
Probability & Random Processes

Sara Pohland

Created: January 29, 2021

Last Modified: December 14, 2023

Contents

I	Introduction to Probability	6
1	Experiments, Models, & Probability	7
1.1	Probability Spaces	7
1.2	Conditional Probability & Bayes' Rule	8
1.3	Independence	9
2	Sequential Experiments	10
2.1	Counting Methods	10
2.2	Independent Trials	10
2.2.1	Binomial Theorem	11
II	Random Variables	12
3	Discrete Random Variables	13
3.1	Discrete Random Variable	13
3.2	Probability Mass Function	13
3.3	Cumulative Distribution Function	13
3.4	Expected Value	14
3.4.1	Linearity of Expectation	15
3.4.2	Jensen's Inequality	15
3.4.3	Tail Sum Formula	15
3.5	Variance & Standard Deviation	16
3.6	Moment Generating Function	17
3.7	Common Discrete Distributions	17
3.7.1	Indicator Random Variable	17
3.7.2	Uniform Random Variable	17
3.7.3	Bernoulli Random Variable	18
3.7.4	Binomial Random Variable	18
3.7.5	Geometric Random Variable	19
3.7.6	Poisson Random Variable	19
3.7.7	Pascal Random Variable	20

4	Continuous Random Variables	21
4.1	Continuous Random Variable	21
4.2	Cumulative Distribution Function	21
4.3	Probability Density Function	22
4.4	Expected Value	22
4.4.1	Linearity of Expectation	23
4.4.2	Jensen's Inequality	23
4.4.3	Tail Sum Formula	23
4.5	Variance & Standard Deviation	23
4.6	Moment Generating Function	24
4.7	Common Continuous Distributions	24
4.7.1	Uniform Random Variable	25
4.7.2	Exponential Random Variable	25
4.7.3	Erlang Random Variable	26
4.7.4	Gaussian Random Variable	26
4.7.5	Standard Normal Random Variable	26
III	Multiple Random Variables	28
5	Multiple Random Variables	29
5.1	Joint Probability Mass Function	29
5.1.1	Two Random Variables	29
5.1.2	Multiple Random Variables	29
5.2	Joint Cumulative Distribution Function	30
5.2.1	Two Random Variables	30
5.2.2	Multiple Random Variables	31
5.3	Joint Probability Density Function	31
5.3.1	Two Random Variables	31
5.3.2	Multiple Random Variables	31
5.4	Expected Value	32
5.4.1	Discrete Random Variables	32
5.4.2	Continuous Random Variables	33
5.5	Variance, Covariance, & Correlation	33
5.5.1	Variance	33
5.5.2	Covariance	34
5.5.3	Correlation Coefficient	34
5.5.4	Correlation	35
5.6	Independent Random Variables	35
5.6.1	Discrete Random Variables	35
5.6.2	Continuous Random Variables	35
5.6.3	Independence Properties	36
5.7	Bivariate Gaussian Random Variables	36

6	Derived Random Variables	38
6.1	Derived Discrete Random Variable	38
6.1.1	Function of One Random Variable	38
6.1.2	Function of Two Random Variables	38
6.1.3	Sums of Independent Random Variables	38
6.2	Derived Continuous Random Variable	39
6.2.1	Function of One Random Variable	39
6.2.2	Function of Two Random Variables	39
6.2.3	Sums of Independent Random Variables	39
7	Conditional Probability Models	40
7.1	Conditioning by an Event	40
7.1.1	Conditional Probability Mass Function	40
7.1.2	Conditional Cumulative Distribution Function	41
7.1.3	Conditional Probability Density Function	41
7.1.4	Conditional Expected Value	42
7.1.5	Conditional Variance	43
7.2	Conditioning by a Random Variable	43
7.2.1	Conditional Probability Mass Function	43
7.2.2	Conditional Cumulative Distribution Function	43
7.2.3	Conditional Probability Density Function	44
7.2.4	Conditional Expected Value	44
7.2.5	Iterated Expectation & Tower Property	45
7.2.6	Conditional Variance	45
IV	Extensions of Random Variables	46
8	Random Vectors	47
8.1	Random Vector	47
8.2	Probability Distributions	47
8.2.1	Probability Mass Functions	47
8.2.2	Cumulative Distribution Functions	48
8.2.3	Probability Density Functions	48
8.3	Properties of Random Vectors	49
8.3.1	Expected Value Vector	49
8.3.2	Covariance & Correlation Matrices	49
8.3.3	Linear Transformation	50
8.4	Gaussian Random Vectors	50
8.4.1	Standard Normal Random Vector	50
9	Concentration Inequalities	52
9.1	Sample Mean	52
9.2	Probability Bounds	53
9.2.1	Markov Inequality	53
9.2.2	Chebyshev Inequality	54

CONTENTS

9.2.3	Chernoff Bound	54
9.3	Law of Large Numbers	55
9.3.1	Convergence of Random Variables	55
9.3.2	Weak Law of Large Numbers	56
9.3.3	Strong Law of Large Numbers	57
9.4	Central Limit Theorem	57

Part I

Introduction to Probability

Chapter 1

Experiments, Models, & Probability

1.1 Probability Spaces

A probability space (Ω, \mathcal{F}, P) is a mathematical construct that allows us to model an experiment. The **sample space** Ω is the set of all possible outcomes/observations of the experiment. The outcomes must be mutually exclusive and collectively exhaustive. The **event space** \mathcal{F} denotes the set of all possible events, where each event is a subset of the sample space containing zero or more outcomes. The **probability law** $P : \mathcal{F} \rightarrow [0, 1]$ assigns a probability between zero and one to each event in the event space.

Example: Consider tossing a fair coin twice. Each toss, we can get either head (H) or tail (T). Suppose we are interested in how many times we get a head.

$$\Omega = \{HH, HT, TH, TT\}$$

$$\mathcal{F} = \{n \text{ heads}, n \in \mathbb{Z}\}$$

$$P(n \text{ heads}) = \begin{cases} 1/4 & \text{if } n = 0 \\ 1/2 & \text{if } n = 1 \\ 1/4 & \text{if } n = 2 \\ 0 & \text{otherwise} \end{cases}$$

A probability law must satisfy the following **probability axioms**:

1. Non-negativity – $P(A) \geq 0, \forall A \in \mathcal{F}$
2. Normalization – $P(\Omega) = 1$
3. Additivity – $P(\cup_i A_i) = \sum_i P(A_i), \forall \text{ disjoint } A_i \in \mathcal{F}$

Probability laws also satisfy several other properties that can be derived from the probability axioms:

1. $P(\emptyset) = 0$
2. $P(A^C) = 1 - P(A)$, $\forall A \in \mathcal{F}$
3. $A \subset B \implies P(A) \leq P(B)$, $\forall A, B \in \mathcal{F}$
4. $P(A \cup B) = P(A) + P(B) - P(A \cap B)$, $\forall A, B \in \mathcal{F}$
5. $P(A \cup B \cup C) = P(A) + P(A^C \cap B) + P(A^C \cap B^C \cap C)$, $\forall A, B, C \in \mathcal{F}$
6. $P\left(\bigcup_{i=1}^n A_i\right) \leq \sum_{i=1}^n P(A_i)$, $\forall A_i \in \mathcal{F}$

1.2 Conditional Probability & Bayes' Rule

The **conditional probability** $P(A|B)$ is the probability that event A occurred given that event B with probability $P(B) > 0$ occurred, which is defined as

$$P(A|B) = \frac{P(A \cap B)}{P(B)}.$$

Notice that, from this definition, we can also say that the probability that event B occurred given that event A with probability $P(A) > 0$ occurred is

$$P(B|A) = \frac{P(A \cap B)}{P(A)}.$$

Now we have two expressions for the probability that both A and B occurred:

$$P(A \cap B) = P(A|B)P(B) \text{ and } P(B \cap A) = P(B|A)P(A).$$

This now allows us to express **Bayes' rule**, which says that given an event B with probability $P(B) > 0$ has occurred, the probability event A occurred is

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

Our definition of conditional probability also allows us to write the **law of total probability**, which says that if the events A_1, \dots, A_n are mutually exclusive and collectively exhaustive, then

$$P(B) = \sum_{i=1}^n P(A_i \cap B) = \sum_{i=1}^n P(B|A_i)P(A_i).$$

Combining the definition of conditional probability with the law of total probability, we get a more general definition of **Bayes' rule**, which says that if the events A_1, \dots, A_n are mutually exclusive and collectively exhaustive, then for any event B such that $P(B) > 0$, we have

$$P(A_j|B) = \frac{P(B|A_j)P(A_j)}{P(B)} = \frac{P(B|A_j)P(A_j)}{\sum_{i=1}^n P(B|A_i)P(A_i)}.$$

1.3 Independence

Two events A and B are said to be **independent** if the occurrence of one event provides no information about the occurrence of the other. For example, the event of drawing an Ace and the event of rolling a 5 are independent because the fact that I drew an Ace does not impact the probability that I roll a 5. More formally, events A and B are independent if and only if

$$P(A|B) = P(A) \text{ and } P(B|A) = P(B).$$

This is equivalent to saying that A and B are independent if and only if

$$P(A \cap B) = P(A)P(B).$$

Three events A_1 , A_2 , and A_3 are **mutually independent** if and only if

1. A_1 and A_2 are independent,
2. A_1 and A_3 are independent,
3. A_2 and A_3 are independent, and
4. $P(A_1 \cap A_2 \cap A_3) = P(A_1)P(A_2)P(A_3)$.

Extending the definition of mutual independence to n events, we can say that a collection of events A_1, \dots, A_n are **mutually independent** if and only if for all collections of events $S \subseteq \{1, \dots, n\}$,

$$P\left(\bigcap_{i \in S} A_i\right) = \prod_{i \in S} P(A_i).$$

We also can define conditional independence, which says that A and B are **conditionally independent** given C if and only if

$$P(A \cap B|C) = P(A|C)P(B|C).$$

Chapter 2

Sequential Experiments

2.1 Counting Methods

Given n total objects, there are $n!$ permutations of those objects. Consider a set of n total objects, and suppose we want to sample $k \leq n$ of those objects. If we choose k out of the n objects without replacement, then there are $(n)_k$ permutations of the objects we choose, where $(n)_k$ is defined as

$$(n)_k = \frac{n!}{(n-k)!}.$$

If we are not interested in the order in which the k objects are chosen, then we say there are $\binom{n}{k}$ ways to choose the k objects, where $\binom{n}{k}$ is defined as

$$\binom{n}{k} = \frac{(n)_k}{k!} = \frac{n!}{k!(n-k)!}.$$

A partition of a set of n objects is a set of m groups each containing n_i objects such that $n = n_1 + \dots + n_m$. The number of ways to choose n_i objects of each of the m groups from the set of n total objects is

$$\binom{n}{n_1, \dots, n_m} = \frac{n!}{n_1! n_2! \dots n_m!}.$$

Now consider the set of n total objects, and suppose we want to sample $k \leq n$ of those objects and replace each object as we sample them. Now there are n^k ways to choose an ordered sample of the objects.

2.2 Independent Trials

Independent trials are identical subexperiments in a sequential experiment. Consider an experiment of n independent trials, where the probability of success in each trial is p . Let an outcome of the experiment be defined as a particular

sequence of successes and failures, resulting in a total of n_1 successes and n_2 failures, where $n_1 + n_2 = n$. The probability of one particular outcome is

$$p^{n_1}(1-p)^{n_2}.$$

There are $\binom{n}{n_1} = \binom{n}{n_2}$ possible outcomes that include n_1 successes and $n_2 = n - n_1$ failures. Therefore, the probability of an experiment of n trials, each with probability of success p , resulting in n_1 successes is

$$\binom{n}{n_1} p^{n_1} (1-p)^{n-n_1}.$$

Now we can generalize this result to experiments with more than two possible outcomes. Consider an experiment which has a sample space $\Omega = \{s_1, \dots, s_m\}$ and the associated probability law $P(s_i) = p_i$. For $n = n_1 + \dots + n_m$ independent trials, the probability of n_i occurrences of s_i for $i \in \{1, \dots, m\}$ is given by

$$\binom{n}{n_1, \dots, n_m} \prod_{i=1}^m p_i^{n_i}.$$

2.2.1 Binomial Theorem

As an aside, it is useful to note a related concept, which is sometimes referred to as the **binomial theorem** or the **binomial expansion**. According to the binomial theorem, we can express $(x + y)^n$ for some integer $n \geq 0$ as

$$(x + y)^n = \sum_{k=0}^n \binom{n}{k} x^k y^{(n-k)}.$$

Part II

Random Variables

Chapter 3

Discrete Random Variables

3.1 Discrete Random Variable

Given an experiment with probability measure P defined on a sample space Ω , a **random variable** is a function that assigns a real number to each outcome. A random variable is a real-valued function $X : \Omega \rightarrow \mathbb{R}$ defined such that

$$\{X = x\} = \{\omega \in \Omega : X(\omega) = x\}.$$

In this formulation, the variable x is called a **realization** of the random variable X . A random variable X is discrete if the range of X , which we denote \mathcal{X} , is a countable set. This set may contain a finite or an infinite number of values.

