

## 24 Decision and Estimation Asymptotics

In this section, we further develop and apply our asymptotic analysis to some basic hypothesis testing and estimation settings of interest.

### 24.1 Asymptotics of Hypothesis Testing

In this section, we employ the large deviation analysis tools we developed so far to characterize the asymptotic behavior of hypothesis testing. We start with the analysis of the likelihood ratio test, study the behavior of both types of errors and explore the geometry of the corresponding probability space. We then discuss the special cases of the classical and Bayesian approaches.

Let  $\mathbf{y}$  be the random variable whose values are used to disambiguate between the two hypotheses in question. Under the null hypothesis  $H_0$ ,  $\mathbf{y}$  is distributed according to the probability distribution  $p_0$ ; under the alternative hypothesis  $H_1$ ,  $\mathbf{y}$  is distributed according to the probability distribution  $p_1$ . Let  $\mathbf{y} = [y_1, \dots, y_N]^T$  be a sequence of  $N$  i.i.d. samples of the random variable  $\mathbf{y}$ . For any constant  $\gamma$ , the log-likelihood ratio test

$$\frac{1}{N} \log \frac{p_1^N(\mathbf{y})}{p_0^N(\mathbf{y})} = \frac{1}{N} \sum_{n=1}^N \log \frac{p_1(y_n)}{p_0(y_n)} \underset{\hat{H}(\mathbf{y})=H_0}{\overset{\hat{H}(\mathbf{y})=H_1}{\geq}} \gamma \quad (1)$$

induces decision regions

$$\mathcal{R}_0 = \left\{ \mathbf{y} \in \mathcal{Y}^N : \frac{1}{N} \sum_{n=1}^N \log \frac{p_1(y_n)}{p_0(y_n)} \leq \gamma \right\} \text{ and } \mathcal{R}_1 = \left\{ \mathbf{y} \in \mathcal{Y}^N : \frac{1}{N} \sum_{n=1}^N \log \frac{p_1(y_n)}{p_0(y_n)} \geq \gamma \right\}. \quad (2)$$

Here we allow a bit of imprecision by including sequences  $\mathbf{y}$  for which the test is met with equality in both sets. This makes our proofs easier, and as we already saw before, does not make much of a difference for the actual test.

The method of types allows us to map sets of observations to sets of probability distributions. Using this mapping, we construct the two corresponding sets on the probability simplex:

$$\mathcal{S}_0 = \left\{ p \in \mathcal{P}^{\mathcal{Y}} : \mathbb{E}_p \left[ \log \frac{p_1(\mathbf{y})}{p_0(\mathbf{y})} \right] \leq \gamma \right\} \quad \text{and} \quad \mathcal{S}_1 = \left\{ p \in \mathcal{P}^{\mathcal{Y}} : \mathbb{E}_p \left[ \log \frac{p_1(\mathbf{y})}{p_0(\mathbf{y})} \right] \geq \gamma \right\}. \quad (3)$$

It is easy to see the connection between the two kinds of sets:

$$\mathcal{R}_i = \{ \mathbf{y} \in \mathcal{Y}^N : \hat{p}(\cdot; \mathbf{y}) \in \mathcal{S}_i \cap \mathcal{P}_N \} \quad \text{for } i = 0, 1. \quad (4)$$

Here we consider the case when  $p_0 \in \mathcal{S}_0$  and  $p_1 \in \mathcal{S}_1$ . Since

$$\mathbb{E}_{p_0} \left[ \log \frac{p_1(y)}{p_0(y)} \right] = -\mathbb{E}_{p_0} \left[ \log \frac{p_0(y)}{p_1(y)} \right] = -D(p_0 \| p_1), \quad (5)$$

$$\mathbb{E}_{p_1} \left[ \log \frac{p_1(y)}{p_0(y)} \right] = D(p_1 \| p_0), \quad (6)$$

the corresponding range of values of the threshold  $\gamma$  is between these two values:

$$-D(p_0 \| p_1) \leq \gamma \leq D(p_1 \| p_0). \quad (7)$$

We discuss the behavior of the test when the threshold is outside this range later in this section.

Using the relationship between the sets of sequences  $\mathcal{R}_i$  and the sets of probability distributions  $\mathcal{S}_i$ , we can characterize the probability of errors in hypothesis testing. Specifically, Sanov's Theorem implies that both the probability of false alarm and the probability of miss decay exponentially with the number of samples  $N$ :

$$P_F = \mathbb{P} \left( \hat{H}(\mathbf{y}) = H_1 \mid H = H_0 \right) \quad (8)$$

$$= P_0\{\mathcal{R}_1\} = P_0\{\mathcal{S}_1 \cap \mathcal{P}_N\} \doteq 2^{-ND(p_0^* \| p_0)}, \quad (9)$$

where  $p_0^*$  is the I-projection of  $p_0$  onto  $\mathcal{S}_1$ , and

$$P_M = 1 - P_D = \mathbb{P} \left( \hat{H}(\mathbf{y}) = H_0 \mid H = H_1 \right) \quad (10)$$

$$= P_1\{\mathcal{R}_0\} = P_1\{\mathcal{S}_0 \cap \mathcal{P}_N\} \doteq 2^{-ND(p_1^* \| p_1)}, \quad (11)$$

where  $p_1^*$  is the I-projection of  $p_1$  onto  $\mathcal{S}_0$ .

Let us consider the geometry of this problem. The linear family of probability distributions

$$\mathcal{L} = \left\{ p \in \mathcal{P}^y : \mathbb{E}_p \left[ \log \frac{p_1(y)}{p_0(y)} \right] = \gamma \right\} \quad (12)$$

forms the boundary between  $\mathcal{S}_0$  and  $\mathcal{S}_1$ . It is easy to see that both  $p_0^*$  and  $p_1^*$  belong to  $\mathcal{L}$ , otherwise we could decrease the corresponding distance between the distribution and its I-projection.

