# Basic concepts in Probability and Statistics

## Peter Sørensen

## 2022-02-09

This compendium summarizes important probability probability concepts, formulas, and distributions widely used in statistical genetics. It is based on the following material:

- https://github.com/wzchen/probability_cheatsheet
- https://en.wikipedia.org/wiki/Algebra_of_random_variables
- https://en.wikipedia.org/wiki/Random_variable

## Probability theory

**Independent Events:**

$A$ and $B$ are independent if knowing whether $A$ occurred gives no information about whether $B$ occurred. More formally, $A$ and $B$ (which have nonzero probability) are independent if and only if one of the following equivalent statements holds:

$$P(A \cap B) = P(A)P(B) \tag{1}$$
$$P(A|B) = P(A) \tag{2}$$
$$P(B|A) = P(B) \tag{3}$$

**Conditional Independence of Events:**

$A$ and $B$ are conditionally independent given $C$ if $P(A \cap B|C) = P(A|C)P(B|C)$. Conditional independence does not imply independence, and independence does not imply conditional independence.

**Joint, Marginal and Conditional Probabilitiy**

- Joint Probability $P(A \cap B)$ or $P(A, B)$ is the joint probability of $A$ and $B$.
- Marginal (Unconditional) Probability $P(A)$ is the marginal probability of $A$.
- Conditional Probability $P(A|B) = P(A, B)/P(B)$ is the conditional probability of $A$, given that $B$ occurred.
- Conditional Probability *is* Probability] $P(A|B)$ is a probability function for any fixed $B$. Any theorem that holds for probability also holds for conditional probability.

**Law of Total Probability (LOTP)**

Let $B_1, B_2, B_3, ...B_n$ be a *partition* of the sample space (i.e., they are disjoint and their union is the entire sample space).

$$P(A) = P(A|B_1)P(B_1) + P(A|B_2)P(B_2) + \cdots + P(A|B_n)P(B_n) \tag{4}$$

$$P(A) = P(A \cap B_1) + P(A \cap B_2) + \cdots + P(A \cap B_n) \tag{5}$$

For **LOTP with extra conditioning**, just add in another event $C$!

$$P(A|C) = P(A|B_1, C)P(B_1|C) + \cdots + P(A|B_n, C)P(B_n|C) \tag{6}$$

$$P(A|C) = P(A \cap B_1|C) + P(A \cap B_2|C) + \cdots + P(A \cap B_n|C) \tag{7}$$

**Bayes' Rule**

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} \tag{8}$$

**Bayes' Rule with with extra conditioning**

$$P(A|B, C) = \frac{P(B|A, C)P(A|C)}{P(B|C)} \tag{9}$$

We can also write

$$P(A|B, C) = \frac{P(A, B, C)}{P(B, C)} = \frac{P(B, C|A)P(A)}{P(B, C)} \tag{10}$$

**Odds Form of Bayes' Rule**

$$\frac{P(A|B)}{P(A^c|B)} = \frac{P(B|A)}{P(B|A^c)} \frac{P(A)}{P(A^c)} \tag{11}$$

The *posterior odds* of $A$ are the *likelihood ratio* times the *prior odds*.

## Random Variables and their Distributions

### Probability Mass Function (PMF)

Gives the probability that a *discrete* random variable takes on the value $x$.

$$p_X(x) = P(X = x) \tag{12}$$

The PMF satisfies

$$p_X(x) \geq 0 \text{ and } \sum_x p_X(x) = 1 \tag{13}$$

### Cumulative Distribution Function (CDF)

Gives the probability that a random variable is less than or equal to $x$.

$$F_X(x) = P(X \leq x) \tag{14}$$

The CDF is an increasing, right-continuous function with

$$F_X(x) \to 0 \text{ as } x \to -\infty \text{ and } F_X(x) \to 1 \text{ as } x \to \infty \tag{15}$$

- Independence: Intuitively, two random variables are independent if knowing the value of one gives no information about the other. Discrete random variables $X$ and $Y$ are independent if for *all* values of $x$ and $y$

$$P(X = x, Y = y) = P(X = x)P(Y = y) \tag{16}$$

### Expected Value and Linearity

The expected Value (a.k.a.~*mean*, *expectation*, or *average*) is a weighted average of the possible outcomes of our random variable. Mathematically, if $x_1, x_2, x_3, \ldots$ are all of the distinct possible values that $X$ can take, the expected value of $X$ is

$$E(X) = \sum_i x_i P(X = x_i) \tag{17}$$

- Linearity: For any random variables $X$ and $Y$, and constants $a, b, c$,

$$E(aX + bY + c) = aE(X) + bE(Y) + c \tag{18}$$

- Same distribution implies same mean: If $X$ and $Y$ have the same distribution, then $E(X) = E(Y)$ and, more generally,

$$E(g(X)) = E(g(Y)) \tag{19}$$

- Conditional Expected Value: is defined like expectation, only conditioned on any event $A$.

