

Topic 4: Parameter Estimation

4.1 Introduction

In many applications we observe data x drawn according to a distribution p with an unknown parameter θ^* , and we want to *learn* or *estimate* such parameter. We write this as

$$x \sim p(x|\theta^*), \quad \theta^* \in \Theta.$$

In other words, we want to find the parameter θ^* that generated our data x .

Example 4.1 (Psychokinesis). Imagine controlling things with your mind. For example, skipping a song by just thinking about it. One way to do this is by putting D sensors in your head (in your headphones, for example). These sensors would record small voltages generated by your brain, store them in a vector $\mathbf{x} \in \mathbb{R}^D$, and send them to a machine (phone, computer, server, etc.). The machine should interpret \mathbf{x} and skip the song if that is what you thought about (see Figure 4.1).

We can expect that whenever you think of skipping a song, \mathbf{x} will be composed of an *skip* signal $\boldsymbol{\mu}_\gg^* \in \mathbb{R}^D$ plus other stuff, i.e.,

$$\mathbf{x} = \boldsymbol{\mu}_\gg^* + \boldsymbol{\eta},$$

where $\boldsymbol{\eta} \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$ represents noise. If we knew $\boldsymbol{\mu}_\gg^*$, we could setup a hypothesis test, similar to Example 3.6:

$$\begin{aligned} H_0 : \mathbf{x} &\sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}) && \Rightarrow \text{do nothing,} \\ H_1 : \mathbf{x} &\sim \mathcal{N}(\boldsymbol{\mu}_\gg^*, \sigma^2 \mathbf{I}) && \Rightarrow \text{skip song.} \end{aligned}$$

However, we do not know $\boldsymbol{\mu}_\gg^*$, and hence we want to *learn* it, or more precisely, *estimate* it. To this end, we can put the sensors in your head, ask you to *think* of skipping a song, and record the response vector \mathbf{x} . We can repeat this experiment N times to obtain i.i.d. samples $\{\mathbf{x}_i\}_{i=1}^N$ according to $\mathcal{N}(\boldsymbol{\mu}_\gg^*, \sigma^2 \mathbf{I})$. We don't know $\boldsymbol{\mu}_\gg^*$, but we know that $\{\mathbf{x}_i\}_{i=1}^N$ were generated according to $\mathcal{N}(\boldsymbol{\mu}_\gg^*, \sigma^2 \mathbf{I})$, so the idea is to find the $\boldsymbol{\mu}_\gg^*$ that most likely generated $\{\mathbf{x}_i\}_{i=1}^N$.

In this example, $p(\mathbf{x}|\theta^*)$ is the gaussian probability density function with parameters $\theta^* = \{\boldsymbol{\mu}_\gg^*, \sigma^2\}$, and we want to estimate $\boldsymbol{\mu}_\gg^*$.

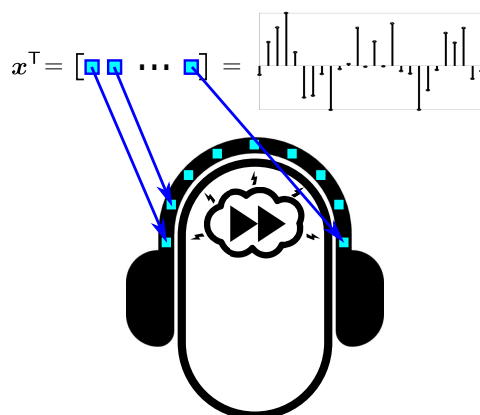


Figure 4.1: Sensors record small voltages generated by your brain and store them in a *signal* vector $\mathbf{x} \in \mathbb{R}^D$. A machine (phone, computer, server, etc.) should interpret \mathbf{x} and *skip* the song if that is what you thought about. See Example 4.1.

Remark 4.1 (Notation). In general, we will use $*$ to denote the *true* parameter that we want to estimate (e.g., θ^*), and $\hat{\cdot}$ to denote an estimator (e.g., $\hat{\theta}$).

4.2 Maximum Likelihood

Again, suppose we are given data $x \sim p(x|\theta^*)$, and we want to estimate θ^* . One way to do this is by finding the parameter $\theta \in \Theta$ that *most likely* generated our data x .

Definition 4.1 (Maximum likelihood estimator (MLE)).

$$\hat{\theta}_{ML} := \arg \max_{\theta \in \Theta} p(x|\theta).$$

Remark 4.2 (Likelihood). The term *likelihood* is often a source of confusion. Recall that in an estimation problem we are given an *instance* of a random variable, i.e., we observe data $x = \mathbf{x}$ drawn according to some probability distribution $p(x|\theta^*)$, and we want to find a parameter $\theta \in \Theta$ that *likely* generated x . The *likelihood* is nothing more than the *probability* that some θ is the parameter that generated x . In other words, the likelihood is $p(x|\theta)$ [evaluated at $x = \mathbf{x}$], and thinking of θ as the variable.

Example 4.2. Let $x_1, \dots, x_N \stackrel{iid}{\sim} \mathcal{N}(\mu^*, \sigma^2)$ with σ^2 known and μ^* unknown. To find $\hat{\mu}_{ML}$ we use elemental techniques from optimization: take the derivative, set to zero, and solve for the desired

parameter. First observe that since $p > 0$,

$$\begin{aligned}\hat{\mu}_{ML} &= \arg \max_{\mu \in \mathbb{R}} p(x_1, \dots, x_N | \mu) = \arg \max_{\mu \in \mathbb{R}} \log p(x_1, \dots, x_N | \mu) = \arg \max_{\mu \in \mathbb{R}} \log \left(\prod_{i=1}^N p(x_i | \mu) \right) \\ &= \arg \max_{\mu \in \mathbb{R}} \sum_{i=1}^N \log p(x_i | \mu) = \arg \max_{\mu \in \mathbb{R}} \sum_{i=1}^N \log \left(\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2} \left(\frac{x_i - \mu}{\sigma} \right)^2} \right) \\ &= \arg \max_{\mu \in \mathbb{R}} \underbrace{\sum_{i=1}^N \log \left(\frac{1}{\sqrt{2\pi}\sigma} \right)}_{\text{constant}} + \sum_{i=1}^N \log \left(e^{-\frac{1}{2} \left(\frac{x_i - \mu}{\sigma} \right)^2} \right) = \arg \max_{\mu \in \mathbb{R}} - \sum_{i=1}^N \frac{1}{2} \left(\frac{x_i - \mu}{\sigma} \right)^2.\end{aligned}$$

Taking derivative with respect to (w.r.t.) μ we have

$$-\frac{\partial}{\partial \mu} \sum_{i=1}^N \frac{1}{2} \left(\frac{x_i - \mu}{\sigma} \right)^2 = \sum_{i=1}^N \frac{x_i - \mu}{\sigma}.$$

Setting to zero and solving for μ we obtain:

$$\hat{\mu}_{ML} = \frac{1}{N} \sum_{i=1}^N x_i.$$

4.3 Bias and Variance

Notice that an estimator $\hat{\theta}$ is a function of the observed data $x = \mathbf{x}$. But since x is a random variable, our estimator will also have some variability. Depending on the particular x that we observe, sometimes our estimator may be better or worse. Hence it is often useful to use the mean (bias) and variance as a measures of the stability of our estimator.

