

# All of Statistic

## 2 - Random Variable

Levi Kassel

May 7, 2023

### Summary

**Random Variable:** A Random variable is a mapping

$$X : \Omega \rightarrow \mathbb{R}$$

that assigns a real number  $X(\omega)$  to each number  $\omega$ .

**Cumulative Distribution Function (CDF):** The cumulative distribution function, or cdf, is the function  $F_X(x) : \mathbb{R} \rightarrow [0, 1]$  defined by:

$$F_X(x) = \mathbb{P}(X \leq x)$$

**Theorem:** Let  $X$  have cdf  $F$  and let  $Y$  have cdf  $G$ . If  $F(x) = G(x)$  for all  $x$ , then  $\mathbb{P}(X \in A) = \mathbb{P}(Y \in A)$  for all  $A$ .

**Theorem:** A function  $F$  mapping the real line to  $[0, 1]$  is a cdf for some probability  $\mathbb{P}$  if and only if  $F$  satisfies the following three conditions:

- $F$  is non-decreasing:  $x_1 < x_2$  implies that  $F(x_1) \leq F(x_2)$
- $F$  is normelized:

$$\lim_{x \rightarrow -\infty} F(x) = 0$$

and:

$$\lim_{x \rightarrow \infty} F(x) = 1$$

- $F$  is right-continuous:  $F(x) = F(x^+)$  for all  $x$ , where:

$$F(x^+) = \lim_{x \rightarrow y, y > x} F(y)$$

**Probability Funtion:**When  $X$  is **discrete**, We define the probability function or probability mass function for  $X$  by

$$f_X(x) = \mathbb{P}(X = x)$$

. Thus,  $f_X(x) \geq 0$  for all  $x \in R$  and  $\sum_i f_X(xi) = 1$ .

**Probability Density Function (PDF):** when  $X$  is **continuous** we define the Probability Density Function  $f_X(x)$  as follow:

- $f_X(x) \geq 0, \forall x$ .
- $\int_{-\infty}^{\infty} f_X(x)dx = 1$
- for every  $a \leq b$ :

$$\mathbb{P}(a < X < b) = \int_a^b f_X(x)dx.$$

- $f_X(x) = F'_X(x)$

**Inverse CDF:** Let  $X$  be a random variable with CDF  $F$ . The inverse CDF or **quantile** function is defined by:

$$F^{-1}(q) = \inf \{ x : F(x) > q \}$$

for  $q \in [0, 1]$ .

## Some Important Discrete Random Variables

**The Point Mass Distribution:**  $X \sim \delta_a$

$$F(x) = \begin{cases} 0, & x < a \\ 1, & x \geq a \end{cases}$$

**The Discrete Uniform Distribution:** Let  $k > 1$  be a given integer. Suppose that  $X$  has probability mass function given by

$$f(x) = \begin{cases} \frac{1}{k}, & \text{for } x = 1, \dots, k \\ 0, & \text{otherwise} \end{cases}$$

**The Bernoulli Distribution:** Let  $X$  represent a binary coin flip. Then  $\mathbb{P}(X = 1) = p$  and  $\mathbb{P}(X = 0) = 1 - p$  for some  $p \in [0, 1]$ . We say that  $X$  has a Bernoulli distribution written  $X \sim \text{Bernoulli}(p)$ . The probability function is

$$f(x) = p^x \cdot (1 - p)^{(1-x)} \text{ for } x \in \{0, 1\}.$$

**The Binomial Distribution:** Flip the coin  $n$  times and let  $X$  be the number of heads. Assume that the tosses are independent.  $X \sim \text{Binomial}(n, p)$

$$f(x) = \begin{cases} \binom{n}{x} \cdot p^x \cdot (1 - p)^{n-x}, & \text{for } x = 0, \dots, n \\ 0, & \text{otherwise} \end{cases}$$

If  $X_1 \sim \text{Binomial}(n_1, p)$  and  $X_2 \sim \text{Binomial}(n_2, p)$  then  $X_1 + X_2 \sim \text{Binomial}(n_1 + n_2, p)$

$n_2, p)$ .

**The Geometric Distribution:**  $X$  has a geometric distribution with parameter  $p \in (0, 1)$ , written  $X \sim \text{Geom}(p)$ , if

$$\mathbb{P}(X = k) = p \cdot (1 - p)^{k-1}, \quad k \geq 1$$

Think of  $X$  as the number of flips needed until the first head when flipping a coin.

**The Poisson Distribution:**  $X$  has a Poisson distribution with parameter  $\lambda$ , written  $X \sim \text{Poisson}(\lambda)$  if:

$$f(x) = e^{-\lambda} \cdot \frac{\lambda^x}{x!}, \quad x \geq 0$$

If  $X_1 \sim \text{Poisson}(\lambda_1)$  and  $X_2 \sim \text{Poisson}(\lambda_2)$  then  $X_1 + X_2 \sim \text{Poisson}(\lambda_1 + \lambda_2)$ .