$$\mathrm{E}(X|A) = \sum_x x P(X = x|A) \tag{20}$$

**Indicator Random Variable**

An indicator random Variable is a random variable that takes on the value 1 or 0. It is always an indicator of some event: if the event occurs, the indicator is 1; otherwise it is 0. They are useful for many problems about counting how many events of some kind occur. Write

$$I_A = \begin{cases} 1 & \text{if } A \text{ occurs,} \\ 0 & \text{if } A \text{ does not occur.} \end{cases}$$

Note that $I_A^2 = I_A, I_A I_B = I_{A \cap B}$, and $I_{A \cup B} = I_A + I_B - I_A I_B$.

- Distribution $I_A \sim \text{Bern}(p)$ where $p = P(A)$.
- Fundamental Bridge The expectation of the indicator for event $A$ is the probability of event $A$: $E(I_A) = P(A)$.

**Variance and Standard Deviation of a Random Variable**

$$\text{Var}(X) = E\left(X - E(X)\right)^2 = E(X^2) - (E(X))^2$$

$$\text{SD}(X) = \sqrt{\text{Var}(X)}$$

**Continuous Random Variables**

A continuous random variable can take on any possible value within a certain interval (for example, $[0, 1]$), whereas a discrete random variable can only take on variables in a list of countable values (for example, all the integers, or the values $1, \frac{1}{2}, \frac{1}{4}, \frac{1}{8}$, etc.) * Do Continuous Random Variables have PMFs? No. The probability that a continuous random variable takes on any specific value is 0. * What's the probability that a CRV is in an interval? Take the difference in CDF values (or use the PDF as described later).

$$P(a \leq X \leq b) = P(X \leq b) - P(X \leq a) = F_X(b) - F_X(a)$$

For $X \sim \mathcal{N}(\mu, \sigma^2)$, this becomes

$$P(a \leq X \leq b) = \Phi\left(\frac{b - \mu}{\sigma}\right) - \Phi\left(\frac{a - \mu}{\sigma}\right) \tag{21}$$

* What is the Probability Density Function (PDF)? The PDF $f$ is the derivative of the CDF $F$.

$$F'(x) = f(x)$$

A PDF is nonnegative and integrates to 1. By the fundamental theorem of calculus, to get from PDF back to CDF we can integrate:

$$F(x) = \int_{-\infty}^{x} f(t)dt \tag{22}$$

To find the probability that a CRV takes on a value in an interval, integrate the PDF over that interval.

$$F(b) - F(a) = \int_a^b f(x)dx \tag{23}$$

Two additional properties of a PDF: it must integrate to 1 (because the probability that a CRV falls in the interval $[-\infty, \infty]$ is 1, and the PDF must always be nonnegative. * How do I find the expected value of a

4

CRV? Analogous to the discrete case, where you sum $x$ times the PMF, for CRVs you integrate $x$ times the PDF.

$$E(X) = \int_{-\infty}^{\infty} xf(x)dx$$

Expected value is *linear*. This means that for *any* random variables $X$ and $Y$ and any constants $a, b, c$, the following is true: %

$$E(aX + bY + c) = aE(X) + bE(Y) + c$$

**Expected value of a function of an random variable**

The expected value of $X$ is defined this way:

$$E(X) = \sum_x xP(X = x) \text{ (for discrete } X)$$

$$E(X) = \int_{-\infty}^{\infty} xf(x)dx \text{ (for continuous } X)$$

The **Law of the Unconscious Statistician (LOTUS)** states that you can find the expected value of a *function of a random variable*, $g(X)$, in a similar way, by replacing the $x$ in front of the PMF/PDF by $g(x)$ but still working with the PMF/PDF of $X$:

$$E(g(X)) = \sum_x g(x)P(X = x) \text{ (for discrete } X)$$

$$E(g(X)) = \int_{-\infty}^{\infty} g(x)f(x)dx \text{ (for continuous } X)$$

- What's a function of a random variable? A function of a random variable is also a random variable. For example, if $X$ is the number of bikes you see in an hour, then $g(X) = 2X$ is the number of bike wheels you see in that hour and $h(X) = \binom{X}{2} = \frac{X(X-1)}{2}$ is the number of *pairs* of bikes such that you see both of those bikes in that hour.