3.2 Probability Mass Function

The frequencies with which a random variable takes on different values is described by its **probability mass function (PMF)**. The PMF, $p_X : \mathcal{X} \rightarrow [0, 1]$, associated with a discrete random variable X is defined as

$$p_X(x) = P(\{X = x\}) = P(\{\omega \in \Omega : X(\omega) = x\}).$$

From the probability axioms, the PMF satisfies the following properties:

1. $p_X(x) \geq 0, \forall x \in \mathcal{X}$
2. $\sum_{x \in \mathcal{X}} p_X(x) = 1$
3. $P(A) = \sum_{x \in A} p_X(x), \forall A \in \mathcal{F}$

3.3 Cumulative Distribution Function

The frequencies with which a random variable takes on different values can also be described by its **cummulative distribution function (CDF)**. The CDF,

$F_X : \mathcal{X} \rightarrow [0, 1]$, associated with a discrete random variable X is defined as

$$F_X(x) = P(\{X \leq x\}) = P(\{\omega \in \Omega : X(\omega) \leq x\}).$$

If we assume the range is given by $\mathcal{X} = \{x_1, x_2, \dots\}$, where $x_1 < x_2 < \dots$, then the CDF of X satisfies the following properties:

1. $F_X(x) = \sum_{x_i \leq x} p_X(x_i)$, $\forall x \in \mathcal{X}$
2. $F_X(-\infty) := \lim_{x \rightarrow -\infty} F_X(x) = 0$
3. $F_X(\infty) := \lim_{x \rightarrow \infty} F_X(x) = 1$
4. $F_X(x') \geq F_X(x)$, $\forall x, x' \in \mathcal{X}$ s.t. $x' \geq x$
5. $p_X(x_i) = F_X(x_i) - F_X(x_i - \epsilon)$, $\forall x_i \in \mathcal{X}$ and some arbitrarily small $\epsilon > 0$
6. $F_X(x) = F_X(x_i)$, $\forall x, x_i, x_{i+1} \in \mathcal{X}$ s.t. $x_i < x < x_{i+1}$
7. $F_X(b) - F_X(a) = P(\{a < X \leq b\})$, $\forall a, b \in \mathcal{X}$

These properties say that the CDF, F_X , starts at zero and increases monotonically to one. Furthermore, there is a discontinuity at each value $x_i \in \mathcal{X}$, and the height of the jump is given by $p_X(x_i)$. Between these discontinuities, the CDF does not change. Finally, the CDF allows us to determine the probability that the random variable X takes on a value between two points.

3.4 Expected Value

For a discrete random variable X with PMF p_X , the **expected value** of X is

$$\mathbb{E}[X] = \sum_{x \in \mathcal{X}} xp_X(x).$$

Given a discrete random variable X with the PMF p_X , $Y = g(X)$ is another discrete random variable with its own PMF, p_Y , which is defined as

$$p_Y(y) = \sum_{x:g(x)=y} p_X(x).$$

This says that $p_Y(y)$ is the sum of probability outcomes $X = x$ for which $Y = y$. For these discrete random variables, the expected value of Y can be defined as

$$\begin{aligned} \mathbb{E}[Y] &= \sum_{y \in \mathcal{Y}} yp_Y(y) = \sum_{y \in \mathcal{Y}} y \left(\sum_{x:g(x)=y} p_X(x) \right) \\ &= \sum_{y \in \mathcal{Y}} \sum_{x:g(x)=y} g(x)p_X(x) = \sum_{x \in \mathcal{X}} g(x)p_X(x). \end{aligned}$$

Now we can see that we compute $\mathbb{E}[Y]$ without needing to compute $p_Y(y)$. This finding is sometimes referred to as the **law of the unconscious statistician**.

3.4.1 Linearity of Expectation

This law of the unconscious statistician leads us to the linearity of expectation, which says that for the discrete random variable X and two constants a and b ,

$$\mathbb{E}[aX + b] = a\mathbb{E}[X] + b.$$

Proof: From the law of the unconscious statistician,

$$\mathbb{E}[aX + b] = \sum_{x \in \mathcal{X}} (ax + b)p_X(x) = a \sum_{x \in \mathcal{X}} xp_X(x) + b \sum_{x \in \mathcal{X}} p_X(x)$$

Using the definition of the expected value for the first term of the sum and the second property of the PMF for the second term, we get the desired result.

3.4.2 Jensen's Inequality

Another useful property of expected value is **Jensen's inequality**, which says that for any discrete random variable X and convex function f ,

$$f(\mathbb{E}[X]) \leq \mathbb{E}[f(X)].$$

Proof: According to the first order condition for convexity (see convex optimization notes), if f is a convex function, then

$$f(y) \geq f(x) + f'(x)(y - x), \quad \forall x, y \in \text{dom}f.$$

If we let $y = X$ and $x = \mathbb{E}[X]$, then this inequality becomes

$$f(X) \geq f(\mathbb{E}[X]) + f'(\mathbb{E}[X])(X - \mathbb{E}[X]).$$

Taking the expectation of both sides and using the linearity of expectation,

$$\begin{aligned} \mathbb{E}[f(X)] &\geq \mathbb{E}[f(\mathbb{E}[X]) + f'(\mathbb{E}[X])(X - \mathbb{E}[X])] \\ &= f(\mathbb{E}[X]) + f'(\mathbb{E}[X])(\mathbb{E}[X] - \mathbb{E}[X]) \\ &= f(\mathbb{E}[X]). \end{aligned}$$

3.4.3 Tail Sum Formula

The tail sum formula for expectation says that for a non-negative, integer-valued random variable X , the expected value is given by

$$\mathbb{E}[X] = \sum_{k=1}^{\infty} P(\{X \geq k\}).$$

Proof: From the definition of the expected value, we know

$$\mathbb{E}[X] = \sum_{x \in \mathcal{X}} xp_X(x) = \sum_{x \in \mathcal{X}} xP(\{X = x\}).$$

If we assume the X is non-negative and integer-valued, then

$$\begin{aligned}\mathbb{E}[X] &= \sum_{x=1}^{\infty} xP(\{X = x\}) = \sum_{x=1}^{\infty} \sum_{k=1}^x P(\{X = x\}) \\ &= \sum_{k=1}^{\infty} \sum_{x=k}^{\infty} P(\{X = x\}) = \sum_{k=1}^{\infty} P(\{X \geq k\}).\end{aligned}$$

3.5 Variance & Standard Deviation

For a discrete random variable X , the **variance** measures the dispersion of its sample values around its expected value, $\mathbb{E}[X]$. The variance of X is defined as

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

We can also express the variance of X using a different formula:

$$\text{Var}(X) = \mathbb{E}[X^2] - (\mathbb{E}[X])^2.$$

Proof: Using the properties of expectation discussed in the previous section, we can show that these expressions are equivalent:

$$\begin{aligned}\text{Var}(X) &= \mathbb{E} \left[(X - \mathbb{E}[X])^2 \right] \\ &= \sum_{x \in \mathcal{X}} (x - \mathbb{E}[X])^2 p_X(x) \\ &= \sum_{x \in \mathcal{X}} (x^2 - 2x\mathbb{E}[X] + (\mathbb{E}[X])^2) p_X(x) \\ &= \sum_{x \in \mathcal{X}} x^2 p_X(x) - 2\mathbb{E}[X] \sum_{x \in \mathcal{X}} x p_X(x) + (\mathbb{E}[X])^2 \sum_{x \in \mathcal{X}} p_X(x) \\ &= \mathbb{E}[X^2] - 2(\mathbb{E}[X])^2 + (\mathbb{E}[X])^2 \\ &= \mathbb{E}[X^2] - (\mathbb{E}[X])^2\end{aligned}$$

For the random variable X and two constants a and b , the variance satisfies

$$\text{Var}(aX + b) = a^2 \text{Var}(X).$$

Proof: We can prove this using the second expression for the variance:

$$\begin{aligned}\text{Var}(aX + b) &= \mathbb{E}[(aX + b)^2] - (\mathbb{E}[aX + b])^2 \\ &= \mathbb{E}[a^2 X^2 + 2abX + b^2] - (\mathbb{E}[aX + b])^2 \\ &= a^2 \mathbb{E}[X^2] + 2ab\mathbb{E}[X] + b^2 - (a\mathbb{E}[X] + b)^2 \\ &= a^2 \mathbb{E}[X^2] + 2ab\mathbb{E}[X] + b^2 - a^2(\mathbb{E}[X])^2 - 2ab\mathbb{E}[X] - b^2 \\ &= a^2 \mathbb{E}[X^2] - a^2(\mathbb{E}[X])^2 \\ &= a^2 \left(\mathbb{E}[X^2] - (\mathbb{E}[X])^2 \right) = a^2 \text{Var}(X)\end{aligned}$$

Another measure of the dispersion of the sample values of a random variable about its expected value is **standard deviation**, which is defined as

$$\sigma_X = \sqrt{\text{Var}(X)}.$$

3.6 Moment Generating Function

Another probability model for the discrete random variable X is the **moment generating function (MGF)**, which is defined for real values of s as

$$\phi_X(s) = \mathbb{E}[e^{sX}] = \sum_{x \in \mathcal{X}} e^{sx} p_X(x).$$

The set of values of s for which $\phi_X(s)$ exists is called the **region of convergence**. MGFs can be used to discuss **moments** of random variables. A discrete random variable X with MGF ϕ_X has the n th moment

$$\mathbb{E}[X^n] = \left. \frac{d^n}{ds^n} \phi_X(s) \right|_{s=0}.$$

3.7 Common Discrete Distributions

There are several common types of discrete random variables with well-known distributions. Some common discrete distributions are discussed in this section.

3.7.1 Indicator Random Variable

An indicator random variable X for an event $A \in \mathcal{F}$ is defined as

$$X = \mathbb{1}_A(x) = \begin{cases} 1 & \text{if } x \in A \\ 0 & \text{otherwise} \end{cases}.$$

The expected value for this random variable is given by

$$\mathbb{E}[X] = P(A).$$

3.7.2 Uniform Random Variable

The PMF of a random variable $X \sim \text{Uniform}(a, b)$ for $a < b$ is given by

$$p_X(x) = \begin{cases} \frac{1}{b-a+1} & \text{if } x \in \{a, a+1, \dots, b-1, b\} \\ 0 & \text{otherwise} \end{cases}.$$

Uniform random variables are used if every possible outcome has the same probability. The expected value and variance for this random variable are

$$\mathbb{E}[X] = \frac{a+b}{2} \quad \text{and} \quad \text{Var}(X) = \frac{(b-a)(b-a+1)}{12}.$$

The moment generating function for this random variable is given by

$$\phi_X(s) = \frac{e^{sk} - e^{s(l+1)}}{1 - e^s}.$$

All these expressions can be derived from their definitions using the PMF of X .

3.7.3 Bernoulli Random Variable

The PMF of a random variable $X \sim \text{Bernoulli}(p)$ for $0 \leq p \leq 1$ is given by

$$p_X(x) = \begin{cases} 1 - p & \text{if } x = 0 \\ p & \text{if } x = 1 \\ 0 & \text{otherwise} \end{cases}.$$

Bernoulli random variables are used to represent the probability of a success or failure in a single binary experiment, where p is the probability of success. The expected value and variance for this random variable are given by

$$\mathbb{E}[X] = p \quad \text{and} \quad \text{Var}(X) = p(1 - p).$$

The moment generating function for this random variable is given by

$$\phi_X(s) = 1 - p + pe^s.$$

All these expressions can be derived from their definitions using the PMF of X .

3.7.4 Binomial Random Variable

The PMF of a random variable $X \sim \text{Binomial}(n, p)$ for $n \in \mathbb{Z}$ and $0 \leq p \leq 1$ is

$$p_X(x) = \begin{cases} \binom{n}{x} p^x (1 - p)^{n-x} & \text{if } x = 0, \dots, n \\ 0 & \text{otherwise} \end{cases}.$$

Binomial random variables are used to represent the probability of some number of successes in n binary trials, where p is the probability of success in any given trial. The expected value and variance for this random variable are given by

$$\mathbb{E}[X] = np \quad \text{and} \quad \text{Var}(X) = np(1 - p).$$

The moment generating function for this random variable is given by

$$\phi_X(s) = (1 - p + pe^s)^n.$$

All these expressions can be derived from their definitions using the PMF of X . If X_1, \dots, X_n are n independent and identically distributed (i.i.d.) (See Section 5.6.1) random variables such that $X_i \sim \text{Bernoulli}(p)$ for $i = \{1, \dots, n\}$, then

$$X = \sum_{i=1}^n X_i.$$

If $X \sim \text{Binomial}(n, p)$ and $Y \sim \text{Binomial}(m, p)$ are independent random variables, then by construction, $X + Y \sim \text{Binomial}(n + m, p)$.

