

Introduction to statistical inference 2

Lauri Mehtätalo
University of Eastern Finland
School of Computing

April 23, 2018

Contents

1	Recap from “Introduction to statistical inference 1”	1
1.1	Random variable	1
1.2	Transformations of random variable	2
1.3	Expected values	3
1.4	Variance	3
1.5	Bivariate random variables	3
1.6	Independence	4
1.7	Covariance	4
1.8	Random vectors	5
1.9	Computing using expected values and variances	7
2	Random samples	9
2.1	Distribution of sum	14
2.2	Sampling from Normally distributed population	15

Chapter 1

Recap from “Introduction to statistical inference 1”

1.1 Random variable

- Random variable is a function from sample space \mathcal{S} of an experiment to sample space of the random variable \mathcal{X} , which is set of real numbers.

$$X : \mathcal{S} \rightarrow \mathcal{X}$$

- The sample space is a set of all possible values random variable can get.
- \mathcal{X} can be
 - an interval of real axis (continuous random variable).

$$\mathcal{X} = [0, 10), \mathcal{X} = [0, 10], \mathcal{X} = (0, 10)$$

- An uncountable set of integers (discrete random variable)

$$\mathcal{X} = \{0, 1, 2, \dots\}$$

- A countable set of integers or real numbers (discrete random variable)

$$\mathcal{X} = \{0, 1\}, \mathcal{X} = \{0, 0.5, 1\}, \mathcal{X} = \{0, 1, \dots, 10\}$$

- Probabilities associated with each value of X are defined by the cumulative distribution function (cdf for short).

$$F_X(x) = P(X \leq x), \text{ where } -\infty < x < \infty$$

Note: $F_X(x)$ is a step function if X is discrete.

Note: $F_X(x)$ is a continuous function if X is continuous.

- $F_X(x)$ or $F(x)$ is cdf, if

- 1) $\lim_{x \rightarrow -\infty} F(x) = 0$ and $\lim_{x \rightarrow \infty} F(x) = 1$

- 2) $F(x)$ is non-decreasing

- 3) $F(x)$ is right-continuous

- Cdf is useful in calculation of any probabilities; for example

$$P(a < X \leq b) = F(b) - F(a)$$

Note: Be careful with $<$ and \leq when working with discrete random variables.

- The probability density function (pdf for short) is defined for continuous random variable as

$$f_X(x) = F'(x) = \frac{dF_X(x)}{dx}, \quad -\infty < x < \infty$$

and

$$\int_{-\infty}^x f_X(t) dt = F_X(x)$$

- The probability mass function (pmf for short) is defined for discrete random variables as

$$f_X(x) = P(X = x)$$

$$F_X(x) = \sum_{k=1}^x f_X(k)$$

1.2 Transformations of random variable

- Consider a monotonic function $g : \mathcal{X} \rightarrow \mathcal{Y}$
- $Y = g(X)$ is also a random variable; function g is called an transformation (muunnos).

- If $g(x)$ is a increasing function of x , then

$$F_Y(y) = F_X(g^{-1}(y))$$

- If $g(x)$ is a decreasing function of x , then

$$F_Y(y) = 1 - F_X(g^{-1}(y))$$

- The pdf of continuous Y is

$$f_Y(y) = F'_Y(y)$$

1.3 Expected values

$$E(X) = \mu_X = \begin{cases} \int_{-\infty}^{\infty} xf(x) dx & \text{if } X \text{ is continuous} \\ \sum_{x \in \mathcal{X}} xf(x) & \text{if } X \text{ is discrete} \end{cases}$$

$$E(g(X)) = \begin{cases} \int_{-\infty}^{\infty} g(x)f(x) dx & \text{if } X \text{ is continuous} \\ \sum_{x \in \mathcal{X}} g(x)f(x) & \text{if } X \text{ is discrete} \end{cases}$$

1.4 Variance

$$\begin{aligned} \sigma_X^2 = \text{Var}(X) &= E(X - \mu_X)^2 \\ &= E(X^2 - 2X\mu_X + \mu_X^2) \\ &= E(X^2) - E(2X\mu_X) + E(\mu_X^2) \\ &= E(X^2) - 2\mu_X \underbrace{E(X)}_{\mu_X} + E(\mu_X^2) \\ &= E(X^2) - \mu_X^2 \\ \text{sd}(X) &= \sqrt{\text{Var}(X)} = \sigma_X \end{aligned}$$