- What's the point? You don't need to know the PMF/PDF of $g(X)$ to find its expected value. All you need is the PMF/PDF of $X$.

**Joint Distributions**

The **joint CDF** of $X$ and $Y$ is

$$F(x, y) = P(X \leq x, Y \leq y)$$

In the discrete case, $X$ and $Y$ have a **joint PMF**

$$p_{X,Y}(x, y) = P(X = x, Y = y).$$

In the continuous case, they have a **joint PDF**

$$f_{X,Y}(x, y) = \frac{\partial^2}{\partial x \partial y} F_{X,Y}(x, y).$$

The joint PMF/PDF must be nonnegative and sum/integrate to 1.

## Conditional Distributions

### Conditioning and Bayes' rule for discrete r.v.s

$$P(Y = y|X = x) = \frac{P(X = x, Y = y)}{P(X = x)} = \frac{P(X = x|Y = y)P(Y = y)}{P(X = x)}$$

### Conditioning and Bayes' rule for continuous r.v.s

$$f_{Y|X}(y|x) = \frac{f_{X,Y}(x, y)}{f_X(x)} = \frac{f_{X|Y}(x|y)f_Y(y)}{f_X(x)}$$

### Hybrid Bayes' rule

$$f_X(x|A) = \frac{P(A|X = x)f_X(x)}{P(A)}$$

## Marginal Distributions

To find the distribution of one (or more) random variables from a joint PMF/PDF, sum/integrate over the unwanted random variables.

### Marginal PMF from joint PMF

$$P(X = x) = \sum_y P(X = x, Y = y)$$

### Marginal PDF from joint PDF

$$f_X(x) = \int_{-\infty}^{\infty} f_{X,Y}(x, y)dy$$

## Independence of Random Variables

Random variables $X$ and $Y$ are independent if and only if any of the following conditions holds: * Joint CDF is the product of the marginal CDFs * Joint PMF/PDF is the product of the marginal PMFs/PDFs * Conditional distribution of $Y$ given $X$ is the marginal distribution of $Y$ Write $X \perp\!\!\!\perp Y$ to denote that $X$ and $Y$ are independent.

## Multivariate LOTUS

LOTUS in more than one dimension is analogous to the univariate LOTUS. For discrete random variables:

$$E(g(X, Y)) = \sum_x \sum_y g(x, y)P(X = x, Y = y)$$

For continuous random variables:

$$E(g(X, Y)) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x, y)f_{X,Y}(x, y)dxdy$$

**Normal Distribution**

Let us say that $X$ is distributed $\mathcal{N}(\mu, \sigma^2)$. We know the following:

- Central Limit Theorem: The Normal distribution is ubiquitous because of the Central Limit Theorem, which states that the sample mean of i.i.d.~r.v.s will approach a Normal distribution as the sample size grows, regardless of the initial distribution.

- Location-Scale Transformation: Every time we shift a Normal r.v.~(by adding a constant) or rescale a Normal (by multiplying by a constant), we change it to another Normal r.v. For any Normal $X \sim \mathcal{N}(\mu, \sigma^2)$, we can transform it to the standard $\mathcal{N}(0,1)$ by the following transformation:

$$Z = \frac{X - \mu}{\sigma} \sim \mathcal{N}(0, 1)$$

- Standard Normal: The Standard Normal, $Z \sim \mathcal{N}(0,1)$, has mean 0 and variance 1. Its CDF is denoted by $\Phi$.

**Multivariate Normal (MVN) Distribution**

A vector $\vec{X} = (X_1, X_2, \ldots, X_k)$ is Multivariate Normal if every linear combination is Normally distributed, i.e., $t_1 X_1 + t_2 X_2 + \cdots + t_k X_k$ is Normal for any constants $t_1, t_2, \ldots, t_k$. The parameters of the Multivariate Normal are the **mean vector** $\vec{\mu} = (\mu_1, \mu_2, \ldots, \mu_k)$ and the **covariance matrix** where the $(i, j)$ entry is $\text{Cov}(X_i, X_j)$.

The Multivariate Normal has the following properties.

