

iid: independently identically distribution

Introduction

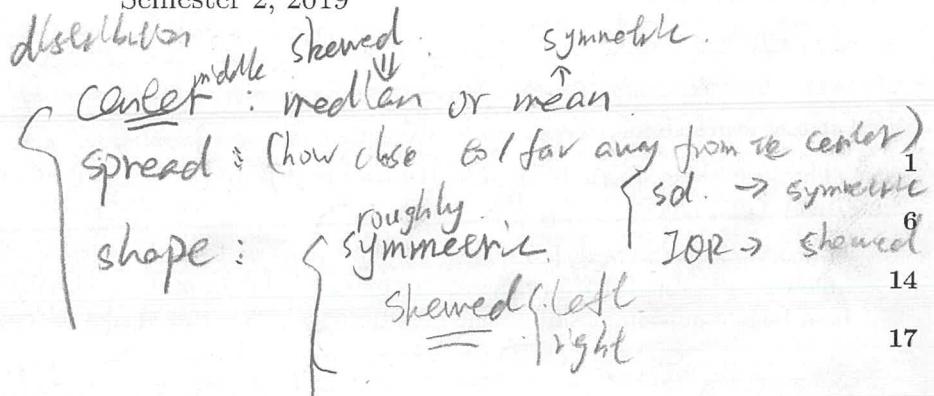
(Module 1)

Statistics (MAST20005) & Elements of Statistics (MAST90058)

Semester 2, 2019

Contents

- 1 Subject information
- 2 Review of probability
- 3 Descriptive statistics
- 4 Basic data visualisations



Aims of this module

- Brief information about this subject
- Brief revision of some prerequisite knowledge (probability)
- Introduce some basic elements of statistics, data analysis and visualisation

1 Subject information

What is statistics?

Let's see some examples...

Examples

- Weather forecasts: Bureau of Meteorology
- Poll aggregation: FiveThirtyEight, The Guardian
- Climate change modelling: Australian Academy of Science
- Discovery of the Higgs Boson (the 'God Particle'): van Dyk (2014)
- Smoking leads to lung cancer: Doll & Hill (1945)
- A/B testing for websites: Google and 41 shades of blue

Tingjin's example

- Real estate price modelling

Damjan's examples

- Genome-wide association studies
- Web analytics
- Lung testing in infants
- Skin texture image analysis
- Wedding 'guestimation'

Goals of statistics

- Answer questions using *data*
- Evaluate *evidence*
- Optimise *study design*
- Make *decisions*

And, importantly:

- Clarify *assumptions*
- Quantify *uncertainty*

Why study statistics?

"The best thing about being a statistician is that you get to play in everyone's backyard." John W. Tukey (1915 2000)

"I keep saying the sexy job in the next ten years will be *statisticians*... The ability to take data – to be able to understand it, to process it, to extract value from it, to visualize it, to communicate it's going to be a hugely important skill in the next decades..." Hal Varian, Google's Chief Economist, Jan 2009

The best job

U.S. News Best Business Jobs in 2019:

1. *Statistician*
2. Mathematician
- ...
6. Actuary

CareerCast (recruitment website) Best Jobs of 2017:

1. *Statistician*
- ...
5. Data Scientist
6. University Professor

Subject overview

Statistics (MAST20005), Elements of Statistics (MAST90058)

These subjects introduce the basic elements of *statistical modelling*, *statistical computation* and *data analysis*. They demonstrate that many commonly used statistical procedures arise as applications of a common theory. They are an entry point to further study of both *mathematical* and *applied* statistics, as well as broader data science.

Students will develop the ability to fit statistical models to data, estimate parameters of interest and test hypotheses. Both classical and Bayesian approaches will be covered. The importance of the underlying *mathematical theory of statistics* and the use of *modern statistical software* will be emphasised.

Joint teaching

MAST20005 and MAST90058 share the same lectures but have separate tutorials and lab classes. The teaching and assessment material for both subjects will overlap significantly.

Subject website (LMS)

- Full information is on the subject website, available through the Learning Management System (LMS).
- Only a brief overview is covered in these notes. Please read all of the info on the LMS as well.
- New material (e.g. problem sets, assignments, solutions) and announcements will appear regularly on the LMS.

Subject structure

- *Lectures:* Three 1-hour lectures per week. Lecture notes/slides will appear on the LMS.
- *Tutorials:* One 1-hour tutorial per week (starting in week 2). Tutorial problems and solutions will appear on the LMS.
- *Computer lab classes:* One 1-hour lab per week (starting in week 2), immediately following the tutorial. Lab notes, exercises and solutions will appear on the LMS.

Computing

- This subject introduces basic statistical computing and programming skills.
- We make extensive use of the R statistical software environment.
- Knowledge of R will be *essential* for some of the tutorial problems, assignment questions and will also be examined.
- We will use the RStudio program as a convenient interface with R.

Textbook

R. Hogg, E. Tanis, and D. Zimmerman. *Probability and Statistical Inference*. 9th Edition, Pearson, 2015.

- This subject is based on Chapters 6 - 9.
- Some of the teaching material is taken from the textbook.
- This textbook is being *phased out* for this subject.
- There are *important differences* between the subject content and the textbook. We will point many of these out, but please ask if unsure.

Assessment

- 3 assignments (20%)
 1. Hand out at the start of week 4, due at the end of week 5
 2. Hand out at the start of week 7, due at the end of week 8
 3. Hand out at the start of week 10, due at the end of week 11
- 45-minute computer lab test held in week 12 (10%)
- 3-hour written examination in the examination period (70%)

Plagiarism declaration

- Everyone must complete the *Plagiarism Declaration Form*
- Do this on the LMS
- Do this ASAP!

Staff contacts

Subject coordinator / Lecturer (stream 2)

Dr Damjan Vukcevic <damjan.vukcevic@unimelb.edu.au>
Room 309, Old Geology South

Lecturer (stream 1)

Dr Tingjin Chu <tingjin.chu@unimelb.edu.au>
Room 104, Peter Hall Building

Tutorial coordinator

Dr Robert Maillardet <rjmail@unimelb.edu.au>
Room G48, Peter Hall Building

See the LMS for details of consultation hours

Online discussion forum (Piazza)

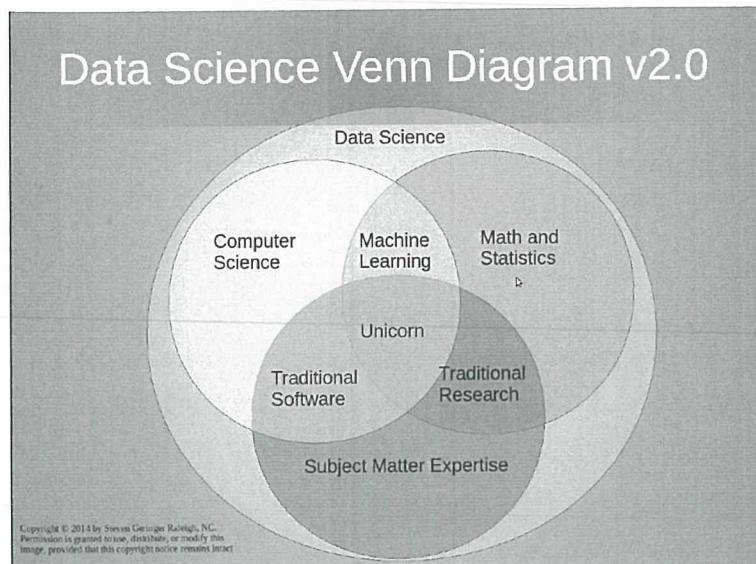
- Access via the LMS
- Post any general questions on the forum
- Do *not* send them by email to staff
- You can answer each others' questions
- Staff will also help to answer questions

Student representatives

Student representatives assist the teaching staff to ensure good communication and feedback from students.

See the LMS to find the contact details of your representatives.

What is Data Science?



Data science is a 'team sport'

Read more at: Data science is inclusive

How to succeed in statistics / data science?

- Get experience with *real data*
- Develop your computational skills, *learn R*

- Understand the *mathematical theory*
- Collaborate with others, *use Piazza*

This subject is challenging

- It is mathematical
 - Manipulating equations
 - Calculus
 - Probability
 - Proofs
- But the ‘real’ world also matters
 - Context can ‘trump’ mathematics
 - More than one correct answer
 - Often uncertain about the answer

Diversity

In 2017: **341 students**

60%	Bachelor of Commerce
24%	Bachelor of Science
6%	Master of Science (Bioinformatics)
10%	8 other degrees/categories

What are your strengths and weaknesses?

Get extra help

- Your classmates
- Piazza
- Textbooks
- Consultation hours
- Oasis

Homework

1. Complete plagiarism declaration on the LMS
2. Log in to Piazza
3. Install RStudio on your computer
4. Start reading lab notes for week 2 (long!)

Tips

The best way to learn statistics is by **solving problems** and ‘**getting your hands dirty**’ with data.

We encourage you to attend all lectures, tutorial and computer labs to get as much practice and feedback as possible.

Good luck!

2 Review of probability

Why probability?

- It forms the mathematical foundation for statistical models and procedures
- Let's review what we know already...

Random variables (notation)

- Random variables (rvs) are denoted by uppercase letters: X, Y, Z , etc.
- Outcomes, or realisations, of random variables are denoted by corresponding lowercase letters: x, y, z , etc.