1.5 Bivariate random variables

- For two discrete random variables, the joint pmf is defined as

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

- For two continuous random variables, we define the joint pdf $f_{X,Y}(x, y)$ as

$$P((X, Y) \in A) = \iint_A f(x, y) dx dy$$

- The expected value for transformation $g(X, Y) : \mathbb{R}^2 \rightarrow \mathbb{R}$ (for example, $g(X, Y) = XY$ or $g(X, Y) = \frac{X}{Y}$) is

$$E(g(X, Y)) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x, y)f(x, y) dx dy$$

if (X, Y) is continuous, and

$$E(g(X, Y)) = \sum_{x, y \in \mathbb{R}^2} g(x, y)f(x, y)$$

if (X, Y) is discrete.

- The marginal pmf / pdf for X are

$$f_X(x) = \sum_{y \in \mathbb{R}} f_{X,Y}(x, y) \quad (\text{pmf})$$

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

and correspondingly for Y .

- The conditional pmf / pdf are both defined as

$$f(y | x) = \frac{f_{X,Y}(x, y)}{f_X(x)} \quad (\text{for both discrete and continuous random variables})$$

and correspondingly for $x | y$.

1.6 Independence

- Random variables are said to be independent ($X \perp\!\!\!\perp Y$) if

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

- For independent random variables, conditional distribution $y | x$ is

$$f(y | x) = f_y(y)$$

regardless of the value of x .

1.7 Covariance

- Covariance measures the linear association between two random variables X and Y

$$\text{cov}(X, Y) = E((X - \mu_X)(Y - \mu_Y))$$

$$= E(XY) - \mu_X\mu_Y$$

$$\text{cov}(X, Y) = \text{cov}(Y, X)$$

$$\text{corr}(X, Y) = \rho_{XY} = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y} \quad \text{where } -1 \leq \rho_{XY} \leq 1$$

Note: $\text{cov}(X, X) = \text{Var}(X)$

Note: If $X \perp\!\!\!\perp Y$, then $\text{cov}(X, Y) = 0$, but if $\text{cov}(X, Y) = 0$, it does not mean X and Y are necessarily independent.

1.8 Random vectors

- Random vectors generalize the bivariate random variables to a n -variate case.

$$\mathbf{X} = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{bmatrix} \quad \text{where } X_i, i = 1, 2, \dots, n \text{ are scalar random variables}$$

- If \mathbf{x} is a continuous random vector, then

$$P(\mathbf{X} \in A) = \int_A \cdots \int f(\mathbf{x}) dx_1 \dots dx_n \quad f : \mathbb{R}^2 \rightarrow \mathbb{R}$$

where $f(\mathbf{x})$ is a joint pdf.

- If \mathbf{x} is a discrete random vector, then

$$P(\mathbf{X} \in A) = \sum \dots \sum f(\mathbf{x})$$

where $f(\mathbf{x})$ is a joint pmf.

- Let $g(\mathbf{x})$ be a transformation $g : \mathbb{R}^2 \rightarrow \mathbb{R}$. Then, the expected value is

$$E(g(\mathbf{X})) = \begin{cases} \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} x f(\mathbf{x}) dx_1 \dots dx_n & \text{if } \mathbf{X} \text{ is continuous} \\ \sum \dots \sum_{\mathbf{X} \in \mathbb{R}^2} g(\mathbf{x}) f(\mathbf{x}) & \text{if } \mathbf{X} \text{ is discrete} \end{cases}$$

- Let us partition the n -variate random vector \mathbf{X} as follows:

$$\mathbf{X} = \begin{pmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{pmatrix}$$

where \mathbf{X}_1 has length k and \mathbf{X}_2 has length $n - k$.