- Any subvector is also MVN.
- If any two elements within an MVN are uncorrelated, then they are independent.
- The joint PDF of a Bivariate Normal $(X, Y)$ with $\mathcal{N}(0,1)$ marginal distributions and correlation $\rho \in (-1, 1)$ is

$$f_{X,Y}(x, y) = \frac{1}{2\pi\tau} \exp\left( -\frac{1}{2\tau^2} (x^2 + y^2 - 2\rho xy) \right),$$

with $\tau = \sqrt{1 - \rho^2}$.

## Covariance and Correlation

- Covariance is the analog of variance for two random variables.

$$\text{Cov}(X, Y) = E\left((X - E(X))(Y - E(Y))\right) = E(XY) - E(X)E(Y)$$

  Note that

$$\text{Cov}(X, X) = E(X^2) - (E(X))^2 = \text{Var}(X)$$

- Correlation is a standardized version of covariance that is always between $-1$ and $1$.

$$\text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)\,\text{Var}(Y)}}$$

- Covariance and Independence: If two random variables are independent, then they are uncorrelated. The converse is not necessarily true (e.g., consider $X \sim \mathcal{N}(0, 1)$ and $Y = X^2$).

$$X \perp\!\!\!\perp Y \longrightarrow \text{Cov}(X, Y) = 0 \longrightarrow E(XY) = E(X)E(Y) \tag{24}$$

- Covariance and Variance:
  The variance of a sum can be found by

$$\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2\,\text{Cov}(X, Y) \tag{25}$$

$$\text{Var}(X_1 + X_2 + \cdots + X_n) = \sum_{i=1}^{n} \text{Var}(X_i) + 2\sum_{i<j} \text{Cov}(X_i, X_j) \tag{26}$$

  If $X$ and $Y$ are independent then they have covariance 0, so

$$X \perp\!\!\!\perp Y \implies \text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y)$$

  If $X_1, X_2, \ldots, X_n$ are identically distributed and have the same covariance relationships (often by **symmetry**), then

$$\text{Var}(X_1 + X_2 + \cdots + X_n) = n\,\text{Var}(X_1) + 2\binom{n}{2}\text{Cov}(X_1, X_2)$$

- Covariance Properties:
  For random variables $W, X, Y, Z$ and constants $a, b$:

$$\text{Cov}(X, Y) = \text{Cov}(Y, X)$$
$$\text{Cov}(X + a, Y + b) = \text{Cov}(X, Y)$$
$$\text{Cov}(aX, bY) = ab\,\text{Cov}(X, Y)$$
$$\text{Cov}(W + X, Y + Z) = \text{Cov}(W, Y) + \text{Cov}(W, Z) + \text{Cov}(X, Y) + \text{Cov}(X, Z)$$

- Correlation: is location-invariant and scale-invariant] For any constants $a, b, c, d$ with $a$ and $c$ nonzero,

$$\text{Corr}(aX + b, cY + d) = \text{Corr}(X, Y)$$

**Conditional Expectation**

**Conditioning on an Event.**  We can find $E(Y|A)$, the expected value of $Y$ given that event $A$ occurred. A very important case is when $A$ is the event $X = x$. Note that $E(Y|A)$ is a *number*. For example: * The expected value of a fair die roll, given that it is prime, is $\frac{1}{3} \cdot 2 + \frac{1}{3} \cdot 3 + \frac{1}{3} \cdot 5 = \frac{10}{3}$.

- Let $Y$ be the number of successes in 10 independent Bernoulli trials with probability $p$ of success. Let $A$ be the event that the first 3 trials are all successes. Then

$$E(Y|A) = 3 + 7p$$

   since the number of successes among the last 7 trials is $\text{Bin}(7, p)$.

- Let $T \sim \text{Expo}(1/10)$ be how long you have to wait until the shuttle comes.  Given that you have already waited $t$ minutes, the expected additional waiting time is 10 more minutes, by the memoryless property. That is, $E(T|T > t) = t + 10$.

| | **Discrete** $Y$ | **Continuous** $Y$ |
|---|---|---|
| | $E(Y) = \sum_y y P(Y = y)$ | $E(Y) = \int_{-\infty}^{\infty} y f_Y(y) dy$ |
| | $E(Y|A) = \sum_y y P(Y = y|A)$ | $E(Y|A) = \int_{-\infty}^{\infty} y f(y|A) dy$ |