## Some Important Continuous Random Variables

**The Uniform Distribution:**  $X$  has a  $\text{Uniform}(a, b)$  distribution, written  $X \sim \text{Uniform}(a, b)$ , if

$$f(x) = \begin{cases} \frac{1}{b-a}, & \text{for } x \in [a, b] \\ 0, & \text{otherwise} \end{cases}$$

$$F(x) = \begin{cases} 0, & x < a \\ \frac{x-a}{b-a}, & \text{for } x \in [a, b] \\ 1, & x > b \end{cases}$$

**Normal (Gaussian):**  $X$  has a Normal (or Gaussian) distribution with parameters  $\mu$  and  $\sigma$ , denoted by  $X \sim N(\mu, \sigma^2)$ , if

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp \left\{ -\frac{1}{2\sigma^2}(x - \mu)^2 \right\}.$$

We say that  $X$  has a standard Normal distribution if  $\mu = 0$  and  $\sigma = 1$  (denoted by  $Z$ ).

The PDF and CDF of a standard Normal are denoted by  $\phi(z)$  and  $\Phi(z)$ .

There is no closed-form expression for  $\Phi$ . Here are some useful facts:

1. if  $X \sim N(\mu, \sigma^2)$ , then  $Z = (X - \mu)/\sigma \sim N(0, 1)$ .
2. if  $Z \sim N(0, 1)$  then  $\mu + \sigma \cdot Z \sim N(\mu, \sigma^2)$ .
3. if  $X \sim N(\mu_i, \sigma_i^2)$ ,  $i = 1, \dots, n$  are independent, then

$$\sum_{i=1}^n X_i \sim N\left(\sum_{i=1}^n \mu_i, \sum_{i=1}^n \sigma_i^2\right)$$

4. if  $X \sim N(\mu, \sigma^2)$ , then:

$$\mathbb{P}(a < X < b) = \mathbb{P}\left(\frac{a - \mu}{\sigma} < Z < \frac{b - \mu}{\sigma}\right) = \Phi\left(\frac{b - \mu}{\sigma}\right) - \Phi\left(\frac{a - \mu}{\sigma}\right)$$

**Exponential Distribution:**  $X$  has an Exponential distribution with parameter  $\beta$ , denoted by  $X \sim Exp(\beta)$ , if

$$f(x) = \frac{1}{\beta} e^{-x/\beta}, \quad x > 0$$

where  $\beta > 0$ .

**Gamma Distribution:** For  $\alpha > 0$ , the Gamma function is defined by

$$\Gamma(\alpha) = \int_0^{\infty} y^{\alpha-1} e^{-y} dy$$

$X$  has a Gamma distribution with parameters  $\alpha$  and  $\beta$ , denoted by  $X \sim Gamma(\alpha, \beta)$ , if

$$f(x) = \frac{1}{\beta^{\alpha} \Gamma(\alpha)} x^{\alpha-1} e^{-x/\beta}, \quad x > 0$$

## Bivariate Distributions

**joint mass function** Given a pair of discrete random variables  $X$  and  $Y$ , define the joint mass function by  $f(x, y) = \mathbb{P}(X = x \text{ and } Y = y)$ .

$f(x, y)$  is a **PDF** for the **continuous** random variables  $(X, Y)$  if:

1.  $f(x, y) \geq 0$  for all  $(x, y)$
2.  $\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x, y) dx dy = 1$
3. for any set  $A \subset \mathbb{R} \times \mathbb{R}$ ,  $\mathbb{P}((X, Y) \in A) = \int \int_A f(x, y) dx dy$

**Marginal Distributions:** If  $(X, Y)$  have joint distribution with mass function  $f_{X,Y}$ , then the marginal mass function for  $X$  is defined by:

For **discrete** variables:

$$f_X(x) = P(X = x) = \sum_y \mathbb{P}(X = x, Y = y) = \sum_y f(x, y)$$

For **continuous** variables:

$$f_X(x) = \int_y f(x, y) dy$$

**Independent Random Variables** Two random variables  $X$  and  $Y$  are independent if, for every  $A$  and  $B$ ,

$$\mathbb{P}(X \in A, Y \in B) = \mathbb{P}(X \in A) \cdot \mathbb{P}(Y \in B)$$

**Theorem:** Let  $X$  and  $Y$  have joint PDF  $f_{X,Y}$ . Then  $X \perp Y$  if and only if  $f_{X,Y}(x,y) = f_X(x)f_Y(y)$  for all values  $x$  and  $y$ . **Theorem:** Suppose that the range of  $X$  and  $Y$  is a (possibly infinite) rectangle. If  $f(x,y) = g(x)h(y)$  for some functions  $g$  and  $h$  (not necessarily probability density functions) then  $X$  and  $Y$  are independent.