- The joint pdf of \mathbf{X} can be written as

$$f(\mathbf{x}) = f(\mathbf{x}_1, \mathbf{x}_2)$$

- The marginal density of \mathbf{X}_1 is

$$f(\mathbf{x}_1) = \int_{\mathbb{R}^{n-k}} f(\mathbf{x}_1, \mathbf{x}_2) d\mathbf{x}_2 \quad \text{if } \mathbf{X} \text{ is continuous}$$

$$f(\mathbf{x}_1) = \sum_{\mathbb{R}^{n-k}} f(\mathbf{x}_1, \mathbf{x}_2) \quad \text{if } \mathbf{X} \text{ is discrete}$$

where $f(\mathbf{x}_1)$ is a k -variate pdf/pmf.

- The conditional pdf/pmf for $\mathbf{X}_2 \mid \mathbf{X}_1$ is

$$f(\mathbf{X}_2 \mid \mathbf{X}_1) = \frac{f(\mathbf{X}_1, \mathbf{X}_2)}{f(\mathbf{X}_1)}$$

- Expected value of random vectors is defined as vector

$$E(\mathbf{x}) = \begin{bmatrix} E(X_1) \\ E(X_2) \\ \vdots \\ E(X_n) \end{bmatrix}_{n \times 1} = \begin{bmatrix} E(\mathbf{X}_1) \\ E(\mathbf{X}_2) \end{bmatrix} = \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}$$

- The variance of a random vector is $n \times n$ symmetric matrix called variance-covariance matrix.

$$\begin{aligned} Var(\mathbf{x})_{n \times n} &= \begin{bmatrix} Var(X_1) & cov(X_1, X_2) & cov(X_1, X_3) & \dots & cov(X_1, X_n) \\ cov(X_2, X_1) & Var(X_2) & cov(X_2, X_3) & \dots & cov(X_2, X_n) \\ cov(X_3, X_1) & cov(X_3, X_2) & Var(X_3) & \dots & cov(X_3, X_n) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ cov(X_n, X_1) & cov(X_n, X_2) & cov(X_n, X_3) & \dots & Var(X_n) \end{bmatrix} \\ &= \begin{bmatrix} Var(\mathbf{x}_1)_{k \times k} & cov(\mathbf{x}_1, \mathbf{x}'_2)_{k \times (n-k)} \\ cov(\mathbf{x}_2, \mathbf{x}'_1)_{(n-k) \times k} & Var(\mathbf{x}_2)_{(n-k) \times (n-k)} \end{bmatrix} \end{aligned}$$

- The correlation matrix is defined as

$$corr(\mathbf{x})_{n \times n} = \begin{bmatrix} 1 & corr(X_1, X_2) & corr(X_1, X_3) & \dots & corr(X_1, X_n) \\ corr(X_2, X_1) & 1 & corr(X_2, X_3) & \dots & corr(X_2, X_n) \\ corr(X_3, X_1) & corr(X_3, X_2) & 1 & \dots & corr(X_3, X_n) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ corr(X_n, X_1) & corr(X_n, X_2) & corr(X_n, X_3) & \dots & 1 \end{bmatrix}$$

- If \mathbf{X} has a n -variate normal distribution, then, with

$$E(\mathbf{x}) = \boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix} \quad \text{and} \quad Var(\mathbf{x})_{n \times n} = \boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_1 & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_2 \end{bmatrix}$$

then $\mathbf{X}_1 \mid \mathbf{X}_2$ has a k -variate normal distribution with

$$E(\mathbf{X}_1 \mid \mathbf{X}_2) = \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_2^{-1}(\mathbf{X}_2 - \boldsymbol{\mu}_2)$$

$$\text{Var}(\mathbf{X}_1 \mid \mathbf{X}_2) = \boldsymbol{\Sigma}_1 - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_2^{-1}\boldsymbol{\Sigma}_{21}$$

Note: $\text{Var}(\mathbf{X}_1 \mid \mathbf{X}_2) \leq \text{Var}(\mathbf{X}_1)$

1.9 Computing using expected values and variances

Let a , b and c be constants, and let X , Y and Z be (scalar) random variables. The following rules hold regardless of the distribution of random variables X , Y and Z .

$$E(c) = c$$

$$E(cX) = cE(X)$$

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

$$E(X + c) = E(X) + c$$

$$E(XY) = E(X)E(Y)$$

Only if $X \perp\!\!\!\perp Y$.

$$E(g(X)) = g(E(X))$$

Only in some special cases, like when $g(X)$ is a linear transformation.

