

STAT2001/STAT3101:

Probability and Inference

Lecturer: Yvo Pokern,
Room 144 Dept. of Statistical Science
Lectures:

- **Medical Sciences 131 A V Hill LT on Thursdays
1pm-2pm**
- **Medical Sciences 131 A V Hill LT on Fridays 2pm-4pm**

Office Hour: Fridays 4:15pm-5:15pm
Email: `y.pokern@ucl.ac.uk`

Lecturecast

Lecturecast recording is active capturing both video of the front of the lecture theatre and audio using the lecturer's lapel microphone.

Lecture Topics

- 1 **Joint (or multivariate) distributions: describe the joint behaviour of more than one random variable.**
- 2 **Transformation of variables: e.g. taking logarithms to get closer to normality, the Body Mass Index**

Lecture Topics

- 1 Joint (or multivariate) distributions: describe the joint behaviour of more than one random variable.**
- 2 Transformation of variables: e.g. taking logarithms to get closer to normality, the Body Mass Index**
- 3 Generating functions: For simplifying calculations for probability distributions, e.g. for sums of independent random variables leading to the Central Limit Theorem**
- 4 Distributions of functions of normally distributed variables: chi-squared, t, F etc.**
- 5 Statistical estimation: how to derive estimators and obtain their properties**

How to use the lecture notes

The lecture notes contain all relevant material for the course, *i.e.* definitions and results as well as the methods required to derive these results.

However, we will work through the examples in detail during the lecture and additional explanations not contained in this booklet will also be given in the lectures. It is therefore essential to attend the lectures and supplement the lecture notes with your own notes. The lecture notes would be woefully incomplete without the weekly exercise sheets that will be handed out separately and discussed in tutorials.

Exercise Sheets

Your solutions to the B-section of exercise sheets 1-8 should be handed in by the deadline stated on the exercise sheet. One randomly selected B-section question will be marked (the same question for everybody) and the 7 best marks out of 8 will constitute the 10% in-course assessed component of this course.

UG Student Handbook:

Plagiarism means attempting to pass off someone else's work as your own, while collusion means passing off joint work as your own unaided effort. Both are unacceptable, particularly in material submitted for examination purposes including exercises done in your own time for in-course assessment.

Tutorials

- **For your tutorial slot: check your online timetable and email, if in serious doubt, contact Karen Leport (Statistical Science, General Office, Room 120)**
- **Attendance at Tutorials is compulsory**
- **Next week's tutorial: section A of exercise sheet 1**
- **Subsequent weeks: section B of previous week's exercise sheet**

Content on Moodle at
<https://moodle.ucl.ac.uk/>

Enrolment key: Rao-Blackwell

- full answers to section A questions on the exercise sheets, very succinct answers to section B questions normally published Thursday mornings at 9:15am
- very succinct answers to section C questions (in time for exam preparation)
- messages, such as rearranged lectures
- general discussion forum moderated by the lecturer
- polls & quizzes to check your understanding and inform lecture progress
- Links to lecturecast material (recorded lectures and additional short explanatory videos for basic maths)

Learning outcomes

At the end of each chapter and of important sections you will find a list of *Learning Outcomes*. These summarize key aspects, and point out what you are expected to be able to do once you have ‘learned’ the material. You can use them to monitor your own progress and to check whether you are well prepared for in–course assessments or the exam. The learning outcomes will be reflected in the examples and exercises given throughout the course.

Revision of Basic Probability

The fundamental idea of probability is that chance can be measured on a scale which runs from zero, which represents *impossibility*, to one, which represents *certainty*.

Sample space, Ω : the set of all outcomes of an experiment (real or hypothetical).

Event, A : a subset of Ω . The elements $\omega \in \Omega$ are called elementary events or outcomes.

Event Space, \mathcal{A} : The family of events whose probability we might be interested in.

Probability measure

P: a mapping from the Event Space \mathcal{A} to $[0, 1]$ such that

- ① **$P(A) \geq 0$**
- ② **$P(\Omega) = 1$**
- ③ **Countable additivity: If A_1, A_2, \dots is a sequence of mutually disjoint sets (i.e. $A_i \cap A_j = \emptyset$, for all $i \neq j$) then**

$$P\left(\bigcup_{i=1}^{\infty} A_i\right) = P(A_1 \cup A_2 \cup \dots) = \sum_{i=1}^{\infty} P(A_i).$$

These conditions are Kolmogorov's axioms of probability.

If Ω is countable (*i.e.* $\Omega = \{\omega_1, \omega_2, \dots\}$) then for P to be a probability measure, the axioms will generally hold for *all* subsets A and mutually disjoint A_i in Ω .

If Ω is uncountable, (like any interval of real numbers) we have to define a ‘suitable’ class of subsets, \mathcal{A} , the event space, for which the axioms hold — in practice this can always be constructed to include all events of interest. For *any* two events A and B we have the addition rule

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

Independence of events

Events A and B are said to be independent if $P(A \cap B) = P(A)P(B)$.

Events A_1, A_2, \dots, A_n are independent if

$$P(A_{i_1} \cap \dots \cap A_{i_k}) = P(A_{i_1}) \dots P(A_{i_k})$$

for all possible choices of k and

$$1 \leq i_1 < i_2 < \dots < i_k \leq n.$$

That is, the product rule must hold for every subclass of the events A_1, \dots, A_n .

Example 1.1

Consider two independent tosses of a fair coin and the events $A =$ ‘first toss is head’, $B =$ ‘second toss is head’, $C =$ ‘different results on two tosses’.

Find the sample space, the probability of an elementary event and the individual probabilities of A , B , and C .

Show that A , B , and C are not independent.



Conditional Probability

Suppose that $P(B) > 0$. Then the conditional probability of A given B , $P(A|B)$ is defined as

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

i.e. the relative weight attached to event A within the restricted sample space B . The conditional probability is undefined if $P(B) = 0$. Note that $P(\cdot|B)$ is a probability measure on B .

Further note that if A and B are independent events then $P(A|B) = P(A)$ and $P(B|A) = P(B)$.

Multiplication Rule

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

Implies

$$\begin{aligned} P(A \cap B) &= P(A|B)P(B) \\ &= P(B|A)P(A) \end{aligned}$$

Note that if $P(B|A) = P(B)$ then we recover the multiplication rule for independent events.

Conditional Independence

Two events A and B are conditionally independent given a third event C if

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

Conditional independence means that once we know that C is true A carries no information on B . Note that conditional independence does not imply marginal independence, nor vice versa.

Example 1.2(1.1 ctd.)

Show that A and B are not conditionally independent given \overline{C} . □

Total Probability

The law of total probability, or partition law follows from the additivity axiom and the definition of conditional probability: suppose that B_1, \dots, B_k are mutually exclusive and exhaustive events (i.e. $B_i \cap B_j = \emptyset$ for all $i \neq j$ and $\cup_i B_i = \Omega$) and let A be any event. Then

$$P(A) = \sum_{j=1}^k P(A \cap B_j) = \sum_{j=1}^k P(A|B_j)P(B_j)$$

Example 1.3

A child gets to throw a fair die. If the die comes up 5 or 6, she gets to sample a sweet from box A which contains 10 chocolate sweets and 20 caramel sweets. If the die comes up 1,2,3 or 4 then she gets to sample a sweet from box B which contains 5 chocolate sweets and 15 caramel sweets. What is the conditional probability she will get a chocolate sweet if the die comes up 5 or 6? What is the conditional probability she will get a chocolate sweet if the die comes up 1,2,3 or 4? What is her probability of getting a chocolate sweet?

Bayes Theorem – Bayesian Statistics

Bayes theorem follows from the law of total probability and the multiplication rule.

Again, let B_1, \dots, B_k be mutually exclusive and exhaustive events and let A be any event with $P(A) > 0$. Then Bayes theorem states that

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

1.2 Revision of random variables (univariate case)

A random variable, X , assigns a real number x to each element ω of the sample space Ω . The probability measure P on Ω then gives rise to a probability distribution for X . More formally, any function $X : \Omega \rightarrow \mathbb{R}$ is called a random variable.

The random variable X may be discrete or continuous.

1.2 Revision of random variables (univariate case)

A random variable, X , assigns a real number x to each element ω of the sample space Ω . The probability measure P on Ω then gives rise to a probability distribution for X . More formally, any function $X : \Omega \rightarrow \mathbb{R}$ is called a random variable.

The random variable X may be discrete or continuous.

The probability measure P on Ω induces a probability distribution for X . In particular, X has (cumulative) distribution function (cdf) $F_X(x) = P(\{\omega : X(\omega) \leq x\})$, which is usually abbreviated to $P(X \leq x)$.

It follows that $F_X(-\infty) = 0$, $F_X(\infty) = 1$, F_X is monotone increasing, and $P(a < X \leq b) = F_X(b) - F_X(a)$.

ANY function?

A random variable, X , assigns a real number x to each element ω of the sample space Ω . The probability measure P on Ω then gives rise to a probability distribution for X . More formally, **any** function $X : \Omega \rightarrow \mathbb{R}$ is called a random variable.

The random variable X may be discrete or continuous.

The probability measure P on Ω induces a probability distribution for X . In particular, X has (cumulative) distribution function (cdf) $F_X(x) = P(\{\omega : X(\omega) \leq x\})$, which is usually abbreviated to $P(X \leq x)$.

It follows that $F_X(-\infty) = 0$, $F_X(\infty) = 1$, F_X is monotone increasing, and $P(a < X \leq b) = F_X(b) - F_X(a)$.

ANY function? **ALL** of Gallia?

ANY function? **ALL** of Gallia? No!

The function must be measurable but it is *very hard* to come up with a non-measurable function...

Example 1.4

Give an example of a random variable whose cdf is right-continuous (it has to be) but not continuous.

Discrete random variables.

X takes only a finite or countably infinite set of values $\{x_1, x_2, \dots\}$. F_X is a step-function, with steps at the x_i of sizes $p_X(x_i) = P(X = x_i)$, and $p_X(\cdot)$ is the *probability mass function (pmf)* of X . (E.g. X = place of horse in race, grade of egg.)

Example 1.5 (1.1 ctd. II)

Consider, for example the random variable X = number of heads obtained on the two tosses. Obtain the pmf and cdf of X . □

Example 1.6

Consider the random variable $X \sim \text{Geo}(p)$ with $P(X = k) = (1 - p)^{k-1}p$ where $k \in \mathbb{N}$. Compute the cdf and sketch it. Is X a discrete or a continuous random variable?

Continuous random variables.

If F_X can be expressed as

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

with f_X non-negative and integrating to one, then X is a continuous random variable.

In this case, X takes values in a non-countable set and

$$P(x < X \leq x + dx) \simeq f_X(x) dx.$$

Thus $f_X(x) dx$ is the probability that X lies in the infinitesimal interval $(x, x + dx)$. Note that the probability that X is *exactly* equal to x is zero for all x (i.e. $P(X = x) = 0$).

Example 1.7

Suppose $f_X(x) = k(2 - x^2)$ on $(-1, 1)$. Calculate k and sketch the pdf. Calculate and sketch the cdf. Is the cdf differentiable? Calculate $P(|X| > 1/2)$.

Example 1.8

The lecturer L recruits a student volunteer V and obtains his/her height in feet and inches, converts to cm denoting the result as $h \in \mathbb{R}$.

$$\begin{aligned} 1\text{in} &= 2.54\text{cm} \\ 1\text{ft} &= 12\text{in} \\ &= 30.12\text{cm} \end{aligned}$$

Example 1.8

The lecturer L recruits a student volunteer V and obtains his/her height in feet and inches, converts to cm denoting the result as $h \in \mathbb{R}$. Then ask for the maximal wager V would be willing to make on each of the following bets involving V's true height H :

1in=2.54cm
1ft=12in
=30.12cm

- ① L pays V £1 if $H \in (h - 15\text{cm}, h + 15\text{cm})$
- ② L pays V £1 if $H \in (h - 5\text{cm}, h + 5\text{cm})$
- ③ L pays V £1 if $H \in (h - 1\text{cm}, h + 1\text{cm})$
- ④ L pays V £1 if $H \in (h - 0.5\text{cm}, h + 0.5\text{cm})$
- ⑤ L pays V £1 if $H \in (h - 0.1\text{cm}, h + 0.1\text{cm})$

Example 1.8

The lecturer L recruits a student volunteer V and obtains his/her height in feet and inches, converts to cm denoting the result as $h \in \mathbb{R}$. Then ask for the maximal wager V would be willing to make on each of the following bets involving V's true height H :

1in=2.54cm
1ft=12in
=30.12cm

- ① L pays V £1 if $H \in (h - 15\text{cm}, h + 15\text{cm})$
- ② L pays V £1 if $H \in (h - 5\text{cm}, h + 5\text{cm})$
- ③ L pays V £1 if $H \in (h - 1\text{cm}, h + 1\text{cm})$
- ④ L pays V £1 if $H \in (h - 0.5\text{cm}, h + 0.5\text{cm})$
- ⑤ L pays V £1 if $H \in (h - 0.1\text{cm}, h + 0.1\text{cm})$

Assuming that fair bets have been offered, calculate the volunteer's subjective probability for each of the intervals.

Expectation.

A distribution has several characteristics that could be of interest such as its shape or skewness. Another one is its expectation, which can be regarded as a summary of the 'average' value of a random variable.

Discrete case:

$$\mathbb{E}(X) = \sum_i x_i p_x(x_i) = \sum_{\omega} X(\omega) P(\{\omega\}).$$

That is, the averaging can be taken over the (distinct) values of X with weights given by the probability distribution p_x , or over the sample space Ω with weights $P(\{\omega\})$.

Continuous case:

$$\mathbb{E}(X) = \int_{-\infty}^{\infty} x f_x(x) dx.$$

Example 1.9

The discrete random variable X has pmf $p_X(k) = \frac{1}{e^\mu - 1} \frac{\mu^k}{k!}$ for $k \in \mathbb{N}$. Compute its expectation.

Functions of a random variable

Let $\phi : \mathbb{R} \rightarrow \mathbb{R}$.

Then the random variable $Y = \phi(X)$ is defined by

$$Y(\omega) \equiv \phi(X)(\omega) = \phi(X(\omega))$$

Since $X : \Omega \rightarrow \mathbb{R}$, it follows that $\phi(X) : \Omega \rightarrow \mathbb{R}$. Thus $Y = \phi(X)$ is also a random variable. Its expectation is

$$\begin{aligned}\mathbb{E}\{\phi(X)\} &= \sum_i \phi(x_i) p_x(x_i). \\ &= \sum_{\omega} \phi(X(\omega)) P(\{\omega\})\end{aligned}$$

The first expression on the right-hand side averages the values of $\phi(x)$ over the distribution of X , whereas the second expression averages the values of $\phi(X(\omega))$ over the probabilities of $\omega \in \Omega$. A third method would be to compute the distribution of Y and average the values of y over the distribution of Y .

Example 1.10 (1.1 ctd III)

Let X be the random variable indicating the number of heads on two tosses. Consider the transformation ϕ with $\phi(0) = \phi(2) = 0$ and $\phi(1) = 1$. Find $\mathbb{E}(X)$ and $\mathbb{E}\{\phi(X)\}$.



Variance

The variance of X is

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

Equivalently $\sigma^2 = \mathbb{E}(X^2) - \{\mathbb{E}(X)\}^2$ (*exercise: prove*). The square root, σ , of σ^2 is called the **standard deviation**.

Variance

The variance of X is

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

Equivalently $\sigma^2 = \mathbb{E}(X^2) - \{\mathbb{E}(X)\}^2$ (*exercise: prove*). The square root, σ , of σ^2 is called the standard deviation.

Example 1.11 (1.1 ctd. IV)

Find $\text{Var}(X)$ and $\text{Var}\{\phi(X)\}$.



Linear functions of X .

The following properties of expectation and variance are easily proved (*exercise/previous notes*):

$$\mathbb{E}(a + bX) = a + b\mathbb{E}(X), \quad \text{Var}(a + bX) = b^2\text{Var}(X)$$

Linear functions of X .

The following properties of expectation and variance are easily proved (*exercise/previous notes*):

$$\mathbb{E}(a + bX) = a + b\mathbb{E}(X), \quad \text{Var}(a + bX) = b^2\text{Var}(X)$$

Example 1.12 (1.1 ctd. V)

Let Y be the excess of heads over tails obtained on the two tosses of the coin. Write down $\mathbb{E}(Y)$ and $\text{Var}(Y)$. □

Linear functions of X .

The following properties of expectation and variance are easily proved (*exercise/previous notes*):

$$\mathbb{E}(a + bX) = a + b\mathbb{E}(X), \quad \text{Var}(a + bX) = b^2\text{Var}(X)$$

Example 1.12 (1.1 ctd. V)

Let Y be the excess of heads over tails obtained on the two tosses of the coin. Write down $\mathbb{E}(Y)$ and $\text{Var}(Y)$. □

Standard distributions: For ease of reference, Appendices 1 and 2 provide definitions of standard discrete and continuous distributions given in earlier courses.

1.3 Joint distributions

Let us first consider the bivariate case. Suppose that the two random variables X and Y share the same sample space Ω (e.g. the height and the weight of an individual). Then we can consider the event

$$\{\omega : X(\omega) \leq x, Y(\omega) \leq y\}$$

and define its probability, regarded as a function of the two variables x and y , to be the joint (cumulative) distribution function of X and Y , denoted by

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

1.3.1 Joint CDF

The joint cumulative distribution function (cdf) has similar *properties* to the univariate cdf. The joint cdf $F_{X,Y}(x, y)$ has the properties

- 1 $F_{X,Y}(-\infty, y) = F_{X,Y}(x, -\infty) = 0$ and $F_{X,Y}(\infty, \infty) = 1$ and
- 2 $F_{X,Y}$ is a nondecreasing function of each of its arguments
- 3 $F_{X,Y}$ is *right-continuous*. That is,
 $F_{X,Y}(x + h, y + k) \rightarrow F_{X,Y}(x, y)$ as $h, k \downarrow 0$ for all x, y .

Note that $F_{X,Y}$ is not necessarily left-continuous as it will have ‘jumps’ if X and/or Y is discrete.

1.3.1 Joint CDF

The joint cumulative distribution function (cdf) has similar *properties* to the univariate cdf. The joint cdf $F_{X,Y}(x, y)$ has the properties

- 1 $F_{X,Y}(-\infty, y) = F_{X,Y}(x, -\infty) = 0$ and $F_{X,Y}(\infty, \infty) = 1$ and
- 2 $F_{X,Y}$ is a nondecreasing function of each of its arguments
- 3 $F_{X,Y}$ is *right-continuous*. That is,
 $F_{X,Y}(x + h, y + k) \rightarrow F_{X,Y}(x, y)$ as $h, k \downarrow 0$ for all x, y .

Note that $F_{X,Y}$ is not necessarily left-continuous as it will have ‘jumps’ if X and/or Y is discrete.

The marginal cdfs of X and Y can be found from

$$F_X(x) = P(X \leq x, Y < \infty) = F_{X,Y}(x, \infty)$$

and

$$F_Y(y) = P(X < \infty, Y \leq y) = F_{X,Y}(\infty, y)$$

We already know in the univariate case that $P(x_1 < X \leq x_2) = F_X(x_2) - F_X(x_1)$. Similarly, we find in the bivariate case that

$$P(x_1 < X \leq x_2, y_1 < Y \leq y_2) = F_{X,Y}(x_2, y_2) - F_{X,Y}(x_1, y_2) - F_{X,Y}(x_2, y_1) + F_{X,Y}(x_1, y_1).$$

We already know in the univariate case that $P(x_1 < X \leq x_2) = F_X(x_2) - F_X(x_1)$. Similarly, we find in the bivariate case that

$$P(x_1 < X \leq x_2, y_1 < Y \leq y_2) = F_{X,Y}(x_2, y_2) - F_{X,Y}(x_1, y_2) - F_{X,Y}(x_2, y_1) + F_{X,Y}(x_1, y_1).$$

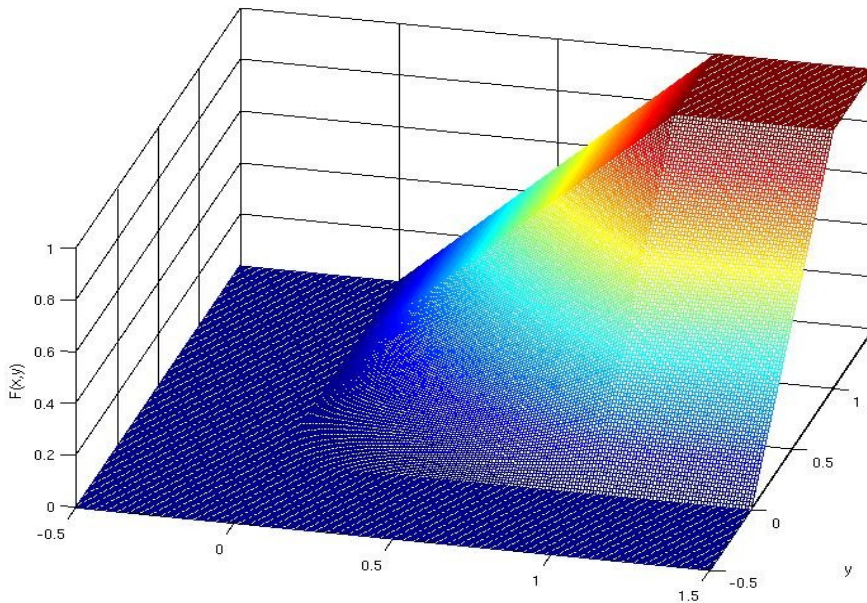
Example 1.13

Consider the function

$$F_{X,Y}(x, y) = x^2y + y^2x - x^2y^2, \quad 0 \leq x \leq 1, 0 \leq y \leq 1.$$

Show that $F_{X,Y}$ has the above properties of CDFs. Find the marginal cdfs of X and Y . Also find $P(0 \leq X \leq \frac{1}{2}, 0 \leq Y \leq \frac{1}{2})$. □

$$F_{x,y}(x, y) = x^2y + y^2x - x^2y^2$$



1.3.2 The discrete case

In many cases of interest, X and Y take only values in a discrete set. Then $F_{X,Y}$ is a step function in each variable separately and we consider the joint probability mass function

$$p_{X,Y}(x_i, y_j) = P(X = x_i, Y = y_j).$$

It is often convenient to represent a discrete bivariate distribution – a joint distribution of two variables – by a **2-way table**. In general, the entries in the table are the *joint probabilities* $p_{X,Y}(x, y)$, while the row and column totals give the *marginal probabilities* $p_X(x)$ and $p_Y(y)$. As always, the total probability is 1.