Definition 4.2 (Bias, variance). The *bias* and *variance* of an estimator $\hat{\theta} \in \mathbb{R}^K$ is defined as

$$\begin{aligned}\text{bias}(\hat{\theta}) &:= \theta^* - \mathbb{E}[\hat{\theta}], \\ \text{var}(\hat{\theta}) &:= \mathbb{E} \left[\|\hat{\theta} - \mathbb{E}[\hat{\theta}]\|_2^2 \right],\end{aligned}$$

where the expectations are w.r.t. the *true* probability distribution of x , i.e., w.r.t. θ^* . If $\text{bias}(\hat{\theta}) = \mathbf{0}$ we say that $\hat{\theta}$ is *unbiased*.

Intuitively, the bias tells us how close we will get to the true parameter on average. The variance tells us how far we expect to be from that average. If an estimator is unbiased but has large variance, it is likely to be inaccurate for our particular sample $x = \mathbf{x}$. Similarly, if an estimator has low variance but is biased, we expect it to be consistently inaccurate (see Figure 4.2). Ideally, we want an unbiased estimator with low variance, meaning that we expect it to be close to the true parameter.

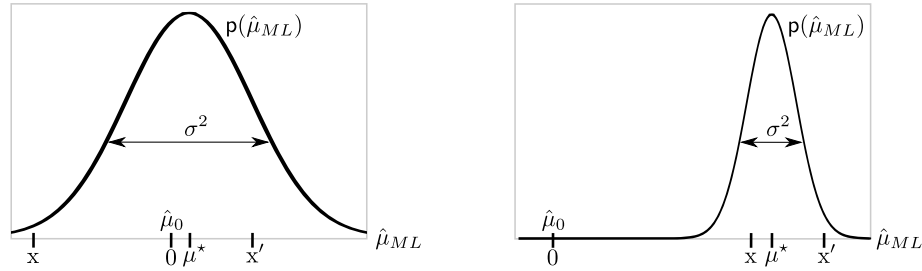


Figure 4.2: An estimator is generally a function of the data. For instance, consider $\hat{\mu}_{ML} = x$ from Example 4.3. Since $x \sim \mathcal{N}(\mu^*, \sigma^2)$, so does $\hat{\mu}_{ML}$. For example, $\hat{\mu}_{ML}$ could take the value x or x' . **Left:** If σ is large, we can expect $\hat{\mu}_{ML}$ to be far from μ^* . In this case the *constant* estimator $\hat{\mu}_0$ (which has zero variance) might be better than $\hat{\mu}_{ML}$. **Right:** If σ is small, we can expect $\hat{\mu}_{ML}$ to be close to μ^* . In this case, $\hat{\mu}_{ML}$ would be better than $\hat{\mu}_0$.

Example 4.3. Let $x \sim \mathcal{N}(\mu^*, \sigma^2)$, and consider two estimators: the constant estimator $\hat{\mu}_0 = 0$ and the MLE $\hat{\mu}_{ML} = x$. Then

$$\begin{aligned} \text{bias}(\hat{\mu}_0) &= \mu^* - \mathbb{E}[\hat{\mu}_0] = \mu^* - \mathbb{E}[0] = \mu^* - 0 = \mu^*, \\ \text{var}(\hat{\mu}_0) &= \mathbb{E}[(\hat{\mu}_0 - \mathbb{E}[\hat{\mu}_0])^2] = \mathbb{E}[(0 - 0)^2] = 0, \\ \text{bias}(\hat{\mu}_{ML}) &= \mu^* - \mathbb{E}[\hat{\mu}_{ML}] = \mu^* - \mathbb{E}[x] = \mu^* - \mu^* = 0, \\ \text{var}(\hat{\mu}_{ML}) &= \mathbb{E}[(\hat{\mu}_{ML} - \mathbb{E}[\hat{\mu}_{ML}])^2] = \mathbb{E}[(x - \mu^*)^2] = \text{var}(x) = \sigma^2. \end{aligned}$$

If μ^* is small, $\hat{\mu}_0$ will always be close to μ^* , regardless of the particular x that we observe. In addition, if σ is big, depending on the particular x that we observe, $\hat{\mu}_{ML}$ could be far from μ^* . In this case $\hat{\mu}_0$ would be better. For similar arguments, if μ^* is big and σ is small, $\hat{\mu}_{ML}$ would be better. See Figure 4.2 to build some intuition.

4.4 Error/Loss and Risk

Example 4.3 shows that some estimators are more accurate than others. As we saw, the bias and variance are good indicators of an estimator's performance. As we will see, these are particular cases of a more general way to measure how good an estimator is: the *risk*. To define it, we first need to introduce the concept of *error*, often also known as *loss*.

Definition 4.3 (Error/loss). An *error* or *loss* function is a mapping $\ell : \Theta \rightarrow \mathbb{R}_+$ that measures the distance between θ^* and $\hat{\theta}$.

Example 4.4. Here are some common loss functions, with $\Theta = \mathbb{R}^K$:

- 0/1 loss:

$$\ell_0(\hat{\theta}) := \mathbb{1}_{\{\hat{\theta} \neq \theta^*\}} = \begin{cases} 0 & \text{if } \hat{\theta} = \theta^*, \\ 1 & \text{if } \hat{\theta} \neq \theta^*. \end{cases}$$

- ℓ_1 loss (absolute error):

$$\ell_1(\hat{\theta}) := \|\hat{\theta} - \theta^*\|_1 = \sum_{k=1}^K |\hat{\theta}_k - \theta_k^*|.$$

- ℓ_2 loss (squared error):

$$\ell_2(\hat{\theta}) := \|\hat{\theta} - \theta^*\|_2^2 = (\hat{\theta} - \theta^*)^\top (\hat{\theta} - \theta^*) = \sum_{k=1}^K (\hat{\theta}_k - \theta_k^*)^2.$$

- ℓ_l log-probability:

$$\ell_l(\hat{\theta}) := -\log p(\mathbf{x}|\hat{\theta}),$$

which measures the distance from θ^* because \mathbf{x} is truly distributed $p(\mathbf{x}|\theta^*)$. We can interpret $\hat{\theta}$ as $\arg \min_{\theta} \ell_l(\theta)$

Recall that an estimator $\hat{\theta}$ is a function of the observed data $x = \mathbf{x}$. But since x is a random variable, our estimator will also have some variability. Depending on the particular \mathbf{x} that we observe, sometimes our estimate may be better or worse. Hence a better indicator of our estimator's performance is the *expected* error/loss, also known as *risk*.

Definition 4.4 (Risk).

$$\text{Risk}(\hat{\theta}) := \mathbb{E} [\ell(\hat{\theta})].$$

Example 4.5. Continuing with Example 4.4,

$$\begin{aligned} \text{Risk}_0(\hat{\theta}) &:= \mathbb{E} [\ell_0(\hat{\theta})] = \mathbb{E} [\mathbb{1}_{\{\hat{\theta} \neq \theta^*\}}] = \mathbb{P}(\hat{\theta} \neq \theta^*), \\ \text{Risk}_1(\hat{\theta}) &:= \mathbb{E} [\ell_1(\hat{\theta})] = \mathbb{E} [\|\hat{\theta} - \theta^*\|_1], \\ \text{Risk}_2(\hat{\theta}) &:= \mathbb{E} [\ell_2(\hat{\theta})] = \mathbb{E} [\|\hat{\theta} - \theta^*\|_2^2] =: \text{MSE}(\hat{\theta}), \\ \text{Risk}_l(\hat{\theta}) &:= \mathbb{E} [\ell_l(\hat{\theta})] = \mathbb{E} [-\log p(\mathbf{x}|\hat{\theta})], \end{aligned}$$

where MSE stands for *mean squared error*.

The next proposition shows that the bias and variance have a tight relation with a particular kind of risk, namely with $\text{Risk}_2 = \text{MSE}$.

