

MATH 323: Probability

Anna Brandenberger

November 24, 2018

This is a transcript of the lectures given by Prof. David Stephens¹ during the fall semester of the 2018-2019 academic year (09-12 2018) for the Probability class (MATH 323). **Subjects covered** are: sample space, events, conditional probability, independence of events, Bayes' Theorem; basic combinatorial probability, random variables, discrete and continuous univariate and multivariate distributions; independence of random variables; inequalities, weak law of large numbers, central limit theorem.

¹ david.stephens@mcgill.ca, BH 1225

1 Basics of Probability

1.1 Review of Set Theory

Definition 1.1. A **set** S is a collection of elements $s \in S$:

- **FINITE:** finite number of elements
- **COUNTABLE:** countably infinite number of elements
- **UNCOUNTABLE:** uncountably infinite number of elements

Definition 1.2. The **power set** $\mathcal{P}(A)$ of A is the set of all subsets of A : $\mathcal{P}(A) = \{B \mid B \subseteq A\}$.

SET OPERATIONS

1. **INTERSECTION:** $s \in A \cap B \iff s \in A \text{ and } s \in B$
extends to A_1, A_2, \dots, A_K (finite): $s \in \bigcap_{k=1}^K A_k \iff s \in A_k \forall k$
2. **UNION:** $s \in A \cup B \iff s \in A \text{ or } s \in B$
extends to A_1, A_2, \dots, A_K (finite²): $s \in \bigcup_{k=1}^K A_k \iff \exists k \text{ s.t. } s \in A_k$
3. **COMPLEMENT:** for $A \subseteq S$, $s \in A' \iff s \in S \text{ but } s \notin A$
4. **SET DIFFERENCE** $A - B \equiv A \setminus B \equiv A \cap B'$
5. **EXCLUSIVE OR (XOR):** $A \oplus B \equiv A \cup B - A \cap B$

Definition 1.3. A_1, A_2, \dots, A_K is called a **partition** of S if these sets are pairwise disjoint $A_j \cap A_k = \emptyset \forall j \neq k$ and exhaustive $\bigcup_{k=1}^K A_k = S$.

Theorem 1.1 (De Morgan's Laws).

$$A' \cap B' = (A \cup B)' \quad A \cup B' = (A \cap B)'$$

Probability is a numerical value representing the chance of a particular event occurring given a particular set of circumstances.

Let $A \subseteq S$:

$$\begin{aligned} A \cap \emptyset &= \emptyset & A \cup \emptyset &= A \\ A \cap S &= A & A \cup S &= S \end{aligned}$$

Intersection \cap and union \cup are commutative and associative.

² technically also extends to a countably infinite number of sets.

Proof. Let $A_1 = A \cup B$, $A_2 = A' \cap B'$. To show that $A'_1 = A_2$, one must show $A_1 \cup A_2 = \emptyset$ and $A_1 \cup A_2 = S$.

$$\begin{aligned} A_1 \cap A_2 &= (A \cup B) \cap A_2 = (A \cap A_2) \cup (B \cap A_2) \\ &= (A \cap A' \cap B') \cup (B \cap A' \cap B') \\ &= \emptyset \cup \emptyset = \emptyset \end{aligned}$$

$$\begin{aligned} A_1 \cup A_2 &= A_1 \cup (A' \cap B') = (A_1 \cup A') \cap (A_1 \cup B') \\ &= (A \cup B \cup A') \cap (A \cup B \cup B') \\ &= S \cap S = S \quad \square \end{aligned}$$

1.2 Sample Spaces and Events

Definition 1.4. The set S of possible outcomes of an experiment is the **sample space** of an experiment, and the individual elements in S are **sample points** or **outcomes**.

Definition 1.5. An **event** A is a collection of sample outcomes $A \subseteq S$. Individual sample outcomes $E_1, E_2, \dots, E_K, \dots \in A$ are termed **simple (elementary) events**. We say that A **occurs** if the outcome $s \in A$. S is the **certain event**; \emptyset the **impossible event**.

Definition 1.6. The **probability function** $P(\cdot)$ assigns numerical values to events: $P : \mathcal{A} \rightarrow \mathbb{R}$, $A \rightarrow p$ (i.e. $P(A) = p$, $A = \mathcal{P}(S)$).

1.3 Axioms and Consequences

Theorem 1.2. Basic probability axioms:

1. $P(A) \geq 0$
2. $P(S) = 1$
3. P is countably additive: if A_1, A_2, \dots s.t. $A_j \cap A_k = \emptyset \forall j \neq k$:

$$P\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} P(A_i)$$

Corollary 1.3. $\forall A \ P(A') = 1 - P(A)$

Corollary 1.4. $P(\emptyset) = 0$

Corollary 1.5. $\forall A \ P(A) \leq 1$

Corollary 1.6. $\forall A \subseteq B, P(A) \leq P(B)$

Corollary 1.7 (General Addition Rule).

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

Using inductive arguments, one can construct a formula for $P(\bigcup_{i=1}^n A_i)$ (c.f. MATH 240) but we will mostly use:

Theorem 1.8 (Boole's Inequality).

$$P\left(\bigcup_{i=1}^n A_i\right) \leq \sum_{i=1}^n P(A_i)$$

Proof. $S = A \cup A'$. By Theorem 1.2(3),
 $P(S) = P(A) + P(A') = 1$
 $\implies P'(A) = 1 - P(A)$ \square

Proof. Simple from Corollary 1.3. \square

Proof. By Theorem 1.2 and 1.3,
 $1 = P(S) = P(A) + P(A') \geq P(A)$ \square

Proof. Take $B = A \cup (A' \cap B)$. Then,
 $P(B) = P(A) + P(A' \cap B) \geq P(A)$ \square

Proof. First, $(A \cup B) = A \cup (A' \cap B)$
 $\implies P(A \cup B) = P(A) + P(A' \cap B)$ [1]
 Then, take $B = (A \cap B) \cup (A' \cap B)$
 $\implies P(B) = P(A \cap B) + P(A' \cap B)$ [2]
 So taking [1] - [2], we get:
 $P(A \cup B) - P(B) = P(A) - P(A \cap B)$
 $P(A \cup B) = P(A) + P(B) - P(A \cap B)$ \square

As a method, one can lay out probabilities in table form such as in Table 1, ensuring that $0 \leq p_* \leq 1$ and probabilities along rows and column add up correctly.

| | A | A' | Total |
|-------|-----------------|------------------|----------|
| B | $p_{A \cap B}$ | $p_{A' \cap B}$ | p_B |
| B' | $p_{A \cap B'}$ | $p_{A' \cap B'}$ | $p_{B'}$ |
| Total | p_A | $p_{A'}$ | p_S |

Table 1: Probabilities laid out in a table.

1.4 Specifying probabilities

Definition 1.7 (Equally likely sample outcomes). Suppose S finite, with sample outcomes E_1, \dots, E_N , then $P(E_i) = \frac{1}{N} \forall i = 1, \dots, N$. For an event $A \in S$, $\exists n \leq N$ s.t. $A = \bigcup_{i=1}^n E_{i_A}$, $i_A \in \{1, \dots, n\}$

$$\implies P(A) = \frac{n}{N} = \frac{\# \text{ sample outcomes in } A}{\# \text{ sample outcomes in } S}$$

Definition 1.8 (Relative frequencies). For a given experiment with a sample space S and interest event A , consider a finite sequence of N repeat experiments and n the number of times that A occurs. The **frequentist definition** of probability, generalizes the previous definition to $P(A) = \lim_{N \rightarrow \infty} \frac{n}{N}$, the relative frequency of appearance of A .

Definition 1.9 (Subjective assessment). For a given experiment with sample space S , the probability of event A is further generalized to a numerical representation of one's own **personal degree of belief** that the actual outcome lies in A , assuming one is **internally consistent**, i.e. rational and coherent.

Examples of experiments with equally likely sample outcomes include fair coins, dice and cards.

1.5 Combinatorial Rules

- MULTIPLICATION PRINCIPLE:** a sequence of k operations i , each with n_i possible outcomes, results in $n_1 \times n_2 \times \dots \times n_k$ possible sequences of outcomes.
- SELECTION PRINCIPLES:** selecting from a finite set can be done:
 - WITH REPLACEMENT:** each selection is independent;
 - WITHOUT REPLACEMENT:** set is depleted by each selection.
- ORDERING:** the order of a sequence of selections can be:
 - IMPORTANT** (312 distinct from 123);
 - UNIMPORTANT** (312 identical to 123).
- DISTINGUISHABLE ITEMS:** objects being selected can be:
 - DISTINGUISHABLE:** individually labelled (e.g. lottery balls);
 - INDISTINGUISHABLE:** labelled according to a type (e.g. colors).

c.f. MATH 240 notes for examples.

Definition 1.10. A **permutation** is an ordered arrangement of distinct objects. The number of ways of ordering an n -set taking r at a time is

$$P_r^n = \frac{n!}{(n-r)!}$$

Definition 1.11. The number of ways of partitioning n distinct elements into k disjoint subsets of sizes n_1, \dots, n_k is the **multinomial coefficient**

$$N = \binom{n}{n_1, \dots, n_k} := \frac{n!}{n_1! \times n_2! \times \dots \times n_k!}$$

a generalization of the **binomial coeff.** $\binom{n}{j} = \frac{n!}{j!(n-j)!} := \binom{n}{j, n-j}$.