Distribution functions

- The cumulative distribution function (cdf) of X is

$$F(x) = \Pr(X \leq x), \quad -\infty < x < \infty$$

- If X is a continuous rv then it has a probability density function (pdf), $f(x)$, that satisfies

$$f(x) = F'(x) = \frac{d}{dx} F(x)$$

$$F(x) = \int_{-\infty}^x f(t) dt$$

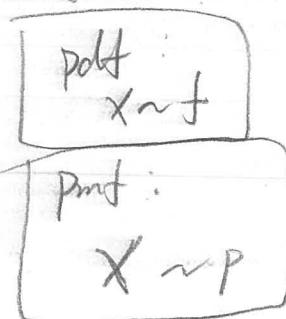
- If X is a discrete rv then it has a probability mass function (pmf),

$$p(x) = \Pr(X = x), \quad x \in \Omega$$

where Ω is a discrete set, e.g. $\Omega = \{1, 2, \dots\}$.

- $\Pr(X > x) = 1 - F(x)$ is called a tail probability of X
- $F(x)$ increases to 1 as $x \rightarrow \infty$ and decreases to 0 as $x \rightarrow -\infty$
- If the rv has a certain distribution with pdf f (or pmf p), we write

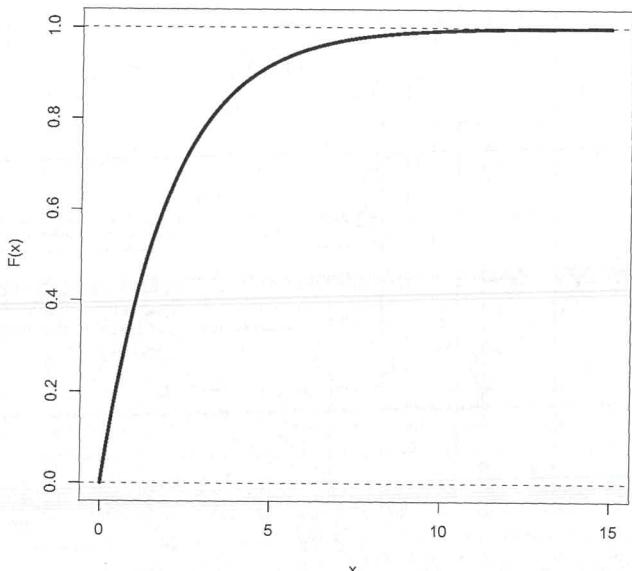
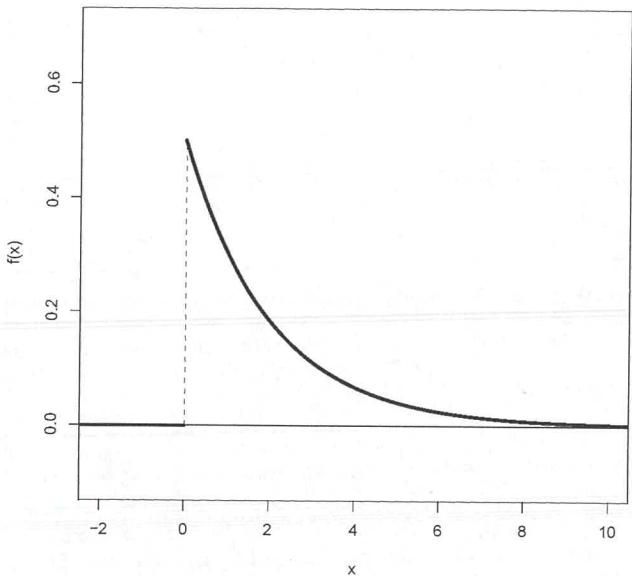
$$X \sim f \quad (\text{or } X \sim p)$$



Example: Unemployment duration

A large group of individuals have recently lost their jobs. Let X denote the length of time (in months) that any particular individual will stay unemployed. It was found that this was well-described by the following pdf:

$$f(x) = \begin{cases} \frac{1}{2}e^{-x/2}, & x \geq 0, \\ 0, & \text{otherwise.} \end{cases}$$



Clearly, $f(x) \geq 0$ for any x and the total area under the pdf is:

$$\Pr(-\infty < X < \infty) = \int_0^\infty \frac{1}{2} e^{-x/2} dx = \frac{1}{2} \left[-2e^{-x/2} \right]_0^\infty = 1.$$

The probability that a person in the population finds a new job within 3 months is:

$$\Pr(0 \leq X \leq 3) = \int_0^3 \frac{1}{2} e^{-x/2} dx = \frac{1}{2} \left[-2e^{-x/2} \right]_0^3 = 0.7769.$$

$$\Pr(0 \leq X \leq 3) = \int_0^3 \frac{1}{2} e^{-x/2} dx = \frac{1}{2} \int_0^3 e^{-x/2} dx = \frac{1}{2} \left[-2e^{-x/2} \right]_0^3 = \frac{1}{2} \left[-2e^{-3/2} + 2 \right]$$

Example: Received calls

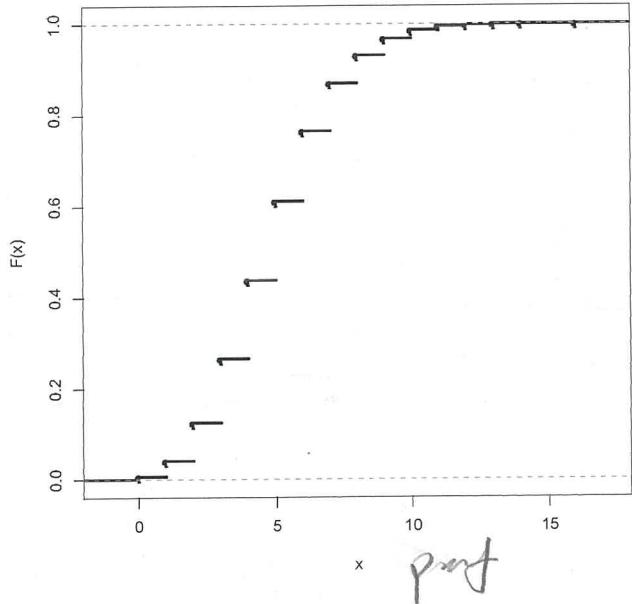
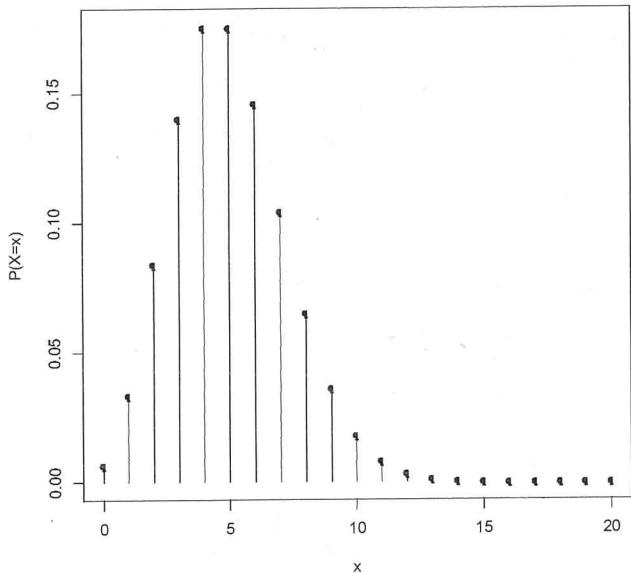
The number of calls received by an office in a given day, X , is well-represented by a pmf with the following expression:

$$p(x) = \frac{e^{-5} 5^x}{x!}, \quad x \in \{0, 1, 2, \dots\},$$

where $x! = 1 \cdot 2 \cdots (x-1) \cdot x$ and $0! = 1$. For example,

$$\Pr(X = 1) = e^{-5} 5 = 0.03368$$

$$\Pr(X = 3) = \frac{e^{-5} 5^3}{3 \cdot 2 \cdot 1} = 0.1403$$



To show that $p(x)$ is a pmf we need to show

$$\sum_{x=0}^{\infty} p(x) = p(0) + p(1) + p(2) + \dots = 1.$$

Since the Taylor series expansion of e^z is $\sum_{i=0}^{\infty} z^i / i!$, we can write

$$\sum_{i=0}^{\infty} \frac{e^{-5} 5^x}{x!} = e^{-5} \sum_{i=0}^{\infty} \frac{5^x}{x!} = e^5 e^{-5} = 1.$$

Moments and variance

- The expected value (or the *first moment*) of a rv is denoted by $\mathbb{E}(X)$ and



$$\mathbb{E}(X) = \sum_{x=-\infty}^{\infty} x p(x)$$

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

(discrete rv)

(continuous rv)

$$\mathbb{E}(x^2) = \mathbb{E}(x^3) =$$

$$\mathbb{E}(X^k) = \sum_{x=-\infty}^{\infty} x^k p(x) \quad \mathbb{E}(x^2) = \int_{-\infty}^{\infty} x^2 f(x) dx$$

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

(discrete rv)

(continuous rv)

- More generally for a function $g(x)$ we can compute

$$\mathbb{E}(g(X)) = \sum_{x=-\infty}^{\infty} g(x) p(x) \quad (\text{discrete rv})$$

$$\mathbb{E}(g(X)) = \int_{-\infty}^{\infty} g(x) f(x) dx \quad (\text{continuous rv})$$

Letting $g(x) = x^k$ gives the moments.

