

Famnit

Probability

The written Corona lectures



Published in the year of our Lord 2021,
when we were all suffering from the Plague
Michaelus Permanus

The Corona lecture notes

PROBABILITY

Mihael Perman

Note:

The following are the Courvoisier lecture notes. They cover most of the material we covered by Zoom except Chapter 2.2 on continuous distributions.

1. Outcomes, events, probabilities

1.1. Outcomes and events

Example: Italian gamblers in the 17th century liked the game where they placed a bet on the outcome of rolling three dice.

○ Popular bets were 9 and 10.

The gamblers had a "theory" that the two popular bets are equivalent in the sense that the probability of winning is the same for both bets. They wrote down two

○ lists :

Sum 9

1 2 6

1 3 5

1 4 4

2 2 5

2 3 4

3 3 3

Sum 10

1 3 6

1 4 5

2 2 6

2 3 5

2 4 4

3 3 4

Based on these two lists the two games were deemed equivalent.

However, gambling experience suggested they were not.

The problem was solved by Galileo Galilei (1564-1642).

He wrote down all possible outcomes.

111	112	113	114	115	116
121	122	123	124	125	126
⋮	⋮	⋮	⋮	⋮	⋮
661	662	663	664	665	666

There are $6^3 = 216$ possible

triplets. Galileo found that

25 sum to 9 and 27 to 10.

Assuming all triplets have the same chance of appearing the problem is solved.

The moral of the story is that we have to write down all possibilities when dealing with an experiment involving chance.

In mathematical language we will talk about the set of all possible outcomes and denote it by Ω .

Examples:

(i) In Galileo's example we have

$$\begin{aligned}\Omega &= \{(i, j, k) : 1 \leq i, j, k \leq 6\} \\ &= \{1, 2, 3, \dots, 6\}^3\end{aligned}$$

(ii) If we toss a coin n times we get a sequence of n heads and tails. In this case

$$\Omega = \{H, T\}^n.$$

(iii) Suppose we arrange n objects in random order. This means that we choose a random permutation or

$$\Omega = S_n = \text{set of all permutations.}$$

(iv) We can think of tossing a coin infinitely many times. In this case

$$\Omega = \{H, T\}^{\mathbb{N}}$$

which is the set of all countably infinite sequences of symbols

H and T .

The next concept is the event.

If we roll three dice an event is, say, that the sum is 9.

An event can either happen or not. But what is an event

mathematically? All the triplets that give a sum of 9 are a subset of $\Omega = \{1, 2, 3, 4, 5, 6\}^3$.

It is plausible to understand events as subsets of Ω .

We will denote events by

A, B, C, \dots

For mathematical reasons denote the family of all events by \mathcal{F} .

We will require the following.

(i) $\Omega \in \mathcal{F}$.

(ii) if $A \in \mathcal{F}$ then $A^c \in \mathcal{F}$.

(iii) if $A_1, A_2, \dots \in \mathcal{F}$, then

$$\bigcup_i A_i \in \mathcal{F}.$$

The union is either finite or countable.

Remark: In mathematics a family of subsets with the above properties is called a σ -algebra.

Remark: In cases of infinite sets Ω not all subsets are necessarily events. For finite Ω we will usually assume that all subsets are events.

In Galileo's example we assumed that all outcomes in Ω are equally likely. The probability of

$A = \{ \text{sum is } 9 \}$ is then $25/216$.

If $B = \{ \text{sum is } 10 \}$ then $P(B) = 27/216$.

We have $A \cap B = \emptyset$ and

$$P(A \cup B) = \frac{25 + 27}{216} = P(A) + P(B).$$

For mathematical reasons it turns out to be better to assign probabilities to events rather than outcomes. The example shows that for disjoint A and B we should have $P(A \cup B) = P(A) + P(B)$.

Any assignment of probabilities should have this property. The mathematical definition is more general.

Definition: Probability is an assignment to every event $A \in \mathcal{F}$ of a real number in such a way that

(i) $0 \leq P(A) \leq 1$, $P(\Omega) = 1$.

(ii) if A_1, A_2, \dots are disjoint we have

$$P\left(\bigcup_i A_i\right) = \sum_i P(A_i)$$

Remark: The sum in (ii) can be finite or infinite.

Remark: The property (ii) is called σ -additivity.

Let us look at some simple consequences of the above definition.

(i) we have $A \cup A^c = \Omega$ and $A \cap A^c = \emptyset$. By additivity

$$1 = P(\Omega) = P(A \cup A^c) = P(A) + P(A^c)$$

so $P(A^c) = 1 - P(A)$

(ii) Let A, B be events. We can write

$$A \cup B = \underbrace{(A \cap B^c) \cup (A \cap B) \cup (A^c \cap B)}_{\text{disjoint}}$$

So

$$P(A \cup B) = P(A \cap B^c) + P(A \cap B) + P(A^c \cap B)$$

But
$$\underbrace{(A \cap B^c) \cup (A \cap B)}_{\text{disjoint}} = A \quad \text{so}$$

$$P(A \cap B^c) + P(A \cap B) = P(A) \Rightarrow$$

$$P(A \cap B^c) = P(A) - P(A \cap B)$$

and similarly

$$P(A^c \cap B) = P(B) - P(A \cap B)$$

Using this in the above expression gives

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

If we have three events we get

$$P(A \cup B \cup C) = P((A \cup B) \cup C)$$

$$= P(A \cup B) + P(C)$$

$$- P((A \cup B) \cap C)$$

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

$$- P((A \cap C) \cup (B \cap C))$$

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

$$- P(A \cap C) - P(B \cap C)$$

$$+ P(A \cap B \cap C)$$

From this we generalize to

Theorem 4.1 (inclusion-exclusion formula)

Let A_1, A_2, \dots, A_n be events.

We have

$$P\left(\bigcup_{i=1}^n A_i\right) = \sum_{i=1}^n P(A_i) - \sum_{1 \leq i < j \leq n} P(A_i \cap A_j)$$

$$+ \sum_{1 \leq i < j < k \leq n} P(A_i \cap A_j \cap A_k)$$

- . . .

$$+ (-1)^{n-1} P(A_1 \cap A_2 \cap \dots \cap A_n)$$

Proof: We know that the formula is valid for $n = 2$. Suppose it is valid for n . We write

$$\bigcup_{i=1}^{n+1} A_i = \bigcup_{i=1}^n A_i \cup A_{n+1} \quad \text{so}$$

$$\begin{aligned} P\left(\bigcup_{i=1}^{n+1} A_i\right) &= P\left(\bigcup_{i=1}^n A_i\right) + P(A_{n+1}) \\ &\quad - \underbrace{P\left(\bigcup_{i=1}^n A_i \cap A_{n+1}\right)} \\ &= P\left(\bigcup_{i=1}^n (A_i \cap A_{n+1})\right). \end{aligned}$$

By the induction assumption the formula is valid for unions of n sets. This means

$$\begin{aligned} P\left(\bigcup_{i=1}^n (A_i \cap A_{n+1})\right) &= \sum_{i=1}^n P(A_i \cap A_{n+1}) \\ &\quad - \sum_{1 \leq i < j \leq n} P(A_i \cap A_j \cap A_{n+1}) \\ &\quad + \dots + (-1)^{n-1} P(A_1 \cap \dots \cap A_{n+1}) \end{aligned}$$

Using this gives the inclusion
exclusion formula for $n+1$ sets,
and the induction step is
completed.

Example: n couples go dancing.
When they are about to leave
the power goes out and each
woman grabs a man at random.
What is the probability that no
woman will grab her man?

In probability language we are
talking about choosing a
random permutation of n
numbers. All permutations
have the same probability $\frac{1}{n!}$.

Figure:

Women	1	2	3	...	i	...	n
Men	3	5	4	...	$2(i)$...	1

Define $A_i = \{ \text{woman } i \text{ grabs her man} \}$.

and $A = \{ \text{no woman grabs her man} \}$

We have

$$A^c = \bigcup_{i=1}^n A_i$$

To use the exclusion-exclusion formula we need the following

probabilities:

$$P(A_i) = \frac{(n-1)!}{n!} = \frac{1}{n}$$

of permutations
sending $i \rightarrow i$.

↑
of all probabilities

$$P(A_i \cap A_j) = \frac{(n-2)!}{n!} = \frac{1}{n(n-1)}$$

⋮

$$P(A_{i_1} \cap \dots \cap A_{i_r}) = \frac{(n-r)!}{n!}$$

The inclusion - exclusion formula gives

$$\begin{aligned} P(A^c) &= \binom{n}{1} \cdot \frac{1}{n} - \binom{n}{2} \cdot \frac{(n-2)!}{n!} \\ &\quad + \binom{n}{3} \frac{(n-3)!}{n!} - \\ &\quad \vdots \\ &\quad + (-1)^n \frac{0!}{n!} \end{aligned}$$

$$= 1 - \frac{1}{2!} + \frac{1}{3!} - \dots + (-1)^{n-1} \cdot \frac{1}{n!}$$

The symbol $\binom{n}{r}$ counts the number of different intersections.

Finally

$$\begin{aligned} P(A) &= 1 - P(A^c) \\ &= \frac{1}{2!} - \frac{1}{3!} + \dots + (-1)^n \cdot \frac{1}{n!} \end{aligned}$$

From Analysis we know

$$e^x = \sum_{k=0}^{\infty} \frac{x^k}{k!} = 1 + \frac{x}{1!} + \frac{x^2}{2!} + \dots$$

Take $x = -1$ to get

$$e^{-1} = \frac{1}{2!} - \frac{1}{3!} + \frac{1}{4!} - \dots$$

- The probability $P(A)$ is a partial sum of the above series which converges fast. We can approximate

$$P(A) \approx e^{-1} = 0.3697$$

○

1.2. Conditional probabilities and independence

Example: Let us return to Galileo's example. We have $\Omega = \{1, 2, 3, 4, 5, 6\}^3$ and all triplets are equally likely. Suppose you know that the first component is 1 but not the other two components. What is your opinion about the probability that the sum is 9? There are 36 triplets of the form $(1, j, k)$.

Of these the triplets

$(1, 2, 6)$

$(1, 3, 5)$

$(1, 4, 4)$

$(1, 5, 3)$

$(1, 6, 2)$

give a sum of 9.

Is it reasonable to assume that given the information that the first component is 1 all the 36 triplets are equally likely?

Yes. So the updated probability

is $5/36$ which is different

from $25/216$. We rewrite

$$5/36 = \frac{5/216}{36/216}$$

and denote $A = \{\text{sum is } 9\}$ and

$B = \{\text{first component is } 1\}$.

We have

$$5/36 = \frac{P(A \cap B)}{P(B)}$$

Definition: The conditional probability of A given B is

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

Remark: If we have additional information about an outcome this usually means that the outcome is in a restricted subset of Ω . In the above example this restricted subset is B.

○ Rewriting the definition we get

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

(if A_1, A_2, \dots, A_n are events we can write

$$P(A_1 \cap \dots \cap A_{n-1} \cap A_n) =$$

$$P(A_n | A_1 \cap \dots \cap A_{n-1}) P(A_1 \cap \dots \cap A_{n-1})$$

Iterating the rule gives

$$P(A_1 \cap \dots \cap A_n)$$

$$= P(A_n | A_1 \cap \dots \cap A_{n-1}) \cdot$$

$$P(A_{n-1} | A_1 \cap \dots \cap A_{n-2})$$

\vdots

$$P(A_2 | A_1) \cdot P(A_1)$$

Definition: A collection $\{H_1, H_2, \dots, H_n\}$ is a partition of Ω if $H_i \cap H_j = \emptyset$ for all $i \neq j$ and $\bigcup_{i=1}^n H_i = \Omega$.

Theorem 1.2 (law of total probabilities)

Let $\{H_1, H_2, \dots, H_n\}$ be a partition and A an event. We have

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

Proof: We write

$$A = A \cap \Omega$$

$$= A \cap \bigcup_{i=1}^n H_i$$

$$= \bigcup_{i=1}^n (A \cap H_i)$$

disjoint events.

It follows that

$$P(A) = \sum_{i=1}^n P(A \cap H_i)$$

$$= \sum_{i=1}^n \frac{P(A \cap H_i)}{P(H_i)} \cdot P(H_i)$$

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

Remark: This formula is useful to compute probabilities. We can often guess conditional probabilities.

Example: In an internet game of chance you have 12 tickets.

1 1 1 1 2 2 3 0 5 5 5 5

The tickets are randomly permuted and turned around so that the player sees

0 0 0 0 0 0 0 0 0 0 0 0

The player then turns around the tickets until the first

0 = STOP. One example is

1 2 0 1 0

The payoff is the sum of digits multiplied by 2 if 0 = DOUBLE is among the visible tickets.

In the above example the payoff is 8.

What is the probability that the player will see the ticket \boxed{D} ?

Define $H_i = \{ \text{first } \boxed{S} \text{ is in position } i \}$
for $i = 1, 2, \dots, 9$. The collection
 $\{H_1, H_2, \dots, H_9\}$ is a partition.

Let $A = \{ \text{we see } \boxed{D} \}$.

First we compute

$$P(H_i) = \frac{8}{12} \cdot \frac{7}{11} \cdot \frac{8-i+2}{12-i+2} \cdot \frac{4}{12-i+1}$$

What about $P(A|H_i)$? (idea?)

If \boxed{S} appears in position i
then the first $i-1$ tickets are
randomly chosen from
 $\boxed{A} \boxed{1} \boxed{1} \boxed{1} \boxed{2} \boxed{2} \boxed{3} \boxed{D}$. So we
choose $i-1$ tickets out of 8

and ask for the (conditional) probability that \textcircled{D} is among the tickets chosen. We have

$$P(A | H_i) = \frac{\binom{7}{i-2}}{\binom{8}{i-1}} \leftarrow \begin{array}{l} \# \text{ of samples} \\ \text{of size } i-1 \\ \text{containing } \textcircled{D}. \end{array}$$

↑
of possible samples

Cancelling we get

$$P(A | H_i) = \frac{i-2}{i-1} = \frac{i-1}{8}$$

Check: For $i=9$ we should get 1.

The rest is adding fractions.

$$\begin{aligned} P(A) &= \sum_{i=1}^9 \frac{8! (12-i)!}{12! (8-i+1)!} \cdot 4 \cdot \frac{(i-1)}{8} \\ &= \frac{8!}{12!} \cdot \frac{1}{2} \sum_{i=1}^9 \frac{(12-i)! (i-1)}{(8-i+1)!} \\ &= \frac{1}{5} \end{aligned}$$

Example : Prisoner's Paradox

Three prisoners are in jail in a dark country. They are all sentenced to death but the ruler will choose one of them at random and pardon him. Here is a conversation between the guard in jail and prisoner A :

A: Guard, you already know who will be pardoned. If you tell me who of the other two will not be pardoned you do not give me any information.

G: If I tell you there will be only two of you left. Your probability of survival is then $\frac{1}{2}$. I do give you some information.

Who is right? To talk about conditional probabilities we need a space of all possible outcomes. Here is a suggestion:

$$\Omega = \left\{ \begin{array}{l} \boxed{AB|B} \quad \frac{1}{3} \\ \boxed{AC|C} \quad \frac{1}{3} \\ \boxed{BC|B} \\ \boxed{BC|C} \end{array} \right\} \quad \frac{1}{3} \rightarrow \begin{array}{l} x/3 \\ (1-x)/3 \end{array}$$

↑

Last letter is what the guard says

First two letters are the wretched prisoners who will be hanged

There is no indication how the last probability of $\frac{1}{3}$ is distributed between the last two outcomes. Let us say

$$\frac{x}{3} \text{ and } \frac{1-x}{3} \text{ for } x \in [0, 1].$$

We compute

$$P(A \text{ survives} \mid \text{Guard says } B)$$

$$= \frac{P((A \text{ survives}) \cap \{\text{Guard says } B\})}{P(\text{Guard says } B)}$$

$$= \frac{x/3}{1/3 + x/3}$$

$$= \frac{x}{1+x}$$

This function has values from 0 to $1/2$ on $[0, 1]$.

Two cases:

(i) if $x = 1/2$ the guard chooses at random when he has the choice. In this case the conditional probability is

$$\frac{1/2}{1 + 1/2} = 1/3 \text{ as before.}$$

(ii) If $x = 0$ then the
conditional probability is 0!

Why?

What if $P(A|B) = P(A)$? Then B does not "tell us anything" about the probability of A. The word we choose is independence. The above equality can be written as

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

by definition. If we have events A, B and C and they are "independent" then $A \cap B$ ought to be independent of C. This leads to the following definition.

Definition:

(i) The events A and B are independent if

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

(ii) Events A_1, A_2, \dots, A_n are independent if for all collections of indices $1 \leq i_1 < i_2 < \dots < i_m \leq n$ we have

$$P(A_{i_1} \cap \dots \cap A_{i_m})$$

$$= P(A_{i_1}) P(A_{i_2}) \dots P(A_{i_m}).$$

Remark: Typically independence is associated with physically different objects, like several

dice, different coins.

Example (Paradox of Chevalier de Méré,
Antoine Gombaud, 1607 - 1684).

Chevalier de Méré considered the
following two games of chance.

(i) You roll a die 4 times. You
win if you see at least one
ace (ace = $\square \cdot$)

(ii) You roll two dice 24 times.
You win if you see at least
one double ace, i.e. $\square \square$

Which of the two games has a
higher probability of winning?

Let us look at the first game.

Define $A_i = \{ \text{the } i\text{-th roll is not } \square \cdot \}$

and $A = \{ \text{we win} \}$.

We have

$$A^c = A_1 \cap A_2 \cap A_3 \cap A_4$$

It is reasonable to assume that subsequent rolls are independent which means that A_1, A_2, A_3, A_4 are independent. It follows

$$\begin{aligned} P(A^c) &= P(A_1 \cap A_2 \cap A_3 \cap A_4) \\ &= P(A_1) \cdot P(A_2) \cdot P(A_3) \cdot P(A_4) \\ &= \frac{5}{6} \cdot \frac{5}{6} \cdot \frac{5}{6} \cdot \frac{5}{6} \\ &= \left(\frac{5}{6}\right)^4 \end{aligned}$$

Finally we get

$$P(A) = 1 - P(A^c) = 1 - \left(\frac{5}{6}\right)^4 \doteq 0.5177$$

For the second game we similarly define

$$A_i = \{ \text{not a double ace on roll } i \},$$
$$i = 1, 2, \dots, 24.$$

$$A = \{ \text{we win} \}$$

We have

$$A^c = A_1 \cap A_2 \cap \dots \cap A_{24}$$

We assume independence and get

$$P(A^c) = P(A_1) P(A_2) \dots P(A_{24})$$

$$= \frac{35}{36} \cdot \frac{35}{36} \cdot \dots \cdot \frac{35}{36}$$

$$= \left(\frac{35}{36} \right)^{24},$$

and finally

$$P(A) = 1 - P(A^c) = 1 - \left(\frac{35}{36} \right)^{24} = 0.4914.$$

Comment: The difference is small but nevertheless important.

Example (Gambler's ruin).

Two gamblers A and B start out with m and n sequins (gold coins) respectively. In each round of the game they toss a coin. If it is heads A gets a coin from B; if it is tails B gets a coin from A.

They play until one of them is left with no coins. What is the probability that A will get all the coins? We assume that tosses are independent and the probability of heads is $p \in (0, 1)$.

Let $A = \{ \text{game over } A \text{ wins} \}$ and denote $p_{m,n} = P(A \text{ wins})$.

Let $H = \{ \text{first toss is heads} \}$.

Then

$$\begin{aligned} P(A) &= \underbrace{P(A|H)} \cdot P(H) + \underbrace{P(A|H^c)} P(H^c) \\ &= p_{m+1, n-1} = p_{m-1, n+1} \end{aligned}$$

We get the recursion

$$p_{m,n} = p \cdot p_{m+1, n-1} + (1-p) p_{m-1, n+1}.$$

Obviously $p_{0, m+n} = 0$ and $p_{m+n, 0} = 1$.

Denote

$$\bar{x}_m = p_{m,n}.$$

We can rewrite the recursion as

$$\bar{x}_m = p \cdot \bar{x}_{m+1} + (1-p) \bar{x}_{m-1}.$$

with $\bar{x}_{m+n} = 1$ and $\bar{x}_0 = 0$.

Define $z := 1-p$. We rewrite

$$\underbrace{(p+z)}_{=1} \bar{x}_m = p \bar{x}_{m+1} + z \bar{x}_{m-1}, \text{ or}$$

$$p (\bar{x}_{m+1} - \bar{x}_m) = z (\bar{x}_m - \bar{x}_{m-1}), \text{ or}$$

$$\bar{x}_{m+1} - \bar{x}_m = \frac{z}{p} (\bar{x}_m - \bar{x}_{m-1})$$

Write

$$\bar{x}_2 - \bar{x}_1 = \frac{z}{p} (\bar{x}_1 - \bar{x}_0)$$

$$\bar{x}_3 - \bar{x}_2 = \frac{z}{p} (\bar{x}_2 - \bar{x}_1)$$

$$= \left(\frac{z}{p}\right)^2 (\bar{x}_1 - \bar{x}_0)$$

⋮

$$\bar{x}_{m+n} - \bar{x}_{m+n-1} = \left(\frac{z}{p}\right)^{m+n-1} (\bar{x}_1 - \bar{x}_0)$$

From this we get by adding

$$\pi_1 \left(1 + \frac{2}{p} + \dots + \left(\frac{2}{p} \right)^{m+n-1} \right)$$

$$= \pi_{m+n}$$

$$= 1$$

It follows that

$$\pi_1 = \frac{1}{1 + \frac{2}{p} + \dots + \frac{2}{p}^{m+n-1}}$$

As a consequence we have

$$\pi_m = p_{m,n} = \frac{1 + \left(\frac{2}{p} \right) + \dots + \left(\frac{2}{p} \right)^{m-1}}{1 + \frac{2}{p} + \dots + \left(\frac{2}{p} \right)^{m+n-1}}$$

How do we know that the game will end?

Let B_k be defined as

$$B_k = \left\{ \text{tossed } (m+n)k+1, (m+n)k+2, \dots \right. \\ \left. \dots, (m+n)(k+1)-1 \right. \\ \left. \text{produce heads} \right\}$$

By independence

$$P(B_k) = p^{m+n}$$

But B_k depend on disjoint blocks of events so they are independent. But

$$\{\text{game ends}\} \supseteq \bigcup_{k=1}^{\infty} B_k$$

We compute

$$\begin{aligned} P\left(\bigcup_{k=1}^{\infty} B_k\right) &= \lim_{r \rightarrow \infty} P\left(\bigcup_{k=1}^r B_k\right) \\ &= \lim_{r \rightarrow \infty} \left(1 - P\left(\bigcap_{k=1}^r B_k^c\right)\right) \\ &= 1 - \lim_{r \rightarrow \infty} P(B_k^c)^r \\ &= 1 - \lim_{r \rightarrow \infty} (1 - p^{m+n})^r \\ &= \underline{1}. \end{aligned}$$

Lemma 1.3 : Let A_1, A_2, \dots be events.

(i) If $A_1 \subseteq A_2 \subseteq \dots$ then

$$P\left(\bigcup_{k=1}^{\infty} A_k\right) = \lim_{n \rightarrow \infty} P(A_n)$$

(ii) If $A_1 \supseteq A_2 \supseteq \dots$ then

$$P\left(\bigcap_{k=1}^{\infty} A_k\right) = \lim_{n \rightarrow \infty} P(A_n)$$

Proof: We only prove (i).

The second assertion follows by de Morgan rules. Write

$$\bigcup_{k=1}^{\infty} A_k = A_1 \cup (A_2 \setminus A_1) \cup ((A_3 \setminus A_1 \cup A_2) \cup \dots$$

The events in the union on the right are disjoint. It follows that

$$P\left(\bigcup_{k=1}^{\infty} A_k\right) = P(A_1) + P(A_2 \setminus A_1) + \dots$$

The infinite series on the right converges and its partial sum is

$$\begin{aligned} P(A_1) + P(A_2 \setminus A_1) + \dots + P(A_n \setminus (A_1 \cup \dots \cup A_{n-1})) \\ = P\left(\bigcup_{k=1}^n A_k\right) = P(A_n). \end{aligned}$$

The assertion follows.

2. Random variables

2.1. Discrete random variables

In Galileo's example we had $\Omega = \{1, 2, 3, 4, 5, 6\}^3$. The outcomes are of the form (i, j, k) . The gamblers, however, were not so interested in the inner structure of the outcome but in the random number that was the sum of dots. We can imagine that random numbers are created through some process that involves chance. In probability we call such random numbers random variables. We denote them by letters X, Y, Z .

Technical note: Formally we understand random variables as functions from Ω to the real numbers \mathbb{R} . We imagine that some invisible "hand" chooses the outcome ω and the random variable X gives the random number $X(\omega)$.

Definition: A random variable X is a function $X: \Omega \rightarrow \mathbb{R}$ such that $X^{-1}((a, b])$ is an event for all $a < b, a, b \in \mathbb{R}$.

Note: The choice of intervals of the form $(a, b]$ is arbitrary. Intervals of the form $(a, b), [a, b]$ do the same.

Definition: A random variable X is discrete if it can only take values in a finite or countable set $\{x_1, x_2, \dots\}$.

We can see that for discrete random variables we can require $X^{-1}(x_i)$ to be an event for all possible values.

This definition is equivalent to the more general one.

To continue with Galileo's example gamblers were interested in the probabilities that the sum is 9 or 10.

We can ask the question for any $k \in \{3, 4, \dots, 18\}$.

A note on notation : We will

write $\{X = x_k\}$ for the event

$X^{-1}(\{x_k\})$. When we write

probabilities $P(\{X = x_k\})$

we will drop the curly

brackets and write $P(X = x_k)$.

We obviously have $\bigcup_{k=3}^{18} \{X = k\} = \Omega$.

Since the events are disjoint
we have

$$\sum_{k=3}^{18} P(X = k) = P(\Omega) = 1.$$

The total probability 1 is
"distributed" among all
possible values of X .

We say that these probabilities
determine the distribution of X .

Definition : Let X be a discrete random variable. The distribution of X is given by the probabilities $P(X = x_k)$ for all possible values of X .

- There is a number of standard distributions in probability.

Binomial distribution

Suppose we toss a coin n times.

Suppose the tosses are independent

- and the probability of heads is $p \in (0, 1)$. Let $X = \#$ of heads in n tosses. This random variable has values $k = 0, 1, 2, \dots, n$.

To describe the distribution we need to compute $P(X = k)$ for all k .

We have $\Omega = \{H, T\}^n$ and

the event $\{X = k\}$ consists of all sequences $H T H H T \dots H$

that contain exactly k heads.

Every such outcome has the

probability $p^k (1-p)^{n-k}$ because

of independence. So we only

need to compute how many

such outcomes there are. But

this is given by $\binom{n}{k}$ because

we need to choose k positions for heads among n positions.

We have

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

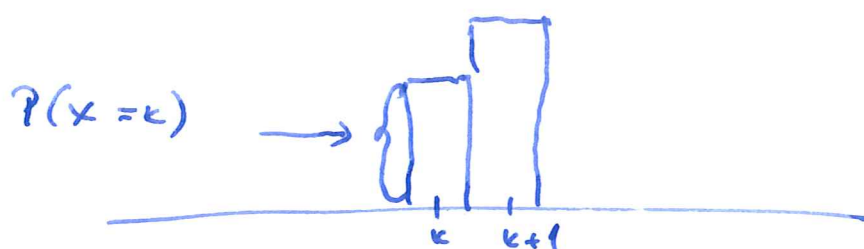
for $k = 0, 1, \dots, n$. We say

that X has binomial distribution with parameters n and p .

Notation: $X \sim \text{Bin}(n, p)$.

The way to visualize a distribution is to draw a histogram. If X is a random variable with integer values we draw a column over a possible value k of X with base 1 and height $P(X=k)$ centered on k .

Figure:



Let us consider $X \sim \text{Bin}(n, p)$.

For $k \geq 1$ we can compute

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

$$= \frac{n-k+1}{k} \cdot \frac{p}{1-p}$$

If $\frac{P(x=k)}{P(x=k-1)} > 1$, then we

have $P(x=k) > P(x=k-1)$, i. e.

the column over k is taller than the column over $k-1$. This happens if

$$\frac{n-k+1}{k} \cdot \frac{p}{1-p} > 1 \Rightarrow$$

$$(n-k+1)p \geq k(1-p) \Rightarrow$$

$$(n+1)p \geq k$$

We have two cases:

1. If $(n+1)p$ is not an integer then the tallest column is over $k = \lfloor (n+1)p \rfloor$.

2. If $(n+1)p$ is an integer then for $k = (n+1)p$ we have

$$\frac{P(x=k)}{P(x=k-1)} = 1$$

This means that the two columns over $k = (u+1)p$ and $(u+1)p-1$ are the tallest and equal.

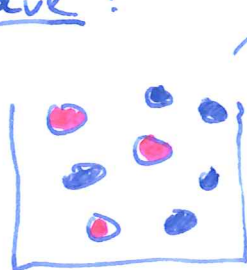
Coin tosses are metaphors for counting "successes" in identical and independent repetitions of the same experiment.

Hyper-geometric distribution

Suppose we have an urn with B black and R red balls.

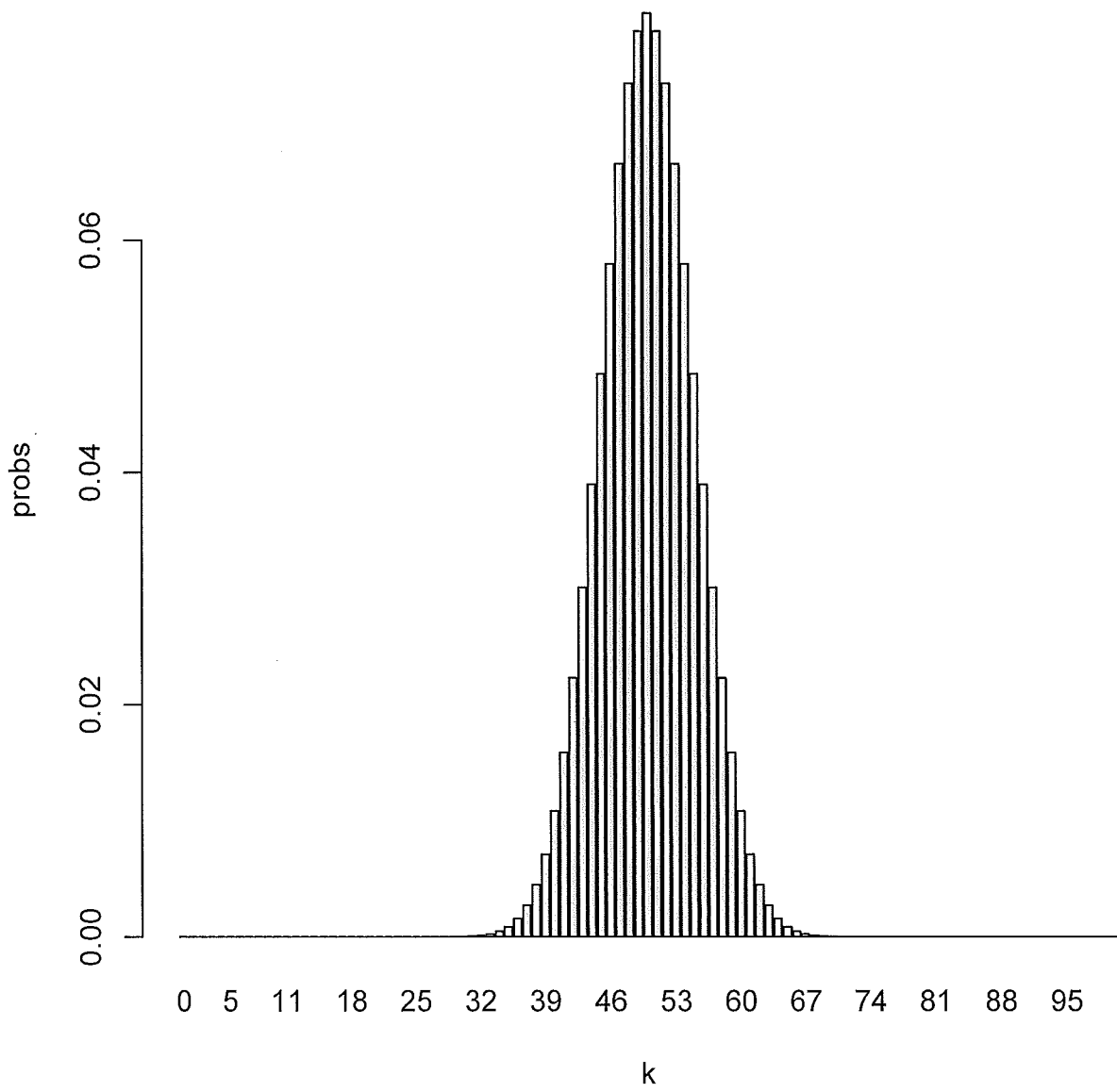
Denote $N = B + R$. Suppose we select $u \leq N$ balls at random from the urn. In mathematical terms this means that all $\binom{N}{u}$ possible selections of u balls are equally likely.