Definition 1.12. The number of **combinations** C_r^n is the number of subsets of size r that can be selected from n objects

$$C_r^n = \binom{n}{r} = n! \cdot P_r^n$$

1.6 Conditional Probability and Independence

Definition 1.13. The **conditional probability** of A given B , where $P(B) > 0$ is

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

Corollary 1.9. $P(A | S) = P(A)$ $P(A | A) = 1$ $P(A | B) \leq 1$

Warning 1. Important distinction between $P(A \cap B)$ and $P(A | B)$:

- $P(A \cap B)$: chance of A and B occurring **RELATIVE TO** S
- $P(A | B)$: chance of A and B occurring **RELATIVE TO** B — non-symmetric: in general, $P(A | B) \neq P(B | A)$.

Theorem 1.10. The conditional probability function satisfies the probability axioms (Theorem 1.2).

Definition 1.14. Two events $A, B \in S$ are **independent** if

$$P(A | B) = P(A) \iff P(A \cap B) = P(A)P(B)$$

Remark 2. Notice then that $P(B | A) = \frac{P(A \cap B)}{P(A)} = \frac{P(A)P(B)}{P(A)} = P(B)$.

Warning 3. Not the same as mutual exclusivity where $P(A \cap B) = 0$!

Given partial knowledge (B happened) i.e. knowing $s \in B \subseteq S$, one can make a more accurate probability assessment of some other event A , as $s \in A \cap B$.

Example 1.1. Given three two-sided cards, one RR, one RB and one BB, one card is selected randomly and one side displayed: it is red (R). Find the probability that the other side is red.

Answer. $S = \{R_1, R_2, B, R, B_1, B_2\}$
 Card 1 (RR) selected $A = \{R_1, R_2\} \implies P(A) = 2/6$
 Red side exposed: $B = \{R_1, R_2, R\} \implies P(B) = 3/6$
 $P(A | B) = \frac{P(A \cap B)}{P(B)} = \frac{2/6}{3/6} = 2/3$

Proof. 1. $P(A | B) = \frac{P(A \cap B)}{P(B)} \geq 0$

2. $P(S | B) = \frac{P(S \cap B)}{P(B)} = \frac{P(B)}{P(B)} = 1$

3. $P(\bigcup_{i=1}^{\infty} (A_i | B)) = \frac{P(\bigcup_{i=1}^{\infty} A_i \cap B)}{P(B)} = \frac{\sum_{i=1}^{\infty} P(A_i \cap B)}{P(B)} = \sum_{i=1}^{\infty} P(A_i | B)$ and since A_i disjoint $\forall i$, (A_i, B) also disjoint. \square

Warning 4. For multiple events $A_1, \dots, A_n \in S$, we can talk about pairwise independence between any A_i, A_j with $i \neq j$, but this does not automatically generalize (c.f. Example 1.2).

Definition 1.15. Events A_1, \dots, A_K are **mutually independent** if, for all subsets \mathcal{I} of $\{1, 2, \dots, K\}$:

$$P\left(\bigcap_{k \in \mathcal{I}} A_k\right) = \prod_{k \in \mathcal{I}} P(A_k)$$

Definition 1.16. Consider events $A_1, A_2, B \in S$ with $P(B) > 0$. A_1 and A_2 are **conditionally independent** given B if

$$\begin{aligned} P(A_1 | A_2 \cap B) &= P(A_1 | B) \\ \iff P(A_1 \cap A_2 | B) &= P(A_1 | B)P(A_2 | B) \end{aligned}$$

Theorem 1.11 (General Multiplication Rule).

$$P(A_1, A_2, \dots, A_K) = P(A_1)P(A_2 | A_1)P(A_3 | A_1 \cap A_2) \cdots P(A_K | A_1 \cap A_2 \cap \cdots \cap A_{K-1})$$

Example 1.2. Rolling two dice with independent outcomes A_1 : first roll outcome is odd; A_2 : second roll outcome is odd and A_3 : total outcome is odd. Then $P(A_1) = P(A_2) = \frac{1}{2}$ and $P(A_1 | A_3) = P(A_2 | A_3) = \frac{1}{2}$. But $P(A_1 \cap A_2 \cap A_3) = 0$, $P(A_1 | A_2 \cap A_3) = 0$, etc.

Proof. Trivial, recursively (Chain Rule for probabilities). \square

We can use **probability trees** to display joint probabilities of multiples, where the junctions are the events and the branches correspond to possible sequences of choices of events. We then *multiply along a branch* to get the joint probability.

Proof. (Theorem 1.12)

$$\begin{aligned} A &= (A \cap B_1) \cup \cdots \cup (A \cap B_n) \\ P(A) &= P(A \cap B_1) + \cdots + P(A \cap B_n) \\ &= P(A | B_1)P(B_1) + \cdots + P(A | B_n)P(B_n) \\ &= \sum_{i=1}^n P(A | B_i)P(B_i) \end{aligned}$$

Proof. (Theorem 1.13). By Definition 1.13,

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

With the stated conditions, we have

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

The second part of the theorem comes from replacing the denominator using Theorem 1.12. \square

1.7 Important Theorems

Theorem 1.12 (Theorem of Total Probability). Given a partition B_1, \dots, B_n of the sample space S into n subsets s.t. $P(B_i) > 0$, then for any event $A \in S$,

$$P(A) = \sum_{i=1}^n P(A | B_i)P(B_i)$$

Theorem 1.13 (Bayes's Theorem). For any $A, B \in S$ where $P(A) > 0$ and $P(B) > 0$:

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

Furthermore, given B_1, B_2, \dots, B_n forming a partition of S :

$$P(B_i | A) = \frac{P(A | B_i)P(B_i)}{P(A)} = \frac{P(A | B_i)P(B_i)}{\sum_{j=1}^n P(A | B_j)P(B_j)}$$

Warning 5. Although $\sum_{i=1}^n P(B_i | A) = 1$, $\sum_{i=1}^n P(A | B_i) \neq 1$ due to the use of different conditioning sets.

Remark 6. Bayes' Theorem is often used to make probability statements concerning an event B that *has not* been observed, given an event A that *has* been observed.

Remark 7 (Prosecutor's Fallacy).

$$P(\text{"Evidence"} \mid \text{"Guilt"}) \neq P(\text{"Guilt"} \mid \text{"Evidence"})$$

Remark 8. Recall that the odds on B are defined as $\frac{P(B)}{P(B')} = \frac{P(B)}{1-P(B)}$; and the conditional odds given A ($P(A) > 0$) are $\frac{P(B|A)}{P(B'|A)} = \frac{P(B|A)}{1-P(B|A)}$. Thus the odds change by a factor of $\frac{P(A|B)}{P(A|B')}$ when one adds the information A .

Corollary 1.14 (Bayes' Theorem for multiple events). Given events $A_1, A_2 \in S$ where $P(A_1 \cap A_2) > 0$ and given B_1, B_2, \dots, B_n forming a partition of S :

$$P(B_i \mid A_1 \cap A_2) = \frac{P(A_1 \cap A_2 \mid B_i)P(B_i)}{\sum_{k=1}^n P(A_1 \cap A_2 \mid B_k)P(B_k)}$$

Example 1.3 (Medical screening). Given events A (positive screening test) and B (person is actually a sufferer). We are interested in $P(B \mid A)$. Suppose $P(B) = p$; $P(A \mid B) = 1 - \alpha$ (true positive rate); $P(A \mid B') = \beta$ (false positive rate) for probabilities p, α, β .

$$\begin{aligned} P(A) &= P(A \mid B)P(B) + P(A \mid B')P(B') \\ &= (1 - \alpha)p + \beta(1 - p) \end{aligned}$$

$$\begin{aligned} P(B \mid A) &= \frac{P(A \mid B)P(B)}{P(A)} \\ &= \frac{(1 - \alpha)p}{(1 - \alpha)p + \beta(1 - p)} \end{aligned}$$

$$\begin{aligned} P(B \mid A') &= \frac{\alpha p}{\alpha p + (1 - \beta)(1 - p)} \\ &\text{by similarly logic.} \end{aligned}$$

2 Random Variables and Probability Distributions

2.1 Random Variables

Definition 2.1. A **random variable** is a mapping $Y : S \rightarrow \mathbb{R}$, $s \mapsto y$ from an arbitrary space S (where the experiment sample outcomes are defined) to the real line \mathbb{R} .

Remark 9. We use uppercase letters X, Y, Z to denote random variables and lowercase x, y, z to denote the real values taken by the random variable.

$Y(\cdot)$ is usually a many-to-one mapping. Also, we usually suppress the dependence

Definition 2.2. A random variable is **discrete** if the image of Y is a countable set $\{y_1, y_2, \dots, y_n, \dots\}$ denoted \mathcal{Y} for positive probabilities.

Corollary 2.1. The probabilities assigned to individual points in the set \mathcal{Y} must satisfy

$$\sum_{y \in \mathcal{Y}} P(Y = y) = 1$$

Proof. follows from the probability axioms. \square

2.2 Probability Mass Function

Definition 2.3. The **probability mass function** (pmf) $p(\cdot)$ for the random variable Y records how probability is distributed across points in \mathbb{R} : $p(y) = P(Y = y)$ where $p(y) > 0$ for $y \in \mathcal{Y}$. It is also denoted $f(y)$, $p_Y(y)$ or $f_Y(y)$. It must:

- (a) specify probabilities, i.e. $0 \leq p(y) \leq 1$
- (b) satisfy Corollary 2.1: $\sum_{y \in \mathcal{Y}} p(y) = 1$

Remark 10. The pmf can be defined pointwise (for a discrete random variable) or in some functional form.

Remark 11. All pmfs take the form $p(y) = cg(y)$ where $g(y)$ also contains information on the set \mathcal{Y} for which $P(y) > 0$.

$$\sum_{y \in \mathcal{Y}} p(y) = 1 \implies \sum_{y \in \mathcal{Y}} g(y) = \frac{1}{c}$$

2.3 Moments: expectation and variance

Definition 2.4. The **expectation** (expected value) $\mathbb{E}[Y]$ of the r.v. Y is $\mathbb{E}[Y] = \sum_{y \in \mathcal{Y}} yp(y)$, the centre of mass of the probability distribution.

