

# 1 SAMPLING DISTRIBUTIONS

In this chapter, we briefly review some concepts from MATH 61.2 regarding the probability distributions of functions of random variables (now called statistics), using techniques like moment-generating functions and the cdf technique. We prove the Lindeberg-Lévy form of the celebrated central limit theorem, as well as derive certain special distributions of sample means, sample variances, and order statistics. Finally, we introduce three new probability distributions, the  $\chi^2$  distribution, the student  $t$  distribution, and the  $\mathcal{F}$  distribution.

## 1.1 Distribution of the mean and the CLT

Given  $n$  random variables,  $X_1, X_2, \dots, X_n$ , a **statistic** is a function  $g(X_1, X_2, \dots, X_n)$  of these random variables. The probability distribution of a statistic is called its sampling distribution.

*Example.* Let  $(X_i)_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} \text{Be}(p)$ . Define  $T_n$  to be their sum,  $T_n = \sum_{i=1}^n X_i$ . Then, we know that  $T_n \sim \text{Bin}(n, p)$ . ■

*Example.* Let  $X_1, X_2 \stackrel{\text{i.i.d.}}{\sim} \mathcal{U}(0, 1)$ . Find the distribution of  $X_1 + X_2$ .

Note that the support for  $X_1 + X_2$  is the unit square in  $\mathbb{R}^2$ . The function  $x_1 + x_2 = t$  corresponds to a line moving diagonally across the square from left to right. We use the cdf technique:

$$F_{X_1+X_2}(t) = \Pr[X_1 + X_2 \leq t] = \iint_R 1 \cdot dx_1 dx_2, \quad R = \{(x_1, x_2) \in [0, 1] : x_1 + x_2 \leq t\}.$$

Notice that we can split the integral into two cases: if (1)  $0 < t < 1$ , then  $R$  is bounded by the two axes and  $x_1 + x_2 = t$ ; whereas if (2)  $1 \leq t < 2$ , then  $R$  is the unit square minus the upper right triangle bounded by  $x_2 = 1, x_1 = 1$ , and  $x_1 + x_2 = t$ . This gives us:

→ when  $0 < t < 1$ :

$$F_{X_1+X_2}(t) = \int_0^t \int_0^{t-x_2} 1 \cdot dx_1 dx_2 = \int_0^t (t - x_2) dx_2 = \left( tx_2 - \frac{x_2^2}{2} \right) \Big|_0^t = t^2 - \frac{t^2}{2} = \frac{t^2}{2}.$$

→ when  $1 \leq t < 2$ :

$$\begin{aligned} F_{X_1+X_2}(t) &= 1 - \int_{t-1}^1 \int_{t-x_2}^1 1 \cdot dx_1 dx_2 = 1 - \int_{t-1}^1 (1 - t + x_2) dx_2 \\ &= 1 - \left( x_2(1 - t) + \frac{x_2^2}{2} \right) \Big|_{t-1}^1 = \frac{1}{2} + (t - 1) - \frac{(t - 1)^2}{2}. \end{aligned}$$

This gives us the cdf of  $X_1 + X_2$ . To find its pdf, we simply differentiate and obtain

$$p_{X_1+X_2}(t) = \begin{cases} t, & \text{when } 0 < t < 1 \\ 2 - t, & \text{when } 1 \leq t < 2 \end{cases}$$

as our desired pdf. ■

Consider  $(X_i)_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mu, \sigma^2)$ . Using moment-generating functions, we can determine the distribution of the so-called **sampling mean**  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ , a common statistic. First recall that  $\mathcal{M}_{\alpha X}(t) = \mathcal{M}_X(\alpha t)$ , and  $\mathcal{M}_{\sum X_i}(t) =$

$\prod_i \mathcal{M}_{X_i}(t)$  for a random variable  $X$ . Now:

$$\begin{aligned}\mathcal{M}_{\bar{X}_n}(t) &= \mathcal{M}_{\frac{1}{n} \sum X_i}(t) = \mathcal{M}_{\sum X_i}\left(\frac{t}{n}\right) = \prod_{i=1}^n \mathcal{M}_{X_i}\left(\frac{t}{n}\right) \\ &= \left\{ \exp \left[ \mu \left( \frac{t}{n} \right) + \frac{1}{2} \sigma^2 \left( \frac{t}{n} \right)^2 \right] \right\}^n = \exp \left[ \mu t + \frac{1}{2} \left( \frac{\sigma^2}{n} \right) t^2 \right], \quad \text{substituting } \frac{t}{n} \text{ into } \mathcal{M}_{X_i}(t)\end{aligned}$$