Example 1.14

Consider three independent tosses of a fair coin. Let $X =$ 'number of heads in first and second toss' and $Y =$ 'number of heads in second and third toss'. Give the probabilities for any combination of possible outcomes of X and Y in a two-way table and obtain the marginal pmfs of X and Y . □

In general, from the joint distribution we can use the law of total probability to obtain the marginal pmf of Y as

$$\begin{aligned} p_Y(y_j) = \mathbf{P}(Y = y_j) &= \sum_{x_i} \mathbf{P}(X = x_i, Y = y_j) \\ &= \sum_{x_i} p_{X,Y}(x_i, y_j). \end{aligned}$$

Similarly, the marginal pmf of X is given by

$$p_X(x_i) = \sum_{y_j} p_{X,Y}(x_i, y_j).$$

The marginal distribution is thus the distribution of just one of the variables.

The joint cdf can be written as

$$F_{X,Y}(x, y) = \sum_{x_i \leq x} \sum_{y_j \leq y} p_{X,Y}(x_i, y_j).$$

Note that there will be jumps in $F_{X,Y}$ at each of the x_i and y_j values.

Random Question to be Marked

A set of English scrabble letters are sorted into three bags: all letters with 1 point are sorted into bag (i), all letters with 0,2 or 3 points are sorted into bag (ii) and all other letters go to bag (iii). A die is thrown: if it comes up 1, bag (i) is sampled from, if it comes up 2 or 3, bag (ii) is sampled from, otherwise bag (iii) is sampled from. Exactly one letter is sampled from the selected bag.

The three events

A_1 One of the letters
B,C,M,N,Q,R,T,Z is sampled.

A_2 One of the letters
A,D,E,F,H,I,G,V,W

A_3 Neither A_1 nor A_2 happens.

are assigned to the exercise questions B1,B2,B3 by audience vote.

What assignment would you like?

Independence

The random variables X and Y , defined on the sample space Ω with probability measure P , are independent if the events

$$\{X = x_i\} \text{ and } \{Y = y_j\}$$

are *independent events*, for all possible values x_i and y_j .
Thus X and Y are independent if,

$$p_{X,Y}(x_i, y_j) = P(X = x_i, Y = y_j) = p_X(x_i)p_Y(y_j)$$

for all x_i, y_j .

This implies that $P(X \in A, Y \in B) = P(X \in A)P(Y \in B)$ for all sets A and B , so that the two events $\{\omega : X(\omega) \in A\}, \{\omega : Y(\omega) \in B\}$ are independent.
(Exercise: prove this.)

NB: If x is such that $p_x(x) = 0$, then $p_{x,y}(x, y_j) = 0$ for all y_j and the above factorisation holds automatically. Thus it does not matter whether we require the factorisation for all *possible* x_i, y_j *i.e.* those with positive probability, or all *real* x, y . (That is, $p_{x,y}(x, y) = p_x(x)p_y(y)$ for all x, y would be an equivalent definition of independence.)

If X, Y are independent then the entries in the two way table are the products of the marginal probabilities. In Example 1.14 we see that X and Y are *not* independent.

Conditional probability distributions

These are defined for random variables by analogy with conditional probabilities of events.

$$\begin{aligned} P(X = x_i \mid Y = y_j) &= \frac{P(X = x_i, Y = y_j)}{P(Y = y_j)} \\ &= \frac{p_{X,Y}(x_i, y_j)}{p_Y(y_j)} \end{aligned}$$

as a function of x_i , for fixed y_j . Then this is a pmf – it is non-negative and

$$\sum_{x_i} \frac{p_{X,Y}(x_i, y_j)}{p_Y(y_j)} = \frac{1}{p_Y(y_j)} \underbrace{\sum_{x_i} p_{X,Y}(x_i, y_j)}_{p_Y(y_j)} = 1,$$

and it gives the probabilities for observing $X = x_i$ given that we already know $Y = y_j$.

We therefore *define* the conditional probability distribution of X given $Y = y_j$ as

$$p_{X|Y}(x_i|y_j) = \frac{p_{X,Y}(x_i, y_j)}{p_Y(y_j)}$$

Conditioning on $Y = y_j$ can be compared to selecting a subset of the population, i.e. only those individuals where $Y = y_j$. The conditional distribution $p_{X|Y}$ of X given $Y = y_j$ then describes the distribution of X within this subgroup.

From the above definition we immediately obtain the multiplication rule for pmfs:

$$p_{X,Y}(x_i, y_j) = p_{X|Y}(x_i|y_j)p_Y(y_j)$$

which can be used to find a bivariate pmf when we known one marginal distribution and one conditional distribution.

Note that if X and Y are *independent* then

$p_{X,Y}(x_i, y_j) = p_X(x_i)p_Y(y_j)$ so that $p_{X|Y}(x_i|y_j) = p_X(x_i)$ i.e. the *conditional* distribution is the same as the *marginal* distribution.

In general, X and Y are independent *if and only if* the conditional distribution of X given $Y = y_j$ is the same as the marginal distribution of X for all y_j . (This condition is equivalent to $p_{X,Y}(x_i, y_j) = p_X(x_i)p_Y(y_j)$ for all x_i, y_j , above). The conditional distribution of Y given $X = x_i$ is defined similarly.

Example 1.15 (1.14 ctd.)

Obtain the conditional pmf of X given $Y = y$. Use this conditional distribution to verify that X and Y are not independent. □

Example 1.16

Suppose that R and N have a joint distribution in which $R|N$ is $\text{Bin}(N, \pi)$ and N is $\text{Poi}(\lambda)$. Show that R is $\text{Poi}(\pi\lambda)$. □

Conditional expectation

Since $p_{X|Y}(x_i|y_j)$ is a probability distribution, it has a mean or expected value:

$$\mathbb{E}(X \mid Y = y_j) = \sum_{x_i} x_i p_{X|Y}(x_i|y_j)$$

which represents the average value of X among outcomes ω for which $Y(\omega) = y_j$. This may also be written $\mathbb{E}_{X|Y}(X \mid Y = y_j)$. We can also regard the conditional expectation $\mathbb{E}(X \mid Y = y_j)$ as the mean value of X in the subgroup characterised by $Y = y_j$.

Example 1.17 (1.14 ctd. II)

Find the conditional expectations $\mathbb{E}(X \mid Y = y)$ for $y = 0, 1, 2$. Plot the graph of the function $\phi(y) = \mathbb{E}[X \mid Y = y]$. What do these values tell us about the relationship between X and Y ? □

In general, what is the relationship between the *unconditional* expectation $\mathbb{E}(X)$ and the *conditional* expectation $\mathbb{E}(X \mid Y = y_j)$?

Example 1.18

Collect the joint distribution of X : gender ($x_1 = M, x_2 = F$) and Y : number of cups of tea drunk today ($y_1 = 0, y_2 = 1, y_3 = 2, y_4 = 3$ or more). Are X and Y independent? What is the expectation of Y ? What is the conditional expectation $Y \mid X = M$ and $Y \mid X = F$? □

Consider the conditional expectation

$\phi(y) = \mathbb{E}_{X|Y}(X|Y = y)$ as a function of y . We may compute its expectation $\mathbb{E}\{\phi(Y)\} = \mathbb{E}_Y\{\mathbb{E}_{X|Y}(X|Y)\}$. But, from the definition of the expectation of a function of Y , we have $\mathbb{E}\{\phi(Y)\} = \sum_{y_j} \phi(y_j)p_Y(y_j)$, so that

$$\mathbb{E}_Y\{\mathbb{E}_{X|Y}(X|Y)\} = \sum_{y_j} \underbrace{\mathbb{E}(X|Y = y_j)}_{\text{function of } y_j} p_Y(y_j)$$

Consider the conditional expectation

$\phi(y) = \mathbb{E}_{X|Y}(X|Y = y)$ as a function of y . We may compute its expectation $\mathbb{E}\{\phi(Y)\} = \mathbb{E}_Y\{\mathbb{E}_{X|Y}(X|Y)\}$. But, from the definition of the expectation of a function of Y , we have $\mathbb{E}\{\phi(Y)\} = \sum_{y_j} \phi(y_j) p_Y(y_j)$, so that

$$\mathbb{E}_Y\{\mathbb{E}_{X|Y}(X|Y)\} = \sum_{y_j} \underbrace{\mathbb{E}(X|Y = y_j)}_{\text{function of } y_j} p_Y(y_j)$$

We show in the lectures that this gives the marginal expectation of $\mathbb{E}(X)$. That is,

$$\mathbb{E}(X) = \mathbb{E}_Y\{\mathbb{E}_{X|Y}(X|Y)\}$$

which is known as the iterated conditional expectation.

The iterated conditional expectation formula is most useful when the conditional distribution of X given $Y = y$ is known and easier to handle than the joint distribution (requiring integration to find the marginal of X if it is not known).

Example 1.19 (1.14 ctd. II)

Verify that $\mathbb{E}(Y) = \mathbb{E}_X\{\mathbb{E}_{Y|X}(Y|X)\}$ in this example.



Example 1.20 (1.16 ctd.)

Find the mean of R using the iterated conditional expectation formula.



Expectation of Functions of two variables

The definition of expectation generalises immediately to functions of two variables, *i.e.* we can compute it from the joint pmf:

$$\begin{aligned}\mathbb{E} \{ \phi(X, Y) \} &= \sum_{\omega} \phi(X(\omega), Y(\omega)) \mathbf{P}(\{\omega\}) \\ &= \sum_{x_i} \sum_{y_j} \phi(x_i, y_j) \mathbf{P}(\{\omega : X(\omega) = x_i, Y(\omega) = y_j\}) \\ &= \sum_{x_i} \sum_{y_j} \phi(x_i, y_j) p_{X,Y}(x_i, y_j)\end{aligned}$$

Iterated Conditional Expectation

Formula:

$$\begin{aligned}\mathbb{E} \{ \phi(\mathbf{X}, \mathbf{Y}) \} &= \sum_{\mathbf{x}_i} \sum_{\mathbf{y}_j} \phi(\mathbf{x}_i, \mathbf{y}_j) \mathbf{p}_{\mathbf{X}|\mathbf{Y}}(\mathbf{x}_i | \mathbf{y}_j) \mathbf{p}_{\mathbf{Y}}(\mathbf{y}_j) \\ &= \sum_{\mathbf{y}_j} \mathbf{p}_{\mathbf{Y}}(\mathbf{y}_j) \underbrace{\sum_{\mathbf{x}_i} \phi(\mathbf{x}_i, \mathbf{y}_j) \mathbf{p}_{\mathbf{X}|\mathbf{Y}}(\mathbf{x}_i | \mathbf{y}_j)}_{\mathbb{E}_{\mathbf{X}|\mathbf{Y}}(\phi(\mathbf{X}, \mathbf{y}_j) | \mathbf{Y}=\mathbf{y}_j)} \\ &= \mathbb{E}_{\mathbf{Y}} \{ \mathbb{E}_{\mathbf{X}|\mathbf{Y}}(\phi(\mathbf{X}, \mathbf{Y}) | \mathbf{Y}) \} .\end{aligned}$$

Taking out what is known (TOK)

$$\mathbb{E}_{X|Y}[\phi(Y)\psi(X, Y)|Y] = \phi(Y)\mathbb{E}_{X|Y}[\psi(X, Y)|Y]$$

This will be shown in lectures for discrete random variables only. It also holds for continuous and mixed random variables, however.

Example 1.21

Consider two random variables X and Y , where the marginal probabilities of Y are $P(Y = 0) = 3/4$, $P(Y = 1) = 1/4$ and the conditional probabilities of X are $P(X = 1|Y = 0) = P(X = 2|Y = 0) = 1/2$ and $P(X = 0|Y = 1) = P(X = 1|Y = 1) = P(X = 2|Y = 1) = 1/3$. Use the iterated conditional expectation formula to find $\mathbb{E}(XY)$.



1.3.3 The continuous case

Now, both X and Y take values in a continuous range and their joint cdf $F_{X,Y}(x, y)$ is differentiable with respect to both x and y . Then $F_{X,Y}$ can be expressed as

$$F_{X,Y}(x, y) = \int_{-\infty}^x \int_{-\infty}^y f_{X,Y}(u, v) dv du$$

where $f_{X,Y}(x, y)$ is the joint probability density function of X and Y . Letting $y \rightarrow \infty$ we get

$$F_X(x) = F_{X,Y}(x, \infty) = \int_{-\infty}^x \left(\int_{-\infty}^{\infty} f_{X,Y}(u, v) dv \right) du$$

But from §1.2 we also know that $F_X(x) = \int_{-\infty}^x f_X(u) du$. It follows that the marginal density function of X is

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

Similarly, Y has marginal density

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

As for the univariate case, we have

$$\begin{aligned} P(x < X \leq x + dx, y < Y \leq y + dy) \\ = \int_x^{x+dx} \int_y^{y+dy} f_{X,Y}(u, v) dv du \simeq f_{X,Y}(x, y) dx dy. \end{aligned}$$

That is, $f_{X,Y}(x, y) dx dy$ is the probability that (X, Y) lies in the infinitesimal rectangle $(x, x + dx) \times (y, y + dy)$. As in the univariate case, $P(X = x, Y = y) = 0$ for all x, y .

Example 1.22

Consider two continuous random variables X and Y with joint density

$$f_{X,Y}(x, y) = \begin{cases} 8xy & 0 \leq x \leq y \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

Sketch the area where $f_{X,Y}$ is positive. Derive the marginal pdfs of X and Y .

Independence

By analogy with the discrete case, two random variables X and Y are said to be independent if their joint density factorises, ie if

$$f_{X,Y}(x, y) = f_X(x)f_Y(y) \text{ for all } x, y.$$

Slogan:

Independence means Factorising

**IF X and Y are independent
THEN $f_{X,Y}(x, y) = f_X(x)f_Y(y)$ for all x, y .**

Equivalent characterisation of independence

Two continuous random variables are independent if and only if there exist functions $g(\cdot)$ and $h(\cdot)$ such that for all (x, y) the joint density factorises as $f_{X,Y}(x, y) = g(x)h(y)$, where g is a function of x only and h is a function of y only.

Proof. If X and Y are independent then simply take $g(x) = f_X(x)$ and $h(y) = f_Y(y)$. For the converse, suppose that $f_{X,Y}(x, y) = g(x)h(y)$ and define

$$G = \int_{-\infty}^{\infty} g(x) dx, \quad H = \int_{-\infty}^{\infty} h(y) dy.$$

Note that both G and H are finite (why?).

Then the marginal densities are $f_X(x) = g(x)H$, $f_Y(y) = Gh(y)$ and either of these equations implies that $GH = 1$. It follows that

$$f_{X,Y}(x, y) = g(x)h(y) = \frac{f_X(x)}{H} \frac{f_Y(y)}{G} = f_X(x)f_Y(y)$$

and so X and Y are independent. □

The advantage of knowing that under independence $f_{X,Y}(x, y) = g(x)h(y)$ is that we don't need to find the marginal densities $f_X(x)$ and $f_Y(y)$ (which would typically involve some integration) to verify independence.

Example 1.23 (1.22 ctd.)

Reminder: Example 1.22

Consider two continuous random variables X and Y with joint density

$$f_{X,Y}(x, y) = \begin{cases} 8xy & 0 \leq x \leq y \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

Sketch the area where $f_{X,Y}$ is positive. Derive the marginal pdfs of X and Y .

Example 1.23

Are X and Y independent?

Conditional distributions

For the conditional distribution, we cannot condition on $Y = y$ in the usual way, as for any arbitrary set A , $P(X \in A \text{ and } Y = y) = P(Y = y) = 0$ when Y is continuous, so that

$$P(X \in A \mid Y = y) = \frac{P(X \in A, Y = y)}{P(Y = y)}$$

is not defined ($0/0$). However, we can consider

$$\frac{P(x < X \leq x + dx, y < Y \leq y + dy)}{P(y < Y \leq y + dy)} \approx \frac{f_{X,Y}(x, y) dx dy}{f_Y(y) dy}$$

and interpret $f_{X,Y}(x, y)/f_Y(y)$ as the conditional density of X given $Y = y$ written as $f_{X|Y}(x \mid y)$.

Note that this *is* a probability density function – it is non-negative and

$$\int_{-\infty}^{\infty} \frac{f_{X,Y}(x, y)}{f_Y(y)} dx = \frac{1}{f_Y(y)} \underbrace{\int_{-\infty}^{\infty} f_{X,Y}(x, y) dx}_{f_Y(y)} = 1.$$

If X and Y are independent then, as before, the conditional density of X given $Y = y$ is just the marginal density of X .

Example 1.24 (1.22 ctd. II)

Give the conditional densities of X given $Y = y$ and of Y given $X = x$ indicating clearly the area where they are positive. Also, find $\mathbb{E}(X|Y = y)$ and $\mathbb{E}(X)$, using the law of iterated conditional expectation for the latter. Compare this with the direct calculation of $\mathbb{E}(X)$. □

Experiment: Question to be marked

Spaghetti Breaking

Measure the length L of a randomly chosen Spaghetti. Break it in two pieces at a random point. Measure the length L_1 of one of the pieces. Convert this length to the question to be marked using the following table:

Experiment: Question to be marked

Spaghetti Breaking

Measure the length L of a randomly chosen Spaghetti. Break it in two pieces at a random point. Measure the length L_1 of one of the pieces. Convert this length to the question to be marked using the following table:

$L_1 \leq 10\text{cm}$	B2
$10\text{cm} \leq L_1 \leq 20\text{cm}$	B3
$20\text{cm} < L_1$	B1

Experiment: Question to be marked

Spaghetti Breaking

Measure the length L of a randomly chosen Spaghetti. Break it in two pieces at a random point. Measure the length L_1 of one of the pieces. Convert this length to the question to be marked using the following table:

$L_1 \leq 10\text{cm}$	B2
$10\text{cm} \leq L_1 \leq 20\text{cm}$	B3
$20\text{cm} < L_1$	B1

Want to know what the odds were? Do exercise sheet 3, question B3!

1.4 Further results on expectations

Expectation of a sum.

Consider the sum $\phi(X) + \psi(Y)$ when X, Y have joint probability mass function $p_{X,Y}(x, y)$. (The continuous case follows similarly, replacing probabilities by probability densities and summations by integrals.)

1.4 Further results on expectations

Expectation of a sum.

Consider the sum $\phi(X) + \psi(Y)$ when X, Y have joint probability mass function $p_{X,Y}(x, y)$. (The continuous case follows similarly, replacing probabilities by probability densities and summations by integrals.) Then

$$\mathbb{E}_{X,Y} [\phi(X) + \psi(Y)] = \sum_{x_i} \sum_{y_j} \{\phi(x_i) + \psi(y_j)\} p_{X,Y}(x_i, y_j)$$

1.4 Further results on expectations

Expectation of a sum.

Consider the sum $\phi(X) + \psi(Y)$ when X, Y have joint probability mass function $p_{X,Y}(x, y)$. (The continuous case follows similarly, replacing probabilities by probability densities and summations by integrals.) Then

$$\begin{aligned}\mathbb{E}_{X,Y} [\phi(X) + \psi(Y)] &= \sum_{x_i} \sum_{y_j} \{\phi(x_i) + \psi(y_j)\} p_{X,Y}(x_i, y_j) \\ &= \sum_{x_i} \phi(x_i) \underbrace{\sum_{y_j} p_{X,Y}(x_i, y_j)}_{p_X(x_i)} \\ &\quad + \sum_{y_j} \psi(y_j) \underbrace{\sum_{x_i} p_{X,Y}(x_i, y_j)}_{p_Y(y_j)}\end{aligned}$$

1.4 Further results on expectations

Expectation of a sum.

Consider the sum $\phi(X) + \psi(Y)$ when X, Y have joint probability mass function $p_{X,Y}(x, y)$. (The continuous case follows similarly, replacing probabilities by probability densities and summations by integrals.) Then

$$\begin{aligned}\mathbb{E}_{X,Y} [\phi(X) + \psi(Y)] &= \sum_{x_i} \sum_{y_j} \{ \phi(x_i) + \psi(y_j) \} p_{X,Y}(x_i, y_j) \\ &= \sum_{x_i} \phi(x_i) \underbrace{\sum_{y_j} p_{X,Y}(x_i, y_j)}_{p_X(x_i)} \\ &\quad + \sum_{y_j} \psi(y_j) \underbrace{\sum_{x_i} p_{X,Y}(x_i, y_j)}_{p_Y(y_j)} \\ &= \mathbb{E}_X [\phi(X)] + \mathbb{E}_Y [\psi(Y)] .\end{aligned}$$

$$\mathbb{E}_{X,Y} [\phi(X) + \psi(Y)] = \mathbb{E}_X [\phi(X)] + \mathbb{E}_Y [\psi(Y)]$$

Note that the subscripts on the E's are unnecessary as there is no possible ambiguity in this equation, and also that this holds regardless of whether or not X and Y are independent.

In particular we have $\mathbb{E}(X + Y) = \mathbb{E}(X) + \mathbb{E}(Y)$. Note the power of this result: there is no need to calculate the probability distribution of $X + Y$ (which may be hard!) if all we need is the mean of $X + Y$.

Expectation of a product.

Now consider $\phi(X)\psi(Y)$. Then

$$\begin{aligned}\mathbb{E}_{X,Y} [\phi(X)\psi(Y)] &= \sum_{x_i} \sum_{y_j} \{\phi(x_i)\psi(y_j)\} p_{X,Y}(x_i, y_j) \\ &= ?\end{aligned}$$

Expectation of a product.