Figure:

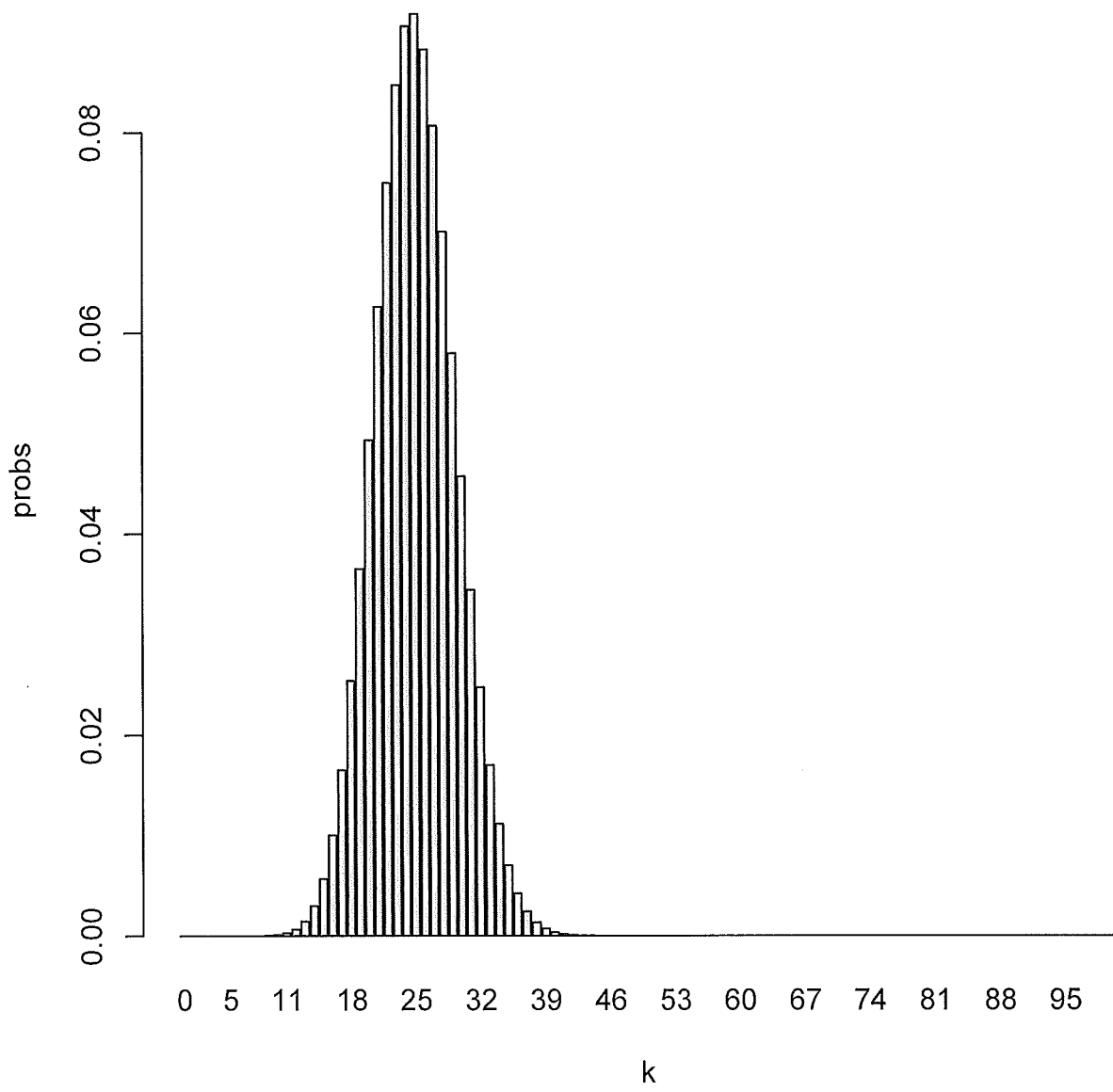


Select u balls at random without replacement.

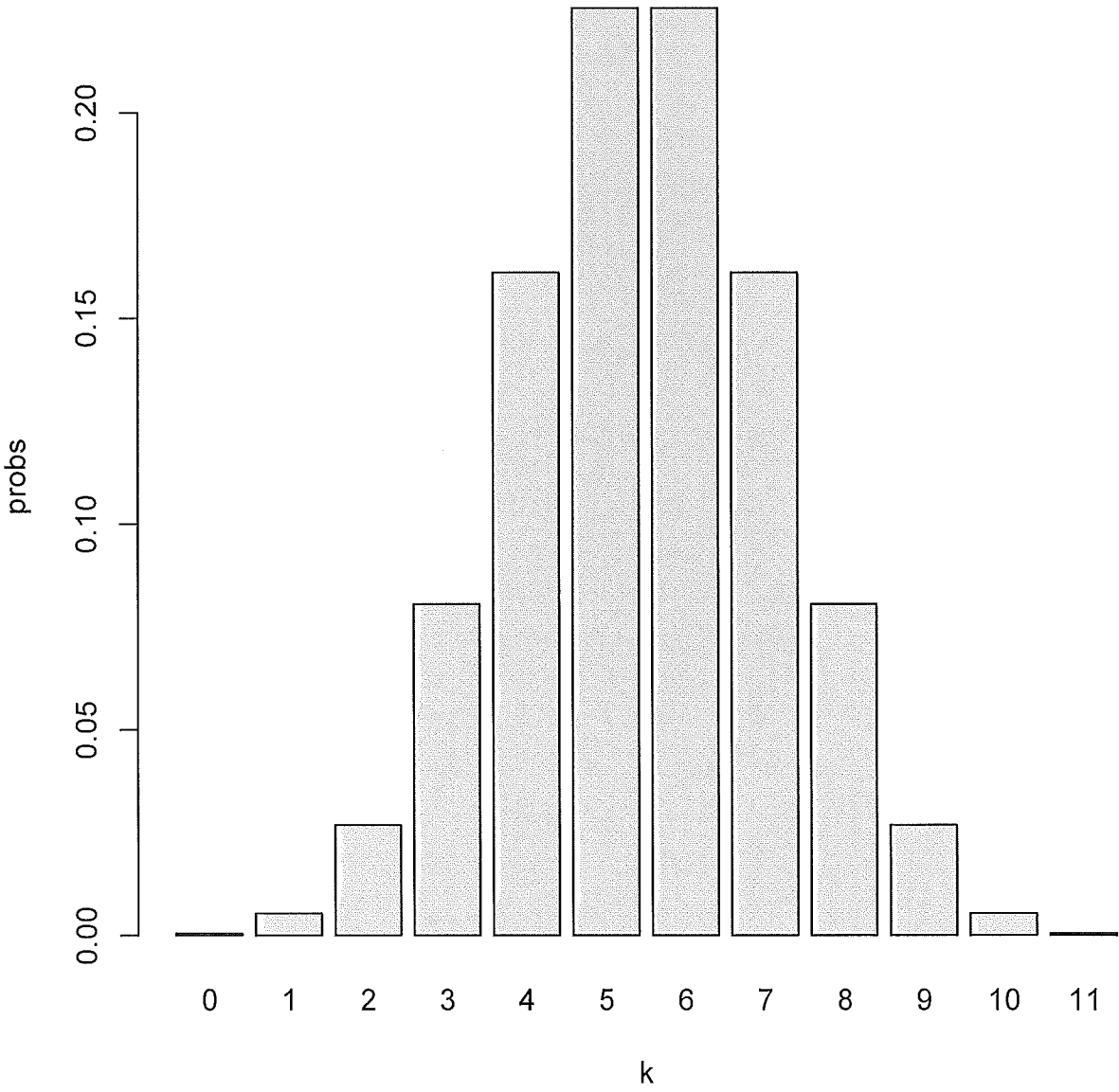
$X \sim \text{Bin}(100, 1/2)$



$X \sim \text{Bin}(100, 1/4)$



$X \sim \text{Bin}(11, 1/2)$



Definition : If

$$P(X=k) = \frac{\binom{B}{k} \binom{R}{n-k}}{\binom{N}{n}} \quad \text{for}$$

$$\max(0, n-R) \leq k \leq \min(n, B)$$

we say that X has the
hyper-geometric distribution with
parameters n, B and $N = B+R$.

Shorthand :

$$X \sim \text{Hypergeom}(n, B, N).$$

Let X = number of black balls among the n selected. X is a random variable with values k that must satisfy

$$\max(0, n-R) \leq k \leq \min(n, B).$$

We have

$$P(X=k) = \frac{\binom{B}{k} \binom{R}{n-k}}{\binom{N}{n}}$$

The denominator is the number of all possible selection and the numerator is the number of selections with exactly k black and $n-k$ white balls. As with the binomial distribution we can calculate

$$\frac{P(X=k)}{P(X=k-1)} = \frac{(B-k+1)}{k} \cdot \frac{(n-k+1)}{R-n+k}$$

After some calculation we find that

$$\frac{P(X=k)}{P(X=k-1)} > 1 \quad \text{if} \quad k < \frac{(B+1)(u+1)}{N+2}$$

Again we have two cases:

1. $\frac{(B+1)(u+1)}{(N+2)}$ is not an integer.

Then $k = \lfloor \frac{(B+1)(u+1)}{N+2} \rfloor$ is the tallest column.

2. $\frac{(B+1)(u+1)}{N+2}$ is an integer.

Then $k = \frac{(B+1)(u+1)}{N+2}$ is still the largest probability but is equal to $P(X=k-1)$.

Example : Lottery.

A lottery ticket has 39 numbers. We cross out m numbers where $m = 8, 9, \dots, 17$.

The winnings depend on the draw.
Each week 7 numbers are drawn.

If all the numbers are among
the ones we crossed out we win
a large amount of money. We

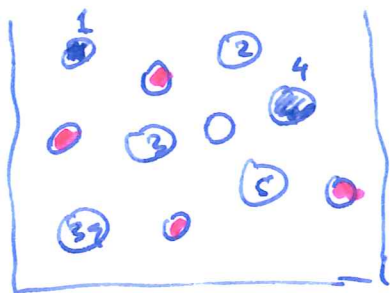
can translate this problem into
a problem involving the

hyper-geometric distribution

I imagine you have balls numbered
1-39. You put them into an

urn. When we cross the
numbers on the lottery
ticket we paint

Figure:

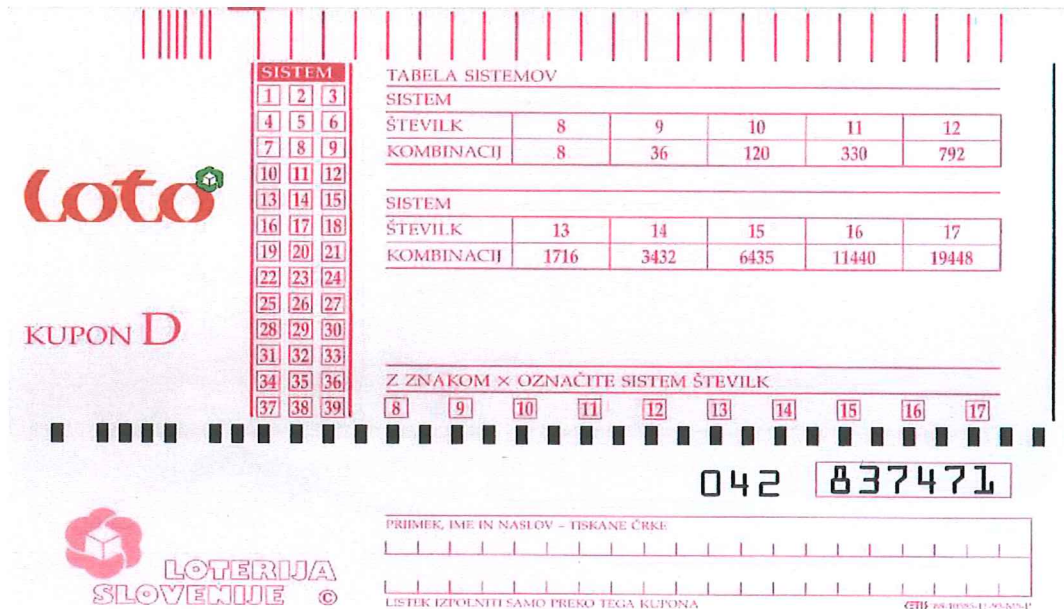


those numbers black.
The others we paint red.
When 7 balls are drawn
the number of correct

guesses is $X \sim \text{Hyper Geom}(7, m, 39)$

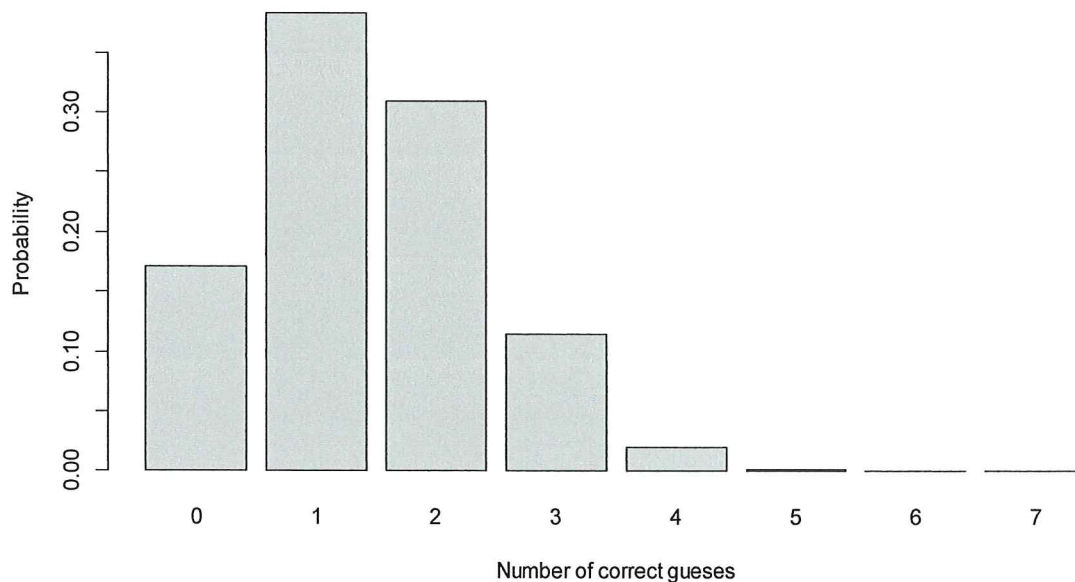
where m is the number of
balls we painted black.

Below are distributions for correct number of guesses in Lottery for $m=8,13,17$. We translated the winning odds in Lottery to a question about the hyper-geometric distribution. The Lottery ticket looks like



On the ticket the player can cross from $m=8$ to $m=17$ numbers. The number of correct guesses is the basis for determining the winnings. The correct guesses are a random variable X . The probability $P(X = 7)$ is the most interesting as it is the probability of jack-pot.

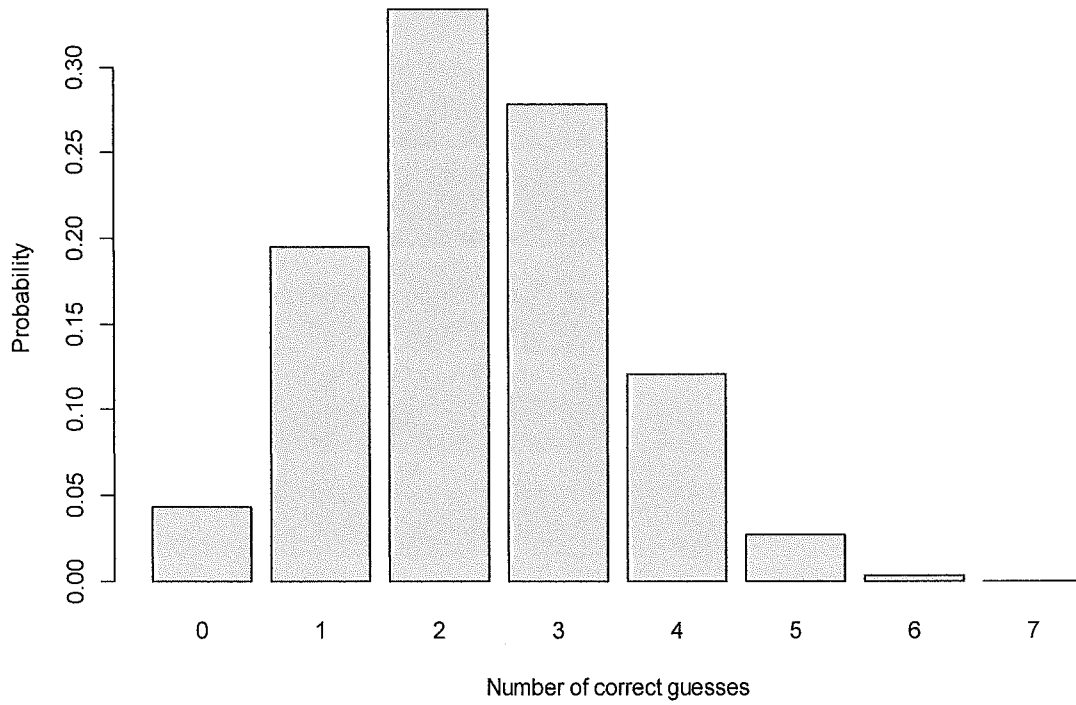
Histogram for Lottery with $m=8$



The numerical values of the above probabilities are:

1.709633e-01 3.829577e-01 3.093120e-01 1.145600e-01 2.045714e-02 1.693005e-03 5.643349e-05 5.201244e-07

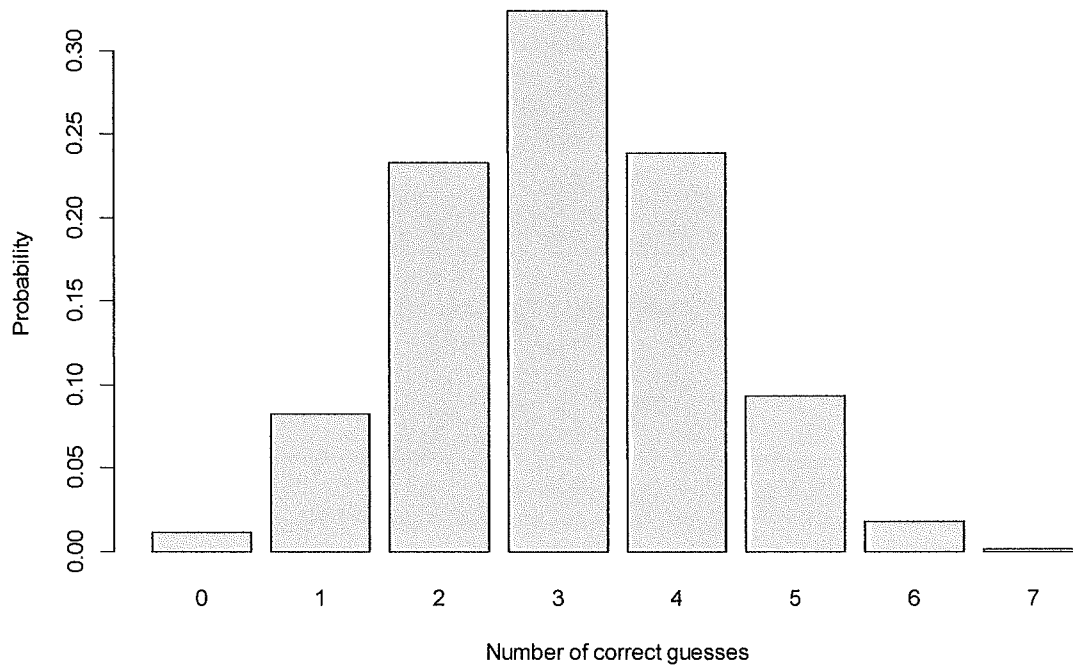
Histogram for Lottery with m=13



The numerical values are:

0.0427672254 0.1945908757 0.3335843584 0.2779869653 0.1208638980 0.0271943770
0.0029007336 0.0001115667

Histogram for Lottery with m=17



Numerical values:

0.011088011 0.082467082 0.232848233 0.323400323 0.238294975 0.092935040 0.017701912
0.001264422

Geometric and negative binomial distribution

Suppose we toss a coin. The tosses are independent and the probability of heads is p .

Let X be the number of tosses until the first heads.

X is a random variable with values in $k = 1, 2, \dots$

Remark: The values of X can be arbitrarily large. The event $\{X = k\}$ happens if we get $\underbrace{TT \dots T}_{(k-1) \text{ times}} H$. By independence this implies

$$P(X = k) = (1-p)^{k-1} \cdot p \quad k = 1, 2, \dots$$

Definition : If

$$P(X=k) = (1-p)^{k-1} \cdot p$$

we say that X has geometric distribution with parameter p .

Shorthand : $X \sim \text{Geom}(p)$.

○ Example : We play roulette and wait for the number 17 to appear. What is the probability that 17 will not appear in the first n tosses?

$$P(X > n) = P(\underbrace{T \dots T}_{n \text{ times}}) = (1-p)^n$$

Because $X \sim \text{Geom}(1/37)$ this means

$$P(X > n) = \left(\frac{36}{37}\right)^n$$

If instead of waiting for heads we wait for the appearance of m heads then X is a random variable with values $k = m, m+1, \dots$

We have

$$P(X = k) = P(\{ \text{exactly } m-1 \text{ heads in first } k-1 \text{ tosses} \} \cap \{ \text{heads of } k\text{-th toss} \})$$

(independence),

$$= P(\text{exactly } m-1 \text{ heads in first } (k-1) \text{ tosses}) \times P(\text{heads on toss } k)$$

$$= \underbrace{\binom{k-1}{m-1} p^{m-1} (1-p)^{(k-1)-(m-1)}}_{\text{by binomial distribution} \times p}$$

$$= \binom{k-1}{m-1} p^m (1-p)^{k-m},$$

$$k = m, m+1, \dots$$

Definition: we say that X has negative binomial distribution with parameters m and p if

$$P(X = k) = \binom{k-1}{m-1} p^m (1-p)^{k-m},$$

$$k = m, m+1, \dots$$

○ Shorthand: $X \sim \text{NegBin}(m, p)$.

Example: The Polish mathematician Stefan Banach (1897 - 1945) was a chain smoker. He always carried two boxes of matches in his pockets. Assume Banach starts with two boxes of n matches. Then he randomly-reaches into the left or right pocket at random with probability $\frac{1}{2}$. At some stage Banach will take the last match from a box but will not notice it. The first time Banach pulls an empty match box from his pockets, the number of matches in the other box is random. Call it X . Possible values for X are $k = 0, 1, \dots, n$. We would like to compute the distribution of X .

Let us define

$A = \{x = k\} \cap \{\text{Banach pulls the empty box from left pocket}\}$

By symmetry and the law of total probabilities we have

$$P(x = k) = 2 \cdot P(A).$$

Let us picture Banach cigarettes.



The event A happens when Banach lights his $(n + (n-k) + 1)$ -st cigarette and at that point has pulled a box from his left pocket exactly $(n+1)$ -st time.

Declare: reach into the left pocket = "success".

Poisson distribution.

Let us look at the binomial distribution $\text{Bin}(n, \frac{\lambda}{n})$ for a given $\lambda > 0$. What happens if n is "large"? Let k be fixed. For $n \geq k$ we have

for $X_n \sim \text{Bin}(n, \frac{\lambda}{n})$ that

$$P(X = k) = \binom{n}{k} \left(\frac{\lambda}{n}\right)^k \left(1 - \frac{\lambda}{n}\right)^{n-k}$$

What happens when $n \rightarrow \infty$?

From Analysis we know that

$$\left(1 + \frac{x}{n}\right)^n \rightarrow e^x \quad \text{for } x \in \mathbb{R}.$$

Rewrite

$$P(X_n = k) = \frac{\lambda^k}{k!} \cdot \frac{n(n-1) \cdots (n-k+1)}{n^k} \times \left(1 - \frac{\lambda}{n}\right)^n \times \left(1 - \frac{\lambda}{n}\right)^{-k}$$

$\begin{matrix} \nearrow 1 \\ \nearrow e^{-\lambda} \\ \nearrow 1 \end{matrix}$

We find that

$$\lim_{n \rightarrow \infty} P(X_n = k) = \frac{\lambda^k}{k!} \cdot e^{-\lambda}.$$

This motivates the following definition

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

for $\lambda > 0$ and $k = 0, 1, 2, \dots$

we say that X has the Poisson distribution with parameter $\lambda > 0$.

Shorthand: $X \sim \text{Po}(\lambda)$.

Remark: From analysis we know that

$$\sum_{k=0}^{\infty} \frac{\lambda^k}{k!} = e^{\lambda}.$$

If $X \sim \text{Po}(\lambda)$ we do get

$$\sum_{k=0}^{\infty} P(X = k) = 1.$$

2.2. Continuous distributions

We can imagine "random numbers" that can take any real number as a value.

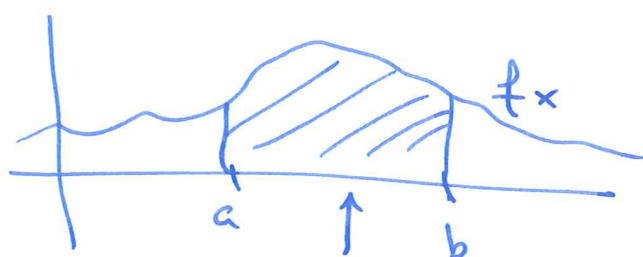
Examples are lifetimes of components, a randomly chosen point on the interval $[0, 1]$, ...

Technically, X is still a function on Ω and we require $X^{-1}((a, b])$ to be an event for all $a < b$.

Definition: The distribution of a random variable is given by the probabilities $P(X \in (a, b])$ for all $a < b$.

The idea of continuous random variables is to describe probabilities $P(X \in (a, b))$ by integrals of a single function.

Figure :



$$\text{Area} = P(a < X \leq b)$$

Definition : The random variable X has continuous distribution if there is a non-negative function $f_x(x)$ called the density such that

$$P(a < X \leq b) = \int_a^b f_x(x) dx$$

for all $a < b$.

Continuous distributions are used in financial modelling and statistics because of their practicality.

There are standard distributions that are often used.

○ Normal distribution

The random variable X has normal distribution with parameters μ and σ if the density is given by

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

Remark: For densities we must have $\int_{-\infty}^{\infty} f_X(x) dx = 1$.

We believe mathematically that f_X is a density.

We will use the notation :

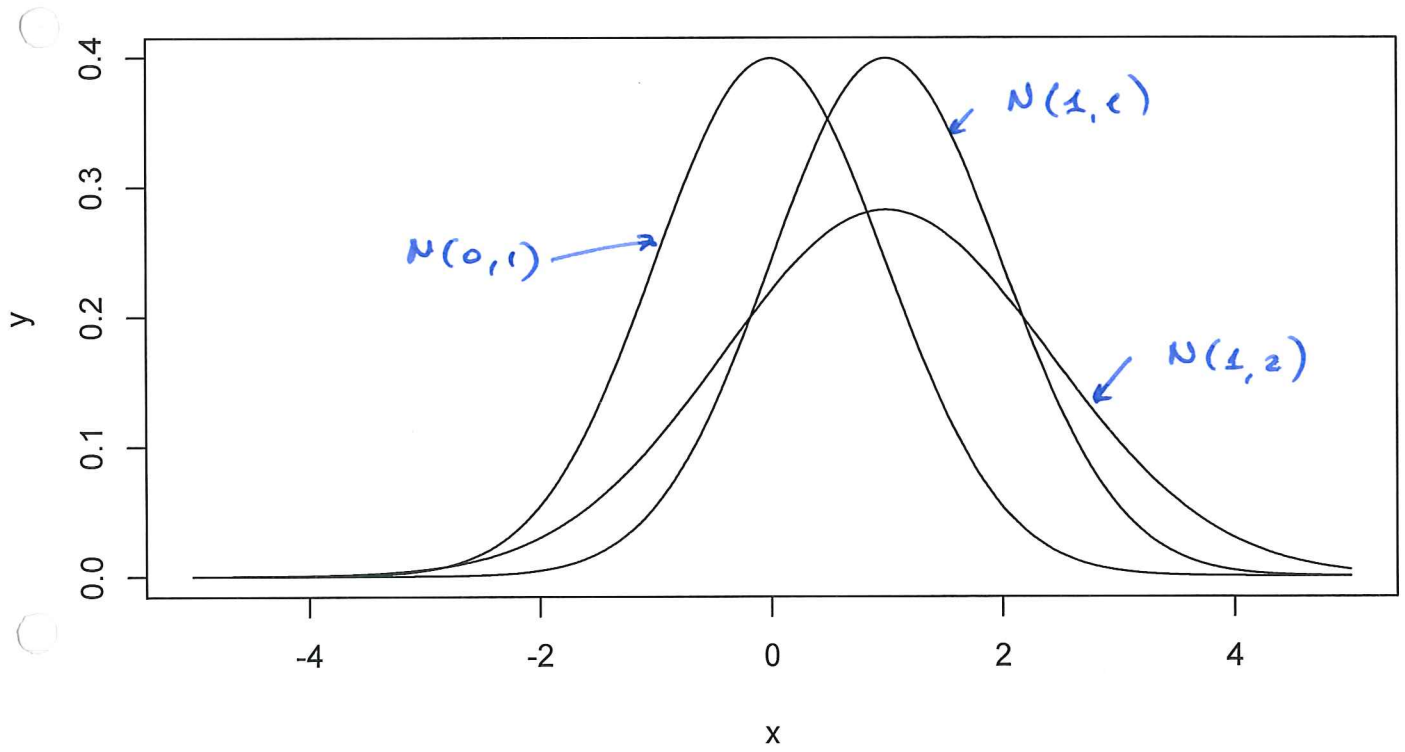
$$X \sim N(\mu, \sigma^2).$$

Remark : The name normal distribution was chosen by the Belgian statistician Adolphe

Quetelet because the histograms of human characteristics such as IQ, height and others are close to normal.

Remark : We will give the right interpretation to parameters μ in σ^2 later.

$X \sim N(0,1), N(1,1), N(1,3)$



Exponential and gamma distribution

The density of the exponential distribution is given by

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

We say that X has the exponential distribution with parameter λ .

Notation : $X \sim \text{exp}(\lambda)$

The exponential distribution is used to model lifetimes of electronic components.

To define the gamma distribution recall the definition of the gamma function.

$$\Gamma(x) = \int_0^{\infty} u^{x-1} \cdot e^{-u} du.$$

The most important properties are:

(i) $\Gamma(x+1) = x \cdot \Gamma(x)$

(ii) $\Gamma(n) = (n-1)!$

(iii) $\Gamma(1/2) = \sqrt{\pi}$.

The density

$$f_X(x) = \begin{cases} \frac{\lambda^a}{\Gamma(a)} x^{a-1} e^{-\lambda x}, & x > 0 \\ 0 & \text{else.} \end{cases}$$

is called the gamma density with parameters a and λ .