$$\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2\text{cov}(X, Y)$$

$$\text{Var}(X - Y) = \text{Var}(X) + \text{Var}(Y) - 2\text{cov}(X, Y)$$

$$\text{Var}(aX) = a^2 \cdot \text{Var}(X)$$

$$\text{cov}(X, Y + Z) = \text{cov}(X, Y) + \text{cov}(X, Z)$$

$$\text{cov}(aX, bY) = ab \cdot \text{cov}(XY)$$

$$\text{Var}(X + a) = \text{Var}(X)$$

$$\text{cov}(X + a, Y + b) = \text{cov}(X, Y)$$

$$E(X) = E_Y [E_{X|Y}(X \mid Y)]$$

$$\text{Var}(X) = E_Y [\text{Var}_{X|Y}(X \mid Y)] + \text{Var}_Y [\text{Var}_{X|Y}(X \mid Y)]$$

Let \mathbf{a} and \mathbf{b} be fixed vectors and \mathbf{X} and \mathbf{Y} random vectors so that the dimensions in the equations match.

$$E(\mathbf{a}'\mathbf{X}) = \mathbf{a}'E(\mathbf{X})$$

$$\text{Var}(\mathbf{a}'\mathbf{X}) = \mathbf{a}'\text{Var}(\mathbf{X})\mathbf{a} \quad \text{Compare to } \text{Var}(aX) = a^2 \cdot \text{Var}(X) = a \cdot \text{Var}(X) \cdot a$$

$$\text{cov}(\mathbf{a}'\mathbf{X}, \mathbf{b}'\mathbf{Y}) = \mathbf{a}'\text{cov}(\mathbf{X}, \mathbf{Y})\mathbf{b}$$

Note: These equations need to be remembered by heart!

Chapter 2

Random samples

Definition 2.1. Random variables X_1, \dots, X_n are called random sample of size n from population $f(x)$, if X_1, \dots, X_n are mutually independent random variables and the marginal pdf/pmf of each X_i is the same function $f(x)$. Alternatively, X_1, \dots, X_n are called independent and identically distributed random variables (i.i.d.) with pdf/pmf $f(x)$.

Note: Sample X_1, \dots, X_n can also be denoted by \mathbf{X} , where $\mathbf{X} = [X_1 \ \dots \ X_n]^T$

- If follows from the mutual independence, of X_1, \dots, X_n that the joint pdf or pmf of \mathbf{X} is

$$f(x_1, \dots, x_n) = f(x_1)f(x_2) \dots f(x_n) = \prod_{i=1}^n f(x_i)$$

- All univariate marginal distributions $f(x_i)$ are the same by definition 2.1.

Example 2.1. Let X_1, \dots, X_n be a random sample from *Exponential*(β) population. X_i specifies the time until failure for n identical cellphones. The exponential pdf is

$$f(x_i) = \frac{1}{\beta} e^{-x_i/\beta}$$

The joint pdf of the sample is

$$f(x_1, x_2, \dots, x_n \mid \beta) = \prod_{i=1}^n \frac{1}{\beta} e^{-x_i/\beta} = \frac{1}{\beta^n} e^{-\frac{1}{\beta} \sum x_i} \quad \text{Recall: } a^b a^c = a^{b+c}$$

What is the probability that all n cellphones last more than 2 years?

$$\begin{aligned}
P(X_1 > 2, \dots, X_n > 2) &= \int_2^\infty \dots \int_2^\infty f(x_1, \dots, x_n) dx_1 \dots dx_n \\
&= \int_2^\infty \dots \int_2^\infty \prod_{i=1}^n \frac{1}{\beta} e^{-x_i/\beta} dx_1 \dots dx_n \\
&= \int_2^\infty \dots \int_2^\infty \prod_{i=2}^n \frac{1}{\beta} e^{-x_i/\beta} \underbrace{\int_2^\infty \frac{1}{\beta} e^{-x_1/\beta} dx_1}_{\substack{\text{Integral over the exponential pdf} \\ \int_2^\infty f(x_1) dx_1 = 1 - F(2) \\ = 1 - (1 - e^{-\frac{1}{\beta^2}}) = e^{-\frac{2}{\beta^2}}} } dx_2 \dots dx_n \\
&= e^{-\frac{2}{\beta}} \underbrace{\int_2^\infty \dots \int_2^\infty \prod_{i=2}^n \frac{1}{\beta} e^{-x_i/\beta} dx_2 \dots dx_n}_{\substack{\text{Identical to original integral except that} \\ \text{there are only } n-1 \text{ terms to be integrated}}} \dots \\
&= e^{-2/\beta} \cdot e^{-2/\beta} \int_2^\infty \dots \int_2^\infty \prod_{i=3}^n \frac{1}{\beta} e^{-x_i/\beta} dx_3 \dots dx_n \\
&= (e^{-2/\beta})^n = e^{-2n/\beta}
\end{aligned}$$