Now consider $\phi(X)\psi(Y)$. Then

$$\begin{aligned}\mathbb{E}_{X,Y} [\phi(X)\psi(Y)] &= \sum_{x_i} \sum_{y_j} \{\phi(x_i)\psi(y_j)\} p_{X,Y}(x_i, y_j) \\ &= ?\end{aligned}$$

If X and Y are independent, then $p_{X,Y}(x_i, y_j) = p_X(x_i) p_Y(y_j)$:

$$\mathbb{E}_{X,Y} [\phi(X)\psi(Y)] = \underbrace{\sum_{x_i} \phi(x_i) p_X(x_i)}_{\mathbb{E}_X[\phi(X)]} \underbrace{\sum_{y_j} \psi(y_j) p_Y(y_j)}_{\mathbb{E}_Y[\psi(Y)]}.$$

Thus, *except for the case where X and Y are independent, we typically have that*

$\mathbb{E}(\text{product}) \neq \text{product of expectations}$

But it is *always* true that

$\mathbb{E}(\text{sum}) = \text{sum of expectations.}$

Slogan:

Independence means Factorising

**IF X and Y are independent
THEN $E_{X,Y} [XY] = E_X [X] E_Y [Y]$**

Covariance

A particular function of interest is the covariance between X and Y . As we will see, this is a measure for the strength of the linear relationship between X and Y . The covariance is defined as

$$\text{Cov}(X, Y) = \mathbb{E} [(X - \mathbb{E}(X))(Y - \mathbb{E}(Y))]$$

An alternative formula for the covariance follows on expanding the bracket, giving

$$\begin{aligned}\text{Cov}(X, Y) &= \mathbb{E} [XY - X\mathbb{E}(Y) - Y\mathbb{E}(X) + \mathbb{E}(X)\mathbb{E}(Y)] \\ &= \mathbb{E}(XY) - \mathbb{E}(X)\mathbb{E}(Y) - \mathbb{E}(X)\mathbb{E}(Y) + \mathbb{E}(X)\mathbb{E}(Y) \\ &= \mathbb{E}(XY) - \mathbb{E}(X)\mathbb{E}(Y)\end{aligned}$$

Note that $\text{cov}(X, X) = \text{var}(X)$, giving the familiar formula $\text{var}(X) = \mathbb{E}(X^2) - \{\mathbb{E}(X)\}^2$.

If X and Y are *independent* then, from above,

$$\mathbb{E}(XY) = \mathbb{E}(X) \mathbb{E}(Y)$$

and it follows that

$$\text{Cov}(X, Y) = 0.$$

***However* in general $\text{cov}(X, Y) = 0 \not\Rightarrow X$ and Y are independent! (Example below)**

Also if $Z = aX + b$, then $E(Z) = aE(X) + b$ and $Z - E(Z) = a\{X - E(X)\}$, so that

$$\begin{aligned}\text{Cov}(Z, Y) &= E\{a(X - E(X))(Y - E(Y))\} \\ &= a\text{Cov}(X, Y).\end{aligned}$$

Using a similar argument we get

$$\text{Cov}(X + Y, W) = \text{Cov}(X, W) + \text{Cov}(Y, W).$$

Exercise: Using the fact that

$\text{Var}(X + Y) = \text{Cov}(X + Y, X + Y)$, derive the general formula $\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2\text{Cov}(X, Y)$.

Correlation:

From above, we see that the covariance varies with the scale of measurement of the variables (lbs/kilos etc), making it difficult to interpret its numerical value. The correlation is a standardised form of the covariance, which is *scale-invariant* and therefore its values are easier to interpret.

The correlation between X and Y is defined by

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

Suppose that $a > 0$. Then $\text{cov}(aX, Y) = a \text{Cov}(X, Y)$ and $\text{Var}(aX) = a^2 \text{Var}(X)$, so it follows that $\text{corr}(aX, Y) = \text{Corr}(X, Y)$, and thus the correlation is *scale-invariant*.

A key result is that

$$-1 \leq \text{Corr}(X, Y) \leq +1$$

for all random variables X and Y . The proof is STAT3101 only and (hopefully) made available via moodle/lecturecast.

Example 1.21 ctd.

Example 1.21

Consider two random variables X and Y , where the marginal probabilities of Y are $P(Y = 0) = 3/4$, $P(Y = 1) = 1/4$ and the conditional probabilities of X are $P(X = 1|Y = 0) = P(X = 2|Y = 0) = 1/2$ and $P(X = 0|Y = 1) = P(X = 1|Y = 1) = P(X = 2|Y = 1) = 1/3$. Use the iterated conditional expectation formula to find $\mathbb{E}(XY)$.

Example 1.24

Find the covariance and correlation of X and Y .

Experiment: Question to be marked

Spaghetti Breaking continued

Measure the length L of a Spaghetti. Break it in half at a randomly point. Measure the length L_1 of the first piece and L_2 of the second piece and compute the realisations u_1, u_2 of the random variables $U_1 = \frac{L_1}{L}$ and $U_2 = \frac{L_2}{L}$.

Question:

- **What is the distribution of U_1 ? Discrete? Continuous? Mixed?**
- **What was the probability of each question (B1,B2,B3)?**
- **What is the correlation of U_1 and U_2 ?**

Example 1.26

Compute the correlation of $X \sim U(-1, 1)$ and $Y = X^2$. Sketch a typical scatter plot of X and Y , e.g. for a sample of size 20. Are X and Y independent?

The conditional variance

Consider random variables X and Y , and the conditional probability distribution of X given $Y = y$. This conditional distribution has a mean, denoted $E(X|Y = y)$, and a variance, $\text{var}(X|Y = y)$. We have already shown that the marginal (unconditional) mean $E(X)$ is related to the conditional mean via the formula

$$E(X) = E_Y\{E_{X|Y}(X|Y)\}.$$

In the lectures we will obtain a similar result for the relation between the marginal and conditional variances. The result is that

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

Example 1.21

Consider two random variables X and Y , where the marginal probabilities of Y are $P(Y = 0) = 3/4$, $P(Y = 1) = 1/4$ and the conditional probabilities of X are $P(X = 1|Y = 0) = P(X = 2|Y = 0) = 1/2$ and $P(X = 0|Y = 1) = P(X = 1|Y = 1) = P(X = 2|Y = 1) = 1/3$. Use the iterated conditional expectation formula to find $\mathbb{E}(XY)$.

Example 1.27 (1.21 ctd.)

Find the conditional variances of X given $Y = 0, 1$.
Compute the marginal variance of X by using the above result.

Example 1.16

Suppose that R and N have a joint distribution in which $R|N$ is $\text{Bin}(N, \pi)$ and N is $\text{Poi}(\lambda)$. Show that R is $\text{Poi}(\pi\lambda)$.

Example 1.28 (1.16 ctd.)

Find the variance of R using the iterated conditional variance formula.

1.5 Standard multivariate distributions

The idea of joint probability distributions extends immediately to more variables, giving general multivariate distributions, i.e. the variables X_1, \dots, X_n have a joint cumulative distribution function

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

and may have a joint probability mass function

$$p_{X_1, \dots, X_n}(x_1, \dots, x_n) = \mathbf{P}(X_i = x_i; i = 1, \dots, n)$$

or joint probability density function

$$f_{X_1, \dots, X_n}(x_1, \dots, x_n),$$

so that a function $\phi(X_1, \dots, X_n)$ has an expectation with respect to this joint distribution etc.

Conditional distributions of a subset of variables given the rest then follow as before; for example, for discrete random variables X_1, X_2, X_3 ,

$$p_{X_1, X_2 | X_3}(x_1, x_2 | x_3) = \frac{p_{X_1, X_2, X_3}(x_1, x_2, x_3)}{p_{X_3}(x_3)}$$

is the conditional pmf of (X_1, X_2) given $X_3 = x_3$. Similarly, the discrete random variables X_1, \dots, X_n are mutually independent if and only if

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

for all x_1, x_2, \dots, x_n . Mutual independence of X_1, \dots, X_n implies mutual independence of the events $\{X_1 \in A_1\}, \dots, \{X_n \in A_n\}$ (exercise: prove).

Finally, we say that X_1 and X_2 are *conditionally independent given X_3* if

$$p_{X_1, X_2 | X_3}(x_1, x_2 | x_3) = p_{X_1 | X_3}(x_1 | x_3) p_{X_2 | X_3}(x_2 | x_3)$$

for all x_1, x_2, x_3 .

These definitions hold for continuous distributions by replacing the pmf by the pdf.

1.5.1 The multinomial distribution

The multinomial distribution is a generalization of the binomial distribution. Suppose that a sample of size n is drawn (*with replacement*) from a population whose members fall into one of $m + 1$ categories. Assume that, for each individual sampled, independently of the rest

$$P(\text{individual is of type } i) = p_i, \quad i = 1, \dots, m + 1$$

where $\sum_{i=1}^{m+1} p_i = 1$. Let N_i be the number of type i individuals in the sample. Note that, since $N_{m+1} = n - \sum_{i=1}^m N_i$, N_{m+1} is determined by N_1, \dots, N_m . We therefore only need to consider the joint distribution of the m random variables N_1, \dots, N_m .

The joint pmf of N_1, \dots, N_m is given by

$$P(N_1 = n_1, \dots, N_m = n_m) = \begin{cases} \frac{n!}{n_1! \dots n_{m+1}!} p_1^{n_1} \dots p_{m+1}^{n_{m+1}}, & n_1 + \dots + n_m \leq n \\ 0 & \text{otherwise.} \end{cases}$$

where $n_{m+1} = n - \sum_{i=1}^m n_i$.

This is the multinomial distribution with index n and parameters p_1, \dots, p_m , where $p_{m+1} = 1 - \sum_{i=1}^m p_i$ (so p_{m+1} is not a 'free' parameter).

$$P(N_1 = n_1, \dots, N_m = n_m) = \begin{cases} \frac{n!}{n_1! \dots n_{m+1}!} p_1^{n_1} \dots p_{m+1}^{n_{m+1}}, & n_1 + \dots + n_m \leq n \\ 0 & \text{otherwise.} \end{cases}$$

To justify the above joint pmf note that we want the probability that the n trials result in exactly n_1 outcomes of the first category, n_2 of the second, \dots , n_{m+1} in the last category. Any specific ordering of these n outcomes has probability $p_1^{n_1} \dots p_{m+1}^{n_{m+1}}$ by the assumption of independent trials, and there are $\frac{n!}{n_1! \dots n_{m+1}!}$ such orderings.

$$P(N_1 = n_1, \dots, N_m = n_m) = \begin{cases} \frac{n!}{n_1! \dots n_{m+1}!} p_1^{n_1} \dots p_{m+1}^{n_{m+1}}, & n_1 + \dots + n_m \leq n \\ 0 & \text{otherwise.} \end{cases}$$

If $m = 1$ the multinomial distribution is just the binomial distribution, i.e. $N_1 \sim \text{Bin}(n, p_1)$, which has mean np_1 and variance $np_1(1 - p_1)$.

$$P(N_1 = n_1, \dots, N_m = n_m) = \begin{cases} \frac{n!}{n_1! \dots n_{m+1}!} p_1^{n_1} \dots p_{m+1}^{n_{m+1}}, & n_1 + \dots + n_m \leq n \\ 0 & \text{otherwise.} \end{cases}$$

Example 1.29

Suppose that a bag contains five red, five black and five yellow balls and that three balls are drawn at random with replacement. What is the probability that there is one of each colour? □

Marginal distribution of N_i .

Clearly N_i can be regarded as the number of successes in n independent Bernoulli trials if we define *success* to be *individual is of type i* . Thus N_i has a binomial distribution, $N_i \sim \text{Bin}(n, p_i)$, with mean np_i and variance $np_i(1 - p_i)$.

Example 1.30

Let N_A , N_B and N_F be the numbers of A grades, B grades and fails respectively amongst a class of 100 students. Suppose that generally 5% of students achieve grade A, 30% grade B and that 5% fail. Write down the joint distribution of N_A , N_B and N_F and find the marginal distribution of N_A . □

Joint distribution of N_i and N_j .

Again we can regard individuals as being one of three types, i , j and $k=\{\text{not } i \text{ or } j\}$. This is the trinomial distribution with probabilities

$$P(N_i = n_i, N_j = n_j) = \begin{cases} \frac{n!}{n_i!n_j!n_k!} p_i^{n_i} p_j^{n_j} p_k^{n_k}, & n_i + n_j \leq n \\ 0 & \text{otherwise} \end{cases}$$

where $n_k = n - n_i - n_j$ and $p_k = 1 - p_i - p_j$. It is intuitively clear that N_i and N_j are dependent and negatively correlated, since a relatively large value of N_i implies a relatively small value of N_j and conversely. We show this as follows. First,

$$\mathbb{E}(N_i N_j) = \sum_{\{n_i, n_j \geq 0, n_i + n_j \leq n\}} \frac{n_i n_j}{n_i! n_j! n_k!} p_i^{n_i} p_j^{n_j} p_k^{n_k}$$

- goal: create a pmf that we know sums to one.

$$\mathbb{E}(N_i N_j)$$

$$= \sum_{\{n_i, n_j \geq 0, n_i + n_j \leq n\}} n_i n_j \frac{n!}{n_i! n_j! n_k!} p_i^{n_i} p_j^{n_j} p_k^{n_k}$$

$$= n(n-1) p_i p_j.$$

$$\sum_{\{n_i-1, n_j-1 \geq 0, n_i+n_j-2 \leq n-2\}} \frac{(n-2)!}{(n_i-1)!(n_j-1)!n_k!} p_i^{n_i-1} p_j^{n_j-1} p_k^{n_k}$$

- goal: create a pmf that we know sums to one.
- Note that we may take $n_i, n_j \geq 1$ in the sum, since if either n_i or n_j is zero then the corresponding term in the sum is zero.

$$\mathbb{E}(N_i N_j)$$

$$= \sum_{\{n_i, n_j \geq 0, n_i + n_j \leq n\}} n_i n_j \frac{n!}{n_i! n_j! n_k!} p_i^{n_i} p_j^{n_j} p_k^{n_k}$$

$$= n(n-1)p_i p_j.$$

$$\sum_{\{n_i-1, n_j-1 \geq 0, n_i+n_j-2 \leq n-2\}} \frac{(n-2)!}{(n_i-1)!(n_j-1)!n_k!} p_i^{n_i-1} p_j^{n_j-1} p_k^{n_k}$$

$$= n(n-1)p_i p_j \cdot 1$$

$$= n(n-1)p_i p_j$$

- goal: create a pmf that we know sums to one.
- Note that we may take $n_i, n_j \geq 1$ in the sum, since if either n_i or n_j is zero then the corresponding term in the sum is zero.
- summation through using known multinomial pmf

Finally

$$\begin{aligned}\text{Cov}(N_i, N_j) &= \mathbb{E}(N_i N_j) - \mathbb{E}(N_i)\mathbb{E}(N_j) \\ &= n(n-1)p_i p_j - (np_i)(np_j) = -np_i p_j\end{aligned}$$

and so

$$\begin{aligned}\text{Corr}(N_i, N_j) &= \frac{-np_i p_j}{\sqrt{np_i(1-p_i)np_j(1-p_j)}} \\ &= -\sqrt{\frac{p_i p_j}{(1-p_i)(1-p_j)}}.\end{aligned}$$

Note that $\text{corr}(N_i, N_j)$ is negative, as anticipated, and also that it does not depend on n .

Random Question to be Marked

A box of sweets contains 4 sweets with black wrapper, 3 with light blue wrapper and 4 with dark blue wrapper. Three sweets are sampled from the box (with replacement). The three events

- A Three sweets of three different colours.**
- B Three sweets of the same colour.**
- C Two sweets black, one of a different colour.**

are assigned to the exercise questions B1,B2,B3 by audience vote. If none of the three events occurs, the experiment is repeated.

What assignment would you like?

Conditional distribution of N_i given $N_j = n_j$.

Given $N_j = n_j$, there are $n - n_j$ remaining independent Bernoulli trials, each with probability of being type i given by

$$P(\text{type } i | \text{not type } j) = \frac{P(\text{type } i)}{P(\text{not type } j)} = \frac{p_i}{1 - p_j}.$$

Thus, given $N_j = n_j$, N_i has a binomial distribution with index $n - n_j$ and probability $\frac{p_i}{1 - p_j}$.

Exercise: Verify this result by using the definition of conditional probability together with the joint distribution of N_i and N_j and the marginal distribution of N_j .

Example 1.31 (1.30 ctd.)

Find the conditional distribution of N_A given $N_F = 10$ and calculate $\text{Corr}(N_A, N_F)$.

Remark:

The multinomial distribution can also be used as a model for *contingency tables*. Let X and Y be discrete variable with a number

of I and J different outcomes, respectively. Then, in a trial of size n , N_{ij} will count the number of outcomes where we observe $X = i$ and $Y = j$. The counts N_{ij} , $i = 1, \dots, I$, $j = 1, \dots, J$, are typically arranged in a contingency table, and from the above considerations we know that their joint distribution is multinomial with parameters n and $p_{ij} = P(X = i, Y = j)$, $i = 1, \dots, I$, $j = 1, \dots, J$.

This leads to the analysis of *categorical data*, for which a question of interest is often ‘are the categories independent?’, *i.e.* is $p_{ij} = p_i p_j$ for all i, j ? Exact significance tests of this hypothesis can be constructed from the multinomial distribution of the entries in the contingency table.

Announcement

Lecturecast video of solution

A video of the fully worked solution of one of the exercise questions is available for you to watch on moodle. You are expected to watch the video **before the tutorial.**

Linear algebra and quadratic equation revision videos have also been made available on lecturecast (links on Moodle) - please use these, especially if you do not know that $\forall \alpha \in \mathbb{R} \forall \mathbf{A} \in \mathbb{R}^{n \times n} : \det(\alpha \mathbf{A}) = \alpha^n \det(\mathbf{A})$.

1.5.3 The multivariate normal distribution

The continuous random variables X and Y are said to have a bivariate normal distribution if they have joint probability density function

$$f_{X,Y}(x,y) = \frac{1}{2\pi\sigma_X\sigma_Y\sqrt{1-\rho^2}} \cdot \exp \left[-\frac{1}{2(1-\rho^2)} \left\{ \left(\frac{x-\mu_X}{\sigma_X} \right)^2 - 2\rho \left(\frac{x-\mu_X}{\sigma_X} \right) \left(\frac{y-\mu_Y}{\sigma_Y} \right) + \left(\frac{y-\mu_Y}{\sigma_Y} \right)^2 \right\} \right]$$

for $-\infty < x, y < \infty$, where

$-\infty < \mu_X, \mu_Y < \infty; \sigma_X, \sigma_Y > 0; \rho^2 < 1$.

The parameters of this distribution are $\mu_X, \mu_Y, \sigma_X^2, \sigma_Y^2$, and ρ . As we will see below, these turn out to be the marginal means, variances, and the correlation of X and Y .

The bivariate normal is widely used as a model for many observed phenomena where dependence is expected, *e.g.* height and weight of an individual, length and width of a petal, income and investment returns. Sometimes the data need to be transformed (*e.g.* by taking logs) before using the bivariate normal.

Marginal distributions. In order to simplify the integrations required to find the marginal densities of X and Y , we set

$$\frac{x - \mu_X}{\sigma_X} = u, \quad \frac{y - \mu_Y}{\sigma_Y} = v.$$

Then, integrating with respect to y , the density of X can be found as

Marginal of a Bivariate Normal

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

- **Marginal Distribution is integral of joint distribution.**

Marginal of a Bivariate Normal

$$\begin{aligned}f_X(x) &= \int_{-\infty}^{\infty} f_{X,Y}(x, y) dy \\&= \int_{-\infty}^{\infty} \frac{1}{2\pi\sigma_X\sigma_Y\sqrt{1-\rho^2}} \\&\quad \cdot \exp\left[-\frac{1}{2(1-\rho^2)}\{u^2 - 2\rho uv + v^2\}\right] \sigma_Y dv\end{aligned}$$

- Marginal Distribution is integral of joint distribution.
- Insert joint bivariate normal density; express in u, v .

Marginal of a Bivariate Normal

$$f_X(x) = \int_{-\infty}^{\infty} \frac{1}{2\pi\sigma_X\sigma_Y\sqrt{1-\rho^2}} \cdot \exp\left[-\frac{1}{2(1-\rho^2)}\left\{u^2 - 2\rho u\mathbf{v} + \mathbf{v}^2\right\}\right] \sigma_Y d\mathbf{v}$$

- Marginal Distribution is integral of joint distribution.
- Insert joint bivariate normal density; express in u, v .
- Take the quadratic and linear term in \mathbf{v} ...

Marginal of a Bivariate Normal

$$\begin{aligned} f_X(x) &= \int_{-\infty}^{\infty} \frac{1}{2\pi\sigma_X\sigma_Y\sqrt{1-\rho^2}} \\ &\quad \cdot \exp\left[-\frac{1}{2(1-\rho^2)}\{u^2 - 2\rho u\mathbf{v} + \mathbf{v}^2\}\right] \sigma_Y d\mathbf{v} \\ &= \frac{1}{\sigma_X\sqrt{2\pi}} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi(1-\rho^2)}} \\ &\quad \cdot \exp\left[-\frac{1}{2(1-\rho^2)}\{(\mathbf{v} - \rho u)^2 + u^2(1-\rho^2)\}\right] d\mathbf{v} \end{aligned}$$

- Marginal Distribution is integral of joint distribution.
- Insert joint bivariate normal density; express in u, v .
- Take the quadratic and linear term in \mathbf{v} ...
- ... and complete the square.

Marginal of a Bivariate Normal

$$\begin{aligned}f_X(x) &= \frac{1}{\sigma_X \sqrt{2\pi}} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi(1-\rho^2)}} \\&\quad \cdot \exp \left[-\frac{1}{2(1-\rho^2)} \left\{ (\mathbf{v} - \rho u)^2 + \mathbf{u}^2(1-\rho^2) \right\} \right] dv \\&= \frac{1}{\sigma_X \sqrt{2\pi}} \exp \left(-\frac{1}{2} \mathbf{u}^2 \right) \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi(1-\rho^2)}} \\&\quad \cdot \exp \left[-\frac{1}{2(1-\rho^2)} \left\{ (\mathbf{v} - \rho u)^2 \right\} \right] dv\end{aligned}$$

- Marginal Distribution is integral of joint distribution.
- Insert joint bivariate normal density; express in u, v .
- Take the quadratic and linear term in \mathbf{v} ...
- ... and complete the square.
- Take term not involving \mathbf{v} outside integral.