3.7.5 Geometric Random Variable

The PMF of a random variable $X \sim \text{Geometric}(p)$ for $0 \leq p \leq 1$ is given by

$$p_X(x) = \begin{cases} p(1-p)^{x-1} & \text{if } x = 1, 2, \dots \\ 0 & \text{otherwise} \end{cases}.$$

Geometric random variables are used to represent to the number of failed trials before the first success, where p is the probability of success in any given trial. The expected value and variance for this random variable are given by

$$\mathbb{E}[X] = \frac{1}{p} \quad \text{and} \quad \text{Var}(X) = \frac{1-p}{p^2}.$$

The moment generating function for this random variable is given by

$$\phi_X(s) = \frac{pe^s}{1 - (1-p)e^s}.$$

All these expressions can be derived from their definitions using the PMF of X . If X_1, \dots, X_n are n independent and identically distributed (i.i.d.) (See Section 5.6.1) random variables such that $X_i \sim \text{Bernoulli}(p)$ for $i = \{1, \dots, n\}$, then

$$X = \min\{k \geq 1 : X_k = 1\}$$

The geometric random variable X is **memoryless**, which means

$$P(\{X > t + s | X > s\}) = P(\{X > t\}), \quad \forall s, t \geq 0.$$

3.7.6 Poisson Random Variable

The PMF of a random variable $X \sim \text{Poisson}(\lambda)$ for $\lambda > 0$ is given by

$$p_X(x) = \begin{cases} \frac{\lambda^x e^{-\lambda}}{x!} & \text{if } x = 0, 1, \dots \\ 0 & \text{otherwise} \end{cases}.$$

Poisson random variables are often used to represent the number of discrete events that occur within a countably infinite interval of time. The expected value and variance for this random variable are given by

$$\mathbb{E}[X] = \lambda \quad \text{and} \quad \text{Var}(X) = \lambda.$$

These expressions can be derived from their definitions using the PMF of X . By construction of the Poisson distribution, if $X \sim \text{Poisson}(\lambda)$ and $Y \sim \text{Poisson}(\mu)$ are independent random variables, then the sum of the random variables has the distribution $X + Y \sim \text{Poisson}(\lambda + \mu)$.

3.7.7 Pascal Random Variable

The PMF of a random variable $X \sim \text{Pascal}(k, p)$ is given by

$$p_X(x) = \begin{cases} \binom{x-1}{k-1} p^k (1-p)^{x-k} & \text{if } x = 1, 2, \dots \\ 0 & \text{otherwise} \end{cases}.$$

The expected value and variance for this random variable are given by

$$\mathbb{E}[X] = \frac{k}{p} \quad \text{and} \quad \text{Var}(X) = \frac{k(1-p)}{p^2}.$$

These expressions can be derived from their definitions using the PMF of X.

Chapter 4

Continuous Random Variables

4.1 Continuous Random Variable

Previously, we said that given an experiment with probability measure P defined on a sample space Ω , a random variable is a function that assigns a real number to each outcome in the sample space of the experiment. We defined a random variable as a real-valued function $X : \Omega \rightarrow \mathbb{R}$ such that

$$\{X = x\} = \{\omega \in \Omega : X(\omega) = x\}.$$

A random variable X is continuous if the range of X , which we denote \mathcal{X} , is not a countable set. Instead, the range consists of one or more intervals of values. Because the probability function is defined over an uncountable range of values, the probability that X takes on a value of exactly x is zero:

$$P(\{X = x\}) = 0.$$

For continuous random variables, we consider the probability that it takes on a value within an interval, rather than taking on an exact value.

4.2 Cumulative Distribution Function

The frequencies with which a continuous random variable takes on different values can be described by its **cumulative distribution function (CDF)**. The CDF, $F_X : \mathcal{X} \rightarrow [0, 1]$, associated with a random variable X is defined as

$$F_X(x) = P(\{X \leq x\}) = P(\{\omega \in \Omega : X(\omega) \leq x\}).$$

The CDF of a continuous random variable satisfies the following properties:

1. $F_X(-\infty) := \lim_{x \rightarrow -\infty} F_X(x) = 0$

2. $F_X(\infty) := \lim_{x \rightarrow \infty} F_X(x) = 1$
3. $F_X(x') \geq F_X(x), \forall x, x' \in \mathcal{X} \text{ s.t. } x' \geq x$
4. $F_X(b) - F_X(a) = P(\{a < X \leq b\})$

Note that these properties also hold if X is a discrete random variable. However, if X is discrete, then the CDF has discontinuities that must satisfy additional properties. For a continuous random variable, the CDF is continuous, so it does not need to satisfy these additional properties.

4.3 Probability Density Function

The frequencies with which a continuous random variable takes on different values can also be described by its **probability density function (PDF)**. The PDF, f_X , of a random variable X is defined as the derivative of its CDF:

$$f_X(x) = \frac{d}{dx}F_X(x).$$

The value of $f_X(x)$ indicates the probability that X is near a sample value x and is a good indication of the likely value of observations. The PDF of the continuous random variable X satisfies the following properties:

1. $f_X(x) \geq 0, \forall x \in \mathcal{X}$
2. $\int_{\mathcal{X}} f_X(x)dx = 1$
3. $F_X(x) = \int_{-\infty}^x f_X(t)dt, \forall x \in \mathcal{X}$
4. $\int_a^b f_X(x)dx = P(\{a < X \leq b\}), \forall a, b \in \mathcal{X}$

4.4 Expected Value

Similar to the discrete case, for a continuous random variable X with the probability density function (PDF) f_X , the **expected value** of X is defined as

$$\mathbb{E}[X] = \int_{\mathcal{X}} x f_X(x)dx.$$

Given a continuous random variable X with the PDF f_X , $Y = g(X)$ is another random variable, but it is not necessarily continuous, so it may not have a well-defined PDF. However, similar to as we showed in the discrete case, the expected value of Y can be defined in terms of the distribution of X as

$$\mathbb{E}[Y] = \int_{\mathcal{X}} g(x)f_X(x).$$

As in the discrete case, we can compute $\mathbb{E}[Y]$ without needing to compute $f_Y(y)$, which is called the **law of the unconscious statistician**.

4.4.1 Linearity of Expectation

In the same way as we saw in the discrete case, the law of the unconscious statistician leads us to the linearity of expectation, which says that for the continuous random variable X and two constants a and b ,

$$\mathbb{E}[aX + b] = a\mathbb{E}[X] + b.$$

Proof: From the law of the unconscious statistician,

$$\mathbb{E}[aX + b] = \int_{\mathcal{X}} (ax + b)f_X(x)dx = a \int_{\mathcal{X}} xf_X(x)dx + b \int_{\mathcal{X}} f_X(x)dx$$

Using the definition of the expected value for the first term of the sum and the second property of the PDF for the second term, we get the desired result.

4.4.2 Jensen's Inequality

Jensen's inequality also holds for continuous random variables and says that for any continuous random variable X and convex function f ,

$$f(\mathbb{E}[X]) \leq \mathbb{E}[f(X)].$$

The proof for this inequality is exactly the same as in the discrete case.

4.4.3 Tail Sum Formula

The tail sum formula for expectation also applies to continuous random variables and says that for a non-negative random variable X , the expected value is

$$\mathbb{E}[X] = \int_0^{\infty} P(\{X \geq t\}) dt = \int_0^{\infty} (1 - F_X(t)) dt.$$

More more information about the tail sum formula, see: https://stat.uiowa.edu/sites/stat.uiowa.edu/files/cae/Lo_Expectation.pdf

4.5 Variance & Standard Deviation

For a continuous random variable X , its **variance** measures the dispersion of sample values of X around its expected value $\mathbb{E}[X]$. The variance of X is

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

As in the discrete case, we can also express the variance as

$$\text{Var}(X) = \mathbb{E}[X^2] - (\mathbb{E}[X])^2.$$

Proof: Using the properties of expectation discussed in the previous section, we can show that these expressions are equivalent:

$$\begin{aligned}\text{Var}(X) &= \mathbb{E} \left[(X - \mathbb{E}[X])^2 \right] \\ &= \int_{\mathcal{X}} (x - \mathbb{E}[X])^2 f_X(x) dx \\ &= \int_{\mathcal{X}} \left(x^2 - 2x\mathbb{E}[X] + (\mathbb{E}[X])^2 \right) f_X(x) dx \\ &= \int_{\mathcal{X}} x^2 f_X(x) - 2\mathbb{E}[X] \int_{\mathcal{X}} x f_X(x) + (\mathbb{E}[X])^2 \int_{\mathcal{X}} f_X(x) \\ &= \mathbb{E}[X^2] - 2(\mathbb{E}[X])^2 + (\mathbb{E}[X])^2 \\ &= \mathbb{E}[X^2] - (\mathbb{E}[X])^2\end{aligned}$$

An important property of variance, which we showed in the discrete case, is that for the continuous random variable X and two constants a and b ,

$$\text{Var}(aX + b) = a^2 \text{Var}(X).$$

The proof for the statement is the same as was shown for the discrete case. Another measure of the dispersion of the sample values of a random variable about its expected value is standard deviation, which is defined as

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

4.6 Moment Generating Function

Another probability model for the continuous random variable X is the **moment generating function (MGF)**, which is defined for real values of s as

$$\phi_X(s) = \mathbb{E}[e^{sX}] = \int_{\mathcal{X}} e^{sx} f_X(x) dx.$$

The set of values of s for which $\phi_X(s)$ exists is called the **region of convergence**. A continuous random variable X with MGF $\phi_X(s)$ has the n th moment

$$\mathbb{E}[X^n] = \left. \frac{d^n}{ds^n} \phi_X(s) \right|_{s=0}.$$

4.7 Common Continuous Distributions

There are several common types of continuous random variables with well-known distributions. Some common continuous distributions are discussed here.

4.7.1 Uniform Random Variable

The PDF of a random variable $X \sim \text{Uniform}(a, b)$ for $a < b$ is given by

$$f_X(x) = \begin{cases} \frac{1}{b-a} & \text{if } a \leq x < b \\ 0 & \text{otherwise} \end{cases}.$$

Uniform random variables are used to model situations in which there is an equal chance of finding an outcome x in any given interval on (a, b) . The CDF for the uniform random variable X is given by

$$F_X(x) = \begin{cases} 0 & \text{if } x \leq a \\ \frac{x-a}{b-a} & \text{if } a < x \leq b \\ 1 & \text{if } x > b \end{cases}.$$

The expected value and variance for this random variable are

$$\mathbb{E}[X] = \frac{a+b}{2} \quad \text{and} \quad \text{Var}(X) = \frac{(b-a)^2}{12}.$$

The moment generating function for this random variable is given by

$$\phi_X(s) = \frac{e^{bs} - e^{as}}{s(b-a)}.$$

All these expressions can be derived from their definitions using the PDF of X .

4.7.2 Exponential Random Variable

The PDF of a random variable $X \sim \text{Exponential}(\lambda)$ for $\lambda > 0$ is given by

$$f_X(x) = \begin{cases} 0 & \text{if } x < 0 \\ \lambda e^{-\lambda x} & \text{if } x \geq 0 \end{cases}.$$

Exponential random variables are used to model the amount of time that passes before an event occurs. The CDF for the exponential random variable X is

$$F_X(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 - e^{-\lambda x} & \text{if } x \geq 0 \end{cases}.$$

The expected value and variance for this random variable are

$$\mathbb{E}[X] = \frac{1}{\lambda} \quad \text{and} \quad \text{Var}(X) = \frac{1}{\lambda^2}.$$

The moment generating function for this random variable is given by

$$\phi_X(s) = \frac{\lambda}{\lambda - s}.$$

All these expressions can be derived from their definitions using the PDF of X . The exponential random variable X is **memoryless**, which means

$$P(\{X > t + s | X > s\}) = P(\{X > t\}), \quad \forall s, t \geq 0.$$

4.7.3 Erlang Random Variable

The PDF of a random variable $X \sim \text{Erlang}(n, \lambda)$ for $\lambda > 0$ and integer $n \geq 1$ is

$$f_X(x) = \begin{cases} 0 & \text{if } x < 0 \\ \frac{\lambda^n x^{n-1} e^{-\lambda x}}{(n-1)!} & \text{if } x \geq 0 \end{cases}.$$

The expected value and variance for this random variable are given by

$$\mathbb{E}[X] = \frac{n}{\lambda} \quad \text{and} \quad \text{Var}(X) = \frac{n}{\lambda^2}.$$

These expressions can be derived from their definitions using the PDF of X . An Erlang random variable with parameters $\lambda > 0$ and $n = 1$ is an exponential random variable with parameter λ .