which is the mgf of a normal random variable with mean  $\mu$  and variance  $\frac{\sigma^2}{n}$ , so  $\bar{X}_n \sim \mathcal{N}(\mu, \frac{\sigma^2}{n})$ . In general, recall also from MATH 61.2 that if  $(X_i)_{i=1}^n$  are independent normal random variables with means  $\mu_i$  and variances  $\sigma_i^2$ , then the distribution of a linear combination of them would obey a normal distribution with a similarly linear combination of their means and variances. Symbolically,

$$(X_i)_{i=1}^n \stackrel{\text{ind}}{\sim} \mathcal{N}(\mu_i, \sigma_i^2) \implies \sum_{i=1}^n a_i X_i \sim \mathcal{N}\left(\sum_{i=1}^n a_i \mu_i, \sum_{i=1}^n a_i^2 \sigma_i^2\right).$$

*Example.* Compute the probability that the sample mean of size 10 taken from a normal population with mean 1 and variance 2 has a value between 1.2 and 3.1.

We know that the sample mean  $\bar{X}_{10} \sim \mathcal{N}(1, \frac{2}{10})$ . Thus,

$$\Pr[1.2 < \bar{X}_{10} < 3.1] = \text{pnorm}(3.1, \text{mean} = 1, \text{sd} = \text{sqrt}(0.2)) - \text{pnorm}(1.2, \text{mean} = 1, \text{sd} = \text{sqrt}(0.2)) \approx 0.3273,$$

using R commands. ■

In practice, though, it might be that we don't know the underlying distribution of a sample  $(X_i)_{i=1}^n$ . This problem arises frequently when dealing with real data, which does not, in general, follow an explicitly given probability distribution. However, we have the following result:

### Theorem 1: Lindeberg-Lévy CLT

Suppose  $(X_i)_{i=1}^n$  are i.i.d. random variables obeying an unknown probability distribution with mean  $\mu$  and variance  $\sigma^2$ . Then,

$$\bar{X}_n \xrightarrow{d} \mathcal{N}(\mu, \frac{\sigma^2}{n}),$$

where  $\xrightarrow{d}$  indicates the distribution approaching  $\mathcal{N}(\mu, \frac{\sigma^2}{n})$  as  $n \rightarrow \infty$ .

*Proof.* To prove this, we will use mgf's. Let  $\mathcal{M}_X(t)$  be the mgf of one of these  $X_i$ 's. To show that  $\mathcal{M}_{\bar{X}_n}(t)$  approaches the mgf of a normal random variable with mean  $\mu$  and variance  $\frac{\sigma^2}{n}$  as  $n \rightarrow \infty$ , it suffices to show that

$$\lim_{n \rightarrow \infty} \mathcal{M}_Z(t) = e^{t^2/2}, \quad \text{where } Z = \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}}, \quad \text{the standardized form of } \bar{X}_n. \quad (1.1)$$

First, we manipulate  $Z$  to make it resemble a linear combination of the  $X_i$ 's:

$$\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} = \frac{\sqrt{n}}{\sigma} \left( \frac{1}{n} \sum_{i=1}^n (X_i - \mu) \right) = \frac{1}{\sigma\sqrt{n}} \sum_{i=1}^n (X_i - \mu).$$

Substituting this into (1.1) gives us

$$\mathcal{M}_Z(t) = \mathcal{M}_{\sum (X_i - \mu)}\left(\frac{t}{\sigma\sqrt{n}}\right) = \prod_{i=1}^n \mathcal{M}_{X_i - \mu}\left(\frac{t}{\sigma\sqrt{n}}\right) = \left[ \mathcal{M}_{X_i - \mu}\left(\frac{t}{\sigma\sqrt{n}}\right) \right]^n.$$