Marginal of a Bivariate Normal

$$\begin{aligned}f_X(x) &= \frac{1}{\sigma_X \sqrt{2\pi}} \exp\left(-\frac{1}{2}u^2\right) \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi(1-\rho^2)}} \\&\quad \cdot \exp\left[-\frac{1}{2(1-\rho^2)}\left\{(\mathbf{v} - \rho u)^2\right\}\right] dv \\&= \frac{1}{\sqrt{2\pi}\sigma_X} \exp\left\{-\frac{1}{2}\left(\frac{x - \mu_X}{\sigma_X}\right)^2\right\}\end{aligned}$$

- Marginal Distribution is integral of joint distribution.
- Insert joint bivariate normal density; express in u, v .
- Take the quadratic and linear term in \mathbf{v} ...
- ... and complete the square.
- Take term not involving \mathbf{v} outside integral.
- Density $N(\rho u, 1 - \rho^2)$ integrates to 1; re-express in x .

Marginal of a Bivariate Normal

We have shown

Marginal of a Bivariate Normal

For (X, Y) bivariate normal with parameters $\mu_X, \mu_Y, \sigma_X, \sigma_Y, \rho$, the marginal distribution of X is normal, with mean μ_X and variance σ_X^2 .

- **By symmetry in X and Y we get that $Y \sim N(\mu_Y, \sigma_Y^2)$ is the marginal distribution of Y .**
- **It will be shown later that the fifth parameter, ρ also has a simple interpretation, for $\rho = \text{Corr}(X, Y)$.**

Conditional distributions

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

- Conditional density is quotient of **joint** and **marginal**.

Conditional distributions

$$\begin{aligned} f_{X|Y}(x|y) &= \frac{f_{X,Y}(x,y)}{f_Y(y)} \\ &= \frac{\sqrt{2\pi}\sigma_Y}{2\pi\sigma_X\sigma_Y\sqrt{1-\rho^2}} \\ &\quad \cdot \exp \left[-\frac{1}{2(1-\rho^2)} \left\{ \left(\frac{x-\mu_X}{\sigma_X} \right)^2 - 2\rho \left(\frac{x-\mu_X}{\sigma_X} \right) \left(\frac{y-\mu_Y}{\sigma_Y} \right) + \right. \right. \\ &\quad \left. \left. + \left(\frac{y-\mu_Y}{\sigma_Y} \right)^2 - (1-\rho^2) \left(\frac{y-\mu_Y}{\sigma_Y} \right)^2 \right\} \right] \end{aligned}$$

- Conditional density is quotient of **joint** and **marginal**.
- Insert **joint** and **marginal** densities.

Conditional distributions

$$\begin{aligned}
 f_{X|Y}(x|y) &= \frac{\sqrt{2\pi}\sigma_Y}{2\pi\sigma_X\sigma_Y\sqrt{1-\rho^2}} \\
 &\cdot \exp \left[-\frac{1}{2(1-\rho^2)} \left\{ \left(\frac{x-\mu_X}{\sigma_X} \right)^2 - 2\rho \left(\frac{x-\mu_X}{\sigma_X} \right) \left(\frac{y-\mu_Y}{\sigma_Y} \right) + \right. \right. \\
 &\quad \left. \left. + \left(\frac{y-\mu_Y}{\sigma_Y} \right)^2 - (1-\rho^2) \left(\frac{y-\mu_Y}{\sigma_Y} \right)^2 \right\} \right] \\
 &= \frac{\sqrt{2\pi}\sigma_Y}{2\pi\sigma_X\sigma_Y\sqrt{1-\rho^2}} \\
 &\cdot \exp \left[-\frac{1}{2\sigma_X^2(1-\rho^2)} \left\{ (x-\mu_X)^2 - 2\rho\frac{\sigma_X}{\sigma_Y}(x-\mu_X)(y-\mu_Y) \right. \right. \\
 &\quad \left. \left. + \rho^2\frac{\sigma_X^2}{\sigma_Y^2}(y-\mu_Y)^2 \right\} \right]
 \end{aligned}$$

- Conditional density is quotient of **joint** and **marginal**.
- Insert **joint** and **marginal** densities.
- Rewrite Expression in $[\cdot]$, factorise σ_X^{-2} .

Conditional distributions

$$\begin{aligned} f_{X|Y}(x|y) &= \frac{\sqrt{2\pi}\sigma_Y}{2\pi\sigma_X\sigma_Y\sqrt{1-\rho^2}} \\ &\cdot \exp \left[-\frac{1}{2\sigma_X^2(1-\rho^2)} \left\{ (x - \mu_X)^2 - 2\rho\frac{\sigma_X}{\sigma_Y}(x - \mu_X)(y - \mu_Y) \right. \right. \\ &\quad \left. \left. + \rho^2\frac{\sigma_X^2}{\sigma_Y^2}(y - \mu_Y)^2 \right\} \right] \\ &= \frac{\sqrt{2\pi}\sigma_Y}{2\pi\sigma_X\sigma_Y\sqrt{1-\rho^2}} \\ &\cdot \exp \left[-\frac{1}{2\sigma_X^2(1-\rho^2)} \left\{ (x - \mu_X) - \rho\frac{\sigma_X}{\sigma_Y}(y - \mu_Y) \right\}^2 \right] \end{aligned}$$

- Conditional density is quotient of **joint** and **marginal**.
- Insert **joint** and **marginal** densities.
- Rewrite Expression in $[\cdot]$, factorise σ_X^{-2} .
- **Summarise** into **one complete square**.

Conditional distributions

$$f_{X|Y}(x|y) = \frac{\sqrt{2\pi}\sigma_Y}{2\pi\sigma_X\sigma_Y\sqrt{1-\rho^2}} \cdot \exp \left[-\frac{1}{2\sigma_X^2(1-\rho^2)} \left\{ (x - \mu_X) - \rho \frac{\sigma_X}{\sigma_Y} (y - \mu_Y) \right\}^2 \right]$$

- Conditional density is quotient of **joint** and **marginal**.
- Insert **joint** and **marginal** densities.
- Rewrite Expression in $[\cdot]$, factorise σ_X^{-2} .
- **Summarise** into **one complete square**.
- So $X|Y = y \sim N(\mu_X + \rho \frac{\sigma_X}{\sigma_Y} (y - \mu_Y), \sigma_X^2(1 - \rho^2))$.

The role of ρ

Note that knowledge of $Y = y$ reduces the variability of X by a factor $(1 - \rho^2)$. The closer the correlation between X and Y , the smaller the conditional variance becomes. Note also that the conditional mean of X is a *linear function* of y . If y is relatively large then the conditional mean of X is also relatively large if $\rho = \text{Corr}(X, Y) > 0$, or is relatively small if $\rho < 0$.

The role of ρ

Note that knowledge of $Y = y$ reduces the variability of X by a factor $(1 - \rho^2)$. The closer the correlation between X and Y , the smaller the conditional variance becomes. Note also that the conditional mean of X is a *linear function* of y . If y is relatively large then the conditional mean of X is also relatively large if $\rho = \text{Corr}(X, Y) > 0$, or is relatively small if $\rho < 0$.

Suppose $\rho = 0$. Then

$$\begin{aligned} f_{X,Y}(x, y) &= \frac{1}{2\pi\sigma_X\sigma_Y} \exp \left[-\frac{1}{2} \left\{ \left(\frac{x - \mu_X}{\sigma_X} \right)^2 + \left(\frac{y - \mu_Y}{\sigma_Y} \right)^2 \right\} \right] \\ &= f_X(x)f_Y(y) \end{aligned}$$

showing that uncorrelated normal variables are independent (remember that this is not true in the general case).

Example 1.32

Let X be the one-year yield of portfolio A and Y be the one-year yield of portfolio B. From past data, the marginal distribution of X is modelled as $N(7, 1)$, whereas the marginal distribution of Y is $N(8, 4)$ (being a more risky portfolio but having a higher average yield). Furthermore, the correlation between X and Y is 0.5. Assuming that X, Y have a bivariate normal distribution, find the conditional distribution of X given that $Y = 9$ and compare this with the marginal distribution of X . Calculate the probability $P(X > 8 | Y = 9)$. □

Linear Algebra Reminder 1

Scalars:

$$\alpha, \beta, \gamma \in \mathbb{R}$$

Vectors:

$$\mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathbb{R}^n, \quad \mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = (x_1, x_2, \dots, x_n)^T$$

Matrices:

$$\mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times n}, \quad \mathbf{A} = \begin{pmatrix} a_{1,1} & a_{1,2} & \dots & a_{1,n} \\ a_{2,1} & a_{2,2} & \dots & a_{2,n} \\ \vdots & & \ddots & \vdots \\ a_{m,1} & \dots & \dots & a_{m,n} \end{pmatrix}$$

Linear Algebra Reminder 2

The usual rules of linear algebra hold:

$$\begin{aligned}(AB)C &= A(BC) & A(B + C) &= AB + AC & (A + B)C &= AC + BC \\ (\alpha + \beta)x &= \alpha x + \beta x & \alpha(x + y) &= \alpha x + \alpha y\end{aligned}$$

Multiplying in component form:

$$(Ax)_i = \sum_{j=1}^n a_{i,j}x_j \quad (AB)_{i,j} = \sum_{k=1}^n a_{i,k}b_{k,j}$$

Transposition

- A^T is the transpose of A with entries $(A^T)_{i,j} = a_{j,i}$,
 $A^T \in \mathbb{R}^{n \times m}$ if $A \in \mathbb{R}^{m \times n}$.
- Transpose to denote the inner product,

$$x^T y = \sum_{i=1}^n x_i y_i,$$

works well in matrix expressions, e.g.

$$x^T A y = (A^T x)^T y$$

since $(A^T)^T = A$.

- transposition and taking inverses:

$$(A^T)^{-1} = (A^{-1})^T =: A^{-T}$$

Symmetric Matrices

A matrix is symmetric if $A^T = A$.

Lemma: Covariance matrices are symmetric.

Proof:

Symmetric Matrices

A matrix is symmetric if $A^T = A$.

Lemma: Covariance matrices are symmetric.

Proof:

$$\begin{aligned}\Sigma_{i,j} = \text{Cov}(X_i, X_j) &= \mathbb{E}[X_i X_j] - \mathbb{E}[X_i]\mathbb{E}[X_j] = \mathbb{E}[X_j X_i] - \mathbb{E}[X_j]\mathbb{E}[X_i] \\ &= \text{Cov}(X_j, X_i) = \Sigma_{j,i}\end{aligned}$$

Determinants

$$\det : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}$$

- $\det(\mathbf{A}) = 0$ if and only if \mathbf{A} is non-invertible (i.e. singular).
- $\det(\mathbf{A}^{-1}) = \frac{1}{\det(\mathbf{A})}$ if \mathbf{A} is invertible
- $\det(\mathbf{A}) = \det(\mathbf{A}^T)$
- $\det(\mathbf{AB}) = \det(\mathbf{A})\det(\mathbf{B})$ for $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$
- The identity matrix has determinant one:

$$\det \begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & & \\ \vdots & & \ddots & \\ 0 & \dots & 0 & 1 \end{pmatrix} = 1$$

Matrix notation:

Define

$$\begin{aligned} X &= \begin{pmatrix} X \\ Y \end{pmatrix} & \mu &= \begin{pmatrix} \mu_X \\ \mu_Y \end{pmatrix} \\ \Sigma &= \begin{pmatrix} \sigma_X^2 & \sigma_{XY} \\ \sigma_{XY} & \sigma_Y^2 \end{pmatrix} & &= \begin{pmatrix} \sigma_X^2 & \rho\sigma_X\sigma_Y \\ \rho\sigma_X\sigma_Y & \sigma_Y^2 \end{pmatrix} \end{aligned}$$

Here X is a random vector, $\mu = \mathbb{E}(X)$ is its mean vector and $\Sigma = \text{Cov}(X)$ is the covariance matrix, or dispersion matrix of X .

Matrix notation:

Define

$$X = \begin{pmatrix} X \\ Y \end{pmatrix} \quad \mu = \begin{pmatrix} \mu_X \\ \mu_Y \end{pmatrix}$$
$$\Sigma = \begin{pmatrix} \sigma_X^2 & \sigma_{XY} \\ \sigma_{XY} & \sigma_Y^2 \end{pmatrix} = \begin{pmatrix} \sigma_X^2 & \rho\sigma_X\sigma_Y \\ \rho\sigma_X\sigma_Y & \sigma_Y^2 \end{pmatrix}$$

Here X is a random vector, $\mu = \mathbb{E}(X)$ is its mean vector and $\Sigma = \text{Cov}(X)$ is the covariance matrix, or dispersion matrix of X .

Then

$$|\Sigma| = \sigma_X^2 \sigma_Y^2 (1 - \rho^2), \quad \Sigma^{-1} = \frac{1}{|\Sigma|} \begin{pmatrix} \sigma_Y^2 & -\rho\sigma_X\sigma_Y \\ -\rho\sigma_X\sigma_Y & \sigma_X^2 \end{pmatrix}$$

Also, writing $x = \begin{pmatrix} x \\ y \end{pmatrix}$,

$$\begin{aligned}(x - \mu)^T \Sigma^{-1} (x - \mu) &= (x - \mu_X, y - \mu_Y) \frac{1}{|\Sigma|} \begin{pmatrix} \sigma_Y^2 & -\rho\sigma_X\sigma_Y \\ -\rho\sigma_X\sigma_Y & \sigma_X^2 \end{pmatrix} \begin{pmatrix} x - \mu_X \\ y - \mu_Y \end{pmatrix} \\&= \frac{1}{\sigma_X^2 \sigma_Y^2 (1 - \rho^2)} \{ (x - \mu_X)^2 \sigma_Y^2 - 2(x - \mu_X)(y - \mu_Y) \rho \sigma_X \sigma_Y \\&\quad + (y - \mu_Y)^2 \sigma_X^2 \} \\&= \frac{1}{1 - \rho^2} \left\{ \left(\frac{x - \mu_X}{\sigma_X} \right)^2 - 2\rho \left(\frac{x - \mu_X}{\sigma_X} \right) \left(\frac{y - \mu_Y}{\sigma_Y} \right) + \left(\frac{y - \mu_Y}{\sigma_Y} \right)^2 \right\}.\end{aligned}$$

It follows that the joint density $f_X(x)$ of X, Y given at the beginning of this section can be written as

$$f_X(x) = \frac{1}{|2\pi\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu) \right\}$$

on noting that $|2\pi\Sigma|^{1/2} = 2\pi|\Sigma|^{1/2}$. The quantity in $\{\cdot\}$ is a quadratic form in $x - \mu$. Note that the above way of writing the joint density resembles much more the univariate density than the explicit formula given at the beginning of the section.

The usefulness of this matrix representation is that the bivariate normal distribution now extends immediately to a general multivariate form, with joint density as given above with

$$\mathbf{X} = (X_1, \dots, X_k)^T, \quad \mathbf{x} = (x_1, \dots, x_k)^T, \quad \boldsymbol{\mu} = (\mu_1, \dots, \mu_k)^T$$

and

$$(\boldsymbol{\Sigma})_{ij} = \text{cov}(X_i, X_j) = \rho_{ij}\sigma_i\sigma_j$$

Further note that, since $\boldsymbol{\Sigma}$ is $k \times k$, we can write

$$|2\pi\boldsymbol{\Sigma}|^{1/2} = (2\pi)^{k/2}|\boldsymbol{\Sigma}|^{1/2}.$$

The usefulness of this matrix representation is that the bivariate normal distribution now extends immediately to a general multivariate form, with joint density as given above with

$$\mathbf{X} = (X_1, \dots, X_k)^T, \quad \mathbf{x} = (x_1, \dots, x_k)^T, \quad \boldsymbol{\mu} = (\mu_1, \dots, \mu_k)^T$$

and

$$(\boldsymbol{\Sigma})_{ij} = \text{cov}(X_i, X_j) = \rho_{ij}\sigma_i\sigma_j$$

Further note that, since $\boldsymbol{\Sigma}$ is $k \times k$, we can write

$$|2\pi\boldsymbol{\Sigma}|^{1/2} = (2\pi)^{k/2}|\boldsymbol{\Sigma}|^{1/2}.$$

For this k -dimensional joint distribution, denoted by $MN(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ or $N_k(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, $\rho_{ij} = \text{Corr}(X_i, X_j)$ and $\text{var}(X_i) = \sigma_i^2$.

It can then be shown that X_i has marginal distribution $N(\mu_i, \sigma_i^2)$, that any two of these variables have a bivariate normal distribution as above, and therefore that the conditional distribution of one variable given the other is also normal.

Example 1.33

Let X_1, X_2, X_3 have a trivariate normal distribution with mean vector (μ_1, μ_2, μ_3) and covariance matrix

$$\Sigma = \begin{bmatrix} a & 0 & 0 \\ 0 & b & 0 \\ 0 & 0 & c \end{bmatrix}.$$

Show that $f_{X_1, X_2, X_3} = f_{X_1} f_{X_2} f_{X_3}$ and give the marginal distributions of X_1, X_2 , and X_3 .



Chapter 2: Transformation of Variables

In this section we will see how to derive the distribution of transformed random variables. This is useful because many statistics applied to data analysis (e.g. test statistics) are transformations of the sample variables.

2.1 Univariate case

Suppose that we have a sample space Ω , a probability function P on Ω , a random variable $X : \Omega \rightarrow \mathbb{R}$, and a function $\phi : \mathbb{R} \rightarrow \mathbb{R}$.

Recall from §1.2 : $Y = \phi(X) : \Omega \rightarrow \mathbb{R}$ is defined by $Y(\omega) = \phi(X)(\omega) = \phi(X(\omega))$. Since $Y = \phi(X)$ is a random variable it also has a probability distribution, which can be determined either directly from P or via the distribution of X .

Discrete case:

$$\begin{aligned}\mathbf{P}(Y = y) &= \mathbf{P}(\{\omega : \phi(\mathbf{X}(\omega)) = y\}) = \sum_{\{\omega : \phi(\mathbf{X}(\omega)) = y\}} \mathbf{P}(\{\omega\}) \\ &= \sum_{\{\mathbf{x} : \phi(\mathbf{x}) = y\}} \mathbf{P}(\{\omega : \mathbf{X}(\omega) = \mathbf{x}\}) \\ &= \sum_{\{\mathbf{x} : \phi(\mathbf{x}) = y\}} p_{\mathbf{x}}(\mathbf{x}).\end{aligned}$$

Discrete case:

$$\begin{aligned} \mathbf{P}(Y = y) &= \mathbf{P}(\{\omega : \phi(X(\omega)) = y\}) = \sum_{\{\omega : \phi(X(\omega)) = y\}} \mathbf{P}(\{\omega\}) \\ &= \sum_{\{x : \phi(x) = y\}} \mathbf{P}(\{\omega : X(\omega) = x\}) \\ &= \sum_{\{x : \phi(x) = y\}} p_x(x). \end{aligned}$$

So, for example

$$\begin{aligned} \mathbf{E}\{Y\} &= \sum_{\omega} \phi(X(\omega)) \mathbf{P}(\{\omega\}) \quad \text{with respect to } \mathbf{P} \text{ on } \Omega \\ &= \sum_x \phi(x) p_x(x) \quad \text{with respect to distribution of } X \\ &= \sum_y y p_{\phi(X)}(y) \quad \text{with respect to distribution of } \phi(X) \end{aligned}$$

Examples

Example 2.1

Consider two independent throws of a fair die. Let X be the sum of the numbers that show up. Give the distribution of X .

Now consider the transformation $Y = (X - 7)^2$. Derive the distribution of Y . □

transformation of univariate continuous random variable: ϕ increasing

In general, suppose that $Y = \phi(X)$ where ϕ is a strictly increasing and differentiable function. Then,

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

transformation of univariate continuous random variable: ϕ increasing

In general, suppose that $Y = \phi(X)$ where ϕ is a strictly increasing and differentiable function. Then,

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

Then, differentiating, Y has density

$$f_Y(y) = f_X(\phi^{-1}(y)) \frac{d}{dy} \phi^{-1}(y) = f_X(x) \frac{dx}{dy} \bigg|_{x=\phi^{-1}(y)},$$

where the index $x = \phi^{-1}(y)$ means that any x in the formula has to be replaced by the inverse $\phi^{-1}(y)$ because $f_Y(y)$ is a function of y .

transformation of univariate continuous random variable: ϕ decreasing

Similarly, if ϕ is decreasing then

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

so that

$$f_Y(y) = -f_X(x) \frac{dx}{dy} \Big|_{x=\phi^{-1}(y)}.$$

transformation of univariate continuous random variable: ϕ decreasing

Similarly, if ϕ is decreasing then

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

so that

$$f_Y(y) = -f_X(x) \left. \frac{dx}{dy} \right|_{x=\phi^{-1}(y)}.$$

In the first case $dy/dx = d\phi(x)/dx$ is positive (since ϕ is increasing), in the second it is negative (since ϕ is decreasing) so either way the transformation formula is

$$f_Y(y) = f_X(x) \left| \frac{dx}{dy} \right|_{x=\phi^{-1}(y)}$$

We can check that the right-hand side of the above formula is a valid pdf as follows. Recall that $\int_{-\infty}^{\infty} f_X(x) dx = 1$. Changing variable to $y = \phi(x)$ we have, for ϕ increasing,

$$1 = \int_{-\infty}^{\infty} \left\{ f_X(x) \frac{dx}{dy} \right\}_{x=\phi^{-1}(y)} dy$$

so that $f_X(x) \left| \frac{dx}{dy} \right|$ is a valid pdf. Similarly for ϕ decreasing.