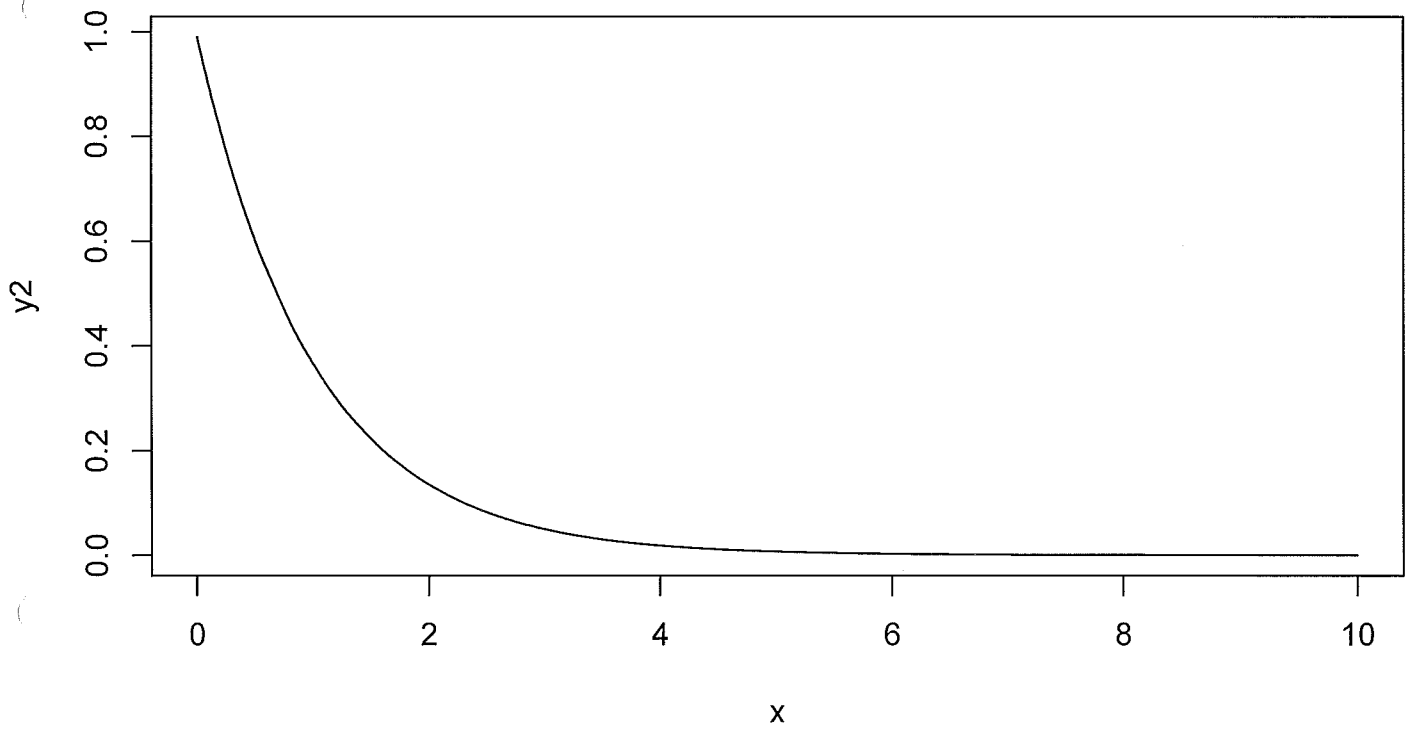
a : shape parameter

λ : scale parameter.

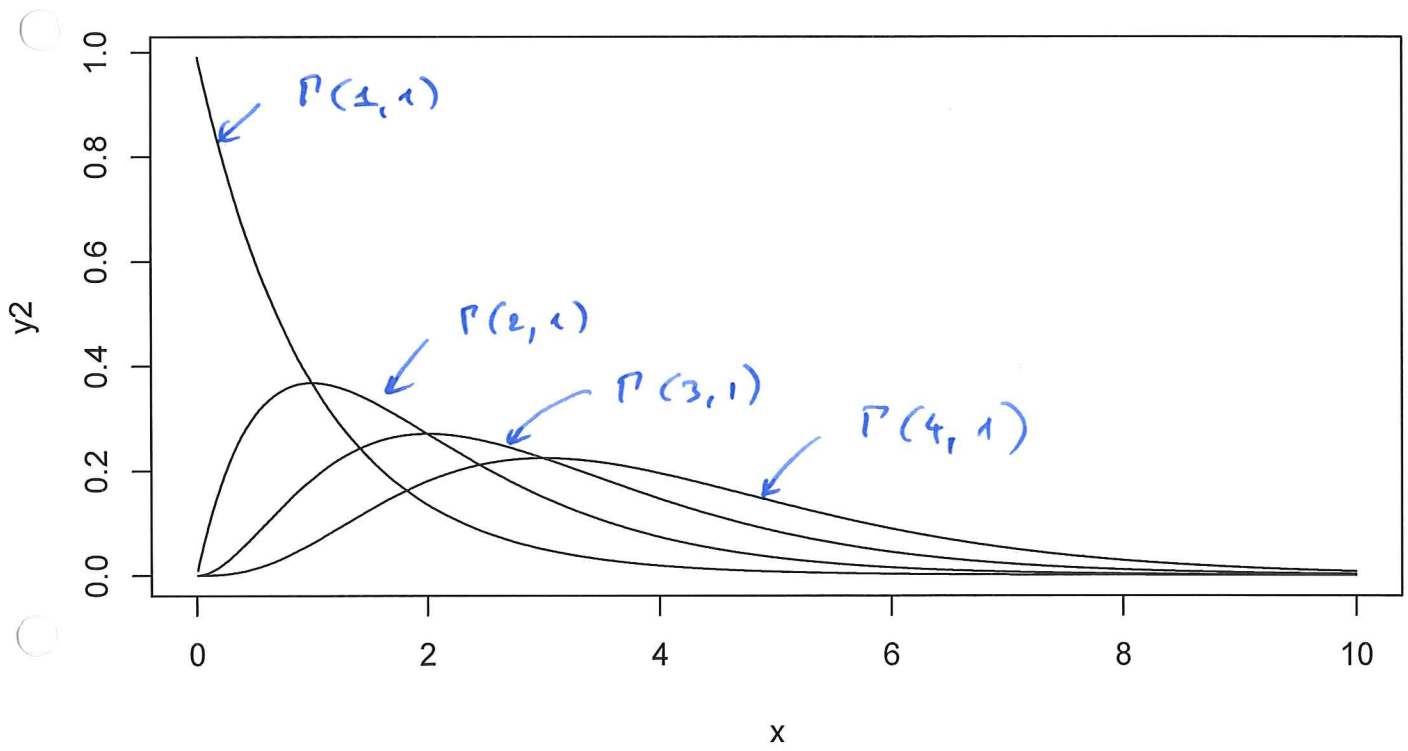
Remark: Sometimes $1/\lambda$ is used instead of λ .

Notation: $X \sim P(a, \lambda)$

Gamma distribution $a=0.5$, $\lambda=0.5$



Various gamma distributions



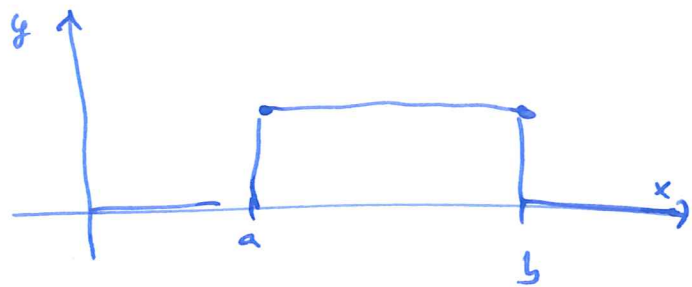
Uniform distribution

The uniform distribution models the choice of a point at random uniformly on the interval (a, b) .

The density is given by

$$f_X(x) = \begin{cases} \frac{1}{b-a} & \text{for } a \leq x \leq b. \\ 0 & , \text{ else.} \end{cases}$$

Figure :



We use the notation :

$$X \sim U(a, b)$$

Most computer generated random numbers are uniform on $[0, 1]$.

2.3. Functions of random variables

Let X be a random variable.

We have that $\{X \leq x\}$ is an event. So the probability is defined.

Definition : The distribution function of X is defined as the function

$$F_X(x) = P(X \leq x)$$

Theorem 2.1 : Let X be a random variable with distribution function F_X .

- (i) F_X is nondecreasing.
- (ii) $\lim_{x \rightarrow \infty} F_X(x) = 1$, $\lim_{x \rightarrow -\infty} F_X(x) = 0$
- (iii) F_X is right continuous.

Proof :

(i) For $x < y$ we have $\{x \leq x\} \subseteq \{x \leq y\}$.

It follows that $P(x \leq x) \leq P(x \leq y)$.

(ii) We have $\Omega = \bigcup_{n=1}^{\infty} \{x \leq n\}$.

The sets in the union are increasing so

$$\begin{aligned} 1 - P(\Omega) &= \lim_{n \rightarrow \infty} P(x \leq n) \\ &= \lim_{n \rightarrow \infty} F_x(n) \end{aligned}$$

The conclusion follows because F_x is nondecreasing. The other limit is proved similarly.

(iii) Fix $x \in \mathbb{R}$. Let $x_n \downarrow x$.

We have $\{x \leq x\} = \bigcap_{n=1}^{\infty} \{x \leq x_n\}$.

The sets $\{x \leq x_n\}$ are decreasing.

It follows

$$P(x \leq x) = \lim_{n \rightarrow \infty} P(x \leq x_n) \quad \text{or}$$

$$F_x(x) = \lim_{n \rightarrow \infty} F_x(x_n)$$

The last statement is equivalent to right continuity.

If X has density $f_X(x)$ then

$$F_X(x) = \int_{-\infty}^x f_X(u) du.$$

Conversely, if

$$F_X(x) = \int_{-\infty}^x g(u) du$$

for all $x \in \mathbb{R}$ and a nonnegative g then g is (a version of)

the density. If g is continuous at x then

$$F_X'(x) = g(x) = f_X(x).$$

Example: Let $X \sim N(0,1)$ i.e.

the density of X is

$$f_X(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$$

Denote $F_X(x)$ by $\Phi(x)$, i.e.

$$\Phi(x) = \int_{-\infty}^x f_X(u) du.$$

Let $Y = aX + b$ for $a > 0$. What is the density of Y ? We have

$$\begin{aligned} P(Y \leq y) &= P(aX + b \leq y) \\ &= P(aX \leq y - b) \\ &= P\left(X \leq \frac{y-b}{a}\right) \\ &= \int_{-\infty}^{\frac{y-b}{a}} \frac{1}{\sqrt{2\pi}} e^{-u^2/2} du \end{aligned}$$

New variable:

$$\frac{y-b}{a} = u \quad \frac{du}{a} = du$$

Example: Let $X \sim N(0,1)$ and

$Y = X^2$. Density of Y ? For $y > 0$

$$\begin{aligned} P(Y \leq y) &= P(X^2 \leq y) \\ &= P(-\sqrt{y} \leq X \leq \sqrt{y}) \\ &= \Phi(\sqrt{y}) - \Phi(-\sqrt{y}) \end{aligned}$$

Comment: In general

$$\boxed{P(a < X \leq b) = F_X(b) - F_X(a)}$$

For continuous random variables
the probabilities $P(X=x) = 0$

for all x and it is not relevant
whether we write $<$ or \leq . We
have

$$\begin{aligned} P(Y \leq y) &= \Phi(\sqrt{y}) - \Phi(-\sqrt{y}) \\ &= \frac{1}{\sqrt{2\pi}} \int_{-\sqrt{y}}^{\sqrt{y}} e^{-u^2/2} du \\ &= \frac{2}{\sqrt{2\pi}} \int_0^{\sqrt{y}} e^{-u^2/2} du \end{aligned}$$

New variable:

$$u^2 = v \Rightarrow$$

$$2u du = dv$$

$$du = \frac{dv}{2\sqrt{v}}$$

$$= \frac{2}{\sqrt{2\pi}} \int_0^y \frac{1}{2\sqrt{v}} e^{-v/2} dv$$

$$= \frac{1}{\sqrt{2\pi}} \int_0^y \frac{1}{\sqrt{v}} e^{-v/2} dv$$

Since $P(Y \geq 0) = 1$ we have

$$f_Y(y) = \begin{cases} \frac{1}{\sqrt{2\pi} \cdot \sqrt{y}} e^{-y/2}, & y > 0 \\ 0, & \text{else.} \end{cases}$$

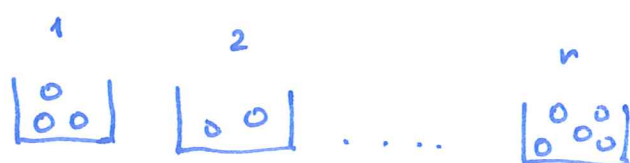
We recognize: $Y \sim \Gamma(1/2, 1/2)$.

3. Multivariate distributions

3.1. Discrete multivariate distributions

Example: Suppose we have r boxes. We are dropping balls into these boxes at random. The probabilities that we hit box k is p_k for $k = 1, 2, \dots, r$; Assume the subsequent drops are independent. There are n balls.

Figure:



We end up with random numbers of balls in boxes. Denote these random numbers by X_1, X_2, \dots, X_r . These random numbers will "in the collective" equal to k_1, k_2, \dots, k_r where $k_i \geq 0$ and $\sum_{i=1}^r k_i = n$. All the random variables X_1, X_2, \dots, X_r simultaneously take a collection

of values. The mathematical objects with several components are vectors.

By analogy we will say that

$\underline{X} = (X_1, X_2, \dots, X_r)$ is a random vector.

The possible values of this random vector are vectors (k_1, k_2, \dots, k_r) with $k_i \geq 0$ and $\sum_{i=1}^r k_i = n$.

For discrete random variables we had that the distribution was given by $P(X=x)$ for all possible x . By analogy

the distribution of the random vector \underline{X} will be given by probabilities $P(\underline{X} = \underline{x})$

where \underline{x} are possible collections/vectors of values. In the above example we need to compute

$$P(\underline{X} = (k_1, k_2, \dots, k_r)) = P(\underbrace{X_1 = k_1, X_2 = k_2, \dots, X_r = k_r}_{\uparrow})$$

This notation means

$$\bigcap_{i=1}^r \{X_i = k_i\}$$

If we want to hit box 1, k_1 times, box 2 k_2 times, ... the possible disjoint ways for this to happen is to get a sequence of hits

$$n_1, n_2, \dots, n_n$$

where k_1 of the n_1, n_2, \dots, n_n are equal to 1, k_2 are equal to 2, ...

The probability of such a sequence of hits is by independence

$$p_1^{k_1} \cdot p_2^{k_2} \cdot \dots \cdot p_r^{k_r}$$

How many sequences of this type are there? we have n positions

$$\begin{matrix} 1 & 2 & & & & & & & n \\ \sqcup & \sqcup & \sqcup & \dots & \dots & \dots & \dots & \dots & \sqcup \end{matrix} \leftarrow \text{Positions}$$

We first choose k_1 positions for 1s. We can do this in $\binom{n}{k_1}$.

Among the $n - k_1$ positions left we choose k_2 positions for $2s$. We can do this in $\binom{n - k_1}{k_2}$ ways.

By the fundamental theorem of combinatorics the total number of possibilities is

$$\begin{aligned} & \binom{n}{k_1} \binom{n - k_1}{k_2} \cdots \binom{n - k_1 - \cdots - k_{v-1}}{k_v} \\ &= \frac{n!}{k_1! (n - k_1)!} \cdot \frac{(n - k_1)!}{k_2! (n - k_1 - k_2)!} \cdots \frac{(n - k_1 - \cdots - k_{v-1})!}{k_v! \cdot 0!} \\ &= \frac{n!}{k_1! \cdot k_2! \cdots k_v!} \end{aligned}$$

All the sequences are disjoint events with the same probabilities. It follows

$$P(X_1 = k_1, \dots, X_v = k_v) = \frac{n!}{k_1! \cdots k_v!} p_1^{k_1} \cdots p_v^{k_v}$$

for $k_i \geq 0$ for $i = 1, 2, \dots, v$ and $\sum_{i=1}^v k_i = n$.

Definition: For a vector with the above distribution we say that it has the multinomial distribution with parameters n and $p = (p_1, p_2, \dots, p_r)$.

Shorthand: $\underline{X} \sim \text{Multinomial}(n, p)$.

Definition: A discrete random vector $\underline{X} = (X_1, X_2, \dots, X_r)$ is a function $\underline{X} : \Omega \rightarrow \{x_1, x_2, \dots, \dots\}$ where $\{x_1, x_2, \dots\}$ is a finite or countable set of possible values, and such that all components X_1, X_2, \dots, X_r are random variables.

Definition: The distribution of a random vector \underline{X} with values in $\{x_1, x_2, \dots\}$ is given by probabilities $P(\underline{X} = \underline{x}_k)$ for all $k = 1, 2, \dots$.

Remark: Typically we will write

$P(x_1 = x_1, \dots, x_r = x_r)$. When the number of components is small we often write $P(x=x, y=y, z=z)$.

Example: Let $N \geq 3$. Choose three

numbers at random ~~from~~ from $\{1, 2, \dots, N\}$ without replacement so that all subsets of three numbers are equally likely. Let x be the smallest of the three numbers, z the largest and y the remaining one.

Example: If we choose 5, 3, 7 ~~we~~ we have $x=3, y=5, z=7$.

What is the distribution of (x, y, z) ?

The possible values are triplets (i, j, k) with $1 \leq i < j < k \leq N$.

We have

$$\begin{aligned} P(X=i, Y=j, Z=k) \\ &= P(\text{we select the subset } \{i, j, k\}) \\ &= \frac{1}{\binom{N}{3}} \end{aligned}$$

What is the distribution of X ?

It has possible values $1, 2, \dots, N-2$.

We notice $\{X=i\} = \cup_{i < j < k \leq N} \{X=i, Y=j, Z=k\}$

$$\begin{aligned} P(X=i) &= \sum_{i < j < k \leq N} P(X=i, Y=j, Z=k) \\ &= \frac{\binom{N-i}{2}}{\binom{N}{3}} \\ &= \frac{(N-i)(N-i-1)}{2} \\ &= \frac{N(N-1)(N-2)}{6} \\ &= \frac{3(N-i)(N-i-1)}{N(N-1)(N-2)} \end{aligned}$$

Definitions:

- (i) The distributions of components of a random vector $\underline{X} = (X_1, X_2, \dots, X_n)$ are called univariate marginal distributions.
- (ii) The distributions of subvectors like (X_1, X_2, \dots, X_s) for $s < n$ are called multivariate marginal distributions.

Example (continuation): What is the distribution of (X, Y) ? We write

$$P\{X=i, Y=j\} = \underbrace{\sum_{k=j+1}^N P\{X=i, Y=j, Z=k\}}_{\text{disjoint events}}$$

We have

$$\begin{aligned} P(X=i, Y=j) &= \sum_{k=j+1}^N P(X=i, Y=j, Z=k) \\ &= \frac{N-j}{\binom{N}{3}} \end{aligned}$$

for $1 \leq i < j < N$.

If $\underline{X} = (X_1, \dots, X_r)$ is a random vector let us write

$$\underline{X}^1 = (X_1, \dots, X_s) \text{ and } \underline{X}^2 = (X_{s+1}, \dots, X_r).$$

Theorem 3.1: Let $\mathcal{R} = \{ \underline{x}_1, \underline{x}_2, \dots \}$ be the set of possible values of \underline{X} .

The marginal distribution of \underline{X}^1 is given by

$$\begin{aligned} P(\underline{X} = \underline{x}^1) &= \sum_{(\underline{x}^1, \underline{x}^2) \in \mathcal{R}} P(\underline{X} = (\underline{x}^1, \underline{x}^2)) \\ &= \sum_{(\underline{x}^1, \underline{x}^2) \in \mathcal{R}} P(\underline{X}^1 = \underline{x}^1, \underline{X}^2 = \underline{x}^2) \end{aligned}$$

Proof: We write

$$\{ \underline{X}^1 = \underline{x}^1 \} = \underbrace{\cup_{(\underline{x}^1, \underline{x}^2) \in \mathcal{R}} \{ \underline{X}^1 = \underline{x}^1, \underline{X}^2 = \underline{x}^2 \}}_{\text{disjoint union.}}$$

It follows that

$$P(\underline{X}^1 = \underline{x}^1) = \sum_{(\underline{x}^1, \underline{x}^2) \in \mathcal{R}} P(\underline{X}^1 = \underline{x}^1, \underline{X}^2 = \underline{x}^2).$$

Independence

For two events A and B we say that they are independent if

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

We would like to define independence for random variables. If X and Y are to be

independent we expect the events $\{X=x\}$ and $\{Y=y\}$ to be independent.

$$\text{So we need } P(X=x, Y=y) = P(X=x)P(Y=y)$$

This is the right intuition. For

the formal definition we generalize to

$$\begin{aligned} P(X \in A, Y \in B) &= \sum_{(x,y) \in A \times B} P(X=x, Y=y) \\ &= \sum_{(x,y) \in A \times B} P(X=x)P(Y=y) \\ &= \left(\sum_{x \in A} P(X=x) \right) \left(\sum_{y \in B} P(Y=y) \right) \\ &= P(X \in A) \cdot P(Y \in B). \end{aligned}$$

Definitions:

(i) Discrete random variables X and Y are independent if

$$P(X \in A, Y \in B) = P(X \in A) P(Y \in B)$$

for any two sets A and B .

(ii) Random variables X_1, X_2, \dots, X_r are independent if

$$P(X_1 \in A_1, \dots, X_r \in A_r) = P(X_1 \in A_1) \dots P(X_r \in A_r)$$

for any sets A_1, A_2, \dots, A_r .

Remark: The second definition is equivalent to saying that

$$P(X_1 = x_1, X_2 = x_2, \dots, X_r = x_r) =$$

$$= P(X_1 = x_1) P(X_2 = x_2) \dots P(X_r = x_r)$$

for all possible values (x_1, \dots, x_r)

of $\underline{X} = (X_1, \dots, X_r)$.

Example: Let $\underline{X} \sim \text{Multinomial}(n, \underline{p})$.

We can easily guess that

$$X_k \sim \text{Bin}(n, p_k) \quad \text{for } k = 1, 2, \dots, r.$$

So

$$P(X_1 = k_1) \dots P(X_r = k_r)$$

$$= \binom{n}{k_1} p_1^{k_1} (1-p_1)^{n-k_1} \dots \binom{n}{k_r} p_r^{k_r} (1-p_r)^{n-k_r}$$

and

$$P(X_1 = k_1, \dots, X_r = k_r) = \frac{n!}{k_1! \dots k_r!} p_1^{k_1} \dots p_r^{k_r}$$

Since $P(X_1 = k_1, \dots, X_r = k_r) \neq P(X_1 = k_1) \dots P(X_r = k_r)$

there is no independence.

Example: Suppose the number of children in a family is Poisson with parameter $\lambda > 0$. Suppose all children are boys or girls with equal probability independently of

each other. Let X be the number of boys and Y the number of girls.

We compute with $N = X + Y$

$$P(X=k, Y=e) = P(X=k, Y=e, N=k+e)$$

$$= P(X=k, Y=e | N=k+e) P(N=k+e)$$

$$= \binom{k+e}{k} \left(\frac{1}{2}\right)^{k+e} \cdot e^{-\lambda} \cdot \frac{\lambda^{k+e}}{(k+e)!}$$

$$= e^{-\lambda/2} \cdot \frac{\left(\frac{\lambda}{2}\right)^k}{k!} \cdot e^{-\lambda/2} \cdot \frac{\left(\frac{\lambda}{2}\right)^e}{e!}$$

On the other hand

$$P(X=k) = \sum_{n=k}^{\infty} P(X=k, N=n)$$

$$= \sum_{n=k}^{\infty} P(X=k | N=n) P(N=n)$$

$$= \sum_{n=k}^{\infty} \binom{n}{k} \left(\frac{1}{2}\right)^n \cdot e^{-\lambda} \cdot \frac{\lambda^n}{n!}$$

$$= \sum_{n=k}^{\infty} \frac{n!}{k! (n-k)!} \left(\frac{1}{2}\right)^n \cdot e^{-\lambda} \cdot \frac{\lambda^n}{n!}$$

$$= \frac{e^{-\lambda} \lambda^k (1/2)^k}{k!} \underbrace{\sum_{n=k}^{\infty} \frac{(1/2)^{n-k}}{(n-k)!}}_{e^{\lambda/2}}$$

$$= \frac{e^{-\lambda/2} (1/2)^k}{k!}$$

We have (the same calculation is valid for girls)

$$P(X=k, Y=e) = P(X=k)P(Y=e)$$

so X, Y are independent.

Theorem 3.2: Suppose X, Y are discrete random variables with values in $\{x_1, x_2, \dots\}$ and $\{y_1, y_2, \dots\}$.

Suppose we have

$$P(X=x, Y=y) = f(x)g(y)$$

for all pairs $(x, y) \in \{x_1, x_2, \dots\} \times \{y_1, y_2, \dots\}$

for some functions $f: \{x_1, x_2, \dots\} \rightarrow \mathbb{R}$

and $g: \{y_1, y_2, \dots\} \rightarrow \mathbb{R}$. Then X

and Y are independent.

Proof: By Theorem 3.1 the marginal distributions are

$$\begin{aligned}P(X=x) &= \sum_y P(X=x, Y=y) \\&= \sum_y f(x) g(y) \\&= f(x) \cdot \underbrace{\sum_y g(y)}_{= c_1}\end{aligned}$$

Similarly

$$P(Y=y) = c_2 \cdot g(y).$$

It follows that

$$\begin{aligned}P(X=x, Y=y) &= \cancel{P(X=x)} \\&= f(x) g(y) \\&= \frac{P(X=x)}{c_1} \cdot \frac{P(Y=y)}{c_2}\end{aligned}$$

To finish the proof we need $c_1 c_2 = 1$.

But

$$\sum_{x,y} P(X=x, Y=y) = 1 \quad \text{and}$$

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

$$= \left(\sum_x P(X=x) \right) \left(\sum_y P(Y=y) \right)$$

$$= 1 \cdot 1.$$

Summing up we get

$$\sum_{x,y} P(X=x, Y=y) = \frac{1}{c_1 c_2} \sum_{x,y} P(X=x)P(Y=y)$$

$$\text{or } 1 = \frac{1}{c_1 \cdot c_2} \cdot 1 \Rightarrow c_1 \cdot c_2 = 1.$$

Definition: Random vectors \underline{X} and \underline{Y} are independent if

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

for all sets A, B .

Remark: The definition is equivalent to

$$P(\underline{x} = \underline{x}, \underline{y} = \underline{y}) = P(\underline{x} = \underline{x}) P(\underline{y} = \underline{y})$$

for all pairs of possible values.

Theorem 3.2 is valid in the following form:

If $P(\underline{x} = \underline{x}, \underline{y} = \underline{y}) = f(\underline{x})g(\underline{y})$ for some functions f, g then $\underline{x}, \underline{y}$ are independent.

3.2. Expected value

Example: In one of on-line games you have 12 tickets

1 1 1 1 2 2 3 10 5 5 5 5

The tickets are turned around and randomly permuted. The player sees

10 10 10 10 10 10 10 10 10 10 10 10

The player then turns tickets from left to right until the ticket

\boxed{S} = STOP appears. Example:

$\boxed{1}$ $\boxed{2}$ \boxed{D} $\boxed{1}$ \boxed{S}

The payout is the sum of all numbers, multiplied by 2 if

\boxed{D} = double appears among the

tickets. In the above example

the payout is 8.

What is the fair price for this game?

Suppose we played this game many times. We can interpret the

payout as a random

variable, X say. Possible

values of X are $\{0, 1, 2, 3, 4, 5, 6, 7,$

$8, 9, 10, 11, 14, 16, 18, 20, 22\}$.

We have denoting possible values of

X by $\{x_1, x_2, \dots, x_{17}\}$:

$$\frac{v_1 + \dots + v_n}{n}$$

$$= \sum_{k=1}^{17} x_k \cdot \underbrace{\frac{\# \text{ of occurrences of } x_k}{n}}_{\approx P(X = x_k)}$$

So the "long term" average will be

$$\sum_{k=1}^{17} x_k P(X = x_k)$$

We will call this average the expected value of a random variable.

Definition : Let X be a discrete random variable with values $\{x_1, x_2, \dots\}$. The expected value

$E(X)$ is defined as

$$E(X) = \sum_{x_k} x_k P(X = x_k)$$

Technical note: We say that X exist if the sum

$$\sum_{x_k} |x_k| \cdot P(X = x_k) \text{ converges.}$$

If f is a function then $Y = f(X)$ is again a discrete random variable. If we "repeat" X we also "repeat" Y . The expectation $E(Y)$ will be approximately

$$\frac{f(x_1) + \dots + f(x_n)}{n} \approx \sum_{x_k} f(x_k) P(X = x_k)$$

by exactly the same argument as before. Formally, we state:

Theorem 3.3: If X is a discrete random variable with values in $\{x_1, x_2, \dots\}$. Let $f: \{x_1, \dots, \dots\} \rightarrow \mathbb{R}$.

We have

$$E[f(X)] = \sum_{x_k} f(x_k) P(X = x_k)$$

Proof: Denote $Y = f(X)$. Possible values are $\{y_1, y_2, \dots\}$. By definition

$$E[f(X)] = E(Y)$$

$$= \sum_{y_e} y_e P(Y = y_e)$$

$$= \sum_{y_e} y_e \sum_{\{x_k: f(x_k) = y_e\}} P(X = x_k)$$

$$= \sum_{y_e} \sum_{\{x_k: f(x_k) = y_e\}} f(x_k) P(X = x_k)$$

$$= \sum_{x_k} f(x_k) P(X = x_k)$$

Technical note: We say that $E(f(X))$

exists if the sum

$$\sum_{x_k} |f(x_k)| P(X = x_k)$$

exists.

Examples :

(i) Let $X \sim \text{Bin}(n, p)$. We compute

$$\begin{aligned} E(X) &= \sum_{k=0}^n k \cdot P(X=k) \\ &= \sum_{k=1}^n k \binom{n}{k} p^k q^{n-k} \quad q = 1-p \\ &= \sum_{k=1}^n n \cdot p \cdot \binom{n-1}{k-1} p^{k-1} q^{(n-1)-(k-1)} \\ &= n \cdot p \underbrace{\sum_{k=1}^n \binom{n-1}{k-1} p^{k-1} q^{(n-1)-(k-1)}}_{= (p+q)^{n-1} = 1} \\ &= n \cdot p \end{aligned}$$

Similarly

$$\begin{aligned} E(X^2) &= \sum_{k=0}^n k^2 \cdot P(X=k) \\ &= \sum_{k=1}^n [k(k-1) + k] \binom{n}{k} p^k q^{n-k} \\ &= \sum_{k=2}^n k(k-1) \binom{n}{k} p^k q^{n-k} \\ &\quad + \underbrace{\sum_{k=1}^n k \binom{n}{k} p^k q^{n-k}}_{= n \cdot p} \end{aligned}$$

$$= \sum_{k=2}^n n(n-1) \binom{n-2}{k-2} \cdot p^2 p^{k-2} q^{(n-2)-(k-2)} + np$$

$$= n(n-1)p^2 \underbrace{\sum_{k=2}^n \binom{n-2}{k-2} p^{k-2} q^{(n-2)-(k-2)}}_{= (p+q)^{n-2} - 1} + np$$

$$= n(n-1)p^2 + np$$

$$= n^2 p^2 + npq$$

(ii) Let $X \sim \text{Neg Bin}(m, p)$.

We have

$$P(X=k) = \binom{k-1}{m-1} p^m q^{k-m}$$

for $k = m, m+1, \dots$. We compute

$$E(X) = \sum_{k=m}^{\infty} k \cdot \binom{k-1}{m-1} p^m \cdot q^{k-m}$$

$$= \sum_{k=m}^{\infty} \binom{(k+1)-1}{(m+1)-1} \cdot m \cdot p^m q^{k-m}$$

$$= \sum_{k=m}^{\infty} \frac{m}{p} \cdot \binom{(k+1)-1}{(m+1)-1} p^{m+1} q^{(k+1)-(m+1)}$$

$$= \frac{m}{p} \cdot \underbrace{\sum_{k=m}^{\infty} \binom{(k+1)-1}{(m+1)-1} p^{m+1} q^{(k+1)-(m+1)}}_{= 1}$$

= 1, because this is the sum of all probabilities in the NegBin($m+1, p$) distribution

$$= \frac{m}{p}$$

In a similar way we find that

$$E(X^2) = \sum_{k=m}^{\infty} k^2 \binom{k-1}{m-1} p^m \cdot q^{k-m}$$

$$= \sum_{k=m}^{\infty} [k(k+1) - k] \binom{k-1}{m-1} p^m q^{k-m}$$

$$= \sum_{k=m}^{\infty} \binom{(k+2)-1}{(m+2)-1} \frac{m(m+1)}{p^2} p^{m+2} q^{k-m}$$

$$- \frac{m}{p}$$

$$= \frac{m(m+1)}{p^2} - \frac{m}{p}$$

$$= \frac{m^2}{p^2} + \frac{m}{p} \left(\frac{1}{p} - 1 \right)$$

$$= \frac{m^2}{p^2} + \frac{m \cdot q}{p^2}$$

(iii) Let $X \sim \text{Hyper Geom}(n, B, N)$.