- The variance of X is defined by

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

and the standard deviation of X is $\text{sd}(X) = \sqrt{\text{var}(X)}$

- "Computational" formula: $\text{var}(X) = \mathbb{E}(X^2) - \{\mathbb{E}(X)^2\}$

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

$$\int_{-\infty}^{\infty} x^2 f(x) dx.$$



Basic properties of expectation and variance

- For any rv X and constant c ,

$$\mathbb{E}(cX) = c\mathbb{E}(X), \quad \text{var}(cX) = c^2 \text{var}(X)$$

- For any two rvs X and Y ,

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

- For any two independent rvs X and Y ,

$$\text{Variance} \quad \text{Independent}$$

$$\text{var}(X + Y) = \text{var}(X) + \text{var}(Y)$$

- More generally,

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

where $\text{cov}(X, Y)$ is the covariance between X and Y

相当于是夹角.

Covariance

协方差

衡量两个变量的总体误差

① 两变量变化趋势一致

cov 为正

- Definition of covariance:

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

- Specifically, for the continuous case

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

where $f(x, y)$ is the bivariate pdf for pair (X, Y) .

- "Computational" formula:

$$\text{cov}(X, Y) = \mathbb{E}(XY) - \mathbb{E}(X)\mathbb{E}(Y)$$

$$\text{cov}(X, Y) = \mathbb{E} \{ (X - \mathbb{E}(X))(Y - \mathbb{E}(Y)) \} = \mathbb{E}(XY) - \mathbb{E}(X)\mathbb{E}(Y)$$

① 若 X, Y independent,

$\Rightarrow \text{cov}(X, Y) = 0$

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

反过来不行

Correlation

- If $\text{cov}(X, Y) > 0$ then X and Y are positively correlated

- If $\text{cov}(X, Y) \leq 0$ then X and Y are negatively correlated

- If $\text{cov}(X, Y) = 0$ then X and Y are uncorrelated

- The correlation between X and Y is defined as:

Pearson 相关系数

$$\rho = \text{cor}(X, Y) = \frac{\text{cov}(X, Y)}{\text{sd}(X)\text{sd}(Y)}, \quad -1 \leq \rho \leq 1$$

② 若 $\text{cov}(X, Y) = 0$

$\Rightarrow X, Y$ independent.

- When $\rho = \pm 1$ then X and Y are perfectly correlated

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

Moment generating functions

mgf

- A moment generating function (mgf) of a rv X is

$$M_X(t) = \mathbb{E}(e^{tX}), \quad t \in (-\infty, \infty)$$

- It enables us to generate moments of X by differentiating at $t = 0$

$$M'_X(0) = \mathbb{E}(X)$$

mgf determines distribution $M_X^{(k)}(0) = \mathbb{E}(X^k), \quad k \geq 1$

- The mgf uniquely determines a distribution. Hence, knowing the mgf is the same as knowing the distribution.

- If X and Y are independent rvs,

$$M_{X+Y}(t) = \mathbb{E}\{e^{t(X+Y)}\} = \mathbb{E}\{e^{tX}\} \mathbb{E}\{e^{tY}\} = M_X(t)M_Y(t)$$

i.e. the mgf of the sum is the product of individual mgfs.

Bernoulli distribution

- X takes on the values 1 (success) or 0 (failure)

- $X \sim \text{Be}(p)$ with pmf

$$P(X) = p^x (1-p)^{1-x}$$

$$p(x) = p^x (1-p)^{1-x}, \quad x \in \{0, 1\}$$

- Properties:

$$\mathbb{E}(X) = p$$

$$\text{var}(X) = p(1-p)$$

$$M_X(t) = pe^t + 1 - p$$

Binomial distribution

n -Bernoulli distribution

- $X \sim \text{Bi}(n, p)$ with pmf

$$p(x) = \binom{n}{x} p^x (1-p)^{n-x}, \quad x \in \{0, 1, \dots, n\}$$

- Properties:

$$\mathbb{E}(X) = np$$

$$\text{var}(X) = np(1-p)$$

$$M_X(t) = (pe^t + 1 - p)^n$$

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

Poisson distribution

: 入: 单位时间随机事件的平均发生次数.

- $X \sim \text{Pn}(\lambda)$ with pmf

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

$$P(X) = e^{-\lambda} \frac{\lambda^x}{x!}$$

- Properties:

$$\mathbb{E}(X) = \text{var}(X) = \lambda$$

$$M_X(t) = e^{\lambda(e^t - 1)}$$

- It arises as an approximation to $\text{Bi}(n, p)$. Letting $\lambda = np$ gives

$$p(x) = \binom{n}{x} p^x (1-p)^{n-x} \approx e^{-\lambda} \frac{\lambda^x}{x!}$$

Poisson \approx Binomial

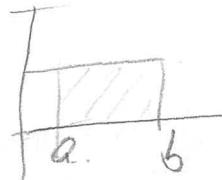
as $n \rightarrow \infty$ and $p \rightarrow 0$.

Uniform distribution

- $X \sim \text{Unif}(a, b)$ with pdf

=

$$f(x) = \frac{1}{b-a}, \quad x \in (a, b)$$



- Properties:

$$\mathbb{E}(X) = \frac{(a+b)}{2}$$

$$\text{var}(X) = \frac{(b-a)^2}{12}$$

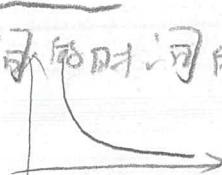
$$M_X(t) = \frac{e^{tb} - e^{ta}}{t(b-a)}$$

- If $b = 1$ and $a = 0$, this is known as the uniform distribution over the unit interval.

Exponential distribution

- $X \sim \text{Exp}(\lambda)$ with pdf

$$f(x) = \lambda e^{-\lambda x}, \quad x \in [0, \infty)$$



- It approximates "time until first success" for independent $\text{Be}(p)$ trials every Δt units of time with $p = \lambda \Delta t$ and $\Delta t \rightarrow 0$

- Properties:

事件以恒定平均速率连续且独立地发生的过程

$$\begin{cases} \mathbb{E}(X) = 1/\lambda \\ \text{var}(X) = 1/\lambda^2 \\ M_X(t) = \frac{\lambda}{\lambda-t} \end{cases}$$



- It is famous for being the only continuous distribution with the memoryless property:

$$\Pr(X > y + x | X > y) = \Pr(X > x), \quad x \geq 0, \quad y \geq 0.$$

Normal distribution

- $X \sim N(\mu, \sigma^2)$ with pdf



$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, \quad x \in (-\infty, \infty), \quad \mu \in (-\infty, \infty), \quad \sigma > 0$$



- It is important in applications because of the Central Limit Theorem (CLT)

- Properties:

$$\mathbb{E}(X) = \mu$$

$$\text{var}(X) = \sigma^2$$

$$M_X(t) = e^{t\mu + t^2\sigma^2/2}$$

- When $\mu = 0$ and $\sigma = 1$ we have the standard normal distribution.

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

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



Quantiles

Let X be a continuous rv. The p th quantile of its distribution is a number π_p such that $p = \Pr(X \leq \pi_p) = F(\pi_p)$. In other words, the area under $f(x)$ to the left of π_p is p :

$$p = \int_{-\infty}^{\pi_p} f(x) dx = F(\pi_p)$$

- π_p is also called the $(100p)$ th percentile
- The 50th percentile (0.5 quantile) is the median, denoted by $m = \pi_{0.5}$
- The 25th and 75th percentiles are the first and third quartiles, denoted by $q_1 = \pi_{0.25}$ and $q_3 = \pi_{0.75}$

Example: Weibull distribution

可靠性分析，寿命推断理论基础

The time X until failure of a certain product has the pdf

$$f(x) = \frac{3x^2}{4} e^{-(x/4)^3}, \quad x \in (0, \infty).$$

The cdf is

$$F(x) = 1 - e^{-(x/4)^3}, \quad x \in (0, \infty)$$

Then $\pi_{0.3}$ satisfies $0.3 = F(\pi_{0.3})$. Therefore,

$$\begin{aligned} & 1 - e^{-(\pi_{0.3}/4)^3} = 0.3 \\ \Rightarrow & \ln(0.7) = -(\pi_{0.3}/4)^3 \\ \Rightarrow & \pi_{0.3} = -4(\ln 0.7)^{1/3} = 2.84. \end{aligned}$$

Law of Large Numbers (LLN)

Consider a collection X_1, \dots, X_n of independent and identically distributed (iid) random variables with $\mathbb{E}(X) = \mu < \infty$, then with probability 1 we have:

$$\frac{1}{n} \sum_{i=1}^n X_i \rightarrow \mu, \quad \text{as } n \rightarrow \infty.$$

The LLN 'guarantees' that long-run averages behave as we expect them to:

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

Central Limit Theorem (CLT)

Consider a collection X_1, \dots, X_n of iid rvs with $\mathbb{E}(X) = \mu < \infty$ and $\text{var}(X) = \sigma^2 < \infty$. Let,

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

Then

$$\frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$$

follows a $N(0, 1)$ distribution as $n \rightarrow \infty$.

This is an extremely important theorem! It provides the 'magic' that will make statistical analysis work.

Example Example

Let X_1, \dots, X_{25} be iid rvs where $X_i \sim \text{Exp}(\lambda = 1/5)$.

Recall that $\mathbb{E}(X) = 5$.

Thus, the LLN implies

$$\bar{X} \rightarrow \mathbb{E}(X) = 5.$$

[25个]

等于 25.

Moreover, since $\text{var}(X) = 1/\lambda^2 = 25$, we have

$$\bar{X} \approx N\left(\frac{1}{\lambda}, \frac{1}{n\lambda^2}\right) = N\left(5, \frac{5^2}{25}\right)$$

Is $n = 25$ large enough?