We can check that the right-hand side of the above formula is a valid pdf as follows. Recall that $\int_{-\infty}^{\infty} f_X(x) dx = 1$. Changing variable to $y = \phi(x)$ we have, for ϕ increasing,

$$1 = \int_{-\infty}^{\infty} \left\{ f_X(x) \frac{dx}{dy} \right\}_{x=\phi^{-1}(y)} dy$$

so that $f_X(x) \left| \frac{dx}{dy} \right|$ is a valid pdf. Similarly for ϕ decreasing.

Example 2.2

Consider $X \sim \text{Uniform}[-\frac{\pi}{2}, \frac{\pi}{2}]$, i.e.

$$f_X(x) = \begin{cases} \frac{1}{\pi} & -\frac{\pi}{2} \leq x \leq \frac{\pi}{2} \\ 0 & \text{otherwise.} \end{cases}$$

Derive the density of $Y = \tan(X)$.



Many-to-one functions

When ϕ is a *many-to-one* function we use the generalised formula $f_Y(y) = \sum f_X(x) \left| \frac{dx}{dy} \right|$, where the summation is over the set $\{x : h(x) = y\}$. That is, we add up the contributions to the density at y from all x values which map to y .

Example 2.3

Suppose that $f_X(x) = 2x$ on $(0, 1)$ and let $Y = (X - \frac{1}{2})^2$. Obtain the pdf of Y . □

Random Question to be Marked

A capacitor of capacitance C is used in a resonant circuit with an inductor of inductivity $L = 10^{-2}\text{H}$. The resonant frequency of the circuit is given by the formula

$$f = \frac{1}{2\pi\sqrt{LC}}$$

Table: Standard normal cdf

x	$\Phi(x)$
0.1	0.5398
0.2	0.5793
0.3	0.6179
0.4	0.6554
0.5	0.6915
0.6	0.7257
0.7	0.7580
0.8	0.7881
0.9	0.8159
1.0	0.8413

A past student project

found that $C \sim N(9.38 \cdot 10^{-8}\text{F}, (1.9 \cdot 10^{-9}\text{F})^2)$.

Event A	Event B	Event C
$f > 5233\text{Hz}$	$5181\text{Hz} < f \leq 5233\text{Hz}$	$f \leq 5181\text{Hz}$

What assignment to the outcomes “question B1”, “question B2”, “Resample” would you like?

2.2 Bivariate case

For the bivariate case we consider two random variables X, Y with joint density $f_{X,Y}(x, y)$. What is the joint density of transformations $U = u(X, Y), V = v(X, Y)$ where $u(\cdot, \cdot)$ and $v(\cdot, \cdot)$ are functions from \mathbb{R}^2 to \mathbb{R} , such as the ratio X/Y or the sum $X + Y$?

In order to use the following generalisation of the method of §2.1, we need to assume that u, v are such that each pair (x, y) defines a unique (u, v) and conversely, so that $u = u(x, y)$ and $v = v(x, y)$ are differentiable and invertible. The formula that gives the joint density of U, V is similar to the univariate case but the derivative, as we used it above, now has to be replaced by the *Jacobian* $J(u, v)$ of this transformation.

The result is that $U = u(X, Y)$, $V = v(X, Y)$ have joint density

$$f_{U,V}(u, v) = f_{X,Y}(x(u, v), y(u, v)) |J(x, y)| \Big|_{\substack{x=x(u,v) \\ y=y(u,v)}}$$

Again, the index $\substack{x=x(u,v) \\ y=y(u,v)}$ means that the x, y have to be replaced by the suitable transformations involving u, v only.

But how do we get the Jacobian $J(x, y)$? It is actually the determinant of the *matrix of partial derivatives* :

$$J(x, y) = \det \left(\frac{\partial(x, y)}{\partial(u, v)} \right) = \det \begin{pmatrix} \frac{\partial x}{\partial u} & \frac{\partial x}{\partial v} \\ \frac{\partial y}{\partial u} & \frac{\partial y}{\partial v} \end{pmatrix}$$

We finally take its absolute value, $|J(x, y)|$.

There are two ways of computing this:

(1) Obtain the inverse transformation

$x = x(u, v)$, $y = y(u, v)$, compute the matrix of partial derivatives $\partial(x, y)/\partial(u, v)$ and then its determinant and absolute value.

(2) Alternatively find the determinant $J(u, v)$ from the matrix of partial derivatives of (u, v) with respect to (x, y) and then its absolute value and invert this.

There are two ways of computing this:

(1) Obtain the inverse transformation

$x = x(u, v)$, $y = y(u, v)$, compute the matrix of partial derivatives $\partial(x, y)/\partial(u, v)$ and then its determinant and absolute value.

(2) Alternatively find the determinant $J(u, v)$ from the matrix of partial derivatives of (u, v) with respect to (x, y) and then its absolute value and invert this.

The two methods are equivalent since

$$\frac{\partial(x, y)}{\partial(u, v)} = \left\{ \frac{\partial(u, v)}{\partial(x, y)} \right\}^{-1}$$

Which way to choose in a specific case will depend on which functions are easier to derive. But note that the inverse transformations $x = x(u, v)$ and $y = y(u, v)$ are required anyway so that the first approach is often preferable.

Example 2.4

Let X and Y be two independent exponential variables with $X \sim \text{Exp}(\lambda)$ and $Y \sim \text{Exp}(\mu)$. Find the distribution of $U = X/Y$ □

Example 2.5

Consider two independent and identically distributed random variables X and Y having a uniform distribution on $[0, 2]$. Derive the joint density of $Z = X/Y$ and $W = Y$, stating the area where this density is positive. Are Z and W independent?

Obtain the marginal density of $Z = X/Y$. □

Sums of random variables

The distribution of a sum $Z = X + Y$ of two (not necessarily independent) random variables X and Y can be derived directly as follows.

In the discrete case note that the marginal distribution of Z is

$$P(Z = z) = \sum_x P(X = x, Z = z) = \sum_x P(X = x, Y = z - x)$$

Sums of random variables

The distribution of a sum $Z = X + Y$ of two (not necessarily independent) random variables X and Y can be derived directly as follows.

In the discrete case note that the marginal distribution of Z is

$$P(Z = z) = \sum_x P(X = x, Z = z) = \sum_x P(X = x, Y = z - x)$$

That is,

$$p_Z(z) = \sum_x p_{X,Y}(x, z - x)$$

Analogously, in the continuous case we get

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

Example 2.6

Let X and Y be two positive random variables with joint pdf

$$f_{X,Y}(x,y) = xye^{-(x+y)}, \quad x, y > 0.$$

Derive and name the distribution of their sum $Z = X + Y$. \square

2.3 Multivariate case

The ideas of §2.2 extend in a straightforward way to the case of more than two variables. The general problem is to find the distribution of $Y = \phi(X)$, where Y is $s \times 1$ and X is $r \times 1$, from the known distribution of X . Here X is the *random vector*

$$X = \begin{pmatrix} X_1 \\ X_2 \\ \cdot \\ \cdot \\ X_r \end{pmatrix}.$$

Discrete case.

$p_Y(y) = \sum p_X(x)$, where the summation is over the set $\{x : \phi(x) = y\}$. That is, we just add up the probabilities of all x -values that give $\phi(x) = y$.

Discrete case.

$p_Y(y) = \sum p_X(x)$, where the summation is over the set $\{x : \phi(x) = y\}$. That is, we just add up the probabilities of all x -values that give $\phi(x) = y$.

Continuous case.

Case (i): ϕ is a one-to-one transformation (so that $s = r$). Then the rule is

$$f_Y(y) = f_X(x(y)) |J(x)|_{x=x(y)}$$

where $J(x) = \det \left(\frac{dx}{dy} \right)$ is the Jacobian of transformation.

Here $\frac{dx}{dy}$ is the matrix of partial derivatives $\left(\frac{dx}{dy} \right)_{ij} = \frac{\partial x_i}{\partial y_j}$.

Case (ii): $s < r$. First transform the s -vector Y to the r -vector Y' , where $Y'_i = Y_i$, $i = 1, \dots, s$, and the other $r - s$ random variables Y'_i , $i = s + 1, \dots, r$, are chosen for convenience. Now find the density of Y' as in case (i) and then integrate out Y'_{s+1}, \dots, Y'_r to obtain the marginal density of Y , as required. (c.f. Examples 2.6 & 2.7 in the bivariate case.)

Case (iii): $s = r$ but $\phi(\cdot)$ is not monotonic. Then there will generally be more than one value of x corresponding to a given y and we need to add the probability contributions from all relevant x s.

Multivariate Transformation: Example

Example 2.7 (linear transformation)

Suppose that $Y = AX$, where A is an $r \times r$ nonsingular matrix. Then $f_Y(y) = f_X(A^{-1}y)|\det(A)|^{-1}$. □

2.4 Approximation of moments

Sometimes we may not need the complete probability distribution of $\phi(X)$, but just the first two moments. Recall that $\mathbb{E}(aX + b) = a\mathbb{E}(X) + b$, so the relation $\mathbb{E}\{\phi(X)\} = \phi(\mathbb{E}(X))$ is true whenever ϕ is a linear function. However, in general if $Y = \phi(X)$ it will not be true that $\mathbb{E}(Y) = \mathbb{E}\{\phi(X)\} = \int \phi(x)f_X(x)dx$ is the same as $\phi(\mathbb{E}(X)) = \phi(\int xf_X(x)dx)$, (or equivalent summations if X is discrete). (c.f. Example 1.1 ctd. in §1.2.)

To find moments of Y we can use the distribution of X , as above. However, the sums or integrals involved may be analytically intractable. In practice an approximate answer may be sufficient.

Intuitively, if X has mean μ_X and X is not very variable, then we would expect $\mathbb{E}(Y)$ to be quite close to $\phi(\mu_X)$. How good is this approximation?

Suppose that $\phi(x)$ is a continuous function of x for which the following Taylor expansion about μ_X exists (which requires the existence of the derivatives of ϕ):

$$\phi(x) = \phi(\mu_X) + (x - \mu_X)\phi'(\mu_X) + \frac{1}{2} (x - \mu_X)^2\phi''(\mu_X) + \dots$$

Suppose that $\phi(x)$ is a continuous function of x for which the following Taylor expansion about μ_X exists (which requires the existence of the derivatives of ϕ):

$$\phi(x) = \phi(\mu_X) + (x - \mu_X)\phi'(\mu_X) + \frac{1}{2} (x - \mu_X)^2\phi''(\mu_X) + \dots$$

Replacing x by X and taking expectations (or, equivalently, multiplying both sides of the above equation by $f_X(x)$ and integrating over x) term by term, we get

$$\begin{aligned}\mathbb{E}\{\phi(X)\} &= \phi(\mu_X) + \phi'(\mu_X) \underbrace{\mathbb{E}(X - \mu_X)}_{=0} + \\ &\quad \frac{1}{2} \phi''(\mu_X) \underbrace{\mathbb{E}\{(X - \mu_X)^2\}}_{\sigma_X^2} + \dots\end{aligned}$$

so that

$$\mathbb{E}(Y) \simeq \phi(\mu_X) + \frac{1}{2} \phi''(\mu_X) \sigma_X^2$$

A usually sufficiently good approximate formula for the variance is based on a first order approximation yielding

$$\begin{aligned}\text{Var}\{\phi(X)\} &= \mathbb{E}\{\phi(X) - \mathbb{E}(\phi(X))\}^2 \\ &\approx \mathbb{E}\{\phi(X) - \phi(\mu_X)\}^2 \approx \mathbb{E}\{(X - \mu_X)^2 (\phi'(\mu_X))^2\} \\ &= \{\phi'(\mu_X)\}^2 \mathbb{E}\{(X - \mu_X)^2\},\end{aligned}$$

**where we have used the approximation $\mathbb{E}(\phi(X)) \simeq \phi(\mu_X)$.
Therefore**

$$\text{var}(Y) \simeq \{\phi'(\mu_X)\}^2 \sigma_X^2$$

A usually sufficiently good approximate formula for the variance is based on a first order approximation yielding

$$\begin{aligned}\text{Var}\{\phi(X)\} &= \mathbb{E}\{\phi(X) - \mathbb{E}(\phi(X))\}^2 \\ &\approx \mathbb{E}\{\phi(X) - \phi(\mu_X)\}^2 \approx \mathbb{E}\{(X - \mu_X)^2 (\phi'(\mu_X))^2\} \\ &= \{\phi'(\mu_X)\}^2 \mathbb{E}\{(X - \mu_X)^2\},\end{aligned}$$

where we have used the approximation $\mathbb{E}(\phi(X)) \simeq \phi(\mu_X)$.
Therefore

$$\text{var}(Y) \simeq \{\phi'(\mu_X)\}^2 \sigma_X^2$$

Example 2.8

Consider a Poisson variable $X \sim \text{Poi}(\mu)$. Find approximations to the expectation and variance of $Y = \sqrt{X}$.



2.5 Order Statistics

Order statistics are a special kind of transformation of the sample variables. Their joint and marginal distributions can be derived by combinatorial considerations.

Suppose that X_1, \dots, X_n are independent with common density f_X . Denote the ordered values by

$X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}$. What is the distribution F_r of $X_{(r)}$?

In particular, $X_{(n)} = \max (X_1, \dots, X_n)$ is the **sample maximum and $X_{(1)} = \min (X_1, \dots, X_n)$ is the **sample minimum**.**

To find the distribution of $X_{(n)}$, note that $\{X_{(n)} \leq x\}$ and $\{\text{all } X_i \leq x\}$ are the same event – and so have the same probability!

Therefore the distribution function of $X_{(n)}$ is

$$\begin{aligned} F_n(x) = P(X_{(n)} \leq x) &= P(\text{all } X_i \leq x) \\ &= P(X_1 \leq x, X_2 \leq x, \dots, X_n \leq x) \\ &= \{F_X(x)\}^n \end{aligned}$$

since the X_i are independent with the same distribution function F_X . Thus

$$F_n(x) = \{F_X(x)\}^n$$

To find the distribution of $X_{(n)}$, note that $\{X_{(n)} \leq x\}$ and $\{\text{all } X_i \leq x\}$ are the same event – and so have the same probability!

Therefore the distribution function of $X_{(n)}$ is

$$\begin{aligned} F_n(x) = P(X_{(n)} \leq x) &= P(\text{all } X_i \leq x) \\ &= P(X_1 \leq x, X_2 \leq x, \dots, X_n \leq x) \\ &= \{F_X(x)\}^n \end{aligned}$$

since the X_i are independent with the same distribution function F_X . Thus

$$F_n(x) = \{F_X(x)\}^n$$

Furthermore, differentiating this expression we see that the density f_n of $X_{(n)}$ is

$$f_n(x) = n\{F_X(x)\}^{n-1}f_X(x)$$

Using a similar argument for $X_{(1)} = \min(X_1, \dots, X_n)$ we see that

$$\begin{aligned} F_1(x) &= P(X_{(1)} \leq x) = P(\text{at least one } X_i \leq x) \\ &= 1 - P(\text{all } X_i > x) = 1 - \{1 - F_X(x)\}^n \end{aligned}$$

so the distribution function of $X_{(1)}$ is

$$F_1(x) = 1 - \{1 - F_X(x)\}^n$$

and, differentiating, the pdf f_1 of $X_{(1)}$ is

$$f_1(x) = n\{1 - F_X(x)\}^{n-1} f_X(x)$$

Consider next the situation for general $1 \leq r \leq n$. For dx sufficiently small we have

$$P(x < X_{(r)} \leq x + dx) = P \left(\begin{array}{l} r - 1 \text{ values } X_i \text{ such that } X_i \leq x, \text{ and} \\ \text{one value in } (x, x + dx], \text{ and} \\ n - r \text{ values such that } X_i > x + dx \end{array} \right)$$

Consider next the situation for general $1 \leq r \leq n$. For dx sufficiently small we have

$$\begin{aligned}
 P(x < X_{(r)} \leq x + dx) &= P \left(\begin{array}{l} r-1 \text{ values } X_i \text{ such that } X_i \leq x, \text{ and} \\ \text{one value in } (x, x + dx], \text{ and} \\ n-r \text{ values such that } X_i > x + dx \end{array} \right) \\
 &= \underbrace{\frac{n!}{(r-1)!(n-r)!}}_{\text{no. of ways of ordering the } r-1, 1 \text{ and } n-r \text{ values}} \{F_X(x)\}^{r-1} f_X(x) dx \{1 - F_X(x + dx)\}^{n-r}
 \end{aligned}$$

Recalling that $f_r(x) = \lim_{dx \rightarrow 0} P(x < X_{(r)} \leq x + dx)/dx$, dividing both sides of the above expression by dx and letting $dx \rightarrow 0$ we obtain the density function of the r th order statistic $X_{(r)}$ as

$$f_r(x) = \frac{n!}{(r-1)!(n-r)!} \{F_X(x)\}^{r-1} \{1 - F_X(x)\}^{n-r} f_X(x)$$

Exercise: Show that this formula gives the previous densities when $r = n$ and $r = 1$.

A hint of extreme value theory

Example 2.9

A village is protected from a river by a dike of height h . The maximum water levels X_i reached by the river in subsequent years $i = 1, 2, 3, \dots$ are modelled as independent following an exponential distribution with mean $\lambda^{-1} = 10$. What is the probability that the village will be flooded (in statistical language this would be called a “threshold exceedance”) at least once in the next 100 years? How high does the dike need to be to make this probability smaller than 0.5?

Activity 2.1

Collect a sample of five student volunteers from the audience and have them measure their heights approximately. Record the heights and compute the mean height. Reorder the volunteers in the sample by height to obtain the first to fifth order statistic, $X_{(1)}, \dots, X_{(5)}$. What do the first, third and fifth order statistic, $X_{(1)}, X_{(3)}, X_{(5)}$ stand correspond to?

Random Question to be Marked

The height difference between the sample maximum $x_{(5)}$ and the sample minimum $x_{(1)}$ in Activity 2.1 in centimetres, compute the remainder under division by 3 and map to the question to be marked using the following table:

Remainder	Question to be marked
$(x_{(5)} - x_{(1)}) \bmod 3 = 0$	Ex5B2
$(x_{(5)} - x_{(1)}) \bmod 3 = 1$	Ex5B3
$(x_{(5)} - x_{(1)}) \bmod 3 = 2$	Ex5B1

3 Generating Functions

The transformation method presented in the previous parts may become tedious when a large number of variables is involved, in particular for transformations of the sample variables when the sample size tends to infinity. Generating functions provide an alternative way of determining a distribution (*e.g.* of sums of random variables).

We consider different generating functions for the discrete and continuous case. For the former, these are the probability generating functions and for the latter the moment generating functions. We further consider joint generating functions and apply these to linear combinations of random variables and finally use them to derive the Central Limit Theorem.

3.1 The probability generating function (pgf)

Suppose that X is a discrete random variable taking values $0, 1, 2, \dots$. Then the probability generating function (pgf) $G(z)$ of X is defined as

$$G(z) \equiv E(z^X)$$

The pgf is a function of particular interest, because it sometimes provides an easy way of determining the distribution of a discrete random variable.

Write $p_i = P(X = i)$, $i = 0, 1, \dots$. Then, by the usual expectation formula,

Coefficients of the Power Series

$$\begin{aligned} G(z) &= E(z^X) = \sum_{i=0}^{\infty} z^i p_i \\ &= p_0 + zp_1 + z^2 p_2 + \dots \end{aligned}$$

Thus $G(z)$ is a *power series* in z , and p_r is the coefficient of z^r .

Coefficients of the Power Series

$$\begin{aligned} G(z) &= E(z^X) = \sum_{i=0}^{\infty} z^i p_i \\ &= p_0 + zp_1 + z^2 p_2 + \dots \end{aligned}$$

Thus $G(z)$ is a *power series* in z , and p_r is the coefficient of z^r . Note that

$$|G(z)| \leq \sum_i |z|^i p_i \leq \sum_i p_i = 1$$

for all $|z| \leq 1$ and that $G(1) = \sum_i p_i = 1$. The sum is therefore convergent for (at least) $|z| \leq 1$.

We know from the theory of Taylor expansions that the r th derivative $G^{(r)}(0) = p_r r!$, yielding an expression for the probability p_r in terms of the r th derivative of G evaluated at $z = 0$:

$$p_r = \frac{G^{(r)}(0)}{r!}, \quad r = 0, 1, 2, \dots$$

In practice it is usually easier to find the power series expansion of G and extract p_r as the coefficient of z^r .

Moments

Whereas the probabilities are related to the derivatives of G at $z = 0$, it turns out that the moments of X are related to the derivatives of G at $z = 1$. To see this, note that

$$G'(z) = p_1 + 2zp_2 + 3z^2p_3 + \cdots = \sum_{i=1}^{\infty} iz^{i-1}p_i$$

so that

$$G'(1) = \sum_{i=1}^{\infty} ip_i = E(X)$$

Moments

Whereas the probabilities are related to the derivatives of G at $z = 0$, it turns out that the moments of X are related to the derivatives of G at $z = 1$. To see this, note that

$$G'(z) = p_1 + 2zp_2 + 3z^2p_3 + \cdots = \sum_{i=1}^{\infty} iz^{i-1}p_i$$

so that

$$G'(1) = \sum_{i=1}^{\infty} ip_i = E(X)$$

Thus, we have

$$E(X) = G'(1)$$

Moments: 2

Further,

$$G''(z) = \sum_{i=2}^{\infty} i(i-1)z^{i-2}p_i$$

so that

$$G''(1) = \sum i(i-1)p_i = E\{X(X-1)\}$$

But we can write

$$\text{var}(X) = E(X^2) - \{E(X)\}^2 = E\{X(X-1)\} + E(X) - \{E(X)\}^2$$

Moments: 2

Further,

$$G''(z) = \sum_{i=2}^{\infty} i(i-1)z^{i-2}p_i$$

so that

$$G''(1) = \sum i(i-1)p_i = E\{X(X-1)\}$$

But we can write

$$\text{var}(X) = E(X^2) - \{E(X)\}^2 = E\{X(X-1)\} + E(X) - \{E(X)\}^2$$

from which we obtain the formula

$$\boxed{\text{var}(X) = G''(1) + G'(1) - \{G'(1)\}^2}$$

Examples

Example 3.1

Let $X \sim \text{Poi}(\mu)$. Find the pgf $G(z) = \mathbb{E}(z^X)$. Also, verify the above formulae for the expectation and variance of X . \square

Example 3.2

Consider the pgf

$$G(z) = (1 - p + pz)^n$$

where $0 < p < 1$ and $n \geq 1$ is an integer. Find the power expansion of $G(z)$ and hence derive the distribution of the random variable X that has this pgf. \square

3.2 The moment generating function (mgf)

Another function of special interest, particularly for continuous variables, is the moment generating function $M(s)$ of X , defined as

$$M(s) \equiv \mathbb{E}(e^{sX}) = \int_{-\infty}^{\infty} e^{sx} f_X(x) dx$$

The moment generating function does not necessarily exist for all $s \in \mathbb{R}$, i.e. the integral might be infinite. However, we assume for the following that $M(s)$ is finite for s in some open interval containing zero.

