$$X \sim f_X(\mathbf{x}; \boldsymbol{\theta}),$$

for some *unknown*  $\boldsymbol{\theta} \in \mathcal{T}$ ; from a sample (or samples) of  $X$ , we want to estimate the latent  $\boldsymbol{\theta}$ .

Here are some stylized examples that can be put into this framework.

**Example:** What is the probability that LeBron James makes a free throw? The model here is that the outcome of each free throw can be captured using a binary-valued random variable  $X_i$ :

$$X_i = 0 \text{ if he misses, } X_i = 1 \text{ if he makes it.}$$

The distribution of these random variables is controlled by a single parameter  $\theta$ :

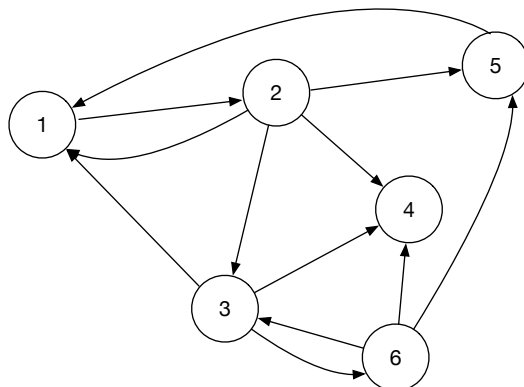
$$X_i = \begin{cases} 1, & \text{with probability } \theta, \\ 0, & \text{with probability } 1 - \theta. \end{cases}$$

---

<sup>1</sup> $f_X(\mathbf{x}; \boldsymbol{\theta})$  should be thought of as a density with argument  $\mathbf{x}$  that is parameterized by some  $\boldsymbol{\theta}$  — different  $\boldsymbol{\theta}$  gives us different probability density functions.

(We are assuming here that the value of  $\theta$  is the same for every free throw.) Given a series of observations  $X_1 = x_1, X_2 = x_2, \dots, X_N = x_N$ , how can we estimate  $\theta$ ?

**Example:** Suppose I represent a social network with  $D$  users with a directed graph like this:



where if there is a directed edge from node  $i$  to node  $j$ , it means user  $j$  is “following” user  $i$ . Suppose that every day, a user (the same user everyday) creates and then shares some “fake news” with his followers. With probability 0.9, a follower finds this news credible and passes it on to their followers; with probability 0.1, a follower flags the news as fake and reports it to you (and does not pass it along). At the end of every day, you have a binary vector whose length is  $D$ ,

$$X_i = \begin{bmatrix} X_i[1] \\ X_i[2] \\ \vdots \\ X_i[D] \end{bmatrix}, \quad X_i[d] = \begin{cases} 1, & \text{if user } d \text{ reported “fake news” on day } i \\ 0, & \text{otherwise} \end{cases}.$$

Given observations  $X_1, X_2, \dots, X_N$  over several days, how can I estimate which user is generating the “fake news”?

**Example:** Suppose that  $X_1, \dots, X_N$  are independent realizations of a  $D$ -dimensional Gaussian random vector with unknown mean  $\boldsymbol{\mu}$  and unknown covariance  $\mathbf{R}$ . Given these  $N$  realizations, how can I estimate  $\boldsymbol{\theta} = (\boldsymbol{\mu}, \mathbf{R})$ ?

In the first example above, the observations are scalars, there is one parameter, and the parameter space  $\mathcal{T} = [0, 1]$  was continuous (all reals between zero and one). In the second, the observations are vectors, there is one parameter, and the parameter space  $\mathcal{T} = \{1, 2, \dots, D\}$  is finite. In the third example, the observations are vectors, the unknown parameters are a vector and a matrix, and the parameter space is  $\mathcal{T} = \mathbb{R}^D \otimes \mathcal{S}_{++}^D$ , where  $\mathcal{S}_{++}^D$  is the set of all  $D \times D$  symmetric positive-definite matrices. In all three of the examples (as in all problems of this type), the parameter(s) induces different distributions on the observed data.