Proposition 4.1 (MSE = bias² + var).

$$\text{MSE}(\hat{\theta}) = \text{bias}(\hat{\theta})^T \text{bias}(\hat{\theta}) + \text{var}(\hat{\theta}).$$

Proof. Write

$$\begin{aligned} \text{MSE}(\hat{\theta}) &= \mathbb{E} \left[\|\hat{\theta} - \theta^*\|_2^2 \right] = \mathbb{E} \left[\underbrace{\|(\hat{\theta} - \mathbb{E}[\hat{\theta}]) + (\mathbb{E}[\hat{\theta}] - \theta^*)\|_2^2}_0 \right] \\ &= \mathbb{E} \left[\left((\hat{\theta} - \mathbb{E}[\hat{\theta}]) + (\mathbb{E}[\hat{\theta}] - \theta^*) \right)^T \left((\hat{\theta} - \mathbb{E}[\hat{\theta}]) + (\mathbb{E}[\hat{\theta}] - \theta^*) \right) \right] \\ &= \mathbb{E} \left[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^T (\hat{\theta} - \mathbb{E}[\hat{\theta}]) + 2(\hat{\theta} - \mathbb{E}[\hat{\theta}])^T (\mathbb{E}[\hat{\theta}] - \theta^*) + (\mathbb{E}[\hat{\theta}] - \theta^*)^T (\mathbb{E}[\hat{\theta}] - \theta^*) \right] \\ &= \mathbb{E} \left[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^T (\hat{\theta} - \mathbb{E}[\hat{\theta}]) \right] + 2\mathbb{E} \left[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^T \underbrace{(\mathbb{E}[\hat{\theta}] - \theta^*)}_{\text{constant}} \right] + \mathbb{E} \left[(\mathbb{E}[\hat{\theta}] - \theta^*)^T \underbrace{(\mathbb{E}[\hat{\theta}] - \theta^*)}_{\text{constant}} \right] \\ &= \mathbb{E} \left[\|\hat{\theta} - \mathbb{E}[\hat{\theta}]\|_2^2 \right] + 2\mathbb{E} \left[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^T (\mathbb{E}[\hat{\theta}] - \theta^*) \right] + \mathbb{E} \left[(\mathbb{E}[\hat{\theta}] - \theta^*)^T (\mathbb{E}[\hat{\theta}] - \theta^*) \right] \\ &= \underbrace{\mathbb{E} \left[\|\hat{\theta} - \mathbb{E}[\hat{\theta}]\|_2^2 \right]}_{\text{var}(\hat{\theta})} + 2 \underbrace{(\mathbb{E}[\hat{\theta}] - \mathbb{E}[\hat{\theta}])^T}_{\mathbf{0}} (\mathbb{E}[\hat{\theta}] - \theta^*) + \underbrace{(\mathbb{E}[\hat{\theta}] - \theta^*)^T}_{\text{bias}(\hat{\theta})} \underbrace{(\mathbb{E}[\hat{\theta}] - \theta^*)}_{\text{bias}(\hat{\theta})}. \end{aligned}$$

□

Example 4.6. Consider Example 4.1. Then

$$\begin{aligned} \hat{\mu}_{ML} &= \arg \max_{\mu \in \mathbb{R}^k} p(\mathbf{x}_1, \dots, \mathbf{x}_N | \mu) = \arg \max_{\mu \in \mathbb{R}^D} \prod_{i=1}^N p(\mathbf{x}_i | \mu) = \arg \max_{\mu \in \mathbb{R}^D} \log \prod_{i=1}^N p(\mathbf{x}_i | \mu) \\ &= \arg \max_{\mu \in \mathbb{R}^D} \sum_{i=1}^N \log p(\mathbf{x}_i | \mu) = \arg \max_{\mu \in \mathbb{R}^D} \sum_{i=1}^N \log \left(\frac{1}{(\sqrt{2\pi}\sigma)^D} e^{-\frac{1}{2\sigma^2} (\mathbf{x}_i - \mu)^T (\mathbf{x}_i - \mu)} \right) \\ &= \arg \max_{\mu \in \mathbb{R}^D} \underbrace{\sum_{i=1}^N \log \frac{1}{(\sqrt{2\pi}\sigma)^D}}_{\text{constant}} - \sum_{i=1}^N \underbrace{\frac{1}{2\sigma^2}}_{\text{constant}} (\mathbf{x}_i - \mu)^T (\mathbf{x}_i - \mu) \\ &= \arg \max_{\mu \in \mathbb{R}^D} - \sum_{i=1}^N (\mathbf{x}_i - \mu)^T (\mathbf{x}_i - \mu) = \arg \max_{\mu \in \mathbb{R}^D} - \sum_{i=1}^N (\mathbf{x}_i^T \mathbf{x}_i - 2\mathbf{x}_i^T \mu + \mu^T \mu) \\ &= \arg \max_{\mu \in \mathbb{R}^D} \sum_{i=1}^N (2\mathbf{x}_i^T \mu - \mu^T \mu). \end{aligned}$$

Now we use our usual tricks: take derivative w.r.t. μ , set to zero and solve for μ . To learn more about how to take derivatives w.r.t. vectors and matrices see *Old and new matrix algebra useful for statistics*

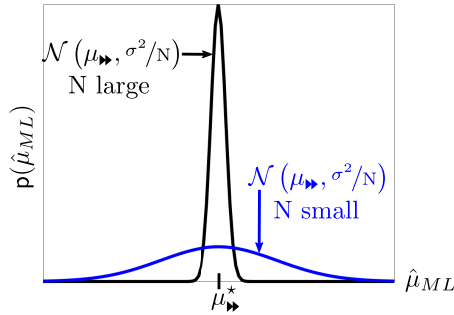


Figure 4.3: Example 4.6 shows that $\hat{\mu}_{ML} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$. This implies that $\hat{\mu}_{ML} \sim \mathcal{N}(\mu_{\star}^*, \frac{\sigma^2}{N} \mathbf{I})$. As N grows, $\text{var}(\hat{\mu}_{ML})$ shrinks, and we expect $\hat{\mu}_{ML}$ to be closer to μ_{\star}^* .

by Thomas P. Minka.

$$\frac{\partial}{\partial \mu} \sum_{i=1}^N (2\mathbf{x}_i^T \mu - \mu^T \mu) = \sum_{i=1}^N \frac{\partial}{\partial \mu} (2\mathbf{x}_i^T \mu - \mu^T \mu) = \sum_{i=1}^N (2\mathbf{x}_i - 2\mu) = 2 \sum_{i=1}^N \mathbf{x}_i - 2N\mu = \mathbf{0}.$$

It follows that

$$\hat{\mu}_{ML} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i.$$

Next notice that

$$\mathbb{E}[\hat{\mu}_{ML}] = \mathbb{E}\left[\frac{1}{N} \sum_{i=1}^N \mathbf{x}_i\right] = \frac{1}{N} \sum_{i=1}^N \mathbb{E}[\mathbf{x}_i] = \frac{1}{N} \sum_{i=1}^N \mu_{\star}^* = \mu_{\star}^*.$$

Hence $\text{bias}(\hat{\mu}_{ML}) = \mathbf{0}$. On the other hand,

$$\text{var}(\hat{\mu}_{ML}) = \text{var}\left(\frac{1}{N} \sum_{i=1}^N \mathbf{x}_i\right) = \frac{1}{N^2} \sum_{i=1}^N \text{var}(\mathbf{x}_i) = \frac{1}{N^2} \sum_{i=1}^N D\sigma^2 = \frac{D\sigma^2}{N}.$$

Hence

$$\text{MSE}(\hat{\mu}_{ML}) = \frac{D\sigma^2}{N}.$$

4.5 Excess Risk/Regret and KL Divergence

Each sample x has a certain probability of being observed, namely $p(x|\theta^*)$. An other way to analyze how good is an estimator is by comparing $p(x|\theta^*)$ against the probability of the same sample under the estimator, i.e., $p(x|\hat{\theta})$. This is the idea behind the concept of *excess risk* or *regret*: to compare $p(x|\theta^*)$ against $p(x|\hat{\theta})$.