Now we construct the exponential family of weighted geometric means of the distributions  $p_0$  and  $p_1$ ,

$$\mathcal{E}_{p_0, p_1}(x) = \left\{ p \in \mathcal{P}^y : p(y; x) = \frac{p_0(y)^{1-x} p_1(y)^x}{Z(x)} \right\} \quad (13)$$

$$= \left\{ p \in \mathcal{P}^y : p(y; x) = \exp \left\{ x \log \frac{p_1(y)}{p_0(y)} - \alpha(x) + \log p_0(y) \right\} \right\} \quad (14)$$

for  $0 \leq x \leq 1$ , and note that it uses  $p_0$  as the base distribution and  $t(y) = \log(p_1(y)/p_0(y))$  as the natural statistic. This exponential family therefore intersects the linear family  $\mathcal{L}$  at  $p_0^*$ , which is the I-projection of the distribution  $p_0$  onto  $\mathcal{L}$ .

The same exponential family has an alternative parameterization that uses  $\tilde{x} = 1 - x$  to index its members:

$$\mathcal{E}_{p_1, p_0}(\tilde{x}) = \left\{ p \in \mathcal{P}^{\mathcal{Y}} : p(y; \tilde{x}) = \frac{p_0(y)^{\tilde{x}} p_1(y)^{1-\tilde{x}}}{Z(\tilde{x})} \right\} \quad (15)$$

$$= \left\{ p \in \mathcal{P}^{\mathcal{Y}} : p(y; \tilde{x}) = \exp \left\{ \tilde{x} \log \frac{p_0(y)}{p_1(y)} - \alpha(\tilde{x}) + \log p_1(y) \right\} \right\} \quad (16)$$

for  $0 \leq \tilde{x} \leq 1$ . Similarly, the linear family  $\mathcal{L}$  has an alternative representation

$$\mathcal{L} = \left\{ p \in \mathcal{P}^{\mathcal{Y}} : \mathbb{E}_p \left[ \log \frac{p_0(y)}{p_1(y)} \right] = -\gamma \right\}. \quad (17)$$

Based on these alternative representations, we conclude that  $p_1^*$ , which is the I-projection of  $p_1$  onto  $\mathcal{L}$ , is also an intersection of  $\mathcal{E}_{p_0, p_1}$  and  $\mathcal{L}$ . This proves that  $p_1^* = p_0^*$ , i.e., the I-projections of  $p_0$  and  $p_1$  onto  $\mathcal{L}$  coincide.

The resulting hypothesis testing information geometry is depicted in Fig. 1.

Applying large deviation analysis we conclude that  $\mathbb{E}_{p(\cdot; x)} [\log(p(y; x)/p_0(y))]$  and  $D(p(\cdot; x) \| p_0)$  are monotonically increasing functions of the parameter  $x$ . Using the alternative representation, we can also show that  $\mathbb{E}_{p(\cdot; x)} [\log(p(y; x)/p_1(y))]$  and  $D(p(\cdot; x) \| p_1)$  are monotonically decreasing functions of the parameter  $x$ . For any member of the exponential family,

$$\mathbb{E}_{p(\cdot; x)} \left[ \log \frac{p_1(y)}{p_0(y)} \right] = \mathbb{E}_{p(\cdot; x)} \left[ \log \frac{p(y; x)}{p_0(y)} \right] - \mathbb{E}_{p(\cdot; x)} \left[ \log \frac{p(y; x)}{p_1(y)} \right] \quad (18)$$

$$= D(p(\cdot; x) \| p_0) - D(p(\cdot; x) \| p_1) \quad (19)$$

is a monotonically increasing function of  $x$ . The distribution  $p_* = p_0^* = p_1^*$  is the member of the exponential family that satisfies

$$\mathbb{E}_{p_*} \left[ \log \frac{p_1(y)}{p_0(y)} \right] = D(p_* \| p_0) - D(p_* \| p_1) = \gamma. \quad (20)$$

We now go back to the expressions for  $P_F$  and  $P_M$  in (9) and (11) and arrive at the following asymptotic characterization:

$$P_F \doteq 2^{-ND(p_* \| p_0)}, \quad P_M \doteq 2^{-ND(p_* \| p_1)}, \quad (21)$$

where

$$p_*(y) = \frac{1}{Z(x^*)} p_0(y)^{1-x^*} p_1(y)^{x^*} \quad (22)$$

with  $x^*$  chosen such that  $D(p_* \| p_0) - D(p_* \| p_1) = \gamma$ .

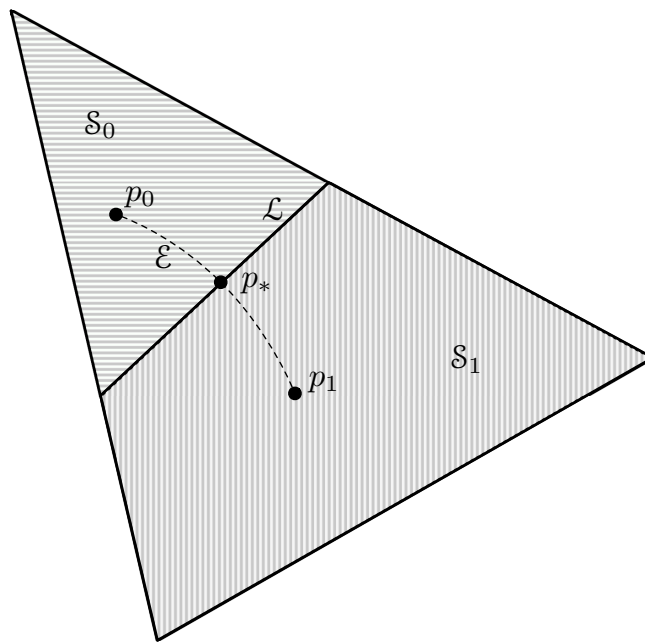


Figure 1: The geometry of hypothesis testing with data distributed i.i.d. according to  $p_0$  or  $p_1$ . The false alarm probability is exponentially small with Chernoff exponent  $D(p_* \| p_0)$ , where  $p_*$  is the I-projection of  $p_0$  onto  $\mathcal{S}_1$ , where the boundary of  $\mathcal{S}_1$  is determined by  $\gamma$ .