A much more simpler solution: notice that $P(X_i > 2) = 1 - F(2) = e^{-1/\beta}$. Because X_i 's are independent, then also events $(X_i > 2)$ are independent event, and so

$$P(\text{All } X_i > 2) = \prod_{i=1}^n P(X_i > 2) = (e^{-2/\beta})^n = e^{-2n/\beta}$$

Illustration See `ExponentialSample.R` for the case where $n = 2$ and $\beta = 3$. Notice that R uses parametrization $\lambda = \frac{1}{\beta}$ for the exponential distribution, so in R, exponential pdf is $f(x | \lambda) = \lambda e^{-\lambda x}$.

Note: Definition 2.1 assumes that X_1, \dots, X_n are independent. In practice, we often do analysis on dependent data. For that use, we could define term “dependent random sample”. Mathematical treatment of dependent sample requires more specific description of dependence structure, e.g. through spatial or temporal autocorrelation models, or explicit models for grouped data.

Dependent data are modeled by assuming that random vectors $\mathbf{X}_1, \dots, \mathbf{X}_n$ are vectors of appropriate length, and there are n independent realizations of them in the data. Quite often $n = 1$ and each replicate to \mathbf{X}_1 , which is a rather long vector.

Example 2.2. Let $\mathbf{X} = [X(u_1) \ X(u_2) \ X(u_3)]^T$ include random variables X at locations u_1, u_2 and u_3 . Assume that \mathbf{X}_1 is normally distributed and the correlation $\rho_{ij} = \text{corr}(X(u_i), X(u_j))$ depends only on the spatial distance $s_{ij} = \|u_i - u_j\|$ between locations u_i and u_j .

The marginal means and variances of $X(u_i)$ are $E(X(u_i)) = \mu$ for all i and $\text{Var}(X(u_i)) = \sigma^2$ for all i . Consider case where $u_1 = (1, 0)$, $u_2 = (0, 0)$, $u_3 = (0, 2)$ and $\rho_{ij} = e^{-s_{ij}/2}$. Correlations and covariances related to different random variables can be seen on table 2.1.

Table 2.1: Correlations between the random variable pairs of example 2.2

Pair	s_{ij}	ρ_{ij}	$cov(X(u_i), X(u_j))$
1,2	1	0.61	$0.61\sigma^2$
1,3	$\sqrt{5}$	0.33	$0.33\sigma^2$
2,3	2	0.37	$0.37\sigma^2$

Let

$$\boldsymbol{\mu} = \begin{bmatrix} \mu \\ \mu \\ \mu \end{bmatrix} \quad \text{and} \quad \boldsymbol{\Sigma} = \sigma^2 \begin{bmatrix} 1 & 0.61 & 0.33 \\ 0.61 & 1 & 0.37 \\ 0.33 & 0.37 & 1 \end{bmatrix}$$

The joint pdf of \mathbf{X} is 3 variate normal distribution with mean $\boldsymbol{\mu}$ and variance

$$f(x_1 \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^3}} |\boldsymbol{\Sigma}|^{-1/2} e^{\frac{1}{2}(\mathbf{x}_1 - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\mathbf{x}_1 - \boldsymbol{\mu})}$$

Note: In random sample, we assume that X_1, \dots, X_n are identically distributed. In the case of random vectors, this is specified by saying that the marginal distributions of the elements of \mathbf{X}_i are identical.

Note: We can also have independent replicates of random vectors, e.g. if we have grouped data where the groups are independent replicates from the process that generates the groups, and the observations within the groups are dependent. This is related to mixed-effects models.

