STAT/QSCI 403: Introduction to Resampling Methods

Spring 2025

Lecture 0: Review on Probability and Statistics

Instructor: Marina Meilă. Course notes courtesy of Yen-Chi Chen

## 0.1 Random Variables

Here we will ignore the formal mathematical definition of a random variable and directly talk about it property. For a random variable X, the *cumulative distribution function* (CDF) of X is

$$P_X(x) = F(x) = P(X \le x).$$

Actually, the distribution of X is completely determined by the CDF F(x), regardless of X being a discrete random variable or a continuous random variable (or a mix of them).

If X is discrete, its probability mass function (PMF) is

$$p(x) = P(X = x).$$

If X is continuous, its probability density function (PDF) is

$$p(x) = F'(x) = \frac{d}{dx}F(x).$$

Moreover, the CDF can be written as

$$F(x) = P(X \le x) = \int_{-\infty}^{x} p(x')dx'.$$

Generally, we write  $X \sim F$  or  $X \sim p$  indicating that the random variable X has a CDF F or a PMF/PDF p.

For two random variables X, Y, their joint CDF is

$$P_{XY}(x,y) = F(x,y) = P(X \le x, Y \le y).$$

The corresponding joint PDF is

$$p(x,y) = \frac{\partial^2 F(x,y)}{\partial x \partial y}.$$

The conditional PDF of Y given X = x is

$$p(y|x) = \frac{p(x,y)}{p(x)},$$

where  $p(x) = \int_{-\infty}^{\infty} p(x, y) dy$  is sometimes called the marginal density function. Note that you can definition the joint PMF and conditional PMF using a similar way.

# 0.2 Expected Value

For a function g(x), the quantity g(X) will also be a random variable and its expected value is

$$\mathbb{E}(g(X)) = \int g(x)dF(x) = \begin{cases} \int_{-\infty}^{\infty} g(x)p(x)dx, & \text{if } X \text{ is continuous} \\ \sum_{x} g(x)p(x), & \text{if } X \text{ is discrete} \end{cases}.$$

When f(x) = x, this reduces to the usual definition of expected value.

Here are some useful properties and quantities related to the expected value:

- $\mathbb{E}(\sum_{j=1}^k c_j g_j(X)) = \sum_{j=1}^k c_j \cdot \mathbb{E}(g_j(X_i)).$
- We often write  $\mu = \mathbb{E}(X)$  as the mean (expectation) of X.
- $\operatorname{Var}(X) = \mathbb{E}((X \mu)^2)$  is the variance of X.
- If  $X_1, \cdots, X_n$  are independent, then

$$\mathbb{E}(X_1 \cdot X_2 \cdots X_n) = \mathbb{E}(X_1) \cdot \mathbb{E}(X_2) \cdots \mathbb{E}(X_n)$$

• If  $X_1, \dots, X_n$  are independent, then

$$\operatorname{Var}\left(\sum_{i=1}^{n} a_i X_i\right) = \sum_{i=1}^{n} a_i^2 \cdot \operatorname{Var}(X_i).$$

• For two random variables X and Y with their mean being  $\mu_X$  and  $\mu_Y$  and variance being  $\sigma_X^2$  and  $\sigma_Y^2$ . The covariance

$$\mathsf{Cov}(X,Y) = \mathbb{E}((X-\mu_x)(Y-\mu_y)) = \mathbb{E}(XY) - \mu_x \mu_y$$

and the (Pearson's) correlation

$$\rho(X,Y) = \frac{\mathsf{Cov}(X,Y)}{\sigma_x \sigma_y}.$$

The conditional expectation of Y given X is the random variable  $\mathbb{E}(Y|X) = g(X)$  such that when X = x, its value is

$$\mathbb{E}(Y|X=x) = \int y p(y|x) dy,$$

where p(y|x) = p(x, y)/p(x).

## 0.3 Common Distributions

### 0.3.1 Discrete Random Variables

**Bernoulli.** If X is a Bernoulli random variable with parameter p, then X = 0 or, 1 such that

$$P(X = 1) = p, \quad P(X = 0) = 1 - p.$$

In this case, we write  $X \sim \mathsf{Ber}(p)$ .

**Binomial.** If X is a binomial random variable with parameter (n, p), then  $X = 0, 1, \dots, n$  such that

$$P(X=k) = \binom{n}{k} p^k (1-p)^{n-k}.$$

In this case, we write  $X \sim Bin(n, p)$ . Note that if  $X_1, \dots, X_n \sim Ber(p)$ , then the sum  $S_n = X_1 + X_2 + \dots + X_n$  is a binomial random variable with parameter (n, p).