Using the expansion $e^{sx} = 1 + sx + \frac{1}{2!} s^2 x^2 + \dots$ we get

$$M(s) = \int_{-\infty}^{\infty} \sum_{n=0}^{\infty} \frac{s^n x^n}{n!} f_x(x) dx = \sum_{n=0}^{\infty} \frac{s^n}{n!} \int_{-\infty}^{\infty} x^n f_x(x) dx$$

integrating term by term (there are no convergence problems here due to assuming finiteness). It follows that

$$M(s) = \sum_{n=0}^{\infty} \frac{s^n}{n!} E(X^n).$$

Thus $M(s)$ is a power series in s and the coefficient of s^n is $E(X^n)/n!$ – hence the name ‘moment generating function’. Again, from the theory of Taylor expansions the r th derivative, $M^{(r)}(0)$, of $M(s)$ at $s = 0$ must therefore equal the r th (raw) moment $E(X^r)$ of X .

In particular we have $M'(0) = E(X)$ and $M''(0) = E(X^2)$, so that

$$E(X) = M'(0)$$

and

$$\text{Var}(X) = M''(0) - \{M'(0)\}^2$$

(Alternatively, and more directly, note that $M'(s) = E(Xe^{sX})$ and $M''(s) = E(X^2 e^{sX})$ and set $s = 0$.)

Note also that $M(0) = E(e^0) = 1$.

It can be shown that if the moment generating function exists on an open interval including zero then it uniquely determines the distribution.

Examples

The pgf tends to be used more for discrete distributions and the mgf for continuous ones, although note that when X takes nonnegative integer values then the two are related by $M(s) = \mathbb{E}(e^{sX}) = \mathbb{E}\{(e^s)^X\} = G(e^s)$.

Example 3.3

Suppose that X has a Gamma distribution with parameters (α, λ) . Find the mgf of X . Use this to derive the expectation and variance of X . □

Example 3.4

Let $X \sim N(\mu, \sigma^2)$ be a standard normal variable. Find the mgf of X . Use this mgf to obtain the expectation and variance of X . □

Linear transformation property

Suppose that $Y = a + bX$ and that we know the mgf $M_X(s)$ of X . What is the mgf of Y ?

We have

$$M_Y(s) = \mathbb{E}(e^{sY}) = \mathbb{E}\{e^{s(a+bX)}\} = e^{sa}\mathbb{E}(e^{sbX}) = e^{as}M_X(bs).$$

We can therefore easily obtain the mgf of any linear function of X from the mgf of X .

Linear transformation property

Suppose that $Y = a + bX$ and that we know the mgf $M_X(s)$ of X . What is the mgf of Y ?

We have

$$M_Y(s) = \mathbb{E}(e^{sY}) = \mathbb{E}\{e^{s(a+bX)}\} = e^{sa}\mathbb{E}(e^{sbX}) = e^{as}M_X(bs).$$

We can therefore easily obtain the mgf of any linear function of X from the mgf of X .

Example 3.5 (3.4 ctd)

Use the mgf of $X \sim N(\mu, \sigma^2)$ to find the distribution of $Y = a + bX$. □

mgf vs. characteristic function

A more general concept is that of the characteristic function. This is defined in a similar way to the mgf and has similar properties but involves complex variables. The main advantage over the moment generating function is that the characteristic function of a random variable always exists. However, we will not consider it here.

3.3 Joint generating functions

So far we have considered the pgf or mgf of a single real variable. The joint distribution of a *collection* of random variables X_1, \dots, X_n can be characterised in a similar way by the *joint* generating functions:

The joint pgf $G(z_1, \dots, z_n)$ of variables X_1, \dots, X_n is a function of n variables, z_1, \dots, z_n , and defined to be

$$G(z_1, \dots, z_n) = E(z_1^{X_1} z_2^{X_2} \dots z_n^{X_n})$$

The joint mgf $M(\mathbf{s}_1, \dots, \mathbf{s}_n)$ is a function of n variables, $\mathbf{s}_1, \dots, \mathbf{s}_n$, and is defined to be

$$M(\mathbf{s}_1, \dots, \mathbf{s}_n) = \mathbb{E}(e^{\mathbf{s}_1 X_1 + \dots + \mathbf{s}_n X_n})$$

These generating functions *uniquely* determine the joint distribution of X_1, \dots, X_n . Note that the mgf may also be written in vector notation as

$$M(\mathbf{s}) = \mathbb{E}(e^{\mathbf{s}^T \mathbf{X}}),$$

where $\mathbf{s} = (\mathbf{s}_1, \dots, \mathbf{s}_n)^T$ and $\mathbf{X} = (X_1, \dots, X_n)^T$.

Independence means Factorise

In both cases we find that if X_1, \dots, X_n are independent random variables then the pgf / mgf are given as the product of the individual pgfs / mgfs. (Recall: $\mathbb{E}(XY) = \mathbb{E}(X)\mathbb{E}(Y)$ when X, Y are independent.)

$$G(z_1, \dots, z_n) = \mathbb{E}(z_1^{X_1} \dots z_n^{X_n}) = \mathbb{E}(z_1^{X_1}) \dots \mathbb{E}(z_n^{X_n})$$

i.e. joint pgf = product of marginal pgfs

$$M(s_1, \dots, s_n) = \mathbb{E}(e^{s_1 X_1} \dots e^{s_n X_n}) = \mathbb{E}(e^{s_1 X_1}) \dots \mathbb{E}(e^{s_n X_n})$$

i.e. joint mgf = product of marginal mgfs.

The above property can be used to characterise independence because it can be shown that the factorisation of the joint mgf holds **if and only if** the variables are independent.

Marginal mgfs

It is straightforward to see that if $M_{X,Y}(s_1, s_2)$ is the joint mgf of X, Y then the marginal mgf of X is given by

$$M_X(s_1) = M_{X,Y}(s_1, 0).$$

(Proof: $\mathbb{E}(e^{s_1 X}) = \mathbb{E}(e^{s_1 X + 0 \cdot Y}) = M_{X,Y}(s_1, 0)$.)

Higher Moments

The joint moment generating function can further be useful to find higher moments of a distribution. More precisely, we can compute $\mathbb{E}(X_i^r X_j^k)$ in the following way

- 1 differentiate $M(s_1, \dots, s_n)$ r times w.r.t. s_i ;
- 2 further differentiate k times w.r.t. s_j ;
- 3 then set all $s_1 = \dots = s_n = 0$.

Following the above steps we get

$$\begin{aligned}\frac{\partial^r M}{\partial s_i^r} &= \mathbb{E}(X_i^r e^{s^T X}) \\ \frac{\partial^{r+k} M}{\partial s_i^r \partial s_j^k} &= \mathbb{E}(X_i^r X_j^k e^{s^T X}),\end{aligned}$$

which gives $\mathbb{E}(X_i^r X_j^k)$ on setting $s = 0$.

Linear transformation property

This is the multivariate generalisation of the univariate transformation property.

Suppose that $Y = a + bX$. Then the mgf of Y is (c.f. Section 3.3)

$$\begin{aligned}M_Y(s) &= \mathbb{E}(e^{s^T Y}) = \mathbb{E}\{e^{s^T(a+bX)}\} \\&= e^{s^T a} \mathbb{E}(e^{bs^T X}) = e^{a^T s} M_X(bs)\end{aligned}$$

Example 3.6: Multivariate normal

Suppose X_1, \dots, X_n are jointly multivariate normally distributed. Then, from equation (1.8) in 1.5.5, the density of $X = (X_1, \dots, X_n)$ is given in matrix notation by

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

where

$$\mathbf{x} = (x_1, \dots, x_n)^T, \quad \mu = (\mu_1, \dots, \mu_n)^T, \quad \Sigma = \text{covariance matrix}$$

ie the $\mu_i = E(X_i)$ are the individual expectations and the $\sigma_{ij} = (\Sigma)_{ij} = \text{cov}(X_i, X_j)$ are the pairwise covariances (variances if $i = j$).

The joint mgf is then (result will be derived next week):

$$M(\mathbf{s}_1, \dots, \mathbf{s}_n) = \exp \left\{ \mathbf{s}^T \boldsymbol{\mu} + \frac{1}{2} \mathbf{s}^T \boldsymbol{\Sigma} \mathbf{s} \right\}$$

***Exercise:* from this derive the mgf of the univariate $N(\mu, \sigma^2)$ distribution.**

Using the multinormal mgf

Example 3.6

Suppose that

$$\begin{pmatrix} X_1 \\ X_2 \end{pmatrix} \sim N_2 \left(\begin{pmatrix} 2 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 & 0.7 \\ 0.7 & 1 \end{pmatrix} \right)$$

Find $E(X_1 X_2)$ (i) directly and (ii) from the joint mgf.



Example 3.5: Multivariate normal

Suppose X_1, \dots, X_n are jointly multivariate normally distributed. Then, from equation (1.8) in §1.5.5, the density of $X = (X_1, \dots, X_n)$ is given in matrix notation by

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

where

$$\mathbf{x} = (x_1, \dots, x_n)^T, \quad \mu = (\mu_1, \dots, \mu_n)^T, \quad \Sigma = \text{covariance matrix}$$

ie the $\mu_i = E(X_i)$ are the individual expectations and the $\sigma_{ij} = (\Sigma)_{ij} = \text{cov}(X_i, X_j)$ are the pairwise covariances (variances if $i = j$).

We obtain the joint mgf as follows. We have

$$\mathbb{E}(\mathbf{e}^{s_1 X_1 + \dots + s_n X_n}) = \int \frac{1}{|2\pi \Sigma|^{1/2}} \cdot \exp \left\{ \mathbf{s}^T \mathbf{x} - \frac{1}{2} (\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu) \right\} d\mathbf{x}$$

(Remember that the integral here represents an n -dimensional integral.) To evaluate the integral we need to complete the square in $\{\cdot\}$. The result (derived in lectures) is that

$$M(\mathbf{s}_1, \dots, \mathbf{s}_n) = \exp \left\{ \mathbf{s}^T \mu + \frac{1}{2} \mathbf{s}^T \Sigma \mathbf{s} \right\}$$

Exercise: from this derive the mgf of the univariate $N(\mu, \sigma^2)$ distribution.

$$\mathbb{E}(X_i X_j)$$

For illustration, we now derive the joint moment $\mathbb{E}(X_i X_j)$. Differentiate first with respect to s_i . Since

$$\mathbf{s}^T \boldsymbol{\mu} = \sum_{k=1}^n \mathbf{s}_k \mu_k, \quad \mathbf{s}^T \boldsymbol{\Sigma} \mathbf{s} = \sum_{k=1}^n \sum_{l=1}^n \mathbf{s}_k \mathbf{s}_l \sigma_{kl}$$

we see that $\partial(\mathbf{s}^T \boldsymbol{\mu}) / \partial s_i = \mu_i$.

Also the terms involving s_i in $\mathbf{s}^T \boldsymbol{\Sigma} \mathbf{s}$ are

$$s_i^2 \sigma_{ii} + 2s_i \sum_{l \neq i} s_l \sigma_{il}$$

giving

$$\partial(\mathbf{s}^T \boldsymbol{\Sigma} \mathbf{s}) / \partial s_i = 2s_i \sigma_{ii} + 2 \sum_{l \neq i} s_l \sigma_{il} = 2 \sum_l s_l \sigma_{il}$$

$$\mathbb{E}(X_i X_j)$$

Therefore

$$\begin{aligned}\mathbb{E}(X_i \mathbf{e}^{\mathbf{s}^T \mathbf{X}}) &= \partial M(\mathbf{s}) / \partial s_i \\ &= \left(\mu_i + \sum_{l=1}^n \sigma_{il} s_l \right) \exp \left\{ \mathbf{s}^T \boldsymbol{\mu} + \frac{1}{2} \mathbf{s}^T \boldsymbol{\Sigma} \mathbf{s} \right\}\end{aligned}$$

$$\mathbb{E}(X_i X_j)$$

Therefore

$$\begin{aligned}\mathbb{E}(X_i e^{\mathbf{s}^T \mathbf{X}}) &= \partial M(\mathbf{s}) / \partial s_i \\ &= \left(\mu_i + \sum_{l=1}^n \sigma_{il} s_l \right) \exp \left\{ \mathbf{s}^T \boldsymbol{\mu} + \frac{1}{2} \mathbf{s}^T \boldsymbol{\Sigma} \mathbf{s} \right\}\end{aligned}$$

Now differentiate again with respect to s_j , $j \neq i$, to give

$$\begin{aligned}\mathbb{E}(X_i X_j e^{\mathbf{s}^T \mathbf{X}}) &= \left\{ \sigma_{ij} + \left(\mu_i + \sum_{l=1}^n \sigma_{il} s_l \right) \left(\mu_j + \sum_{l=1}^n \sigma_{jl} s_l \right) \right\} \\ &\quad \cdot \exp \left\{ \mathbf{s}^T \boldsymbol{\mu} + \frac{1}{2} \mathbf{s}^T \boldsymbol{\Sigma} \mathbf{s} \right\}.\end{aligned}$$

Setting $\mathbf{s} = \mathbf{0}$ now gives $\mathbb{E}(X_i X_j) = \sigma_{ij} + \mu_i \mu_j$ and therefore $\text{cov}(X_i, X_j) = \mathbb{E}(X_i X_j) - \mu_i \mu_j = \sigma_{ij}$, as claimed in Section 1.5.2.

Using the multinormal mgf

Example 3.7

Suppose that

$$\begin{pmatrix} X_1 \\ X_2 \end{pmatrix} \sim N_2 \left(\begin{pmatrix} 2 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 & 0.7 \\ 0.7 & 1 \end{pmatrix} \right)$$

Find $E(X_1 X_2)$ (i) directly and (ii) from the joint mgf.



Using the multinormal mgf

Example 3.7

Suppose that

$$\begin{pmatrix} X_1 \\ X_2 \end{pmatrix} \sim N_2 \left(\begin{pmatrix} 2 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 & 0.7 \\ 0.7 & 1 \end{pmatrix} \right)$$

Find $E(X_1 X_2)$ (i) directly and (ii) from the joint mgf. □

Remark.

The multivariate normal density of X_1, \dots, X_n is only valid when Σ is a non-singular matrix, which can be shown to hold if and only if no exact linear relationships exist between X_1, \dots, X_n . However, the multivariate normal distribution can still be defined when Σ is singular. Note in particular that the joint mgf is valid even when Σ is singular.

3.4 Linear combinations of random variables

We will now use the above methods to derive (properties of) the distribution of linear combinations of random variables.

Let X_1, \dots, X_n be the original variables. A linear combination is defined by

$$Y = a_1 X_1 + \dots + a_n X_n,$$

for any real-valued constants a_1, \dots, a_n . A popular linear combination is for example the sample mean $Y = \bar{X}$, for which $a_i = 1/n$, $i = 1, \dots, n$.

Let us first find the expectation and variance of the linear combination Y in terms of the moments of the X_i . The methods from § 1.4 can be used for this purpose. First

$$E(Y) = a_1 E(X_1) + \cdots + a_n E(X_n)$$

regardless of whether or not the X_i are independent.

Let us first find the expectation and variance of the linear combination Y in terms of the moments of the X_i . The methods from § 1.4 can be used for this purpose. First

$$E(Y) = a_1 E(X_1) + \cdots + a_n E(X_n)$$

regardless of whether or not the X_i are independent. For the variance we have

$$\text{var}(Y) = \text{cov}\left(\sum_i a_i X_i, \sum_j a_j X_j\right)$$

Let us first find the expectation and variance of the linear combination Y in terms of the moments of the X_j . The methods from § 1.4 can be used for this purpose. First

$$E(Y) = a_1 E(X_1) + \cdots + a_n E(X_n)$$

regardless of whether or not the X_j are independent. For the variance we have

$$\begin{aligned} \text{var}(Y) &= \text{cov}\left(\sum_i a_i X_i, \sum_j a_j X_j\right) \\ &= \sum_i \sum_j a_i a_j \text{cov}(X_i, X_j) \end{aligned}$$

Let us first find the expectation and variance of the linear combination Y in terms of the moments of the X_i . The methods from § 1.4 can be used for this purpose. First

$$E(Y) = a_1 E(X_1) + \cdots + a_n E(X_n)$$

regardless of whether or not the X_i are independent. For the variance we have

$$\begin{aligned}\text{var}(Y) &= \text{cov}\left(\sum_i a_i X_i, \sum_j a_j X_j\right) \\ &= \sum_i \sum_j a_i a_j \text{cov}(X_i, X_j) \\ &= \sum_i a_i^2 \text{var}(X_i) + \sum_{i \neq j} a_i a_j \text{cov}(X_i, X_j).\end{aligned}$$

Let us first find the expectation and variance of the linear combination Y in terms of the moments of the X_i . The methods from § 1.4 can be used for this purpose. First

$$E(Y) = a_1 E(X_1) + \cdots + a_n E(X_n)$$

regardless of whether or not the X_i are independent. For the variance we have

$$\begin{aligned}\text{var}(Y) &= \text{cov}\left(\sum_i a_i X_i, \sum_j a_j X_j\right) \\ &= \sum_i \sum_j a_i a_j \text{cov}(X_i, X_j) \\ &= \sum_i a_i^2 \text{var}(X_i) + \sum_{i \neq j} a_i a_j \text{cov}(X_i, X_j).\end{aligned}$$

In particular, we see that if the X_i are *independent* then

$$\text{var}(\sum a_i X_i) = \sum_i \text{var}(a_i X_i) = \sum a_i^2 \text{var}(X_i)$$

but not otherwise.

In vector notation, if we write $Y = a^T X$, where $X = (X_1, \dots, X_n)^T$, $a = (a_1, \dots, a_n)^T$ then this relations becomes

$$\mathbb{E}(Y) = a^T \mu, \quad \text{var}(Y) = a^T \Sigma a,$$

where $\mu = \mathbb{E}(X)$, $\Sigma = \text{Cov}(X)$.

In vector notation, if we write $Y = a^T X$, where $X = (X_1, \dots, X_n)^T$, $a = (a_1, \dots, a_n)^T$ then this relations becomes

$$\mathbb{E}(Y) = a^T \mu, \quad \text{var}(Y) = a^T \Sigma a,$$

where $\mu = \mathbb{E}(X)$, $\Sigma = \text{Cov}(X)$.

Now we want to find out about the actual distribution of Y . If we have the joint distribution of X_1, \dots, X_n we could proceed by transformation similar to §2.2, but this could be very tedious if n is large. Instead, let us explore an approach based on the joint pgf / mgf of the X_i . (The result below for the mgf is actually just a special case of the earlier linear transformation property of a joint mgf.)

If Y is discrete we find for its pgf

$$G_Y(z) = E(z^Y) = E(z^{a_1 X_1 + \dots + a_n X_n}) = E(z^{a_1 X_1} z^{a_2 X_2} \dots z^{a_n X_n}),$$

which is the joint pgf of X_1, \dots, X_n evaluated at $z_i = z^{a_i}$, i.e. $G(z^{a_1}, \dots, z^{a_n})$.

If Y is discrete we find for its pgf

$$G_Y(z) = E(z^Y) = E(z^{a_1 X_1 + \dots + a_n X_n}) = E(z^{a_1 X_1} z^{a_2 X_2} \dots z^{a_n X_n}),$$

which is the joint pgf of X_1, \dots, X_n evaluated at $z_i = z^{a_i}$, i.e. $G(z^{a_1}, \dots, z^{a_n})$.

Similarly, if Y is continuous its mgf is given as

$$M_Y(s) = E(e^{sY}) = E(e^{s(a_1 X_1 + \dots + a_n X_n)}) = E(e^{sa_1 X_1 + \dots + sa_n X_n}),$$

which is the joint mgf of X_1, \dots, X_n evaluated at $s_i = sa_i$, i.e. $M(sa_1, \dots, sa_n)$.

So, an alternative to the transformation method is to obtain the *joint* pgf or mgf and use this to derive the pgf or mgf of Y . Of course, we still have to get from there to the probability mass function or density if needed — but often the generating function of the univariate Y is known to belong to a specific distribution.

So, an alternative to the transformation method is to obtain the *joint* pgf or mgf and use this to derive the pgf or mgf of Y . Of course, we still have to get from there to the probability mass function or density if needed — but often the generating function of the univariate Y is known to belong to a specific distribution.

A simplification is available if X_1, \dots, X_n are independent. In this case we have

$$G_Y(z) = E(z^{a_1 X_1} z^{a_2 X_2} \dots z^{a_n X_n}) = \prod_{i=1}^n E(z^{a_i X_i}) = \prod_{i=1}^n G_{X_i}(z^{a_i}),$$

which is the product of the individual pgfs $G_{X_i}(z^{a_i})$ of the X_i evaluated at $z_i = z^{a_i}$.

In the case of independent continuous random variables the mgfs factorise similarly:

$$M_Y(\mathbf{s}) = \mathbb{E}(e^{sa_1 X_1 + \dots + sa_n X_n}) = \prod_{i=1}^n \mathbb{E}(e^{sa_i X_i}) = \prod_{i=1}^n M_{X_i}(sa_i),$$

which is the product of the individual mgfs $M_{X_i}(sa_i)$ of the X_i evaluated at $s_i = sa_i$.