Equivalently, we can compare $\log p(x|\theta^*)$ against $\log p(x|\hat{\theta})$. This is done thinking about cases where we have i.i.d. samples x_1, \dots, x_N , so that $\log p(\mathbf{x}|\cdot)$ factors nicely into $\sum_{i=1}^N \log p(x_i|\cdot)$. Finally, since samples are random variables, it is better to compare the expected values of these quantities, which is precisely what

the *regret* compares.

Definition 4.5 (Excess risk/regret).

$$\text{Regret}(\hat{\theta}) := \mathbb{E}[\log p(\mathbf{x}|\theta^*)] - \mathbb{E}[\log p(\mathbf{x}|\hat{\theta})].$$

The next proposition shows that the regret is the *Kullback-Leibler divergence* (KL divergence).

Proposition 4.2 (Regret = KL divergence).

$$\text{Regret}(\hat{\theta}) = D\left(p(\mathbf{x}|\theta^*) \parallel p(\mathbf{x}|\hat{\theta})\right) := \int p(\mathbf{x}|\theta^*) \log \frac{p(\mathbf{x}|\theta^*)}{p(\mathbf{x}|\hat{\theta})} d\mathbf{x}.$$

Proof.

$$\mathbb{E}[\log p(\mathbf{x}|\theta^*)] - \mathbb{E}[\log p(\mathbf{x}|\hat{\theta})] = \mathbb{E}[\log p(\mathbf{x}|\theta^*) - \log p(\mathbf{x}|\hat{\theta})] = \mathbb{E}\left[\log \frac{p(\mathbf{x}|\theta^*)}{p(\mathbf{x}|\hat{\theta})}\right].$$

□

D is a well-studied function. Proposition 4.2 allows us to translate from **Regret** to D and use all the machinery and results known for D.

Example 4.7. Suppose $x_1, \dots, x_N \stackrel{iid}{\sim} \mathcal{N}(\mu^*, \sigma^{*2})$. Then the regret of any estimator $\hat{\theta} = \{\hat{\mu}, \hat{\sigma}\}$ is the well-known KL divergence of two gaussian distributions:

$$\text{Regret}(\hat{\theta}) = \log\left(\frac{\hat{\sigma}}{\sigma^*}\right) + \frac{\sigma^{*2} + (\mu^* - \hat{\mu})^2}{2\hat{\sigma}^2} - \frac{1}{2}.$$

Here is another example of the usefulness of Proposition 4.2:

Proposition 4.3 (Information inequality).

$$D\left(p(\mathbf{x}|\theta^*) \parallel p(\mathbf{x}|\hat{\theta})\right) \geq 0$$

with equality if and only if $p(\mathbf{x}|\theta^*) = p(\mathbf{x}|\hat{\theta})$.

Proof. See Theorem 2.6.3 in *Elements of Information Theory* by Cover and Thomas, second edition. □

Proposition 4.3 implies that (as intuition suggests) in expectation, the true distribution $\mathbf{p}(\mathbf{x}|\theta)$ is more likely to have produced \mathbf{x} than *any* other distribution $\mathbf{p}(\mathbf{x}|\hat{\theta})$. See *Elements of Information Theory* by Cover and Thomas to learn more about KL-divergence, entropy and other useful information theory functions and results.

4.6 Estimators of Functions

Often we want to estimate a function g of θ^* .

Example 4.8. Suppose we observe $x_1, \dots, x_N \stackrel{iid}{\sim} \text{Poisson}(\lambda^*)$, with $\lambda^* \in \mathbb{R}$ unknown. We want to know the probability γ that a new sample x is larger than λ^* . Here $\gamma = g(\lambda^*) = \mathbb{P}(x \geq \lambda^*)$.

The next theorem shows that the MLE of a function is the function of the MLE.

Theorem 4.1 (Invariance of the MLE). Let $x_1, \dots, x_N \stackrel{iid}{\sim} \mathbf{p}(\mathbf{x}|\theta^*)$. Let $\gamma^* = g(\theta^*)$ for some function $g: \Theta \rightarrow \Gamma$. Then the MLE of γ^* , defined as

$$\hat{\gamma}_{ML} := \arg \max_{\gamma \in \Gamma} \left(\max_{\theta \in g^{-1}(\gamma)} \mathbf{p}(\mathbf{x}|\theta) \right),$$

is given by $\hat{\gamma}_{ML} = g(\hat{\theta}_{ML})$. Here g^{-1} denotes the inverse image of g , i.e., $g^{-1}(\gamma) = \{\theta \in \Theta : g(\theta) = \gamma\}$.

Think of γ^* as a parameter (that is a function of θ^*) that we also want to estimate.

Proof. Γ is defined as the range of g , so even if g is not one-to-one, the sets $\{g^{-1}(\gamma)\}_{\gamma \in \Gamma}$ form a partition of Θ (see Figure 4.4 to build some intuition). Therefore,

$$\bigcup_{\gamma \in \Gamma} g^{-1}(\gamma) = \Theta,$$

which in turns implies

$$\max_{\gamma \in \Gamma} \left(\max_{\theta \in g^{-1}(\gamma)} \mathbf{p}(\mathbf{x}|\theta) \right) = \max_{\theta \in \Theta} \mathbf{p}(\mathbf{x}|\theta).$$

□

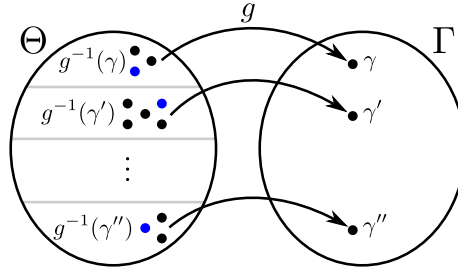


Figure 4.4: In Theorem 4.1, since $\{g^{-1}(\gamma)\}_{\gamma \in \Gamma}$ forms a partition of Θ , taking the max over Θ is the same as first taking the max of over each individual $g^{-1}(\gamma)$, highlighted in blue, and then taking the max over these maximums. More precisely, $\max_{\theta \in \Theta} \cdot = \max_{\gamma \in \Gamma} (\max_{\theta \in g^{-1}(\gamma)} \cdot)$.