4.7.4 Gaussian Random Variable

The PDF of a random variable $X \sim N(\mu, \sigma^2)$ for real number μ and $\sigma > 0$ is

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

The graph of $f_X(x)$ has a bell shape and is centered at μ . The parameter σ reflects the width of the bell such that smaller values of σ correspond to a narrow bell shape with a high peak and larger values of σ correspond to a wide bell with a low peak. The height of the peak is given by $1/\sqrt{2\pi\sigma^2}$. The expected value and variance for this random variable are given by

$$\mathbb{E}[X] = \mu \quad \text{and} \quad \text{Var}(X) = \sigma^2.$$

The moment generating function for this random variable is given by

$$\phi_X(s) = \exp\left(s\mu + \frac{1}{2}s^2\sigma^2\right).$$

All these expressions can be derived from their definitions using the PDF of X .

4.7.5 Standard Normal Random Variable

The standard normal random variable $Z \sim N(0, 1)$ is the Gaussian random variable with zero mean and unit variance, which has the PDF

$$f_Z(z) = \frac{1}{\sqrt{2\pi}} e^{-z^2/2}.$$

To prove that this is actually a valid PDF, first notice that this function is non-negative for all values of z . Now we need to show $\int_{-\infty}^{\infty} f_Z(z) dz = 1$. To do

so, we will consider the square of the given integral:

$$\begin{aligned}
 \left(\int_{-\infty}^{\infty} f_Z(z) dz \right)^2 &= \left(\int_{-\infty}^{\infty} f_Z(u) du \right) \left(\int_{-\infty}^{\infty} f_Z(v) dv \right) \\
 &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_Z(u) f_Z(v) dudv \\
 &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \left(\frac{1}{\sqrt{2\pi}} e^{-u^2/2} \right) \left(\frac{1}{\sqrt{2\pi}} e^{-v^2/2} \right) dudv \\
 &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \frac{1}{2\pi} e^{-(u^2+v^2)/2} dudv \\
 &= \int_0^{2\pi} \int_0^{\infty} \frac{1}{2\pi} e^{-r^2/2} r dr d\theta \\
 &= \int_0^{2\pi} \frac{1}{2\pi} d\theta \int_0^{\infty} e^{-r^2/2} r dr \\
 &= \frac{1}{2\pi} (2\pi) \int_0^{\infty} e^{-t} dt \\
 &= 1 \quad \checkmark
 \end{aligned}$$

The CDF for the standard normal random variable Z is given by

$$\phi(z) := F_Z(z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-u^2/2} du.$$

The CDF for Z satisfies an important symmetrical property:

$$\phi(-z) = 1 - \phi(z).$$

The standard normal complementary CDF is given by

$$Q(z) = 1 - \phi(z) = \frac{1}{\sqrt{2\pi}} \int_z^{\infty} e^{-u^2/2} du.$$

To transform a Gaussian random variable $X \sim N(\mu, \sigma^2)$ into a standard normal random variable $Z \sim N(0, 1)$, we can define Z as

$$Z = \frac{X - \mu}{\sigma}.$$

We can then define the probability distributions of the Gaussian random variable X with mean μ and standard deviation σ in terms of the standard normal CDF:

$$F_X(x) = \Phi\left(\frac{x - \mu}{\sigma}\right) \quad \text{and} \quad P(\{a < X \leq b\}) = \Phi\left(\frac{b - \mu}{\sigma}\right) - \Phi\left(\frac{a - \mu}{\sigma}\right).$$

Part III

Multiple Random Variables

Chapter 5

Multiple Random Variables

5.1 Joint Probability Mass Function

5.1.1 Two Random Variables

Let X and Y be two discrete random variables defined on the same probability space, (Ω, \mathcal{F}, P) . Suppose the range of X and Y are \mathcal{X} and \mathcal{Y} respectively. Their **joint probability mass function (PMF)** describes the frequencies of their joint outcomes. The joint PMF of X and Y is defined as

$$\begin{aligned} p_{XY}(x, y) &= P(\{X = x, Y = y\}) \\ &= P(\{\omega \in \Omega : X(\omega) = x, Y(\omega) = y\}) \\ &= P(\{\omega \in \Omega : X(\omega) = x\} \cap \{\omega \in \Omega : Y(\omega) = y\}). \end{aligned}$$

The joint probability mass function satisfies the following properties:

1. $p_{XY}(x, y) \geq 0, \forall x \in \mathcal{X}, y \in \mathcal{Y}$
2. $\sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p_{XY}(x, y) = 1$
3. $P(A) = \sum_{(x, y) \in A} p_{XY}(x, y), \forall A \in \mathcal{F}$

For random variables X and Y with joint PMF p_{XY} , the **marginal PMFs** are probability models for the individual random variables and are defined as

$$p_X(x) = \sum_{y \in \mathcal{Y}} p_{XY}(x, y) \quad \text{and} \quad p_Y(y) = \sum_{x \in \mathcal{X}} p_{XY}(x, y).$$

5.1.2 Multiple Random Variables

We can also expand the definition of the joint probability mass function (PMF) to n random variables. Let X_1, \dots, X_n be n discrete random variables defined

on the same probability space, (Ω, \mathcal{F}, P) . Suppose the range of X_i is \mathcal{X}_i for $i = 1, \dots, n$. For the random variables X_1, \dots, X_n , the joint PMF is

$$p_{X_1 \dots X_n}(x_1, \dots, x_n) = P(\{X_1 = x_1, \dots, X_n = x_n\}).$$

For n random variables, the joint PMF satisfies the following properties:

1. $p_{X_1 \dots X_n}(x_1, \dots, x_n) \geq 0, \forall x_1 \in \mathcal{X}_1, \dots, x_n \in \mathcal{X}_n$
2. $\sum_{x_1 \in \mathcal{X}_1} \dots \sum_{x_n \in \mathcal{X}_n} p_{X_1 \dots X_n}(x_1, \dots, x_n) = 1$
3. $P(A) = \sum_{(x_1, \dots, x_n) \in A} p_{X_1 \dots X_n}(x_1, \dots, x_n), \forall A \in \mathcal{F}$

Consider four random variables X_1, X_2, X_3 , and X_4 with joint PMF $p_{X_1 X_2 X_3 X_4}$. The marginal PMFs can be computed in the following way:

$$p_{X_2 X_3 X_4}(x_2, x_3, x_4) = \sum_{x_1 \in \mathcal{X}_1} p_{X_1 X_2 X_3 X_4}(x_1, x_2, x_3, x_4)$$

$$p_{X_1 X_4}(x_1, x_4) = \sum_{x_2 \in \mathcal{X}_2} \sum_{x_3 \in \mathcal{X}_3} p_{X_1 X_2 X_3 X_4}(x_1, x_2, x_3, x_4)$$

5.2 Joint Cumulative Distribution Function

5.2.1 Two Random Variables

Let X and Y be two random variables defined on the same probability space, (Ω, \mathcal{F}, P) . Suppose the range of X and Y are \mathcal{X} and \mathcal{Y} respectively. Their **joint cumulative distribution function (CDF)** describes the frequencies of their joint outcomes. The joint CDF of X and Y is defined as

$$\begin{aligned} F_{XY}(x, y) &= P(\{X \leq x, Y \leq y\}) \\ &= P(\{\omega \in \Omega : X(\omega) \leq x, Y(\omega) \leq y\}) \\ &= P(\{\omega \in \Omega : X(\omega) \leq x\} \cap \{\omega \in \Omega : Y(\omega) \leq y\}). \end{aligned}$$

The joint cumulative distribution function (CDF) satisfies several properties:

1. $0 \leq F_{XY}(x, y) \leq 1, \forall x \in \mathcal{X}, y \in \mathcal{Y}$
2. $F_{XY}(\infty, \infty) := \lim_{x \rightarrow \infty} \lim_{y \rightarrow \infty} F_{XY}(x, y) = 1$
3. $F_{XY}(x, -\infty) := \lim_{y \rightarrow -\infty} F_{XY}(x, y) = 0, \forall x \in \mathcal{X}$
4. $F_{XY}(-\infty, y) := \lim_{x \rightarrow -\infty} F_{XY}(x, y) = 0, \forall y \in \mathcal{Y}$
5. $F_X(x) = F_{XY}(x, \infty) := \lim_{y \rightarrow \infty} F_{XY}(x, y), \forall x \in \mathcal{X}$
6. $F_Y(y) = F_{XY}(\infty, y) := \lim_{x \rightarrow \infty} F_{XY}(x, y), \forall y \in \mathcal{Y}$
7. $F_{XY}(x_1, y_1) \leq F_{XY}(x_2, y_2), \forall x_1, x_2 \in \mathcal{X}, y_1, y_2 \in \mathcal{Y} : x_1 \leq x_2, y_1 \leq y_2$
8. $P(\{x_1 < X \leq x_2, y_1 < Y \leq y_2\}) =$
 $F_{XY}(x_2, y_2) - F_{XY}(x_2, y_1) - F_{XY}(x_1, y_2) + F_{XY}(x_1, y_1),$
 $\forall x_1, x_2 \in \mathcal{X}, y_1, y_2 \in \mathcal{Y} \text{ s.t. } x_1 \leq x_2 \text{ and } y_1 \leq y_2$

5.2.2 Multiple Random Variables

We can also expand the definition of the joint cumulative distribution function (CDF) to n random variables. Let X_1, \dots, X_n be n random variables defined on the same probability space, (Ω, \mathcal{F}, P) . The joint CDF is defined as

$$F_{X_1 \dots X_n}(x_1, \dots, x_n) = P(\{X_1 \leq x_1, \dots, X_n \leq x_n\}).$$

The joint CDF satisfies similar properties as in the case of two random variables.

5.3 Joint Probability Density Function

5.3.1 Two Random Variables

Let X and Y be two continuous random variables defined on the same probability space, (Ω, \mathcal{F}, P) . Suppose the range of X and Y are \mathcal{X} and \mathcal{Y} respectively. The **joint probability density function (PDF)** of the continuous random variables X and Y is defined in terms of their joint CDF as

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

The joint PDF of random variables X and Y satisfies the following properties:

1. $f_{XY}(x, y) \geq 0, \forall x \in \mathcal{X}, y \in \mathcal{Y}$
2. $\int_{\mathcal{X}} \int_{\mathcal{Y}} f_{XY}(x, y) dx dy = 1$
3. $F_{XY}(x, y) = \int_{-\infty}^x \int_{-\infty}^y f_{XY}(u, v) du dv$
4. $P(A) = \iint_A f_{XY}(x, y) dx dy, \forall A \in \mathcal{F}$

For random variables X and Y with joint PDF f_{XY} , the **marginal PDFs** are probability models for the individual random variables and are defined as

$$f_X(x) = \int_{-\infty}^{\infty} f_{XY}(x, y) dy \quad \text{and} \quad f_Y(y) = \int_{-\infty}^{\infty} f_{XY}(x, y) dx.$$

5.3.2 Multiple Random Variables

We can also expand the definition of the joint probability density function (PDF) to n random variables. Let X_1, \dots, X_n be n continuous random variables defined on the same probability space, (Ω, \mathcal{F}, P) . Suppose the range of X_i is \mathcal{X}_i for $i = 1, \dots, n$. For the random variables X_1, \dots, X_n , the joint PDF is

$$f_{X_1 \dots X_n}(x_1, \dots, x_n) = \frac{\partial^n}{\partial x_1 \dots \partial x_n} F_{X_1 \dots X_n}(x_1, \dots, x_n)$$

For n random variables, the joint PDF satisfies the following properties:

1. $f_{X_1 \dots X_n}(x_1, \dots, x_n) \geq 0, \forall x_1 \in \mathcal{X}_1, \dots, x_n \in \mathcal{X}_n$
2. $\int_{\mathcal{X}_1} \dots \int_{\mathcal{X}_n} f_{X_1 \dots X_n}(x_1, \dots, x_n) dx_1 \dots dx_n = 1$
3. $F_{X_1 \dots X_n}(x_1, \dots, x_n) = \int_{-\infty}^{x_1} \dots \int_{-\infty}^{x_n} f_{X_1 \dots X_n}(u_1, \dots, u_n) du_1 \dots du_n$
4. $P(A) = \int \dots \int_A f_{X_1 \dots X_n}(x_1, \dots, x_n) dx_1 \dots dx_n, \forall A \in \mathcal{F}$

Consider four random variables X_1, X_2, X_3 , and X_4 with joint PDF is $f_{X_1 X_2 X_3 X_4}$. The marginal PDFs can be computed in the following way:

$$f_{X_2 X_3 X_4}(x_2, x_3, x_4) = \int_{\mathcal{X}_1} f_{X_1 X_2 X_3 X_4}(x_1, x_2, x_3, x_4) dx_1$$

$$f_{X_1 X_4}(x_1, x_4) = \int_{\mathcal{X}_2} \int_{\mathcal{X}_3} f_{X_1 X_2 X_3 X_4}(x_1, x_2, x_3, x_4) dx_2 dx_3$$