Let us agree that $\binom{a}{b} = 0$

if $b > a$ or $b < 0$. We compute

$$E(X) = \sum_k k \cdot \frac{\binom{B}{k} \binom{R}{n-k}}{\binom{N}{n}}$$

$$= \sum_k \frac{B \binom{B-1}{k-1} \cdot \binom{R}{(n-1)-(k-1)} \cdot n}{\binom{N-1}{n-1} \cdot N}$$

$$= n \cdot \frac{B}{N} \cdot \underbrace{\sum_k \frac{\binom{B-1}{k-1} \binom{R}{(n-1)-(k-1)}}{\binom{N-1}{n-1}}}_{= 1, \text{ because}}$$

this is the sum of all probs. in

$\text{Hyper Geom}(n-1, B-1, N-1)$
distribution.

$$= n \cdot \frac{B}{N}$$

The most important theoretical property of expectation is linearity.

Theorem 3.4 : Let X, Y be discrete random variables.

We have

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

Proof : Denote $Z = aX + bY$. Z is a discrete random variable with values $\{z_1, z_2, \dots\}$. We have

$$\begin{aligned} E(Z) &= \sum_{z_m} z_m \cdot P(Z = z_m) \\ &= \sum_{z_m} z_m \cdot \sum_{\{(x_k, y_l) : ax_k + by_l = z_m\}} P(X = x_k, Y = y_l) \\ &= \sum_{z_m} \sum_{-} (ax_k + by_l) P(-) \end{aligned}$$

$$= \sum_{x_k, y_l} (ax_k + by_l) P(X=x_k, Y=y_l)$$

$$= a \cdot \sum_{x_k, y_l} x_k P(X=x_k, Y=y_l)$$

$$+ b \cdot \sum_{x_k, y_l} y_l P(X=x_k, Y=y_l)$$

$$= a \cdot \sum_{x_k} x_k P(X=x_k)$$

$$+ b \cdot \sum_{y_l} y_l P(Y=y_l)$$

$$= E(X) + E(Y)$$

Technical note : We assume that

$E(X)$ and $E(Y)$ exist. In this case $E(aX + bY)$ exist as well.

Remark : We have derived that

$$E(X) = \sum_{x_k, y_l} x_k P(X=x_k, Y=y_l)$$

A consequence of Theorem 3.4 is that linearity is valid for more general linear combinations.

If X_1, X_2, \dots, X_v are random variables such that $E(X_k)$ exists then

$$E \left[\sum_{k=1}^v a_k X_k \right] = \sum_{k=1}^v a_k E(X_k).$$

Finally, we state

Theorem 3.5: Let \underline{X} be a discrete random vector in \mathbb{R}^n and let $f: \mathbb{R}^n \rightarrow \mathbb{R}$ be a function. We have

$$E[f(\underline{X})] = \sum_{\underline{x}_k} f(\underline{x}_k) P(\underline{X} = \underline{x}_k)$$

Proof: The proof is identical to the proof of Theorem 3.3.

Example: Let $\underline{X} \sim \text{Multinomial}(n, p)$.

What is $E(X_k \cdot X_l)$? We know that

$X_k + X_l \sim \text{Bin}(n, p_k + p_l)$ so

$$E[(X_k + X_l)^2] = n(p_k + p_l)(1 - p_k - p_l) + n^2(p_k + p_l)^2$$

$$E[X_k^2 + 2X_k X_l + X_l^2]$$

$$E(X_k^2) + 2E(X_k X_l) + E(X_l^2)$$

$$= n p_k (1 - p_k) + n^2 p_k^2$$

$$+ 2 E(X_k X_l)$$

$$+ n p_l (1 - p_l) + n^2 p_l^2$$

This is an equation for $E(X_k X_l)$

from which we compute

$$E(X_k X_l) = -n p_k p_l + n^2 p_k p_l$$

Definition: A random variable X with values in $\{0, 1\}$ is called an indicator or a Bernoulli random variable. We denote $p = P(X=1)$

Shorthand: $X \sim \text{Bernoulli}(p)$.

By definition

$$E(X) = 0 \cdot P(X=0) + 1 \cdot P(X=1) = p$$

Remark: Since $X: \Omega \rightarrow \{0, 1\}$ we can denote $A = \{X=1\}$ which is an event.

Every indicator is associated with an event A . We will write

I_A or 1_A for the indicator of A i.e. the random variable X , for which $X(\omega) = 1$ if $\omega \in A$ and 0 else.

- In many cases complicated random variables can be written as linear combinations of ~~more~~ ~~complicated~~ simpler random variables. Expectations can then be computed in simpler ways using linearity.

Example: Let us return to the first example.

$\overset{1}{\boxed{1}}$ $\overset{2}{\boxed{1}}$ $\overset{3}{\boxed{1}}$ $\overset{4}{\boxed{1}}$ $\overset{1}{\boxed{2}}$ $\overset{2}{\boxed{2}}$ $\boxed{3}$ $\boxed{10}$ $\boxed{15}$ $\boxed{15}$ $\boxed{15}$ $\boxed{15}$

Label the tickets with 1 from 1 to 4, and the tickets with 2 from 1 to 2.

We have

$$X = \sum_{i=1}^4 1_{A_{1,i}} + 2 \sum_{i=1}^4 1_{A_{2,i}} \\ + 2 \sum_{i=1}^2 1_{B_{1,i}} + 4 \sum_{i=1}^2 1_{B_{2,i}} \\ + 3 1_{C_1} + 6 1_{C_2}$$

By symmetry

$$P(A_{1,i}) = P(B_{1,i}) = P(C_1)$$

and

$$P(A_{2,i}) = P(B_{2,i}) = P(C_2).$$

This means that

$$E(X) = 11 \cdot P(A_{1,1}) + 22 \cdot P(A_{2,1})$$

We compute $P(A_{1,1})$ and $P(A_{2,1})$

by noticing that if we only

look at tickets

$\boxed{1}$ $\boxed{10}$ $\boxed{5}$ $\boxed{5}$ $\boxed{5}$ $\boxed{5}$ among the 12

permutated tickets they too are randomly permutated. We say that the induced permutation is random. It follows that $A_{1,1}$ happens if we see

$\boxed{1} \boxed{3} * * * *$

The probability is

$$\frac{1}{6} \times \frac{4}{5} = \frac{2}{15}$$

The event $A_{2,1}$ happens if we see

$\boxed{1} \boxed{2} * * * *$ or $\boxed{2} \boxed{1} * * * *$.

The probability is

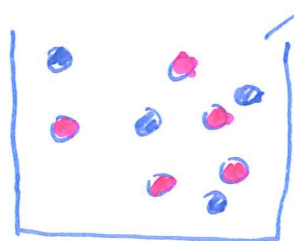
$$2 \cdot \frac{1}{6} \cdot \frac{1}{5} = \frac{1}{15}$$

Finally,

$$E(x) = 11 \cdot \frac{2}{15} + 22 \cdot \frac{1}{15} = \frac{44}{15}$$

$$= 2.93$$

Example: The hyper-geometric distribution is created by selecting balls out of a box.



Select n balls at random
 $X = \#$ of black balls

We can imagine that balls are selected one by one at random until we have n balls. Define

$$I_k = \begin{cases} 1, & \text{if the } k\text{-th ball is black.} \\ 0, & \text{else,} \end{cases}$$

for $k = 1, 2, \dots, n$. We have

$$\begin{aligned} E(X) &= E(\underbrace{I_1 + \dots + I_n}_{= X}) \\ &= E(I_1) + \dots + E(I_n) \end{aligned}$$

$$= P(I_1=1) + P(I_2=1) + \dots + P(I_n=1)$$

But the k -th ball is equally likely to be any of the N balls.

We are asking this question before the selection process begins.

This means that

$$P(I_1=1) = P(I_2=1) = \dots = P(I_n=1).$$

But

$$P(I_1=1) = P(\text{first ball selected is black})$$

$$= \frac{B}{N}.$$

It follows that

$$E(X) = n \cdot \frac{B}{N}.$$

Comment: The idea to write X as a linear combination of indicators is called the method of indicators.

3.3. Joint continuous distributions

For a continuous random variable X with density f_X we have

$$P(a \leq X \leq b) = \int_a^b f_X(x) dx.$$

More generally, for a set A we can say

$$\begin{aligned} P(X \in A) &= \int_A f_X(x) dx \\ &= \int_{-\infty}^{\infty} f_X(x) \chi_A(x) dx, \end{aligned}$$

where χ_A is the characteristic function of the set A . This last form has an easy extension to \mathbb{R}^2 , \mathbb{R}^3 or \mathbb{R}^4 .

For \mathbb{R}^2 we can say that

$$P(\underline{X} \in A) = \iint_A f_{\underline{X}}(x, y) dx dy$$

for an appropriate non-negative function. In \mathbb{R}^3

we have for $\underline{x} = (x_1, x_2, x_3)$

$$P(\underline{x} \in A) = \int\int\int_A f_{\underline{x}}(x_1, x_2, x_3) dx_1 dx_2 dx_3.$$

In probability we will write single integrals even in higher dimensions.

If $A \subseteq \mathbb{R}^n$ we will write

$$\int\int\int\int_A f(x_1, x_2, \dots, x_n) dx_1 dx_2 \dots dx_n$$

$$= \int_A f(\underline{x}) d\underline{x}$$

Definition: If for a random vector

\underline{x} we have

$$P(\underline{x} \in A) = \int_A f_{\underline{x}}(\underline{x}) d\underline{x}$$

for a non-negative function $f_{\underline{x}}: \mathbb{R}^n \rightarrow \mathbb{R}$ and all (reasonable) sets A we say that \underline{x} has continuous distribution with density $f_{\underline{x}}$.

Technical note: In more dimensions in general we say that the distribution of \underline{x} is described by probabilities $P(\underline{x} \in A)$ for all reasonable sets $A \subseteq \mathbb{R}^n$.

"Reasonable" means all sets that are formed from open sets by complements, countable unions and countable intersections. Such sets are called Borel sets.

Example: Let (x, y) be a random vector with density $f_{x,y}(x, y)$ given by

$$f_{x,y}(x, y) = \frac{1}{2\pi\sqrt{1-\rho^2}} e^{-\frac{x^2 - 2\rho xy + y^2}{2(1-\rho^2)}}$$

for $\rho \in (-1, 1)$. Let us check that

$f_{x,y}$ is a density. This means

that it is non-negative and integrates to 1. We know that

$$\frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{\infty} e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx = 1,$$

because the latter is the integral of the normal density. We integrate

$$\int_{\mathbb{R}^2} f_{X,Y}(x,y) dx dy =$$

$$= \frac{1}{2\pi\sqrt{1-\rho^2}} \int_{\mathbb{R}^2} e^{-\frac{x^2-2\rho xy+y^2}{2(1-\rho^2)}} dx dy$$

$$= \frac{1}{2\pi\sqrt{1-\rho^2}} \int_{-\infty}^{\infty} dx \cdot \int_{-\infty}^{\infty} e^{-\frac{x^2-2\rho xy+y^2}{2(1-\rho^2)}} dy$$

This is called Fubini's theorem.

$$= (*)$$

We note

$$\begin{aligned} & (x^2 - 2\rho xy + y^2) / (1 - \rho^2) \\ &= [(y - \rho x)^2 + (1 - \rho^2)x^2] / (1 - \rho^2) \\ &= \frac{(y - \rho x)^2}{1 - \rho^2} + x^2 \end{aligned}$$

$$\begin{aligned} (*) &= \frac{1}{2\pi\sqrt{1-\rho^2}} \int_{-\infty}^{\infty} e^{-\frac{x^2}{2}} dx \\ &\quad \cdot \underbrace{\int_{-\infty}^{\infty} e^{-\frac{(y-\rho x)^2}{2(1-\rho^2)}} dy}_{= \sqrt{2\pi} \cdot \sqrt{1-\rho^2}} \end{aligned}$$

$$= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-x^2/2} dx$$

$$= 1.$$

Let us continue with this example and compute $P(x \geq 0, y \geq 0)$.

In other words

$$P(x \geq 0, Y \geq 0) = P((x, Y) \in [0, \infty)^2)$$

By definition

$$P((x, Y) \in [0, \infty)^2)$$

$$= \int_{[0, \infty)^2} f_{x, Y}(x, y) dx dy$$

$$= \frac{1}{2\pi\sqrt{1-\rho^2}} \int_{[0, \infty)^2} e^{-\frac{x^2 - 2\rho xy + y^2}{2(1-\rho^2)}} dx dy$$

$$= \frac{1}{2\pi\sqrt{1-\rho^2}} \int_0^{\infty} e^{-x^2/2} dx \cdot \underbrace{\int_0^{\infty} e^{-\frac{(y-\rho x)^2}{2(1-\rho^2)}} dy}_{\text{New variable:}}$$

New variable:

$$\frac{y-\rho x}{\sqrt{1-\rho^2}} = u$$

$$= \frac{1}{2\pi\sqrt{1-\rho^2}} \int_0^{\infty} e^{-x^2/2} dx \cdot$$

$$\int_{-\frac{\rho x}{\sqrt{1-\rho^2}}}^{\infty} e^{-u^2/2} du \cdot \sqrt{1-\rho^2}$$

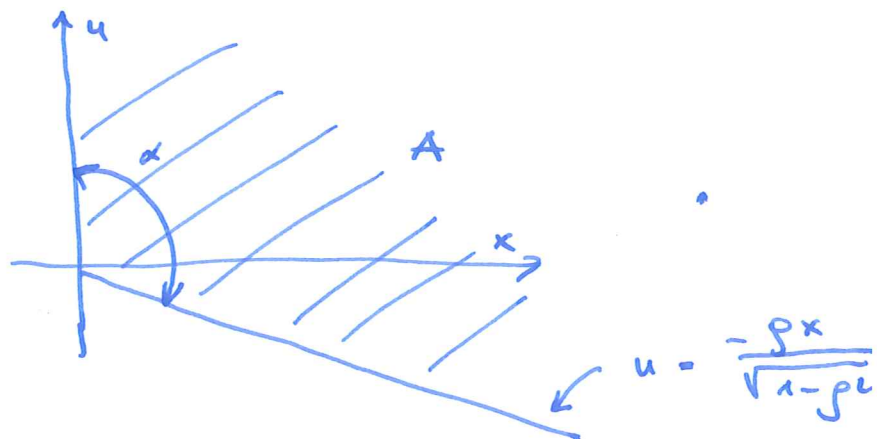
This last integral is the integral of the function

$$f(x, u) = \frac{1}{2\pi} e^{-\frac{x^2 + u^2}{2}} \quad \text{over the set}$$

$$A = \left\{ (x, u) : x \geq 0, u \geq -\frac{\rho x}{\sqrt{1-\rho^2}} \right\}$$

by Fubini.

Figure :



We observe :

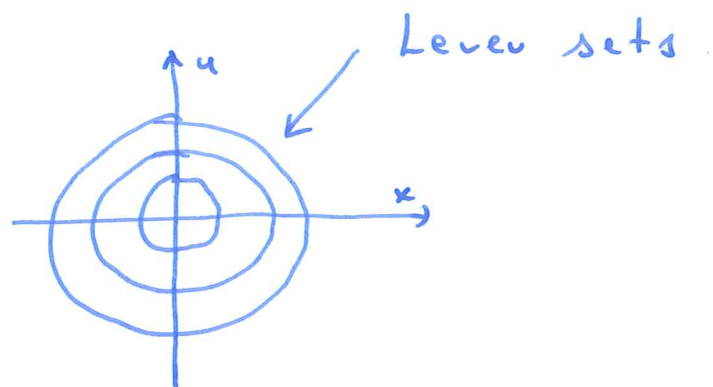
(i) $f(x, u)$ integrates to 1.

We get this by taking $\rho = 0$ in the previous example.

(ii) The function

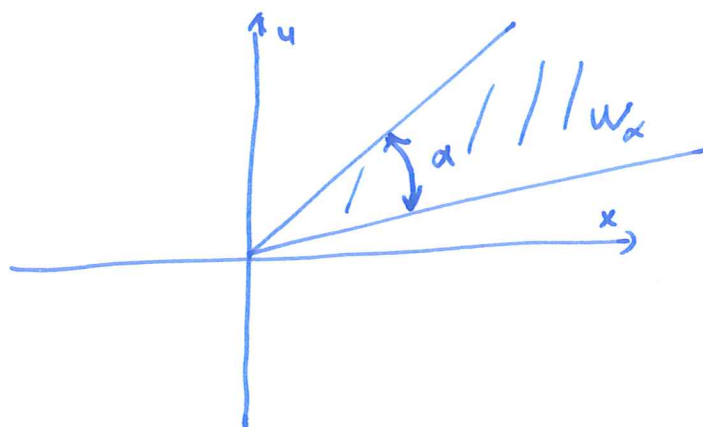
$f(x, u)$ is rotationally symmetric.

Figure :



The integral over a wedge of angle α is proportional to the angle

Figure :



$$\int_{W_\alpha} f(x, u) dx du = \frac{\alpha}{2\pi}$$

In our case the angle α equals

$$\alpha = \frac{\pi}{2} + \arctan\left(\frac{\rho}{\sqrt{1-\rho^2}}\right).$$

Finally we have

$$P(X \geq 0, Y \geq 0) =$$

$$= \frac{1}{4} + \frac{1}{2\pi} \arctan\left(\frac{\rho}{\sqrt{1-\rho^2}}\right)$$

Marginal distributions

Suppose (X, Y) has density $f_{X,Y}(x, y)$.

What is the density of X ?

We know from the ~~1st~~ 2nd Chapter that if for any $a < b$

$$P(a \leq X \leq b) = \int_a^b g(x) dx \quad \text{+ this}$$

implies that $g(x) = f_X(x)$.

we compute

$$P(a \leq x \leq b) = P(a \leq x \leq b, Y \in \mathbb{R})$$

$$= P((x, Y) \in [a, b] \times \mathbb{R})$$

$$= \int_{[a, b] \times \mathbb{R}} f_{X, Y}(x, y) dx dy$$

$$= \int_a^b dx \cdot \underbrace{\int_{-\infty}^{\infty} f_{X, Y}(x, y) dy}$$

This is a function
of x , say $g(x)$

$$= \int_a^b g(x) dx.$$

Theorem 3.6: Let (X, Y) have
the density $f_{X, Y}(x, y)$. We have

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

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

Proof: Done already.

Comments

- (i) The two formulae in Theorem 3.6 are called formulae for marginal densities.
- (ii) A rigorous statement must include the assumptions that $x \mapsto f_{x,y}(x,y)$ and $y \mapsto f_{x,y}(x,y)$ are Riemann integrable, and that f_x and f_y are Riemann integrable.

Example :

$$f_{x,y}(x,y) = \frac{1}{2\pi\sqrt{1-\rho^2}} e^{-\frac{x^2 - 2xy + y^2}{2(1-\rho^2)}}$$

We have

$$f_x(x) = \int_{-\infty}^{\infty} f_{x,y}(x,y) dy.$$

$$= \int_{-\infty}^{\infty} \frac{1}{2\pi\sqrt{1-\rho^2}} e^{-\frac{x^2}{2}} \cdot e^{-\frac{(y-\rho x)^2}{2(1-\rho^2)}} dy$$

$$= \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$$

$$\cdot \underbrace{\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}\sqrt{1-\rho^2}} e^{-\frac{(y-\rho x)^2}{2(1-\rho^2)}} dy}_{= 1, \text{ because it is the integral of the } N(\rho x, 1-\rho^2) \text{ dist}}$$

$$= \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{x^2}{2}}.$$

Conclusion : $X \sim N(0, 1)$.

Theorem 3.6 has a more general version.

Theorem 3.6a: Let $(\underline{x}, \underline{y})$ be a random vector with $\underline{x} \in \mathbb{R}^b$ and $\underline{y} \in \mathbb{R}^2$ with density $f_{\underline{x}, \underline{y}}(\underline{x}, \underline{y})$. Then

$$f_{\underline{x}}(\underline{x}) = \int_{\mathbb{R}^2} f_{\underline{x}, \underline{y}}(\underline{x}, \underline{y}) d\underline{y}$$

$$f_{\underline{y}}(\underline{y}) = \int_{\mathbb{R}^b} f_{\underline{x}, \underline{y}}(\underline{x}, \underline{y}) d\underline{x}$$

Proof: Same as before.

Example: If (x, y, z) has density $f_{x, y, z}(x, y, z)$ then

$$f_{x, y}(x, y) = \int_{-\infty}^{\infty} f_{x, y, z}(x, y, z) dz$$

and

$$f_x(x) = \int_{\mathbb{R}^2} f_{x, y, z}(x, y, z) dy dz.$$

Independence

In general we say that x, y are independent if

$$P(x \in A, y \in B) = P(x \in A) \cdot P(y \in B).$$

If (x, y) has density $f_{x, y}(x, y)$

this means that for $A = [a, b]$

and $B = [c, d]$ we have

$$\begin{aligned} & \int_{[a, b] \times [c, d]} f_{x, y}(x, y) dx dy \\ & \quad = P(x \in [a, b], y \in [c, d]) \\ & \quad = \left(\int_a^b f_x(x) dx \right) \cdot \left(\int_c^d f_y(y) dy \right) \\ & \quad \quad \quad \underbrace{\hspace{10em}}_{P(a \leq x \leq b)} \quad \quad \quad \underbrace{\hspace{10em}}_{P(c \leq y \leq d)} \end{aligned}$$

Fubini.

$$= \int_{[a, b] \times [c, d]} f_x(x) \cdot f_y(y) dx dy$$

We borrow a statement from
Analysis 2: if for functions
 $f(x,y)$ and $g(x,y)$ we have:

(i) For all rectangles $Q = [a,b] \times [c,d]$
the functions are Riemann
integrable.

(ii)

$$\int_Q f(x,y) dx dy = \int_Q g(x,y) dx dy$$

for all Q

then $f(x,y) = g(x,y)$.

Technical note: in fact f, g
can differ but only on a set
of measure 0.

If x, y are independent we
have

$$\int_Q f_{x,y}(x,y) dx dy = \int_Q f_x(x) f_y(y) dx dy$$

Theorem 3.7 : Let (X, Y) have density $f_{X, Y}(x, y)$. The random variables X and Y are independent if and only if $f_{X, Y}(x, y) = f_X(x) f_Y(y)$.

Proof : If $f_{X, Y}(x, y) = f_X(x) \cdot f_Y(y)$

then

$$\underbrace{P((X, Y) \in A \times B)}_{\int_{A \times B} f_{X, Y}(x, y) dx dy} = P(X \in A, Y \in B)$$

Fubini

$$= \underbrace{\int_A f_X(x) dx}_{P(X \in A)} \cdot \underbrace{\int_B f_Y(y) dy}_{P(Y \in B)}$$

Independence follows.

If X, Y are independent we proved above that $f_{X, Y}(x, y) = f_X(x) \cdot f_Y(y)$.

Theorem 3.7 has a more general version.

Theorem 3.7a : Let $\underline{X}, \underline{Y}$ be continuous random vectors with density $f_{\underline{X}, \underline{Y}}(\underline{x}, \underline{y})$. The vectors \underline{X} and \underline{Y} are independent if

and only if $f_{\underline{X}, \underline{Y}}(\underline{x}, \underline{y}) = f_{\underline{X}}(\underline{x}) \cdot f_{\underline{Y}}(\underline{y})$.

Proof : Same as above.

Theorem 3.8 : Let (X, Y) have density $f_{X, Y}(x, y)$. If

$$f_{X, Y}(x, y) = \cancel{g(x)} \cancel{h(y)} g(x) h(y)$$

for nonnegative functions g and h then X, Y are independent.

Proof: By the formulae for marginal density we have

$$\begin{aligned}f_x(x) &= \int_{-\infty}^{\infty} f_{x,y}(x,y) dy \\&= \int_{-\infty}^{\infty} g(x) h(y) dy \\&= g(x) \cdot \underbrace{\int_{-\infty}^{\infty} h(y) dy}_{= c_1}\end{aligned}$$

Similarly

$$f_y(y) = h(y) \cdot \underbrace{\int_{-\infty}^{\infty} g(x) dx}_{c_2}$$

It follows

$$f_{x,y}(x,y) = \frac{f_x(x)}{c_1} \cdot \frac{f_y(y)}{c_2}$$

We need to prove that $c_1 \cdot c_2 = 1$.

Integrate both sides over \mathbb{R}^2 .

We get

$$1 = \int_{\mathbb{R}^2} f_{x,y}(x,y) dx dy$$

$$= \int_{\mathbb{R}^2} \frac{1}{c_1 c_2} f_x(x) f_y(y) dx dy$$

$$= \frac{1}{c_1 c_2} \int_{\mathbb{R}^2} f_x(x) f_y(y) dx dy$$

Fubini

$$= \frac{1}{c_1 c_2} \int_{-\infty}^{\infty} f_x(x) dx \cdot \int_{-\infty}^{\infty} f_y(y) dy$$

$$= \frac{1}{c_1 \cdot c_2} \cdot 1 \cdot 1$$

It follows $c_1 \cdot c_2 = 1$.

Example 1 Let

$$f_{x,y}(x,y) = \frac{1}{2\pi \sqrt{1-\rho^2}} e^{-\frac{x^2 - 2\rho xy + y^2}{2(1-\rho^2)}}$$

We computed

$$f_X(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2} \quad \text{and}$$

$$f_Y(y) = \frac{1}{\sqrt{2\pi}} e^{-y^2/2}$$

We see that

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

only when $\rho = 0$.

3.4. Functions of random vectors

Discrete case

In the discrete case we will only look at integer valued random variables. If X, Y are two such variables then $Z = X + Y$ is an integer valued

random variable. We have

$$\{Z = u\} = \underbrace{\bigcup_{k \in \mathbb{Z}} \{X = k, Y = u - k\}}_{\text{disjoint union}}$$

We have

$$P(Z = u) = \sum_{k \in \mathbb{Z}} P(X = k, Y = u - k)$$

Special cases:

(i) if X, Y are non-negative we have

$$P(Z = u) = \sum_{k=0}^u P(X = k, Y = u - k)$$

(ii) if X, Y are independent then

$$P(Z = u) = \sum_{k \in \mathbb{Z}} P(X = k) P(Y = u - k)$$

Examples: (i) Let X, Y be independent and $X \sim P_0(\mu), Y \sim P_0(\lambda)$. Let $Z = X + Y$. By the above formula

$$P(Z = u) = \sum_{k=0}^u P(X=k) \cdot P(Y=u-k)$$

$$= \sum_{k=0}^u \frac{e^{-\mu} \mu^k}{k!} \cdot \frac{e^{-\lambda} \lambda^{u-k}}{(u-k)!}$$

$$= \frac{e^{-(\lambda+\mu)}}{u!} \sum_{k=0}^u \underbrace{\frac{u!}{k!(u-k)!}}_{= \binom{u}{k}} \mu^k \lambda^{u-k}$$

$$= \frac{e^{-(\lambda+\mu)}}{u!} (\lambda + \mu)^u.$$

Conclusion: $Z = X + Y \sim P_0(\lambda + \mu)$.

(ii) Let X, Y be independent and have the Polya distribution.

This means that

$$P(X=k) = \frac{\beta^a (a)_k}{k! (1+\beta)^{a+k}} \quad k=0,1,\dots$$

$$P(Y=l) = \frac{\beta^b (b)_l}{l! (1+\beta)^{b+l}} \quad l=0,1,\dots$$

Here $(a)_0 = 1$ and

$$(a)_k = a(a+1)\dots(a+k-1)$$

is the Pochhammer symbol.

Let $Z = X + Y$. We are looking for the distribution of Z . By the formula we have

$$\begin{aligned} P(Z=n) &= \sum_{k=0}^n P(X=k)P(Y=n-k) \\ &= \frac{\beta^{a+b}}{(1+\beta)^{a+b+n}} \sum_{k=0}^n \frac{(a)_k (b)_{n-k}}{k! (n-k)!} \\ &= \frac{\beta^{a+b}}{(1+\beta)^{a+b+n} \cdot n!} \sum_{k=0}^n \binom{n}{k} (a)_k (b)_{n-k}. \end{aligned}$$

The last formula is similar to the binomial formula.

To prove it we will use a few facts from Analysis:

(i) The gamma function is defined as

$$\Gamma(x) = \int_0^{\infty} u^{x-1} \cdot e^{-u} du, \quad x > 0$$

Integration by parts gives

$$\Gamma(x+1) = x \Gamma(x) \quad \text{and}$$

as a consequence

$$\Gamma(a+n) = (a+n-1)(a+n-2) \cdots a \cdot \Gamma(a)$$

We can write

$$(a)_n = \frac{\Gamma(a+n)}{\Gamma(a)}$$

(ii) The Beta function is defined as

$$B(p, q) = \int_0^1 u^{p-1} (1-u)^{q-1} du, \quad p, q > 0$$

The connection between Γ and B functions is given by Euler:

$$B(p, q) = \frac{\Gamma(p) \cdot \Gamma(q)}{\Gamma(p+q)}$$

We compute

$$\sum_{k=0}^n \binom{n}{k} (a)_k (b)_{n-k}$$

$$= \sum_{k=0}^n \binom{n}{k} \frac{\Gamma(a+k)}{\Gamma(a)} \cdot \frac{\Gamma(b+n-k)}{\Gamma(b)}$$

$$= \frac{\Gamma(a+b+n)}{\Gamma(a)\Gamma(b)} \sum_{k=0}^n \binom{n}{k} \frac{\Gamma(a+k)\Gamma(b+n-k)}{\Gamma(a+b+n)}$$

$$= \frac{\Gamma(a+b+n)}{\Gamma(a)\Gamma(b)} \sum_{k=0}^n \binom{n}{k} B(a+k, b+n-k)$$

$$= \frac{\Gamma(a+b+u)}{\Gamma(a)\Gamma(b)} \sum_{k=0}^n \binom{n}{k} \int_0^1 u^{a+k-1} (1-u)^{b+n-k-1} du$$

$$= \frac{\Gamma(a+b+u)}{\Gamma(a)\Gamma(b)} \int_0^1 u^{a-1} (1-u)^{b-1} \underbrace{\sum_{k=0}^n \binom{n}{k} u^k (1-u)^{n-k}}_{=1} du$$

def.

$$= \frac{\Gamma(a+b+u)}{\Gamma(a)\Gamma(b)} B(a, b)$$

Euler

$$= \frac{\Gamma(a+b+u)}{\Gamma(a)\Gamma(b)} \cdot \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$$

$$= \frac{\Gamma(a+b+u)}{\Gamma(a+b)}$$

$$= (a+b)_n$$

Finally we have

$$P(Z=u) = \frac{\beta^{a+b} (a+b)_n}{n! (1+\beta)^{a+b+n}}, \quad n=0, 1, \dots$$

Example: Suppose X, Y are independent and $X \sim \text{Bin}(m, p)$, $Y \sim \text{Bin}(n, p)$. We expect

$Z = X + Y \sim \text{Bin}(m+n, p)$. The

formal proof:

$$\begin{aligned}
 P(Z=l) &= \sum_{k=0}^l P(X=k, Y=l-k) \\
 &= \sum_{k=\max(0, n-l)}^{\min(l, m)} \binom{m}{k} p^k q^{m-k} \binom{n}{l-k} p^{l-k} q^{n-l+k} \\
 &= \sum_{k=\max(0, n-l)}^{\min(l, m)} \binom{m}{k} \binom{n}{l-k} \underbrace{p^l q^{m+n-l}}_{\text{does not depend on } k}
 \end{aligned}$$

The sum is computed by the following combinatorial argument:

Suppose we need to choose l elements from the union of sets with m and n elements.

This can be done in $\binom{m+n}{l}$ ways. We can count in another way: we choose k elements

from the first set and $l-k$ from the other. This is possible for $k \geq \max(0, n-l)$ and $l \leq \min(l, m)$. This splits all the choices in disjoint subsets so

$$\binom{m+n}{l} = \sum_{k=\max(0, n-l)}^{\min(l, m)} \binom{m}{k} \binom{n}{l-k}.$$

Finally

$$P(Z=l) = \binom{m+n}{l} p^l q^{m+n-l}$$

Continuous case

The most important formula is the transformation formula.