With a probabilistic model in place for the observations, given a sample<sup>2</sup>  $X_1 = \mathbf{x}_1, X_2 = \mathbf{x}_2, X_N = \mathbf{x}_N$ , the **likelihood** of a particular set of parameters  $\boldsymbol{\theta}$  is

$$L(\boldsymbol{\theta}; \mathbf{x}_1, \dots, \mathbf{x}_N) = f_{X_1, \dots, X_N}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N; \boldsymbol{\theta}).$$

We introduce this new notation to emphasize that the likelihood should be thought of as a function of  $\boldsymbol{\theta}$ . The maximum likelihood estimation is simply the parameters that maximize  $L(\boldsymbol{\theta}; \cdot)$ :

$$\hat{\boldsymbol{\theta}}_{\text{mle}} = \arg \max_{\boldsymbol{\theta} \in \mathcal{T}} L(\boldsymbol{\theta}; \mathbf{x}_1, \dots, \mathbf{x}_N).$$

Because many times the joint distribution of  $X_1, \dots, X_N$  involves multiplying a bunch of functions together (especially when the observations are independent), it is often convenient to work with the

---

<sup>2</sup>As in the two examples, the observations might be vectors or scalars.

## log likelihood

$$\ell(\boldsymbol{\theta}; \mathbf{x}_1, \dots, \mathbf{x}_N) = \log f_{X_1, \dots, X_N}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N; \boldsymbol{\theta}).$$

It should be absolutely clear<sup>3</sup> that  $L(\boldsymbol{\theta}; \cdot)$  and  $\ell(\boldsymbol{\theta}; \cdot)$  are maximized at the same place, so

$$\hat{\boldsymbol{\theta}}_{\text{mle}} = \arg \max_{\boldsymbol{\theta} \in \mathcal{T}} \ell(\boldsymbol{\theta}; \mathbf{x}_1, \dots, \mathbf{x}_N) = \arg \min_{\boldsymbol{\theta} \in \mathcal{T}} -\ell(\boldsymbol{\theta}; \mathbf{x}_1, \dots, \mathbf{x}_N).$$

Let's see how this works.

**Example:** In the LeBron James freethrow example above, given  $X_1 = x_1, \dots, X_N = x_N$ , and assuming each of the trials are independent of one another, we have

$$L(\theta; x_1, \dots, x_N) = \prod_{n=1}^N \theta^{x_n} (1 - \theta)^{1-x_n}.$$

Notice that the expression inside the product above is  $\theta$  if  $x_n = 1$ , and  $1 - \theta$  if  $x_n = 0$ . With

$$S_N = \sum_{n=1}^N x_n,$$

the expression above becomes

$$L(\theta; x_1, \dots, x_N) = \theta^{S_N} (1 - \theta)^{N-S_N},$$

---

<sup>3</sup>If it is not absolutely clear to you, then you have a new homework question: Let  $\phi : \mathbb{R} \rightarrow \mathbb{R}$  be a monotonically increasing function on  $\mathbb{R}$ , and  $f : \mathbb{R}^D \rightarrow \mathbb{R}$  be a function on  $\mathbb{R}^D$ . Show that  $\arg \max_{\mathbf{v}} \phi(f(\mathbf{v})) = \arg \max_{\mathbf{v}} f(\mathbf{v})$ .

and so

$$\ell(\theta; x_1, \dots, x_N) = S_N \log \theta + (N - S_N) \log(1 - \theta).$$

The first and second derivatives of  $\ell$  are

$$\begin{aligned} \frac{d\ell(\theta; \cdot)}{d\theta} &= \frac{S_N}{\theta} - \frac{N - S_N}{1 - \theta} \\ \frac{d^2\ell(\theta; \cdot)}{d\theta^2} &= -\frac{S_N}{\theta^2} - \frac{N - S_N}{(1 - \theta)^2}. \end{aligned}$$

Since  $0 \leq S_N \leq N$ , the second derivative is  $\leq 0$  for all  $\theta \in \mathcal{T} = [0, 1]$ , and so we can find the maximizer by setting the first derivative equal to zero. This yields

$$\hat{\theta}_{\text{mle}} = \frac{S_N}{N} = \frac{1}{N} \sum_{n=1}^N x_n.$$

**Example:** Suppose that  $X$  is a scalar random variable distributed uniformly on the interval  $[a, b]$ ,  $X \sim \text{Uniform}([a, b])$ , where  $a$  and  $b$  are unknown. Given  $X_1 = x_1, X_2 = x_2, \dots, X_N = x_N$ , what is the MLE for  $\boldsymbol{\theta} = (a, b)$ ?