Example 4.9. Continuing with Example 4.8, we can use Theorem 4.1 to first compute $\hat{\lambda}_{ML}$ and then obtain $\hat{\gamma}_{ML} = g(\hat{\lambda}_{ML}) = P(X > \hat{\lambda}_{ML})$.

$$\begin{aligned} \hat{\lambda}_{ML} &= \arg \max_{\lambda \in \mathbb{R}} \prod_{i=1}^N P(x = x_i | \lambda) = \arg \max_{\lambda \in \mathbb{R}} \log \prod_{i=1}^N P(x = x_i | \lambda) = \arg \max_{\lambda \in \mathbb{R}} \sum_{i=1}^N \log P(x = x_i | \lambda) \\ &= \arg \max_{\lambda \in \mathbb{R}} \sum_{i=1}^N \log \frac{\lambda^{x_i} e^{-\lambda}}{x_i!} = \arg \max_{\lambda \in \mathbb{R}} \sum_{i=1}^N \log \lambda^{x_i} - \sum_{i=1}^N \lambda - \underbrace{\sum_{i=1}^N \log x_i!}_{\text{constant}} \\ &= \arg \max_{\lambda \in \mathbb{R}} \sum_{i=1}^N x_i \log \lambda - N\lambda = \arg \max_{\lambda \in \mathbb{R}} \log \lambda \sum_{i=1}^N x_i - N\lambda. \end{aligned}$$

It follows that

$$\frac{\partial}{\partial \lambda} \left(\log \lambda \sum_{i=1}^N x_i - N\lambda \right) = \frac{1}{\lambda} \sum_{i=1}^N x_i - N.$$

Setting to zero and solving for λ we have:

$$\hat{\lambda}_{ML} = \frac{1}{N} \sum_{i=1}^N x_i.$$

It follows by Theorem 4.1 that

$$\hat{\gamma}_{ML} = g(\hat{\lambda}_{ML}) = P(x > \hat{\lambda}_{ML}) = P\left(x > \frac{1}{N} \sum_{i=1}^N x_i\right).$$

4.7 Asymptotics

Example 4.6 shows that the larger N , i.e., the more *training* samples we have, the closer $\hat{\mu}_{ML}$ will be to the true *skip* signal μ_{\star}^* that we want to estimate (see Figure 4.3). This is a desirable property in an estimator called *consistency*.

Definition 4.6 (Consistency). We say $\hat{\theta}$ is *consistent* if $\hat{\theta} \rightarrow \theta^*$ as $N \rightarrow \infty$. More precisely, $\hat{\theta}$ is *consistent w.r.t. Risk* if $\text{Risk}(\hat{\theta}) \rightarrow 0$ as $N \rightarrow \infty$.

Example 4.10. $\hat{\mu}_{ML}$ from Example 4.6 is consistent w.r.t. Risk = MSE, because

$$\text{MSE}(\hat{\mu}_{ML}) = \text{E} [\|\hat{\mu}_{ML} - \mu_{\star}^*\|_2^2] = \frac{D\sigma^2}{N} \xrightarrow{N \rightarrow \infty} 0,$$

i.e., $\hat{\mu}_{ML} \rightarrow \mu_{\star}^*$ as $N \rightarrow \infty$.

As mentioned before, an estimator is a function of the data we observe, and hence it is a random variable. The next theorem states that if we observe a large number N of i.i.d. samples, then regardless of the distribution of these data, $\hat{\theta}_{ML}$ will be distributed normal, centered at θ^* , and with variance decreasing with N , implying consistency.

Theorem 4.2 (Asymptotic distribution of the MLE). Let $\mathbf{x}_1, \dots, \mathbf{x}_N \stackrel{iid}{\sim} p(\mathbf{x}|\theta^*)$, with $\theta^* \in \mathbb{R}^K$. Let

$$L(\theta) := \sum_{i=1}^N \log p(\mathbf{x}_i|\theta).$$

Suppose $\frac{\partial L(\theta)}{\partial \theta_k}$ and $\frac{\partial^2 L(\theta)}{\partial \theta_k \partial \theta_\ell}$ exist for every $k, \ell \in \{1, \dots, K\}$. Then

$$\hat{\theta}_{ML} \stackrel{N \rightarrow \infty}{\sim} \mathcal{N}\left(\theta^*, \frac{1}{N} \mathbf{I}_{\theta^*}^{-1}\right),$$

where \mathbf{I}_{θ^*} is the *Fisher-information matrix*, whose elements are defined as:

$$[\mathbf{I}_{\theta^*}]_{k\ell} := -\text{E} \left[\frac{\partial^2 L(\theta)}{\partial \theta_k \partial \theta_\ell} \Big|_{\theta=\theta^*} \right].$$

Proof. The proof follows as a consequence of the central limit theorem, which says that if x_1, \dots, x_N are i.i.d. with mean μ and variance σ^2 , then $\frac{1}{\sqrt{N}} \sum_{i=1}^N x_i$ is asymptotically distributed $\mathcal{N}(\mu, \sigma^2)$. For a detailed proof see Theorem 10.1.12 in *Statistical Inference* by George Casella and Roger L. Berger, second edition. \square

Example 4.11 (Exam-type question). We want to find out which of two treatments is more effective. So we give treatment k to n_k people, and record the number x_k of people that react favorably to treatment k ($k = 1, 2$). We can model this as $x_k \sim \text{Binomial}(n_k, p_k^*)$, where x_1 is independent of x_2 . Then the difference $\delta^* = p_1^* - p_2^*$ would be a good indicator of which treatment is more effective, and by how much. To estimate δ^* (which is a function of p_1^* and p_2^*), we can first estimate p_1^* and p_2^* , and

then use Theorem 4.1. First notice that,

$$\hat{p}_k = \arg \max_{p_k \in [0,1]} P(x_k = x_k | p_k) = \arg \max_{p_k \in [0,1]} \log P(x_k = x_k | p_k) = \arg \max_{p_k \in [0,1]} L(p_k).$$

Next we will use our usual tricks: take derivative, set to zero and solve for the desired parameters.

$$\begin{aligned} \frac{\partial L(p_k)}{\partial p_k} &= \frac{\partial}{\partial p_k} \log \left(\binom{n_k}{x_k} p_k^{x_k} (1-p_k)^{n_k-x_k} \right) = \frac{\partial}{\partial p_k} \left(\underbrace{\log \binom{n_k}{x_k}}_{\text{constant}} + \log p_k^{x_k} + \log(1-p_k)^{n_k-x_k} \right) \\ &= \frac{\partial}{\partial p_k} (x_k \log p_k + (n_k - x_k) \log(1-p_k)) = \frac{x_k}{p_k} - \frac{n_k - x_k}{1-p_k}. \end{aligned}$$

Setting to zero,

$$0 = \frac{x_k(1-p_k) - (n_k - x_k)p_k}{p_k(1-p_k)} = x_k - \cancel{x_k p_k} - n_k p_k + \cancel{x_k p_k} = x_k - n_k p_k,$$

and solving for p_k we obtain $\hat{p}_k = \frac{x_k}{n_k}$. It follows by Theorem 4.1 that

$$\hat{\delta}_{ML} = \hat{p}_1 - \hat{p}_2 = \frac{x_1}{n_1} - \frac{x_2}{n_2}.$$

Theorems 4.1 and 4.2 imply that $\hat{\delta}_{ML} \rightarrow \delta^*$ as $n_1, n_2 \rightarrow \infty$ (can you show this?). But we can also use them to quantify how accurate our estimator $\hat{\delta}_{ML}$ will be. For example, if we only test $n_k = 1$ people per treatment, our estimator might not be very reliable. We want to know how many people we need to test to guarantee that $\hat{\delta}_{ML}$ is pretty close to the true δ^* . More precisely, we want

$$P \left(|\hat{\delta}_{ML} - \delta^*| \geq \beta \right) \leq \alpha.$$

Recall that the sum of binomials is only binomial if they share the same p . Hence we do not know the distribution of $\hat{\delta}_{ML}$. Fortunately, we can use Theorem 4.2 to know its asymptotic distribution. To this end, we first need to compute $I_{p_k^*}$, so write:

$$\begin{aligned} \frac{\partial^2 L(p_k)}{\partial^2 p_k} &= \frac{\partial}{\partial p_k} \left(\frac{x_k}{p_k} - \frac{n_k - x_k}{1-p_k} \right) = x_k \frac{\partial p_k^{-1}}{\partial p_k} - (n_k - x_k) \frac{\partial (1-p_k)^{-1}}{\partial p_k} \\ &= -x_k p_k^{-2} - (n_k - x_k)(1-p_k)^{-2} = -\frac{x_k}{p_k^2} - \frac{n_k - x_k}{(1-p_k)^2}. \end{aligned}$$