Let's now deal with the expectation of a function $g(Y)$.

Definition 2.5. Let $g(\cdot)$ real-valued function and $p(y)$ the pmf of Y :

Set of values of $g(Y)$: $\mathcal{G} = \{g_i, i = 1, \dots, m : g_i = g(y), \text{ some } y \text{ in } \mathcal{Y}\}$

$$\mathbb{E}[g(Y)] = \sum_{y \in \mathcal{Y}} g(y)p(y)$$

Example 2.1 (Discrete Uniform Distribution).

Suppose $\mathcal{Y} = \{y_1, \dots, y_n\}$, $p(y) = \frac{1}{n}$. Then $\mathbb{E}[Y] = \frac{1}{n} \sum_{i=1}^n y_i = \bar{y}$. Here $p(y) = c$ is constant and the number of elements in \mathcal{Y} is $\frac{1}{c}$.

Definition 2.6. The **variance** $\mathbb{V}[Y]$ of Y is a measure of how dispersed the probability distribution is, defined as

$$\begin{aligned}\mathbb{V}[Y] &= \mathbb{E}[(Y - \mu)^2] \text{ where } \mu = \mathbb{E}[Y] \\ &= \sum_y (y - \mu)^2 p(y) \geq 0\end{aligned}$$

Definition 2.7. The **standard deviation** of Y is $\sqrt{\mathbb{V}[Y]}$.

Theorem 2.2. The expectation $\mathbb{E}[Y]$ satisfies linearity.

Remark 12. $\mathbb{V}[Y]$ can be infinite even if $\mathbb{E}[Y]$ is bounded.

Theorem 2.3. $\mathbb{E}[(Y - \mu)^2] = \mathbb{E}[Y^2] - \mu^2$

Corollary 2.4. $\mathbb{V}[aY + b] = a^2 \mathbb{V}[Y]$

2.4 Binomial distribution

Considering an experiment with two outcomes 'success' and 'failure', define the random variable Y to be $Y = 1$ for 'success' and $Y = 0$ for 'failure', with $P(1) = p$, $P(0) = 1 - p$.

Definition 2.8. The **Bernoulli** distribution is defined with pmf

$$p_Y(y) = p^y q^{1-y} = p^y (1-p)^{1-y} \quad \forall y \in \{0, 1\}, \text{ given } p \in (0, 1)$$

Theorem 2.5. The expectation of the Bernoulli dist. is $\mu = \mathbb{E}[Y] = p$.

Theorem 2.6. The variance of the Bernoulli dist. is $\mathbb{V}[Y] = p(1-p)$.

We now extend the definition to a sequence of multiple experiments with two outcomes, and have the r.v. Y record the number of 'successes': $Y = y \Leftrightarrow$ "y successes in n trials". Consider a sequence with y successes and $n - y$ failures: the probability of any such sequence is $p^y q^{n-y}$. The number of possible such sequences is $C_y^n = \binom{n}{y}$.

Definition 2.9. The **Binomial** distribution is defined with pmf

$$p_y(y) \equiv P(Y = y) = \binom{n}{y} p^y q^{n-y} \quad \forall y \in \{0, 1, 2, \dots, n\}$$

Theorem 2.7. The expectation of the binomial dist. is $\mathbb{E}[Y] = np$.

Proof. (Theorem 2.2)

Let $g(Y) = ag_1(Y) + bg_2(Y)$ for separate $g_1(\cdot)$ and $g_2(\cdot)$.

$$\begin{aligned}\mathbb{E}[g(Y)] &= \mathbb{E}[ag_1(Y) + bg_2(Y)] \\ &= \sum_y (ag_1(y) + bg_2(y))p(y) \\ &= a \sum_y g_1(y)p(y) + b \sum_y g_2(y)p(y) \\ &= a\mathbb{E}[g_1(Y)] + b\mathbb{E}[g_2(Y)] \quad \square\end{aligned}$$

Proof. (Theorem 2.3)

$$\begin{aligned}\mathbb{E}[(Y - \mu)^2] &= \mathbb{E}[Y^2 + 2Y\mu + \mu^2] \\ &= \mathbb{E}[Y^2] - 2\mu\mathbb{E}[Y] + \mu^2 \\ &= \mathbb{E}[Y^2] - 2\mu^2 + \mu^2 \\ &= \mathbb{E}[Y^2] - \mu^2 \quad \square\end{aligned}$$

Proof. (Corollary 2.4)

$$\begin{aligned}\mathbb{V}[aY + b] &= \mathbb{E}[(aY + b - \mathbb{E}[aY + b])^2] \\ &= \mathbb{E}[(aY + b - (a\mathbb{E}[Y] + b))^2] \\ &= \mathbb{E}[(aY + b - a\mu - b)^2] \\ &= \mathbb{E}[a^2(Y - \mu)^2] \\ &= a^2 \mathbb{V}[Y] \quad \square\end{aligned}$$

Proof. (Theorem 2.5)

$$\begin{aligned}\mathbb{E}[Y] &= \sum_{i=0}^1 yp(y) = 0p(0) + 1p(1) \\ &= 0 \cdot q + 1 \cdot p = p \quad \square\end{aligned}$$

Proof. (Theorem 2.6)

$$\begin{aligned}\mathbb{V}[Y] &= \mathbb{E}[Y^2] - \mu^2 \\ &= \sum_{i=0}^1 y^2 p(y) - \mu^2 = 1^2 p(1) \\ &= p - p^2 = p(1-p) = pq \quad \square\end{aligned}$$

Proof. (Theorem 2.7)

$$\begin{aligned}\mathbb{E}[Y] &= \sum_{y=0}^n y \binom{n}{y} p^y q^{n-y} \\ &= n \sum_{y=1}^n \frac{(n-1)!}{(y-1)!(n-y)!} p^y q^{n-y} \\ &= n \sum_{j=0}^{n-1} \frac{(n-1)!}{j!((n-1)-j)!} p^{j+1} q^{n-j-1} \\ &= np \sum_{j=0}^{n-1} \frac{(n-1)!}{j!((n-1)-j)!} p^j q^{n-1-j} \\ &= np(p+q)^{n-1} = np \quad \square\end{aligned}$$

Theorem 2.8. The variance of the binomial dist. is $\mathbb{V}[Y] = npq$.

Proof.

$$\begin{aligned}
 \mathbb{E}[Y^2] &= \mathbb{E}[Y(Y-1)] + \mathbb{E}[Y] \text{ (by linearity)} \\
 \mathbb{E}[Y(Y-1)] &= \sum_{y=0}^n y(y-1) \binom{n}{y} p^y q^{n-y} = \sum_{y=2}^n y(y-1) \frac{n!}{y!(n-y)!} p^y q^{n-y} \\
 &= n(n-1) \sum_{y=2}^n \frac{(n-2)!}{(y-2)!(n-y)!} p^y q^{n-y} \\
 &= n(n-1) \sum_{j=0}^{n-2} \frac{(n-2)!}{j!(n-2-j)!} p^{j+2} q^{n-j-2} \\
 &= n(n-1)p^2 \sum_{j=0}^{n-2} \frac{(n-2)!}{j!(n-2-j)!} p^j q^{n-j-2} \\
 &= n(n-1)p^2(p+q)^{n-2} = n(n-1)p^2 \\
 \mathbb{V}[Y] &= \mathbb{E}[Y^2] - \mu^2 = n(n-1)p^2 + np - (np)^2 \\
 &= np - np^2 = npq
 \end{aligned}$$

Remark 13. Notice how the Bernoulli distribution is a special case of the binomial distribution with $n = 1$.

2.5 Geometric Distribution

Consider independent binary trials, where one counts the number of trials carried out until the first success. Let Y be the r.v. recording the number of trials carried out up to and including the first success: $Y = y \Leftrightarrow "y-1 \text{ failures followed by } 1 \text{ success}"$.

Definition 2.10. The **geometric** distribution is defined with pmf

$$p(y) \equiv P(Y = y) = q^{y-1}p \quad \forall y \in \{1, 2, \dots\}$$

Remark 15. The geometric distribution can also be written with pmf

$$p(y) \equiv P(Y^* = y) = q^y p \quad \forall y \in \{0, 1, 2, \dots\}$$

where the r.v. Y^* then records the number of failures observed before the first success: $Y^* = Y - 1$.

Theorem 2.9. The expectation of the geometric dist. is $\mathbb{E}[Y] = \frac{1}{p}$.

Theorem 2.10. The variance of the geometric dist. is $\mathbb{V}[Y] = \frac{1-p}{p^2}$.

Definition 2.11. The **cumulative distribution function** (cdf) of a distribution is defined as $F(y) = P(Y \leq y)$.