**Answer:**

**Example:** Suppose we observe  $X_1 = \mathbf{x}_1, \dots, X_N = \mathbf{x}_N$ , where the  $X_n$  are independent and identically distributed Gaussian random vectors in  $\mathbb{R}^D$ :

$$X_n \sim \text{Normal}(\boldsymbol{\mu}, \mathbf{R}).$$

From these observations, we want to estimate  $\boldsymbol{\mu}$  and  $\mathbf{R}$ . The MLE is the solution to

$$\underset{\boldsymbol{\mu} \in \mathbb{R}^D, \mathbf{R} \in \mathcal{S}_{++}^D}{\text{maximize}} \quad \prod_{n=1}^N (2\pi)^{-D/2} (\det \mathbf{R})^{-1/2} \exp(-(\mathbf{x}_n - \boldsymbol{\mu})^T \mathbf{R}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) / 2).$$

Taking the log of the likelihood function, this is equivalent to

$$\underset{\boldsymbol{\mu} \in \mathbb{R}^D, \mathbf{R} \in \mathcal{S}_{++}^D}{\text{maximize}} \quad \frac{-ND}{2} \log(2\pi) + \frac{N}{2} \log \det \mathbf{R}^{-1} - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \mathbf{R}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}).$$

The first terms doesn't depend on either  $\boldsymbol{\mu}$  or  $\mathbf{R}$ , so we can drop it, leaving us with

$$\underset{\boldsymbol{\mu} \in \mathbb{R}^D, \mathbf{R} \in \mathcal{S}_{++}^D}{\text{maximize}} \quad \frac{N}{2} \log \det \mathbf{R}^{-1} - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \mathbf{R}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}).$$

It turns out that we can tease out an estimate for  $\boldsymbol{\mu}$  from the expression above that is independent of  $\mathbf{R}$ . The maximizer of the second term above is the same as solving the minimization program<sup>4</sup>

$$\underset{\boldsymbol{\mu} \in \mathbb{R}^D}{\text{minimize}} \quad \sum_{n=1}^N \|\mathbf{R}^{-1/2} (\mathbf{x}_n - \boldsymbol{\mu})\|_2^2.$$

---

<sup>4</sup>We know that the expression  $\mathbf{R}^{-1/2}$  makes sense thanks to the fact that  $\mathbf{R} \in \mathcal{S}_{++}^D$  and the Sylvester theorem.

This is a least-squares problem equivalent to  $\text{minimize}_{\boldsymbol{\mu}} \|\mathbf{y} - \mathbf{A}\boldsymbol{\mu}\|_2^2$  with

$$\mathbf{A} = \begin{bmatrix} \mathbf{R}^{-1/2} \\ \mathbf{R}^{-1/2} \\ \vdots \\ \mathbf{R}^{-1/2} \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} \mathbf{R}^{-1/2} \mathbf{x}_1 \\ \mathbf{R}^{-1/2} \mathbf{x}_2 \\ \vdots \\ \mathbf{R}^{-1/2} \mathbf{x}_N \end{bmatrix}.$$

Thus

$$\hat{\boldsymbol{\mu}}_{\text{mle}} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n,$$

and so the estimator for the mean does not depend on the covariance. We can now estimate  $\mathbf{R}$  as the solution to

$$\underset{\mathbf{R} \in \mathcal{S}_{++}^D}{\text{maximize}} \quad \frac{N}{2} \log \det \mathbf{R}^{-1} - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \hat{\boldsymbol{\mu}})^T \mathbf{R}^{-1} (\mathbf{x}_n - \hat{\boldsymbol{\mu}}).$$