5.4 Expected Value

5.4.1 Discrete Random Variables

Given a two discrete random variables X and Y with joint PMF p_{XY} , we can define a new random variable $W = g(X, Y)$. The expected value of W is

$$\mathbb{E}[W] = \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} g(x, y) p_{XY}(x, y).$$

Using this definition of the expected value, we can show that the expectation of a sum of two discrete random variables is the sum of its expected values:

$$\mathbb{E}[X + Y] = \mathbb{E}[X] + \mathbb{E}[Y].$$

Proof: We can directly prove this property for two discrete random variables:

$$\begin{aligned} \mathbb{E}[X + Y] &= \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} (x + y) p_{XY}(x, y) \\ &= \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} x p_{XY}(x, y) + \sum_{y \in \mathcal{Y}} \sum_{x \in \mathcal{X}} y p_{XY}(x, y) \\ &= \sum_{x \in \mathcal{X}} x \left(\sum_{y \in \mathcal{Y}} p_{XY}(x, y) \right) + \sum_{y \in \mathcal{Y}} y \left(\sum_{x \in \mathcal{X}} p_{XY}(x, y) \right) \\ &= \sum_{x \in \mathcal{X}} x p_X(x) + \sum_{y \in \mathcal{Y}} y p_Y(y) \\ &= \mathbb{E}[X] + \mathbb{E}[Y] \end{aligned}$$

5.4.2 Continuous Random Variables

Given a two continuous random variables X and Y with joint PDF f_{XY} , we can define a new random variable $W = g(X, Y)$. The expected value of W is

$$\mathbb{E}[W] = \int_{\mathcal{X}} \int_{\mathcal{Y}} g(x, y) f_{XY}(x, y) dx dy.$$

As for the discrete case, the expectation of a sum of two continuous random variables is the sum of its expected values:

$$\mathbb{E}[X + Y] = \mathbb{E}[X] + \mathbb{E}[Y].$$

Proof: We can directly prove this for two continuous random variables:

$$\begin{aligned} \mathbb{E}[X + Y] &= \int_{\mathcal{X}} \int_{\mathcal{Y}} (x + y) f_{XY}(x, y) dx dy \\ &= \int_{\mathcal{X}} \int_{\mathcal{Y}} x f_{XY}(x, y) dx dy + \int_{\mathcal{X}} \int_{\mathcal{Y}} y f_{XY}(x, y) dx dy \\ &= \int_{\mathcal{X}} x \left(\int_{\mathcal{Y}} f_{XY}(x, y) dy \right) dx + \int_{\mathcal{Y}} y \left(\int_{\mathcal{X}} f_{XY}(x, y) dx \right) dy \\ &= \int_{\mathcal{X}} x f_X(x) dx + \int_{\mathcal{Y}} y f_Y(y) dy \\ &= \mathbb{E}[X] + \mathbb{E}[Y] \end{aligned}$$

5.5 Variance, Covariance, & Correlation

Previously, we considered the variance as a measure of the dispersion of sample values for a single random variable. We can also consider the variance of multiple random variables. This leads to two related metrics: covariance and correlation.

5.5.1 Variance

For any two random variables X and Y , the variance of their sum can be expressed in terms of the variances of the individual variables as

$$\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2\mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])].$$

Proof: We can use the definition of the variance to show this property:

$$\begin{aligned}
 \text{Var}(X + Y) &= \mathbb{E} \left[(X + Y - \mathbb{E}[X + Y])^2 \right] \\
 &= \mathbb{E} \left[(X + Y - \mathbb{E}[X] - \mathbb{E}[Y])^2 \right] \\
 &= \mathbb{E} \left[X^2 + Y^2 + (\mathbb{E}[X])^2 + (\mathbb{E}[Y])^2 + 2XY - 2X\mathbb{E}[X] - 2Y\mathbb{E}[X] - 2X\mathbb{E}[Y] - 2Y\mathbb{E}[Y] + 2\mathbb{E}[X]\mathbb{E}[Y] \right] \\
 &= \mathbb{E} \left[(X - \mathbb{E}[X])^2 + (Y - \mathbb{E}[Y])^2 + 2(X - \mathbb{E}[X])(Y - \mathbb{E}[Y]) \right] \\
 &= \mathbb{E} \left[(X - \mathbb{E}[X])^2 \right] + \mathbb{E} \left[(Y - \mathbb{E}[Y])^2 \right] + 2\mathbb{E} \left[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y]) \right] \\
 &= \text{Var}(X) + \text{Var}(Y) + 2\mathbb{E} \left[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y]) \right]
 \end{aligned}$$

5.5.2 Covariance

The last term in the equation above is called the **covariance** and is denoted:

$$\sigma_{XY} := \text{Cov}(X, Y) = \mathbb{E} \left[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y]) \right].$$

The covariance describes how the pair of random variables, X and Y , vary together. If the covariance is positive, then as X increases, Y generally increases as well. If the covariance is negative, then as X increases, Y generally decreases. If the covariance is zero, then X and Y are said to be **uncorrelated**. The covariance satisfies the following useful linearity property:

$$\begin{aligned}
 \text{Cov} \left(\sum_{j=1}^n A_{ij} X_j + b_i, \sum_{l=1}^n A_{kl} X_l + b_k \right) \\
 &= \text{Cov} \left(\sum_{j=1}^n A_{ij} X_j, \sum_{l=1}^n A_{kl} X_l \right) \\
 &= \sum_{j=1}^n \sum_{l=1}^n A_{ij} A_{kl} \text{Cov}(X_j, X_l) \\
 &= \sum_{j=1}^n A_{ij} A_{kj} \text{Var}(X_j) + \sum_{1 \leq j < l \leq n} A_{ij} A_{kl} \text{Cov}(X_j, X_l)
 \end{aligned}$$

5.5.3 Correlation Coefficient

The **correlation coefficient** is a normalized version of the covariance, which indicates the relationship between two random variables, regardless of the measurement units. While the unit of the covariance is the product of the units of X and Y , the correlation coefficient is dimensionless and is unaffected by scale changes. The correlation coefficient for two random variables X and Y is

$$\rho_{XY} := \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}} = \frac{\sigma_{XY}}{\sigma_X \sigma_Y}.$$

The correlation coefficient, ρ_{XY} , is restricted to values between -1 and 1 . If ρ_{XY} is positive, there is a positive correlation between X and Y . Similarly, if it is negative, then there is a negative correlation. If the correlation coefficient is zero, then there is no correlation. If the magnitude of ρ_{XY} is close to zero, then the random variables are only weakly correlated. If the magnitude of ρ_{XY} is close to one, then the random variables are highly correlated. If it is equal to one, then there is a linear relationship between the two variables.

5.5.4 Correlation

Another metric used to describe the relationship between two random variables is **correlation**. The correlation between random variables X and Y is

$$r_{XY} := \mathbb{E}[XY] = \text{Cov}(X, Y) + \mathbb{E}[X]\mathbb{E}[Y].$$

If the correlation between two random variables, X and Y , is zero, then X and Y are said to be **orthogonal**.

5.6 Independent Random Variables

5.6.1 Discrete Random Variables

Two discrete random variables X and Y , with corresponding ranges \mathcal{X} and \mathcal{Y} , are considered to be **independent** if and only if

$$p_{XY}(x, y) = p_X(x)p_Y(y), \quad \forall x \in \mathcal{X}, y \in \mathcal{Y}.$$

Similarly, n discrete random variables X_1, \dots, X_n , with corresponding ranges $\mathcal{X}_1, \dots, \mathcal{X}_n$, are considered to be independent if and only if

$$p_{X_1 \dots X_n}(x_1, \dots, x_n) = \prod_{i=1}^n p_{X_i}(x_i), \quad \forall x_i \in \mathcal{X}_i, i = 1, \dots, n.$$

These variables are said to be **independent and identically distributed (i.i.d.)** if and only if (1) the random variables are independent and (2) all of the random variables have the same PMF, which we denote p_X , meaning

$$p_{X_1 \dots X_n}(x_1, \dots, x_n) = \prod_{i=1}^n p_X(x_i).$$

5.6.2 Continuous Random Variables

Two continuous random variables X and Y , with corresponding ranges \mathcal{X} and \mathcal{Y} , are considered to be **independent** if and only if

$$f_{XY}(x, y) = f_X(x)f_Y(y), \quad \forall x \in \mathcal{X}, y \in \mathcal{Y}.$$

Similarly, n continuous random variables X_1, \dots, X_n , with corresponding ranges $\mathcal{X}_1, \dots, \mathcal{X}_n$, are considered to be independent if and only if

$$f_{X_1 \dots X_n}(x_1, \dots, x_n) = \prod_{i=1}^n f_{X_i}(x_i), \quad \forall x_i \in \mathcal{X}_i, i = 1, \dots, n.$$

These variables are said to be **independent and identically distributed (i.i.d.)** if and only if (1) the random variables are independent and (2) all of the random variables have the same PDF, which we denote f_X , meaning

$$f_{X_1 \dots X_n}(x_1, \dots, x_n) = \prod_{i=1}^n f_X(x_i)$$

5.6.3 Independence Properties

For two independent variables X and Y , the following properties hold:

1. $\mathbb{E}[g(X)h(Y)] = \mathbb{E}[g(X)]\mathbb{E}[h(Y)]$
2. $r_{XY} = \mathbb{E}[XY] = \mathbb{E}[X]\mathbb{E}[Y]$
3. $\text{Cov}(X, Y) = 0$
4. $\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y)$

Notice that all independent random variables are uncorrelated, but they are not necessarily orthogonal. Also note that while independence implies that two random variables are uncorrelated, two uncorrelated random variables are not necessarily independent.

5.7 Bivariate Gaussian Random Variables

One special class of a pair of random variables is called bivariate Gaussian random variables. Let X be a Gaussian random variable with mean μ_X and standard deviation σ_X , and let Y be a Gaussian random variable with mean μ_Y and standard deviation σ_Y . Assume that the correlation coefficient between X and Y is given by ρ_{XY} . If X and Y form a pair of bivariate Gaussian random variables, then their joint PDF can be expressed as the following:

$$f_{XY}(x, y) = \frac{1}{2\pi\sigma_X\sigma_Y\sqrt{1-\rho_{XY}^2}} \exp\left(-\frac{\left(\frac{x-\mu_X}{\sigma_X}\right)^2 + \left(\frac{y-\mu_Y}{\sigma_Y}\right)^2 - \frac{2\rho_{XY}(x-\mu_X)(y-\mu_Y)}{\sigma_X\sigma_Y}}{2(1-\rho_{XY}^2)}\right).$$

Previously, we said two independent variables are uncorrelated but two uncorrelated random variables are not necessarily independent. For bivariate Gaussian random variables, they are uncorrelated if and only if they are independent.