Suppose the vector (x, y)

has density $f_{x,y}(x, y)$. We

form a new vector (u, v) by

$$\Phi(x, y) = (\Phi_1(x, y), \Phi_2(x, y)).$$

Example:

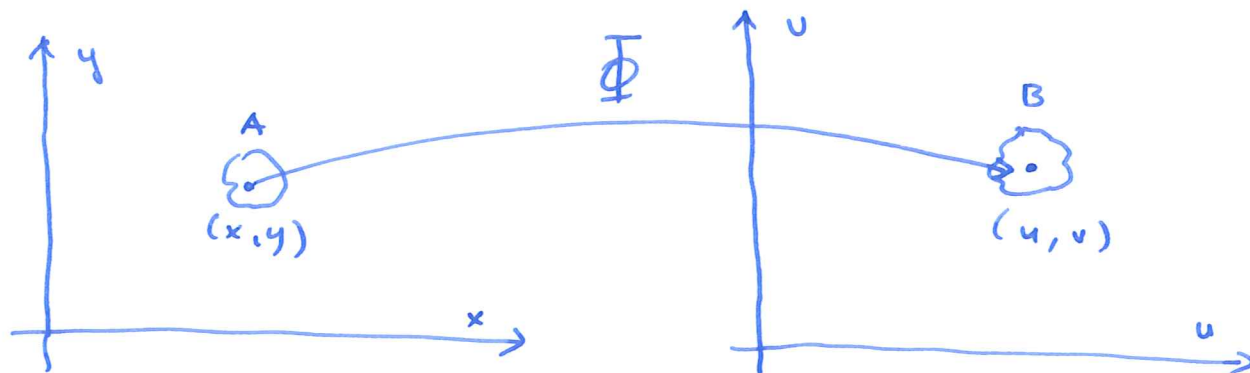
$$\Phi(x, y) = \left(\frac{x}{x+y}, x+y \right)$$

$$(u, v) = \left(\frac{x}{x+y}, x+y \right).$$

Question: What is the density

$f_{u,v}(u, v)$ of (u, v) ?

Idea :



By definition for "small" A and B

$$P((x, y) \in A) \approx f_{x, y}(x, y) |A|$$

$$P((u, v) \in B) \approx f_{u, v}(u, v) \cdot |B|$$

If Φ is bijective then

$$P((x, y) \in A) = P((u, v) \in B)$$

if $B = \Phi(A)$. So

$$f_{x, y}(x, y) |A| = f_{u, v}(u, v) \cdot |B|$$

or

$$f_{u, v}(u, v) \approx f_{x, y}(x, y) \cdot \frac{|A|}{|B|}$$

But from Analysis 2 we know

$$\frac{|A|}{|B|} \approx |\mathcal{J}_{\Phi^{-1}}(u,v)|.$$

Theorem 3.9 (transformation formula). Let (x,y) be

○ a vector with density $f_{x,y}(x,y)$.

Suppose $P((x,y) \in \mathcal{U}) = 1$ for

an open set \mathcal{U} . Let

○ $\Phi: \mathcal{U} \rightarrow \mathcal{I}$ be a bijective

map which is continuously

partially differentiable. Let

$$(u,v) = \Phi(x,y).$$

The the density $f_{u,v}(u,v)$ is

$$f_{u,v}(u,v) = f_{x,y}(\Phi^{-1}(u,v))$$

$$\cdot |\mathcal{J}_{\Phi^{-1}}(u,v)|$$

where $J\Phi^{-1}$ is the Jacobian
determinant of Φ^{-1} .

Proof: Let $B \subseteq \mathcal{P}$. We compute

$$P((U, V) \in B)$$

$$= P((X, Y) \in \Phi^{-1}(B))$$

$$= \int_{\Phi^{-1}(B)} f_{X, Y}(x, y) dx dy$$

$$= (*)$$

New variable: $(x, y) = \Phi^{-1}(u, v)$

$$dx dy = |J\Phi^{-1}(u, v)| du dv$$

$$(*) = \int_B f_{X, Y}(\Phi^{-1}(u, v)) |J\Phi^{-1}(u, v)| du dv.$$

Comment: We used the formula for a new variable in double integrals.

Example: Let X, Y be independent with $X \sim P(a, \lambda)$

and $Y \sim P(b, \lambda)$. This means

$$f_X(x) = \frac{\lambda^a}{\Gamma(a)} x^{a-1} \cdot e^{-\lambda x}, \quad x > 0$$

$$f_Y(y) = \frac{\lambda^b}{\Gamma(b)} y^{b-1} e^{-\lambda y}, \quad y > 0$$

By independence

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

let

$$\Phi(x,y) = \left(\frac{x}{x+y}, x+y \right)$$

for $x, y > 0$. We can take

$$U = (0, \infty)^2 \quad \text{and} \quad J = (0, 1) \times (0, \infty).$$

Φ is bijective and continuously differentiable. To find Φ^{-1} we need to solve equations

$$\frac{x}{x+y} = u, \quad x+y = v.$$

We get

$$x = u \cdot v$$

$$y = v - x = v - u \cdot v$$

$$= v(1-u)$$

This means

$$\Phi^{-1}(u, v) = (uv, v(1-u)),$$

We compute

$$\begin{aligned} J_{\Phi^{-1}}(u, v) &= \det \begin{pmatrix} v & u \\ -v & 1-u \end{pmatrix} \\ &= v \end{aligned}$$

The density $f_{u,v}(u, v)$ is given by

$$f_{u,v}(u, v) = f_{x,y}(uv, v(1-u)) \cdot$$

$$|J_{\Phi^{-1}}(u, v)| =$$

$$= f_x(uv) f_y(v(1-u)) \cdot v$$

$$= \frac{\lambda^a}{\Gamma(a)} (uv)^{a-1} \cdot e^{-\lambda uv} \cdot \frac{\lambda^b}{\Gamma(b)} [v(1-u)]^{b-1} \cdot e^{-\lambda v(1-u)}$$

$\cdot v$

$$= \frac{\lambda^{a+b}}{\Gamma(a)\Gamma(b)} u^{a-1} (1-u)^{b-1} \cdot v^{a+b-1} \cdot e^{-\lambda v}$$

for $(u, v) \in (0, 1) \times (0, \infty)$.

We note:

(i) u, v are independent

(ii) $f_u(u) = \text{const.} \cdot u^{a-1} (1-u)^{b-1}$

$$f_v(v) = \text{const.} \cdot v^{a+b-1} \cdot e^{-\lambda v}$$

It follows $u = \frac{x}{x+y} \sim \text{Beta}(a, b)$

and $V = x+y \sim \Gamma(a+b, \lambda)$.

Example : Suppose (x, y)
has density $f_{x, y}(x, y)$. Let

$$\Phi(x, y) = (x, x+y) = (x, z)$$

What is the density $f_{x, z}(x, z)$?

By the transformation formula

$$f_{x, z}(x, z) = f_{x, y}(x, z-x) \cdot |J_{\Phi^{-1}}(x, z)|$$

But $\Phi^{-1}(x, z) = (x, z-x) \Rightarrow$

$$|J_{\Phi^{-1}}(x, z)| = 1.$$

We have

$$f_{x, z}(x, z) = f_{x, y}(x, z-x)$$

The density of z is the marginal density of $f_{X,Z}(x,z)$.

We have

$$f_z(z) = \int_{-\infty}^{\infty} f_{X,Y}(x, z-x) dx$$

If X, Y are independent we get

$$f_z(z) = \int_{-\infty}^{\infty} f_X(x) f_Y(z-x) dx$$

Comment: The above formula is known as convolution in Analysis.

Example : Let X, Y be independent with $X \sim N(\mu, \sigma^2)$ and $Y \sim N(\nu, \tau^2)$. What is the

density of $Z = X + Y$. Assume

first $\mu = \nu = 0$ and $\sigma^2 + \tau^2 = 1$.

In this case

$$f_Z(z) = \int_{-\infty}^{\infty} f_X(x) f_Y(z-x) dx$$

$$= \frac{1}{2\pi \cdot \sigma \tau} \int_{-\infty}^{\infty} e^{-\frac{x^2}{2\sigma^2}} \cdot e^{-\frac{(z-x)^2}{2\tau^2}} dx$$

$$= \frac{1}{2\pi \sigma \tau} \int_{-\infty}^{\infty} e^{-x^2 \left[\frac{1}{2\sigma^2} + \frac{1}{2\tau^2} \right]}$$

$$\cdot e^{\frac{xz}{\tau^2}} \cdot e^{-\frac{z^2}{2\tau^2}} dx$$

$$= \frac{1}{2\pi \sigma \tau} \int_{-\infty}^{\infty} e^{-\frac{x^2}{2\sigma^2 \tau^2} + \frac{xz}{\tau^2}} \cdot e^{-\frac{z^2}{2\tau^2}} dx$$

$$= \frac{1}{2\pi\sigma\tau} \int_{-\infty}^{\infty} e^{-\frac{1}{2\sigma^2\tau^2} (x - z\sigma^2)^2} \cdot e^{-\frac{z^2\sigma^2}{2\tau^2}} dx$$

$$= \underbrace{\frac{1}{\sqrt{2\pi\sigma\tau}} \int_{-\infty}^{\infty} e^{-\frac{(x - z\sigma^2)^2}{2\sigma^2\tau^2}} dx}_{= 1}$$

$$\cdot \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2} \left[\frac{1}{\tau^2} - \frac{\sigma^2}{\tau^2} \right]} = 1$$

$$= \frac{1}{\sqrt{2\pi}} e^{-z^2/2}$$

Conclusion: $Z \sim N(0, 1)$.

We know: if $X \sim N(\mu, \sigma^2)$ then $aX + b \sim N(a\mu + b, a^2\sigma^2)$.

In general: $X \sim N(\mu, \sigma^2)$, $Y \sim N(\nu, \tau^2)$

$$X + Y = \sqrt{\sigma^2 + \tau^2} \cdot X'$$

$$\left(\underbrace{\frac{X - \mu}{\sqrt{\sigma^2 + \tau^2}}}_{X'} + \underbrace{\frac{Y - \nu}{\sqrt{\sigma^2 + \tau^2}}}_{Y'} \right) + \mu + \nu$$

We have $X' \sim N(0, \frac{\sigma^2}{\sigma^2 + \tau^2})$ in

$Y' \sim N(0, \frac{\tau^2}{\sigma^2 + \tau^2})$. The

expression $X' + Y' \sim N(0, 1)$.

It follows that

$$X + Y \sim N(\mu + \nu, \sigma^2 + \tau^2)$$

Example: Let X, Y be

independent standard normal.

Let $Z = \frac{Y}{X}$. Density of Z ?

Define

$$\Phi(x, y) = (x, \frac{y}{x})$$

$$\Phi^{-1}(x, z) = (x, xz) \Rightarrow$$

$$J_{\Phi^{-1}}(x, z) = \det \begin{pmatrix} 1 & 0 \\ z & x \end{pmatrix} = x$$

The density of (x, z) is

$$\begin{aligned} f_{x,z}(x, z) &= f_x(x) f_y(xz) \cdot |x| \\ &= \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{x^2}{2}} \cdot \frac{1}{\sqrt{2\pi}} e^{-\frac{(xz)^2}{2}} \cdot |x| \end{aligned}$$

We get the density of z as the
○ marginal density

$$f_z(z) = \int_{-\infty}^{\infty} \frac{1}{2\pi} \cdot e^{-\frac{x^2(1+z^2)}{2}} \cdot |x| dx$$

$$= \frac{1}{\pi} \int_0^{\infty} e^{-\frac{x^2(1+z^2)}{2}} x \cdot dx$$

$$= \frac{1}{\pi(1+z^2)} \left(-e^{-\frac{x^2(1+z^2)}{2}} \right) \Big|_0^{\infty}$$

$$= \frac{1}{\pi(1+z^2)}$$

Example: Let X, Y be independent with $X \sim P(a, \lambda)$ and $Y \sim P(b, \lambda)$. Let $Z = X + Y$. We established that $Z \sim P(a+b, \lambda)$ but will ~~do~~ it again using convolution.

$$\begin{aligned}
 f_Z(z) &= \int_0^z f_X(x) f_Y(z-x) dx \\
 &= \int_0^z \frac{\lambda^a}{\Gamma(a)} x^{a-1} e^{-\lambda x} \\
 &\quad \cdot \frac{\lambda^b}{\Gamma(b)} (z-x)^{b-1} e^{-\lambda(z-x)} dx \\
 &= \frac{\lambda^{a+b}}{\Gamma(a) \Gamma(b)} \cdot e^{-\lambda z} \\
 &\quad \cdot \int_0^z x^{a-1} (z-x)^{b-1} dx
 \end{aligned}$$

New variable: $x = z \cdot u$
 $dx = z \cdot du$

$$= \frac{\lambda^{a+b}}{\Gamma(a) \Gamma(b)} \cdot e^{-z} \cdot \int_0^1 u^{a-1} \cdot (1-u)^{b-1} \cdot z^{a+b-1} du$$

$$= \frac{\lambda^{a+b}}{\Gamma(a)\Gamma(b)} e^{-z} \cdot z^{a+b-1} B(a, b)$$

The result is a density which means that it integrates to 1.

But we know that

$$\frac{\lambda^{a+b}}{\Gamma(a+b)} \int_0^{\infty} z^{a+b-1} e^{-\lambda z} dz = 1.$$

This means that

$$\frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \cdot B(a, b) = 1$$

which is Euler's identity!

We have used probability to derive Euler's identity.

Theorem 3.6 has a more general version.

Theorem 3.6 a: Let \underline{x} be a random vector with density

$f_{\underline{x}}(\underline{x})$. Assume $P(\underline{x} \in \mathcal{U}) = 1$

for some open set $\mathcal{U} \subseteq \mathbb{R}^n$ and

let $\Phi: \mathcal{U} \rightarrow \mathcal{S} \subseteq \mathbb{R}^n$ be a bijective

map between \mathcal{U} and \mathcal{S} such

that Φ and Φ^{-1} are continuously

partially differentiable. Let

$\underline{y} = \Phi(\underline{x})$. Then \underline{y} has

the density

$$f_{\underline{y}}(\underline{y}) = f_{\underline{x}}(\Phi^{-1}(\underline{y})) \cdot |\mathcal{J}\Phi^{-1}(\underline{y})|.$$

Proof: Same as before.

Example: Let $\underline{x} = (x_1, x_2, \dots, x_r)$ such that x_1, x_2, \dots, x_r are independent and $x_k \sim N(0, 1)$ for all $k = 1, 2, \dots, r$; Let \underline{A} be an invertible matrix. Define

$$\underline{\Phi}(\underline{x}) = \underline{A}\underline{x} + \underline{\mu} \quad \text{for } \underline{\mu} \in \mathbb{R}^n.$$

We have $\underline{\Phi}^{-1}(\underline{y}) = \underline{A}^{-1}(\underline{y} - \underline{\mu})$

and $J_{\underline{\Phi}^{-1}}(\underline{y}) = \det(\underline{A}^{-1}) = \frac{1}{\det(\underline{A})}$

The transformation formula gives for $\underline{y} = \underline{A}\underline{x} + \underline{\mu}$

$$f_{\underline{y}}(\underline{y}) = f_{\underline{x}}(\underline{\Phi}^{-1}(\underline{y})) \cdot |J_{\underline{\Phi}^{-1}}(\underline{y})|.$$

We have

$$f_{\underline{x}}(\underline{x}) = \prod_{k=1}^r f_{x_k}(x_k)$$

$$= \frac{1}{(2\pi)^{v/2}} \cdot e^{-\frac{1}{2} \sum_{k=1}^v x_k^2}$$

$$= \frac{1}{(2\pi)^{v/2}} \cdot e^{-\frac{1}{2} \underline{x}^T \underline{x}}$$

It follows

$$f_{\underline{y}}(\underline{y}) = \frac{1}{(2\pi)^{v/2} |\det(\underline{A})|}$$

$$\times e^{-\frac{1}{2} [\underline{A}^{-1}(\underline{x} - \underline{\mu})]^T [\underline{A}^{-1}(\underline{x} - \underline{\mu})]}$$

$$= \frac{1}{(2\pi)^{v/2} |\det(\underline{A})|}$$

$$\times e^{-\frac{1}{2} (\underline{x} - \underline{\mu})^T (\underline{A}^{-1})^T \underline{A}^{-1} (\underline{x} - \underline{\mu})}$$

Denote $\underline{\Sigma} = \underline{A} \cdot \underline{A}^T$. We have

$$\underline{\Sigma}^{-1} = (\underline{A}^T)^{-1} \cdot \underline{A}^{-1} = (\underline{A}^{-1})^T (\underline{A}^{-1})$$

We have

$$\frac{1}{|\det(\underline{A})|} = \frac{1}{\sqrt{\det(\underline{\Sigma})}}$$

and

$$f_{\underline{y}}(\underline{y}) = \frac{1}{(2\pi)^{n/2} \sqrt{\det(\underline{\Sigma})}}$$

$$\times e^{-\frac{1}{2}(\underline{x} - \underline{\mu})^T \underline{\Sigma}^{-1}(\underline{x} - \underline{\mu})}.$$

Comment: The above density is called the multivariate normal density with parameters $\underline{\mu} \in \mathbb{R}^n$ and $\underline{\Sigma}$ ($r \times r$).

Example : Let \underline{x} have density

$$f_{\underline{x}}(\underline{x}) = \frac{1}{(2\pi)^{n/2} \sqrt{\det \underline{\Sigma}}} e^{-\frac{1}{2} (\underline{x} - \underline{\mu})^T \underline{\Sigma}^{-1} (\underline{x} - \underline{\mu})}$$

If $\underline{x} = (x_1, x_2, \dots, x_n)$ what is

the distribution of $\underline{x}^{(1)} = (x_1, x_2, \dots, x_p)$,

$p < n$? Denote $\underline{x} = \begin{pmatrix} \underline{x}^{(1)} \\ \underline{x}^{(2)} \end{pmatrix}$,

$$\underline{\mu} = \begin{pmatrix} \underline{\mu}^{(1)} \\ \underline{\mu}^{(2)} \end{pmatrix} \quad \text{and} \quad \underline{\Sigma} = \begin{pmatrix} \underline{\Sigma}_{11} & \underline{\Sigma}_{12} \\ \underline{\Sigma}_{21} & \underline{\Sigma}_{22} \end{pmatrix}.$$

$\underline{\mu}^{(1)}$ is a p -dimensional vector,

$\underline{\Sigma}_{11}$ ($p \times p$), $\underline{\Sigma}_{21}$ ($p \times q$), $\underline{\Sigma}_{12}$ ($q \times p$),

$\underline{\Sigma}_{22}$ ($q \times q$). Define $\underline{\Phi} : \mathbb{R}^n \rightarrow \mathbb{R}^n$

by

$$\underline{\Phi}(\underline{x}) = \begin{pmatrix} \underline{x}^{(1)} \\ \underline{x}^{(2)} - \underline{\Sigma}_{21} \underline{\Sigma}_{11}^{-1} \underline{x}^{(1)} \end{pmatrix}.$$

$\Phi(\underline{x})$ is a linear map

$$\Phi(\underline{x}) = \underbrace{\begin{pmatrix} \underline{I}_r & 0 \\ -\underline{\Sigma}_{21} \underline{\Sigma}_{11}^{-1} & \underline{I}_r \end{pmatrix}}_{\underline{A}} \underline{x}$$

Since the matrix \underline{A} is lower triangular we have

$$\begin{aligned} D\Phi &= \underline{A} \Rightarrow \int \Phi(\underline{x}) = \underline{1} \\ &\Rightarrow \int \Phi^{-1}(\underline{y}) = \underline{1}. \end{aligned}$$

We have

$$\Phi^{-1}(\underline{y}) = \begin{pmatrix} y^{(1)} \\ y^{(2)} + \underline{\Sigma}_{21} \underline{\Sigma}_{11}^{-1} y^{(1)} \end{pmatrix}$$

It follows that

$$f_{\underline{Y}}(\underline{y}) = f_{\underline{X}}(\Phi^{-1}(\underline{y})) \cdot \underline{1}$$

We need some linear algebra.

Suppose \underline{A} , \underline{B} are invertible matrices. Write

$$\underline{A} = \begin{pmatrix} \underline{A}_{11} & \underline{A}_{12} \\ \underline{A}_{21} & \underline{A}_{22} \end{pmatrix} \quad \text{and} \quad \underline{B} = \begin{pmatrix} \underline{B}_{11} & \underline{B}_{12} \\ \underline{B}_{21} & \underline{B}_{22} \end{pmatrix}$$

where \underline{A}_{ij} and \underline{B}_{ij} are of the same dimension. If $\underline{A} \cdot \underline{B} = \underline{I}$ we have

$$\underline{A}_{11} \underline{B}_{11} + \underline{A}_{12} \underline{B}_{21} = \underline{I}$$

$$\underline{A}_{11} \underline{B}_{12} + \underline{A}_{12} \underline{B}_{22} = 0$$

For simplicity we assume $\underline{\mu} = 0$.

We need to compute

$$[\underline{\Phi}^{-1}(y)]^T \cdot \underline{\Sigma}^{-1} [\underline{\Phi}^{-1}(y)].$$

In matrix form this means

$$y^T \begin{pmatrix} \underline{I}_p & \underline{\Sigma}_{11}^{-1} \underline{\Sigma}_{21} \\ \underline{0} & \underline{I}_p \end{pmatrix} \underline{\Sigma}^{-1} \begin{pmatrix} \underline{I}_p & \underline{0} \\ \underline{\Sigma}_{21} \underline{\Sigma}_{11}^{-1} & \underline{I}_q \end{pmatrix}$$

Denote $\underline{A} = \underline{\Sigma}^{-1}$. From

$$\begin{pmatrix} \underline{A}_{11} & \underline{A}_{12} \\ \underline{A}_{21} & \underline{A}_{22} \end{pmatrix} \begin{pmatrix} \underline{\Sigma}_{11} & \underline{\Sigma}_{12} \\ \underline{\Sigma}_{21} & \underline{\Sigma}_{22} \end{pmatrix} = \underline{I}_n$$

we have

$$\underline{A}_{11} \underline{\Sigma}_{11} + \underline{A}_{21} \underline{\Sigma}_{21} = \underline{I}_p$$

$$\underline{A}_{11} \underline{\Sigma}_{12} + \underline{A}_{12} \underline{\Sigma}_{22} = \underline{0}$$

$$\underline{A}_{21} \underline{\Sigma}_{11} + \underline{A}_{22} \underline{\Sigma}_{21} = \underline{0}$$

We compute

$$\begin{pmatrix} \underline{A}_{11} & \underline{A}_{12} \\ \underline{A}_{21} & \underline{A}_{22} \end{pmatrix} \begin{pmatrix} \underline{I}_p & \underline{0} \\ \underline{\Sigma}_{21} \underline{\Sigma}_{11}^{-1} & \underline{I}_q \end{pmatrix}$$

$$= \begin{pmatrix} \underline{A}_{11} + \underline{A}_{12} \underline{\Sigma}_{21} \underline{\Sigma}_{11}^{-1}, & \underline{A}_{12} \\ \underline{A}_{21} + \underline{A}_{22} \underline{\Sigma}_{21} \underline{\Sigma}_{11}^{-1}, & \underline{A}_{22} \end{pmatrix}$$

$$= \begin{pmatrix} \underline{\Sigma}_{11}^{-1} & \underline{A}_{12} \\ \underline{0} & \underline{A}_{22} \end{pmatrix}$$

Continue to get

$$\begin{pmatrix} \underline{I}_p & \underline{\Sigma}_{11}^{-1} \underline{\Sigma}_{12} \\ 0 & \underline{I}_2 \end{pmatrix} \begin{pmatrix} \underline{\Sigma}_{11}^{-1} & \underline{A}_{12} \\ 0 & \underline{A}_{22} \end{pmatrix}$$

$$= \begin{pmatrix} \underline{\Sigma}_{11}^{-1} & \underline{A}_{12} + \underline{\Sigma}_{11}^{-1} \underline{\Sigma}_{12} \underline{A}_{22} \\ - & \underline{A}_{22} \end{pmatrix}$$

But

$$\begin{pmatrix} \underline{\Sigma}_{11} & \underline{\Sigma}_{12} \\ \underline{\Sigma}_{21} & \underline{\Sigma}_{22} \end{pmatrix} \begin{pmatrix} \underline{A}_{11} & \underline{A}_{12} \\ \underline{A}_{21} & \underline{A}_{22} \end{pmatrix} = \underline{I}_n$$

gives

$$\underline{\Sigma}_{11} \underline{A}_{12} + \underline{\Sigma}_{12} \underline{\Sigma}_{22} = 0, \text{ so}$$

$$\underline{\Sigma}_{11} (\underline{A}_{12} + \underline{\Sigma}_{11}^{-1} \underline{\Sigma}_{12} \underline{A}_{22}) = 0$$

The linear equations give

$$\underline{A}_{22} = (\underline{\Sigma}_{22} - \underline{\Sigma}_{12} \underline{\Sigma}_{11}^{-1} \underline{\Sigma}_{21})$$

(see Appendix)

So we have

$$[\bar{\Phi}^{-1}(y)]^T \underline{\Sigma}^{-1} [\bar{\Phi}^{-1}(y)]$$

$$= y^T \begin{pmatrix} \underline{\Sigma}_{11}^{-1} & 0 \\ 0 & (\underline{\Sigma}_{22} - \underline{\Sigma}_{21} \underline{\Sigma}_{11}^{-1} \underline{\Sigma}_{12})^{-1} \end{pmatrix} y$$

$$= [y^{(1)}]^T \underline{\Sigma}_{11}^{-1} y^{(1)}$$

$$+ [y^{(2)}]^T (\underline{\Sigma}_{22} - \underline{\Sigma}_{21} \underline{\Sigma}_{11}^{-1} \underline{\Sigma}_{12})^{-1} y^{(2)}$$

Comment: in general replace

y by $y - A$. So we have

$$f_{\underline{\Sigma}}(y) = f(y^{(1)}) \cdot g(y^{(2)}).$$

This means that

$$\underline{y}^{(1)} = (x_1, \dots, x_r)$$

$$\underline{y}^{(2)} = \underline{x}^{(2)} - \underline{\Sigma}_{21} \underline{\Sigma}_{11}^{-1} \underline{x}^{(1)}$$

are independent vectors.

Appendix: if we have

$$\begin{pmatrix} \underline{\Sigma}_{11} & \underline{\Sigma}_{12} \\ \underline{\Sigma}_{21} & \underline{\Sigma}_{22} \end{pmatrix} \begin{pmatrix} \underline{A}_{11} & \underline{A}_{12} \\ \underline{A}_{21} & \underline{A}_{22} \end{pmatrix} = \underline{I}_n$$

then

$$\underline{\Sigma}_{11} \underline{A}_{11} + \underline{\Sigma}_{12} \underline{A}_{22} = \underline{I}_p$$

$$\underline{\Sigma}_{11} \underline{A}_{12} + \underline{\Sigma}_{12} \underline{A}_{22} = \underline{0}$$

$$\underline{\Sigma}_{21} \underline{A}_{11} + \underline{\Sigma}_{22} \underline{A}_{21} = \underline{0}$$

$$\underline{\Sigma}_{21} \underline{A}_{12} + \underline{\Sigma}_{22} \underline{A}_{22} = \underline{I}_q$$

We have a system of 4 linear equations with 4 unknowns.

Multiply the second equation with $\underline{\Sigma}_{11}^{-1}$ from the left to get

$$\underline{A}_{12} + \underline{\Sigma}_{11}^{-1} \underline{\Sigma}_{12} \underline{A}_{22} = \underline{0}$$

Insert this into the last equation to get

$$- \underline{\Sigma}_{21} \underline{\Sigma}_{11}^{-1} \underline{\Sigma}_{12} \underline{A}_{22} + \underline{\Sigma}_{22} \underline{A}_{22} = \underline{I}_p$$

We have

$$\underline{A}_{22} = \left(\underline{\Sigma}_{22} - \underline{\Sigma}_{12} \underline{\Sigma}_{11}^{-1} \underline{\Sigma}_{21} \right)^{-1}$$

This result is known as the inversion lemma.

Remark: Invertibility follows from the fact that the product is \underline{I}_p .

3.5. Conditional distributions

In elementary probability

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

If X is a discrete random variable with values $\{x_1, x_2, \dots\}$

the distribution is given by

the probabilities $P(X = x_k)$.

If we have additional information

in the sense that the event

B has happened our opinion

about the probabilities of

event $\{X = x_k\}$ change to the

conditional probabilities

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

We can verify easily that

$$\sum_{x_k} P(X = x_k | B) = 1.$$

This observation motivates the definition of conditional probabilities and distributions.

○ Definition: Let X be a discrete random variable with values in $\{x_1, x_2, \dots\}$.

The conditional distribution of X given B with $P(B) > 0$

○ is given by

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

Comment: In most cases the event B is of the form

$B = \{Y = y_e\}$ for some random variable Y .

Example: Let X, Y be independent
with $X \sim \text{Bin}(m, p)$ and

$Y \sim \text{Bin}(u, p)$. Let $Z = X + Y$.

We know that $Z \sim \text{Bin}(m+u, p)$.

The conditional distribution of

X given $\{Z = r\}$ is given by

$$P(X=k | Z=r) = \frac{P(X=k, Z=r)}{P(Z=r)}$$

$$= \frac{P(X=k, Y=r-k)}{P(Z=r)}$$

$$= \frac{P(X=k) P(Y=r-k)}{P(Z=r)} \quad \text{indep}$$

$$= \frac{\binom{m}{k} p^k q^{m-k} \cdot \binom{u}{r-k} p^{r-k} q^{u-r+k}}{\binom{m+u}{r} p^r q^{m+u-r}}$$

$$= \frac{\binom{m}{k} \binom{n}{r-k}}{\binom{m+n}{r}}$$

for $k \leq \min(m, r)$ and
 $k \geq \max(0, r-n)$.

We recognize the hypergeometric distribution. We write

$$X | Z=r \sim \text{Hypergeom}(r, m, m+n).$$

Definition: Let \underline{X} be a discrete random vector with values $\{x_1, x_2, \dots\}$. Let B be an event. The conditional distribution of \underline{X} given B with $P(B) > 0$ is given by conditional probabilities

$$P(\underline{X} = \underline{x}_k | B) = \frac{P(\{X = x_k\} \cap B)}{P(B)}.$$

As before in most cases B is of the form $B = \{Y = y_k\}$ for some random vector \underline{Y} .

Example: let $\underline{X} = (X_1, \dots, X_r)$

be multinomial with parameters n and $p = (p_1, p_2, \dots, p_r)$. Let $s < r$.

What is the conditional distribution of (X_1, X_2, \dots, X_s)

given $Y = X_1 + X_2 + \dots + X_s = m$.

Denote $Z = X_{s+1} + \dots + X_r$. We know

that $Y \sim \text{Bin}(n, p_1 + \dots + p_s)$. We

compute for $k_1 + \dots + k_s = m$

$$P(X_1 = k_1, \dots, X_s = k_s \mid Y = m)$$

$$= \frac{P(X_1 = k_1, \dots, X_s = k_s, Z = n - m)}{P(Z = n - m)}$$

$$= \frac{n!}{k_1! \dots k_s! (n-m)!}$$

$$\times p_1^{k_1} \dots p_s^{k_s} (1-p_1-\dots-p_s)^{n-m}$$

$$/ \binom{n}{m} (p_1+\dots+p_s)^m (1-p_1-\dots-p_s)^{n-m}$$

$$= \frac{m!}{k_1! \dots k_s!}$$

$$\times \frac{p_1^{k_1} \dots p_s^{k_s}}{(p_1+\dots+p_s)^m}$$

= (*)

We denote: $\tilde{p}_k = \frac{p_k}{(p_1+\dots+p_s)}$

for $k = 1, 2, \dots, s$. We have:

$$* = \frac{m!}{k_1! \dots k_s!} \tilde{p}_1^{k_1} \dots \tilde{p}_s^{k_s}$$

Conclusion: (X_1, X_2, \dots, X_s)

has the multinomial distribution

with parameters $n-m$ and

$\hat{p} = (\hat{p}_1, \dots, \hat{p}_s)$. We write

$\underline{x}' = (x_1, \dots, x_s)$ and

$\underline{x}' \mid x_1 + \dots + x_s = m \sim \text{Multinom}(m, \hat{p})$.