Note: Definition 2.1 specifies sampling from an infinite population (or population model): the sampling procedure does not change the population.

The population may also be finite, like numbers in hat. In that case, the sampling can be made with replacement: whenever a number has been drawn, the value is recorded but the number is put back to the hat. This procedure, which is useful e.g. in Bootstrapping, fulfils the conditions of definition 2.1.

If we do the sampling without replacement (which often makes much sense), then the conditions of definition 2.1 are not fulfilled: each draw changes the population by removing one unit from the finite population. This leads to so called design-based inference, which is covered in the literature of sampling theory.

The difference to definition 2.1 is important especially if the sample size n is large compared to the size of population, but may be unimportant if $N \gg n$. In this course, we are considering only sampling according to definition 2.1.

Random samples can be summarized to a well defined summary called statistic.

Definition 2.2. Let $\mathbf{X} = X_1, \dots, X_n$ be a random sample of size n from population $f(x)$ and let $T(\mathbf{X})$ be a real-valued or vector-valued function, whose domain includes

the sample space of \mathbf{X} . The random variable $Y = T(\mathbf{X})$ is called a statistic (statistikka, tunnusluku). The probability distribution of Y is called the sampling distribution of Y (otantajakauma).

Examples of statistics

Sample minimum $\min(\mathbf{X})$

Sample maximum $\max(\mathbf{X})$

Sample median $\text{median}(\mathbf{X})$

Also other like sample mean, variance, standard deviation, quartiles etc.

Example 2.3. Let \mathbf{X} be a random sample of size 5 from *Exponential*(2) population. R script `statistics.R` illustrates the distributions of $\min(\mathbf{X})$, $\max(\mathbf{X})$, $\text{median}(\mathbf{X})$ and $\text{mean}(\mathbf{X})$ through simulation. Script repeats the sampling from Exponential distribution $M = 10000$ times and illustrates the distribution of the above mentioned sample statistic when $n = 5$.

Definition 2.3. The sample mean is the statistic defined by

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

Definition 2.4. The sample variance is the statistic defined by

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

- The sample standard deviation is the statistic defined by $S = \sqrt{S^2}$.
- The observed values of these random variables are denoted by \bar{X} , s^2 and s .

Theorem 2.1. Let X_1, \dots, X_n be any numbers. Then

$$\begin{aligned} \text{a) } \min_a \sum_{i=1}^n (x_i - a)^2 &= \sum_{i=1}^n (x_i - \bar{x})^2 \\ \text{b) } (n-1)s^2 &= \sum_{i=1}^n (x_i - \bar{x})^2 = \underbrace{\sum_{i=1}^n x_i^2 - n\bar{x}^2}_{\substack{\text{Compare to} \\ \text{Var}(\mathbf{X}) = E(\mathbf{X}^2) - (E(\mathbf{X}))^2}} \end{aligned}$$

Lemma 2.1. Let X_1, \dots, X_n be a random sample from a population and the $g(X)$ be a function such that $E(g(X_i))$ and $\text{Var}(g(X_i))$ exist. Then

$$E \left[\sum_{i=1}^n g(X_i) \right] = n \cdot E(g(X_i))$$

and

$$Var\left(\sum_{i=1}^n g(X_i)\right) = n \cdot Var(g(X_i))$$

Proof

$$E\left(\sum g(X_i)\right) = \sum E(g(X_i)) = *$$

Because X_i 's are identically distributed

$$* = nE(g(X_i))$$

$$\begin{aligned} Var\left(\sum g(X_i)\right) &= E\left[\sum g(X_i) - E\left(\sum g(X_i)\right)\right]^2 \\ &= E\left\{\sum \left[\underbrace{g(X_i) - E(g(X_i))}_{t_i}\right]\right\}^2 = * \end{aligned}$$