CHAPTER 5. MULTIPLE RANDOM VARIABLES

Another important property of bivariate Gaussian random variables is that a pair of linear combinations of bivariate Gaussian random variables forms another pair of bivariate Gaussian random variables. If X and Y form a pair of bivariate Gaussian random variables with parameters $(\mu_X, \mu_Y, \sigma_X, \sigma_Y, \rho_{XY})$, and the random variables W_1 and W_2 are given by the linearly independent equations $W_1 = a_1X + b_1Y$ and $W_2 = a_2X + b_2Y$, then W_1 and W_2 form a pair of bivariate Gaussian random variables such that

$$(1) \mathbb{E}[W_i] = a_i\mu_X + b_i\mu_Y,$$

$$(2) \text{Var}(W_i) = a_i^2\sigma_X^2 + b_i^2\sigma_Y^2 + 2a_ib_i\rho_{XY}\sigma_X\sigma_Y, \text{ and}$$

$$(3) \text{Cov}(W_1, W_2) = a_1a_2\sigma_X^2 + b_1b_2\sigma_Y^2 + (a_1b_2 + a_2b_1)\rho_{XY}\sigma_X\sigma_Y$$

Chapter 6

Derived Random Variables

6.1 Derived Discrete Random Variable

6.1.1 Function of One Random Variable

Given a discrete random variable X with PMF p_X and some function g , the derived random variable $Y = g(X)$ has the PMF given by

$$p_Y(y) = \sum_{x:g(x)=y} p_X(x).$$

If g is an invertible function, then we can express the PMF as

$$p_Y(y) = p_X(g^{-1}(y)).$$

6.1.2 Function of Two Random Variables

We can extend the definition given previously for functions of two random variables. Given discrete random variables X and Y with PMFs p_X and p_Y respectively, the derived random variable $W = g(X, Y)$ has the PMF given by

$$p_W(w) = \sum_{(x,y):g(x,y)=w} p_{XY}(x, y).$$

6.1.3 Sums of Independent Random Variables

If X and Y are discrete random variables, then their sum is another discrete random variable, $W = X + Y$. If X and Y are independent with PMFs p_X and p_Y respectively, then the PMF of W is the convolution of p_X and p_Y :

$$p_W(w) = (p_X * p_Y)(w) = \sum_{k=-\infty}^{\infty} p_X(k)p_Y(w - k).$$

6.2 Derived Continuous Random Variable

6.2.1 Function of One Random Variable

Given a continuous random variable X with CDF F_X and some function g , the derived random variable $Y = g(X)$ has the CDF given by

$$F_Y(y) = P(\{Y \leq y\}) = P(\{g(X) \leq y\}).$$

If g is an invertible function, then we can express the CDF as

$$F_Y(y) = P(\{X \leq g^{-1}(y)\}) = F_X(g^{-1}(y)).$$

Given the CDF of the derived random variable Y , we can find its probability density function (PDF) by computing the derivative of the CDF:

$$f_Y(y) = \frac{d}{dy}F_Y(y).$$

If g is a linear function with slope a and intercept b , then $Y = aX + b$. For this special case, we can use the following identities to derive the CDF and PDF Y :

$$F_Y(y) = F_X\left(\frac{y-b}{a}\right) \quad \text{and} \quad f_Y(y) = \frac{1}{a}f_X\left(\frac{y-b}{a}\right).$$

6.2.2 Function of Two Random Variables

We can extend the definition given previously for functions of two random variables. Given continuous random variables X and Y with CDFs F_X and F_Y respectively, the derived random variable $W = g(X, Y)$ has the CDF given by

$$F_W(w) = \iint_{(x,y):g(x,y) \leq w} f_{XY}(x,y) dx dy.$$

Again, the probability density function (PDF) of W is then given by

$$f_W(w) = \frac{d}{dw}F_W(w).$$

6.2.3 Sums of Independent Random Variables

If X and Y are continuous random variables, then their sum is another continuous random variable, $W = X + Y$. If X and Y are independent with PDFs f_X and f_Y respectively, then the PDF of W is the convolution of f_X and f_Y :

$$f_W(w) = (f_X * f_Y)(w) = \int_{-\infty}^{\infty} f_X(\tau)f_Y(w - \tau)d\tau$$

Chapter 7

Conditional Probability Models

7.1 Conditioning by an Event

7.1.1 Conditional Probability Mass Function

For a discrete random variable X with PMF p_X and an event A with probability $P(A) > 0$, the **conditional probability mass function (PMF)** is defined as

$$p_{X|A}(x) = P(\{X = x|A\}) = \begin{cases} \frac{p_X(x)}{P(A)} & \text{if } x \in A \\ 0 & \text{otherwise} \end{cases}.$$

The conditional PMF satisfies the following properties:

1. $p_{X|A}(x) \geq 0, \forall x \in A$
2. $\sum_{x \in A} p_{X|A}(x) = 1$
3. $P(B|A) = \sum_{x \in B} p_{X|A}(x)$

For a discrete random variable X resulting from an experiment with partition A_1, \dots, A_m , we can express the PMF of X in terms of its conditional PMF as

$$p_X(x) = \sum_{i=1}^m p_{X|A_i}(x)P(A_i).$$

For two discrete random variables X and Y with the joint PMF p_{XY} and an event A with probability $P(A) > 0$, the **conditional joint PMF** is given by

$$p_{XY|A}(x, y) = P(\{X = x, Y = y|A\}) = \begin{cases} \frac{p_{XY}(x, y)}{P(A)} & \text{if } (x, y) \in A \\ 0 & \text{otherwise} \end{cases}.$$

7.1.2 Conditional Cumulative Distribution Function

For a random variable X with CDF F_X and event A with probability $P(A) > 0$, the **conditional cumulative distribution function (CDF)** is defined as

$$F_{X|A}(x) = P(\{X \leq x|A\}) = \begin{cases} \frac{F_X(x)}{P(A)} & \text{if } x \in A \\ 0 & \text{otherwise} \end{cases}.$$

For two random variables X and Y with joint CDF F_{XY} and an event A with probability $P(A) > 0$, the **conditional joint CDF** is given by

$$F_{XY|A}(x, y) = P(\{X \leq x, Y \leq y|A\}) = \begin{cases} \frac{F_{XY}(x, y)}{P(A)} & \text{if } x \in A \\ 0 & \text{otherwise} \end{cases}.$$

7.1.3 Conditional Probability Density Function

For a continuous random variable X with PDF f_X and conditional CDF $F_{X|A}$ and an event A with probability $P(A) > 0$, the **conditional probability density function (PDF)** can be derived from the conditional CDF as

$$f_{X|A}(x) = \frac{d}{dx} F_{X|A}(x) = \begin{cases} \frac{f_X(x)}{P(A)} & \text{if } x \in A \\ 0 & \text{otherwise} \end{cases}.$$

The conditional PDF satisfies the following properties:

1. $f_{X|A}(x) \geq 0, \forall x \in A$
2. $\int_A f_{X|A}(x) dx = 1$
3. $P(B|A) = \int_B p_{X|A}(x) dx$

For a continuous random variable X resulting from an experiment with partition A_1, \dots, A_m , we can express the PDF of X in terms of its conditional PDF as

$$f_X(x) = \sum_{i=1}^m f_{X|A_i}(x) P(A_i).$$

For two continuous random variables X and Y with joint PDF f_{XY} and conditional joint CDF $F_{XY|A}$ and an event A with probability $P(A) > 0$, the **conditional joint PDF** is given by

$$f_{XY|A}(x, y) = \frac{\partial^2}{\partial x \partial y} F_{XY|A}(x, y) = \begin{cases} \frac{f_{XY}(x, y)}{P(A)} & \text{if } (x, y) \in A \\ 0 & \text{otherwise} \end{cases}.$$

7.1.4 Conditional Expected Value

For a *discrete* random variable X and an event A with probability $P(A) > 0$, the **conditional expected value** of X is defined as

$$\mathbb{E}[X|A] = \sum_{x \in \mathcal{X}} xp_{X|A}(x).$$

If X is a *discrete* random variable, Y is a random variable defined as $Y = g(X)$, and A is an event, then the conditional expected value of Y is given by

$$\mathbb{E}[Y|A] = \mathbb{E}[g(X)|A] = \sum_{x \in \mathcal{X}} g(x)p_{X|A}(x).$$

If X and Y are two *discrete* random variables, W is defined as $W = g(X, Y)$, and A is an event, then the conditional expected value of W is given by

$$\mathbb{E}[W|A] = \mathbb{E}[g(X, Y)|A] = \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} g(x, y)p_{XY|A}(x, y).$$

For a *discrete* random variable X resulting from an experiment with partition A_1, \dots, A_m , the expected value of X can be expressed as

$$\mathbb{E}[X] = \sum_{i=1}^m \mathbb{E}[X|A_i]P(A_i).$$

For a *continuous* random variable X and an event A with probability $P(A) > 0$, the **conditional expected value** of X is defined as

$$\mathbb{E}[X|A] = \int_{\mathcal{X}} xf_{X|A}(x)dx.$$

If X is a *continuous* random variable, Y is defined as $Y = g(X)$, and A is an event, then the conditional expected value of Y is given by

$$\mathbb{E}[Y|A] = \mathbb{E}[g(X)|A] = \int_{\mathcal{X}} g(x)f_{X|A}(x)dx.$$

If X and Y are two *continuous* random variables, W is defined as $W = g(X, Y)$, and A is an event, then the conditional expected value of W is given by

$$\mathbb{E}[W|A] = \mathbb{E}[g(X, Y)|A] = \int_{\mathcal{Y}} \int_{\mathcal{X}} g(x, y)f_{XY|A}(x, y)dxdy.$$

For a *continuous* random variable X resulting from an experiment with partition A_1, \dots, A_m , the expected value of X can be expressed as

$$\mathbb{E}[X] = \sum_{i=1}^m \mathbb{E}[X|A_i]P(A_i).$$

7.1.5 Conditional Variance

For a random variable X (discrete or continuous) and an event A with probability $P(A) > 0$, the **conditional variance** of X is defined as

$$\text{Var}(X|A) = \mathbb{E} \left[\left(X - \mathbb{E}[X|A] \right)^2 \middle| A \right] = \mathbb{E}[X^2|A] - \left(\mathbb{E}[X|A] \right)^2.$$

The **conditional standard deviation** is defined in terms of the variance as

$$\sigma_{X|A} = \sqrt{\text{Var}(X|A)}.$$

7.2 Conditioning by a Random Variable

7.2.1 Conditional Probability Mass Function

For two discrete random variables X and Y , we can define two **conditional probability mass functions (PMFs)**, which are given by

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

and

$$p_{Y|X}(y|x) = P(\{Y = y|X = x\}) = \frac{p_{XY}(x, y)}{p_X(x)}.$$

The **conditional PMF theorem** says that the joint PMF can be expressed as

$$p_{XY}(x, y) = p_{X|Y}(x|y)p_Y(y) = p_{Y|X}(y|x)p_X(x).$$

If X and Y are independent random variables, then

$$p_{X|Y}(x|y) = p_X(x) \quad \text{and} \quad p_{Y|X}(y|x) = p_Y(y)$$

7.2.2 Conditional Cumulative Distribution Function

For two random variables X and Y , we can define two **conditional cumulative distribution functions (CDFs)**, which are given by

$$F_{X|Y}(x|y) = P(\{X \leq x|Y \leq y\})$$

and

$$F_{Y|X}(y|x) = P(\{Y \leq y|X \leq x\}).$$

7.2.3 Conditional Probability Density Function

For two continuous random variables X and Y , two **conditional probability density functions (PDFs)** are defined in terms of the conditional CDFs as

$$f_{X|Y}(x|y) = \frac{d}{dx}F_{X|Y}(x|y) = \frac{f_{XY}(x, y)}{f_Y(y)}$$

and

$$f_{Y|X}(y|x) = \frac{d}{dy}F_{Y|X}(y|x) = \frac{f_{XY}(x, y)}{f_X(x)}.$$

The **conditional PDF theorem** says that the joint PDF can be expressed as

$$f_{XY}(x, y) = f_{X|Y}(x|y)f_Y(y) = f_{Y|X}(y|x)f_X(x).$$

If X and Y are independent random variables, then

$$f_{X|Y}(x|y) = f_X(x) \quad \text{and} \quad f_{Y|X}(y|x) = f_Y(y).$$

7.2.4 Conditional Expected Value

If X and Y are two *discrete* random variables, then the **conditional expected value** of X given that $Y = y$ is defined as

$$\mathbb{E}[X|Y = y] = \sum_{x \in \mathcal{X}} xp_{X|Y}(x|y).$$

If X and Y are two *discrete* random variables and W is defined as $W = g(X, Y)$, then the conditional expected value of W given $Y = y$ is defined as

$$\mathbb{E}[W|Y = y] = \mathbb{E}[g(X, Y)|Y = y] = \sum_{x \in \mathcal{X}} g(x, y)p_{X|Y}(x|y).$$

If X and Y are two *continuous* random variables, then the **conditional expected value** of X given that $Y = y$ is defined as

$$\mathbb{E}[X|Y = y] = \int_{\mathcal{X}} xf_{X|Y}(x|y)dx.$$

If X and Y are two *continuous* random variables and W is defined as $W = g(X, Y)$, then the conditional expected value of W given $Y = y$ is defined as

$$\mathbb{E}[W|Y = y] = \mathbb{E}[g(X, Y)|Y = y] = \int_{\mathcal{X}} g(x, y)f_{X|Y}(x|y)dx.$$

If the random variables X and Y are independent, then

$$\mathbb{E}[X|Y = y] = E[X], \quad \forall y \in \mathcal{Y} \quad \text{and} \quad \mathbb{E}[Y|X = x] = E[Y], \quad \forall x \in \mathcal{X}.$$

7.2.5 Iterated Expectation & Tower Property

The conditional expected value $\mathbb{E}[X|Y]$ is a function of the random variable Y such that if $Y = y$, then $\mathbb{E}[X|Y] = \mathbb{E}[X|Y = y]$. Because $\mathbb{E}[X|Y]$ is a function of a random variable, it is also a random variable. To calculate the expected value of $g(X)$, we can compute the **iterated expected value**:

$$\mathbb{E}[g(X)] = \mathbb{E}[\mathbb{E}[g(X)|Y]].$$

If X is a discrete random variable, then this expected value is given by

$$\mathbb{E}[g(X)] = \sum_{y \in \mathcal{Y}} \mathbb{E}[g(X)|Y = y] p_Y(y).$$

If X is a continuous random variable, then this expected value is given by

$$\mathbb{E}[g(X)] = \int_{\mathcal{Y}} \mathbb{E}[g(X)|Y = y] f_Y(y) dy.$$

Similar to the iterated expected value, we also have the **tower property**, which says that for any two random variables X and Y and any function f ,

$$\mathbb{E}[f(Y)X] = \mathbb{E}[f(Y)\mathbb{E}[X|Y]].$$

7.2.6 Conditional Variance

If X and Y are two random variables, the variance of X given that $Y = y$ is

$$\begin{aligned} \text{Var}(X|Y = y) &= \mathbb{E} \left[\left(X - \mathbb{E}[X|Y = y] \right)^2 \middle| Y = y \right] \\ &= \mathbb{E}[X^2|Y = y] - (\mathbb{E}[X|Y = y])^2. \end{aligned}$$