### 24.1.1 Classical Hypothesis Testing

In the classical setting, we place an upper limit on one of the error terms and minimize the other one. Here we constrain  $P_M$  and minimize  $P_F$ . The analysis for bounding  $P_F$  and minimizing  $P_M$  is virtually identical.

Eq. (21) implies that to keep  $P_M$  from decaying exponentially with  $N$ , we need to set  $\gamma = D(p_1 \| p_0)$ , implying  $p_* = p_1$ . In this case,

$$P_F \doteq 2^{-ND(p_1 \| p_0)}, \quad P_M \doteq 1, \quad (23)$$

i.e.,  $P_M$  decays subexponentially while the Chernoff exponent for  $P_F$  is equal to  $D(p_1 \| p_0)$ . A variant of this fact is often referred to as Stein's Lemma: when  $P_M \leq \alpha$ , the fastest rate of decay for  $P_F$  is equal to  $D(p_1 \| p_0)$ .

This setup is a bit unnatural as we fight the exponential behavior of the errors by growing the threshold on the likelihood ratio exponentially in the number of samples  $N$ . Since  $\gamma$  is equal to the logarithm of the threshold in the likelihood ratio test, we effectively commit to growing the threshold  $\eta = 2^{N\gamma} = 2^{ND(p_1 \| p_0)}$  exponentially with  $N$  to force one of the errors to decay subexponentially.

### 24.1.2 Bayesian Hypothesis Testing

In the Bayesian setting, our goal is to minimize the expected cost

$$\mathbb{E}[C] = C_{01}P_0P_F + C_{10}P_1P_M \doteq C_{01}P_02^{-ND(p_* \| p_0)} + C_{10}P_12^{-ND(p_* \| p_1)} \quad (24)$$

$$\doteq 2^{-N \min\{D(p_* \| p_0), D(p_* \| p_1)\}}. \quad (25)$$

To maximize the rate of the cost decay, we must maximize the smallest of  $D(p_* \| p_0)$  and  $D(p_* \| p_1)$ . Since  $D(p(\cdot; x) \| p_0)$  monotonically increases with  $x$  and  $D(p(\cdot; x) \| p_1)$  monotonically decreases with  $x$ , the highest rate of cost decay in (25) is achieved for  $p_*$  that satisfies  $D(p_* \| p_0) = D(p_* \| p_1)$ . This correspond to  $\gamma = 0$ , and  $\eta = 2^{N\gamma} = 1$ .

We also note that as  $N \rightarrow \infty$ , the effect of the prior and of the costs is eliminated. In this regime, the Bayesian likelihood ratio test reduces to the ML likelihood ratio test. The probability of error decays exponentially, with the Chernoff exponent equal to  $D(p_* \| p_0) = D(p_* \| p_1)$ .

### 24.1.3 Special Case: Extreme Values of Threshold

Before concluding this section, we analyze the special case when both probability distributions  $p_0$  and  $p_1$  belong to the same set  $\mathcal{S}_1$ . This corresponds to placing the threshold in the LRT so low that the typical sequences with respect to  $p_0$  are still assigned to the decision region  $\mathcal{R}_1$ , i.e.,

$$\gamma < \mathbb{E}_{p_0} \left[ \log \frac{p_1(y)}{p_0(y)} \right] = -D(p_0 \| p_1). \quad (26)$$

Similar to the analysis earlier in this section, we can show that the I-projections of both distributions onto  $\mathcal{S}_0$  coincide:

$$p_* = \arg \min_{p \in \mathcal{S}_0} D(p||p_1) = \arg \min_{p \in \mathcal{S}_0} D(p||p_0). \quad (27)$$

Then Sanov's theorem implies that the probability that a sequence  $\mathbf{y}$  of  $N$  i.i.d. samples of  $\mathbf{y}$  has a type that belongs to  $\mathcal{S}_0$  decays exponentially with  $N$ , whether the data is generated by  $p_0$  or by  $p_1$ :

$$P_M = P_1\{\mathcal{S}_0 \cap \mathcal{P}_N\} \doteq 2^{-ND(p_*||p_1)}, \quad (28)$$

$$1 - P_F = P_0\{\mathcal{S}_0 \cap \mathcal{P}_N\} \doteq 2^{-ND(p_*||p_0)}. \quad (29)$$

In other words, the probability of miss  $P_M$  decays exponentially, but the probability of false alarm  $P_F$  approaches 1 as  $N \rightarrow \infty$ .

The other special case of an exceedingly high threshold leads to a virtually identical analysis. If  $\gamma > D(p_1||p_0)$ ,  $N \rightarrow \infty$ , both the probability of false alarm and the probability of true detection decay exponentially.

## 24.2 Asymptotics of Parameter Estimation

As we saw earlier, there are many ways to evaluate the quality of estimators. The analysis of their asymptotic behavior is richer still. To develop our results, we begin by developing the concepts of convergence we will need for our development.

### 24.2.1 Convergence of Random Sequences

In this section, we define some notions of convergence of sequences of random variables that we need for our development.

**Definition 1.** *A sequence of random variables  $z_1, z_2, \dots$  is said to converge almost surely (or with probability one) to a random variable  $z$  if*

$$\mathbb{P}\left(\lim_{N \rightarrow \infty} z_N = z\right) = 1. \quad (30)$$

*We express this convergence using the notation  $z_N \xrightarrow{\text{a.s.}} z$ .*

By comparison, recall that the weak law of large numbers (WLLN) focused on convergence in probability, defined as follows.

**Definition 2.** *A sequence of random variables  $z_1, z_2, \dots$  is said to converge in probability (or with high probability) to a random variable  $z$  if for every  $\epsilon > 0$  we have*

$$\lim_{N \rightarrow \infty} \mathbb{P}(|z_N - z| < \epsilon) = 1. \quad (31)$$

*We express this convergence using the notation  $z_N \xrightarrow{\text{P}} z$ .*

Almost-sure convergence is stronger than convergence in probability, as the following alternative characterization emphasizes.<sup>1</sup>

**Fact 1** (Alternative Characterization of Almost Sure Convergence). *For a collection of random variables  $z, z_1, z_2, \dots$ , we have  $z_N \xrightarrow{\text{a.s.}} z$  if and only if for every  $\epsilon > 0$ , there exists an  $N_0$  such that*

$$\mathbb{P}(|z_N - z| < \epsilon \text{ for all } N > N_0) = 1.$$

In other words, for an almost-sure convergent sequence, when a sample path gets within  $\epsilon$  of  $z$ , it remain within that tolerance forever thereafter (with probability one).