Remark 14. Notice that

$$\sum_{y=1}^{\infty} q^{y-1}p = p \sum_{j=0}^{\infty} q^j = p \frac{1}{1-q} = 1$$

Proof. (of Theorem 2.9)

$$\begin{aligned}
 \mathbb{E}[Y] &= \sum_{y=0}^{\infty} yq^{y-1}p = p \sum_{y=1}^{\infty} yq^{y-1} \\
 &= p \frac{1}{(1-q)^2} = p \frac{1}{p^2} = \frac{1}{p} \quad \square
 \end{aligned}$$

Proof. (of Theorem 2.10)

$$\begin{aligned}
 \mathbb{E}[Y(Y-1)] &= \sum_{y=0}^{\infty} y(y-1)pq^y \\
 &= pq \sum_{y=2}^{\infty} y(y-1)q^{y-2} = pq \frac{d^2}{dq^2} \left(\sum_{y=0}^{\infty} q^y \right) \\
 &= pq \frac{d^2}{dq^2} \left(\frac{1}{1-q} \right) = pq \frac{2}{(1-q)^3} = 2 \frac{q}{p^2} \\
 \mathbb{V}[Y] &= \mathbb{E}[Y(Y-1)] + \mathbb{E}[Y] - (\mathbb{E}[Y])^2 \\
 &= 2 \frac{1-p}{p^2} + \frac{1}{p} - \frac{1}{p^2} = \frac{1-p}{p^2}
 \end{aligned}$$

□

Theorem 2.11. The cdf of the geometric distribution is

$$F(y) = P(Y \leq y) = 1 - (1 - p)^y \quad y \in \{1, 2, \dots\}$$

Proof.

$$\begin{aligned} P(Y > y) &= \sum_{j=y+1}^{\infty} P(Y = j) = \sum_{j=y+1}^{\infty} q^{j+1} p = pq^y \sum_{j=y+1}^{\infty} q^{j-y-1} = \frac{pq^y}{1-q} \\ &= q^y \quad \forall y \in \{1, 2, \dots\} \\ P(Y \leq y) &= 1 - P(Y > y) = 1 - (1 - p)^y \quad \forall y \in \{1, 2, \dots\} \quad \square \end{aligned}$$

We can define all 'interval' probabilities such as $P(Y \leq a)$, $P(a \leq Y \leq b)$, etc. via $F(y)$, since \forall distributions,

$$P(Y \in \mathcal{X}) = \sum_{y \in \mathcal{X}} p(y)$$

2.6 Negative Binomial Distribution

Consider independent binary trials, where one counts the number of trials carried out until the r th success. Let r.v. Y record the number of trials carried out until the r th success: $Y = y \Leftrightarrow A \cap B$: $A = "r - 1$ successes in the first $y - 1$ trials" and $B = "success on the y th trial".$

Definition 2.12. The **negative binomial** distribution has pmf

$$p(y) \equiv P(Y = y) = \binom{y-1}{r-1} p^r q^{y-r} \quad \forall y \in \{r, r+1, \dots\}$$

Remark 16. As for the geometric distribution, the negative binomial dist. can also be written with pmf

$$p(y) = P(Y^* = y) = \binom{y+r-1}{r-1} p^r q^y \quad \forall y \in \{0, 1, 2, \dots\}$$

where the r.v. Y^* records the number of failures observed up to the r th success: $Y^* = Y - r$.

Theorem 2.12. The expectation of the negative binomial distribution is $\mathbb{E}[Y] = \frac{r}{p}$.

Theorem 2.13. The variance of the negative binomial distribution is $\mathbb{V}[Y] = \frac{r(1-p)}{p^2}$.

Exercise 2.2. Show for the negative binomial distribution that

$$\sum_{y=r}^{\infty} p(y) = 1$$

Proofs left to the reader.

2.7 Hypergeometric Distribution

Consider a finite set of N items with r type I and $N - r$ type II. We select n items without replacement. Let r.v. $Y = y \Leftrightarrow "y$ type I items".

Definition 2.13. The **hypergeometric** distribution is defined with pmf

$$p(y) \equiv P(Y = y) = \frac{\binom{r}{y} \times \binom{N-r}{n-y}}{\binom{N}{n}} \quad \forall y \in \mathcal{Y}$$

$\binom{r}{y}$: number of ways of selecting y type I from the r available.
 $\binom{N-r}{n-y}$: number of ways of selecting $n - y$ type II from the $N - r$ available.
 $\binom{N}{n}$: total number of ways of selecting n items from N .

Remark 17. An equivalent expression for $p(y)$ is

$$p(y) = \frac{\binom{n}{y} \binom{N-n}{r-y}}{\binom{N}{r}}$$

2.8 Poisson distribution

Consider the binomial experiment setting, and take the limiting case where we allow n to get larger and p to get smaller, but keeping $n \times p$ constant: i.e. take $p = \frac{\lambda}{n}$ and let n large. Then,

$$\begin{aligned} p(y) &= \binom{n}{y} p^y (1-p)^{n-y} = \frac{n!}{y!(n-y)!} \left(\frac{\lambda}{n}\right)^y \left(1 - \frac{\lambda}{n}\right)^{n-y} \\ &= \frac{\lambda^y}{y!} \left(1 - \frac{\lambda}{n}\right)^n \underbrace{\frac{n!}{(n-y)! n^y} \left(1 - \frac{\lambda}{n}\right)^{-y}}_{\rightarrow 1 \text{ as } n \rightarrow \infty} \rightarrow \frac{\lambda^y}{y!} e^{-\lambda} \end{aligned}$$

Definition 2.14. The **Poisson** distribution is defined with pmf

$$p(y) = \frac{\lambda^y}{y!} e^{-\lambda} \quad \forall y \in \{0, 1, 2, \dots\}$$

Theorem 2.14. The expectation of the Poisson dist. is $\mathbb{E}[Y] = \lambda$.

Theorem 2.15. The variance of the Poisson dist. is $\mathbb{V}[Y] = \lambda$.

Remark 19. The Poisson distribution corresponds to a model where incidents occur independently at a constant rate λ .

2.9 Moments and Moment-Generating Functions

Definition 2.15. Let $\mu = \mathbb{E}[Y]$. We define

$$\textbf{kth moment: } \mu'_k = \mathbb{E}[Y^k] = \sum_y y^k p(y)$$

$$\textbf{kth central moment: } \mu_k = \mathbb{E}[(Y - \mu)^k] = \sum_y (y - \mu)^k p(y)$$

Definition 2.16. **Generating functions** $g(t)$ for the sequence $\{a_j\} \in \mathbb{R}$ is the function defined by the sum

$$g(t) = \sum_j a_j t^j$$

Definition 2.17. **Exponential generating function:** $g(t) = \sum_j a_j \frac{t^j}{j!}$.

Remark 18. When N and r are large, we can approximate the hypergeometric to a binomial distribution

$$p(y) \approx \binom{n}{y} \left(\frac{r}{N}\right)^y \left(1 - \frac{r}{N}\right)^{n-y}$$

λ : rate at which successes occur in a continuous approximation to the experiment.

An example use case would be: "If Mike makes an average number of 2 sales a day, what are the odds that he makes 3 sales tomorrow?"

Proof. (Theorem 2.14)

$$\begin{aligned} \mathbb{E}[Y] &= \sum_{y=0}^{\infty} y p(y) = \sum_{y=0}^{\infty} y \frac{\lambda^y}{y!} e^{-\lambda} \\ &= e^{-\lambda} \sum_{y=1}^{\infty} \frac{\lambda^y}{(y-1)!} \\ &= e^{-\lambda} \lambda \sum_{y=1}^{\infty} \frac{\lambda^{y-1}}{(y-1)!} \\ &= e^{-\lambda} \lambda e^{\lambda} = \lambda \quad \square \end{aligned}$$

Proof. (Theorem 2.15)

$$\begin{aligned} \mathbb{E}[Y(Y-1)] &= \sum_{y=1}^{\infty} y(y-1) \frac{\lambda^y}{y!} e^{-\lambda} \\ &= e^{-\lambda} \lambda^2 \sum_{y=2}^{\infty} \frac{\lambda^{y-2}}{(y-2)!} = e^{-\lambda} \lambda^2 e^{\lambda} \\ &= \lambda^2 \\ \mathbb{V}[Y] &= \mathbb{E}[Y(Y-1)] + \mathbb{E}[Y] - \mu^2 \\ &= \lambda^2 + \lambda - \lambda^2 = \lambda \quad \square \end{aligned}$$

Example 2.3. The generating function for $\{a_j\} = \binom{n}{j} \quad \forall j = 0, 1, \dots$ is

$$g(t) = (1+t)^n = \sum_{j=0}^n \binom{n}{j} t^j$$

Definition 2.18. For pmf $p(y)$, the **moment generating function** (mgf) $m(t)$ is the exponential generating function for the moments of $p(y)$, requiring the sum to be finite in the neighbourhood of $t = 0$:

$$m(t) = \sum_{j=0}^{\infty} \mu'_j \frac{t^j}{j!}$$

Theorem 2.16. The moment generating function can be written

$$m(t) = \mathbb{E}[e^{tY}]$$

Theorem 2.17. The k th derivative of $m(t)$ evaluated at $t = 0$ gives

$$m^{(k)}(0) = \frac{d^k}{dt^k} m(0) = \mu'_k$$

Theorem 2.18 (Uniqueness). Suppose $m_1(t)$ and $m_2(t)$ are two mgfs for pmfs $p_1(t)$ and $p_2(t)$. Then:

$$m_1(t) = m_2(t) \quad \forall |t| \leq b \text{ for some } b \iff p_1(y) = p_2(y) \quad \forall y$$

Remark 20. It is difficult to see the use for moment generating functions right now, but it can be used to identify distributions, since two random variables with the same mgf share the same distribution (Theorem 2.18), and it can be used to calculate the moments of a random variables.