(1) For the continuous case the intuitive idea is that we will define conditional densities. If (X, Y) has density $f_{X,Y}(x,y)$ then the conditional density of Y given $X=x$ should be proportional to the function

$$y \mapsto f_{X,Y}(x,y)$$

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

Example : let

$$f_{X,Y}(x,y) = \frac{1}{2\pi\sqrt{1-\rho^2}} \times e^{-\frac{x^2 - 2\rho xy + y^2}{2(1-\rho^2)}}$$

for $|\rho| < 1$. We know that $X \sim N(0,1)$ i.e.

$$f_X(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$$

We write

$$f_{X,Y}(x,y) = \frac{1}{2\pi\sqrt{1-\rho^2}} e^{-x^2/2} \times e^{-\frac{(y-\rho x)^2}{2(1-\rho^2)}}$$

It follows that

$$f_{Y|X=x}(y) \\ = \frac{1}{\sqrt{2\pi}\sqrt{1-\rho^2}} e^{-\frac{(y-\rho x)^2}{2(1-\rho^2)}}$$

We notice

$$Y|X=x \sim N(\rho x, 1-\rho^2)$$

The definition has a vector version.

Definition: Let $(\underline{x}, \underline{y})$ have

density $f_{\underline{x}, \underline{y}}(\underline{x}, \underline{y})$. Assume $f_{\underline{x}}(\underline{x}) > 0$. The conditional density of \underline{y} given $\{\underline{x} = \underline{x}\}$ is given by

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

Example : Let $\underline{X} = (\underline{X}^{(1)}, \underline{X}^{(2)}) \sim N(\underline{\mu}, \underline{\Sigma})$.

What is $f_{\underline{X}^{(2)} | \underline{X}^{(1)} = \underline{x}^{(1)}}(\underline{x}^{(2)})$?

Direct calculation is difficult but we found out that $\underline{X}^{(1)}$ and

$\underline{X}^{(2)} - \underline{\Sigma}_{21} \underline{\Sigma}_{11}^{-1} \underline{X}^{(1)}$ are independent

vectors. If we write $\underline{Y} = \underline{X}^{(2)} - \underline{\Sigma}_{21} \underline{\Sigma}_{11}^{-1} \underline{X}^{(1)}$

we know that

$$f_{\underline{Y}}(\underline{y}) = \frac{1}{(2\pi)^{2/2} \sqrt{\det(\underline{\Sigma}_{22} - \underline{\Sigma}_{21} \underline{\Sigma}_{11}^{-1} \underline{\Sigma}_{12})}}$$

$$\times \exp\left(-\frac{1}{2} \left(\underline{y} - \underline{\mu}^{(2)} + \underline{\Sigma}_{21} \underline{\Sigma}_{11}^{-1} \underline{\mu}^{(1)}\right)^T \left(\underline{\Sigma}_{22} - \underline{\Sigma}_{21} \underline{\Sigma}_{11}^{-1} \underline{\Sigma}_{12}\right)^{-1} \left(\underline{y} - \underline{\mu}^{(2)} + \underline{\Sigma}_{21} \underline{\Sigma}_{11}^{-1} \underline{\mu}^{(1)}\right)\right)$$

But

$$\begin{pmatrix} \underline{X}^{(1)} \\ \underline{X}^{(2)} \end{pmatrix} = \begin{pmatrix} \underline{\mu}^{(1)} & 0 \\ \underline{\Sigma}_{21} \underline{\Sigma}_{11}^{-1} & \underline{I}_2 \end{pmatrix} \begin{pmatrix} \underline{X}^{(1)} \\ \underline{Y} \end{pmatrix}$$

The Jacobian of this is 1
 so we can write

$$f_{\underline{x}^{(2)}}(\underline{x}^{(2)}) = f_{\underline{x}^{(1)}}(\underline{x}^{(1)})$$

$$\times f_{\underline{y}}(\underline{x}^{(2)} - \underline{\Sigma}_{21} \underline{\Sigma}_{11}^{-1} \underline{x}^{(1)})$$

Now it is easy to divide by

$f_{\underline{x}^{(1)}}(\underline{x}^{(1)})$. We get

$$f_{\underline{x}^{(2)}} | \underline{x}^{(1)} = \underline{x}^{(1)}(\underline{x}^{(2)})$$

$$= f_{\underline{y}}(\underline{x}^{(2)} - \underline{\Sigma}_{12} \underline{\Sigma}_{11}^{-1} \underline{x}^{(1)}).$$

Using the form of $f_{\underline{y}}$ we

find:

$$\underline{x}^{(2)} | \underline{x}^{(1)} = \underline{x}^{(1)} \sim N\left(\underline{\mu}^{(2)} + \underline{\Sigma}_{21} \underline{\Sigma}_{11}^{-1} (\underline{x}^{(1)} - \underline{\mu}^{(1)}), \underline{\Sigma}_{22} - \underline{\Sigma}_{21} \underline{\Sigma}_{11}^{-1} \underline{\Sigma}_{12}\right)$$

4. Expectation and variance

4.1. Expectation in general

For a discrete random variable we defined

$$E(X) = \sum_{x_k} x_k P(X = x_k)$$

$$E[f(X)] = \sum_{x_k} f(x_k) P(X = x_k)$$

For a discrete random vector we have

$$E[f(\underline{X})] = \sum_{\underline{x}_k} f(\underline{x}_k) P(\underline{X} = \underline{x}_k)$$

We need to extend the notion of expectation to continuous random variables and vectors.

Definitions:

(i) Let X have density $f_X(x)$.

We define

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

$$E[f(X)] = \int_{-\infty}^{\infty} f(x) f_X(x) dx$$

Technical note: we say

that $E(X)$ exists if the integral $\int_{-\infty}^{\infty} |x| f_X(x) dx$ converges and similarly for $E[f(X)]$.

(ii) Let the random vector \underline{X} have density $f_{\underline{X}}(\underline{x})$. We define

$$E[f(\underline{X})] = \int_{\mathbb{R}^n} f(\underline{x}) f_{\underline{X}}(\underline{x}) d\underline{x}$$

Examples :

(i) $X \sim N(\mu, \sigma^2)$

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

$$= \int_{-\infty}^{\infty} x \cdot \frac{1}{\sqrt{2\pi} \sigma} \cdot e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx$$

= (*) New variable: $\frac{x-\mu}{\sigma} = u$

$$= \frac{1}{\sqrt{2\pi}} \cdot \int_{-\infty}^{\infty} (\sigma u + \mu) e^{-\frac{u^2}{2}} du$$

$$= \mu \cdot \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-\frac{u^2}{2}} du$$

$$= \mu$$

because $\int_{-\infty}^{\infty} u \cdot e^{-\frac{u^2}{2}} du = 0$

(odd function).

We continue

$$E(x^2) = \int_{-\infty}^{\infty} x^2 \cdot f_X(x) dx$$

$$= \frac{1}{\sqrt{2\pi} \sigma} \int_{-\infty}^{\infty} x^2 \cdot e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx$$

$$= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} (\sigma u + \mu)^2 e^{-u^2/2} du$$

$$= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} (\sigma^2 u^2 + \mu^2) \cdot e^{-u^2/2} du$$

(the middle term = 0)

$$= \sigma^2 \cdot \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} u^2 \cdot e^{-u^2/2} du + \mu^2 = (*)$$

We integrate by parts

$$\int_{-\infty}^{\infty} u^2 \cdot e^{-u^2/2} du =$$

$$= \int_{-\infty}^{\infty} u \cdot u \cdot e^{-u^2/2} du$$

$$= -u \cdot e^{-u^2/2} \Big|_{-\infty}^{\infty} + \int_{-\infty}^{\infty} e^{-u^2/2} du$$

$$= \sqrt{2\pi}$$

$$(*) = \sigma^2 + \mu^2$$

Definition : Let X be a random variable. Let $\mu = E(X)$.

We call $E(X^m)$ the m -th moment of X and $E[(X-\mu)^m]$ the m -th central moment of X .

Example : Let $X \sim \Gamma(a, \lambda)$.

We compute the m -th moment of X as

$$\begin{aligned} E(X^m) &= \int_0^{\infty} x^m \cdot f_X(x) dx \\ &= \int_0^{\infty} \frac{\lambda^a}{\Gamma(a)} \cdot x^m \cdot x^{a-1} \cdot e^{-\lambda x} dx \\ &= \frac{\lambda^a}{\Gamma(a)} \int_0^{\infty} x^{m+a-1} e^{-\lambda x} dx \\ &= \frac{\lambda^a}{\Gamma(a)} \cdot \frac{\Gamma(m+a)}{\lambda^{m+a}} \end{aligned}$$

$$= \frac{P(m+a)}{P(a) \lambda^m}$$

$$= \frac{(a)_m}{\lambda^m}$$

where $(a)_m = a(a+1) \cdots (a+m-1)$.

Example: Let (X, Y) have

the density $f_{X,Y}(x, y) =$

$$= \frac{1}{2\pi\sqrt{1-\rho^2}} e^{-\frac{x^2 - 2\rho xy + y^2}{2(1-\rho^2)}}$$

with $|\rho| < 1$. We have

$$E(XY) = \int_{\mathbb{R}^2} xy \cdot f_{X,Y}(x, y) dx dy$$

$$= \frac{1}{2\pi\sqrt{1-\rho^2}} \int_{\mathbb{R}^2} xy \cdot e^{-\frac{x^2 - 2\rho xy + y^2}{2(1-\rho^2)}} dx dy$$

Fubini

$$= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} x \cdot e^{-x^2/2} dx \times$$

$$\times \frac{1}{\sqrt{2\pi}\sqrt{1-\rho^2}} \int_{-\infty}^{\infty} y \cdot e^{-\frac{(y-\rho x)^2}{2(1-\rho^2)}} dy$$

The last integral is ρx .

We computed it in the first example. We have

$$\begin{aligned} E(x^2) &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \rho x^2 \cdot e^{-x^2/2} dx \\ &= \rho \cdot E(x^2) \\ &= \rho \end{aligned}$$

If $X \sim N(0, 1)$ then $E(x^2) = \sigma^2 + \mu^2 = \sigma^2$.

Theorem 4.1: Let (X, Y) have density $f_{X,Y}(x,y)$. We have

$$E[\alpha X + \beta Y] = \alpha E(X) + \beta E(Y)$$

Proof: We have

$$E[\alpha X + \beta Y]$$

$$= \int_{\mathbb{R}^2} (\alpha x + \beta y) f_{X,Y}(x,y) dx dy$$

$$= \alpha \cdot \int_{\mathbb{R}^2} x f_{X,Y}(x,y) dx dy$$

$$+ \beta \int_{\mathbb{R}^2} y f_{X,Y}(x,y) dx dy$$

$$= \alpha \cdot E(X) + \beta E(Y)$$

Remarks:

(i) The same proof works

with $\alpha f(x,y) + \beta g(x,y)$ i.e.

$$E[\alpha f(x,y) + \beta g(x,y)]$$

$$= \alpha E[f(x,y)] + \beta E[g(x,y)].$$

(ii) The theorem has a vector version:

$$E \left[\sum_{k=1}^r \alpha_k X_k \right] = \sum_{k=1}^r \alpha_k E(X_k).$$

The expectation is always linear.

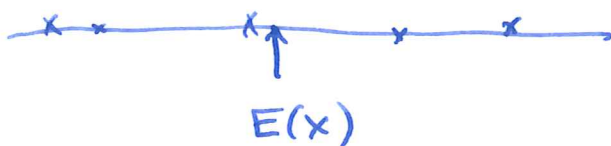
4.2 Variance and covariance

We motivated the expectation as a "long term average".

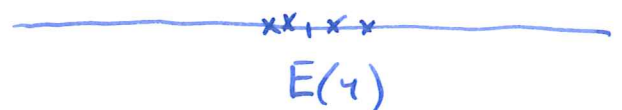
We would like to devise a measure of dispersion of a random variable X .

Figure :

(i)



(ii)



In the figure we have "repeated" values of X and Y . The values of X are more "dispersed". Why do we say this? On average the values of X are further

- away from $E(X)$. Repetitions to the right and to the left contribute an equal amount to dispersion so we take absolute distances. However,
- Gauss chose the square.

His choice was motivated by mathematical considerations.

Denote v_1, v_2, \dots, v_n the repetitions of X .

The dispersion according to Gauss is

$$\frac{(v_1 - E(x))^2 + \dots + (v_n - E(x_n))^2}{n}$$

If we take $f(x) = (x - E(x))^2$
we see that the above average
 $\approx E[(x - E(x))^2]$.

Definition: The variance of the random variable X is given by

$$\text{var}(X) = E[(X - E(X))^2]$$

We compute

$$\begin{aligned} E[(X - E(X))^2] &= \\ &= E[X^2 - 2E(X) \cdot X + E(X)^2] \end{aligned}$$

$$\begin{aligned} \text{Var.} &= E(x^2) - 2E(x) \cdot E(x) + E(x)^2 \\ &= E(x^2) - [E(x)]^2 \end{aligned}$$

Alternative form:

$$\text{var}(x) = E(x^2) - [E(x)]^2$$

Examples:

(i) $X \sim \text{Bin}(n, p)$

We know

$$E(x^2) = npq + n^2 p^2 \text{ and}$$

$$E(x) = n \cdot p$$

$$\text{var}(x) = E(x^2) - E(x)^2$$

$$= npq$$

(ii) $X \sim \text{Neg Bin}(m, p)$

We know

$$E(X) = \frac{m}{p}$$

$$E(X^2) = \frac{m \cdot q}{p^2} + \frac{m^2}{p^2}$$

$$\text{var}(X) = \frac{m \cdot q}{p^2}$$

(iii) $X \sim N(\mu, \sigma^2)$

We know

$$E(X) = \mu$$

$$E(X^2) = \sigma^2 + \mu^2$$

$$\begin{aligned} \text{var}(X) &= E(X^2) - [E(X)]^2 \\ &= \sigma^2 \end{aligned}$$

Comment: If $X \sim N(\mu, \sigma^2)$ the two parameters have a nice

interpretation. They are the expectation and the variance.

What about the variance of sums? We compute

$$\text{var}(x+y) = E[(x+y)^2] - [E(x+y)]^2$$

$$= E(x^2 + 2xy + y^2)$$

$$- (E(x) + E(y))^2$$

l.u.

$$= \underbrace{E(x^2)} + 2E(xy) + \underbrace{E(y^2)}$$

$$- \underbrace{E(x)^2} - 2E(x)E(y) - \underbrace{E(y)^2}$$

$$= \text{var}(x) + \text{var}(y)$$

$$+ 2[E(xy) - E(x)E(y)]$$

There is no reason for the term in square brackets to be 0.

Definition : Let X, Y be random variables. The quantity

$$E(XY) - E(X) \cdot E(Y)$$

is called the covariance of X and Y and denoted by

○ $\text{cov}(X, Y)$.

Remark : An application of linearity gives that

$$\text{cov}(X, Y) = E[(X - E(X))(Y - E(Y))].$$

○ Theorem 4.2 : Let X_1, \dots, X_r and Y_1, Y_2, \dots, Y_s be random variables.

We have

$$\begin{aligned} \text{cov}\left(\sum_{k=1}^n \alpha_k X_k, \sum_{l=1}^s \beta_l Y_l\right) \\ = \sum_{k=1}^n \sum_{l=1}^s \alpha_k \beta_l \text{cov}(X_k, Y_l) \end{aligned}$$

Proof: We compute

$$E \left[\left(\sum_{k=1}^r \alpha_k X_k \right) \left(\sum_{e=1}^s \beta_e Y_e \right) \right]$$

$$= E \left[\sum_{k=1}^r \sum_{e=1}^s \alpha_k \beta_e X_k Y_e \right]$$

$$\stackrel{\text{lin.}}{=} \sum_{k=1}^r \sum_{e=1}^s \alpha_k \beta_e E(X_k Y_e)$$

On the other hand

$$E \left(\sum_{k=1}^r \alpha_k X_k \right) \cdot E \left(\sum_{e=1}^s \beta_e Y_e \right)$$

$$\stackrel{\text{lin.}}{=} \sum_{k=1}^r \sum_{e=1}^s \alpha_k \beta_e E(X_k) E(Y_e)$$

We subtract and get the result.

Remark: The property is called bilinearity.

The definitions give us further properties of covariances that follow from definitions:

- (i) $\text{var}(aX) = a^2 \text{var}(X)$
- (ii) $\text{cov}(X, X) = \text{var}(X)$
- (iii) $\text{cov}(X, Y) = \text{cov}(Y, X)$
- (iv) $\text{cov}(\alpha X, \beta Y) = \alpha\beta \text{cov}(X, Y)$.

Theorem 4.3 : Let X_1, \dots, X_r

be random variables. We have

$$\text{var}\left(\sum_{k=1}^r \alpha_k X_k\right)$$

$$= \sum_{k=1}^r \alpha_k^2 \text{var}(X_k)$$

$$+ \sum_{\substack{k, l=1 \\ k \neq l}}^r \alpha_k \alpha_l \text{cov}(X_k, X_l).$$

Proof : Follows directly from Theorem 4.2.

Special case: Let X, Y be discrete and independent.

Then

$$\begin{aligned} E(X \cdot Y) &= \sum_{x,y} x \cdot y P(X=x, Y=y) \\ &= \sum_{x,y} x \cdot y P(X=x) P(Y=y) \\ &= \left(\sum_x x P(X=x) \right) \cdot \left(\sum_y y P(Y=y) \right) \\ &= E(X) E(Y). \end{aligned}$$

This means $\text{cov}(X, Y) = 0$.

A similar calculation is valid for continuous X, Y .

Remark: For independent X, Y and functions f, g we have

$$E[f(X)g(Y)] = E(f(X))E(g(Y)).$$

with the same proof.

As a consequence for independent
 X_1, X_2, \dots, X_n we have

$$\text{var}(X_1 + \dots + X_n) = \text{var}(X_1) + \dots + \text{var}(X_n)$$

Examples : (i) If $\underline{X} = (X_1, \dots, X_r)$
is multivariate normal we have

$$E(X_k X_l) = -np_k p_l + n^2 p_k p_l$$

and $E(X_k) = np_k$ and $E(X_l) = np_l$.

We have

$$\text{cov}(X_k, X_l) = -n p_k p_l$$

(ii) If

$$f_{X,Y}(x,y) = \frac{1}{2\pi\sqrt{1-\rho^2}} e^{-\frac{x^2 - 2\rho xy + y^2}{2(1-\rho^2)}}$$

we have

$$E(X) = \rho \quad \text{and}$$

$$E(Y) = \rho$$

so

$$\text{cov}(X, Y) = \rho$$

Method of indicators :

If we can write a random variable X as a sum of indicators we can in many cases compute variances by Theorem 4.3.

○ If $I \sim \text{Bernoulli}(p)$ then

$$E(I^2) = E(I) = p \quad \text{so}$$

$$\text{var}(I) = E(I^2) - E(I)^2$$

$$= p - p^2$$

$$= p(1-p)$$

○ If I, J are indicators then

$$\text{cov}(I, J) = E(I \cdot J) - E(I)E(J)$$

$$= P(I=1, J=1)$$

$$- P(I=1)P(J=1)$$

Example: Let $X \sim \text{HyperGeom}(n, B, N)$.

We wrote

$$X = I_1 + \dots + I_n.$$

We have

$$\begin{aligned} \text{var}(X) &= \sum_{k=1}^n \text{var}(I_k) \\ &\quad + \sum_{\substack{k, l=1 \\ k \neq l}}^n \text{cov}(I_k, I_l) \end{aligned}$$

We know: $I_k \sim \text{Bernoulli}\left(\frac{B}{N}\right)$

$$\text{so } \text{var}(I_k) = \frac{B}{N} \cdot \left(1 - \frac{B}{N}\right).$$

We have

$$\begin{aligned} P(I_1 = 1, I_2 = 1) &= P(\text{1st \& 2nd ball} \\ &\quad \text{black}) \end{aligned}$$

$$= \frac{B}{N} \cdot \frac{B-1}{N-1}.$$

It follows

$$\begin{aligned}\text{Cov}(I_1, I_2) &= \frac{B}{N} \cdot \frac{B-1}{N-1} - \left(\frac{B}{N}\right)^2 \\ &= \frac{B}{N} \left[\frac{(B-1)N - B(N-1)}{N(N-1)} \right] \\ &= \frac{B}{N} \left[\frac{-N + B}{N(N-1)} \right] \\ &= -\frac{B}{N} \left(1 - \frac{B}{N}\right) \cdot \frac{1}{N-1}\end{aligned}$$

But by symmetry (I_k, I_l) has the same distribution as (I_1, I_2) .

So all covariances are the same.

We have

$$\begin{aligned}\text{var}(\underline{X}) &= n \frac{B}{N} \left(1 - \frac{B}{N}\right) + n(n-1) \times (-1) \times \\ &\quad \times \frac{B}{N} \left(1 - \frac{B}{N}\right) \cdot \frac{1}{N-1} \\ &= n \frac{B}{N} \left(1 - \frac{B}{N}\right) \left(1 - \frac{n-1}{N-1}\right) \\ &= n \cdot \frac{B}{N} \left(1 - \frac{B}{N}\right) \frac{N-n}{N-1}\end{aligned}$$

4.3. Conditional expectation

Idea: If X is a discrete random variable then

$$E[f(X)] = \sum_{x_k} f(x_k) P(X=x_k).$$

The expectation is computed using the

distribution. We can use the same

idea with the conditional

distribution $P(X=x_k | B)$ for $P(B) > 0$.

Definition: Let X be a discrete

random variable. The conditional

expectation of X given B is

given by

$$E(X | B) = \sum_{x_k} x_k \cdot P(X=x_k | B)$$

and

$$E(f(X) | B) = \sum_{x_k} f(x_k) P(X=x_k | B).$$

Technical note: We understand the

existence of $E(X | B)$ the same way

as for usual expectations.

In most cases B will be of the form $B = \{Y = ye\}$ for some random variable Y .

Example: Players A and B get 5 cards each from a well shuffled deck of cards. Let X be the number of aces of A and Y the number of aces of B. We have

$Y|X=k \sim \text{HyperGeom}(5, 4-k, 47)$. It

follows that

$$E(Y|X=k) = 5 \cdot \frac{4-k}{47}$$

We know that for $Z \sim \text{HyperGeom}(n, B, N)$

$$\text{we have } \text{var}(Z) = n \cdot \frac{B}{N} \left(1 - \frac{B}{N}\right) \frac{N-n}{N-1}.$$

so

$$E(Z^2) = \text{var}(Z) + n^2 \cdot \frac{B^2}{N^2}$$

We have

$$E(Y^2 | X=k)$$

$$= 5 \cdot \frac{4-k}{47} \left(1 - \frac{4-k}{47}\right) \cdot \frac{47-5}{47-1} \\ + 5^2 \cdot \frac{(4-k)^2}{47^2}$$

There is an alternative way to write the conditional expectation.

If X is discrete and B is an event we compute

$$E[\underbrace{X \cdot \mathbb{1}_B}_{\uparrow}] = \sum_{x_k} x_k \cdot P(\{X=x_k\} \cap B) \\ = (*)$$

This random variable has value x_k with probability $P(\{X=x_k\} \cap B)$ and possibly 0 with probability

$$P(B \cup \{X=0\})$$

$$\begin{aligned}
 (*) &= \sum_{x_k} x_k \frac{P(\{X = x_k\} \cap B)}{P(B)} \cdot P(B) \\
 &= P(B) \cdot E(X|B)
 \end{aligned}$$

We have

$$E(X|B) = \frac{E(X \cdot 1_B)}{P(B)}$$

$$E[f(X)|B] = \frac{E[f(X) \cdot 1_B]}{P(B)}.$$

Theorem 4.4: Let $\{H_1, H_2, \dots\}$ be

a partition of $(\Omega, \mathcal{F}, \mathbb{P})$. Let X

be a discrete random variable with

$E(|X|) < \infty$. We have

$$E(X) = \sum_k E(X|H_k) \cdot P(H_k)$$

Proof: We compute

$$\sum_k E(X | H_k) \cdot P(H_k)$$

$$= \sum_k \left(\sum_{x_e} x_e P(X=x_e | H_k) \right) P(H_k)$$

$$= \sum_{x_e} x_e \underbrace{\sum_k P(X=x_e | H_k) P(H_k)}_{= P(X=x_e)}$$

$$= \sum_{x_e} x_e P(X=x_e)$$

$$= E(X)$$

We used the law of total probabilities.

The statement is the law of

total expectations.

Example: We toss a coin until we get r consecutive heads.

Tosses are independent and the

probability of heads is p . Let

X be the number of tosses

needed.

Example: if $r = 4$ and we get

H T T H H T H T T H T H H H H

$$X = 15$$

We want $E(X)$. Let

$H_k = \{\text{the first T appears in position } k\}$.

We have

$$E(X | H_k) = r \quad \text{if } k = r+1, \dots$$

and

$$E(X | H_k) = k + E(X) \quad \text{if } k = 1, \dots, r;$$

The law of total expectation gives

$$\begin{aligned} E(X) &= \sum_{k=1}^{\infty} E(X | H_k) P(H_k) \\ &= \sum_{k=1}^r (k + E(X)) P(H_k) \\ &\quad + \sum_{k=r+1}^{\infty} r \cdot P(H_k) \end{aligned}$$

This last expression is a linear equation for $E(x)$. We compute

$$(i) \sum_{k=r+1}^{\infty} r P(H_k) =$$

$$= r \cdot \sum_{k=r+1}^{\infty} p^{k-1} \cdot q$$

$$= r \cdot p^r \cdot q \cdot \sum_{k=0}^{\infty} p^k$$

$$= r \cdot p^r \cdot q \cdot \frac{1}{1-p}$$

$$= r \cdot p^r$$

$$(ii) \sum_{k=1}^r P(H_k) = q \cdot \sum_{k=1}^r p^{k-1}$$

$$= q \cdot \frac{1-p^r}{1-p}$$

$$= 1-p^r$$

$$(iii) \sum_{k=1}^r k \cdot P(H_k) = \sum_{k=1}^r k \cdot p^{k-1} \cdot q$$

$$= \frac{d}{dp} (p + p^2 + \dots + p^r) \cdot q$$

$$= (*)$$

$$\begin{aligned}
 (*) &= \frac{d}{dp} \left(\frac{p(1-p^n)}{1-p} \right) \cdot 2 \\
 &= \frac{[(1-p^n) - r p p^{n-1}](1-p) + p(1-p^n)}{(1-p)^2} \cdot 2 \\
 &= \frac{1-p^n - r p^n + r p^{n+1}}{(1-p)^2} \cdot 2 \\
 &= \frac{1-p^n - r p^n + r p^{n+1}}{2}
 \end{aligned}$$

Rewrite

$$\begin{aligned}
 E(x) &= \frac{1-p^n - r p^n + r p^{n+1}}{2} \\
 &\quad + E(x) \cdot (1-p^n) \\
 &\quad + r \cdot p^n
 \end{aligned}$$

We have

$$\begin{aligned}
 E(x) &= r + \frac{1-p^n - r p^n + r \cdot p^{n+1}}{2 \cdot p^n} \\
 &= r + \frac{1-p^n + r p^n (p-1)}{2 \cdot p^n} \\
 &= \frac{1-p^n}{2 \cdot p^n}
 \end{aligned}$$

$$E[f(\underline{x})] = \sum_k E[f(\underline{x}) | H_k] P(H_k)$$

The proof is identical to the proof we had before.

Definition: We define

$$\text{var}(x | B) = E(x^2 | B) - [E(x | B)]^2$$

and

$$\text{cov}(x, y | B) = E(xy | B) - E(x | B)E(y | B).$$

For continuous random variables we use conditional densities

to compute conditional expectations.

We have:

Definition: Let X have

conditional density $f_{Y|X=x}(y)$

given $X=x$.

We define

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

$$E[f(Y) | X=x] = \int_{-\infty}^{\infty} f(y) f_{Y|X=x}(y) dy.$$

Technical note: We define

existence the way existence is defined for usual expectations.

Comment: The same definition holds for vectors. We define

$$E[f(\underline{Y}) | \underline{X} = \underline{x}]$$

$$= \int_{\mathbb{R}^2} f(\underline{y}) f_{\underline{Y}|\underline{X}=\underline{x}}(\underline{y}) d\underline{y}.$$

Definition: We define

$$\text{var}(Y | \underline{X} = \underline{x}) = E(Y^2 | \underline{X} = \underline{x}) - \left[E(Y | \underline{X} = \underline{x}) \right]^2$$

and

$$\text{cov}(Y_1, Y_2 | \underline{x} = \underline{x})$$

$$= E(Y_1 \cdot Y_2 | \underline{x} = \underline{x})$$

$$- E(Y_1 | \underline{x} = \underline{x}) E(Y_2 | \underline{x} = \underline{x})$$

Example: let

$$f_{X,Y}(x,y) = \frac{1}{2\pi\sqrt{1-\rho^2}} e^{-\frac{x^2 - 2\rho xy + y^2}{2(1-\rho^2)}}$$

We have computed that

$$f_{Y|X=X}(y) = \frac{1}{2\pi\sqrt{1-\rho^2}} e^{-\frac{(y-\rho x)^2}{2(1-\rho^2)}}$$

or $Y|X=x \sim N(\rho x, 1-\rho^2)$. From this

we have

$$E(Y|X=x) = \rho x \quad \text{and}$$

$$\text{var}(Y|X=x) = 1-\rho^2.$$

5. Generating functions

5.1. Definitions and basic properties

The idea of generating functions comes from analysis and combinatorics. If c_0, c_1, \dots is a sequence of complex numbers then we can define the power series

$$G(s) = \sum_{k=0}^{\infty} c_k \cdot s^k \quad \text{for } s \in \mathbb{C}.$$

We know from analysis that such power series converge for $|s| < R$ where R is the radius of convergence. Analysis further gives that

$$\frac{1}{R} = \limsup_{n \rightarrow \infty} \sqrt[n]{|c_n|}.$$

If $|c_n| \leq 1$ for all n then

$$\frac{1}{R} \leq 1 \Rightarrow R \geq 1.$$

In this chapter we will only look at non-negative integer valued random variables.

○ Definition: Let X be a random variable with values $0, 1, 2, \dots$

We define the generating function of X , denoted by $G_X(s)$ as the power series

$$G_X(s) = \sum_{k=0}^{\infty} P(X=k) \cdot s^k$$

Comments:

(i) The idea is to "pack" up the distribution in a function.

(ii) Since $\sum_{k=0}^{\infty} P(X=k) = 1$.

the power series is dominated by $P(X=k)$ for $|s| \leq 1$ and converges uniformly to a continuous function.

Examples:

(i) if $X \sim \text{Bin}(n, p)$ we have

$$\begin{aligned} G_X(s) &= \sum_{k=0}^n P(X=k) \cdot s^k \\ &= \sum_{k=0}^n \binom{n}{k} p^k q^{n-k} \cdot s^k \\ &= \sum_{k=0}^n \binom{n}{k} (ps)^k q^{n-k} \\ &= (ps + q)^n. \end{aligned}$$

(ii) if $X \sim \text{Po}(\lambda)$ we have

$$\begin{aligned} G_X(s) &= \sum_{k=0}^{\infty} P(X=k) s^k \\ &= \sum_{k=0}^{\infty} \frac{e^{-\lambda} \cdot \lambda^k}{k!} \cdot s^k \\ &= (*) \end{aligned}$$

$$\begin{aligned}
 (*) &= e^{-\lambda} \cdot \sum_{k=0}^{\infty} \frac{(\lambda s)^k}{k!} \\
 &= e^{-\lambda} \cdot e^{\lambda s} \\
 &= e^{-\lambda(1-s)}
 \end{aligned}$$

(iii) Let $X \sim \text{Neg Bin}(m, p)$.

From analysis we have that
for $|x| < 1$

$$(1+x)^a = \sum_{k=0}^{\infty} \binom{a}{k} x^k \quad \text{where}$$

$$\binom{a}{k} = \frac{a(a-1)\cdots(a-k+1)}{k!}$$

The above formula is known as the Newton formula. Replace x by $-x$ and let $a = -r$ for some integer $r > 0$.