It is also useful to note that almost-sure convergence follows when the convergence in probability is sufficiently fast, i.e., when  $\mathbb{P}(|z_N - z| > \epsilon)$  decays sufficiently strongly with  $N$ . For example, we have the following fact

**Fact 2.** *For a collection of random variables  $z, z_1, z_2, \dots$ , we have  $z_N \xrightarrow{\text{a.s.}} z$  if for every  $\epsilon > 0$  we have*

$$\sum_{n=1}^{\infty} \mathbb{P}(|z_n - z| > \epsilon) < \infty.$$

In contrast, a useful notion of convergence that is *weaker* than convergence in probability is convergence in *distribution*.

**Definition 3.** *A sequence of random variables  $z_1, z_2, \dots$  is said to converge in distribution (or in law) to a random variable  $z$  if*

$$\lim_{N \rightarrow \infty} P_{z_N}(z) = P_z(z) \tag{32}$$

*for all  $z$  at which  $P_z(\cdot)$  is continuous. We express this convergence using the notation  $z_N \xrightarrow{d} z$ .*

Although we omit a proof, an equivalent characterization of convergence in distribution is that

$$\lim_{N \rightarrow \infty} \mathbb{E}[g(z_N)] = \mathbb{E}[g(z)], \quad \text{for all bounded, continuous functions } g(\cdot).$$

Moreover, convergence in distribution implies that for some events  $\mathcal{A} \subset \mathbb{R}$  for which probabilities can be evaluated, we have

$$\mathbb{P}_{p_{z_N}}[\mathcal{A}] = \int_{\mathcal{A}} p_{z_N}(c) \, dc \quad \rightarrow \quad \mathbb{P}_{p_z}[\mathcal{A}] = \int_{\mathcal{A}} p_z(c) \, dc, \quad \text{as } N \rightarrow \infty. \tag{33}$$

---

<sup>1</sup>As an illustration of this distinction, there are examples of sequences  $z_1, z_2, \dots$  that converge in probability to  $z$  but with probability one *no* individual sequence converges, i.e.,  $\mathbb{P}(\lim_{N \rightarrow \infty} z_N = z) = 0$ .

However, in general, convergence in distribution certainly does *not* imply (33) for *all* such events  $\mathcal{A}$ . Thus Definition 3 can be viewed as a relatively weak form of convergence in distribution. To obtain such a result requires a stronger notion of convergence in distribution, such as convergence in *divergence*.<sup>2</sup>

**Definition 4.** A sequence of random variables  $\mathbf{z}_1, \mathbf{z}_2, \dots$  is said to converge in divergence (or strongly in law) to a random variable  $\mathbf{z}$  if

$$\lim_{N \rightarrow \infty} D(p_{\mathbf{z}_N} \| p_{\mathbf{z}}) = 0. \quad (34)$$

We express this convergence using the notation  $\mathbf{z}_N \xrightarrow{D} \mathbf{z}$ .

It is this stronger notion of convergence in distribution that is often needed in inference applications.<sup>3</sup> To confirm that it is indeed stronger, it suffices to note that for arbitrary densities  $p$  and  $q$  we have, for an arbitrary  $y$ ,

$$\begin{aligned} |P(y) - Q(y)| &= \left| \int_{-\infty}^y p(b) \, db - \int_{-\infty}^y q(b) \, db \right| \\ &= \left| \int_{-\infty}^y [p(b) - q(b)] \, db \right| \\ &\leq \int_{-\infty}^y |p(b) - q(b)| \, db \\ &\leq \int_{-\infty}^{+\infty} |p(b) - q(b)| \, db \\ &\leq \sqrt{D(p \| q) / (2 \log e)}, \end{aligned}$$

where to obtain the last inequality we have used (a tight version of) Pinsker's inequality for continuous alphabets. Hence if  $p, p_1, p_2, \dots$  are such that  $D(p_N \| p) \rightarrow 0$ , then  $|P_N(y) - P(y)| \rightarrow 0$  as well.

The following theorem asserts that continuous transformations of random variables preserve their mode of convergence.

**Theorem 1** (Continuous Mapping Theorem). *If  $\mathbf{z}, \mathbf{z}_1, \mathbf{z}_2, \dots$  are defined over  $\mathbb{R}^k$  and if  $g: \mathbb{R}^k \mapsto \mathbb{R}$  is a continuous function, then:*

$$\mathbf{z}_N \xrightarrow{d} \mathbf{z} \quad \Rightarrow \quad g(\mathbf{z}_N) \xrightarrow{d} \mathbf{z} \quad (35)$$

$$\mathbf{z}_N \xrightarrow{p} \mathbf{z} \quad \Rightarrow \quad g(\mathbf{z}_N) \xrightarrow{p} \mathbf{z} \quad (36)$$

$$\mathbf{z}_N \xrightarrow{\text{a.s.}} \mathbf{z} \quad \Rightarrow \quad g(\mathbf{z}_N) \xrightarrow{\text{a.s.}} \mathbf{z} \quad (37)$$

More generally the theorem holds whenever the set of points at which  $g(\cdot)$  is discontinuous has probability zero under  $p_{\mathbf{z}}$ .

<sup>2</sup>Another example of such a sufficiently stronger notion is convergence in *density*, i.e.,  $p_{\mathbf{z}_N}(z) \rightarrow p_{\mathbf{z}}(z)$  as  $N \rightarrow \infty$ , for each  $z$ .

<sup>3</sup>Recall that it is what we needed in our approximation results on exponential families.



### 24.2.2 Convergence of Sample-Averages

Equipped with the preceding notions of convergence, we now discuss how sample averages of i.i.d. random variables behave with respect to these notions, which our analysis of estimators will rely on.