Using the easily checked fact that  $\mathbf{w}^T \mathbf{S} \mathbf{w} = \text{trace}(\mathbf{S} \mathbf{w} \mathbf{w}^T)$  for vectors  $\mathbf{w}$  and sym+def matrices  $\mathbf{S}$ , and the fact that trace is a linear operator, and the fact that an inverse of a sym+def matrix is again sym+def, this is equivalent to

$$\underset{\mathbf{S} \in \mathcal{S}_{++}^D}{\text{maximize}} \quad \log \det \mathbf{S} - \text{trace}(\mathbf{S} \boldsymbol{\Sigma}),$$

where  $\boldsymbol{\Sigma}$  is the sample covariance,

$$\boldsymbol{\Sigma} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \hat{\boldsymbol{\mu}})(\mathbf{x}_n - \hat{\boldsymbol{\mu}})^T.$$

It is a fact that  $\log \det$  is concave in its matrix argument  $\mathbf{S}$  over  $\mathcal{S}_{++}^D$ , and its gradient is

$$\nabla \log \det \mathbf{S} = \mathbf{S}^{-1}.$$

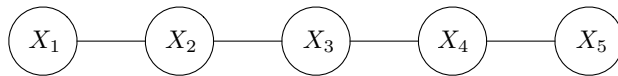
The function  $\text{trace}(\mathbf{S}\mathbf{\Sigma})$  is linear in  $\mathbf{S}$ , and hence both concave and convex; its gradient is

$$\nabla \text{trace}(\mathbf{S}\mathbf{\Sigma}) = \mathbf{\Sigma}.$$

So setting the gradient equal to zero yields:

$$\hat{\mathbf{S}} = \mathbf{\Sigma}^{-1}, \quad \text{or} \quad \hat{\mathbf{R}}_{\text{mle}} = \mathbf{\Sigma} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \hat{\boldsymbol{\mu}}_{\text{mle}})(\mathbf{x}_n - \hat{\boldsymbol{\mu}}_{\text{mle}})^{\text{T}}.$$

**Example:** Suppose that we have some information about the covariance matrix  $\mathbf{R}$ . One for this information could take is the conditional independence structure. For instance, if this were quantified with the graph



we would know that certain entries in the inverse covariance matrix were zero. In this case, we would set up a constrained optimization program

$$\begin{aligned} & \underset{\mathbf{S} \in \mathcal{S}_{++}^5}{\text{maximize}} \quad \log \det \mathbf{S} - \text{trace}(\mathbf{S}\mathbf{\Sigma}) \quad \text{subject to} \\ & S[1, 3] = 0, \quad S[1, 4] = 0, \quad S[1, 5] = 0, \quad S[2, 4] = 0, \quad S[2, 5] = 0 \\ & S[3, 1] = 0, \quad S[3, 5] = 0, \quad S[4, 1] = 0, \quad S[4, 2] = 0. \end{aligned}$$

This program does not have an explicit solution, but it is a concave program with linear constraints, so there is an established methodology for solving it.



**Example:** (This will be on the homework next week.) Let  $\mathbf{A}$  be a fixed  $M \times N$  matrix with full column rank. Suppose we observe

$$Y = \mathbf{A}\boldsymbol{\theta} + Z, \quad Z \sim \text{Normal}(\mathbf{0}, \boldsymbol{\Sigma}).$$

Given  $Y = \mathbf{y} \in \mathbb{R}^M$ , find the MLE for  $\boldsymbol{\theta} \in \mathbb{R}^N$ .

**Example:** (This will also be on the homework next week.) Let  $Z[1], \dots, Z[N]$  be a sequence of independent Gaussian random variables with mean 0 and variance 1. You observe the random vector  $X$  in  $\mathbb{R}^N$  that is generated through the autoregressive process

$$X[k] = \begin{cases} Z[1], & k = 1 \\ aX[k-1] + Z[k], & k > 1. \end{cases}$$

Given  $X = \mathbf{x}$ , find the MLE for  $a \in \mathbb{R}$ .

**Example:** Suppose that  $X$  is a two-sided Laplacian random variable, meaning that

$$f_X(x; \theta) = \frac{1}{2}e^{-|x-\theta|}.$$

We observe  $N$  independent realizations of  $X$ ,  $X_1 = x_1, X_2 = x_2, \dots, X_N = x_N$ . What is the MLE for  $\theta$ ?