**Poisson.** If X is a Poisson random variable with parameter  $\lambda$ , then  $X = 0, 1, 2, 3, \cdots$  and

$$P(X=k) = \frac{\lambda^k e^{-\lambda}}{k!}.$$

In this case, we write  $X \sim \mathsf{Poi}(\lambda)$ .

#### 0.3.2 Continuous Random Variables

**Uniform.** If X is a uniform random variable over the interval [a, b], then

$$p(x) = \frac{1}{b-a}I(a \le x \le b),$$

where I(statement) is the indicator function such that if the statement is true, then it outputs 1 otherwise 0. Namely, p(x) takes value  $\frac{1}{b-a}$  when  $x \in [a, b]$  and p(x) = 0 in other regions. In this case, we write  $X \sim \text{Uni}[a, b]$ .

**Normal.** If X is a normal random variable with parameter  $(\mu, \sigma^2)$ , then

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

In this case, we write  $X \sim N(\mu, \sigma^2)$ .

**Exponential.** If X is an exponential random variable with parameter  $\lambda$ , then X takes values in  $[0,\infty)$  and

$$p(x) = \lambda e^{-\lambda x}.$$

In this case, we write  $X \sim \mathsf{Exp}(\lambda)$ . Note that we can also write

$$p(x) = \lambda e^{-\lambda x} I(x \ge 0).$$

### 0.4 Useful Theorems

We write  $X_1, \dots, X_n \sim F$  when  $X_1, \dots, X_n$  are IID (independently, identically distributed) from a CDF F. In this case,  $X_1, \dots, X_n$  is called a *random sample*.

For a sequence of random variables  $Z_1, \dots, Z_n, \dots$ , we say  $Z_n$  converges in probability to a fixed number  $\mu$  if for any  $\epsilon > 0$ ,

$$\lim_{n \to \infty} P(|Z_n - \mu| > \epsilon) = 0$$

and we will write

 $Z_n \xrightarrow{P} \mu.$ 

In other words,  $Z_n$  converges in probability implies that the distribution is concentrating at the targeting point.

Let  $F_1, \dots, F_n, \dots$  be the corresponding CDFs of  $Z_1, \dots, Z_n, \dots$ . For a random variable Z with CDF F, we say  $Z_n$  converges in distribution to Z if for every x,

$$\lim_{n \to \infty} F_n(x) = F(x).$$

In this case, we write

$$Z_n \xrightarrow{D} Z.$$

Namely, the CDF's of the sequence of random variables converge to a the CDF of a fixed random variable.

**Theorem 0.1 (Weak) Law of Large Number.** Let  $X_1, \dots, X_n \sim F$  and  $\mu = \mathbb{E}(X_1)$ . If  $\mathbb{E}|X_1| < \infty$ , then the sample average

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$$

converges in probability to  $\mu$ . i.e.,

$$\bar{X}_n \stackrel{P}{\to} \mu.$$

**Theorem 0.2 Central Limit Theorem.** Let  $X_1, \dots, X_n \sim F$  and  $\mu = \mathbb{E}(X_1)$  and  $\sigma^2 = Var(X_1) < \infty$ . Let  $\bar{X}_n$  be the sample average. Then

$$\sqrt{n}\left(\frac{X_n-\mu}{\sigma}\right) \xrightarrow{D} N(0,1).$$

Note that N(0,1) is also called standard normal random variable.

### 0.5 Estimators and Estimation Theory

Let  $X_1, \dots, X_n \sim F$  be a random sample. Here we can interpret F as the population distribution we are sampling from (that's why we are generating data from this distribution). Any numerical quantity (or even non-numerical quantity) of F that we are interested in is called the **parameter of interest**. For instance, the parameter of interest can be the mean of F, the median of F, standard deviation of F, first quartile of F, ... etc. The parameter of interest can even be  $P(X \ge t) = 1 - F(t) = S(t)$ . The function S(t) is called the *survival function*, which is a central topic in biostatistics and medical research.

When we know (or assume) that F is a certain distribution with some parameters, then the parameter of interest can be the parameter describing that distribution. For instance, if we assume F is an exponential distribution with an unknown parameter  $\lambda$ . Then this unknown parameter  $\lambda$  might be the parameter of interest.

Most of the statistical analysis is concerned with the following question:

"given the parameter of interest, how can I use the random sample to infer it?"