**Theorem 2** (Strong Law of Large Numbers (SLLN)). *Let  $w_1, \dots, w_N$  be a set of i.i.d. random variables with mean  $\mu$  and  $\mathbb{E}[|w_n|] < \infty$ . Then*

$$\frac{1}{N} \sum_{n=1}^N w_n \xrightarrow{\text{a.s.}} \mu, \quad \text{as } N \rightarrow \infty, \quad (38)$$

*i.e., the sample-average converges almost surely to  $\mu$ .*

By comparison, the weak law states that

$$\frac{1}{N} \sum_{n=1}^N w_n \xrightarrow{\text{p}} \mu.$$

Note that the strong law does not require stronger conditions than what we required in our statement of the weak law, in an earlier installment of the notes.

For analyzing estimators, we need a more general version of the law of large numbers, which is called the *uniform* law.

**Theorem 3** (Uniform Law of Large Numbers (ULLN)). *Let  $w_1, \dots, w_N$  be a set of i.i.d. random variables, let  $\Theta$  be a compact parameter set, and let  $g$  be a function such that: i)  $g(w; \theta)$  is continuous in  $\theta \in \Theta$ , and ii) there exists  $g_+(w) > 0$  such that  $\mathbb{E}[g_+(w)] < \infty$  and  $|g(w; \theta)| \leq g_+(w)$  for all  $w$  and  $\theta \in \Theta$ . Then*

$$\max_{\theta \in \Theta} \left| \frac{1}{N} \sum_{n=1}^N g(w_n; \theta) - \mu(\theta) \right| \xrightarrow{\text{a.s.}} 0, \quad \text{as } N \rightarrow \infty, \quad (39)$$

*where  $\mu(\theta) = \mathbb{E}[g(w; \theta)]$ , i.e., the sample-average almost surely converges uniformly<sup>4</sup> to its mean.*

This is the *strong* uniform law of large numbers, corresponding to almost sure convergence (i.e., with probability one). There is also a weak version, where the convergence is in probability (i.e., with high probability).

The following celebrated theorem establishes that the distribution of a sample-average converges, in an appropriate sense, to a Gaussian.

---

<sup>4</sup>As a (nonrandom) illustration of the kind of convergence phenomena this result is trying to avoid, consider a function  $f_N: [0, 1] \mapsto [0, 1]$  for  $N = 1, 2, \dots$ , where  $f_N(x) = \mathbb{1}_{1/N^2 \leq x \leq 1/N}$ . Then  $f_N(x) \rightarrow 0$  as  $n \rightarrow \infty$  for each  $x \in [0, 1]$ , i.e.,  $f_N(x)$  converges *pointwise* to  $f(x) \equiv 0$ . However  $\max_{x \in [0, 1]} |f_N(x) - f(x)| = \max_{x \in [0, 1]} |f_N(x)| = 1$  for all  $N \geq 1$ , so  $f_N(x)$  does *not* converge uniformly to  $f(x) \equiv 0$ .

**Theorem 4** (Central Limit Theorem (CLT)). *Let  $w_1, \dots, w_N$  be a set of i.i.d. random variables with mean  $\mu$  and variance  $\sigma^2 < \infty$ . Then<sup>5</sup>*

$$\sqrt{N} \left[ \frac{1}{N} \sum_{n=1}^N w_n - \mu \right] \xrightarrow{d} \mathcal{N}(0, \sigma^2), \quad \text{as } N \rightarrow \infty. \quad (40)$$

Note that (40) expresses that for large  $N$  the sample average behaves like

$$\frac{1}{N} \sum_{n=1}^N w_n \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{N}\right).$$

As an aside, a theorem due to Pólya establishes that if  $z_N \xrightarrow{d} z$  and  $P_z(\cdot)$  is continuous, then the convergence is uniform, i.e.,

$$\lim_{N \rightarrow \infty} \max_z |P_{z_N}(z) - P_z(z)| = 0.$$

Hence, the convergence in the CLT is uniform. A useful bound on the (specifically,  $1/\sqrt{N}$ ) rate of this uniform convergence is provided by what is referred to as the Berry-Esseen Theorem.

There is also a stronger CLT corresponding to convergence in divergence, which is particularly useful for inference applications.

**Theorem 5** (Strong Central Limit Theorem (SCLT)). *Let  $w_1, \dots, w_N$  be a set of i.i.d. random variables with mean  $\mu$  and variance  $\sigma^2 < \infty$ . Then*

$$e_N \triangleq \sqrt{N} \left[ \frac{1}{N} \sum_{n=1}^N w_n - \mu \right] \xrightarrow{D} \mathcal{N}(0, \sigma^2), \quad \text{as } N \rightarrow \infty \quad (41)$$

*if and only if  $D(p_{e_N} \| \mathcal{N}(0, \sigma^2)) < \infty$  for some  $N$ .*

### 24.2.3 Asymptotics of ML Estimates

Let  $\{p_y(\cdot; a) : a \in \mathcal{X}\}$  be a parameterized family of distributions,  $y$  be a random variable distributed according to a family member  $p_y(\cdot; x)$ . We start by constructing the normalized log-likelihood of  $N$  i.i.d. samples drawn from an arbitrary member of the family  $p_y(\cdot; a)$ , viewed as a function of  $a$ :

$$\ell_y^N(a) = \frac{1}{N} \sum_{n=1}^N \log p_y(y_n; a). \quad (42)$$

---

<sup>5</sup>Here we adopt a slight abuse of notation, since we really should write  $\xrightarrow{d} z \sim \mathcal{N}(0, \sigma^2)$ . However, we will use this convenient alternate notation interchangeably in the sequel, since the risk of confusion is minimal.

Since  $\ell_{\mathbf{y}}^N(a)$  depends on the random variables  $y_1, \dots, y_N$ , it is also a random variable. At the same time, we can view it as a (probabilistic) function of  $a$ . It is easy to see that

$$\mathbb{E}_{p_y(\cdot; x)} [\log p_y(y; a)] = -D(p_y(\cdot; x) \| p_y(\cdot; a)) - H(p_y(\cdot; x)). \quad (43)$$

The LLN therefore implies that

$$\ell_{\mathbf{y}}^N(a) \rightarrow g(a) \triangleq -D(p_y(\cdot; x) \| p_y(\cdot; a)) - H(p_y(\cdot; x)). \quad (44)$$

In other words, with probability one the likelihood function  $L_N(a; \mathbf{y})$  we observe for any particular sequence  $\mathbf{y}$  will approach  $g(a)$  as  $N \rightarrow \infty$ . Since the maximum of the function  $g(a)$  is achieved at the true value of the parameter  $x$ , we expect the ML estimate of  $x$  to be close to  $x$ .