A simulation exercise

Generate $B = 1000$ samples of size n . For each sample compute \bar{x} . The continuous curve is the normal $N(5, 5^2/n)$ distribution prescribed by the CLT.

Sample 1: $x_1^{(1)}, \dots, x_n^{(1)} \rightarrow \bar{x}^{(1)}$

Sample 2: $x_1^{(2)}, \dots, x_n^{(2)} \rightarrow \bar{x}^{(2)}$

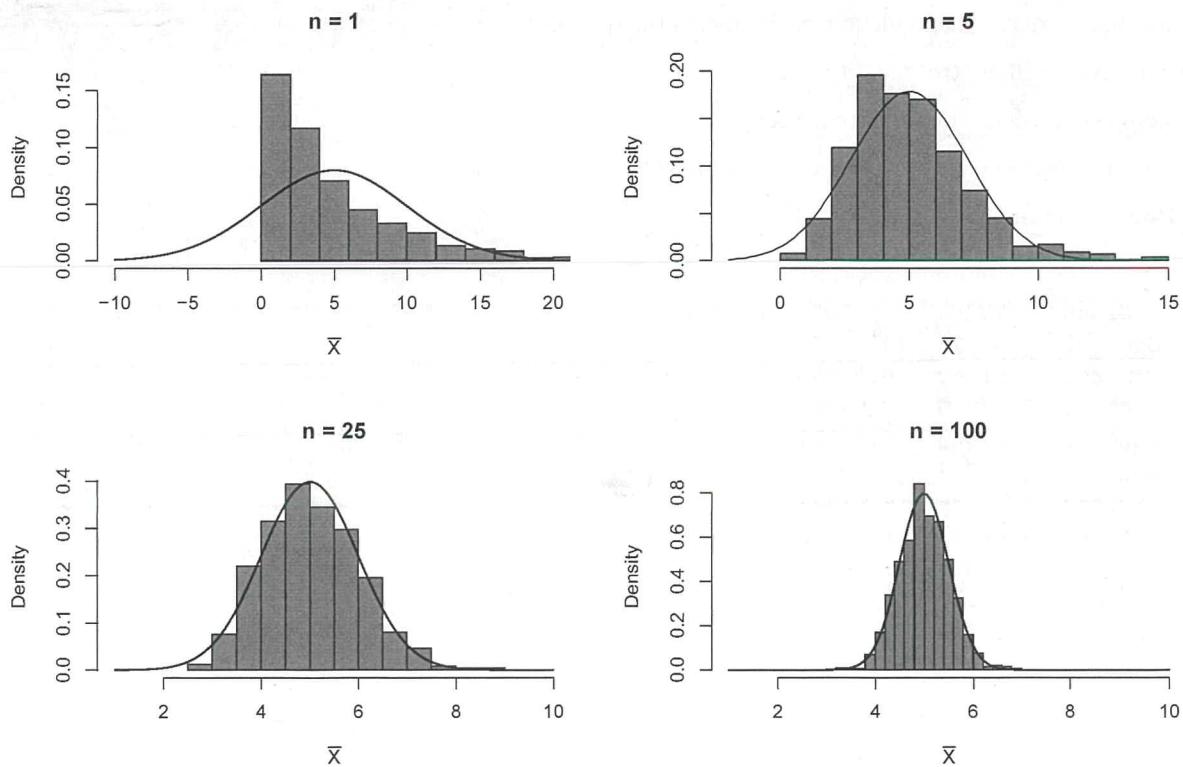
⋮

Sample B: $x_1^{(B)}, \dots, x_n^{(B)} \rightarrow \bar{x}^{(B)}$

Then represent the distribution of $\{\bar{x}^{(b)}, b = 1, \dots, B\}$ by a histogram.

A simulation exercise

The distribution of \bar{X} approaches the theoretical distribution (CLT). Moreover it will be more and more concentrated around μ (LLN). To see this, note that $\text{var}(\bar{X}) = \sigma^2/n \rightarrow 0$ as $n \rightarrow \infty$.



$$1 = \int_0^2 ax^3 dx = \left[\frac{a}{4} x^4 \right]_0^2 = 2a$$

Challenge problem

Let X_1, X_2, \dots, X_{25} be iid rvs with pdf $f(x) = ax^3$ where $0 < x < 2$.

1. What is the value of a ? $a = \frac{1}{4}$

2. Calculate $\mathbb{E}(X_1)$ and $\text{var}(X_1)$.

3. What is an approximate value of $\Pr(\bar{X} < 1.5)$?

$$4a - 0 \Rightarrow 1$$

$$a = \frac{1}{4}$$

$$\boxed{3} \quad \mathbb{E}(X_1) = \int_0^2 x \cdot f(x) dx$$

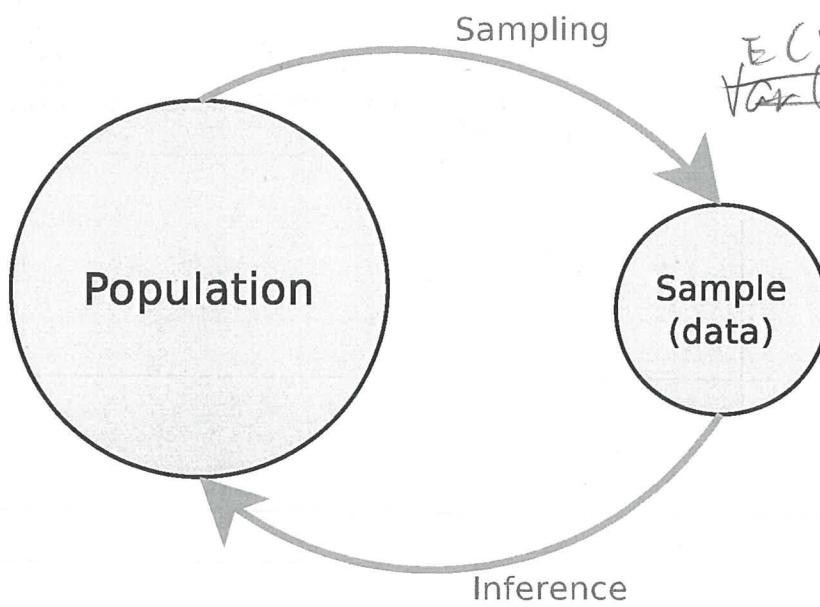
$$= \int_0^2 x \cdot \frac{1}{4} x^3 dx$$

$$= \int_0^2 \frac{1}{4} x^4 dx = \left[\frac{1}{20} x^5 \right]_0^2$$

$$= \frac{32}{20} - 0 = \frac{8}{5}$$

3 Descriptive statistics

Statistics: the big picture



$$\mathbb{E}(X_1^2) = \int_0^2 x^2 \cdot f(x) dx$$

$$= \int_0^2 x^2 \cdot \frac{1}{4} x^3 dx$$

$$= \int_0^2 \frac{1}{4} x^5 dx$$

$$= \boxed{4} \left[\frac{1}{20} x^6 \right]_0^2$$

$$= \frac{1}{20} \cdot 2^6 = \frac{16}{6} = \frac{8}{3}$$

$$\text{Var}(X_1) = \mathbb{E}(X_1^2) - (\mathbb{E}(X_1))^2$$

$$= \frac{8}{3} - \left(\frac{8}{5} \right)^2 = \frac{8}{75}$$

Example: Stress and cancer

- An experiment gives independent measurements on 10 mice
- Mice are divided in control and stress groups
- The biologist considers two different proteins:
 - Vascular endothelial growth factor C (VEGFC)
 - Prostaglandin-endoperoxide synthase 2 (COX2)

Mouse	Group	VEGFC	COX2
1	Control	0.96718	14.05901
2	Control	0.51940	6.92926
3	Control	0.73276	0.02799
4	Control	0.96008	6.16924
5	Control	1.25964	7.32697
6	Stress	4.05745	6.45443
7	Stress	2.41335	12.95572
8	Stress	1.52595	13.26786
9	Stress	6.07073	55.03024
10	Stress	5.07592	29.92790

$$\bar{X} = \frac{\sum X_i}{n}$$

$$\star \quad Z = \frac{\bar{X} - \mu}{\sqrt{\frac{\sigma^2}{n}}} \approx N(0, 1)$$

$$\Pr(\bar{X} < 1.5) = \Pr(Z < \frac{1.5 - \mu}{\sqrt{\frac{\sigma^2}{n}}})$$

$$= \Pr(Z < -1.53) =$$

$$\approx \Phi(-1.53) = 0.063$$

by computer

Data & sampling

- The data are numbers:

x_1, \dots, x_n

$t \in \mathbb{R}$

- The model for the data is a random sample, that is a sequence of iid rvs:

X_1, X_2, \dots, X_n

$\sim \mathcal{B}$

This model is equivalent to random selection from a hypothetical infinite population.)

- The goal is to use the data to learn about the distribution of the random variables (and, therefore, the population).

Statistic function of sample rvs. sample → population realisations

- A statistic $T = \phi(X_1, \dots, X_n)$ is a function of the sample and its realisation is denoted by $t = \phi(x_1, \dots, x_n)$.

- Note: the word "statistic" can also be used to refer to both the realisation, t , as well as the random variable, T . Sometime need to be more specific about which one is meant.

- A statistic has two purposes:

- Describe or summarise the sample descriptive statistics
- Estimate the distribution generating the sample inferential statistics

- A statistic can be both descriptive and inferential, it depends on how you wish to use/interpret it (see later)

- We now introduce some commonly used descriptive statistics...