Let  $\theta = \theta(F)$  be the parameter of interest and let  $\hat{\theta}_n$  be a statistic (a function of the random sample  $X_1, \dots, X_n$ ) that we use to estimate  $\theta$ . In this case,  $\hat{\theta}_n$  is called an *estimator*. For an estimator, there are two important quantities measuring its quality. The first quantity is the **bias**:

$$\operatorname{Bias}(\hat{\theta}_n) = \mathbb{E}(\hat{\theta}_n) - \theta,$$

which captures the systematic deviation of the estimator from its target. The other quantity is the **variance**  $Var(\hat{\theta}_n)$ , which measures the size of stochastic fluctuation.

**Example.** Let  $X_1, \dots, X_n \sim F$  and  $\mu = \mathbb{E}(X_1)$  and  $\sigma^2 = \mathsf{Var}(X)$ . Assume the parameter of interest is the population mean  $\mu$ . Then a natural estimator is the sample average  $\hat{\mu}_n = \bar{X}_n$ . Using this estimator, then

$$\mathbf{bias}(\hat{\mu}_n) = \mu - \mu = 0, \quad \mathsf{Var}(\hat{\mu}_n) = \frac{\sigma^2}{n}.$$

Therefore, when  $n \to \infty$ , both bias and variance converge to 0. Thus, we say  $\hat{\mu}_n$  is a **consistent** estimator of  $\mu$ . Formally, an estimator  $\hat{\theta}_n$  is called a consistent estimator of  $\theta$  if  $\hat{\theta}_n \xrightarrow{P} \theta$ .

The following lemma is a common approach to prove consistency:

**Lemma 0.3** Let  $\hat{\theta}_n$  be an estimator of  $\theta$ . If  $\mathbf{bias}(\hat{\theta}_n) \to 0$  and  $\mathbf{Var}(\hat{\theta}_n) \to 0$ , then  $\hat{\theta}_n \xrightarrow{P} \theta$ . *i.e.*,  $\hat{\theta}_n$  is a consistent estimator of  $\theta$ .

In many statistical analysis, a common measure of the quality of the estimator is the mean square error (MSE), which is defined as

$$\mathsf{MSE}(\hat{\theta}_n) = \mathsf{MSE}(\hat{\theta}_n, \theta) = \mathbb{E}\left((\hat{\theta}_n - \theta)^2\right)$$

By simple algebra, the MSE of  $\hat{\theta}_n$  equals

$$\begin{split} \mathsf{MSE}(\hat{\theta}_n, \theta) &= \mathbb{E}\left((\hat{\theta}_n - \theta)^2\right) \\ &= \mathbb{E}\left((\hat{\theta}_n - \mathbb{E}(\hat{\theta}_n) + \mathbb{E}(\hat{\theta}_n) - \theta)^2\right) \\ &= \underbrace{\mathbb{E}\left((\hat{\theta}_n - \mathbb{E}(\hat{\theta}_n))^2\right)}_{=\mathsf{Var}(\hat{\theta}_n)} + 2\underbrace{\mathbb{E}\left(\hat{\theta}_n - \mathbb{E}(\hat{\theta}_n)\right)}_{=0} \cdot (\mathbb{E}(\hat{\theta}_n) - \theta) + \left(\underbrace{\mathbb{E}(\hat{\theta}_n) - \theta}_{=\mathsf{bias}(\hat{\theta}_n)}\right)^2 \\ &= \mathsf{Var}(\hat{\theta}_n) + \mathsf{bias}^2(\hat{\theta}_n). \end{split}$$

Namely, the MSE of an estimator is the variance plus the square of bias. This decomposition is also known as the *bias-variance tradeoff* (or bias-variance decomposition). By the Markov inequality,

$$\mathsf{MSE}(\hat{\theta}_n, \theta) \to 0 \Longrightarrow \hat{\theta}_n \xrightarrow{P} \theta.$$

i.e., if an estimator has MSE converging to 0, then it is a consistent estimator. The convergence of MSE is related to the  $L_2$  convergence in probability theory.

Note that we write  $\theta = \theta(F)$  for the parameter of interest because  $\theta$  is a quantity derived from the population distribution F. Thus, we may say that the parameter of interest  $\theta$  is a 'functional' (function of function; the input is a function, and the output is a real number).

• : There are two common methods of finding an estimator: the first one is called the MLE (maximum likelihood estimator), the other one is called the MOM (method of moments)<sup>1</sup>. You can google these two terms and you will find lots of references about them.

Question to think: if the parameter of interest is  $F(x) = P(X \le x)$ , what will be the estimator of it?

<sup>&</sup>lt;sup>1</sup>https://en.wikipedia.org/wiki/Method\_of\_moments\_(statistics) and MIT open course