The following theorem makes this framework of reasoning precise establishes conditions under which this indeed happens.

**Theorem 6.** *Consider a model family  $\{p_y(\cdot; a) : a \in \mathcal{X}\}$  for a parameter set  $\mathcal{X} \subset \mathbb{R}$ , and let  $y_1, \dots, y_N$  be i.i.d. according to the model  $p_y(\cdot; x)$  for some  $x \in \mathcal{X}$ . Then if:*

1.  $\mathcal{X}$  is compact;
2.  $p_y(y; a)$  is continuous in  $a \in \mathcal{X}$  for all  $y$  and continuous in  $y$  in a neighborhood of  $a = x$ ;
3. there exists an  $\ell_y^+$  such that  $\mathbb{E} [|\ell_y^+|] < \infty$  and

$$\tilde{\ell}_y(a) \triangleq \log \frac{p_y(y; a)}{p_y(y; x)} \leq \tilde{\ell}_y^+, \quad \text{for all } a \in \mathcal{X} \text{ and } y; \text{ and}$$

4. the model is identifiable in the sense that  $x' = x$  if  $p_y(y; x') = p_y(y; x)$  (almost everywhere);

then the ML estimator  $\hat{x}_N(y_1, \dots, y_N)$  satisfies

$$\hat{x}_N(y_1, \dots, y_N) \xrightarrow{\text{a.s.}} x, \quad (45)$$

i.e.,  $\hat{x}_N(y_1, \dots, y_N)$  converges to  $x$  as  $N \rightarrow \infty$  with probability one.

When (45) holds for an estimator, we say that the estimator is *consistent*. More precisely, we say it is *strongly* consistent, to distinguish this from the scenario where  $\hat{x}_N \xrightarrow{\text{P}} x$ , in which case the estimator is weakly consistent. Thus, Theorem 6 establishes that the ML estimator is consistent.

*Proof.* First, for  $a \in \mathcal{X}$  let

$$\tilde{\ell}_{\mathbf{y}}^N(a) \triangleq \frac{1}{N} \sum_{n=1}^N \tilde{\ell}_{y_n}(a).$$

Then

$$\hat{x}_N(y_1, \dots, y_N) = \arg \max_{a \in \mathcal{X}} \tilde{\ell}_{\mathbf{y}}^N(a), \quad (46)$$

and

$$\mathbb{E} \left[ \tilde{\ell}_{\mathbf{y}}^N(a) \right] = -D(p_y(\cdot; x) \| p_y(\cdot; a)) \leq 0. \quad (47)$$

Next, note that

$$\tilde{\ell}_{\mathbf{y}}^N(\hat{x}_N) = \max_{a \in \mathcal{X}} \tilde{\ell}_{\mathbf{y}}^N(a) \geq \tilde{\ell}(x) = 0. \quad (48)$$

Now for  $\rho > 0$  define the compact set

$$\mathcal{X}_\rho \triangleq \{a \in \mathcal{X} : |a - x| \geq \rho\},$$

and let

$$\bar{D}_\rho = \max_{a \in \mathcal{X}_\rho} -D(p_y(\cdot; x) \| p_y(\cdot; a)) < 0, \quad (49)$$

where the strict inequality follows from the fact that  $x \notin \mathcal{X}_\rho$ . Finally, let  $x_\rho \in \mathcal{X}_\rho$  denote the value  $a$  at which the maximum in (49) is achieved.

Then when the model is such that the ULLN can be applied for the parameter set  $\mathcal{X}_\rho$ , we know that for any choice of  $\epsilon > 0$ , there exists an  $N_0$  such that for all  $N > N_0$ ,

$$|\tilde{\ell}_{\mathbf{y}}^N(a) + D(p_y(\cdot; x) \| p_y(\cdot; a))| < \epsilon, \quad \text{for all } a \in \mathcal{X}_\rho, \quad (50)$$

with probability one. So let us choose, say,  $\epsilon = -\bar{D}_\rho/2$ , so (50) implies

$$\tilde{\ell}_{\mathbf{y}}^N(a) < -\bar{D}_\rho/2 - D(p_y(\cdot; x) \| p_y(\cdot; a)), \quad \text{for all } a \in \mathcal{X}_\rho. \quad (51)$$

Hence

$$\begin{aligned} \max_{a \in \mathcal{X}_\rho} \tilde{\ell}_{\mathbf{y}}^N(a) &\leq -\bar{D}_\rho/2 + \max_{a \in \mathcal{X}_\rho} -D(p_y(\cdot; x) \| p_y(\cdot; a)) \\ &= -\bar{D}_\rho/2 + \bar{D}_\rho = \bar{D}_\rho/2 \\ &< 0. \end{aligned} \quad (52)$$

Comparing (52) with (48) we conclude that  $\hat{x}_N(y_1, \dots, y_N) \notin \mathcal{X}_\rho$  as  $N \rightarrow \infty$  with probability one, i.e.,  $|\hat{x}_N(y_1, \dots, y_N) - x| < \rho$ . Since  $\rho > 0$  is arbitrary, we conclude  $\hat{x}_N(y_1, \dots, y_N) \xrightarrow{\text{a.s.}} x$ .  $\square$

Not only does the ML estimator  $\hat{x}_N$  tend toward the correct value  $x$  for large  $N$ , but its fluctuations about  $x$  are Gaussian in the large  $N$  regime. In particular, a consistent ML estimator behaves for large  $N$  like

$$\hat{x}_N \sim \mathcal{N}\left(x, \frac{1}{NJ_y(x)}\right). \quad (53)$$

The following theorem more precisely characterizes this behavior under suitable conditions (that, for example, allow the needed ULLN to be applied).