Hence

$$\begin{aligned} I_{p_k^*} &= -E \left[\frac{\partial^2 L(p_k)}{\partial^2 p_k} \Big|_{p_k=p_k^*} \right] = -E \left[\left(-\frac{x_k}{p_k^2} - \frac{n_k - x_k}{(1-p_k)^2} \right) \Big|_{p_k=p_k^*} \right] \\ &= E \left[\frac{x_k}{p_k^{*2}} + \frac{n_k - x_k}{(1-p_k^*)^2} \right] = \frac{E[x_k]}{p_k^{*2}} + \frac{n_k - E[x_k]}{(1-p_k^*)^2} \\ &= \frac{n_k p_k}{p_k^{*2}} + \frac{n_k - n_k p_k}{(1-p_k^*)^2} = \frac{n_k}{p_k^*} + \frac{n_k(1-p_k)}{(1-p_k^*)^2} \\ &= \frac{n_k}{p_k^*} + \frac{n_k}{1-p_k^*} = n_k \frac{1-p_k^* + p_k^*}{p_k^*(1-p_k^*)} = \frac{n_k}{p_k^*(1-p_k^*)}, \end{aligned}$$

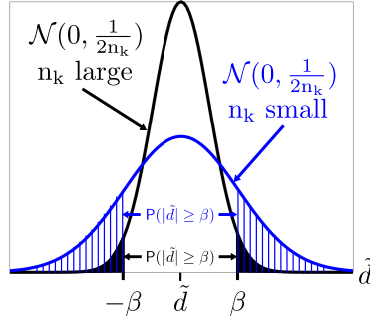


Figure 4.5: In Example 4.11, $\mathbb{P}(|\hat{\delta}_{ML} - \delta^*| \geq \beta) \leq \mathbb{P}(|\tilde{\delta}| \geq \beta)$, where $\tilde{\delta} \sim \mathcal{N}(0, \frac{1}{2n_k})$. As n_k grows, $\text{var}(\tilde{\delta})$ shrinks, and hence so does $\mathbb{P}(|\tilde{\delta}| \geq \beta)$. We want to find the n_k to guarantee that $\mathbb{P}(|\tilde{\delta}| \geq \beta) \leq \alpha$. This will also imply that $\mathbb{P}(|\hat{\delta}_{ML} - \delta^*| \geq \beta) \leq \alpha$.

which implies

$$\mathbf{I}_{\mathbf{p}_k^*}^{-1} = \frac{\mathbf{p}_k^*(1 - \mathbf{p}_k^*)}{n_k}.$$

By Theorem 4.2, $\hat{\mathbf{p}}_k \stackrel{N \rightarrow \infty}{\sim} \mathcal{N}\left(\mathbf{p}_k^*, \frac{1}{N} \frac{\mathbf{p}_k^*(1 - \mathbf{p}_k^*)}{n_k}\right)$, with $N = 1$ (because we only observe one x_1 and one x_2). Since sums of gaussians are gaussians, it follows that

$$\hat{\delta}_{ML} \stackrel{N \rightarrow \infty}{\sim} \mathcal{N}\left(\mathbf{p}_1^* - \mathbf{p}_2^*, \frac{\mathbf{p}_1^*(1 - \mathbf{p}_1^*)}{n_1} + \frac{\mathbf{p}_2^*(1 - \mathbf{p}_2^*)}{n_2}\right).$$

Now suppose for simplicity that we will test the same patients for each treatment, i.e., $n_1 = n_2$. Then this simplifies to

$$\hat{\delta}_{ML} \stackrel{N \rightarrow \infty}{\sim} \mathcal{N}\left(\delta^*, \frac{\mathbf{p}_1^*(1 - \mathbf{p}_1^*) + \mathbf{p}_2^*(1 - \mathbf{p}_2^*)}{n_k}\right).$$

Next observe that since $p(1 - p) \leq 1/4$ for every $p \in [0, 1]$, the asymptotic variance of $\hat{\delta}_{ML}$ is bounded by $\frac{1}{2n_k}$. Letting $\tilde{\delta}$ be a $\mathcal{N}(0, \frac{1}{2n_k})$ distributed random variable,

$$\mathbb{P}\left(|\hat{\delta}_{ML} - \delta^*| \geq \beta\right) \leq \mathbb{P}\left(|\tilde{\delta}| \geq \beta\right) = 2Q_{0, \frac{1}{2n_k}}(\beta) = 2Q(\sqrt{2n_k}\beta).$$

See Figure 4.5 to build some intuition. If we want $2Q(\sqrt{2n_k}\beta) \leq \alpha$, then we need to test

$$n_k \geq \left(\frac{Q^{-1}(\alpha/2)}{\sqrt{2}\beta}\right)^2 \tag{4.1}$$

patients per treatment.

4.8 Experiments

It is often good to run some experiments to verify that our findings are correct. For instance, we can run a simulation of Example 4.11.

Recall that the ultimate goal in Example 4.11 is to determine which of two treatments is better, and by how

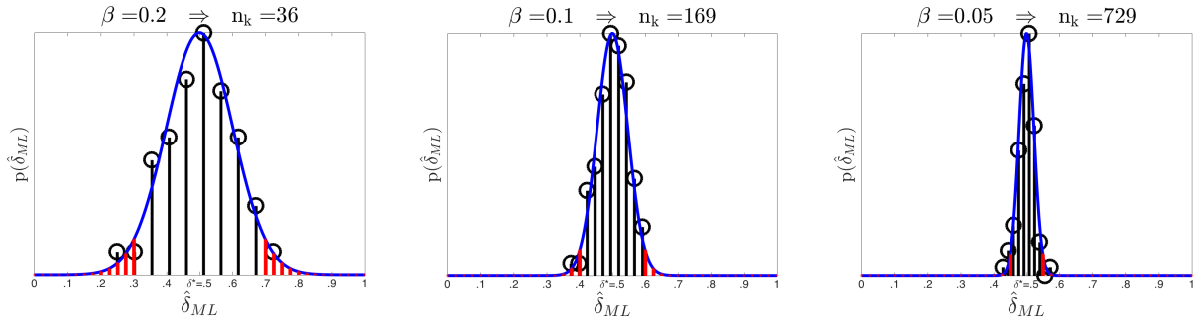


Figure 4.6: Results of the simulation of Example 4.11. In **black** is the histogram of $\hat{\delta}_{ML}$ over $T = 100$ trials. In **blue** is the asymptotic distribution of $\hat{\delta}_{ML}$. In **red** is the probability α that $|\hat{\delta}_{ML} - \delta^*| \geq \beta$. Notice that as n_k grows, the variance of $\hat{\delta}_{ML}$ decreases, and so $\hat{\delta}_{ML}$ gets closer and closer to δ^* . The code for this simulation is in the appendix.

much. This is determined by δ^* , so our goal is to estimate δ^* . To this end, we need to know how many patients n_k we need to test on each treatment to be $(1 - \alpha)\%$ sure that our estimate $\hat{\delta}_{ML}$ will be within β error of δ^* , i.e., we want $P(|\hat{\delta}_{ML} - \delta^*| \geq \beta) \leq \alpha$. Using Theorems 4.1 and 4.2 we found that n_k needs to be as in (4.1). We can run a simulation to verify whether this is accurate as follows.