With t_i , we can get following

$$\begin{aligned} E\left(\sum t_i \sum t_i\right) &= E\left(\sum t_i^2 + \sum_{i \neq j} t_i t_j\right) \\ &= \sum E(t_i^2) + \sum_{i \neq j} E(t_i t_j) \end{aligned}$$

which we can expand back to original equation

$$\begin{aligned} * &= \underbrace{\sum E([g(X_i) - E(g(X_i))])^2}_{\sum E(t_i^2)} + \underbrace{\sum_{i \neq j} E[(g(X_i) - E(g(X_i)))(g(X_j) - E(g(X_j)))]}_{\sum_{i \neq j} E(t_i t_j)} \\ &= \sum Var(g(X_i)) + \sum_{i \neq j} cov(g(X_i), g(X_j)) \\ &= \underbrace{n \cdot Var(X_1)}_{\text{Because } X_i \text{'s are identically distributed}} + n(n-1) \cdot \underbrace{0}_{\text{Because } X_i \text{'s are independent}} = n \cdot Var[g(X_1)] \end{aligned}$$

Note: The proof of $E(\sum g(X_i))$ did not utilize the assumption of independence. Therefore it is valid also for dependent samples. The poof of $Var(\sum g(X_i))$ used the assumption of independence. Therefore it is not valid for dependent samples.

- For dependent samples, we get

$$Var\left(\sum g(X_i)\right) = \sum_{i=1}^n \sum_{j=1}^n cov[g(X_i), g(X_j)]$$

- Proof left as an exercise (recall, that $Var(X_i) = cov(X_i, X_i)$).

$$\begin{aligned}
\text{Var}(\mathbf{x}) &= \begin{bmatrix} \text{Var}(X_1) & \text{cov}(X_1, X_2) & \dots & \text{cov}(X_1, X_n) \\ \text{cov}(X_2, X_1) & \text{Var}(X_2) & \dots & \text{cov}(X_2, X_n) \\ \vdots & \vdots & \ddots & \vdots \\ \text{cov}(X_n, X_1) & \text{cov}(X_n, X_2) & \dots & \text{Var}(X_n) \end{bmatrix} \\
\text{Var}(\sum X_i) &= \sum \sum \text{cov}(X_i, X_j)
\end{aligned}$$

- That is, the sum of all elements of matrix $\text{Var}(\mathbf{X})$.

Note: This extends the rule

$$\text{Var}(X_1 + X_2) = \text{Var}(X_1) + \text{Var}(X_2) + 2 \cdot \text{cov}(X_1, X_2)$$

to general sum $\sum X_i$.

Theorem 2.2. Let X_1, \dots, X_n be a random sample from a population with mean μ and variance $\sigma^2 < \infty$. Then

- $E(\bar{X}) = \mu$
- $\text{Var}(\bar{X}) = \frac{\sigma^2}{n} \quad E(\bar{X} - \mu)^2$
- $E(S^2) = \sigma^2$

- Proof for a and b are familiar from ISI1 course. Proof of c applies the theorem 2.1 and is left as an exercise.

Note: Because $E(\bar{X}) = \mu$, the statistic \bar{X} is said to be an unbiased estimator (harhaton estimaattori) of population mean μ . Also S^2 is an unbiased estimator of population variance σ^2 .

Example 2.4. Illustrate these results through simulation. Intuitively, mean of an observed value of a statistic over a large number of samples should be close to expected value of the statistic in question. Therefore, (a), (b) and (c) of theorem 2.2 can be demonstrated by simulating M samples of a fixed size n and exploring how the means of the sample values of \bar{x} , $(\bar{x} - \mu)^2$ and s^2 behave as $M \rightarrow \infty$.

R-script `demonstratebias.R` implements this by assuming that the population has the $\text{Uniform}(0, 10)$ distribution and $n = 10$.

2.1 Distribution of sum

Theorem 2.3. (Convolution formula)

If X and Y are independent, continuous random variables with pdf's $f_X(x)$ and $f_Y(y)$, then the pdf of the sum $Z = x + Y$ is

$$f_Z(z) = \int_{-\infty}^{\infty} f_X(w)f_Y(z-w) dw$$

Proof: See Casella Berger, p 215-216.

Illustration: See example 1.20 of `notes.pdf`

Note: To compute the complete distribution of a sum $Z = \sum_{i=1}^n X_i$, where X_i 's do not need to be identically distributed, but they are independent, we need to apply theorem 2.3 iteratively. E.g. to find the distribution of $X_1 + X_2 + X_3$, you may first find the distribution of $Z = X_1 + X_2$ using theorem 2.3 and thereafter use theorem 2.3 again to find the pdf of $Z + X_3$. This is not trivial in very general case, but there are easier ways to do this, e.g. when X_i 's are identically distributed and independent (i.i.d).