$t \in \mathbb{R}$

{ descriptive
inferential

{ $T \in \text{Random variable}$
 $t = \phi(x_1, \dots, x_n)$ [realisations]

Moment statistics

$$\text{Sample mean} = \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{23.59}{10} = 2.359$$

$$\text{Sample variance} = s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = 3.98761$$

$$\text{Sample standard deviation} = s = \sqrt{3.98761} = 1.9969$$

These are 'sample' or 'empirical' versions of moments of a random variable

Empirical means 'derived from the data'

Order statistics

sort

Arrange the sample x_1, \dots, x_n in order of increasing magnitude and define:

$$x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$$

Then $x_{(k)}$ is the kth order statistic.

Special cases:

- $x_{(1)}$ is the sample minimum
- $x_{(n)}$ is the sample maximum

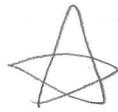
For the example data,

$$x_{(1)} = 0.52, \quad x_{(2)} = 0.73, \quad \dots, \quad x_{(10)} = 6.07$$

What is $x_{(3.25)}$?

Let it be 0.25 of the way from $x_{(3)}$ to $x_{(4)}$,

$$\begin{aligned}x_{(3.25)} &= x_{(3)} + 0.25 \cdot (x_{(4)} - x_{(3)}) \\&= 0.96 + 0.25 \cdot (0.97 - 0.96) \\&= 0.9625\end{aligned}$$



In other words, define it via linear interpolation.

Exercise: verify that $x_{(7.75)} = 3.6480$

$$x_{(7.75)} = x_{(7)} + 0.75 \cdot (x_{(8)} - x_{(7)})$$

Why do this? It allows us to define...

Sample quantiles

General definition ('Type 7' quantiles):

$$P \rightarrow P_r$$

Type 7:

$$\hat{\pi}_p = x_{(k)}, \text{ where } k = 1 + (n-1)p$$

Special cases:

$$k = 1 + (n-1)x_{0.5} = 5.5$$

$$\text{Sample median} = \hat{\pi}_{0.5} = x_{(5.5)} = \frac{1.26 + 1.53}{2} = 1.395$$

$$\text{Sample 1st quartile} = \hat{\pi}_{0.25} = x_{(3.25)} = 0.9625$$

$$\text{Sample 3rd quartile} = \hat{\pi}_{0.75} = x_{(7.75)} = 3.6480$$

Also:

IQR

$$\text{Interquartile range} = \hat{\pi}_{0.75} - \hat{\pi}_{0.25} = 2.685$$

$\hat{\pi}_{0.25}$ and $\hat{\pi}_{0.75}$ contain about 50% of the sample between them

Note: Type 7 quantiles are the default in R, but there are many alternatives! Don't worry too much about the differences between them. We will discuss this in a bit more detail later in the semester.

Some descriptive statistics in R

```
> x <- round(VEGFC, digit = 2)

> x
[1] 0.97 0.52 0.73 0.96 1.26 4.06 2.41 1.53 6.07 5.08

> sort(x) # order statistics
[1] 0.52 0.73 0.96 0.97 1.26 1.53 2.41 4.06 5.08 6.07

> summary(x) # sample mean & sample quantiles
   Min. 1st Qu. Median Mean 3rd Qu. Max.
0.5200 0.9625 1.3950 2.3590 3.6475 6.0700

> var(x) # sample variance
[1] 3.98761

> sd(x) # sample standard deviation
[1] 1.9969

> IQR(x) # interquartile range
[1] 2.685
```

Frequency statistics

Can also define empirical versions of pdf, pmf, cdf

Will see in the next section...

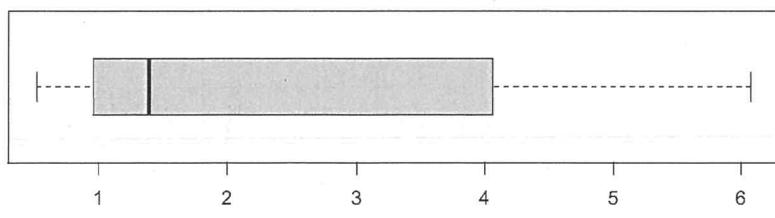
4 Basic data visualisations

Box plot

Graphical summary of data from a single variable

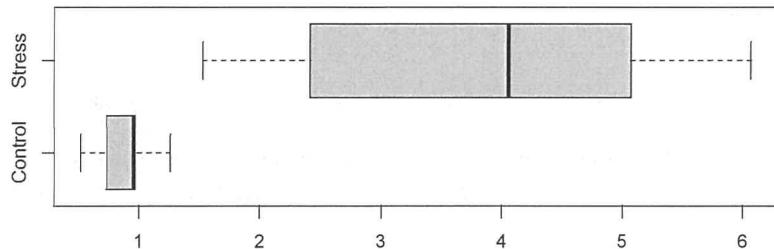
Main box: $\hat{\pi}_{0.25}, \hat{\pi}_{0.5}, \hat{\pi}_{0.75}$

'Whiskers': $x_{(1)}, x_{(n)}$ (but R does something more complicated, see tutorial problems)



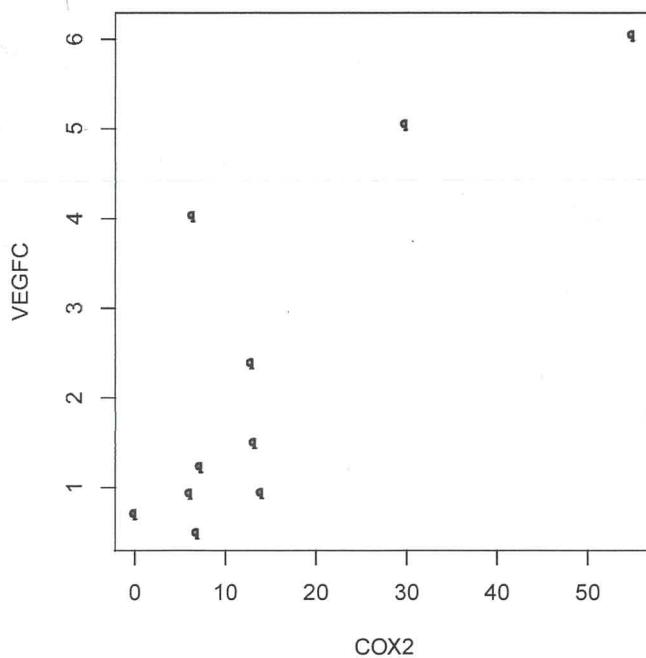
Convenient way of comparing data from different groups

Example: VEGFC (Stress vs Control)



Scatter plot

For comparing data from two variables (usually continuous)



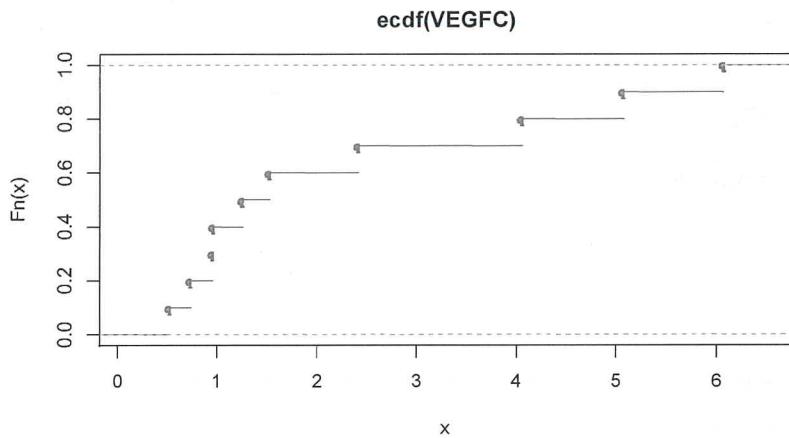
Empirical cdf

The sample cdf, or empirical cdf, is defined as

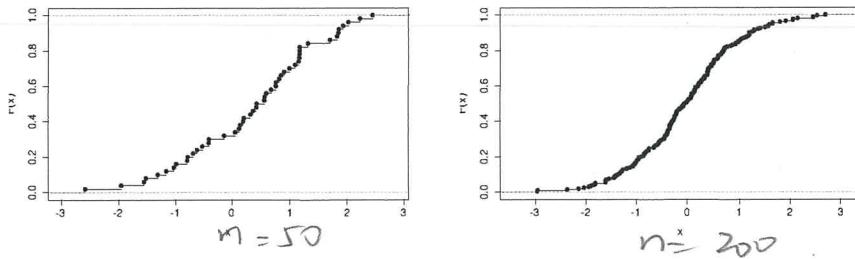
How many $\hat{F}(x) = \frac{1}{n} \sum_{i=1}^n I(x_i \leq x)$ are smaller than interval value
where $I(\cdot)$ is the indicator function ($I(x_i \leq x)$ has value 1 if $x_i \leq x$ and value 0 if $x_i > x$).

For example, for the previous data,

$$\hat{F}(2) = \frac{1}{10} \sum_{i=1}^{10} I(x_i \leq 2) = \frac{6}{10} = 0.6$$



It has the form of a discrete cdf. However, it will approximate the cdf of a continuous variable if the sample size is large. The following diagram shows cdfs based on $n = 50$ and $n = 200$ observations sampled from a standard normal distribution, $N(0, 1)$.