In the case of independent continuous random variables the mgfs factorise similarly:

$$M_Y(\mathbf{s}) = \mathbb{E}(e^{s a_1 X_1 + \dots + s a_n X_n}) = \prod_{i=1}^n \mathbb{E}(e^{s a_i X_i}) = \prod_{i=1}^n M_{X_i}(s a_i),$$

which is the product of the individual mgfs $M_{X_i}(s a_i)$ of the X_i evaluated at $s_i = s a_i$.

Example 3.8 (3.2 ctd.)

Consider independent random variables X_1, \dots, X_n with $X_i \sim \text{Bin}(m_i, p)$, i.e. they have different numbers of trials m_i but the same success probability p . Find the pgf and the distribution of $Y = \sum X_i$.

Example 3.9 (3.3 ctd.)

Let X_1, \dots, X_n be independent with $X_i \sim \text{Gam}(\alpha_i, \lambda)$. Find the mgf and the distribution of $Y = \sum X_i$.

Example 3.8

Reminder: $Y = \sum_{i=1}^n a_i X_i$

If Y is discrete we find for its pgf

$$G_Y(z) = E(z^Y) = E(z^{a_1 X_1 + \dots + a_n X_n}) = E(z^{a_1 X_1} z^{a_2 X_2} \dots z^{a_n X_n}),$$

which is the joint pgf of X_1, \dots, X_n evaluated at $z_i = z^{a_i}$, i.e. $G(z^{a_1}, \dots, z^{a_n})$.

Example 3.8

Reminder: $Y = \sum_{i=1}^n a_i X_i$

If Y is discrete we find for its pgf

$$G_Y(z) = E(z^Y) = E(z^{a_1 X_1 + \dots + a_n X_n}) = E(z^{a_1 X_1} z^{a_2 X_2} \dots z^{a_n X_n}),$$

which is the joint pgf of X_1, \dots, X_n evaluated at $z_i = z^{a_i}$, i.e. $G(z^{a_1}, \dots, z^{a_n})$.

Example 3.8 (3.2 ctd.)

Consider independent random variables X_1, \dots, X_n with $X_i \sim \text{Bin}(m_i, p)$, i.e. they have different numbers of trials m_i but the same success probability p . Find the pgf and the distribution of $Y = \sum X_i$.

Example 3.9

Reminder: $Y = \sum_{i=1}^n a_i X_i$

In the case of independent continuous random variables X_i the mgfs factorise similarly:

$$M_Y(\mathbf{s}) = \mathbb{E}(e^{sa_1 X_1 + \dots + sa_n X_n}) = \prod_{i=1}^n \mathbb{E}(e^{sa_i X_i}) = \prod_{i=1}^n M_{X_i}(sa_i),$$

which is the product of the individual mgfs $M_{X_i}(sa_i)$ of the X_i evaluated at $s_i = sa_i$.

Example 3.9

Reminder: $Y = \sum_{i=1}^n a_i X_i$

In the case of independent continuous random variables X_i the mgfs factorise similarly:

$$M_Y(\mathbf{s}) = \mathbb{E}(e^{sa_1 X_1 + \dots + sa_n X_n}) = \prod_{i=1}^n \mathbb{E}(e^{sa_i X_i}) = \prod_{i=1}^n M_{X_i}(sa_i),$$

which is the product of the individual mgfs $M_{X_i}(sa_i)$ of the X_i evaluated at $s_i = sa_i$.

Example 3.9

Let X_1, \dots, X_n be independent with $X_i \sim \text{Gam}(\alpha_i, \lambda)$. Find the mgf and the distribution of $Y = \sum X_i$.

Announcement: Peer Learning Sessions

- Full solution and marking scheme
- Dept. of Statistical Science, Tues, Dec 1st, 6pm-7pm
- Room 102: Question Ex4C2, Room 116: Ex4C3, Room B17: Ex3C2

Vote for week 9's section C question you want solved on Moodle!

course in statistics, such as STAT1004 and STAT1005 or MATH7501. STAT2001 is of intermediate level and mostly intended for 2nd year statistics as well as Mathematics and Statistics students. STAT3101 is of advanced level and mostly intended for 3rd year mathematics and quantitatively oriented natural sciences students. The courses STAT2001 and STAT3101 share lectures but there are differences in the assessments.

My office is Room 144 in the Department of Statistical Science. Office hours Tuesdays 3pm-3:55pm and Fridays 4:15pm-5:30pm and by email appointment for those who cannot make this time.



News forum



General Discussion Forum

Discussion forum for general questions concerning STAT2001/3101 including questions on homework assignments as long as they do not give away the answer. Moderated by the lecturer on a weekly basis.

Topic 1

Lectures



Lecture times



Topic to be explained again in lecture?



2014/15 Current Lecture Notes, now including chapter 4 564.2KB PDF document

Moodle Hot Question:

Submit and Vote topic you want re-explained in the lecture!

Example 3.10 (3.6 ctd): Multivariate normal

Suppose X_1, \dots, X_n are jointly multivariate normally distributed. Consider again the linear transformation $Y = a_1 X_1 + \dots + a_n X_n$, or $Y = \mathbf{a}^T \mathbf{X}$ in vector notation, where $\mathbf{a} = (a_1, \dots, a_n)^T$, $\mathbf{X} = (X_1, \dots, X_n)^T$. It follows from the general result that the mgf of Y is

$$M_Y(s) = M_X(\mathbf{s}\mathbf{a}) = \mathbb{E}\{\exp(\mathbf{s}\mathbf{a}^T \mathbf{X})\} = \exp\left(\mathbf{s}\mathbf{a}^T \mu + \frac{1}{2} \mathbf{s}^2 \mathbf{a}^T \Sigma \mathbf{a}\right)$$

from chapter 3.4. By comparison with the univariate mgf (see Example 3.4) we see that

$$\boxed{Y = \mathbf{a}^T \mathbf{X} \sim N(\mathbf{a}^T \mu, \mathbf{a}^T \Sigma \mathbf{a})}$$

We have seen earlier that for any random vector X we have $E(a^T X) = a^T \mu$, $\text{var}(a^T X) = a^T \Sigma a$, so the importance of the foregoing result is that any linear combination of normal variables is itself normally distributed even if the variables are correlated (and thus not independent).

Remark: Positive semi-definite and positive-definite

Definition:

A matrix $\Sigma \in \mathbb{R}^{n \times n}$ is called positive semi-definite if for all vectors $a \in \mathbb{R}^n$ we have $a^T \Sigma a \geq 0$. It is called positive definite if for all vectors $a \in \mathbb{R}^n \setminus \{0\}$ we have $a^T \Sigma a > 0$.

Remark: Positive semi-definite and positive-definite

Definition:

A matrix $\Sigma \in \mathbb{R}^{n \times n}$ is called **positive semi-definite** if for all vectors $a \in \mathbb{R}^n$ we have $a^T \Sigma a \geq 0$. It is called **positive definite** if for all vectors $a \in \mathbb{R}^n \setminus \{0\}$ we have $a^T \Sigma a > 0$.

Result:

All covariance matrices are positive semi-definite.

Remark: Positive semi-definite and positive-definite

Definition:

A matrix $\Sigma \in \mathbb{R}^{n \times n}$ is called **positive semi-definite** if for all vectors $a \in \mathbb{R}^n$ we have $a^T \Sigma a \geq 0$. It is called **positive definite** if for all vectors $a \in \mathbb{R}^n \setminus \{0\}$ we have $a^T \Sigma a > 0$.

Result:

All covariance matrices are positive semi-definite.

Sketch Proof:

From the last lecture, we know that for

$Y = a^T X = \sum_{i=1}^n a_i X_i$ with X_1, \dots, X_n being a random vector of size n and any $a \in \mathbb{R}^n$ we have $\text{Var}(Y) = a^T \Sigma a$.

Therefore, $a^T \Sigma a = \text{Var}(a^T X) \geq 0$ because all variances are non-negative.

Example 3.11 (3.7 ctd.)

Example 3.7

Suppose that

$$\begin{pmatrix} X_1 \\ X_2 \end{pmatrix} \sim N_2 \left(\begin{pmatrix} 2 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 & 0.7 \\ 0.7 & 1 \end{pmatrix} \right)$$

Find $E(X_1 X_2)$ (i) directly and (ii) from the joint mgf. □

Example 3.11 (3.7 ctd.)

Use the joint mgf of X_1 and X_2 to find the distribution of $Y = X_1 - X_2$. Is the covariance matrix positive semi-definite? Is it positive definite?

The Central Limit Theorem

If X_1, X_2, \dots are i.i.d. with mean μ and variance $\sigma^2 < \infty$, then

$$\frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \xrightarrow{d} N(0, 1)$$

as $n \rightarrow \infty$. This means that

$$P(a < \frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \leq b) \approx \Phi(b) - \Phi(a)$$

The Central Limit Theorem

If X_1, X_2, \dots are i.i.d. with mean μ and variance $\sigma^2 < \infty$, then

$$\frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \xrightarrow{d} N(0, 1)$$

as $n \rightarrow \infty$. This means that

$$P\left(a < \frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \leq b\right) \approx \Phi(b) - \Phi(a)$$

The Central Limit Theorem is amazing: From **next to no assumptions** (independence and equal means and variances) we arrive at a **phenomenally strong conclusion**: The sample mean asymptotically follows a Gaussian distribution. Regardless of which particular distribution the X_n follow, the resulting distribution is always Gaussian!

Sample Mean and CLT

Thus we can compute probabilities about \bar{X}_n from

$$\mathbf{P}(a < \bar{X}_n \leq b) \approx \Phi\left(\frac{\sqrt{n}(b - \mu)}{\sigma}\right) - \Phi\left(\frac{\sqrt{n}(a - \mu)}{\sigma}\right).$$

NB. There are many generalisations of the theorem where the assumptions of independence, common distribution or finite variance of the X_i are relaxed. More general proofs typically make use of the characteristic function mentioned earlier.

Central Limit Theorem

Example 3.12

Let X_1, \dots, X_n be an i.i.d. sample of exponential variables, i.e. $X_i \sim \text{Exp}(\lambda)$. Find formulae to approximate the probabilities $P(\bar{X}_n \leq x)$ and $P(\sum X_i \leq x)$.

Central Limit Theorem

Example 3.13

Normal approximation to the binomial

Consider a sequence of independent Bernoulli trials with constant probability of success π , so that $P(X_i = 1) = \pi$, $i = 1, 2, \dots$. Use the Central Limit Theorem with $\mu = \pi$ and $\sigma^2 = \pi(1 - \pi)$ to derive the normal approximation to the binomial distribution and perform a continuity correction.

Sketch Proof of CLT 1

- We additionally assume X_i have MGF defined on an open interval containing zero.
- To start, suppose $\mathbb{E}(X) = \mu = 0$. This implies $\mathbb{E}(\bar{X}_n) = 0$ and $\text{var}(\bar{X}_n) = \sigma^2/n$.
- Use standardised variable $Z_n = \sqrt{n} \cdot \bar{X}_n/\sigma$, $\mathbb{E}[Z_n] = 0$ and $\text{Var}(Z_n) = 1$.

Sketch Proof of CLT 1

- Use standardised variable $Z_n = \sqrt{n} \cdot \bar{X}_n / \sigma$, $\mathbb{E}[Z_n] = 0$ and $\text{Var}(Z_n) = 1$.

Note that Z_n is a linear combination of the X_i , since

$$Z_n = \frac{\sqrt{n} \sum_{i=1}^n X_i}{\sigma n} = \sum_{i=1}^n \frac{X_i}{\sigma \sqrt{n}}.$$

Thus, taking coefficients $a_i = 1/(\sigma \sqrt{n})$, the mgf of Z_n is

$$M_Z(s) = \prod_{i=1}^n M_{X_i} \left(\frac{s}{\sigma \sqrt{n}} \right) = \left\{ M_X \left(\frac{s}{\sigma \sqrt{n}} \right) \right\}^n$$

since the X_i all have the same distribution and thus the same mgf M_X .

Sketch Proof of CLT 2

Taylor series expansion of M_X about 0:

$$M_X(t) = M_X(0) + tM'_X(0) + \frac{t^2}{2}M''_X(0) + \varepsilon_t, \quad (1)$$

where ε_t is a term for which we know that $\varepsilon_t/t^2 \rightarrow 0$ for $t \rightarrow 0$, denoted $\varepsilon_t = o(t^2)$.

To make use of the above Taylor series expansion, note that (recalling that $\mu = 0$)

$$\begin{aligned} M_X(0) &= \mathbb{E}(e^{0X}) = 1 \\ M'_X(0) &= E(X_i) = 0 \\ M''_X(0) &= E(X_i^2) = \sigma^2. \end{aligned}$$

Inserting these values in (1) and replacing t by $s/(\sigma\sqrt{n})$ we find that...

Sketch Proof of CLT 3

Inserting $M_X(0)$, $M'_X(0)$, $M''_X(0)$ in (1) and replacing t by $s/(\sigma\sqrt{n})$ we find that

$$M_X(t) = M_X(0) + tM'_X(0) + \frac{t^2}{2}M''_X(0) + \varepsilon_t \quad (1)$$

turns into

$$M_X\left(\frac{s}{\sigma\sqrt{n}}\right) = 1 + 0 + \frac{1}{2}\sigma^2\left(\frac{s}{\sigma\sqrt{n}}\right)^2 + o\left(\frac{1}{n}\right) = 1 + \frac{s^2}{2n} + o\left(\frac{1}{n}\right).$$

Sketch Proof of CLT 4

We know the MGF M_Z of Z_n is given as the n th power of the MGF on the last slide and we find that

$$\begin{aligned} M_Z(s) &= \left\{ M_X \left(\frac{s}{\sigma\sqrt{n}} \right) \right\}^n = \left\{ 1 + \frac{s^2}{2n} + o\left(\frac{1}{n}\right) \right\}^n \\ &= \left(1 + \frac{\frac{1}{2}s^2 + \delta_n}{n} \right)^n \longrightarrow e^{\frac{1}{2}s^2} \end{aligned}$$

as $n \rightarrow \infty$, since $\delta_n \rightarrow 0$. (Recall: $(1 + \frac{a_n}{n})^n \rightarrow e^a$ as $n \rightarrow \infty$ provided $\lim_{n \rightarrow \infty} a_n = a$)

Sketch Proof of CLT 4

We know the MGF M_Z of Z_n is given as the n th power of the MGF on the last slide and we find that

$$\begin{aligned} M_Z(s) &= \left\{ M_X \left(\frac{s}{\sigma\sqrt{n}} \right) \right\}^n = \left\{ 1 + \frac{s^2}{2n} + o \left(\frac{1}{n} \right) \right\}^n \\ &= \left(1 + \frac{\frac{1}{2}s^2 + \delta_n}{n} \right)^n \longrightarrow e^{\frac{1}{2}s^2} \end{aligned}$$

as $n \rightarrow \infty$, since $\delta_n \rightarrow 0$. (Recall: $(1 + \frac{a_n}{n})^n \rightarrow e^a$ as $n \rightarrow \infty$ provided $\lim_{n \rightarrow \infty} a_n = a$)

Punch-Line

The limiting mgf is the one that we know as belonging to the standard normal distribution (*cf.* Example 3.4).

Sketch Proof of CLT: Finish

The above proves the claim in the case $\mu = 0$.

In the case that $\mathbb{E}(X_i) = \mu \neq 0$ define $Y_i = X_i - \mu$. Then $\mathbb{E}(Y_i) = 0$ so the result already proved gives

$Z_n = \sqrt{n}(\bar{X} - \mu)/\sigma = \sqrt{n}\bar{Y}/\sigma \xrightarrow{d} N(0, 1)$ as required. \square

Sums of Uniforms

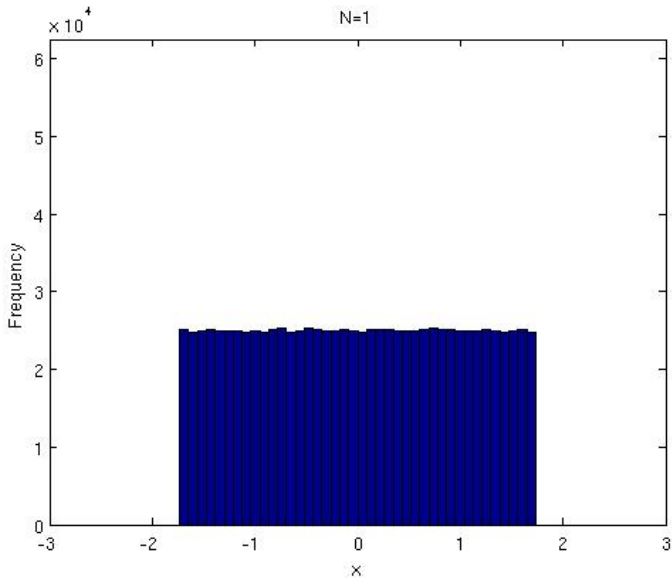
As an illustration of the CLT, consider the normalised sum of uniforms:

$$\sqrt{\frac{12}{N}} \sum_{i=1}^N U_i$$

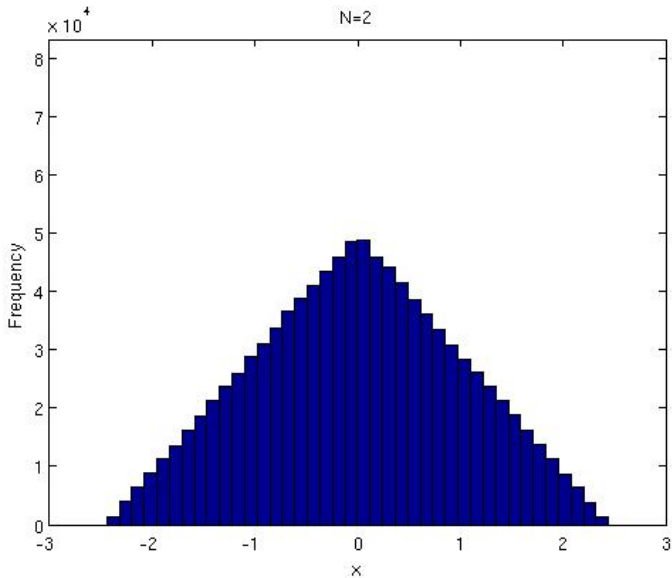
where $U_i \sim \text{Unif}(-0.5, 0.5)$.

- **For $N = 1$ the distribution is uniform (obvious)**
- **For $N = 2$ the distribution is triangle (do the convolution or derive it via mgfs)**
- **For $N > 2$ the calculation becomes lengthy**
- **As $N \rightarrow \infty$, the distribution approaches $\mathcal{N}(0, 1)$.**

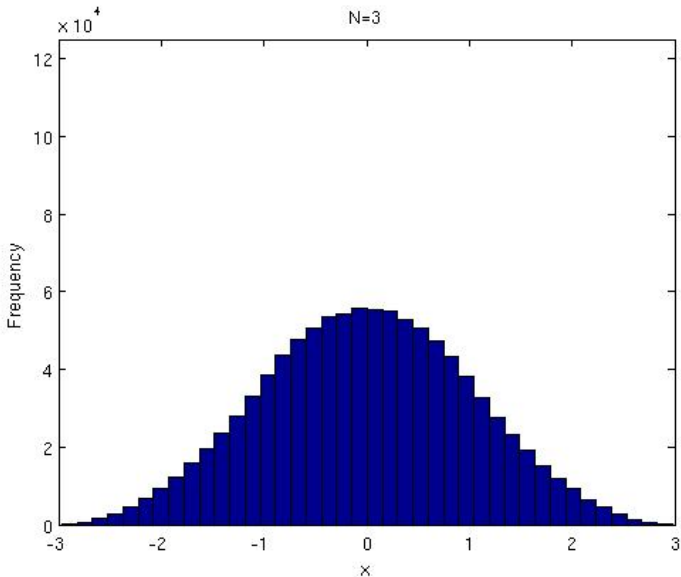
Sum of Uniforms: N=1



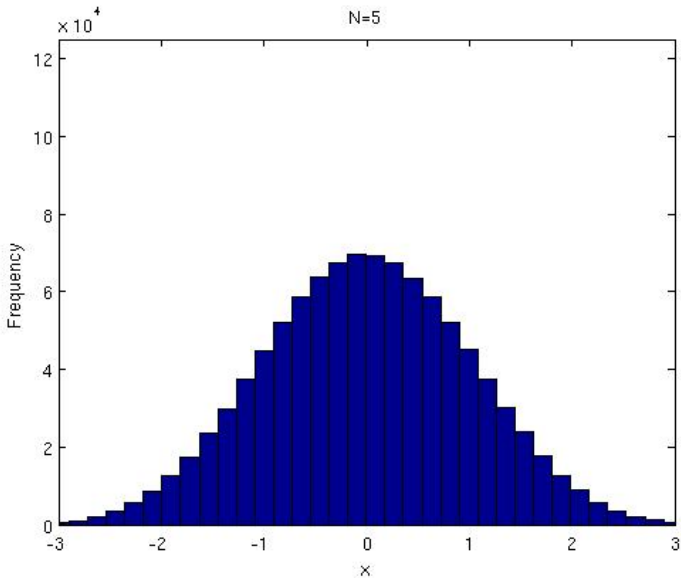
Sum of Uniforms: N=2



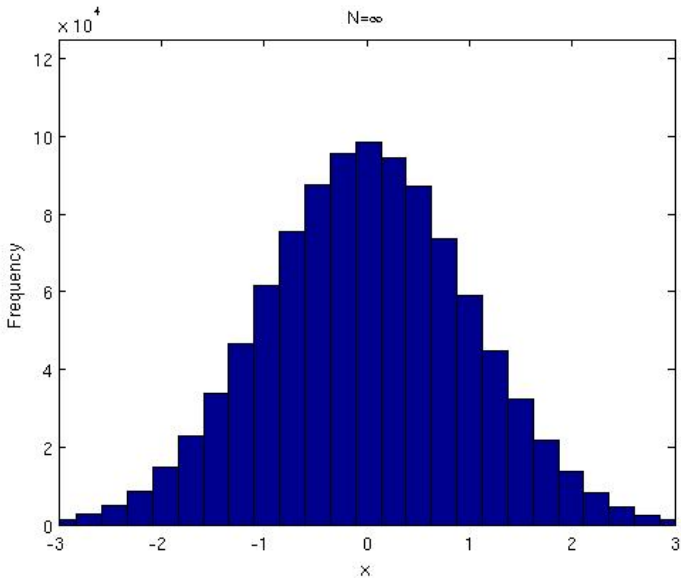
Sum of Uniforms: N=3



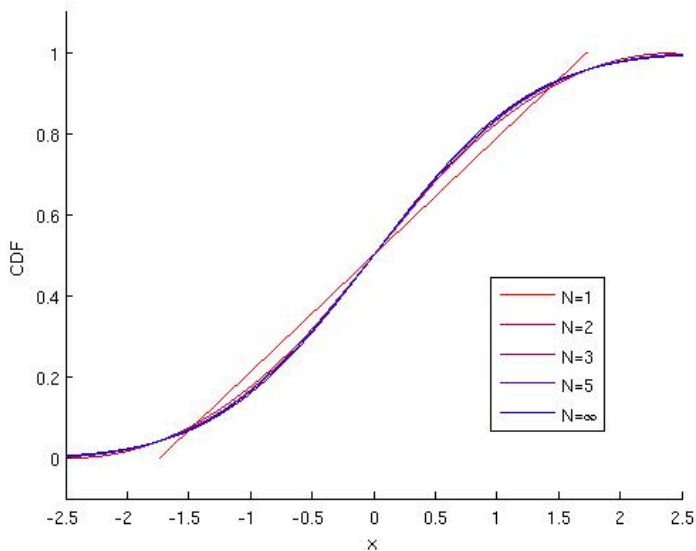
Sum of Uniforms: N=5



Sum of Uniforms: $N = \infty$



Sum of Uniforms: the sequence of CDFs



Central Limit Theorem

Example 3.14

Sums of i.i.d. uniformly distributed random variables

Let $X_i \sim U(0, 1)$ be independent random variables.