Proof. (Theorem 2.16)

$$\begin{aligned} m(t) &= \sum_{j=0}^{\infty} \mu'_j \frac{t^j}{j!} = \sum_{j=0}^{\infty} \sum_y y^j p(y) \frac{t^j}{j!} \\ &= \sum_y \left(\sum_{j=0}^{\infty} \frac{(yt)^j}{j!} \right) p(y) \\ &= \sum_y e^{ty} p(y) \quad \square \end{aligned}$$

Proof. (Theorem 2.17)

Simply write out the terms. \square

Proof. (Theorem 2.18)

Beyond the scope of this course. \square

2.10 Probability-Generating Functions

Definition 2.19. The **probability-generating function** (pgf) for random variable Y with $P(Y = i) = y_i$ is

$$G(t) = \mathbb{E}[t^Y] = \sum_{i=0}^{\infty} p_i t^i$$

Definition 2.20. The k th **factorial moment** $\mu_{[k]}$ is defined as

$$\mu_{[k]} = \mathbb{E}[Y(Y-1) \cdots (Y-k+1)]$$

Remark 21. Note that

$$G^{(k)}(t) = \frac{d^{(k)}(t)}{dt^k} \left(\sum_{i=0}^{\infty} p_i t^i \right) = \sum_{i=k}^{\infty} i(i-1) \cdots (i-k+1) p_i t^{i-k}$$

$$G^{(k)}(1) = \sum_{i=k}^{\infty} i(i-1) \cdots (i-k+1) p_i = \mu_{[k]}$$

Corollary 2.19. $G(t) \equiv \mathbb{E}[t^Y] = \mathbb{E}[\exp(Y \ln t)] = m(\ln t)$

2.11 Continuous Random Variables

Theorem 2.20. Properties of the cdf $F(Y)$ (c.f. Definition 2.11):

1. "starts at zero": $F(-\infty) \equiv \lim_{y \rightarrow -\infty} F(y) = 0$
2. "ends at one": $F(\infty) \equiv \lim_{y \rightarrow \infty} F(y) = 1$
3. "non-decreasing in between": $y_1 < y_2 \implies F(y_1) \leq F(y_2)$

Definition 2.21. A r.v. Y with cdf $F(y)$ is called **continuous** if $F(y)$ is continuous for all real y . i.e. $\forall y$,

$$\lim_{\varepsilon \rightarrow 0^+} F(y) = \lim_{\varepsilon \rightarrow 0^-} F(y) = F(y)$$

Warning 22. If $F(y)$ is continuous, then the probability attached to any individual point in the space is null: $P(Y = y) = 0 \forall y \in \mathbb{R}$.

2.12 Probability Density Function

Definition 2.22. For a continuous cdf $F(y)$, the **probability density function** (pdf) $f(y)$ is the first derivative of $F(\cdot)$ at y when it exists:

$$f(y) = \frac{dF(y)}{dt} \implies F(y) = \int_{-\infty}^y f(t)dt$$

Theorem 2.21. Properties of the pdf $f(y)$:

1. non-negative: $f(y) \geq 0 \quad -\infty < y < \infty$
2. integrates to 1: $\int_{-\infty}^{\infty} f(y)dy = 1$

Proof. Follows from probability axioms - properties from Theorem 2.20. \square

Warning 23. The probability density function does NOT represent probabilities; it does not integrate to 1.

We can compute the interval probability $P(a \leq Y \leq b)$ as

$$P(a \leq Y \leq b) = P(Y \leq b) - P(Y \leq a) = \int_a^b f(t)dt$$

ALL THE DEFINITIONS AND THEOREMS shown for the discrete case still hold for the continuous case:

1. EXPECTATION: for continuous r.v. Y and a function $g(Y)$,

$$\begin{aligned} \mathbb{E}[Y] &= \int_{-\infty}^{\infty} yf(y)dy = \int_y yf(y)dy = \mu \\ \mathbb{E}[g(Y)] &= \int_{-\infty}^{\infty} g(y)f(y)dy \quad (\text{both: only if integral is finite}) \end{aligned}$$

2. LINEARITY OF EXPECTATION (Theorem 2.2) still holds.

3. VARIANCE: for continuous r.v. Y ,

$$\mathbb{V}[Y] = \mathbb{E}[(Y - \mu)^2] = \int_{-\infty}^{\infty} (y - \mu)^2 f(y)dy$$

4. MOMENTS, CENTRAL MOMENTS AND FACTORIAL MOMENTS:

$$\mu'_k = \mathbb{E}[Y^k] = \int_{-\infty}^{\infty} y^k f(y) dy$$

$$\mu_k = \mathbb{E}[(Y - \mu)^k] = \int_{-\infty}^{\infty} (y - \mu)^k f(y) dy$$

$$\mu_{[k]} = \mathbb{E}[Y(Y-1)\cdots(Y-k+1)] = \int_{-\infty}^{\infty} y \cdots (y-k+1) f(y) dy$$

5. MOMENT GENERATING FUNCTION $m(t)$: for $|t| \leq b$, $b > 0$

$$m(t) = \mathbb{E}[e^{tY}] = \int_{-\infty}^{\infty} e^{ty} f(y) dy$$

6. FACTORIAL MOMENT GENERATING FUNCTION $G(t)$: for $1 - b \leq t \leq 1 + b$, $b > 0$,

$$G(t) = \mathbb{E}[t^Y] = \int_{-\infty}^{\infty} t^y f(y) dy$$

2.13 Distributions for continuous RVs

Definition 2.23. The **continuous uniform distribution** has a pdf that is constant on some interval (θ_1, θ_2) , $\theta_1 < \theta_2$:

$$f(y) = \frac{1}{\theta_2 - \theta_1} \quad \forall \theta_1 \leq y \leq \theta_2 \text{ and zero otherwise.}$$

By direct calculation, we can compute that the cdf $F(y)$ is

$$F(y) = \begin{cases} 0 & y \leq \theta_1 \\ \frac{y - \theta_1}{\theta_2 - \theta_1} & \theta_1 \leq y \leq \theta_2 \\ 1 & y \geq \theta_2 \end{cases}$$

and also find $\mathbb{E}[Y] = \frac{\theta_1 + \theta_2}{2}$ and $\mathbb{V}[Y] = \frac{(\theta_2 - \theta_1)^2}{12}$.

Definition 2.24. The **Normal (Gaussian) distribution** has pdf defined with parameters σ and μ by the pdf

$$f(y) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2\sigma^2}(y - \mu)^2\right) \quad \forall -\infty < y < \infty$$

It can be shown that $\mathbb{E}[Y] = \mu$ and $\mathbb{V}[Y] = \sigma^2$.

The cdf $F(y) = \int_{-\infty}^y \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2\sigma^2}(t - \mu)^2\right) dt$ which cannot be computed analytically.

Instead, if Y has Normal distribution with parameters μ and σ , and we consider the r.v. $Z = \frac{Y - \mu}{\sigma}$, then $\forall z \in \mathbb{R}$,

$$\begin{aligned} F_Z(z) &= P(Z \leq z) = P(\sigma Z + \mu \leq \sigma z + \mu) = P(Y \leq \sigma z + \mu) \\ &= F_Y(\sigma z + \mu) = \int_{-\infty}^{\sigma z + \mu} \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2\sigma^2}(t - \mu)^2\right) dt \\ &= \int_{-\infty}^z \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{u^2}{2}\right) du \quad \text{substituting } u = (t - \mu)/\sigma \end{aligned}$$

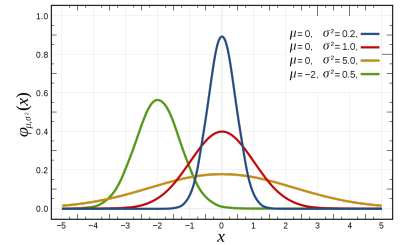


Figure 1: Normal distribution for various values of μ and σ .

Definition 2.25. The **Gamma distribution** with parameters $\alpha, \beta > 0$ has pdf

$$f(y) = \begin{cases} 0 & y < 0 \\ \frac{1}{\beta^\alpha \Gamma(\alpha)} y^{\alpha-1} e^{-y/\beta} & 0 \leq y < \infty \end{cases}$$

where $\Gamma(\alpha) = \int_0^\infty t^{\alpha-1} e^{-t} dt$ which satisfies $\Gamma(\alpha) = (\alpha - 1)\Gamma(\alpha - 1)$, and for $n \in \mathbb{N}$, $\Gamma(n) = (n - 1)!$ with $\Gamma(1) = 1$.

It can be shown that $\mathbb{E}[Y] = \alpha\beta$ and $\mathbb{V}[Y] = \alpha\beta^2$.

SPECIAL CASES of the Gamma Distribution are

1. CHI-SQUARE DISTRIBUTION μ degrees of freedom: $\alpha = \mu/2, \beta = 2$
2. EXPONENTIAL DISTRIBUTION: $\alpha = 1, f(y) = \frac{1}{\beta} e^{-y/\beta} \quad 0 \leq y < \infty$

The moment-generating function of the Gamma distribution is

$$m(t) = \mathbb{E}[e^{tY}] = \int_0^\infty e^{ty} \frac{1}{\beta^\alpha \Gamma(\alpha)} y^{\alpha-1} e^{-y/\beta} dy = \left(\frac{1}{1 - \beta t} \right)^\alpha$$

Definition 2.26. The **beta distribution** is suitable for an r.v. defined on $[0, 1]$, with pdf

$$f(y) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} y^{\alpha-1} (1 - y)^{\beta-1} = \frac{1}{B(\alpha, \beta)} y^{\alpha-1} (1 - y)^{\beta-1} \quad 0 \leq y \leq 1$$

where $B(\alpha, \beta)$ is the **Beta** function $B(\alpha, \beta) = \int_0^1 y^{\alpha-1} (1 - y)^{\beta-1} dy$.

It can be found that $\mathbb{E}[Y^k] = \frac{1}{B(\alpha, \beta)} B(\alpha + k, \beta)$ for $k = 1, 2, \dots$

Remark 24. A special case of the beta function with $\alpha = \beta = 1$ has $B(1, 1) = \frac{\Gamma(1)\Gamma(1)}{\Gamma(2)} = \frac{1 \cdot 1}{1} = 1$ and results in the pdf

$$f(y) = 1 \quad \forall 0 \leq y \leq 1 \quad \text{and zero otherwise}$$

the continuous uniform distribution with $\theta_1 = 0, \theta_2 = 1$.