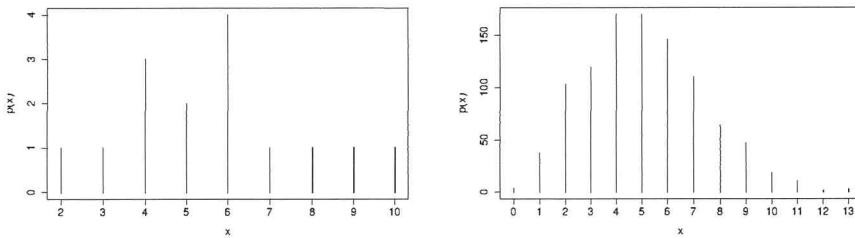
Empirical pmf

If the underlying variable is discrete we use the pmf corresponding to the sample cdf \hat{F}

$$\hat{p}(x) = \frac{1}{n} \sum_{i=1}^n I(x_i = x)$$

How many values are equal to the value.

For example, the following shows $\hat{p}(x)$ of size $n = 15$ from $Pn(5)$ (left) and the true pmf $p(x)$ of $Pn(5)$ (right)



Histograms

Histograms and smoothed pdfs

If the underlying variable is continuous we would prefer to obtain an approximation of the pdf. There are several approaches that can be used.

1. Histogram, \hat{f}_h (h is the bin length). First divide the entire range of values into a series of small intervals (bins) and then count how many values fall into each interval. For interval $[a, b)$, where $b - a = h$, draw a rectangle with height:

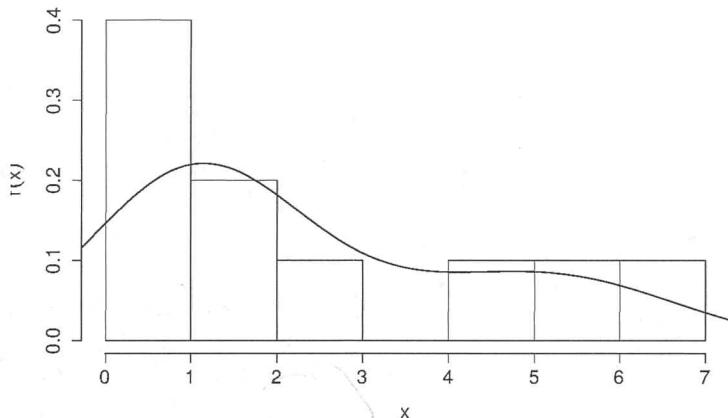
$$\hat{f}_h(x) = \frac{1}{hn} \sum_{i=1}^n I(a \leq x_i < b)$$

2. Smoothed pdf, \hat{f}_h (h is the 'bandwidth parameter'),

$$\hat{f}_h(x) = \frac{1}{hn} \sum_{i=1}^n K\left(\frac{x_i - x}{h}\right),$$

where $K(\cdot)$ is the kernel (a non-negative function that integrates to 1 and with mean zero) and h is a parameter that controls the level of smoothing.

Example: VEGFC

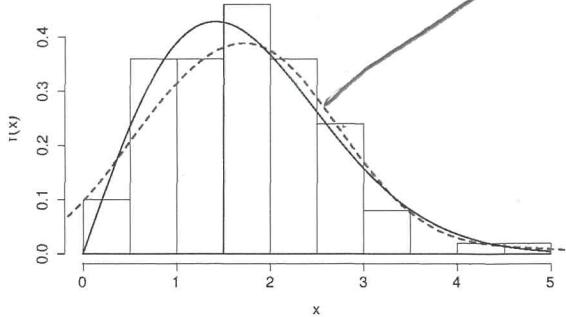


Simulated data

Consider $n = 100$ observations from the Weibull distribution with pdf

$$f(x) = \frac{1}{2}xe^{-(x/2)^2}, x > 0$$

True density (solid black curve), smoothed pdf (red dashed curve)



Quantile-quantile (QQ) plots

- For comparing the similarity of two probability distributions
- We plot their quantiles against each other (as a scatter plot)
- Typically, we compare data against a theoretical distribution
- The points in the plot are $(\hat{\pi}_p, \pi_p)$
- Note:** some people plot these the other way around: $(\pi_p, \hat{\pi}_p)$
- One axis shows the data, written here as sample quantiles:

$$\hat{\pi}_p = x_{(k)}, \text{ where } p = \frac{k}{n+1} \quad (\text{'Type 6' quantiles})$$

for $k = 1, \dots, n$

- Other axis shows corresponding quantiles for a theoretical distribution:

$$\pi_p = F^{-1}(p) = F^{-1}\left(\frac{k}{n+1}\right)$$

- The points in the plot therefore are,

$$\left\{ x_{(k)}, F^{-1}\left(\frac{k}{n+1}\right) \right\}.$$

Example: VEGFC

Is the sample from an exponential distribution with cdf $F(x) = 1 - e^{-\lambda x}$?

Sample quantiles (data):

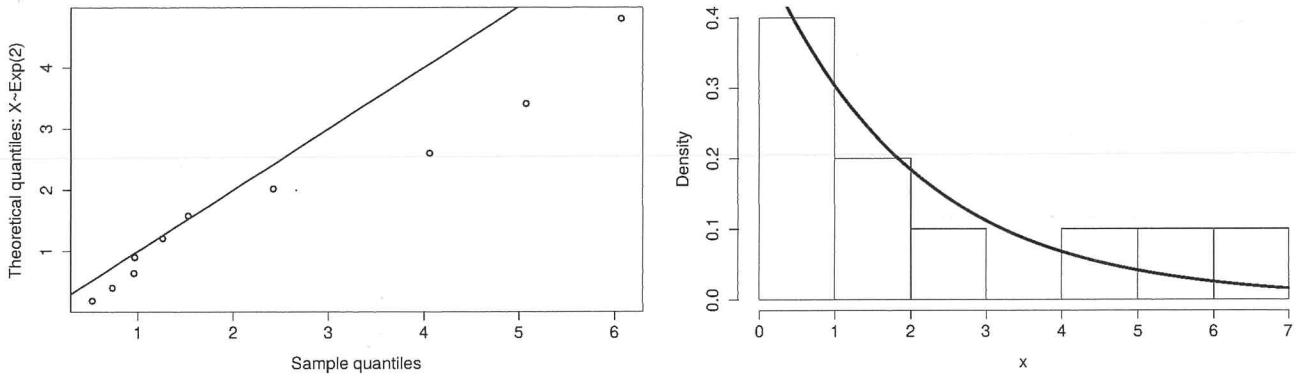
$$x_{(1)} = 0.52, x_{(2)} = 0.73, \dots, x_{(10)} = 6.07$$

Theoretical quantiles:

$$F^{-1}(p) = -\ln(1-p)/\lambda \quad (\text{e.g. set } \lambda = 0.5)$$

$$1/(10+1) = 0.09, 2/(10+1) = 0.18, \dots, 10/(10+1) = 0.91$$

$$F^{-1}(0.09) = 0.19, F^{-1}(0.18) = 0.40, \dots, F^{-1}(0.91) = 4.80$$



The right tail of the sample does not quite match the theoretical model (tail of the sample distribution is heavier).

Normal QQ plots

Normal QQ plots

If $X \sim N(\mu, \sigma^2)$, then $X = \mu + \sigma Z$, where $Z \sim N(0, 1)$. Therefore, if the normal model is correct

$$x_{(k)} \approx \mu + \sigma \Phi^{-1} \left(\frac{k}{n+1} \right)$$

where $\Phi(z) = P(Z \leq z)$ is the standard normal cdf.

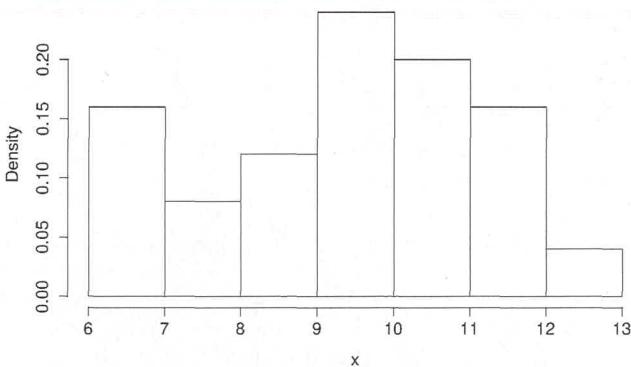
So, if we plot the points

$$\left(x_{(k)}, \Phi^{-1} \left(\frac{k}{n+1} \right) \right), \quad k = 1, \dots, n$$

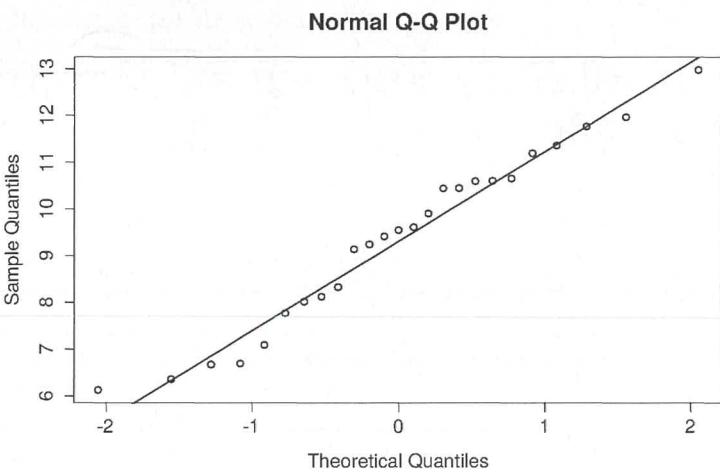
the result should be a straight line with intercept μ and slope σ . The values $\Phi^{-1}(k/(n+1))$ are called **normal scores**.