**Theorem 7.** *Consider a model family  $\{p_y(\cdot; a) : a \in \mathcal{X}\}$  for a parameter set  $\mathcal{X} \subset \mathbb{R}$ , and let  $y_1, \dots, y_N$  be i.i.d. according to the model  $p_y(\cdot; x)$  for some  $x \in \mathcal{X}$ . If*

1.  $\mathcal{X}$  is open;
2.  $\partial^2 p_y(y; a) / \partial a^2$  is continuous in  $a \in \mathcal{X}$  for all  $y$ , and we have the usual regularity condition

$$\int \frac{\partial}{\partial a} p_y(y; a) dy = 0;$$

3. there exists  $\ddot{\ell}_y^+$  such that  $\mathbb{E} [\ddot{\ell}_y^+; x] < \infty$  and

$$|\ddot{\ell}_y(a)| \leq \ddot{\ell}_y^+, \quad \text{for all } a \text{ in a neighborhood of } x \text{ and all } y,$$

where

$$\ell_y(a) \triangleq \log p_y(y; a); \tag{54}$$

4.  $J_y(x) > 0$ ; and
5. the ML estimator  $\hat{x}_N = \hat{x}_N(y_1, \dots, y_N)$  is (strongly) consistent;

then

$$\sqrt{N}(\hat{x}_N - x) \xrightarrow{d} \mathbf{N}(0, 1/J_y(x)). \tag{55}$$

*Proof.* Let

$$\ell_{\mathbf{y}}^N(a) \triangleq \frac{1}{N} \sum_{n=1}^N \ell_{y_n}(a), \tag{56}$$

which we note is a sample-average of i.i.d. random variables. Then by Taylor's Theorem<sup>6</sup> we can write the score function as

$$\dot{\ell}_{\mathbf{y}}^N(a) = \dot{\ell}_{\mathbf{y}}^N(x) + (a - x) \ddot{\ell}_{\mathbf{y}}^N(a') \tag{58}$$

for some  $a'$  between  $a$  and  $x$ . Let us choose  $a = \hat{x}_N$ , in which case, since  $\hat{x}_N$  is achieved at a stationary point of the likelihood function (because  $\mathcal{X}$  is open), we have

$$0 = \dot{\ell}_{\mathbf{y}}^N(\hat{x}_N) = \dot{\ell}_{\mathbf{y}}^N(x) + (\hat{x}_N - x) \ddot{\ell}_{\mathbf{y}}^N(x'), \quad \text{for some } x' \text{ between } \hat{x}_N \text{ and } x,$$

---

<sup>6</sup>The special case of Taylor's Theorem we need states that if a function  $f$  is differentiable, then for any  $t$  and  $t_0$  we can write

$$f(t) = f(t_0) + (t - t_0) \dot{f}(t') \tag{57}$$

for some  $t'$  between  $t$  and  $t_0$ .

whence

$$\sqrt{N}(\hat{x}_N - x) = \frac{\sqrt{N}\dot{\ell}_{\mathbf{y}}^N(x)}{-\ddot{\ell}_{\mathbf{y}}^N(x')}. \quad (59)$$

Now since, since

$$\dot{\ell}_{\mathbf{y}}^N(a) = \frac{1}{N} \sum_{n=1}^N \dot{\ell}_{\mathbf{y}_n}(a)$$

is also the sample-average of i.i.d. random variables for any  $a$ , and since, as we saw in our development of the Cramér-Rao bound,

$$\mathbb{E} \left[ \dot{\ell}_{\mathbf{y}}^N(x) \right] = 0 \quad \text{and} \quad \mathbb{E} \left[ (\dot{\ell}_{\mathbf{y}}^N(x))^2 \right] = J_{\mathbf{y}}(x)$$

under our regularity conditions, via the CLT we conclude that the numerator of the right-hand side of (59) satisfies

$$\sqrt{N}\dot{\ell}_{\mathbf{y}}^N(x) \xrightarrow{d} \mathcal{N}(0, J_{\mathbf{y}}(x)).$$

Turning next to the denominator of the right-hand side of (59), we note that if

$$-\ddot{\ell}_{\mathbf{y}}^N(x') \xrightarrow{\text{a.s.}} J_{\mathbf{y}}(x) \quad (60)$$

then by the Continuous Mapping Theorem we obtain (55), since  $\text{var}[z/J_{\mathbf{y}}(x)] = (\text{var } z)/J_{\mathbf{y}}(x)^2$  for any random variable  $z$ .

Thus, it remains only to show (60). Proceeding, since<sup>7</sup>  $\mathbb{E}[\ddot{\ell}_{\mathbf{y}}(a); x]$  is continuous in  $a$  and since, as we also showed in our development of the Cramér-Rao bound,  $\mathbb{E}[\ddot{\ell}_{\mathbf{y}}(x); x] = -J_{\mathbf{y}}(x)$ , for any  $\epsilon > 0$  there exists a  $\rho > 0$  such that

$$\left| \mathbb{E}[\ddot{\ell}_{\mathbf{y}}(a); x] + J_{\mathbf{y}}(x) \right| < \epsilon/2$$

whenever  $a \in \bar{\mathcal{S}}_{\rho} \triangleq \{a' : |a' - x| \leq \rho\}$ .

Now since

$$\ddot{\ell}_{\mathbf{y}}^N(a) = \frac{1}{N} \sum_{n=1}^N \ddot{\ell}_{\mathbf{y}_n}(a)$$

is also a sample-average of i.i.d. random variables, by the ULLN there exists  $n_1$  such that for all  $n > n_1$

$$\max_{a \in \bar{\mathcal{S}}_{\rho}} \left| \ddot{\ell}_{\mathbf{y}}(a) - \mathbb{E}[\ddot{\ell}_{\mathbf{y}}(a); x] \right| < \epsilon/2 \quad (61)$$

Moreover, since  $\hat{x}_N$  is (strongly) consistent, there exists  $n_2$  such that for all  $n > n_2$  we have  $\hat{x}_N \in \bar{\mathcal{S}}_{\rho}$  and, in turn,  $x' \in \bar{\mathcal{S}}_{\rho}$ , since  $x'$  lies between  $\hat{x}_N$  and  $x$ .