First generate random vectors $\mathbf{y}_k \in \mathbb{R}^{n_k}$ with i.i.d. Bernoulli(p_k^*) entries. The i^{th} entry in \mathbf{y}_k simulates whether patient i reacted favorably to treatment k . Then x_k is simply the sum of the entries in \mathbf{y}_k . From our results in Example 4.11 we know that $\hat{\delta}_{ML} = \frac{x_1}{n_1} - \frac{x_2}{n_2}$.

We can repeat this experiment T trials, plot a histogram of $\hat{\delta}_{ML}$, and verify whether $\hat{\delta}_{ML}$ truly is within β of δ^* on $(1 - \alpha)\%$ of the trials. Furthermore, we can compare it with the asymptotic distribution obtained by Theorem 4.2. Figure 4.6 shows some results. The code for this simulation is in the appendix.

4.9 Cramer-Rao Lower Bound

Theorem 4.2 tells us that $\hat{\theta}_{ML} \rightarrow \theta^*$ at a rate of $1/\sqrt{N}$ (because $\text{var}(\hat{\theta}_{ML}) = E[(\hat{\theta}_{ML} - \theta^*)^2] = \frac{1}{N} \mathbf{I}_{\theta^*}^{-1}$). Intuitively, this means that we need about N samples to be within $1/\sqrt{N}$ error. This raises the question: is there a better estimator? For example, one that achieves $1/N$ error with the same number of samples? The next theorem shows that this is *not* the case. More precisely, it shows that there exists *no* estimator with a faster convergence rate than $1/\sqrt{N}$.

Theorem 4.3 (Cramer-Rao Lower Bound). Let $\mathbf{x}_1, \dots, \mathbf{x}_N \stackrel{iid}{\sim} p(\mathbf{x}|\theta^*)$, with $\theta^* \in \mathbb{R}^K$. Suppose $\frac{\partial L(\theta)}{\partial \theta_k}$ and $\frac{\partial^2 L(\theta)}{\partial \theta_k \partial \theta_\ell}$ exist for every $k, \ell \in \{1, \dots, K\}$. Let $\hat{\theta}$ be an unbiased estimator of θ^* . Define the *error covariance matrix* as

$$\hat{\mathbf{C}} := E[(\hat{\theta} - \theta^*)(\hat{\theta} - \theta^*)^T].$$

Then

$$\hat{\mathbf{C}} \succeq \frac{1}{N} \mathbf{I}_{\theta^*}^{-1},$$

where \succeq means that the eigenvalues of $\hat{\mathbf{C}} - \frac{1}{N} \mathbf{I}_{\theta^*}^{-1}$ are ≥ 0

Proof. We will prove this for $D = K = 1$, i.e., $x_1, \dots, x_N, \theta \in \mathbb{R}$ are scalars. The general result follows by similar arguments. Since $\hat{\theta}$ is unbiased, we know

$$0 = \mathbb{E}[\hat{\theta} - \theta^*] = \int (\hat{\theta} - \theta^*) p(\mathbf{x}|\theta^*) d\mathbf{x}.$$

Taking derivatives w.r.t. θ we have:

$$\begin{aligned} 0 &= \frac{\partial}{\partial \theta} \int (\hat{\theta} - \theta) p(\mathbf{x}|\theta) d\mathbf{x} \Big|_{\theta=\theta^*} = \int \frac{\partial}{\partial \theta} \underbrace{(\hat{\theta} - \theta)}_u \underbrace{p(\mathbf{x}|\theta)}_v d\mathbf{x} \Big|_{\theta=\theta^*} \\ &= \underbrace{\int -p(\mathbf{x}|\theta) d\mathbf{x} \Big|_{\theta=\theta^*}}_{-1} + \int (\hat{\theta} - \theta) \frac{\partial p(\mathbf{x}|\theta)}{\partial \theta} d\mathbf{x} \Big|_{\theta=\theta^*}. \end{aligned}$$

Recall that

$$\frac{\partial \log p(\mathbf{x}|\theta)}{\partial \theta} = \frac{1}{p(\mathbf{x}|\theta)} \frac{\partial p(\mathbf{x}|\theta)}{\partial \theta}.$$

Plugging this into the previous equation, we get

$$\begin{aligned} 1 &= \int (\hat{\theta} - \theta) \frac{\partial \log p(\mathbf{x}|\theta)}{\partial \theta} p(\mathbf{x}|\theta) d\mathbf{x} \Big|_{\theta=\theta^*} = \mathbb{E} \left[(\hat{\theta} - \theta) \cdot \frac{\partial \log p(\mathbf{x}|\theta)}{\partial \theta} \Big|_{\theta=\theta^*} \right] \\ &\leq \underbrace{\sqrt{\mathbb{E}[(\hat{\theta} - \theta^*)^2]}}_{\sqrt{\text{var}(\hat{\theta})}} \sqrt{\mathbb{E} \left[\left(\frac{\partial \log p(\mathbf{x}|\theta)}{\partial \theta} \right)^2 \Big|_{\theta=\theta^*} \right]}, \end{aligned} \quad (4.2)$$

where the last step follows by the Cauchy-Schwartz inequality: for two random vectors \mathbf{x}, \mathbf{y} , we define their inner product as $\langle \mathbf{x}, \mathbf{y} \rangle := \mathbb{E}[\mathbf{x}^\top \mathbf{y}]$ and their norm as $\|\mathbf{x}\|_2 := \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle} = \sqrt{\mathbb{E}[\mathbf{x}^\top \mathbf{x}]}$. Then the Cauchy-Schwartz inequality ($|\langle \mathbf{x}, \mathbf{y} \rangle| \leq \|\mathbf{x}\| \|\mathbf{y}\|$) becomes:

$$|\mathbb{E}[\mathbf{x}^\top \mathbf{y}]| \leq \sqrt{\mathbb{E}[\mathbf{x}^\top \mathbf{x}] \mathbb{E}[\mathbf{y}^\top \mathbf{y}]}.$$

We will now show that the second term in (4.2) is $-\sqrt{\mathbf{I}_{\theta^*}}$. Write:

$$\begin{aligned} \frac{\partial^2 \log p(\mathbf{x}|\theta)}{\partial \theta^2} &= \frac{\partial}{\partial \theta} \left(\frac{\partial \log p(\mathbf{x}|\theta)}{\partial \theta} \right) = \frac{\partial}{\partial \theta} \left(\frac{1}{p(\mathbf{x}|\theta)} \frac{\partial p(\mathbf{x}|\theta)}{\partial \theta} \right) \\ &= -\frac{1}{p(\mathbf{x}|\theta)^2} \frac{\partial p(\mathbf{x}|\theta)}{\partial \theta} \frac{\partial p(\mathbf{x}|\theta)}{\partial \theta} + \frac{1}{p(\mathbf{x}|\theta)} \frac{\partial^2 p(\mathbf{x}|\theta)}{\partial \theta^2} \\ &= -\underbrace{\left(\frac{1}{p(\mathbf{x}|\theta)} \frac{\partial p(\mathbf{x}|\theta)}{\partial \theta} \right)^2}_{\frac{\partial \log p(\mathbf{x}|\theta)}{\partial \theta}} + \frac{1}{p(\mathbf{x}|\theta)} \frac{\partial^2 p(\mathbf{x}|\theta)}{\partial \theta^2}. \end{aligned}$$