However, recall that in any case, the moments are easy

$$E(\sum X_i) = \sum E(X_i)$$

$$Var(\sum X_i) = \sum Var(X_i) \quad (\text{If } X_i\text{'s are independent})$$

2.2 Sampling from Normally distributed population

Theorem 2.4. Let $X_{n \times 1}$ be random sample of size n from a $N(\mu, \sigma^2)$ distribution and let $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ and $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$. Then

- \bar{X} and S^2 are independent random variables.
- \bar{X} has $N(\mu, \frac{\sigma^2}{n})$
- $(n-1)S^2/\sigma^2$ has chi-squared distribution with parameter $n-1$. Parameter $n-1$ is commonly called the “degrees of freedom”.

Note: The chi-squared distribution with p degrees of freedom has the pdf

$$f_X(x) = \frac{1}{\underbrace{\Gamma(p/2)}_{\text{Gamma -function}}} X^{p/2-1} e^{-x/2}$$

- Gamma -function: $\Gamma(z) = \int_0^{\infty} x^{z-1} e^{-x} dx 2^{p/2}$

Lemma 2.2. We use notation χ_p^2 for a chi-squared random variables with p degrees of freedom

- If Z is a $N(0, 1)$ random variable, then $Z^2 \sim \chi_1^2$, that is, the square of a standard normal random variable is a chi-squared random variable.

- If X_1, \dots, X_n are independent and $X_i \sim \chi_{p_i}^2$ then $\sum X_i \sim \chi_{\sum p_i}^2$. That is, independent chi-squared random variables add to a chi-squared random variable, and the degrees of freedom also add.

Example 2.5. Illustrate the result of theorem 2.4 using R-script. In file `normalchi.R`, we implement the following:

Repeat M times

1. Generate random variables $X_i \sim N(0, 1)$, where $n = 1, \dots, 9$
2. Compute $Z = \sum x_i^2$ for each sample
3. Plot the histogram of the M obtained values of Z and compare to the distribution χ_n^2 .

Note: If \bar{X}_n is a random sample from $N(\mu, \sigma^2)$ population, then the random variable

$$Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$$

where $\mu = E(\bar{X})$ and $\sigma/\sqrt{n} = \text{Var}(\bar{X})$.

- This transformation could be used to make inference on μ using the observed \bar{x} , if σ^2 is known.
- However, that is not the case, and therefore we use $\frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$ which can also be written as

$$\frac{(\bar{X} - \mu)/(\sigma/\sqrt{n})}{\sqrt{S^2/\sigma^2}}$$

where $(\bar{X} - \mu)/(\sigma/\sqrt{n}) \sim N(0, 1)$ and $\sqrt{S^2/\sigma^2} \sim \chi_{n-1}^2$.

Theorem 2.5. Let \bar{X} be a random sample from $N(\mu, \sigma^2)$. The random variable

$$T = \frac{\bar{X} - \mu}{S/\sqrt{n}}$$

has Student's t-distribution with parameter (n-1) "degrees of freedom".

In general, a random variable T has Student's t-distribution with p degrees of freedom ($T \sim t_p$) if T has the pdf

$$f_T(t) = \frac{\Gamma(\frac{p+1}{2})}{\Gamma(\frac{p}{2})} \frac{1}{(p\pi)^{1/2}} \frac{1}{(1 + t^2/p)^{(p+1)/2}}$$

Proof: See Casella Berger, p.223-224.

Illustration: File `normalchi.R`

- We use samples of sizes 3, 4, \dots , 100 and $M = 10000$.

- Plot empirical histogram of $\frac{\bar{X} - \mu}{s/\sqrt{n}}$
- Compare to t_{n-1} distribution and to $N(0, \sigma^2)$.
- $\mu = 5$ and $\sigma^2 = 3^2$.
- Because $S^2 \rightarrow \sigma^2$ as n increases, the difference between $N(0, 1)$ and t_{n-1} becomes meaningless as $n \rightarrow \infty$.

Note: t_p has only p moments. Especially $E(T_p) = 0$ if $p > 1$.