Example: simulated data

Consider 25 observations from $X \sim N(10, 2)$. The histogram is not very helpful:



But the QQ plot is much clearer:



Point estimation

(Module 2)

Statistics (MAST20005) & Elements of Statistics (MAST90058)

Contents

- 1 Estimation & sampling distributions
- 2 Estimators
- 3 Method of moments
- 4 Maximum likelihood estimation

Aims of this module

- Introduce the main elements of statistical inference and estimation, especially the idea of a sampling distribution $\bar{X} = \frac{1}{3}(X_1 + X_2 + X_3)$
- Show the simplest type of estimation: that of a single number
- Show some general approaches to estimation, especially the method of maximum likelihood

1 Estimation & sampling distributions

Motivating example

On a particular street, we measure the time interval (in minutes) between each car that passes:

2.55 2.13 3.18 5.94 2.29 2.41 8.72 3.71

We believe these follow an exponential distribution:

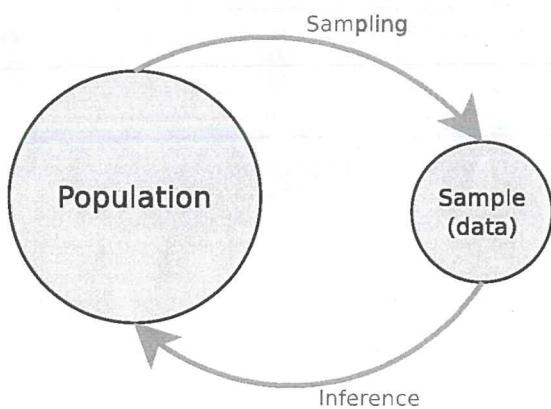
$$X_i \sim \text{Exp}(\lambda)$$

What can we say about λ ?

Can we approximate it from the data?

Yes! We can do it using a statistic. This is called estimation.

Statistics: the big picture



We want to start learning how to do inference. First, we need a good understanding of the 'sampling' part.

Distributions of statistics

Consider sampling from $X \sim \text{Exp}(\lambda = 1/5)$.

Convenient simplification: set $\theta = 1/\lambda$. This makes $\mathbb{E}(X) = \theta$ and $\text{var}(X) = \theta^2$.

Note: There are two common parameterisations,

$$f_X(x) = \lambda e^{-\lambda x}, \quad x \in [0, \infty)$$

$$f_X(x) = \frac{1}{\theta} e^{-\frac{1}{\theta} x}, \quad x \in [0, \infty)$$

λ is called the *rate parameter* (relates to a *Poisson process*)

Be clear about which is being used!

Take a large number of samples, each of size $n = 100$:

1.	1.84	1.19	11.73	5.64	17.98	0.26	...
2.	2.67	7.15	5.99	1.03	0.65	3.18	...
3.	16.99	2.15	2.60	5.40	3.64	2.01	...
4.	2.21	1.54	4.27	5.29	3.65	0.83	...
5.	12.24	1.59	2.56	1.38	5.72	0.69	...
							...

Then calculate some statistics (\bar{x} , $x_{(1)}$, $x_{(n)}$, etc.) for each one:

	Min.	Median	Mean	Max.
1.	0.02	4.10	5.17	23.96
2.	0.16	4.48	5.84	39.90
3.	0.17	3.39	4.38	15.61
4.	0.03	3.73	5.43	34.02
5.	0.01	3.12	4.71	19.94
				...

As we continue this process, we get some information on the distributions of these statistics.

Sampling distribution (definition)

Recall that any statistic $T = \phi(X_1, \dots, X_n)$ is a random variable.

The *sampling distribution* of a statistic is its probability distribution, given an assumed population distribution and a sampling scheme (e.g. random sampling).

Sometimes we can determine it exactly, but often we might resort to simulation.

In the current example, we know that:

$$\begin{aligned} X_{(1)} &\sim \text{Exp}(100\lambda) \\ \sum X_i &\sim \text{Gamma}(100, \lambda) \end{aligned}$$

How to estimate?

Suppose we want to estimate θ from the data. What should we do?

Reminder:

- Population mean, $\mathbb{E}(X) = \theta = 5$
- Population variance, $\text{var}(X) = \theta^2 = 5^2$
- Population standard deviation, $\text{sd}(X) = \theta = 5$

Can we use the sample mean, \bar{X} , as an estimate of θ ? Yes!

Can we use the sample standard deviation, S , as an estimate of θ ? Yes!

Will these statistics be good estimates? Which one is better? Let's see...

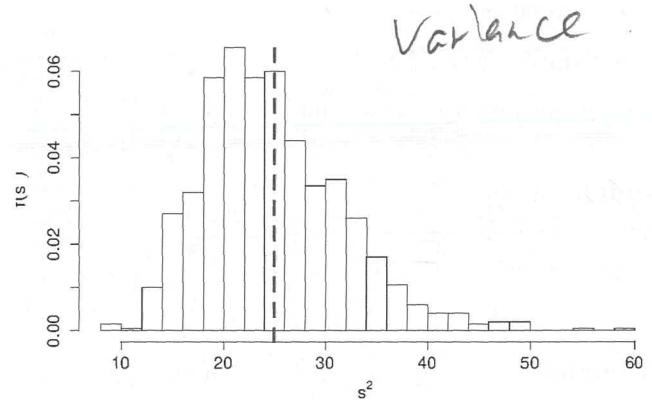
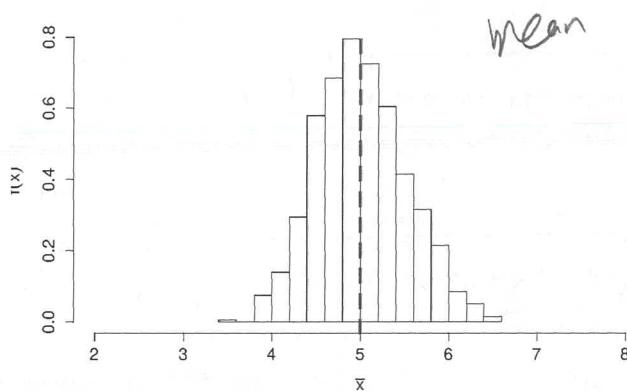
We need to know properties of their sampling distributions, such as their mean and variance.

Note: we are referring to the distribution of the statistic T rather than the population distribution from which we draw samples, X .

For example, it is natural to expect that:

- $E(\bar{X}) \approx \mu$ (sample mean \approx population mean)
- $E(S^2) \approx \sigma^2$ (sample variance \approx population variance)

Let's see for our example:



Left: distribution of \bar{X} . Right: distribution of S^2 . Vertical dashed lines: true values, $E(X) = 5$ and $\text{var}(X) = 5^2$.

- Should we use \bar{X} or S to estimate θ ? Which one is the better estimator?
- We would like the sample distribution of the estimator to be as close as possible to the true value $\theta = 5$.
- In practice, for any given dataset, we don't know which estimate is the closest, since we don't know the true value.
- We should use the one that is more likely to be the closest.
- Simulation: consider 250 samples of size $n = 100$ and compute:

$$\bar{x}_1, \dots, \bar{x}_{250},$$

$$s_1, \dots, s_{250}$$

```
> summary(x.bar)
  Min. 1st Qu. Median Mean 3rd Qu. Max.
3.789 4.663 4.972 5.015 5.365 6.424
> sd(x.bar)
[1] 0.4888185

> summary(s)
  Min. 1st Qu. Median Mean 3rd Qu. Max.
3.502 4.473 4.916 5.002 5.512 7.456
> sd(s)
[1] 0.7046119
```

From our simulation, $sd(\bar{X}) \approx 0.49$ and $sd(S) \approx 0.70$. So, in this case it looks like \bar{X} is superior to S .

2 Estimators

Definitions

- A parameter is a quantity that describes the population distribution, e.g. μ and σ^2 for $N(\mu, \sigma^2)$

- The parameter space is the set of all possible values that a parameter might take, e.g. $-\infty < \mu < \infty$ and $0 \leq \sigma < \infty$.
- An estimator (or point estimator) is a statistic that is used to estimate a parameter. It refers specifically to the random variable version of the statistic, e.g. $T = u(X_1, \dots, X_n)$.
- An estimate (or point estimate) is the observed value of the estimator for a given dataset. In other words, it is a realisation of the estimator, e.g. $t = u(x_1, \dots, x_n)$, where x_1, \dots, x_n is the observed sample (data).
- 'Hat' notation: If T is an estimator for θ , then we usually refer to it by $\hat{\theta}$ for convenience.

Examples

We will now go through a few important examples:

- Sample mean
- Sample variance
- Sample proportion

In each case, we assume a sample of iid rvs, X_1, \dots, X_n , with mean μ and variance σ^2 .

Sample mean

$$\bar{X} = \frac{1}{n} (X_1 + X_2 + \dots + X_n) = \frac{1}{n} \sum_{i=1}^n X_i$$

Properties:

- $\mathbb{E}(\bar{X}) = \mu$
- $\text{var}(\bar{X}) = \frac{\sigma^2}{n}$

Also, the Central Limit Theorem implies that usually:

$$\bar{X} \approx N\left(\mu, \frac{\sigma^2}{n}\right)$$

Often used to estimate the population mean, $\hat{\mu} = \bar{X}$.

Sample variance

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

Properties:

- $\mathbb{E}(S^2) = \sigma^2$
- $\text{var}(S^2) = (\text{a messy formula})$

Often used to estimate the population variance, $\hat{\sigma}^2 = S^2$.