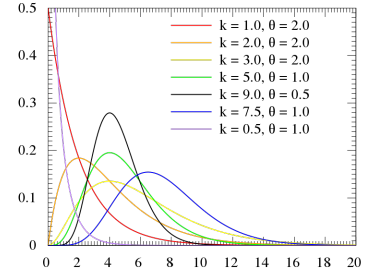


Figure 2: Gamma distribution for various values of α and β , which are respectively the shape and scale parameters.

The Gamma distribution is to the Exponential distribution what the Binomial is to the Bernoulli distribution.

e.g use case of the Gamma dist.: "How long does it take for a bucket to fill up if left out in the rain?"

The exponential distribution is known as the waiting-time distribution.

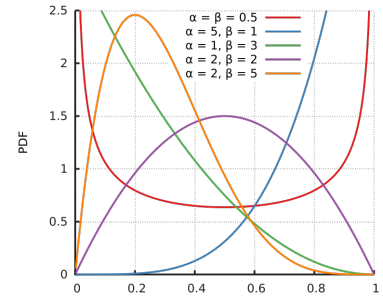


Figure 3: Beta distribution for various values of α and β .

3 Probability Calculation Methods

3.1 Transformations in One Dimension

Given r.v. Y defined on \mathcal{Y} (the region on which Y is nontrivial), we are interested in the transformed r.v. U (defined on \mathcal{U}), by

$$U = h(Y) \text{ where } h(\cdot) : \mathcal{Y} \rightarrow \mathbb{R}$$

DISCRETE CASE

Suppose $h(\cdot)$ maps values from (countable) $\{y_1, y_2, \dots, y_n, \dots\}$ to (countable) $\{u_1, u_2, \dots, u_n, \dots\}$ with $u_i = h(y_i)$.

- if $h(\cdot)$ is 1-1, then identify the set $\{u_i\}_{i=1}^\infty$ that $\{y_i\}_{i=1}^\infty$ map onto, and for each i , assign $p_Y(y_i)$ to u_i to obtain $p_U(u_i)$.
- if $h(\cdot)$ is not 1-1, we must identify the set \mathcal{U} and compute for each case the probability assigned to the values in the set.

CONTINUOUS CASE

We again want to see how values in the set \mathcal{Y} are mapped to a set \mathcal{U} . Consider $\mathcal{B} \subseteq \mathcal{U}$ the set of values that points in $\mathcal{A} \subseteq \mathcal{Y}$ can be mapped to under the transformation $h(\cdot)$:

$$\mathcal{B} = \{u : u = h(y) \mid y \in \mathcal{A}\} \implies P(Y \in \mathcal{A}) = P(U \in \mathcal{B})$$

So to compute $F_U(\cdot)$ or $f_U(\cdot)$, $\mathcal{B} \subseteq \mathcal{U}$, so fix \mathcal{B} . If $h(\cdot)$ is a 1-1 transformation, then the inverse transformation $h^{-1}(\cdot)$ can usually be computed explicitly. We know $P(U \in \mathcal{B}) = P(Y \in h^{-1}(\mathcal{B}))$, and

$$\textbf{Theorem 3.1.} \quad f_U(u) = f_Y(h^{-1}(u)) \left| \frac{dh^{-1}(u)}{du} \right|$$

Proof. Two cases depending on $h(\cdot)$

- 1-1 increasing:

$$\begin{aligned} F_U(u) &= P(U \leq u) = P(h(Y) \leq u) = P(Y \leq h^{-1}(u)) = F_Y(h^{-1}(u)) \\ f_U(u) &= \frac{dF_U(u)}{du} = \frac{dF_Y(h^{-1}(u))}{du} = f_Y(h^{-1}(u)) \underbrace{\frac{dh^{-1}(u)}{du}}_{\geq 0} \end{aligned}$$

- 1-1 decreasing:

$$\begin{aligned} F_U(u) &= P(U \leq u) = P(h(Y) \leq u) = P(Y \geq h^{-1}(u)) = 1 - F_Y(h^{-1}(u)) \\ f_U(u) &= \frac{dF_U(u)}{du} = -\frac{dF_Y(h^{-1}(u))}{du} = -f_Y(h^{-1}(u)) \underbrace{\frac{dh^{-1}(u)}{du}}_{\leq 0} \quad \square \end{aligned}$$

We typically deal with simple transformations, e.g.

$$\begin{aligned} U &= 2Y & U &= \exp(-Y) \\ U &= \frac{Y}{Y+1} & U &= Y^2 \end{aligned}$$

Recall that for a continuous r.v. Y ,

$$P(Y = y) = 0 \quad \forall y \in \mathbb{R}$$

The identity $P(U \in \mathcal{B}) = P(Y \in h^{-1}(\mathcal{B}))$ can be re-written as

$$\int_{\mathcal{B}} f_U(u) du = \int_{h^{-1}(\mathcal{B})} f_Y(y) dy$$

Example 3.1. If $U = cY + d$ for $c > 0$ and d ,

$$\begin{aligned} m_U(t) &= \mathbb{E}[e^{tU}] = \mathbb{E}[e^{t(cY+d)}] \\ &= \mathbb{E}[e^{(ct)Y} e^{dt}] = e^{dt} \mathbb{E}[e^{(ct)Y}] \\ \implies m_U(t) &= e^{dt} m_Y(ct) \end{aligned}$$

3.2 Techniques for sums of random variables

DISCRETE RANDOM VARIABLES

For discrete random variables Y_1 and Y_2 independent taking values on \mathcal{Y}_1 and \mathcal{Y}_2 respectively. Let $Y = Y_1 + Y_2$.

$$\begin{aligned}
 P(Y = y) &= P\left(\bigcup_{y_1 \in \mathcal{Y}_1} (Y = y) \cap (Y_1 = y_1)\right) \\
 &= \sum_{y_1 \in \mathcal{Y}_1} P((Y = y) \cap (Y_1 = y_1)) \\
 &= \sum_{y_1 \in \mathcal{Y}_1} P((Y_2 = y - y_1) \cap (Y_1 = y_1)) \\
 &= \sum_{y_1 \in \mathcal{Y}_1} P(Y_2 = y - y_1)P(Y_1 = y_1) \\
 &= \sum_{y_1 \in \mathcal{Y}_1^*} P(Y_2 = y - y_1)P(Y_1 = y_1)
 \end{aligned}$$

Theorem 3.2. If Y_1 and Y_2 are independent Poisson variables with parameters λ_1 and λ_2 , then $Y = Y_1 + Y_2$ also has a Poisson distribution, with parameter $\lambda_1 + \lambda_2$.

CONTINUOUS RANDOM VARIABLES

If Y_1 and Y_2 are independent continuous random variables, then the sum random variable $Y = Y_1 + Y_2$ is also continuous and we aim to compute $F_Y(y) = P(Y \leq y)$. By definition, with $f_Y(y)$ pdf of Y :

$$F_Y(y) = P(Y \leq y) = \int_{-\infty}^y f_Y(t) dt \text{ where } (Y \leq y) = (Y_1 + Y_2 \leq y)$$

From the discrete case above, we analogize to:

Definition 3.1. Convolution formula for the cdf of $Y = Y_1 + Y_2$:

$$f_Y(y) = \int_{-\infty}^{\infty} f_{Y_2}(y - y_1) f_{Y_1}(y_1) dy_1$$

where $f_Y(\cdot)$, $f_{Y_1}(\cdot)$ and $f_{Y_2}(\cdot)$ are the pdf's for Y , Y_1 and Y_2 resp.

Proof? .-

The important components of this calculation are:

- the partition of the event $Y = y$ according to possible values of Y_1
- the theorem of total probability
- the independence of Y_1 and Y_2

where

$$\mathcal{Y}_1^* = \{y_1 : y_1 \in \mathcal{Y}_1 \text{ and } y - y_1 \in \mathcal{Y}_2\}$$

Proof.

$$\begin{aligned}
 P(Y = y) &= \sum_{y_1=0}^{y-y_1} P(Y_1 = y_1)P(Y_2 = y - y_1) \\
 &= \sum_{y_1=0}^{y-y_1} \frac{\lambda_1^{y_1} e^{-\lambda_1}}{y_1!} \frac{\lambda_2^{y-y_1} e^{-\lambda_2}}{(y-y_1)!} \\
 &= e^{-(\lambda_1+\lambda_2)} \lambda_2^y \sum_{y_1=0}^{y-y_1} \left(\frac{\lambda_1}{\lambda_2}\right)^{y_1} \frac{1}{y_1!(y-y_1)!} \\
 &= \frac{e^{-(\lambda_1+\lambda_2)} \lambda_2^y}{y!} \left(1 + \frac{\lambda_1}{\lambda_2}\right)^y \\
 &= \frac{e^{-(\lambda_1+\lambda_2)} (\lambda_1 + \lambda_2)^y}{y!} \quad y = 0, 1, 2, \dots
 \end{aligned}$$

Poisson dist. with param. $\lambda_1 + \lambda_2$. \square

We can extend this result to more independent r.v.'s using the convolution formula recursively, e.g.

$$Y = Y_1 + Y_2 + Y_3 = (Y_1 + Y_2) + Y_3$$

4 Multivariate Distributions

We now consider probability specifications for more multiple random variables where there is not necessarily independence. We consider the vector random variable (Y_1, Y_2, \dots, Y_n) .