We get

$$\begin{aligned}(1-x)^{-r} &= \sum_{k=0}^{\infty} \binom{-r}{k} \cdot (-x)^k \\ &= \sum_{k=0}^{\infty} \frac{(-r)(-r-1)\dots(-r-k+1)}{k!} (-1)^k \cdot x^k \\ &= \sum_{k=0}^{\infty} \frac{r(r+1)\dots(r+k-1)}{k!} \cdot x^k \\ &= \sum_{k=0}^{\infty} \frac{(r-1)! r(r+1)\dots(r+k-1)}{(r-1)! \cdot k!} x^k \\ &= \sum_{k=0}^{\infty} \frac{(r+k-1)!}{(r-1)! \cdot k!} x^k \\ &= \sum_{k=0}^{\infty} \binom{r+k-1}{r-1} x^k.\end{aligned}$$

We compute

$$\begin{aligned}G_X(s) &= \sum_{k=m}^{\infty} P(X=k) \cdot s^k \\ &= \sum_{k=m}^{\infty} \binom{k-1}{m-1} \cdot p^m \cdot q^{k-m} \cdot s^k \\ &= (*)\end{aligned}$$

$$(*) = \sum_{l=0}^{\infty} \binom{m+l-1}{m-1} p^m \cdot q^l s^{m+l}$$

$$= p^m \cdot s^m \cdot \sum_{l=0}^{\infty} \binom{m+l-1}{m-1} s^l \cdot q^l$$

$$= \frac{p^m \cdot s^m}{(1-qs)^m}$$

$$= \left(\frac{ps}{1-qs} \right)^m$$

(iv) The computation in previous example gives

$$(1-x)^{-a} = \sum_{k=0}^{\infty} \frac{(a)_k}{k!} \cdot x^k \quad |x| < 1.$$

Let X have the Po'lya distribution

$$P(X=k) = \frac{\beta^a (a)_k}{k! (1+\beta)^{a+k}}$$

We have

$$\begin{aligned}G_X(s) &= \sum_{k=0}^{\infty} P(X=k) \cdot s^k \\&= \sum_{k=0}^{\infty} \frac{\beta^a (a)_k}{k! (1+\beta)^{a+k}} \cdot s^k \\&= \sum_{k=0}^{\infty} \frac{\beta^a (a)_k}{k! (1+\beta)^a} \left(\frac{s}{1+\beta}\right)^k \\&= \frac{\beta^a}{(1+\beta)^a} \cdot \left(1 - \frac{s}{1+\beta}\right)^{-a} \\&= \left(\frac{\beta}{1+\beta-s}\right)^a\end{aligned}$$

Theorem 5.1: Let X be a nonnegative integer valued random variable and let $G_X(s)$ be its generating function.

Then $G_X(s)$ uniquely determines the distribution of X .

Proof: Since $G_X(s)$ converges for $|s| < 1$ we have

$$G_X^{(n)}(0) = n! P(X=n).$$

Theorem 6.2: Let X be an integer valued random variable with generating function $G_X(s)$.

(i)

$$E(X) = \lim_{s \uparrow 1} G_X'(s)$$

(ii)

$$E[X(X-1)\dots(X-m+1)]$$

$$= \lim_{s \uparrow 1} G_X^{(m)}(s)$$

Proof: Let $\varepsilon > 0$ and assume first that $E(X) < \infty$.

There is a N_ε such that for $n \geq N_\varepsilon$ we have $\sum_{k=n}^{\infty} k P(X=k) < \varepsilon$.

This means that

$$E(X) - \sum_{k=0}^{N_\varepsilon-1} k P(X=k) < \varepsilon.$$

Since all the coefficients in the power series are non-negative we

have that for $s \in (0, 1)$

$$\sum_{k=0}^{N_\varepsilon-1} k P(X=k) s^{k-1} \leq G'_X(s) \leq E(X).$$

As $s \uparrow 1$ we have

$$\sum_{k=0}^{N_\varepsilon-1} k P(X=k) \leq \lim_{s \uparrow 1} G'_X(s) \leq E(X)$$

The limit exists because $G_X(s)$ is nondecreasing on $(0, 1)$. But the above means that

$$E(X) - \varepsilon \leq \lim_{s \uparrow 1} G'_X(s) \leq E(X)$$

for all $\varepsilon > 0$.

If $E(X) = \infty$ we have that

$$\lim_{n \uparrow \infty} G'_X(s) \geq \sum_{k=1}^n k P(X=k)$$

for any finite n . This implies

that $\lim_{s \uparrow 1} G'_X(s) = \infty$.

(ii) The proof is similar.

Theorem 5.3: Let X, Y be independent.

Then

$$G_{X+Y}(s) = G_X(s) \cdot G_Y(s)$$

Proof: We can write

$$E(s^X) = \sum_{k=0}^{\infty} s^k \cdot P(X=k) = G_X(s)$$

Then

$$\begin{aligned} E(s^{X+Y}) &= E(s^X \cdot s^Y) \\ &\stackrel{\text{indp}}{=} E(s^X) \cdot E(s^Y) \end{aligned}$$

$$= G_X(s) \cdot G_Y(s)$$

Comment: This is the most important property of generating functions.

By extension we have for

○ independent X_1, X_2, \dots, X_r

$$G_{X_1 + X_2 + \dots + X_r}(s) = G_{X_1}(s) \cdot G_{X_2}(s) \cdots G_{X_r}(s)$$

Examples:

○ (i) X, Y independent $X \sim \text{Bin}(m, p)$ and $Y \sim \text{Bin}(n, p)$. We have

$$\begin{aligned} G_{X+Y}(s) &= G_X(s) \cdot G_Y(s) \\ &= (ps+q)^m \cdot (ps+q)^n \\ &= (ps+q)^{m+n} \end{aligned}$$

This last function is the generating function of the Bin $(m+n, p)$ distribution. It follows, by uniqueness, $X+Y \sim \text{Bin}(m+n, p)$.

(ii) Let X, Y be independent and

$$P(X=k) = \frac{\beta^a (a)_k}{k! (1+\beta)^{a+k}}, \quad k=0,1,\dots$$

$$P(Y=l) = \frac{\beta^b (b)_l}{l! (1+\beta)^{b+l}}, \quad l=0,1,\dots$$

We have

$$\begin{aligned} G_{X+Y}(s) &= G_X(s) \cdot G_Y(s) \\ &= \left(\frac{\beta}{1+\beta-s} \right)^a \cdot \left(\frac{\beta}{1+\beta-s} \right)^b \\ &= \left(\frac{\beta}{1+\beta-s} \right)^{a+b} \end{aligned}$$

Conclusion:

$$P(X+Y=k) = \frac{\beta^{a+b} (a+b)_k}{k! (1+\beta)^{a+b}}, \quad k=0, 1, \dots$$

Comment: We have computed the distribution of $x+y$ before but the above is much more elegant.

(iii) Suppose X_1, X_2, \dots, X_r are independent and $X_i \sim \text{Geom}(p)$. We have

$$\begin{aligned} G_{X_i}(s) &= \sum_{i=0}^{\infty} s^i \cdot 2^{i-1} \cdot p \\ &= ps \sum_{i=0}^{\infty} (2s)^i \\ &= \frac{ps}{1-2s} \end{aligned}$$

It follows

$$G_{X_1 + \dots + X_r}(s) = \left(\frac{ps}{1-2s} \right)^r$$

Conclusion: $X_1 + \dots + X_r \sim \text{Neg Bin}(r, p)$.

5.2. Branching processes

In applications of probability we often calculate sums of a random number of random variables. Let

X_1, X_2, \dots be random variables and

N an integer valued nonnegative

random variable. We need to

define $X_1 + X_2 + \dots + X_N$. Formally

we define

$$X = \sum_{k=1}^{\infty} X_k \cdot \mathbb{1}(N \geq k)$$

Comment:

For a fixed $\omega \in \Omega$

we have $N(\omega) < \infty$ and so the

sum is finite because only

a finitely many terms are $\neq 0$.

We will write

$$X = X_1 + X_2 + \dots + X_N.$$

Theorem 5.4 : let N, X_1, X_2, \dots be independent, X_1, X_2, \dots equally distributed non-negative integer valued random variables and N non-negative integer valued. Let $X = X_1 + X_2 + \dots + X_N$. Then

$$G_X(s) = G_N(G_{X_1}(s)).$$

Proof: We use the formula for the total expectation.

$$\begin{aligned} G_X(s) &= E(s^X) \\ &= \sum_{k=0}^{\infty} E(s^X | N=k) P(N=k) \\ &= \sum_{k=0}^{\infty} E(s^{X_1 + \dots + X_k} | N=k) P(N=k) \\ &\stackrel{\text{indep}}{=} \sum_{k=0}^{\infty} E(s^{X_1 + \dots + X_k}) P(N=k) \\ &= (*) \end{aligned}$$

$$\begin{aligned}
 \phi(x) &= \sum_{k=0}^{\infty} [G_{X_1}(s)]^k \cdot P(N=k) \\
 &= G_N(G_{X_1}(s)).
 \end{aligned}$$

Example : A hen lays N eggs. A chick hatches from each egg with probability p independent of all other eggs. Suppose $N \sim Po(\lambda)$. What is the distribution of the number of chicks? In mathematical notation we are asking about the distribution of $I_1 + I_2 + \dots + I_N$ where I_1, I_2, \dots are independent with $I_k \sim \text{Bernoulli}(p)$ and independent of N . We have

$$G_X(s) = G_N(G_{X_1}(s)).$$

$$\begin{aligned}
 G_{I_1}(s) &= s^0 P(I_1=0) + s^1 \cdot P(I_1=1) \\
 &= q + ps
 \end{aligned}$$

We have

$$\begin{aligned}G_X(s) &= G_N(p+qs) \\ &= e^{-\lambda(1-p-qs)} \\ &= e^{-\lambda(p+qs)} \\ &= e^{-\lambda p(1-s)}\end{aligned}$$

This last function is the generating function of the $P(\lambda p)$ distribution so $X \sim Po(\lambda p)$.

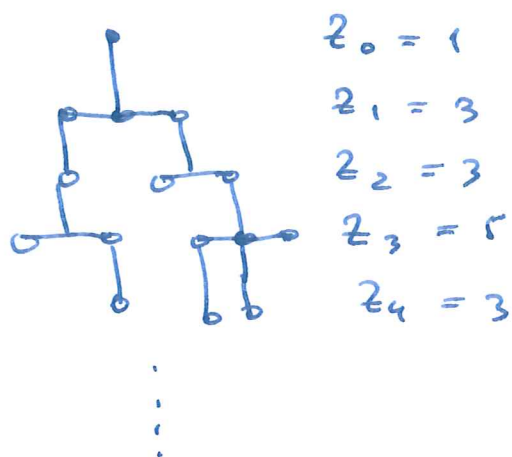
Branching processes

In 1874 Sir Francis Galton (1822-1911) asked the following question:

suppose you take an English aristocrat. He will have a random number of sons. His sons will have a random number of sons, ...

The problem is to determine the probability that the family tree will die out.

Figure: A possible family tree



The problem was solved by Galton and Watson in 1875 (F. Galton, H.W.

Watson, *Proceedings of the Royal Society*, (1875)

On the probability of the extinction of families, *Journal of the Royal Anthropological Institute* 4, (1875) 138-144) using generating functions.

To solve the problem mathematically we need a few additional assumptions:

- (i) Generations are simultaneous.
- (ii) Each individual has sons independently of all the others.
- (iii) The random number of sons has the same distribution for all individuals.

The above assumptions imply the following mathematical formulation:

Let $\{\xi_{n,k}\}_{n \geq 1, k \geq 1}$ be independent, equally distributed non-negative integer valued random variables with generating function G .

We define

$$Z_0 = 1 \quad \text{and recursively}$$

$$Z_{n+1} = \xi_{n+1,1} + \xi_{n+1,2} + \dots + \xi_{n+1,Z_n}$$

The sequence Z_0, Z_1, \dots of random variables is called the branching process.

The above means that Z_u individuals in the u -th generation have randomly many offspring.

The random variable Z_u depends on $\xi_{m,k}$ for $m \leq u$ so it is independent of

$$\xi_{u+1,1}, \xi_{u+1,2}, \dots$$

Denote $G_u(s) = G_{Z_u}(s)$. By Theorem 5.4 we have

$$G_{u+1}(s) = G_u(G(s)).$$

By definition $G_1(s) = G(s)$ and by the above recursion

$$G_2(s) = G_1(G(s)) = (G \circ G)(s)$$

$$G_3(s) = G_2(G(s)) = (G \circ G \circ G)(s)$$

\vdots

$$G_u(s) = (G \circ G \circ \dots \circ G)(s).$$

Since composition is associative we have

$$G_{n+1}(s) = G(G_n(s))$$

Let $A = \{\text{the family tree dies out}\}$.

The family tree dies out if one of the generations is empty so

$$A = \bigcup_{n=1}^{\infty} \{z_n = 0\}.$$

But $\{z_1 = 0\} \subseteq \{z_2 = 0\} \subseteq \dots$

In the first chapter we proved that for $A_1 \subseteq A_2 \subseteq \dots$ we have

$$P\left(\bigcup_{n=1}^{\infty} A_n\right) = \lim_{n \rightarrow \infty} P(A_n).$$

Denote $\eta = P(A)$. We have

$$\eta = P(A) = \lim_{n \rightarrow \infty} P(z_n = 0).$$

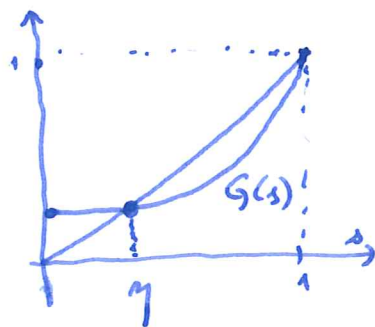
But $P(z_n = 0) = G_n(0)$.

Theorem 5.5 : The probability η satisfies the equation $\eta = G(\eta)$ and is the smallest solution of the above equation on $[0, 1]$.

Comments :

- (i) If $\eta = G(\eta)$ we say that η is a fixed point of G .
- (ii) Since $G(1) = 1$ there is always at least one fixed point on $[0, 1]$. The set of fixed points is compact so it contains a smallest point.

Figure :



Proof : $G(s)$ is continuous on $[0, 1]$.

So we have

$$\begin{aligned}\eta &= \lim_{n \rightarrow \infty} G_{n+1}(0) = \lim_{n \rightarrow \infty} G(G_n(0)) \\ &= G\left(\lim_{n \rightarrow \infty} G_n(0)\right) = G(\eta).\end{aligned}$$

So η is a fixed point. To

prove that η is the smallest

fixed point let $\bar{\eta}$ be a fixed point on $[0,1]$. We have

$$0 \leq \bar{\eta}$$

Since G is nondecreasing on $[0,1]$ it follows

$$G(0) \leq G(\bar{\eta}) = \bar{\eta}$$

$$G(G(0)) \leq G(\bar{\eta}) = \bar{\eta}$$

\vdots

$$(G \circ \dots \circ G)(0) = G_n(0) \leq \bar{\eta}$$

So

$$\lim_{n \rightarrow \infty} G_n(0) = \eta \leq \bar{\eta}.$$

This means that any fixed point \bar{y} is $\geq y$ which proves the theorem.

Example: Suppose every individual has 0, 1, 2, 3 sons with probability $1/4$ each. This means that

$$G(s) = \frac{1 + s + s^2 + s^3}{4}$$

We need all solutions of

$$G(s) = s \Leftrightarrow 1 - 3s + s^2 + s^3 = 0$$

We know that $s = 1$ is a solution so we can factor

$$1 - 3s + s^2 + s^3 = (s-1)(s^2 + 2s - 1)$$

The solutions are

$$s = 1$$

$$s = -1 + \sqrt{2}$$

$$s = -1 - \sqrt{2}$$

The smallest

fixed point

on $[0, 1]$ is $-1 + \sqrt{2}$

$$\approx 0.4142.$$

Example: Suppose $G(s) = \frac{p}{1-qs}$.

What is $G_n(s)$?

$$\begin{aligned}G_2(s) &= G(G(s)) \\&= \frac{p}{1-q \cdot \frac{p}{1-qs}} \\&= \frac{p(1-qs)}{1-pq-qs}\end{aligned}$$

$$\begin{aligned}G_3(s) &= \frac{p}{1-q \cdot \frac{p(1-qs)}{1-pq-qs}} \\&= \frac{p(1-pq-qs)}{1-pq-qs-pq(1-qs)} \\&= \frac{p(1-pq-qs)}{1-2pq-qs(1-pq)}\end{aligned}$$

We see that

$$G_n(s) = \frac{a_n - b_n s}{c_n - d_n s}$$

$$G_{n+1}(s) = \frac{a_n - b_n \cdot \frac{p}{1-qs}}{c_n - d_n \cdot \frac{p}{1-qs}}$$

Multiplying out we get

$$a_{n+1} = a_n - p \cdot b_n$$

$$b_{n+1} = q a_n$$

We have $f_0(x) = 1 \Rightarrow a_0 = 0, b_0 = -1$

Write in matrix form

$$\begin{pmatrix} a_{n+1} \\ b_{n+1} \end{pmatrix} = \begin{pmatrix} 1 & -p \\ q & 0 \end{pmatrix} \begin{pmatrix} a_n \\ b_n \end{pmatrix}.$$

Iteration gives

$$\begin{aligned} \begin{pmatrix} a_n \\ b_n \end{pmatrix} &= \begin{pmatrix} 1 & -p \\ q & 0 \end{pmatrix}^n \begin{pmatrix} a_0 \\ b_0 \end{pmatrix} \\ &= \begin{pmatrix} 1 & -p \\ q & 0 \end{pmatrix}^n \begin{pmatrix} 0 \\ -1 \end{pmatrix} \end{aligned}$$

We need to find the power of the matrix. Suppose $p \neq q$.

We can check by multiplication that

$$\begin{aligned} \underbrace{\begin{pmatrix} p & 1 \\ q & 1 \end{pmatrix}}_A \begin{pmatrix} p & 0 \\ 0 & q \end{pmatrix} \underbrace{\begin{pmatrix} \frac{1}{p-q} & -\frac{1}{p-q} \\ -\frac{q}{p-q} & \frac{p}{p-q} \end{pmatrix}}_{A^{-1}} \\ = \begin{pmatrix} 1 & -p \\ q & 0 \end{pmatrix} \end{aligned}$$

We have diagonalized the matrix

$\begin{pmatrix} 1 & -p \\ \ell & 0 \end{pmatrix}$. This means

$$\begin{aligned} \begin{pmatrix} 1 & -p \\ \ell & 0 \end{pmatrix}^n &= \begin{pmatrix} p & 1 \\ \ell & 1 \end{pmatrix} \begin{pmatrix} p^n & 0 \\ 0 & \ell^n \end{pmatrix} \begin{pmatrix} 1 & -1 \\ -\ell & p \end{pmatrix} \cdot \frac{1}{p-\ell} \\ &= \begin{pmatrix} p^{n+1} - \ell^{n+1} & -p^{n+1} + p\ell^n \\ \ell p^n - \ell^{n+1} & -\ell p^n + p\ell^n \end{pmatrix} \cdot \frac{1}{p-\ell} \end{aligned}$$

We find:

$$\begin{aligned} \begin{pmatrix} 1 & -p \\ \ell & 0 \end{pmatrix}^n \begin{pmatrix} 0 \\ -1 \end{pmatrix} &= \\ &= \frac{1}{p-\ell} \begin{pmatrix} p(\ell^n - p^n) \\ \ell \ell (p^{n-1} - \ell^{n-1}) \end{pmatrix} = \begin{pmatrix} a_n \\ b_n \end{pmatrix} \end{aligned}$$

For c_n, d_n the procedure is the same except that $\begin{pmatrix} c_0 \\ d_0 \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$

and

$$\begin{pmatrix} c_n \\ d_n \end{pmatrix} = \begin{pmatrix} p^{n+1} - \ell^{n+1} \\ \ell (p^n - \ell^n) \end{pmatrix} \cdot \frac{1}{p-\ell}$$

Finally,

$$G_n(s) = \frac{p(p^n - q^n) - qs(p^{n-1} - q^{n-1})}{p^{n+1} - q^{n+1} - qs(p^n - q^n)}$$

We find

$$\begin{aligned} P(z_n = 0) &= G_n(0) \\ &= \frac{p(p^n - q^n)}{p^{n+1} - q^{n+1}} \end{aligned}$$

We get: if $p > q$, then

$$\lim_{n \rightarrow \infty} P(z_n = 0) = 1$$

if $p < q$

$$\lim_{n \rightarrow \infty} P(z_n = 0) = \frac{p}{q} < 1.$$

Comment: The fixed points

satisfy $\frac{p}{1 - qs} = s \Rightarrow$

$$qs^2 - s + p = (s-1)(qs-p) = 0 \Rightarrow$$

There is a fixed point in $(0,1)$ other than $\frac{1}{2}$ if $p < q$.

Comment : For $p = q = \frac{1}{2}$ we get

$$G_n(s) = \frac{n - (n-1)s}{n+1 - ns}$$

and

$$G_n(0) = \frac{n}{n+1} \rightarrow 1, \text{ as } n \rightarrow \infty.$$

Theorem 5.6: Let z_0, z_1, z_2, \dots be a branching process. Let $\mu = E(z_1)$.

- (i) If $\mu < 1$ then $\eta = 1$.
- (ii) If $\mu > 1$ then $\eta \in (0, 1)$.
- (iii) If $\mu = 1$ and $G(s) \neq s$ then $\eta = 1$.

Comment: The case $G(s) = s$ is

uninteresting.

The proof of the theorem has more to do with analysis than probability.

Proof: Note that $\mu = \lim_{s \uparrow 1} G'(s)$.

The function $G'(s)$ is nondecreasing

on $(0, 1)$.

(i) If $\mu < 1$, then $G'(s) \leq \mu < 1$ for all $s \in (0, 1)$. If $G(\bar{\eta}) = \bar{\eta}$

for $\bar{\eta} \in (0, 1)$ then by

Lagrange theorem

$$\begin{aligned} G(1) - G(\bar{\eta}) &= G'(\xi)(1 - \bar{\eta}) \quad \text{for } \xi \in (\bar{\eta}, 1). \\ \text{"} & \\ 1 - \bar{\eta} & \end{aligned}$$

But $G'(\xi) < 1$. So we have a contradiction.

(ii) If $\mu > 1$ there is a $\delta > 0$ such that $G'(s) > 1$ for $s \in (1-\delta, 1)$. By Lagrange for $s \in (1-\delta, 1)$ we have

$$G(1) - G(s) = G'(\xi)(1-s) \geq (1-s)$$

for $\xi \in (s, 1)$. This implies

$$1 - G(s) > 1 - s \Rightarrow G(s) < s.$$

On the other hand $G(0) \geq 0$.

For $s \in (1-\delta, 0)$ we have that

$$F(0) = G(0) - 0 \geq 0$$

$$F(s) = G(s) - s < 0.$$

There must be a zero of F on $(0, s) \subset (0, 1) \Rightarrow \eta \in (0, 1)$.

(iii) If $g(s) \neq 1$ then either $g(s) = 1$
in which case $\eta = 1$ or g is
strictly convex on $(0, 1)$, or
 $g'(s)$ is strictly increasing on $(0, 1)$.

If $g(\bar{\eta}) = \bar{\eta}$ for some $\bar{\eta} \in (0, 1)$
this would imply

$$\begin{aligned}g(1) - g(\bar{\eta}) &= 1 - \bar{\eta} \\ &= g'(\xi)(1 - \bar{\eta})\end{aligned}$$

for some $\xi \in (\bar{\eta}, 1)$. This means

$g'(\xi) = 1$. But $g'(s)$ is strictly
increasing meaning $\lim_{s \rightarrow 1} g'(s) > 1$.

A contradiction.

Paujer recursion

If $X = X_1 + X_2 + \dots + X_N$ we know that

$$G_X(s) = G_N(G_{X_1}(s)).$$

In principle we get $P(X=k)$

by expanding the right side into power series. But this is often difficult and recursive formulae are needed. This problem is often dealt with in insurance.

Definition: The random variable N is of Paujer class if

$$P(N=u) = \left(a + \frac{b}{u}\right) P(N=u-1)$$

for $u = 1, 2, \dots$

Examples: (i) Take $a = 0$ and $b > 0$.

$$\text{We get } P(N=u) = \frac{b}{u} P(N=1) \Rightarrow$$

$$P(N=u) = e^{-b} \cdot \frac{b^u}{u!} \Rightarrow N \sim Po(b).$$

(ii) Suppose $N \sim \text{Bin}(M, p)$. We have
compu test

$$\frac{P(N=u)}{P(N=u-1)} = \frac{M-u+1}{u} \cdot \frac{p}{q}$$

$$= \left(-\frac{p}{q} + \frac{(M+1)p}{q \cdot u} \right)$$

We take $a = -\frac{p}{q}$, $b = \frac{(M+1)p}{q}$

We see that $P(N=M+1) = 0$.

We compute

$$P(N=1) = \left(-\frac{p}{q} + \frac{(M+1)p}{q} \right) P(N=0)$$

$$= \frac{p}{q} M \cdot P(N=0)$$

$$P(N=2) = \left(-\frac{p}{q} + \frac{(M+1)p}{2q} \right) P(N=1)$$

$$= \frac{p}{q} \left(-1 + \frac{M+1}{2} \right) P(N=1)$$

$$= \frac{p}{q} \cdot \frac{M-1}{2} \cdot \frac{p}{q} \cdot \frac{M}{1}$$

Continuing we get

$$P(N=u) = \left(\frac{p}{z}\right)^u \cdot \frac{M(M-1)\cdots(M-u+1)}{u!} P(N=0)$$
$$= \left(\frac{p}{z}\right)^u \binom{M}{u} \cdot P(N=0)$$

Because all probabilities must add to 1 we have

$$\left(1 + \frac{p}{z}\right)^M P(N=0) = 1 \Rightarrow$$

$$P(N=0) = z^{-M}$$

or

$$P(N=u) = \binom{M}{u} p^u z^{M-u}$$

Conclusion: N is in the Poisson class.

Theorem 5.7: For $|s| < 1$ the generating function of N satisfies

$$(1-as) G_N'(s) = (a+b) G_N(s).$$

Proof: From the recursion equation we have

$$P(N=n) \cdot s^n = \left(a + \frac{b}{n}\right) P(N=n-1) \cdot s^{n-1}$$

Sum both sides over $n = 1, 2, \dots$

We get

$$G_N(s) - G_N(0)$$

$$= a \sum_{n=1}^{\infty} P(N=n-1) s^{n-1}$$

$$+ b \cdot \sum_{n=1}^{\infty} \frac{s^{n-1}}{n} P(N=n-1)$$

$$= as G_N(s) + b \sum_{n=0}^{\infty} \left(\int_0^1 u^n du \right) P(N=n)$$

$$= as G_N(s) + b \cdot \int_0^1 \left(\sum_{n=0}^{\infty} u^n P(N=n) \right) du$$

$$= as G_N(s) + b \int_0^1 G_N(u) du.$$

It is legitimate to interchange summation and integration because the sum converges uniformly on $[0, 1]$.

Take derivatives to get

$$G'_N(s) = a G_N(s) + as G'_N(s) + b G_N(s).$$

Rearranging gives the equation.

If $X = X_1 + X_2 + \dots + X_N$ where X_1, X_2 are independent equally distributed we get

$$G_X(s) = G_N(G_{X_1}(s))$$

Taking derivatives we get

$$G'_X(s) = G'_N(G_{X_1}(s)) \cdot G'_{X_1}(s)$$

Multiply both sides by

$$1 - aG_{X_1}(G_{X_1}(s))$$

and use Theorem 5.7. We get

$$\begin{aligned} G'_X(s) (1 - a G_{X_1}(s)) \\ = (a+b) G_X(s) G'_{X_1}(s) \end{aligned}$$

Denote $P(N=n) = p_n$ for $n=0, 1, \dots$

Denote $P(X=r) = g_r$ and

$P(X_1=k) = f_k$ for $k=0, 1, \dots$

We have that $X=0$ if either
 $N=0$ or $N>0$ and $X_1+\dots+X_N=0$.

It follows

$$P(X=0) = P(N=0) + \sum_{n=1}^{\infty} P(N=n) \cdot f_0^n$$

$$= G_N(f_0)$$

In our notation

$$P(X=0) = g_0 = G_N(f_0).$$

From Analysis we know that

$$\left(\sum_{k=0}^{\infty} a_k x^k \right) \left(\sum_{k=0}^{\infty} b_k x^k \right) = \sum_{k=0}^{\infty} c_k x^k$$

with

$$c_k = \sum_{i=0}^k a_i b_{k-i}.$$

Comment: This is called the Cauchy product of power series.

In the formulae connecting the generating functions

Equate coefficients for s^n .

$$h: (n+1)g_{n+1} = a \sum_{k=0}^n f_k \cdot (n+1-k)g_{n+1-k}$$

$$\begin{aligned} 2: (a+b) \sum_{k=0}^n (k+1) f_{k+1} \cdot g_{n-k} \\ = (a+b) \sum_{k=1}^{n+1} k f_k \cdot g_{n-k+1} \end{aligned}$$

Rearranging we get

$$(n+1)g_{n+1} - a f_0 (n+1)g_{n+1}$$

$$= a \sum_{k=1}^n f_k (n+1-k) g_{n+1-k} +$$

$$(a+b) \sum_{k=0}^n (k+1) f_{k+1} g_{n-k}$$

$$= a \sum_{k=1}^n f_k (n+1-k) g_{n+1-k}$$

$$+ (a+b) \sum_{k=1}^{n+1} k f_k \cdot g_{n+1-k}$$

$$= a \cdot \sum_{k=1}^{n+1} f_k (n+1-k) g_{n+1-k}$$

$$+ (a+b) \sum_{k=1}^{n+1} k f_k g_{n+1-k}$$

Divide by $(n+1)(1-af_0)$ to get

$$g_{n+1} = \frac{1}{1-af_0} \sum_{k=1}^{n+1} \left(a + \frac{bk}{n+1} \right) f_k g_{n+1-k}$$

Theorem 5.8 (Paujer recursion)

We have

$$g_{n+1} = \frac{1}{1-af_0} \sum_{k=1}^{n+1} \left(a + \frac{bk}{n+1} \right) f_k g_{n+1-k}$$

G. The central limit theorem

We will be interested in the distribution of sums $S_n = X_1 + X_2 + \dots + X_n$.

For some types of distributions we know the answer but in general this question is difficult to answer.

Moreover, in statistics we need to approximate such distributions even if we do not exactly know the distributions of X_1, X_2, \dots .

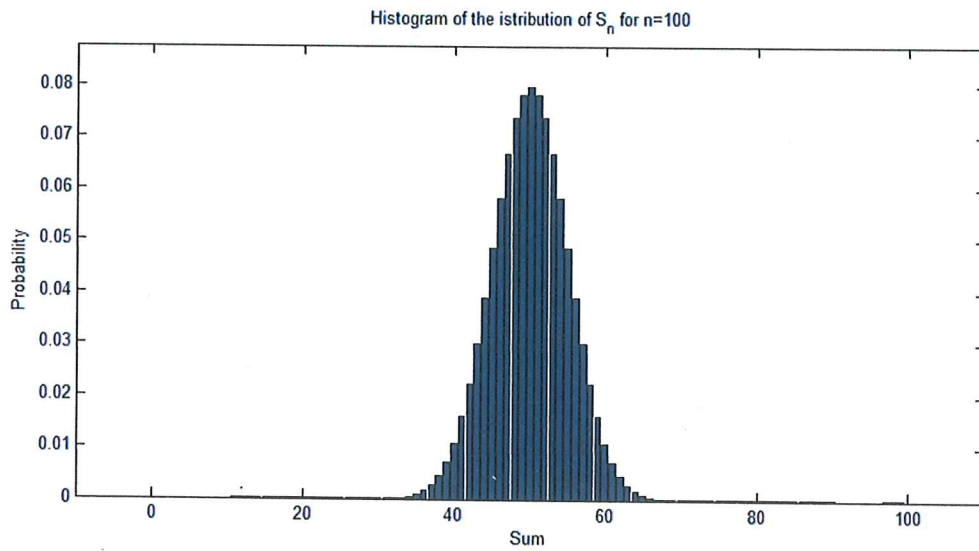
The setup we will look at will be: X_1, X_2, \dots are independent, equally distributed random variables.

We denote $S_n = X_1 + X_2 + \dots + X_n$.

Let us look at examples of distributions of S_n for simple distributions.