**Conditioning on a Random Variable:**  We can also find $E(Y|X)$, the expected value of $Y$ given the random variable $X$. This is *a function of the random variable $X$*. It is *not* a number except in certain special cases such as if $X \perp\!\!\!\perp Y$. To find $E(Y|X)$, find $E(Y|X = x)$ and then plug in $X$ for $x$. For example: * If $E(Y|X = x) = x^3 + 5x$, then $E(Y|X) = X^3 + 5X$. * Let $Y$ be the number of successes in 10 independent Bernoulli trials with probability $p$ of success and $X$ be the number of successes among the first 3 trials. Then $E(Y|X) = X + 7p$. * Let $X \sim \mathcal{N}(0, 1)$ and $Y = X^2$. Then $E(Y|X = x) = x^2$ since if we know $X = x$ then we know $Y = x^2$. And $E(X|Y = y) = 0$ since if we know $Y = y$ then we know $X = \pm\sqrt{y}$, with equal probabilities (by symmetry). So $E(Y|X) = X^2, E(X|Y) = 0$.

**Properties of Conditional Expectation:**

- $E(Y|X) = E(Y)$ if $X \perp\!\!\!\perp Y$
- $E(h(X)W|X) = h(X)E(W|X)$ (**taking out what's known**) \ In particular, $E(h(X)|X) = h(X)$. \*$E(E(Y|X)) = E(Y)$ (**Adam's Law**, a.k.a.~Law of Total Expectation)

**Adam's Law (a.k.a.~Law of Total Expectation):**  The aaw of total expectation can also be written in a way that looks analogous to LOTP. For any events $A_1, A_2, \ldots, A_n$ that partition the sample space,

$$E(Y) = E(Y|A_1)P(A_1) + \cdots + E(Y|A_n)P(A_n) \tag{27}$$

For the special case where the partition is $A, A^c$, this says

$$E(Y) = E(Y|A)P(A) + E(Y|A^c)P(A^c) \tag{28}$$

**Eve's Law (a.k.a.~Law of Total Variance)]**

$$\text{Var}(Y) = E(\text{Var}(Y|X)) + \text{Var}(E(Y|X))$$

**Law of Large Numbers (LLN)**

Let $X_1, X_2, X_3 \ldots$ be i.i.d.~with mean $\mu$. The **sample mean** is

$$\bar{X}_n = \frac{X_1 + X_2 + X_3 + \cdots + X_n}{n}$$

The **Law of Large Numbers** states that as $n \to \infty$, $\bar{X}_n \to \mu$ with probability 1. For example, in flips of a coin with probability $p$ of Heads, let $X_j$ be the indicator of the $j$th flip being Heads. Then LLN says the proportion of Heads converges to $p$ (with probability 1).

**Central Limit Theorem (CLT)**

**Approximation using CLT**  We use $\dot{\sim}$ to denote *is approximately distributed*. We can use the **Central Limit Theorem** to approximate the distribution of a random variable $Y = X_1 + X_2 + \cdots + X_n$ that is a sum of $n$ i.i.d. random variables $X_i$. Let $E(Y) = \mu_Y$ and $\text{Var}(Y) = \sigma_Y^2$. The CLT says

$$Y \dot{\sim} \mathcal{N}(\mu_Y, \sigma_Y^2)$$

If the $X_i$ are i.i.d.~with mean $\mu_X$ and variance $\sigma_X^2$, then $\mu_Y = n\mu_X$ and $\sigma_Y^2 = n\sigma_X^2$. For the sample mean $\bar{X}_n$, the CLT says

$$\bar{X}_n = \frac{1}{n}(X_1 + X_2 + \cdots + X_n) \dot{\sim} \mathcal{N}(\mu_X, \sigma_X^2/n)$$

**Asymptotic Distributions using CLT**  We use $\xrightarrow{D}$ to denote *converges in distribution to* as $n \to \infty$. The CLT says that if we standardize the sum $X_1 + \cdots + X_n$ then the distribution of the sum converges to $\mathcal{N}(0,1)$ as $n \to \infty$:

$$\frac{1}{\sigma\sqrt{n}}(X_1 + \cdots + X_n - n\mu_X) \xrightarrow{D} \mathcal{N}(0,1)$$

In other words, the CDF of the left-hand side goes to the standard Normal CDF, $\Phi$. In terms of the sample mean, the CLT says

$$\frac{\sqrt{n}(\bar{X}_n - \mu_X)}{\sigma_X} \xrightarrow{D} \mathcal{N}(0,1)$$