**Conditional Distributions:** The conditional probability mass function is:  
For **discrete** variables:

$$f_{X|Y} = \mathbb{P}(X = x|Y = y) = \frac{\mathbb{P}(X = x, Y = y)}{\mathbb{P}(Y = y)} = \frac{f_{X,Y}(x,y)}{f_Y(y)}, \quad \text{if } f_Y(y) > 0$$

For **continuous** variables:

$$\mathbb{P}(X \in A|Y = y) = \int_A f_{X|Y}(x|y)dx$$

**Multivariate Distributions and iid Samples** If  $X_1, \dots, X_n$  are independent and each has the same marginal distribution with CDF  $F$ , we say that  $X_1, \dots, X_n$  are iid (independent and identically distributed) and we write:

$$X_1, \dots, X_n \sim F$$

**Multinomial** The multivariate version of a Binomial is called a Multinomial. We say that  $X$  has a *Multinomial*( $n, p$ ) distribution written  $X \sim \text{Multinomial}(n, p)$ . The probability function is:

$$f(x) = \binom{n}{x_1, \dots, x_n} \cdot p^{x_1} \dots p^{x_n}$$

**Lemma:** Suppose that  $X \sim \text{Multinomial}(n, p)$  where  $X = (X_1, \dots, X_k)$  and  $p = (p_1, \dots, p_k)$ . The marginal distribution of  $X_j$  is *Binomial*( $n, p_j$ ).

**Multivariate Normal** let:

$$\mathbf{Z} = \begin{pmatrix} z_1 \\ \cdot \\ \cdot \\ \cdot \\ z_k \end{pmatrix}$$

where  $Z_1, \dots, Z_k \sim N(0, 1)$  are independent. The density of  $\mathbf{Z}$  is :

$$f(z) = \prod_{i=1}^k f(z_i) = \frac{1}{2\pi^{k/2}} \exp \left\{ -\frac{1}{2} \sum_{j=1}^k z_j^2 \right\} = \frac{1}{2\pi^{k/2}} \exp \left\{ -\frac{1}{2} \mathbf{z}^T \mathbf{z} \right\}$$

More generally, a vector  $\mathbf{X}$  has a multivariate Normal distribution, denoted by  $\mathbf{X} \sim N(\mu, \Sigma)$ , if it has density

$$f(x; \mu, \Sigma) = \frac{1}{2\pi^{k/2} |\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu) \right\}$$

**Theorem:** If  $Z \sim N(0, I)$  and  $X = \mu + \Sigma^{-1/2} \cdot Z$  then  $X \sim N(\mu, \Sigma)$ . Conversely, if  $X \sim N(\mu, \Sigma)$ , then  $\Sigma^{1/2}(X - \mu) \sim N(0, I)$

**Theorem:** Let  $X \sim N(\mu, \Sigma)$ . Then:

1. The marginal distribution of  $X_a$  is  $X_a \sim N(\mu_a, \Sigma_{aa})$ .
2. The conditional distribution of  $X_b$  given  $X_a = x_a$  is:

$$X_b|X_a = x_a \sim N(\mu_b + \Sigma_{ba}\Sigma_{aa}^{-1} \cdot (x_a - \mu_a), \Sigma_{bb} - \Sigma_{ba}\Sigma_{aa}^{-1}\Sigma_{ab}).$$

3. If  $a$  is a vector then  $a^T X \sim N(a^T \mu, a^T \Sigma a)$ .
4.  $V = (X\mu)^T \Sigma^{-1} (X\mu) \sim \chi_k^2$ .

## Transformations of Random Variables

Let  $Y = r(X)$  be a function of  $X$ , for example,  $Y = X^2$  or  $Y = e^X$ . We call  $Y = r(X)$  a transformation of  $X$ .

How do we compute the pdf and cdf of  $Y$ ?

**Discrete case:** The mass function of  $Y$  is given by:

$$f_Y(y) = \mathbb{P}(Y = y) = \mathbb{P}(r(X) = y) = \mathbb{P}(\{x; r(x) = y\}) = \mathbb{P}(X \in r^{-1}(y)).$$

**Continuous case:**

- (a) For each  $y$ , find the set  $A_y = \{x : r(x) \leq y\}$ .
- (b) Find the CDF:

$$F_Y(y) = \mathbb{P}(Y \leq y) = \mathbb{P}(r(X) \leq y) = \mathbb{P}(\{x; r(x) \leq y\}) = \int_{A_y} f_X(x) dx$$

- (c) The PDF is  $f_Y(y) = F_Y'(y)$ .

## Transformations of Several Random Variables

$X$  and  $Y$  are given random variables, we might want to know the distribution of  $X/Y$ ,  $X + Y$ ,  $\max(X, Y)$  or  $\min(X, Y)$ . Let  $Z = r(X, Y)$  be the function of interest. The steps for finding  $f_Z$  are the same as before:

- (a) For each  $z$ , find the set  $A_z = \{(x, y) : r(x, y) \leq z\}$ .
- (b) Find the CDF:

$$F_Z(z) = \mathbb{P}(Z \leq z) = \mathbb{P}(r(X, Y) \leq z) = \mathbb{P}(\{(x, y); r(x, y) \leq z\}) = \int \int_{A_z} f_{X,Y}(x, y) dx dy.$$

- (c) Then  $f_Z(z) = F_Z'(z)$ .