For the second term write:

$$\mathbb{E} \left[\frac{1}{p(\mathbf{x}|\theta)} \frac{\partial^2 p(\mathbf{x}|\theta)}{\partial \theta^2} \right] = \int \frac{1}{p(\mathbf{x}|\theta)} \frac{\partial^2 p(\mathbf{x}|\theta)}{\partial \theta^2} p(\mathbf{x}|\theta) d\mathbf{x} = \int \frac{\partial^2 p(\mathbf{x}|\theta)}{\partial \theta^2} d\mathbf{x} = \frac{\partial^2}{\partial \theta^2} \underbrace{\int p(\mathbf{x}|\theta) d\mathbf{x}}_1 = 0.$$

It follows that

$$\begin{aligned} \mathbb{E} \left[\left(\frac{\partial \log p(\mathbf{x}|\theta)}{\partial \theta} \right)^2 \Big|_{\theta=\theta^*} \right] &= \mathbb{E} \left[\frac{\partial^2 \log p(\mathbf{x}|\theta)}{\partial \theta^2} \Big|_{\theta=\theta^*} \right] = \mathbb{E} \left[\frac{\partial^2 \log \prod_{i=1}^N p(x_i|\theta)}{\partial \theta^2} \Big|_{\theta=\theta^*} \right] \\ &= \mathbb{E} \left[\frac{\partial^2 \sum_{i=1}^N \log p(x_i|\theta)}{\partial \theta^2} \Big|_{\theta=\theta^*} \right] = N \mathbb{E} \left[\frac{\partial^2 \log p(x_i|\theta)}{\partial \theta^2} \Big|_{\theta=\theta^*} \right] = -N \mathbf{I}_{\theta^*} \end{aligned}$$

Plugging this in (4.2) and taking squares, we have:

$$\text{var}(\hat{\theta}) \geq \frac{1}{N} \mathbf{I}_{\theta^*}^{-1},$$

as desired. □

4.10 Mixtures

Theorem 4.3 shows that the MLE is optimal in the sense that it achieves *the best* convergence rate to the true parameter. However, it is not always easy, or even possible, to derive or compute the MLE.

Example 4.12. Suppose we observe x_1, \dots, x_N from a *mixture* of K gaussians (see Figure 4.7):

$$x_i \stackrel{iid}{\sim} p(\mathbf{x}|\boldsymbol{\theta}^*) = \sum_{k=1}^K \rho_k^* \mathcal{N}(\mu_k^*, \sigma_k^{*2}) = \sum_{k=1}^K \rho_k^* \frac{1}{\sqrt{2\pi}\sigma_k^*} e^{-\frac{1}{2}\left(\frac{x-\mu_k^*}{\sigma_k^*}\right)^2}.$$

where $\rho_1, \dots, \rho_K \geq 0$ and $\sum_{k=1}^K \rho_k = 1$. ρ_k is the probability that a sample corresponds to the k^{th} gaussian. In this case $\boldsymbol{\theta}^* = [\theta_k^* \cdots \theta_k^*]^T$, where $\theta_k^* = \{\rho_k^*, \mu_k^*, \sigma_k^*\}$, and computing the MLE

$$\arg \max_{\boldsymbol{\theta} \in \Theta} p(\mathbf{x}|\boldsymbol{\theta}) = \arg \max_{\boldsymbol{\theta} \in \Theta} \log p(\mathbf{x}|\boldsymbol{\theta}).$$

is not so easy because the log of a sum does not factor nicely, i.e., $p(\mathbf{x}|\boldsymbol{\theta})$ is no longer convex.

In cases like these where it is not easy to compute the MLE. Estimating mixtures is an active field of research, as they are good models for classification. So, how would you do it? Give it some thought. Maybe you come up with a great idea and you'd become famous! ;)

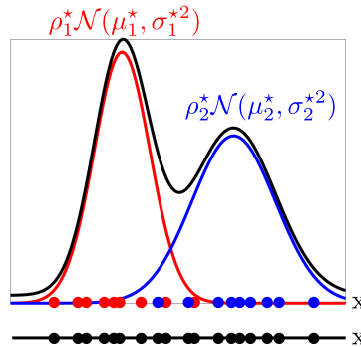


Figure 4.7: See Example 4.12. In a *mixture* each sample x_i is drawn according to one out of K distributions that we want to estimate. The challenge is that we don't know which distribution generated each x_i , and so we don't know which x_i 's should be used to estimate the k^{th} parameter. In this illustration we only observe the *black* points. If we knew which one is red and which one is blue, we could use only the red to estimate $\{\mu_1^*, \sigma_1^*\}$ and only the blue to estimate $\{\mu_2^*, \sigma_2^*\}$. The challenge is that we don't know the colors, and so we also have to estimate which distribution each point belongs to. Here $\rho_k^* \in [0, 1]$ models the fraction of points that correspond to the k^{th} distribution, which also has to be estimated. So, how would you do it? Give it some thought. Maybe you come up with a great idea and you'd become famous! ;)

A Code for Simulation of Example 4.11

```

1 clear all; close all; clc;
2
3 % ===== Code to simulate the asymptotic distribution of delta.hat =====
4 % ===== See Example 4.11.
5
6 T = 100;           % Number of trials.
7 beta = 0.1;       % Error margin that we want.
8 alpha = 0.05;     % Allowed probability of having bigger error than we want.
9 p1.star = 0.75;   % True probability of reacting favorably to treatment 1.
10 p2.star = 0.25;   % True probability of reacting favorably to treatment 1.
11 delta.star = p1.star-p2.star; % Parameter we want to estimate.
12 delta.hat = zeros(T,1); % Estimates of delta.star over trials.
13
14 % Number of patients to test each treatment, to guarantee desired accuracy.
15 nk = ceil(norminv(alpha/2)/(sqrt(2)*beta))^2;
16
17 for t=1:T,
18     % Yk contains the outcomes of the patients given treatment k,
19     % 1=reacted favorably, 0=not favorably.
20     Y1 = rand([nk,1]) < p1.star;
21     Y2 = rand([nk,1]) < p2.star;
22
23     % Xk is the number of patients that reacted favorably to treatment k.
24     X1 = sum(Y1);
25     X2 = sum(Y2);
26
27     % hat_pk is the estimate of pk.star.
28     p1.hat = X1 / nk;
29     p2.hat = X2 / nk;
30
31     delta.hat(t) = p1.hat - p2.hat;
32 end
33
34 % Plot results.
35 figure(1);
36 axes('Box','on');
37 hold on;
38 [h,rangeh] = hist(delta.hat);
39 stem(rangeh,h,'k','LineWidth',4,'MarkerSize',20);
40
41 % Asymptotic distribution of delta.hat.
42 range = 0:.001:1;
43 mu_delta = delta.star;
44 sigma_delta = sqrt((p1.star*(1-p1.star) + p2.star*(1-p2.star))/nk);
45 asymptotic_dist = normpdf(range,mu_delta,sigma_delta);
46 plot(range,asymptotic_dist/max(asymptotic_dist)*max(h),'b','LineWidth',4);
47
48 % Highlight beta and the probability alpha.
49 bounds.alpha = [0:.025:delta.star-beta,delta.star+beta:0.025:1];
50 region.alpha = normpdf(bounds.alpha,mu_delta,sigma_delta);
51 stem(bounds.alpha,region.alpha/max(asymptotic_dist)*max(h),'r',...
52     'LineWidth',5,'MarkerSize',1);
53
54 % Make figure look sexy.
55 axis tight;
56 set(gca,'XTick',[0:.1:1],'xticklabel',...
57     {'0','.1','.2','.3','.4','\delta*=.5','.6','.7','.8','.9','1'},...
58     'fontsize',15);
59 set(gca,'YTick',[],'yticklabel',[]);
60 ylabel('p$({\hat{\delta}}_{ML})$', 'Interpreter','latex','fontsize',25);

```