4.1 Discrete joint distributions

Definition 4.1. For Y_1, \dots, Y_n be discrete random variables on \mathbb{Z} . The **joint (discrete bivariate for $n = 2$) probability mass function** is

$$p(y_1, \dots, y_n) = P\left(\bigcap_{i=1}^n (Y_i = y_i)\right) \equiv P(Y_1 = y_1, \dots, Y_n = y_n)$$

We can think of $p(y_1, y_2)$ as specifying probabilities in a table with Y_1 and Y_2 values in columns / rows.

We can find the joint pmf for various situations using combinatorial arguments.

Theorem 4.1. The joint pmf has the basic properties:

$$0 \leq p(y_1, \dots, y_n) \leq 1 \quad \forall y_i \ i \in [1, n]$$

$$\sum_{i=1}^n \sum_{y_i=-\infty}^{\infty} p(y_1, \dots, y_n) = \sum_{y_1=-\infty}^{\infty} \cdots \sum_{y_n=-\infty}^{\infty} p(y_1, \dots, y_n) = 1$$

4.2 Marginal CDFs and PDFs

Definition 4.2. Given a joint pmf for Y_1 and Y_2 : $P_{Y_1, Y_2}(\cdot, \cdot)$, the **marginal probability mass function** for Y_i is merely the pmf P_{Y_i} of Y_i , which can be found via the Theorem of Total Probability (1.12).

$$p_{Y_1}(y_1) = P(Y_1 = y_1) = \sum_{y_2=-\infty}^{\infty} p_{Y_1, Y_2}(y_1, y_2)$$

If $p_{Y_1, Y_2}(y_1, y_2)$ specifies the values in a probability table, we compute the marginal pmf for Y_1 and Y_2 by summing respectively across the rows and down the columns of the table.

Definition 4.3. The **conditional mass function** for Y_2 given $Y_1 = y_1$ is

$$p_{Y_2|Y_1}(y_1, y_2) = P(Y_2 = y_2 | Y_1 = y_1) = \frac{P_{Y_1, Y_2}(y_1, y_2)}{P_{Y_1}(y_1)}$$

Recall the fundamental result

$$p_{Y_1, Y_2}(y_1, y_2) = p_{Y_1}(y_1)p_{Y_2|Y_1}(y_2 | y_1) \text{ when } p_{Y_1}(y_1) > 0$$

Definition 4.4. Discrete r.v.'s Y_1, \dots, Y_n are **independent** if

$$p_{Y_1, \dots, Y_n}(y_1, \dots, y_n) = \prod_{i=1}^n p_{Y_i}(y_i)$$

i.e. $p_{Y_i|Y_j}(y_i | y_j) = p_{Y_i} \ \forall i, j \in [1, n]$

| | Y_2 | | | $p_{Y_1}(\cdot)$ |
|------------------|-------|--------------|---|------------------|
| | 1 | 2 | 3 | |
| Y_1 | 1 | | | $p_{Y_1}(2)$ |
| | 2 | | | |
| | 3 | | | |
| $p_{Y_2}(\cdot)$ | | $p_{Y_2}(2)$ | | 1 |

Definition 4.5. The **joint cumulative distribution function cdf** $F_{Y_1, Y_2}(y_1, y_2)$ and **joint probability density function pdf** $f_{Y_1, Y_2}(y_1, y_2)$

$$F_{Y_1, Y_2}(y_1, y_2) = P((Y_1 \leq y_1) \cap (Y_2 \leq y_2)) \quad \forall (y_1 \in \mathbb{R}, y_2 \in \mathbb{R})$$

$$= \int_{-\infty}^{y_1} \int_{-\infty}^{y_2} f_{Y_1, Y_2}(t_1, t_2) dt_2 dt_1$$

Theorem 4.2. The joint cdf has the following properties:

$$\lim_{y_1 \rightarrow -\infty} \lim_{y_2 \rightarrow -\infty} F_{Y_1, Y_2}(y_1, y_2) = 0 \quad \lim_{y_1 \rightarrow \infty} \lim_{y_2 \rightarrow \infty} F_{Y_1, Y_2}(y_1, y_2) = 1$$

$$F_{Y_1, Y_2}(y_1, y_2) \leq F_{Y_1, Y_2}(y_1 + c, y_2) \text{ and } F_{Y_1, Y_2}(y_1, y_2 + c) \quad \forall y_1, y_2, c > 0$$

$$\lim_{y_1 \rightarrow \infty} F_{Y_1, Y_2}(y_1, y_2) = F_{Y_2}(y_2) \quad \lim_{y_2 \rightarrow \infty} F_{Y_1, Y_2}(y_1, y_2) = F_{Y_1}(y_1)$$

Theorem 4.3. By the probability axioms, the joint pdf

- is non-negative: $f_{Y_1, Y_2}(y_1, y_2) \geq 0$
- integrates to 1: $\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_{Y_1, Y_2}(y_1, y_2) dy_2 dy_1 = 1$

Definition 4.6. The **marginal probability density function pdf** $f_{Y_i}(y_i)$ for each fixed y_i , $i \in \{1, 2\}$ is

$$f_{Y_i}(y_i) = \int_{-\infty}^{\infty} f_{Y_1, Y_2}(y_1, y_2) dy_2$$

Definition 4.7. The **conditional probability density function pdf** is

$$f_{Y_2 | Y_1}(y_2 | y_1) = \frac{f_{Y_1, Y_2}(y_1, y_2)}{f_{Y_1}(y_1)}$$

As in the discrete case, we have the **CHAIN RULE FACTORIZATION**

$$f_{Y_1, Y_2}(y_1, y_2) = f_{Y_1}(y_1) f_{Y_2 | Y_1}(y_2 | y_1) = f_{Y_2}(y_2) f_{Y_1 | Y_2}(y_1 | y_2)$$

4.3 Independence and Correlation

Definition 4.8. Continuous r.v.'s (Y_1, \dots, Y_n) are **independent** if

$$f_{Y_1, \dots, Y_n}(y_1, \dots, y_n) = \prod_{i=1}^n f_{Y_i}(y_i)$$

Remark 25. An easy way to prove that two r.v.'s are not independent is to check if $\mathcal{Y}_{12} = \mathcal{Y}_1 \times \mathcal{Y}_2$ where \mathcal{Y}_{12} is the support of the joint pdf.

To compute $f_{Y_1, Y_2}(y_1, y_2)$ from $F_{Y_1, Y_2}(y_1, y_2)$, use

$$f_{Y_1, Y_2}(y_1, y_2) = \frac{\partial^2}{\partial y_1 \partial y_2} (F_{Y_1, Y_2}(y_1, y_2))$$

where the partial differentiations are commutative given the second derivatives are continuous at that point.

The conditional pdf $f_{Y_1, Y_2}(y_1, y_2)$ is a pdf in y_2 for every fixed y_1 , so it is non-negative and integrates to 1 over y_2 .

In the bivariate case, we have

$$f_{Y_2 | Y_1}(y_2 | y_1) = f_{Y_2}(y_2)$$

$$f_{Y_1 | Y_2}(y_1 | y_2) = f_{Y_1}(y_1)$$

Extending the idea of expectation to multiple variables Y_1, \dots, Y_n , we get the following, respectively discrete and continuous, expressions

$$\begin{aligned}\mathbb{E}[g(y_1, \dots, y_n)] &= \sum_{i=-\infty}^{\infty} \cdots \sum_{i=-\infty}^{\infty} g(y_1, \dots, y_n) f_{Y_1, \dots, Y_n}(y_1, \dots, y_n) \\ \mathbb{E}[g(y_1, \dots, y_n)] &= \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} g(y_1, \dots, y_n) f_{Y_1, \dots, Y_n}(y_1, \dots, y_n) dy_1 \cdots dy_n\end{aligned}$$

and linearity of expectation (Theorem 2.2) still holds.

Definition 4.9. The **covariance** $\text{Cov}[Y_1, Y_2]$ between Y_1 and Y_2 with expectations μ_1 and μ_2 is

$$\text{Cov}[Y_1, Y_2] = \mathbb{E}[(Y_1 - \mu_1)(Y_2 - \mu_2)] = \mathbb{E}[Y_1 Y_2] - \mu_1 \mu_2$$

Definition 4.10. The **correlation** $\text{Corr}[Y_1, Y_2]$ between Y_1 and Y_2 is

$$\text{Corr}[Y_1, Y_2] = \frac{\text{Cov}[Y_1, Y_2]}{\sqrt{\mathbb{V}[Y_1] \mathbb{V}[Y_2]}}$$

Theorem 4.4. $-1 \leq \text{Corr}[Y_1, Y_2] \leq 1$

Definition 4.11. Y_1, Y_2 **uncorrelated** if $\text{Cov}[Y_1, Y_2] = \text{Corr}[Y_1, Y_2] = 0$.

Theorem 4.5. For Y_1 and Y_2 independent,

$$\mathbb{E}[g_1(Y_1)g_2(Y_2)] = \mathbb{E}[g_1(Y_1)]\mathbb{E}[g_2(Y_2)]$$

Theorem 4.6. If Y_1 and Y_2 independent then they are uncorrelated.

Warning 26. Uncorrelation doesn't imply independence!

For multiple variables, we consider pairwise covariance between all of them and put the results into a matrix called covariance matrix.