Compute the approximate distribution of $Y = \sum_{n=1}^{100} X_n$ and the probability that $Y < 60$.



4 Distributions of Functions of Normally Distributed Variables

These distributional results are needed for applications in statistics (such as estimation and hypothesis testing) in the context of samples of normally distributed variables. The Central Limit Theorem implies that the normal distribution arises in, or is a good approximation to, many practical situations.

4.1 The chi-squared (χ^2) distribution

Preliminaries. Recall the density of the $\text{Gam}(\alpha, \lambda)$ distribution,

$$f_X(x) = \frac{\lambda^\alpha x^{\alpha-1} e^{-\lambda x}}{\Gamma(\alpha)}, \quad x > 0$$

its mean and variance $E(X) = \alpha/\lambda$, $\text{Var}(X) = \alpha/\lambda^2$, its mgf $\{\lambda/(\lambda - s)\}^\alpha$ (see example 3.3) and the additive property: if X_1, \dots, X_n are independent $\text{Gam}(\alpha, \lambda)$ random variables then $X_1 + \dots + X_n$ is $\text{Gam}(n\alpha, \lambda)$.

We start by showing that if X has the standard normal distribution then X^2 has the gamma distribution with index $1/2$ and scale parameter $1/2$.

The moment generating function of X^2 is given by

$$\begin{aligned} E(e^{sX^2}) &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2}x^2 + sx^2 \right\} dx \\ &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2}x^2(1 - 2s) \right\} dx \end{aligned}$$

The moment generating function of X^2 is given by

$$\begin{aligned} E(e^{sX^2}) &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}x^2 + sx^2\right\} dx \\ &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}x^2(1 - 2s)\right\} dx \\ &= (1 - 2s)^{-1/2}, \end{aligned}$$

where the final line follows by comparison with the integral of the density of a normal variable with mean 0 and variance $(1 - 2s)^{-1}$.

By comparison with the gamma mgf, we see that X^2 has a gamma distribution with index and scale parameter each having the value 1/2. This distribution is also known as the **χ^2 distribution with one degree of freedom**. (Thus $\chi_1^2 \equiv \text{Gam}(\frac{1}{2}, \frac{1}{2})$.)

Example 4.1

Verify the result that $X^2 \sim \text{Gam}(1/2, 1/2)$ by using the transformation $\phi(x) = x^2$ on a mean zero normal random variable $X \sim N(0, 1)$.

Now let X_1, \dots, X_ν be independent standard normal variables, where ν is a positive integer. Then it follows from the additive property of the gamma distribution that their **sum of squares** $X_1^2 + \dots + X_\nu^2$ has the gamma distribution $\text{Gam}(\frac{\nu}{2}, \frac{1}{2})$ with index $\nu/2$ and scale parameter $1/2$.

This distribution is also known as the **χ^2 distribution with ν degrees of freedom** and is written as χ_ν^2 . Its mgf is therefore $\{ \frac{1}{2} / (\frac{1}{2} - s) \}^{\nu/2} = (1 - 2s)^{-\nu/2}$. (Thus $\chi_\nu^2 \equiv \text{Gam}(\frac{\nu}{2}, \frac{1}{2})$.)

It further follows from the mean and variance of the gamma distribution that the mean and variance of $U \sim \chi_\nu^2$ are (verify)

$$\mathbb{E}(U) = \nu, \quad \text{Var}(U) = 2\nu$$

Example 4.2

Verify the above expectation and variance of χ_ν^2 by considering $U = \sum_{i=1}^{\nu} X_i^2$ with independent $X_i \sim N(0, 1)$.

The χ^2 distribution

pdf

The pdf of the χ^2 distribution with ν degrees of freedom ($\nu > 0$) is

$$f(u) = \frac{u^{\frac{1}{2}\nu-1} e^{-\frac{1}{2}u}}{2^{\frac{1}{2}\nu} \Gamma(\frac{1}{2}\nu)}$$

for $u > 0$.

This can be verified by comparison with the $\text{Gam}(\frac{\nu}{2}, \frac{1}{2})$ density (see Appendix 2).

Application to sampling distributions

We know that if X has a normal distribution with mean μ and variance σ^2 then the standardised variable $(X - \mu)/\sigma$ has the standard normal distribution. It follows that if X_1, \dots, X_n are independent normal variables, all with mean μ and variance σ^2 , then

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

because the left-hand side is the sum of squares of n independent standard normal random variables.

Since $\sigma^2 = E(X - \mu)^2$, when μ is known the sample average $S_\mu^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2$ is an intuitively natural estimator of σ^2 . The above result gives the sampling distribution of S_μ^2 as $nS_\mu^2 \sim \sigma^2 \chi_n^2$.

Mean and Variance of χ^2

We can now deduce the mean and variance of S_μ^2 from those of the χ_n^2 distribution:

$$\mathbb{E}(S_\mu^2) = \frac{\sigma^2}{n} \mathbb{E}\left(\frac{nS_\mu^2}{\sigma^2}\right) = \frac{\sigma^2}{n} \times n = \sigma^2$$

$$\text{Var}(S_\mu^2) = \frac{\sigma^4}{n^2} \text{Var}\left(\frac{nS_\mu^2}{\sigma^2}\right) = \frac{\sigma^4}{n^2} \times 2n = \frac{2\sigma^4}{n}$$

Note that the expectation formula is generally true, even if the X_i are not normal. However, the variance formula is only true for normal distributions.

When μ is unknown, it is natural to estimate it by \bar{X} . We already know that $E(\bar{X}) = \mu$ and $\text{var}(\bar{X}) = \frac{\sigma^2}{n}$. Since a linear combination of normal variables is again normally distributed, we deduce that

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

However, in order to be able to use this result for statistical inference when σ^2 is also unknown, we need to estimate σ^2 . An intuitively natural estimator of σ^2 when μ is unknown is the **sample variance** $S^2 = \frac{1}{n-1} \sum (X_i - \bar{X})^2$. The reason for the factor $\frac{1}{n-1}$ rather than $\frac{1}{n}$ is that S^2 is unbiased for σ^2 , as we will see in the next section. In order to deduce the sampling distribution of S^2 we need to find the distribution of the sum of squares $\sum (X_i - \bar{X})^2$.

4.3.1 The distribution of $\sum (X_i - \bar{X})^2$

Let X_1, \dots, X_n be independent, normally distributed variables with common mean μ and variance σ^2 . We will see that $\sum_{i=1}^n (X_i - \bar{X})^2 / \sigma^2$ has the χ^2 distribution with $n - 1$ degrees of freedom.

4.3.1 The distribution of $\sum (X_i - \bar{X})^2$

Let X_1, \dots, X_n be independent, normally distributed variables with common mean μ and variance σ^2 . We will see that $\sum_{i=1}^n (X_i - \bar{X})^2 / \sigma^2$ has the χ^2 distribution with $n - 1$ degrees of freedom.

Result 4.1:

\bar{X} and $S^2 = \sum (X_i - \bar{X})^2 / (n - 1)$ are independent.

Result 4.2:

The sampling distribution of S^2 is $\frac{(n-1)}{\sigma^2} S^2 \sim \chi_{n-1}^2$.

Derivation of Result 4.1

To find the joint mgf, we relate a linear combination of \bar{X} and the U_i to a linear combination of the original X_i as follows:

$$s_0 \bar{X} + \sum_{i=1}^n s_i U_i = \sum_{i=1}^n a_i X_i,$$

where $a_i = \frac{s_0}{n} + s_i - \bar{s}$ and $\bar{s} = \sum_{i=1}^n s_i / n$ (exercise: derive this).

Derivation of Result 4.1

To find the joint mgf, we relate a linear combination of \bar{X} and the U_i to a linear combination of the original X_i as follows:

$$s_0 \bar{X} + \sum_{i=1}^n s_i U_i = \sum_{i=1}^n a_i X_i,$$

where $a_i = \frac{s_0}{n} + s_i - \bar{s}$ and $\bar{s} = \sum_{i=1}^n s_i / n$ (exercise: derive this).

The joint mgf of \bar{X} and the U_i is now given by

$$M_{\bar{X}, U_1, \dots, U_n}(s_0, s_1, \dots, s_n) = \mathbb{E}(e^{s_0 \bar{X} + \sum s_i U_i}) = \mathbb{E}(e^{\sum a_i X_i}),$$

which is the joint mgf of X_1, \dots, X_n evaluated at the above choice of the a_i .

Derivation of Result 4.1

As X_1, \dots, X_n is an independent normal sample we know that this joint mgf (*c.f.* Example 3.6) is given by

$$\mathbb{E}(e^{\sum a_i X_i}) = \prod_i M_{X_i}(a_i) = \exp \left\{ \left(\sum a_i \right) \mu + \frac{1}{2} \left(\sum a_i^2 \right) \sigma^2 \right\}$$

But $\sum_{i=1}^n a_i = s_0$ and $\sum_{i=1}^n a_i^2 = \frac{s_0^2}{n} + \sum_{i=1}^n (s_i - \bar{s})^2$, giving

$$M_{\bar{X}, U_1, \dots, U_n}(s_0, s_1, \dots, s_n) = \exp \left\{ s_0 \mu + \frac{1}{2n} s_0^2 \sigma^2 + \frac{1}{2} \sigma^2 \sum_{i=1}^n (s_i - \bar{s})^2 \right\}$$

Derivation of Result 4.1

Therefore the joint mgf $M_{\bar{X}, U_1, \dots, U_n}$ *factorises* into

$$\exp \left\{ s_0 \mu + \frac{1}{2n} s_0^2 \sigma^2 \right\} \quad \text{and} \quad \exp \left\{ \frac{1}{2} \sigma^2 \sum_{i=1}^n (s_i - \bar{s})^2 \right\}.$$

This factorisation implies that

- (1) \bar{X} and U_1, \dots, U_n are independent;
- (2) the marginal mgf of \bar{X} is $\exp(s_0 \mu + \frac{1}{2} s_0^2 \sigma^2 / n)$
 $\implies \bar{X} \sim N(\mu, \sigma^2 / n)$ (as we already know).

Example 4.3

Show that \bar{X} and $U_i = X_i - \bar{X}$ are independent by computing their covariance.

4.3.1 The distribution of $\sum (X_i - \bar{X})^2$

Let X_1, \dots, X_n be independent, normally distributed variables with common mean μ and variance σ^2 . We will see that $\sum_{i=1}^n (X_i - \bar{X})^2 / \sigma^2$ has the χ^2 distribution with $n - 1$ degrees of freedom.

Result 4.1:

\bar{X} and $S^2 = \sum (X_i - \bar{X})^2 / (n - 1)$ are independent.

Result 4.2:

The sampling distribution of S^2 is $\frac{(n-1)}{\sigma^2} S^2 \sim \chi_{n-1}^2$.

Derivation of Result 4.2

$$\sum (X_i - \mu)^2 = \sum (X_i - \bar{X})^2 + n(\bar{X} - \mu)^2$$

so that

$$\frac{\sum (X_i - \mu)^2}{\sigma^2} = \frac{\sum (X_i - \bar{X})^2}{\sigma^2} + \frac{(\bar{X} - \mu)^2}{\sigma^2/n}.$$

But, since $(X_i - \mu)/\sigma$ are independent $N(0, 1)$ and $(\bar{X} - \mu)/(\sigma/\sqrt{n})$ is $N(0, 1)$, we have

$$\frac{\sum (X_i - \mu)^2}{\sigma^2} \sim \chi_n^2 \quad \text{and} \quad \frac{(\bar{X} - \mu)^2}{\sigma^2/n} \sim \chi_1^2$$

Derivation of Result 4.2

$$\sum (X_i - \mu)^2 = \sum (X_i - \bar{X})^2 + n(\bar{X} - \mu)^2$$

so that

$$\frac{\sum (X_i - \mu)^2}{\sigma^2} = \frac{\sum (X_i - \bar{X})^2}{\sigma^2} + \frac{(\bar{X} - \mu)^2}{\sigma^2/n}.$$

But, since $(X_i - \mu)/\sigma$ are independent $N(0, 1)$ and $(\bar{X} - \mu)/(\sigma/\sqrt{n})$ is $N(0, 1)$, we have

$$\frac{\sum (X_i - \mu)^2}{\sigma^2} \sim \chi_n^2 \quad \text{and} \quad \frac{(\bar{X} - \mu)^2}{\sigma^2/n} \sim \chi_1^2$$

Using mgfs, it follows that

$$\frac{\sum (X_i - \bar{X})^2}{\sigma^2} \sim \chi_{n-1}^2$$

since $(\bar{X} - \mu)^2/(\sigma^2/n)$ and $\sum (X_i - \bar{X})^2/\sigma^2$ are indep.

These results give the sampling distribution of S^2 to be $(n-1)S^2 \sim \sigma^2 \chi_{n-1}^2$. In particular, we can deduce the sampling mean and variance of S^2 to be

$$\mathbb{E}(S^2) = \frac{\sigma^2}{n-1} \mathbb{E}\left(\frac{(n-1)S^2}{\sigma^2}\right) = \frac{\sigma^2}{n-1} \times n-1 = \sigma^2$$

$$\text{Var}(S^2) = \frac{\sigma^4}{(n-1)^2} \text{Var}\left(\frac{(n-1)S^2}{\sigma^2}\right) = \frac{2\sigma^4}{n-1}$$

recalling that the mean and variance of χ_{n-1}^2 are $n-1$, $2(n-1)$ respectively.

As for the case where μ is known, the **unbiasedness** property of S^2 is generally true, even if the X_i are not normal, but the variance formula is only true for normal distributions.

4.3.2 Student's t distribution

If X_1, \dots, X_n are independent normally distributed variables with common mean μ and variance σ^2 then \bar{X} is normally distributed with mean μ and variance σ^2/n . Recall (STAT1004 & STAT1005, or MATH7502) that if σ^2 is known, then we may test the hypothesis $\mu = \mu_0$ by examining the statistic

$$Z = \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}}.$$

Z is a linear transformation of a normal variable and hence is also normally distributed. When $\mu = \mu_0$ we see that $Z \sim N(0, 1)$ so we conduct the test by computing Z and referring to the $N(0, 1)$ distribution.

Now if σ^2 is unknown, it is intuitively reasonable to estimate it by S^2 (recalling that $\mathbb{E}(S^2) = \sigma^2$) and use the statistic

$$T = \frac{\bar{X} - \mu_0}{S/\sqrt{n}}.$$

However, in order to conduct the test we need to know the distribution of this statistic when $\mu = \mu_0$.

Note that we can write T as

$$T = \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}} \frac{\sigma}{S} = \frac{Z}{\sqrt{U/(n-1)}},$$

where $U = \sum_i (X_i - \bar{X})^2 / \sigma^2$. From the above results we have $Z \sim N(0, 1)$, $U \sim \chi_{n-1}^2$ and Z, U are independent random variables. (Note that the distribution of T does not depend on μ_0 and σ^2 , but only on the known number n of observations and is therefore suitable as a test statistic.)

We can now find the probability distribution of T by the transformation method that was described in §2.2.

Alternatively, we can derive it from the F distribution as shown below. The distribution of T , denoted by t_{n-1} , is known as **Student's t distribution with $n - 1$ degrees of freedom**.

The general description of the t distribution is as follows. Suppose that Z has a standard normal distribution, U has a χ^2 distribution with ν degrees of freedom, and Z, U are independent random variables. Then

$$T = \frac{Z}{\sqrt{U/\nu}}$$

has the Student's t distribution with ν degrees of freedom, denoted by t_ν .

T has probability density function

$$f_T(t) = \frac{\Gamma((\nu + 1)/2)}{\Gamma(\nu/2)} \frac{1}{\sqrt{\nu\pi}} \left(1 + \frac{t^2}{\nu}\right)^{-(\nu+1)/2}$$

for $-\infty < t < \infty$.

The distribution of T is symmetrical about 0, so that $E(T) = 0$ ($\nu > 1$). It can further be shown that the variance of T is $\nu/(\nu - 2)$ for $\nu > 2$.

Example 4.4

Write down the pdf of a t_1 -distribution and identify the distribution by name (look back at Example 2.2). Why is $\nu > 1$ needed for the expected value to be zero?

STAT3101 only: What happens as the sample size ν becomes large? Derive the limit of the pdf in this case.

Example 2.2

Consider $X \sim \text{Uniform}[-\frac{\pi}{2}, \frac{\pi}{2}]$, i.e.

$$f_X(x) = \begin{cases} \frac{1}{\pi} & -\frac{\pi}{2} \leq x \leq \frac{\pi}{2} \\ 0 & \text{otherwise.} \end{cases}$$

Derive the density of $Y = \tan(X)$.



4.4 The F distribution

Now suppose that we have two *independent* samples of observations: X_1, \dots, X_m are independent normally distributed variables with common mean μ_X and variance σ_X^2 , while Y_1, \dots, Y_n are independent and normally distributed with common mean μ_Y and variance σ_Y^2 . Suppose that we wish to test the hypothesis that the variances σ_X^2 and σ_Y^2 are equal. If $\sigma_X^2 = \sigma_Y^2$ then $\sigma_X^2/\sigma_Y^2 = 1$ and it is natural to examine the ratio of the two sample variances, S_X^2/S_Y^2 , and compare its value with 1.

Here

$$S_X^2 = \sum_{i=1}^m (X_i - \bar{X})^2 / (m - 1), \quad S_Y^2 = \sum_{i=1}^n (Y_i - \bar{Y})^2 / (n - 1).$$

Now, since $(m - 1)S_X^2/\sigma_X^2 \sim \chi_{m-1}^2$, and $(n - 1)S_Y^2/\sigma_Y^2 \sim \chi_{n-1}^2$, we can write

$$\frac{S_X^2/\sigma_X^2}{S_Y^2/\sigma_Y^2} = \frac{U/(m - 1)}{V/(n - 1)}$$

where $U = \sum_i (X_i - \bar{X})^2/\sigma_X^2 \sim \chi_{m-1}^2$,
 $V = \sum_i (Y_i - \bar{Y})^2/\sigma_Y^2 \sim \chi_{n-1}^2$ and U and V are independent variables. When $\sigma_X^2 = \sigma_Y^2$, the left hand side is just the ratio S_X^2/S_Y^2 .

The above considerations motivate the following general description of the F distribution. Suppose that $U \sim \chi_{\alpha}^2$, $V \sim \chi_{\beta}^2$ and U, V are independent. Then

$$W = \frac{U/\alpha}{V/\beta}$$

has the F distribution with (α, β) degrees of freedom, denoted $F_{\alpha, \beta}$.

The density of $W \sim F_{a,b}$ is most easily written using the **Beta function** $B(a, b) \equiv \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$ as

$$f_W(w) = \frac{\frac{\alpha}{\beta} \left(\frac{\alpha w}{\beta} \right)^{\frac{\alpha}{2}-1}}{B\left(\frac{\alpha}{2}, \frac{\beta}{2}\right) \left(1 + \frac{\alpha w}{\beta} \right)^{\frac{\alpha+\beta}{2}}}$$

for $w \geq 0$.

Note that, from the definition of the F distribution, we have

$$\frac{1}{W} = \frac{V/\beta}{U/\alpha} \sim F_{\beta,\alpha}$$

i.e. if $W \sim F_{\alpha,\beta}$ then $1/W \sim F_{\beta,\alpha}$.

It can be shown that, for all α and for $\beta > 2$,

$$E(W) = \beta/(\beta - 2).$$

We can now apply this result to our test statistic discussed earlier. Under the hypothesis that $\sigma_X^2 = \sigma_Y^2$ we have (taking $\alpha = m - 1$, $\beta = n - 1$)

$$\boxed{\frac{S_X^2}{S_Y^2} \sim F_{m-1,n-1}}$$

This is the sampling distribution of the variance ratio. Note that this distribution is free from σ_X and σ_Y .

Example 4.5

Suppose that X , Y , U are independent random variables such that $X \sim N(2, 9)$, $Y \sim t_4$ and $U \sim \chi_3^2$. Give four functions of the above variables that have the following distributions:

- 1 χ_1^2
- 2 χ_4^2
- 3 t_3
- 4 $F_{1,4}$