The conditional variance $\text{Var}(X|Y)$ is a function of the random variable Y such that if $Y = y$, then $\text{Var}(X|Y) = \text{Var}(X|Y = y)$. Because $\text{Var}(X|Y)$ is a function of a random variable, it is also a random variable. Using iterated expectation, we can write the **minimum mean-square error (MMSE)** of X given Y :

$$\mathbb{E}[\text{Var}(X|Y)] = \mathbb{E} \left[\left(X - \mathbb{E}[X|Y] \right)^2 \right].$$

The **law of total variance** says that the variance of X can be expressed as

$$\text{Var}(X) = \mathbb{E}[\text{Var}(X|Y)] + \text{Var}(\mathbb{E}[X|Y]).$$

Rearranging this equation, we see that we can express the **minimum mean-square error (MMSE)** of X given Y as

$$\mathbb{E}[\text{Var}(X|Y)] = \text{Var}(X) - \text{Var}(\mathbb{E}[X|Y]).$$

Part IV

Extensions of Random Variables

Chapter 8

Random Vectors

8.1 Random Vector

A **random vector** with n dimensions is a concise representation of a set of n random variables. Often we express the n -dimensional random vector \mathbf{X} as

$$\mathbf{X} = \begin{bmatrix} X_1 \\ \vdots \\ X_n \end{bmatrix},$$

where each element X_i is a random variable. A realization or sample value of the random vector \mathbf{X} is another vector \mathbf{x} , which can be expressed as

$$\mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix},$$

where the i th element x_i is a sample value of the random variable X_i .

8.2 Probability Distributions

8.2.1 Probability Mass Functions

The probability mass function (PMF) for a discrete random vector \mathbf{X} is

$$p_{\mathbf{X}}(\mathbf{x}) = p_{X_1 \dots X_n}(x_1, \dots, x_n).$$

If \mathbf{X} is a discrete random vector and $\mathbf{W} = g(\mathbf{X})$, the PMF of \mathbf{W} is given by

$$p_{\mathbf{W}}(\mathbf{w}) = \sum_{\mathbf{x}:g(\mathbf{x})=\mathbf{w}} p_{\mathbf{X}}(\mathbf{x}).$$

The joint PMF for a pair of discrete random vectors, \mathbf{X} and \mathbf{Y} , is defined as

$$p_{\mathbf{X}, \mathbf{Y}}(\mathbf{x}, \mathbf{y}) = p_{X_1 \dots X_n Y_1 \dots Y_m}(x_1, \dots, x_n, y_1, \dots, y_m).$$

The random vectors \mathbf{X} and \mathbf{Y} are said to be independent if and only if

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

8.2.2 Cumulative Distribution Functions

The cumulative distribution function (CDF) for a random vector \mathbf{X} is

$$F_{\mathbf{X}}(\mathbf{x}) = F_{X_1 \dots X_n}(x_1, \dots, x_n).$$

If \mathbf{X} is a discrete random vector and $\mathbf{W} = g(\mathbf{X})$, the CDF of \mathbf{W} is given by

$$F_{\mathbf{W}}(\mathbf{w}) = \sum_{\mathbf{x}: g(\mathbf{x}) \leq \mathbf{w}} p_{\mathbf{X}}(\mathbf{x}).$$

If \mathbf{X} is a continuous random vector and $\mathbf{W} = g(\mathbf{X})$, the CDF of \mathbf{W} is given by

$$F_{\mathbf{W}}(\mathbf{w}) = \int_{g(\mathbf{x}) \leq \mathbf{w}} f_{\mathbf{X}}(\mathbf{x}) d\mathbf{x}.$$

The joint CDF for a pair of random vectors, \mathbf{X} and \mathbf{Y} , is defined as

$$F_{\mathbf{X}, \mathbf{Y}}(\mathbf{x}, \mathbf{y}) = F_{X_1 \dots X_n Y_1 \dots Y_m}(x_1, \dots, x_n, y_1, \dots, y_m).$$

8.2.3 Probability Density Functions

The probability density function (PDF) for a continuous random vector \mathbf{X} is

$$f_{\mathbf{X}}(\mathbf{x}) = f_{X_1 \dots X_n}(x_1, \dots, x_n).$$

If \mathbf{X} is a continuous random vector and the random variable $\mathbf{W} = g(\mathbf{X})$ has the CDF $F_{\mathbf{W}}$, the PDF of \mathbf{W} can be determined from the CDF of \mathbf{W} :

$$f_{\mathbf{W}}(\mathbf{w}) = \nabla_{\mathbf{w}} F_{\mathbf{W}}(\mathbf{w}).$$

The joint PDF for a pair of continuous random vectors, \mathbf{X} and \mathbf{Y} , is defined as

$$f_{\mathbf{X}, \mathbf{Y}}(\mathbf{x}, \mathbf{y}) = f_{X_1 \dots X_n Y_1 \dots Y_m}(x_1, \dots, x_n, y_1, \dots, y_m).$$

The random vectors \mathbf{X} and \mathbf{Y} are said to be independent if and only if

$$f_{\mathbf{X}, \mathbf{Y}}(\mathbf{x}, \mathbf{y}) = f_{\mathbf{X}}(\mathbf{x})f_{\mathbf{Y}}(\mathbf{y}).$$

8.3 Properties of Random Vectors

8.3.1 Expected Value Vector

The expected value of a random vector is a column vector whose elements are the expected values of the individual random variables, which we express as

$$\boldsymbol{\mu}_{\mathbf{X}} = \mathbb{E}[\mathbf{X}] = \begin{bmatrix} \mathbb{E}[X_1] \\ \vdots \\ \mathbb{E}[X_n] \end{bmatrix}.$$

If \mathbf{X} is a discrete random vector, its expected value is determined by its PMF:

$$\mathbb{E}[\mathbf{X}] = \sum_{\mathbf{x} \in \mathcal{X}} \mathbf{x} p_{\mathbf{X}}(\mathbf{x}).$$

If \mathbf{X} is a continuous random vector, its expected value is defined by its PDF:

$$\mathbb{E}[\mathbf{X}] = \int_{\mathcal{X}} \mathbf{x} f_{\mathbf{X}}(\mathbf{x}) d\mathbf{x}.$$

If \mathbf{X} is a discrete random vector and $\mathbf{Y} = g(\mathbf{X})$, then \mathbf{Y} has the expected value

$$\mathbb{E}[\mathbf{Y}] = \mathbb{E}[g(\mathbf{X})] = \sum_{\mathbf{x} \in \mathcal{X}} g(\mathbf{x}) p_{\mathbf{X}}(\mathbf{x}).$$

If \mathbf{X} is a continuous random vector and $\mathbf{Y} = g(\mathbf{X})$, \mathbf{Y} has the expected value

$$\mathbb{E}[\mathbf{Y}] = \mathbb{E}[g(\mathbf{X})] = \int_{\mathcal{X}} g(\mathbf{x}) f_{\mathbf{X}}(\mathbf{x}) d\mathbf{x}.$$

When the components of \mathbf{X} are all independent, the following property holds:

$$\mathbb{E}[g_1(X_1) \dots g_n(X_n)] = \prod_{i=1}^n \mathbb{E}[g_i(X_i)].$$

8.3.2 Covariance & Correlation Matrices

The covariance of an n -dimensional random vector \mathbf{X} is an $n \times n$ matrix $\boldsymbol{\Sigma}_{\mathbf{X}}$ whose (i, j) th element is $\text{Cov}(X_i, X_j)$. In vector notation, we write this as

$$\boldsymbol{\Sigma}_{\mathbf{X}} = \mathbb{E} \left[(\mathbf{X} - \boldsymbol{\mu}_{\mathbf{X}})(\mathbf{X} - \boldsymbol{\mu}_{\mathbf{X}})^T \right].$$

From this definition of the covariance, we can see that the **covariance matrix** $\boldsymbol{\Sigma}_{\mathbf{X}}$ is always positive semidefinite (PSD), meaning $\mathbf{v}^T \boldsymbol{\Sigma}_{\mathbf{X}} \mathbf{v} \geq 0$ for all $\mathbf{v} \in \mathbb{R}^n$.

The correlation of an n -dimensional random vector \mathbf{X} is an $n \times n$ matrix $\mathbf{R}_{\mathbf{X}}$ whose (i, j) th element is $\mathbb{E}[X_i X_j]$. In vector notation, we write this as

$$\mathbf{R}_{\mathbf{X}} = \mathbb{E}[\mathbf{X} \mathbf{X}^T].$$

For a random vector \mathbf{X} with **correlation matrix** $\mathbf{R}_{\mathbf{X}}$, **covariance matrix** $\boldsymbol{\Sigma}_{\mathbf{X}}$, and **expected value vector** $\boldsymbol{\mu}_{\mathbf{X}}$, we can write the following relationship:

$$\mathbf{R}_{\mathbf{X}} = \boldsymbol{\Sigma}_{\mathbf{X}} + \boldsymbol{\mu}_{\mathbf{X}} \boldsymbol{\mu}_{\mathbf{X}}^T.$$

8.3.3 Linear Transformation

Consider the n -dimensional random vector \mathbf{X} with correlation matrix \mathbf{R}_X , covariance matrix Σ_X , and expected value vector $\boldsymbol{\mu}_X$. If \mathbf{A} is an $m \times n$ matrix and \mathbf{b} is an m -dimensional vector, then the random vector $\mathbf{Y} = \mathbf{A}\mathbf{x} + \mathbf{b}$ has the following mean, correlation matrix, and covariance matrix:

$$\begin{aligned}\boldsymbol{\mu}_Y &= \mathbf{A}\boldsymbol{\mu}_X + \mathbf{b} \\ \mathbf{R}_Y &= \mathbf{A}\mathbf{R}_X\mathbf{A}^T + (\mathbf{A}\boldsymbol{\mu}_X)\mathbf{b}^T + \mathbf{b}(\mathbf{A}\boldsymbol{\mu}_X)^T + \mathbf{b}\mathbf{b}^T \\ \Sigma_Y &= \mathbf{A}\Sigma_X\mathbf{A}^T\end{aligned}$$

8.4 Gaussian Random Vectors

The random vector \mathbf{X} is a Gaussian random vector with expected value vector $\boldsymbol{\mu}_X$ and covariance matrix Σ_X if and only if its PDF can be expressed as

$$f_{\mathbf{X}}(\mathbf{x}) = \frac{1}{(2\pi)^{n/2}|\Sigma_X|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_X)^T \Sigma_X^{-1}(\mathbf{x} - \boldsymbol{\mu}_X)\right).$$

The Gaussian random vector \mathbf{X} has independent components if and only if all of its components are uncorrelated, meaning $\text{Cov}(X_i, X_j) = 0$ for all $i \neq j \in [1, n]$. Therefore, a Gaussian random vector \mathbf{X} has independent components if and only if Σ_X is a diagonal matrix, whose diagonal elements are $\sigma_i^2 = \text{Var}(X_i)$, i.e.

$$\Sigma_X = \text{diag}(\sigma_1^2, \dots, \sigma_n^2) = \begin{bmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_n^2 \end{bmatrix}.$$

Consider an n -dimensional Gaussian random variable \mathbf{X} with expected value $\boldsymbol{\mu}_X$ and covariance Σ_X . If \mathbf{A} is an $m \times n$ matrix and \mathbf{b} is an m -dimensional vector, then the random vector $\mathbf{Y} = \mathbf{A}\mathbf{X} + \mathbf{b}$ is an m -dimensional Gaussian random vector with mean vector and covariance matrix given by

$$\boldsymbol{\mu}_Y = \mathbf{A}\boldsymbol{\mu}_X + \mathbf{b} \quad \text{and} \quad \Sigma_Y = \mathbf{A}\Sigma_X\mathbf{A}^T.$$

8.4.1 Standard Normal Random Vector

The n -dimensional standard normal random vector \mathbf{Z} is the n -dimensional Gaussian random vector with expected value vector $\boldsymbol{\mu}_Z = \mathbf{0}_n$ and covariance matrix $\Sigma_Z = \mathbf{I}_n$. This means each component of \mathbf{Z} is a random variable Z_i with an expected value of zero and variance of one. Furthermore, the covariance between all of the random variables is zero, which implies that all of the random variables are uncorrelated, and thus independent.