---

<sup>7</sup>Recall our notational convention is such that  $\mathbb{E}[f(\mathbf{y}); x] = \mathbb{E}_{p_{\mathbf{y}}(\cdot; x)}[f(\mathbf{y})]$  for a given  $f$ .

Hence, for all  $n > \max(n_1, n_2)$  we have

$$\begin{aligned}
\left| -\ddot{\ell}_{\mathbf{y}}(x') - J_{\mathbf{y}}(x) \right| &= \left| \ddot{\ell}_{\mathbf{y}}(x') + J_{\mathbf{y}}(x) \right| \\
&= \left| \ddot{\ell}_{\mathbf{y}}(x') - \mathbb{E} \left[ \ddot{\ell}_{\mathbf{y}}(x'); x \right] + \mathbb{E} \left[ \ddot{\ell}_{\mathbf{y}}(x'); x \right] + J_{\mathbf{y}}(x) \right| \\
&\leq \left| \ddot{\ell}_{\mathbf{y}}(x') - \mathbb{E} \left[ \ddot{\ell}_{\mathbf{y}}(x'); x \right] \right| + \left| \mathbb{E} \left[ \ddot{\ell}_{\mathbf{y}}(x'); x \right] + J_{\mathbf{y}}(x) \right| \\
&\leq \underbrace{\max_{a \in \mathcal{S}_\rho} \left| \ddot{\ell}_{\mathbf{y}}(a) - \mathbb{E} \left[ \ddot{\ell}_{\mathbf{y}}(a); x \right] \right|}_{\leq \epsilon/2} + \underbrace{\left| \mathbb{E} \left[ \ddot{\ell}_{\mathbf{y}}(x'); x \right] + J_{\mathbf{y}}(x) \right|}_{\leq \epsilon/2} \\
&\leq \epsilon,
\end{aligned}$$

which establishes (60) as desired.  $\square$

We observe that Theorem 7 further establishes that the variance of the ML estimator asymptotically satisfies the Cramér-Rao bound in a particular sense. More generally, we have the following definition.

**Definition 5.** Let  $y_1, \dots, y_N$  be i.i.d. according to the model  $p_{\mathbf{y}}(\cdot; x)$  for some  $x \in \mathcal{X}$ . An estimator  $\hat{x}_N = \hat{x}_N(y_1, \dots, y_N)$  is asymptotically efficient if

$$\text{var} \left[ \sqrt{N}(\hat{x}_N - x) \right] \rightarrow 1/J_{\mathbf{y}}(x), \quad \text{as } N \rightarrow \infty. \quad (62)$$

Hence, Theorem 7 establishes that the ML estimator is *asymptotically efficient*.

It is straightforward to extend the results of this section to the case of vector parameters. In such analysis, the inverse of the Fisher information matrix, for example, determines the asymptotic covariance structure of the estimator.

#### 24.2.4 Mismatched ML Estimation

In the preceding analysis, the true model lies in the parameterized class, i.e.  $x \in \mathcal{X}$ . However, in practice, we cannot always know we have chosen  $\mathcal{X}$  large enough such that this is the case. Let us examine what happens when  $x \notin \mathcal{X}$ .

If the data is generated by a probability distribution  $q$ , the SLLN implies

$$\frac{1}{N} \sum_{n=1}^N \log \frac{p_{\mathbf{y}}(y_n; x)}{q(y_n)} \xrightarrow{\text{a.s.}} -D(q(\cdot) \| p_{\mathbf{y}}(\cdot; x)), \quad (63)$$

therefore

$$p_{\mathbf{y}}(\mathbf{y}; x) = \prod_{n=1}^N p_{\mathbf{y}}(y_n; x) \doteq e^{-ND(q(\cdot) \| p_{\mathbf{y}}(\cdot; x))} \prod_{n=1}^N q(y_n). \quad (64)$$

If  $q$  is not in the family  $p_y(\cdot; x)$ , the ML estimate is equal to the value of the parameter  $x^*$  that minimizes  $D(q(\cdot) \| p_y(\cdot; x))$ . Since  $D(q(\cdot) \| p_y(\cdot; x))$  is equal to the log-loss penalty for making inference based on  $p_y(\cdot; x)$  when the observations are generated by  $q$ , ML estimation achieves asymptotically optimal solution with respect to the log-loss function.

### 24.2.5 Asymptotics of MAP Estimation

In the Bayesian case, we model the hidden parameter as a random variable  $x$  distributed according to the probability distribution  $p_x(\cdot)$ . To understand the asymptotic behavior of the MAP estimates, we consider the log ratio of the posterior probability distribution of  $x$  and the posterior probability of the true parameter value  $x_0$ , given a sequence  $\mathbf{y}$  of  $N$  i.i.d. samples of the random variable  $y$ :

$$\frac{1}{N} \log \frac{p_{x|\mathbf{y}}(x|\mathbf{y})}{p_{x|\mathbf{y}}(x_0|\mathbf{y})} = \frac{1}{N} \log \frac{p_x(x)}{p_x(x_0)} + \frac{1}{N} \sum_{n=1}^N \log \frac{p_y(y_n|x)}{p_y(y_n|x_0)}. \quad (65)$$

As  $N \rightarrow \infty$ , the first term vanishes while the second term does not. Instead, it converges to  $g(x) = -D(p_y(\cdot|x_0) \| p_y(\cdot|x))$ . The estimates we construct by maximizing the posterior probability distribution become increasingly close to those constructed based on maximizing the likelihood. This shows again that asymptotically, the importance of priors is diminished. However, choosing a prior might affect the rate of convergence, as we will see in the next section when we analyze model capacity.

## 24.3 Further Reading

The following texts provide useful supplementary reading:

Cover and Thomas, *Elements of Information Theory*, provide an in-depth discussion and analysis of typical sets, the method of types and hypothesis testing.

Serfling, *Approximation Theorems of Mathematical Statistics*, offers a detailed analysis of asymptotic behavior of parameter estimation and hypothesis testing.

Ferguson, *A Course in Large Sample Theory*, has a careful, self-contained treatment of several aspects of estimation asymptotics.