Since we are working with exponents to the  $n$ , it might be helpful to take the logarithm of this expression. Hence, notice that to prove (1.1), it further suffices to show

$$\lim_{n \rightarrow \infty} \ln \mathcal{M}_Z(t) = \frac{t^2}{2}.$$

As a final simplification, we can make a substitution letting  $h := \frac{t}{\sigma\sqrt{n}}$ , so that  $n = \frac{t^2}{\sigma^2 h^2}$  and  $h \rightarrow 0$  as  $n \rightarrow \infty$ . Then:

$$\begin{aligned} \lim_{n \rightarrow \infty} \ln \mathcal{M}_Z(t) &= \lim_{n \rightarrow \infty} \ln \left[ \mathcal{M}_{X_i - \mu} \left( \frac{t}{\sigma\sqrt{n}} \right) \right]^n = \lim_{n \rightarrow \infty} n \cdot \ln \left( \mathcal{M}_{X_i - \mu} \left( \frac{t}{\sigma\sqrt{n}} \right) \right) \\ &= \lim_{h \rightarrow 0} \frac{t^2}{\sigma^2 h^2} \cdot \ln \left( \mathcal{M}_{X_i - \mu}(h) \right) = \frac{t^2}{\sigma^2} \lim_{h \rightarrow 0} \frac{\ln \left( \mathcal{M}_{X_i - \mu}(h) \right)}{h^2} \quad \text{form } \frac{0}{0}, \text{ since } \mathcal{M}_X(0) = 1 \text{ for all r.v.'s } X \\ &= \frac{t^2}{\sigma^2} \lim_{h \rightarrow 0} \frac{\frac{1}{\mathcal{M}_{X_i - \mu}(h)} \cdot \mathcal{M}'_{X_i - \mu}(h)}{2h} = \frac{t^2}{\sigma^2} \lim_{h \rightarrow 0} \frac{\mathcal{M}'_{X_i - \mu}(h)}{2h \cdot \mathcal{M}_{X_i - \mu}(h)} \quad \text{form } \frac{0}{0}, \text{ since } \mathcal{M}'_{X_i - \mu}(0) = \mathbb{E}[X_i - \mu] = 0 \\ &\quad \mathbb{E}[(X_i - \mu)^2] = \text{var}[X_i] = \sigma^2 \\ &= \frac{t^2}{\sigma^2} \lim_{h \rightarrow 0} \frac{\underbrace{\mathcal{M}''_{X_i - \mu}(h)}_{2} \cdot \underbrace{\mathcal{M}'_{X_i - \mu}(h)}_0}{2 \cdot \mathcal{M}_{X_i - \mu}(h) + 2h \cdot \mathcal{M}'_{X_i - \mu}(h)} = \frac{t^2}{\sigma^2} \cdot \frac{\sigma^2}{2} = \frac{t^2}{2}, \end{aligned}$$

as desired. □

*Example.* Let  $(X_i)_{i=1}^{64} \stackrel{\text{i.i.d.}}{\sim} \mathcal{U}(0, 1)$ . Find  $k \in \mathbb{R}$  such that  $\Pr \left[ \sqrt[64]{\prod_{i=1}^{64} X_i} \geq k \right] = 0.10$ .

Taking the logarithm of the expression inside the probability function gives us

$$\Pr \left[ \frac{1}{64} \sum_{i=1}^{64} \ln X_i \geq \ln k \right] = 0.10,$$

which we can recognize as the sampling mean of random variables  $(\ln X_i)_{i=1}^{64}$ . Ideally, then, we can apply the CLT to approximate the value of  $k$  that will make this statement true, since this will give us an approximation of the distribution of  $\overline{\ln X}_{64}$ , but first we need to find  $\mathbb{E}[\ln X_i]$  and  $\text{var}[\ln X_i]$ . We use the cdf technique. Note that since  $X_i$  takes on values between 0 and 1,  $\ln X_i$  will take on values less than 0. Hence

$$F_{\ln X_i}(t) = \Pr[\ln X_i \leq t] = \Pr[X_i \leq e^t] = \int_0^{e^t} 1 \cdot dx_i, t < 0,$$

and so  $p_{\ln X_i}(t) = e^t$ , with  $t < 0$ . From here we can see that  $\mathbb{E}[\ln X_i] = -1$  and  $\text{var}[\ln X_i] = 1$ . Then, by the CLT, we have  $\frac{1}{64} \sum_{i=1}^{64} \ln X_i \xrightarrow{d} \mathcal{N}(-1, \frac{1}{64})$ . Finally, we compute

$$\begin{aligned} \Pr \left[ \frac{1}{64} \sum_{i=1}^{64} \ln X_i \geq \ln k \right] &= 0.10 \\ 1 - \Pr \left[ \frac{1}{64} \sum_{i=1}^{64} \ln X_i \leq \ln k \right] &= 0.10 \\ \ln k &= \text{qnorm}(0.9, \text{mean} = -1, \text{sd} = 1/8) \\ k &\approx e^{-0.84}, \end{aligned}$$

as desired. ■