Proof. in textbook. \square

Proof. (Theorem 4.5)

$$\begin{aligned}\mathbb{E}[g_1(Y_1)g_2(Y_2)] &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g_1(y_1)g_2(y_2)f_{Y_1, Y_2}(y_1, y_2)dy_1dy_2 \\ &= \int_{-\infty}^{\infty} g_1(y_1)f_{Y_1}(y_1)dy_1 \int_{-\infty}^{\infty} g_2(y_2)f_{Y_2}(y_2)dy_2 \\ &= \mathbb{E}[g_1(Y_1)]\mathbb{E}[g_2(Y_2)]\end{aligned}\quad \square$$

Proof. (Theorem 4.6)

$$\begin{aligned}\text{Cov}[Y_1, Y_2] &= \mathbb{E}[Y_1 Y_2] - \mu_1 \mu_2 \\ &= \mathbb{E}[Y_1]\mathbb{E}[Y_2] - \mu_1 \mu_2 \\ &= 0\end{aligned}\quad \square$$

4.4 Transformations of random variables

LINEAR COMBINATIONS

Let Y_1, \dots, Y_n with $\mathbb{E}[Y_i] = \mu_i$ and X_1, \dots, X_m with $\mathbb{E}[X_i] = \xi_i$, define $U_1 = \sum_{i=1}^n a_i Y_i$ and $U_2 = \sum_{j=1}^m b_j X_j$. We get the following:

Theorem 4.7.

$$\begin{aligned}\mathbb{E}[U_1] &= \sum_{i=1}^n a_i \mathbb{E}[Y_i] = \sum_{i=1}^n a_i \mu_i \\ \mathbb{V}[U_1] &= \sum_{i=1}^n a_i^2 \mathbb{V}[Y_i] + 2 \sum_{j=2}^n \sum_{i=1}^{j-1} a_i a_j \text{Cov}[Y_i, Y_j] \\ \text{Cov}[U_1, U_2] &= \sum_{i=1}^n \sum_{j=1}^m a_i b_j \text{Cov}[Y_i, X_j]\end{aligned}$$

Back to sums, consider $Y = Y_1 + Y_2$ and recall the convolution formula (Definition 3.1). An easier method than computing convolutions is to recall the moment generating function for $U = g(Y)$ (Definition 2.18) $m_U(t) = \mathbb{E}[e^{tU}] = \mathbb{E}[e^{tg(Y)}]$. If Y_1 and Y_2 are independent and $Y = a_1g_1(Y_1) + a_2g_2(Y_2)$:

$$m_Y(t) = \mathbb{E}[e^{t(a_1g_1(Y_1) + a_2g_2(Y_2))}] = \mathbb{E}[e^{ta_1g_1(Y_1)}] \mathbb{E}[e^{ta_2g_2(Y_2)}]$$

Theorem 4.8. For $Y = Y_1 + Y_2$ where Y_1, Y_2 independent,

$$m_Y(t) = m_{Y_1}(t)m_{Y_2}(t)$$

Proof. $m_Y(t) = \mathbb{E}[e^{t(Y_1+Y_2)}] = \mathbb{E}[e^{tY_1}] \mathbb{E}[e^{tY_2}] = m_{Y_1}(t)m_{Y_2}(t)$ \square

GENERAL TRANSFORMATIONS

Approaching from a 'change of variable' perspective:

Theorem 4.9. Consider the bijective transformation

$$(Y_1, Y_2) \longrightarrow (U_1, U_2)$$

where we can write

$$Y_1 = h_1^{-1}(U_1, U_2) \text{ and } Y_2 = h_2^{-1}(U_1, U_2)$$

So for $(u_1, u_2) \in \mathbb{R}^2$,

$$f_{U_1, U_2}(u_1, u_2) = f_{Y_1, Y_2}(h_1^{-1}(u_1, u_2), h_2^{-1}(u_1, u_2)) |J(u_1, u_2)|$$

$$\text{where } J(u_1, u_2) = \det \begin{bmatrix} \frac{\partial h_1^{-1}(u_1, u_2)}{\partial u_1} & \frac{\partial h_1^{-1}(u_1, u_2)}{\partial u_2} \\ \frac{\partial h_2^{-1}(u_1, u_2)}{\partial u_1} & \frac{\partial h_2^{-1}(u_1, u_2)}{\partial u_2} \end{bmatrix} \text{ is the Jacobian.}$$

Proof. c.f. textbook. \square

Example 4.1 (Sum of Binomial r.v.'s). For Y_1, Y_2 independent Binomial r.v.'s with parameters (n_1, p) and (n_2, p) we have for $Y = Y_1 + Y_2$

$$\begin{aligned} m_Y(t) &= m_{Y_1}(t)m_{Y_2}(t) \\ &= (pe^t + q)^{n_1}(pe^t + q)^{n_2} \\ &= (pe^t + q)^{n_1+n_2} \end{aligned}$$

Notice this is the mgf of the Binomial distribution $(n_1 + n_2, p)$.

Example 4.2 (Sum of Poisson r.v.'s). Similarly to the previous example, for Y_1 and Y_2 independent Poisson r.v.'s with parameters λ_1 and λ_2 to find

$$m_Y(t) = \exp((\lambda_1 + \lambda_2)(e^t - 1))$$

the mgf of the Poisson distribution with parameter $\lambda_1 + \lambda_2$.

Example 4.3 (Sum of Normal r.v.'s). Similarly for Y_1 and Y_2 independent Normal r.v.'s with parameters (μ_1, σ_1^2) and (μ_2, σ_2^2) , we get

$$m_Y(t) = \exp((\mu_1 + \mu_2)t + (\sigma_1^2 + \sigma_2^2)t^2/2)$$

the mgf of the Normal distribution with parameters $(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$.

5 Probability Inequalities and Theorems

Theorem 5.1 (Markov's Inequality). Suppose $g(y)$ function of a continuous r.v. Y such that $\mathbb{E}[|g(Y)|] < \infty$. For any $k > 0$ let

$$\mathcal{A}_k = \{y \in \mathbb{R} : |g(y)| \geq k\}$$

then we have

$$P(|g(Y)| \geq k) \leq \frac{\mathbb{E}[|g(Y)|]}{k}$$

Proof. We want to prove that $\mathbb{E}[|g(Y)|] \geq kP(|g(Y)| \geq k)$

$$\begin{aligned} \mathbb{E}[|g(Y)|] &= \int_{-\infty}^{\infty} |g(y)| f_Y(y) dy = \int_{\mathcal{A}_k} |g(y)| f_Y(y) dy + \int_{\mathcal{A}_k^c} |g(y)| f_Y(y) dy \\ &\geq \int_{\mathcal{A}_k} |g(y)| f_Y(y) dy \geq \int_{\mathcal{A}_k} k f_Y(y) dy = kP(|g(Y)| \geq k) \quad \square \end{aligned}$$

Theorem 5.2 (Chebychev's Inequality). For any distribution with finite mean μ and variance σ^2 , we have that $\forall c > 0$,

$$P(\mu - c\sigma < Y < \mu + c\sigma) \geq 1 - \frac{1}{c^2}$$

Proof. Let $g(t) = (t - \mu)^2$ in Theorem 5.1:

$$P((t - \mu)^2 < k) \geq 1 - \frac{\mathbb{E}[(t - \mu)^2]}{k}$$

Recall $\mathbb{E}[(Y - \mu)^2] = \sigma^2$ so choose $k = c^2\sigma^2$

$$P((t - \mu)^2 < c^2\sigma^2) = P\left(\sqrt{(t - \mu)^2} < \sqrt{c^2\sigma^2}\right) = P(|Y - \mu| < c\sigma)$$

$$P(|Y - \mu| < c\sigma) \geq 1 - \frac{\mathbb{E}[(t - \mu)^2]}{c^2\sigma^2} = 1 - \frac{1}{c^2} \quad \square$$

Theorem 5.3 (Sample mean). Suppose Y_1, \dots, Y_n independent r.v.'s that have the same distribution, with expectation μ and variance σ^2 . Let $\bar{Y}_n = \frac{1}{n} \sum_{i=1}^n Y_i$ be the sample mean r.v. Then $\forall \varepsilon > 0$,

$$P(|\bar{Y}_n - \mu| \geq \varepsilon) \leq \frac{\sigma^2}{n\varepsilon^2}$$

Definition 5.1. A sequence of r.v.'s $Y_1, Y_2, \dots, Y_n, \dots$ converges in probability to $c \in \mathbb{R}$ as $n \rightarrow \infty$, written as $Y_n \xrightarrow{p} c$, if $\forall \varepsilon > 0$,

$$\lim_{n \rightarrow \infty} P(|Y_n - c| < \varepsilon) = 1$$

Markov's Inequality holds for discrete r.v.'s as well (just replace all the integrals with sums in the proof).

It can be written equivalently as

$$P(|g(Y)| < k) \geq 1 - \frac{\mathbb{E}[|g(Y)|]}{k}$$

Theorem 5.2 bounds how much probability can be present in the tails of the distribution.

It can be written equivalently as

$$P(|Y - \mu| \geq c\sigma) \leq \frac{1}{c^2}$$

Proof. Recall that for $U = \sum_{i=1}^n Y_i$ we have $\mathbb{E}[U] = n\mu$ and $\mathbb{V}[U] = n\sigma^2$.

$$\mathbb{E}[\bar{Y}_n] = \frac{1}{n} \mathbb{E}\left[\sum_{i=1}^n Y_i\right] = \frac{n\mu}{n} = \mu$$

$$\mathbb{V}[\bar{Y}_n] = \frac{1}{n^2} \mathbb{V}\left[\sum_{i=1}^n Y_i\right] = \frac{n\sigma^2}{n^2} = \frac{\sigma^2}{n}$$

$$P\left(|\bar{Y}_n - \mu| \geq \frac{c\sigma}{\sqrt{n}}\right) \leq \frac{1}{c^2} \quad \forall c$$

Let $\varepsilon = \frac{c\sigma}{\sqrt{n}}$ such that $c = \frac{\sqrt{n}\varepsilon}{\sigma}$

$$\implies P(|\bar{Y}_n - \mu| \geq \varepsilon) \leq \frac{\sigma^2}{n\varepsilon^2}$$

where ε arbitrary. \square