We can easily transform between general Gaussian random vectors and standard normal random vectors. Assume \mathbf{X} is an n -dimensional Gaussian random vector

with expected value $\boldsymbol{\mu}_{\mathbf{X}}$ and covariance $\boldsymbol{\Sigma}_{\mathbf{X}}$. Because the covariance matrix is positive semidefinite, we can define an $n \times n$ matrix \mathbf{A} such that $\mathbf{A}\mathbf{A}^T = \boldsymbol{\Sigma}_{\mathbf{X}}$. The random vector $\mathbf{Z} = \mathbf{A}^{-1}(\mathbf{X} - \boldsymbol{\mu}_{\mathbf{X}})$ is an n -dimensional standard normal random vector. Similarly, given the n -dimensional standard normal random vector \mathbf{Z} , an invertible $n \times n$ matrix \mathbf{A} and an n -dimensional vector \mathbf{b} , the random vector $\mathbf{X} = \mathbf{A}\mathbf{Z} + \mathbf{b}$ is an n -dimensional Gaussian random vector with expected value $\boldsymbol{\mu}_{\mathbf{X}} = \mathbf{b}$ and covariance $\boldsymbol{\Sigma}_{\mathbf{X}} = \mathbf{A}\mathbf{A}^T$.

Chapter 9

Concentration Inequalities

9.1 Sample Mean

Let X_1, \dots, X_n be independent and identically distributed random samples drawn from a population with variance σ^2 and expected value μ . The **sample mean** of this set of n samples is denoted $M_n(X)$ and is defined as

$$M_n(X) = \frac{1}{n} \sum_{i=1}^n X_i.$$

Note that the sample mean is a random variable because it is a function of random variables. Its expected value and variance can be expressed as

$$\mathbb{E}[M_n(X)] = \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n X_i\right] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[X_i] = \frac{1}{n} \sum_{i=1}^n \mu = \frac{1}{n}(n)\mu = \mu$$

and

$$\text{Var}(M_n(X)) = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(X_i) = \frac{1}{n^2} \sum_{i=1}^n \sigma^2 = \frac{1}{n}\sigma^2.$$

The expected value of the sample mean is simply the mean of the population, and the variance of the sample mean is the variance of the population divided by the number of samples. As the number of samples approaches infinity, the variance of the sample mean approaches zero. Therefore, the sample mean converges to the population mean as the number of samples approaches infinity.

Now suppose that X_1, \dots, X_n are correlated random samples drawn from a population with variance σ^2 , expected value μ , and correlation coefficient $\rho > 0$. For this case, the expected value of the sample mean would not change. However,

the variance for the sample mean would now be given by

$$\begin{aligned}
 \text{Var}(M_n(X)) &= \text{Var}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n^2} \text{Var}\left(\sum_{i=1}^n X_i\right) \\
 &= \frac{1}{n^2} \text{Cov}\left(\sum_{i=1}^n X_i, \sum_{i=1}^n X_i\right) \\
 &= \frac{1}{n^2} \left(\sum_{i=1}^n \text{Var}(X_i) + \sum_{i=1}^n \sum_{j \neq i} \text{Cov}(X_i, X_j) \right) \\
 &= \frac{1}{n^2} \left(\sum_{i=1}^n \sigma^2 + \sum_{i=1}^n \sum_{j \neq i} \rho \sigma^2 \right) \\
 &= \frac{1}{n^2} (n\sigma^2 + n(n-1)\rho\sigma^2) \\
 &= \frac{1 + (n-1)\rho}{n} \sigma^2.
 \end{aligned}$$

Notice that the variance of the sample mean no longer approaches zero as the number of samples approaches infinity. Therefore, the sample mean may not converge to the population mean as the number of samples approaches infinity.

9.2 Probability Bounds

9.2.1 Markov Inequality

Theorem: The **Markov inequality** says that for a non-negative random variable X and a positive constant $c > 0$, the following inequality holds:

$$P(\{X \geq c\}) \leq \frac{\mathbb{E}[X]}{c}.$$

Proof: From the definition of the indicator function, we can write:

$$\begin{aligned}
 \mathbf{1}\{X \geq c\} &= \begin{cases} 1 & \text{if } X \geq c \\ 0 & \text{if } X < c \end{cases} \\
 c \cdot \mathbf{1}\{X \geq c\} &= \begin{cases} c & \text{if } X \geq c \\ 0 & \text{if } X < c \end{cases}
 \end{aligned}$$

Now notice that because X is a non-negative random variable,

$$c \cdot \mathbf{1}\{X \geq c\} \leq X.$$

Because the expected value function is monotonically increasing,

$$\mathbb{E}[c \cdot \mathbf{1}\{X \geq c\}] \leq \mathbb{E}[X]$$

From the linearity of expectation, we can see that

$$c \cdot \mathbb{E}[\mathbf{1}\{X \geq c\}] \leq \mathbb{E}[X].$$

Because c is strictly positive, we can divide both sides by c , so

$$\mathbb{E}[\mathbf{1}\{X \geq c\}] \leq \frac{\mathbb{E}[X]}{c}.$$

Finally, from our discussion of common discrete distributions, we know the expected value of an indicator function, which leads us to the desired result:

$$P(\{X \geq c\}) \leq \frac{\mathbb{E}[X]}{c}.$$

Extension: Note that there are several other inequalities that can be derived from the Markov inequality. One useful extension says that for a random variable X with finite variance and $\mathbb{E}[X] = 0$,

$$P(\{X \geq c\}) \leq \frac{\mathbb{E}[X^2]}{\mathbb{E}[X^2] + c^2}, \quad \forall c \geq 0.$$

9.2.2 Chebyshev Inequality

Theorem: The **Chebyshev inequality** says that for any arbitrary random variable X with expected value μ_X and variance $\text{Var}(X)$ and for an arbitrary positive constant $c > 0$, the following inequality holds:

$$P(\{|X - \mu_X| \geq c\}) \leq \frac{\text{Var}(X)}{c^2}.$$

Proof: This inequality comes directly from the Markov inequality.

9.2.3 Chernoff Bound

Theorem: The **Chernoff bound** says that that for any arbitrary random variable X with moment generating function $\phi_X(s)$ and any constant c ,

$$P(\{X \geq c\}) \leq \min_{s \geq 0} e^{-sc} \phi_X(s).$$

Proof: This inequality comes directly from the Markov inequality.

It's interesting to consider the implications of the Chernoff bound for Gaussian random variables. For the Gaussian random variable $X \sim N(\mu, \sigma^2)$, it says

$$\begin{aligned} P(\{X \geq c\}) &\leq \min_{s \geq 0} \left\{ e^{-sc} \cdot \exp\left(\frac{s\mu + s^2\sigma^2}{2}\right) \right\} \\ &= \min_{s \geq 0} \left\{ \exp\left(\frac{s(\mu - 2c) + s^2\sigma^2}{2}\right) \right\}. \end{aligned}$$

To find the optimal value, \hat{s} , that minimizes this objective, we can take the derivative of the objective with respect to s , plug in \hat{s} , and set the expression equal to zero. Following this approach, we get the following result:

$$\begin{aligned} \exp\left(\frac{\hat{s}(\mu - 2c) + \hat{s}^2\sigma^2}{2}\right)\left(\frac{\mu}{2} - c + \sigma^2\hat{s}\right) &= 0 \\ \left(\frac{\mu}{2} - c + \sigma^2\hat{s}\right) &= 0 \\ \hat{s} &= \frac{2c - \mu}{2\sigma^2} \end{aligned}$$

Plugging this expression for \hat{s} into our previous inequality, we get

$$\begin{aligned} P(\{X \geq c\}) &\leq \min_{s \geq 0} \left\{ \exp\left(\frac{s(\mu - 2c) + s^2\sigma^2}{2}\right) \right\} \\ &= \exp\left(\frac{\hat{s}(\mu - 2c) + \hat{s}^2\sigma^2}{2}\right) \\ &= \exp\left(\frac{-\frac{(\mu - 2c)^2}{2\sigma^2} + \frac{(\mu - 2c)^2}{4\sigma^2}}{2}\right) \\ &= \exp\left(-\frac{(\mu - 2c)^2}{8\sigma^2}\right). \end{aligned}$$

Because the Gaussian random variable X is symmetric, we can also write

$$P(\{|X| \geq c\}) \leq 2 \exp\left(-\frac{(\mu - 2c)^2}{8\sigma^2}\right).$$

9.3 Law of Large Numbers

9.3.1 Convergence of Random Variables

There are three commonly used modes of convergence for random variables:

1. Almost Sure Convergence

For a sequence of random variables X_1, X_2, \dots , we say that X_n converges to the random variable X almost surely if

$$P\left(\left\{\lim_{n \rightarrow \infty} X_n = X\right\}\right) = 1.$$

We can equivalently write this expression as

$$P\left(\left\{\omega \in \Omega : \lim_{n \rightarrow \infty} X_n(\omega) = X(\omega)\right\}\right) = 1.$$

This means that the values of X_n approach the value of X , in the sense that events for which X_n does not converge to X have probability zero.

2. Convergence in Probability

For a sequence of random variables X_1, X_2, \dots , we say that X_n converges to the random variable X in probability if for every constant $\epsilon > 0$,

$$\lim_{n \rightarrow \infty} P(\{|X_n - X| > \epsilon\}) = 0.$$

3. Convergence in Distribution

For a sequence of random variables X_1, X_2, \dots with corresponding CDFs F_{X_1}, F_{X_2}, \dots , we say that X_n converges to the random variable X with CDF F_X in distribution if for every $x \in \mathbb{R}$ where the CDFs are continuous,

$$\lim_{n \rightarrow \infty} F_{X_n}(x) = F_X(x).$$

Note that almost sure convergence implies convergence in probability, and convergence in probability implies convergence in distribution. These implications do not hold in the opposite direction.

9.3.2 Weak Law of Large Numbers

Suppose that the random variable X represents a population from which samples are drawn. Let μ_X be the expected value of X , σ^2 be the variance, and $M_n(X)$ be the sample mean. The **weak law of large numbers** for a finite number of samples says that for any positive constant $c > 0$, the following holds:

$$(1) P(\{|M_n(X) - \mu_X| \geq c\}) \leq \frac{\sigma^2}{nc^2}$$

and

$$(2) P(\{|M_n(X) - \mu_X| < c\}) \geq 1 - \frac{\sigma^2}{nc^2}.$$

These inequalities come directly from the Chebyshev inequality. From our definitions of convergence, the weak law of large numbers says that the empirical mean $M_n(X)$ converges to the true mean μ_X in probability.

Taking the limit of these inequalities as n goes to infinity, we get the weak law of large numbers for an infinite number of samples:

$$\lim_{n \rightarrow \infty} P(\{|M_n(X) - \mu_X| \geq c\}) = 0$$

$$\lim_{n \rightarrow \infty} P(\{|M_n(X) - \mu_X| < c\}) = 1$$

This finding aligns with our expectations. If $M_n(X)$ is the average of a very large sample of independent and identically distributed random variables, then we expect that any given realization of $M_n(X)$ is very close to the expected value μ_X . As the number of samples goes to infinity, we expect that any realization from $M_n(X)$ will be equal to the expected value of the individual random variables. Therefore, as the number of samples goes to infinity, $P(\{|M_n(X) - \mu_X| > 0\}) = 0$.

9.3.3 Strong Law of Large Numbers

The **strong law of large numbers** says that for any random variable X with expected value μ_X and sample mean $M_n(X)$, the sample mean converges to the true mean μ_X almost surely. Note that the strong law of large numbers implies the weak law because almost sure convergence implies convergence in probability. This law does not say that anomalies cannot happen for a very large number n , but with probability one, these anomalies will not happen.

9.4 Central Limit Theorem

Let X_1, \dots, X_n be independent and identically distributed random variables with variance $\text{Var}(X_i) = \sigma^2 < \infty$ and expected value $\mathbb{E}[X_i] = \mu$. Now consider a new random variable S_n , which is defined as the following sum:

$$S_n = \frac{1}{\sqrt{n\sigma^2}} \left(\sum_{i=1}^n X_i - n\mu \right).$$

The **central limit theorem** says that S_n converges to the standard normal random variable $Z \sim \mathcal{N}(0, 1)$ in distribution, meaning that for all $x \in \mathbb{R}$,

$$P(\{S_n \leq x\}) \rightarrow \phi(x).$$

The central limit theorem is useful because it says that when independent random variables are summed up, their properly normalized sum tends toward a standard normal distribution, even if the original variables themselves are not normally distributed. This implies that probabilistic and statistical methods that work for normal distributions can be applied to many problems involving other types of distributions, assuming a large enough sample size.