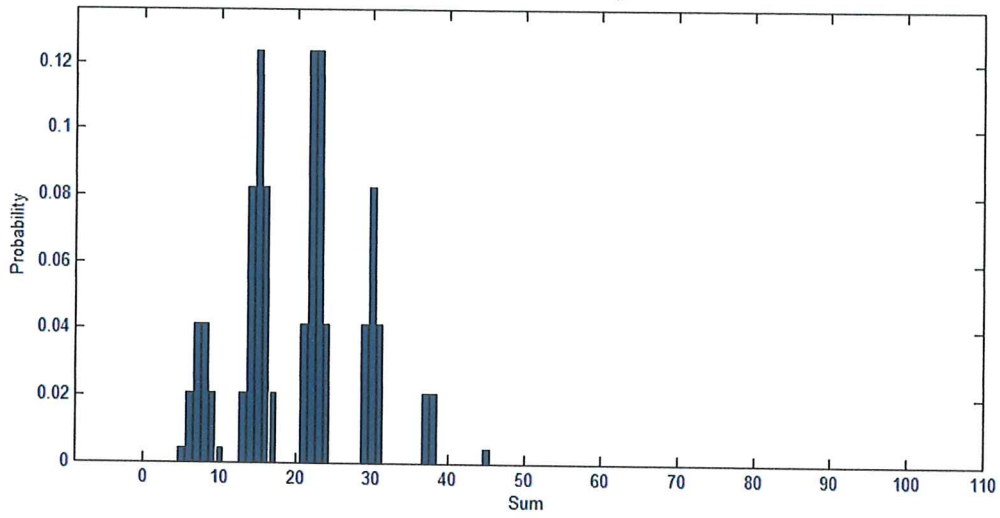
We will look at a few examples of distributions of S_n for different distributions of X_1 and different n .

1. Let $P(X_1 = 0) = P(X_1 = 1) = \frac{1}{2}$. Take $n = 100$. Let $S_n = X_1 + \dots + X_n$. The histogram of the distribution of S_n is:

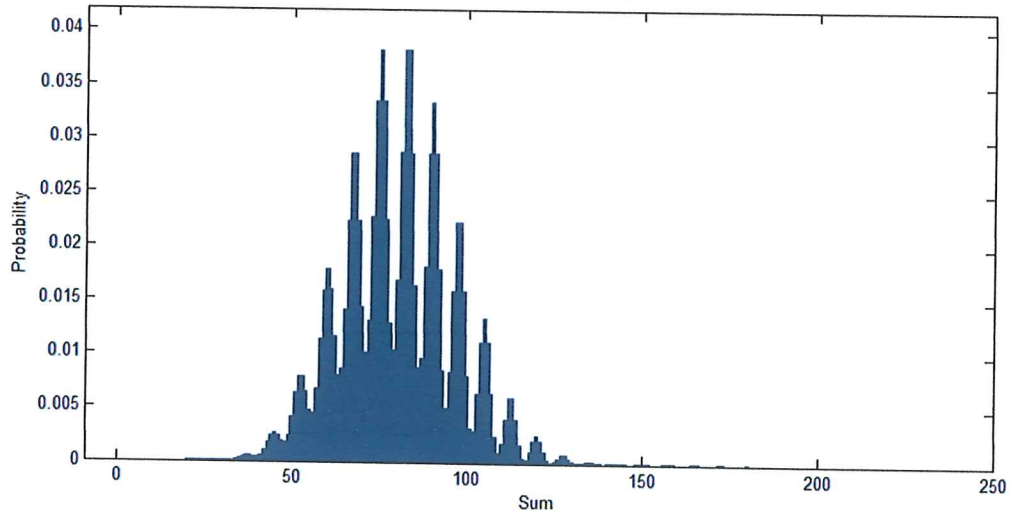


2. Let $P(X_1 = 1) = P(X_1 = 2) = P(X_1 = 9) = \frac{1}{3}$. Let $n = 5, 20, 50, 200$.

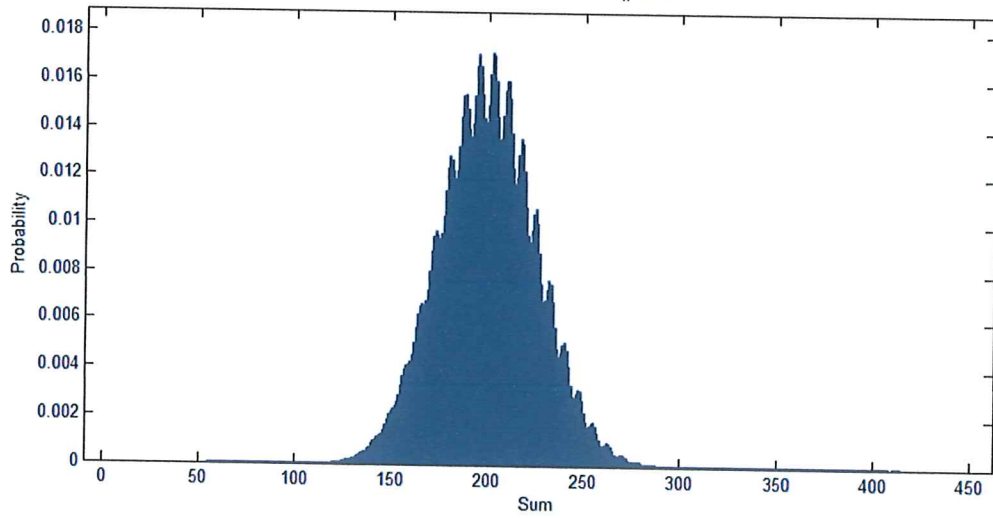
Histogram of the istribution of S_n for $n=5$

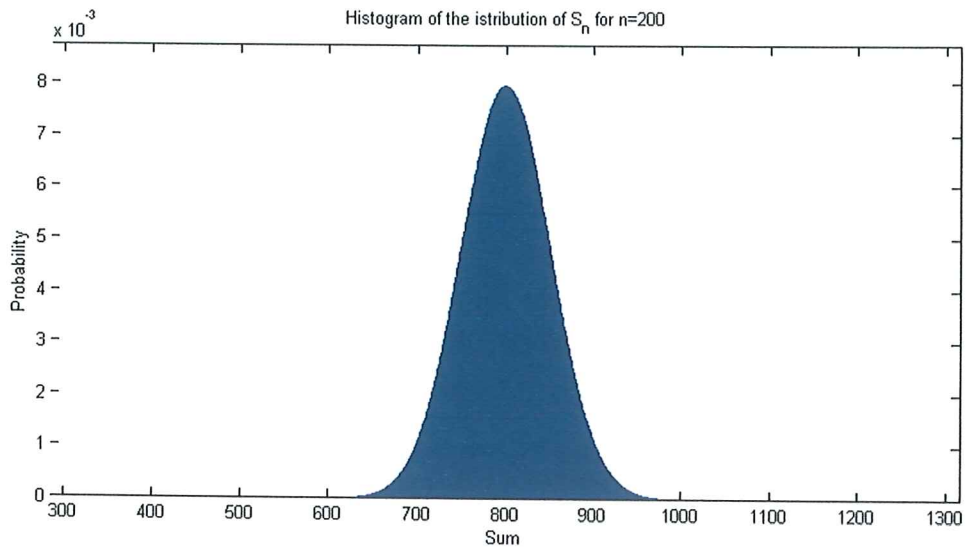


Histogram of the istribution of S_n for $n=20$

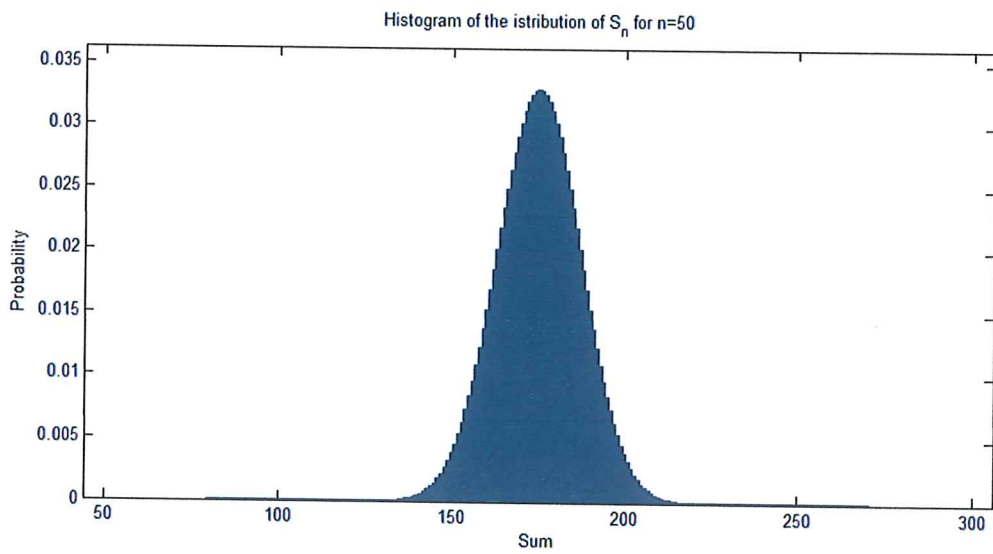


Histogram of the istribution of S_n for $n=50$

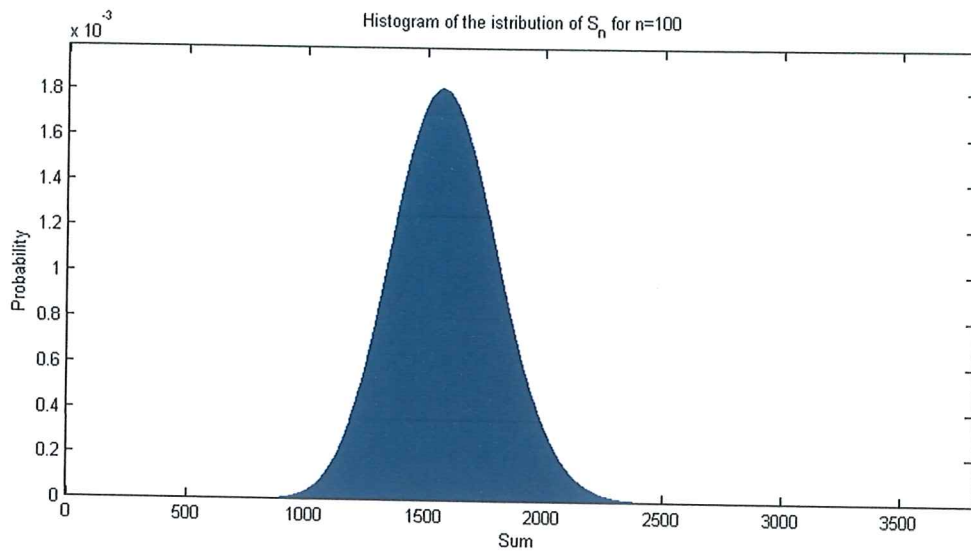




3. Take $P(X_1 = k) = 1/6$ for $k = 1, 2, \dots, 6$. Take $n = 50$.



4. Take $P(X_1 = 2^k) = 1/7$ for $k = 0, 1, 2, 3, 4, 5, 6$. Take $n = 100$.



From the examples we infer that the distribution of S_n is similar to the normal distribution. This is not clear in itself. But if we accept this observation we need to find the normal distribution that fits the histogram of S_n well.

Idea: We match the top of the two distributions which means that the mean of the normal distribution will be $E(S_n)$. We match the dispersion by choosing the second parameter to be $\text{var}(S_n)$.

Denote: $E(X_1) = \mu$ $\text{var}(X_1) = \sigma^2$

$$E(S_n) = n E(X_1) = n\mu$$

$$\text{var}(S_n) = n \text{var}(X_1) = n\sigma^2$$

To approximate the distribution
we can say

$$P(a \leq S_n \leq b) \approx \frac{1}{\sqrt{2\pi} \cdot \tau} \int_a^b e^{-\frac{(x-\nu)^2}{2\tau^2}} dx$$

The area of columns in the histogram between a and b ~~are~~ is exactly $P(a \leq S_n \leq b)$. We superimpose a curve closely following the histogram and replace the area of columns with the integral under the curve.

To turn the above into a mathematical theorem we will reformulate.

Take $a = \nu + \alpha \cdot \tau$ and $b = \nu + \beta \cdot \tau$.

We compute

$$P(a \leq S_u \leq b)$$

$$= P(\nu + \alpha \cdot \tau \leq S_u \leq \nu + \beta \cdot \tau)$$

$$\approx \frac{1}{\sqrt{2\pi} \tau} \int_{\nu + \alpha \cdot \tau}^{\nu + \beta \cdot \tau} e^{-\frac{(x - \nu)^2}{2\tau^2}} dx$$

New variable: $\frac{x - \nu}{\tau} = u$

$$= \frac{1}{\sqrt{2\pi}} \int_{\alpha}^{\beta} e^{-u^2/2} du.$$

On the other hand we have

$$P(\nu + \alpha \cdot \tau \leq S_u \leq \nu + \beta \cdot \tau)$$

$$= P\left(\alpha \leq \frac{S_u - \nu}{\tau} \leq \beta\right)$$

$$= P\left(\alpha \leq \frac{S_u - E(S_u)}{\sqrt{\text{var}(S_u)}} \leq \beta\right)$$

Definition: The expression

$$\tilde{S}_n = \frac{S_n - E(S_n)}{\sqrt{\text{var}(S_n)}}$$

is called the standardized sum.

We noticed that the approximation

is "better" if n is "large". We

expect the mathematical form

to include limits.

Theorem 6.1 (central limit theorem)

Let X_1, X_2, \dots be independent

equally distributed random

variables with $E(X_i) = \mu$ and

$\text{var}(X_i) = \sigma^2 < \infty$. Let $S_n = X_1 + \dots + X_n$.

For any $\alpha < \beta$ we have

$$\lim_{n \rightarrow \infty} P\left(\alpha \leq \frac{S_n - E(S_n)}{\sqrt{\text{var}(S_n)}} \leq \beta\right)$$

$$= \Phi(\beta) - \Phi(\alpha)$$

where Φ is the distribution function of the standard normal distribution.

Comments:

(i) we will prove the theorem in several steps. It is true as it is formulated but we will impose the additional assumption $E(|X_i|^3) < \infty$.

(ii) It is enough to prove

$$\lim_{n \rightarrow \infty} P\left(\frac{S_n - E(S_n)}{\sqrt{\text{var}(S_n)}} \leq \beta\right) = \Phi(\beta)$$

for $\beta \in \mathbb{R}$.

(iii) In the limit we get equality.

For finite n we use the limit as an approximation.

To prove the central limit theorem we need the following result.

Theorem 6.2 (Lindeberg-Bergman)

Let X_1, X_2, \dots, X_n be independent and such that $\text{var}(X_1 + X_2 + \dots + X_n) = 1$ and $E(X_1 + \dots + X_n) = 0$. Assume that $E(|X_k|^3) < \infty$ for all $k = 1, 2, \dots, n$.

Let f be a three times continuously differentiable function such that $|f(x)|, |f'(x)|, |f''(x)|, |f'''(x)| \leq M$ for some $M < \infty$ and all $x \in \mathbb{R}$.

Let $S_n = X_1 + X_2 + \dots + X_n$. Then

$$|E(f(S_n)) - E(f(Z))|$$

$$\leq \frac{1}{6} M \left(1 + \sqrt{\frac{8}{\pi}}\right) E(|X_1|^3 + \dots + |X_n|^3).$$

for $Z \sim N(0, 1)$.

Proof: Without loss of generality we can assume $E(X_k) = 0$ for all $k = 1, 2, \dots, n$. Let Z_1, Z_2, \dots, Z_n be independent and independent of X_1, X_2, \dots, X_n and such that $Z_k \sim N(0, \text{var}(X_k))$, $k = 1, 2, \dots, n$.

Since $\text{var}(X_1) + \dots + \text{var}(X_n) = 1$ by assumption we have that

$$Z = Z_1 + Z_2 + \dots + Z_n \sim N(0, 1).$$

Define

$$a_1 = E[f(Z_1 + Z_2 + \dots + Z_n)] - E[f(X_1 + Z_2 + \dots + Z_n)]$$

$$a_2 = E[f(X_1 + Z_2 + \dots + Z_n)] - E[f(X_1 + X_2 + \dots + Z_n)]$$

⋮

$$a_n = E[f(X_1 + \dots + X_{n-1} + Z_n)] - E[f(X_1 + X_2 + \dots + X_n)]$$

By triangle inequality we have

$$|E[f(X_1 + \dots + X_n)] - E[f(Z_1 + \dots + Z_n)]| \leq \sum_{k=1}^n |a_k|$$

By Taylor we have

$$f(x+h) - f(x) = f'(x) \cdot h + \frac{1}{2} f''(x) h^2 + r$$

where $r = \frac{1}{6} f'''(\xi) h^3$ for some

ξ between x and $x+h$. By our

assumption $|r| \leq \frac{1}{6} \cdot M \cdot |h|^3$.

○ Define

$$Y_1 = z_2 + z_3 + \dots + z_n$$

$$Y_2 = x_1 + z_3 + \dots + z_n$$

$$Y_3 = x_1 + x_2 + z_4 + \dots + z_n$$

$$Y_n = x_1 + x_2 + \dots + x_{n-1}$$

Note that Y_k is independent of (x_k, z_k) for all $k = 1, 2, \dots, n$.

We use Taylor's expansion around Y_k to get

$$E \left[f(x_1 + \dots + x_{k-1} + z_k + \dots + z_n) \right]$$

$$= E \left[f(y_k) + f'(y_k) z_k + \frac{1}{2} f''(y_k) z_k^2 + R_k \right]$$

and

$$E \left[f(x_1 + \dots + x_k + z_{k+1} + \dots + z_n) \right]$$

$$\stackrel{(1)}{=} E \left[f(y_k) + f'(y_k) x_k + \frac{1}{2} f''(y_k) x_k^2 + \tilde{R}_k \right]$$

Subtracting we get

$$a_k = E \left[f'(y_k)(z_k - x_k) + \frac{1}{2} f''(y_k)(z_k^2 - x_k^2) \right. \\ \left. + R_k - \tilde{R}_k \right]$$

$$\stackrel{(1)}{=} E \left[f'(y_k)(z_k - x_k) \right] \\ + E \left[\frac{1}{2} f''(y_k)(z_k^2 - x_k^2) \right] \\ + E \left[R_k - \tilde{R}_k \right]$$

By independence

$$\begin{aligned} E [f'(Y_k) (Z_k - X_k)] \\ &= E [f'(Y_k)] \underbrace{E [Z_k - X_k]}_{= 0 \text{ by assumption}} \\ &= 0 \end{aligned}$$

and

$$\begin{aligned} E [f''(Y_k) (Z_k^2 - X_k^2)] \\ &= E [f''(Y_k)] \underbrace{E [Z_k^2 - X_k^2]}_{= 0 \text{ by assumption}} \\ &= 0. \end{aligned}$$

We are left with

$$a_k = E [R_k - \tilde{R}_k]$$

Since $-x \leq |x|$ we have $|E(x)| \leq E(|x|)$

so

$$|a_k| \leq E [|R_k - \tilde{R}_k|]$$

But

$$|R_k| \leq \frac{1}{6} M \cdot |X_k|^3$$

$$|\tilde{R}_k| \leq \frac{1}{6} M |Z_k|^3$$

so

$$|R_k - \tilde{R}_k| \leq \frac{1}{6} M (|X_k|^3 + |Z_k|^3).$$

It follows

$$E[|R_k - \tilde{R}_k|] \leq \frac{1}{6} \cdot M (E(|X_k|^3) + E(|Z_k|^3)).$$

A standard calculation gives that for $Z \sim N(0, \sigma^2)$ we have

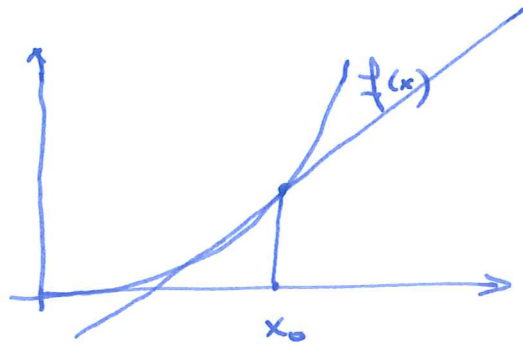
$$E(|Z|^3) = \sqrt{\frac{8}{\pi}} \cdot \sigma^3$$

We then have

$$\begin{aligned} E(|Z_k|^3) &= \sqrt{\frac{8}{\pi}} \text{var}(X_k)^{3/2} \\ &= \sqrt{\frac{8}{\pi}} \cdot E(X_k^2)^{3/2} \\ &= \sqrt{\frac{8}{\pi}} \cdot E(|X_k|^2)^{3/2} \end{aligned}$$

Since $f(x) = x^{3/2}$ is convex,
the function is above its tangent.

Figure :



We have for $x_0 > 0$

$$f(x) = x^{3/2} \geq \underbrace{f'(x_0)(x - x_0) + f(x_0)}_{\text{tangent}}$$

Take $x_0 = E(|X_k|^2)$. We have

$$(|X_k|^2)^{3/2} \geq f'(x_0)(|X_k|^2 - x_0) + x_0^{3/2}$$

Taking expectations we get

$$\begin{aligned} E(|X_k|^3) &\geq E(|X_k|^2)^{3/2} \\ &= (\text{var}(X_k))^{3/2} \end{aligned}$$

Taking all inequalities we get

$$E(|x_k|^3) + E(|z_k|^3)$$

$$\leq \left(1 + \sqrt{\frac{8}{\pi}}\right) E(|x_k|^3)$$

Finally,

$$\sum_{k=1}^n |a_k| \leq \sum_{k=1}^n \frac{1}{6} \cdot M \left(1 + \sqrt{\frac{8}{\pi}}\right) E(|x_k|^3) \quad \square$$

The inequality is valid for arbitrary x_1, x_2, \dots, x_n provided they are independent. If

x_1, x_2, \dots, x_n are independent and equally distributed define

$$x_k' = \frac{x_k - E(x_k)}{\sqrt{\text{var}(S_n)}}$$

for $S_n = x_1 + x_2 + \dots + x_n$. Note that

$$x_1' + x_2' + \dots + x_n' = \frac{S_n - E(S_n)}{\sqrt{\text{var}(S_n)}} = \tilde{S}_n$$

We have that X_1', X_2', \dots, X_n'
 are independent, $E(X_k') = 0$
 and $\text{var}(X_1' + \dots + X_n') = 1$.

Theorem 6.2 implies

$$|E[f(\hat{S}_n)] - E[f(z)]|$$

$$\leq \frac{1}{6} M \left(1 + \sqrt{\frac{8}{\pi}}\right) \cdot n \cdot E(|X_1'|^3)$$

But

$$E(|X_1'|^3) = E\left(\frac{|X_1 - E(X_1)|^3}{\sqrt{n} \sqrt{\text{var}(X_1)}}\right)$$

$$= \frac{1}{n^{3/2} \text{var}(X_1)^{3/2}} E(|X_1 - E(X_1)|^3)$$

If $\gamma = E(|X_1 - E(X_1)|^3) < \infty$ we have

$$|E[f(\hat{S}_n)] - E[f(z)]|$$

$$\leq \frac{1}{6} \cdot M \left(1 + \sqrt{\frac{8}{\pi}}\right) \cdot \frac{1}{\sqrt{n}} \cdot \frac{\gamma}{\text{var}(X_1)^{3/2}}$$

$$\rightarrow 0, \text{ as } n \rightarrow \infty.$$

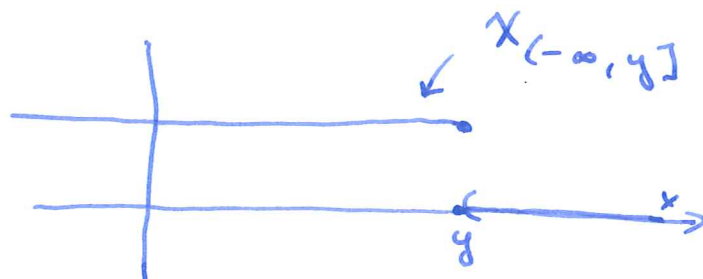
We will prove the central limit theorem under the additional assumption that $\gamma = E(|X_1 - E(X_1)|^3) < \infty$.

Comment: If $E(|X_1|^3) < \infty$ then $\gamma < \infty$.

Proof: Let $\varepsilon > 0$. The distribution function $\bar{\Phi}$ of $Z \sim N(0,1)$ is continuous so for a fixed $\beta \in \mathbb{R}$ there is a $\delta > 0$ such that $|\bar{\Phi}(x) - \bar{\Phi}(\beta)| < \varepsilon$ for $|x - \beta| < \delta$.

Denote by $X_{(-\infty, y]}$ the indicator function of the interval $(-\infty, y]$.

Figure:



Analysis 1 gives: there are functions $f^{-\varepsilon}$ and f^{ε} with values on $[0, 1]$ such that:

(i) $f^{-\varepsilon}, f^{\varepsilon}$ are three times continuously differentiable.

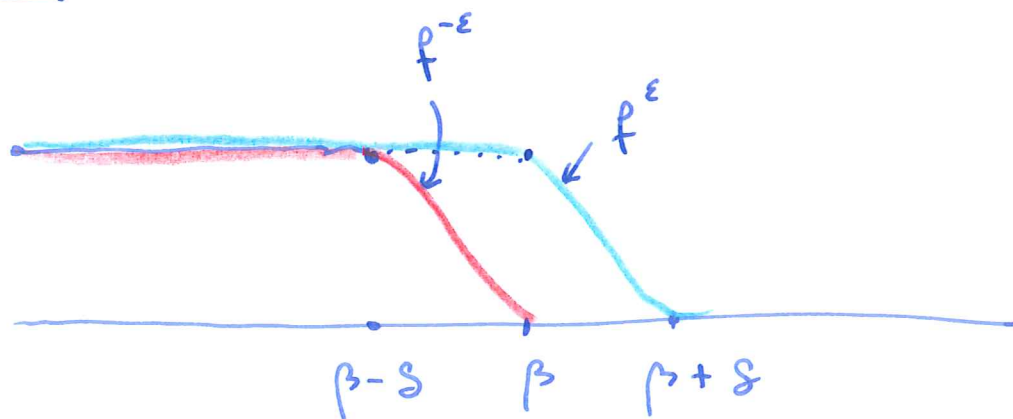
(ii) derivatives up to the third are bounded by $M < \infty$

(iii)

$$\chi_{(-\infty, \beta - \delta]} \leq f^{-\varepsilon} \leq \chi_{(-\infty, \beta]}$$

$$\leq f^{\varepsilon} \leq \chi_{(-\infty, \beta + \delta]}$$

Figure:



We have

$$\begin{aligned} E(\chi_{(-\infty, y]}(\tilde{S}_n)) \\ = P(\tilde{S}_n \leq y) \end{aligned}$$

and similarly

$$E(\chi_{(-\infty, y]}(z)) = P(z \leq y) = \Phi(y).$$

Inequalities in (iii) imply

$$\begin{aligned} \Phi(\beta - \delta) &\leq E(\chi^{-\varepsilon}(z)) \\ &\leq \Phi(\beta) \leq E(\chi^{\varepsilon}(z)) \\ &\leq \Phi(\beta + \delta) \end{aligned}$$

$$\text{But } E(\chi^{-\varepsilon}(\tilde{S}_n)) \rightarrow E(\chi^{-\varepsilon}(z))$$

$$E(\chi^{\varepsilon}(\tilde{S}_n)) \rightarrow E(\chi^{\varepsilon}(z))$$

as $n \rightarrow \infty$.

For sufficiently large n
we will have

$$\Phi(\beta - \delta) - \varepsilon \leq \mathbb{P}(\tilde{S}_n \leq \beta) \leq \Phi(\beta + \delta) + \varepsilon.$$

or

$$\Phi(\beta) - 2\varepsilon \leq \mathbb{P}(\tilde{S}_n \leq \beta) \leq \Phi(\beta) + 2\varepsilon.$$

⊖ This proves that

$$\lim_{n \rightarrow \infty} \mathbb{P}(\tilde{S}_n \leq \beta) = \Phi(\beta).$$

Examples: Typically we want to

⊖ estimate probabilities of the form

$$\mathbb{P}(a \leq S_n \leq b).$$

we compute

$$\mathbb{P}(a \leq S_n \leq b)$$

$$= \mathbb{P}(a - E(S_n) \leq S_n - E(S_n) \leq b - E(S_n))$$

$$= P \left(\underbrace{\frac{a - E(S_n)}{\sqrt{\text{var}(S_n)}}}_{\alpha} \leq \frac{S_n - E(S_n)}{\sqrt{\text{var}(S_n)}} \leq \underbrace{\frac{b - E(S_n)}{\sqrt{\text{var}(S_n)}}}_{\beta} \right)$$

$$= P \left(\alpha \leq \frac{S_n - E(S_n)}{\sqrt{\text{var}(S_n)}} \leq \beta \right)$$

$$\stackrel{\text{CLT}}{\approx} \Phi(\beta) - \Phi(\alpha)$$

⊖

(i) Let X_1, X_2, \dots be independent and $X_k \sim \text{Bernoulli}(p)$. We know that $S_n = X_1 + X_2 + \dots + X_n \sim \text{Bin}(n, p)$

so $E(S_n) = n \cdot p$ and $\text{var}(S_n) = npq$.

○ Assume $n = 10,000$ and $p = \frac{1}{2}$

and $a = 4950$ and $b = 5050$. We have

$$P(4950 \leq S_n \leq 5050)$$

$$= P \left(\frac{4950 - 5000}{50} \leq \frac{S_n - E(S_n)}{\sqrt{\text{var}(S_n)}} \leq \frac{5050 - 5000}{50} \right)$$

$$\approx P(-1 \leq Z \leq 1)$$

Statistical software gives

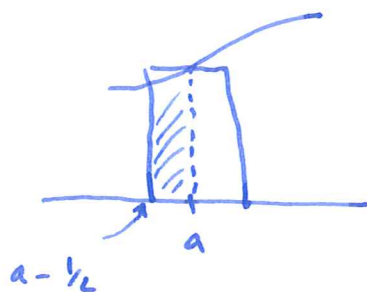
$$P(-1 \leq z \leq 1) = \Phi(1) - \Phi(-1)$$

$$= 0.6827$$

The exact probability is 0.6875.

If x_1, x_2, \dots are integer valued we can improve the approximation by changing a to $a - 1/2$ and b to $b + 1/2$.

Figure:



Changing a to $a - 1/2$ adds the "half" of the column over a . This correction is called correction for continuity.

Using this correction we get

$$P(4950 \leq S_n \leq 5500)$$

$$\approx \Phi(1.0100) - \Phi(-1.0100)$$

$$\doteq 0.6875.$$

This is accurate to 4 decimals!

(ii) Let X_1, X_2, \dots be independent

$$\text{and } P(X_i = 1) = P(X_i = 2) = P(X_i = 9) = 1/3$$

Let $n = 300$. We find

$$E(S_{300}) = 300 \cdot 4 = 1200$$

$$\text{var}(S_{300}) = 300 \cdot \frac{38}{3} = 3800$$

We approximate

$$P(1,100 \leq S_{300} \leq 1,300)$$

$$= P\left(\frac{1100 - 1200}{\sqrt{3800}} \leq \frac{S_{300} - E(S_{300})}{\sqrt{3800}} \leq \frac{1300 - 1200}{\sqrt{3800}}\right)$$

$$\approx \Phi(1.622) - \Phi(-1.622)$$

$$\doteq 0.8952$$

The exact probability using the fast Fourier transform turns out to be 0.8970. If we include the continuity correction we get 0.8970!

Can we say anything about the accuracy approximation? The answer is yes but the proof is demanding.

Theorem 6.3 (Berry-Esséen) Let $\gamma = E(|X_1 - E(X_1)|^3)$ and keep all the assumptions of Theorem 6.1.

Then

$$\sup_{x \in \mathbb{R}} \left| P\left(\frac{S_n - E(S_n)}{\sqrt{\text{var}(S_n)}} \leq x\right) - \Phi(x) \right| \leq \frac{C \cdot \gamma}{\sqrt{n} \text{var}(X_1)^{3/2}}$$

where $C < 0.4748$.

For a proof see

Shvertsova, I., On the accuracy
of the normal approximation
for sums of independent symmetric
random variables, Dokl. Akad. Nauk
443 (2012), no. 6, 671-676.

⊙ THE END ⊙