Sample proportion

For a discrete random variable, we might be interested in how often a particular value appears. Counting this gives the sample frequency:

$$\text{freq}(a) = \sum_{i=1}^n I(X_i = a)$$

Let the population proportion be $p = \Pr(X = a)$. Then we have:

$$\text{freq}(a) \sim \text{Bi}(n, p)$$

Divide by the sample size to get the *sample proportion*. This is often used as an estimator for the population proportion:

$$\hat{p} = \frac{\text{freq}(a)}{n} = \frac{1}{n} \sum_{i=1}^n I(X_i = a)$$

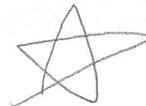
For large n , we can approximate this with a *normal distribution*:

$$\hat{p} \approx N\left(p, \frac{p(1-p)}{n}\right)$$

Note:

- The sample pmf and the sample proportion are the same, both of them estimate the probability of a given event or set of events.
- The pmf is usually used when the interest is in many different events/values, and is written as a function, e.g. $\hat{p}(a)$.
- The proportion is usually used when only a single event is of interest (getting heads for a coin flip, a certain candidate winning an election, etc.).

Examples for a normal distribution



If the sample is drawn from a normal distribution, $X_i \sim N(\mu, \sigma^2)$, we can derive *exact* distributions for these statistics.

Sample mean:

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

$$S \sim (\bar{X}, \frac{\sigma^2}{n})$$

Sample variance:

$$S^2 \sim \frac{\sigma^2}{n-1} \chi_{n-1}^2$$

$$\mathbb{E}(S^2) = \sigma^2, \quad \text{var}(S^2) = \frac{2\sigma^4}{n-1}$$



χ_k^2 is the chi-squared distribution with k degrees of freedom. (more details in Module 3)

Bias

Consider an estimator $\hat{\theta}$ of θ .

- If $\mathbb{E}(\hat{\theta}) = \theta$, the estimator is said to be *unbiased*.
- The *bias* of the estimator is, $\mathbb{E}(\hat{\theta}) - \theta$

Examples:

- The sample variance is unbiased for the population variance, $\mathbb{E}(S^2) = \sigma^2$. (problem 5 in week 3 tutorial)
- What if we divide by n instead of $n-1$ in the denominator?

Transformations and biasedness

$$\mathbb{E}\left(\frac{n-1}{n} S^2\right) = \frac{n-1}{n} \sigma^2 < \sigma^2$$

⇒ biased!

In general, if $\hat{\theta}$ is unbiased for θ , then it will usually be the case that $g(\hat{\theta})$ is biased for $g(\theta)$.

Unbiasedness is not preserved under transformations.

Challenge problem

Is the sample standard deviation, $S = \sqrt{S^2}$, biased for the population standard deviation σ ?

$$\mathbb{E}(S^2) = \sigma^2$$

$$\sqrt{\mathbb{E}(S^2)} = \sigma$$

$$\mathbb{E}(S) = ?$$

$$\begin{aligned} \text{Var}(S) &= \mathbb{E}(S^2) - \mathbb{E}(S)^2 > 0 \quad \downarrow \\ 5 \Rightarrow \mathbb{E}(S^2) &> \mathbb{E}(S)^2 \\ \sigma^2 &= \sqrt{\mathbb{E}(S^2)} > \mathbb{E}(S) \Rightarrow \sigma > \mathbb{E}(S). \end{aligned}$$

biased

Choosing between estimators

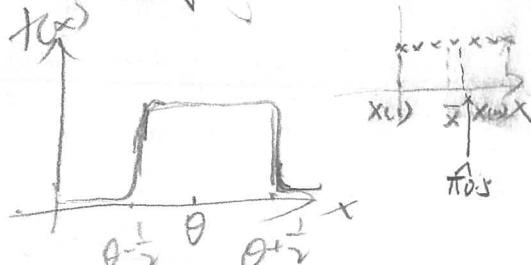
- Evaluate and compare the sampling distributions of the estimators.
- Generally, prefer estimators that have **smaller bias** and **smaller variance** (and it can vary depending on the aim of your problem).
- Sometimes, we only know asymptotic properties of estimators (will see examples later).

Note: this approach to estimation is referred to as *frequentist* or *classical* inference. The same is true for most of the techniques we will cover. We will also learn about an alternative approach, called *Bayesian* inference, later in the semester.

Challenge problem (uniform distribution)

Take a random sample of size n from the uniform distribution with pdf:

$$f(x) = 1 \quad \left(\theta - \frac{1}{2} < x < \theta + \frac{1}{2} \right)$$



Can you think of some estimators for θ ? What is their bias and variance?

Challenge problem (boundary problem)

Take a random sample of size n from the shifted exponential distribution, with pdf:

$$f(x) = e^{-(x-\theta)} \quad (x > \theta)$$

Equivalently:

$$X_i \sim \theta + \text{Exp}(1)$$

Can you think of some estimators for θ ? What is their bias and variance?

Coming up with (good) estimators?

How can we do this for any given problem?

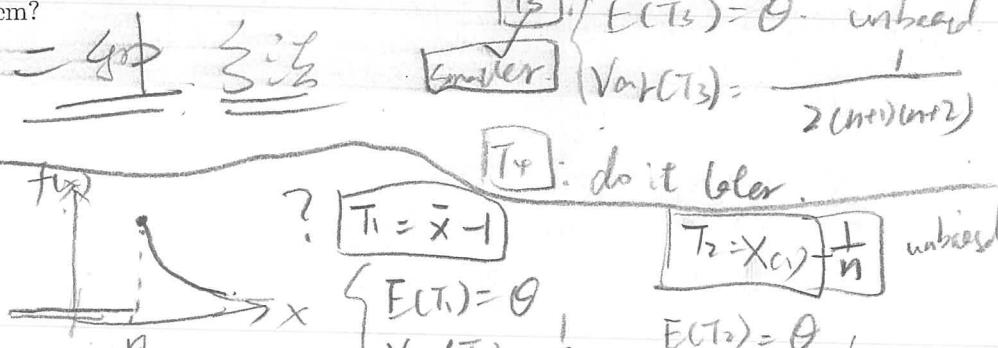
We will cover two general methods:

- Method of moments
- Maximum likelihood

3 Method of moments

Method of moments (MM)

- Idea:
 - Make the population distribution resemble the empirical (data) distribution...
 - ... by equating theoretical moments with sample moments
 - Do this until you have enough equations, and then solve them
- Example: if $E(\bar{X}) = \theta$, then the method of moments estimator of θ is \bar{X} .
- General procedure (for r parameters):
 - X_1, \dots, X_n i.i.d. $f(x | \theta_1, \dots, \theta_r)$.
 - k th moment is $\mu_k = E(X^k)$
 - k th sample moment is $M_k = \frac{1}{n} \sum X_i^k$
 - Set $\mu_k = M_k$, for $k = 1, \dots, r$ and solve for $(\theta_1, \dots, \theta_r)$.
- Alternative: Can use the variance instead of the second moment (sometimes more convenient).



$$E(X^K) \approx \frac{1}{n} \sum X^K$$

Remarks

- An intuitive approach to estimation
- Can work in situations where other approaches are too difficult
- Usually biased
- Usually not optimal (but may suffice)
- Note: some authors use a 'bar' ($\bar{\theta}$) or a 'tilde' ($\tilde{\theta}$) to denote MM estimators rather than a 'hat' ($\hat{\theta}$). This helps to distinguish different estimators when comparing them to each other.

Example: Geometric distribution

- Sampling from: $X \sim \text{Geom}(p)$

- The first moment:

$$E(X) = \sum_{x=1}^{\infty} xp(1-p)^{x-1} = \frac{1}{p}$$

- The MM estimator is obtained by solving

$$\bar{X} = \frac{1}{p}$$

which gives

$$\tilde{p} = \frac{1}{\bar{X}}$$

Example: Normal distribution

- Sampling from: $X \sim N(\mu, \sigma^2)$
- Population moments: $E(X) = \mu$ and $E(X^2) = \sigma^2 + \mu^2$
- Sample moments: $M_1 = \bar{X}$ and $M_2 = \frac{1}{n} \sum X_i^2$
- Equating them:

$$\bar{X} = \mu \quad \text{and} \quad \frac{1}{n} \sum X_i^2 = \sigma^2 + \mu^2$$

Solving these gives:

$$\tilde{\mu} = \bar{X} \quad \text{and} \quad \tilde{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$$

Note:

- This is not the usual sample variance!
- $\tilde{\sigma}^2 = \frac{n-1}{n} S^2$
- This one is biased, $E(\tilde{\sigma}^2) = \frac{n-1}{n} \sigma^2 \neq \sigma^2$.

Example: Gamma distribution

- Sampling from: $X \sim \text{Gamma}(\alpha, \theta)$
- The pdf is:

$$f(x | \alpha, \theta) = \frac{1}{\Gamma(\alpha)\theta^\alpha} x^{\alpha-1} \exp\left(\frac{-x}{\theta}\right)$$

- Population moments: $E(X) = \alpha\theta$ and $\text{var}(X) = \alpha\theta^2$
- Sample moments: $M = \bar{X}$ and $S^2 = \frac{1}{n-1} \sum (X_i - \bar{X})^2$

Equating them:

$$\bar{X} = \alpha\theta \quad \text{and} \quad S^2 = \alpha\theta^2$$

Solving these gives:

$$\tilde{\theta} = \frac{S^2}{\bar{X}} \quad \text{and} \quad \tilde{\alpha} = \frac{\bar{X}^2}{S^2}$$

Note:

- This is an example of using S^2 instead of M_2

4 Maximum likelihood estimation

Method of maximum likelihood (ML)

- Idea: find the 'most likely' explanation for the data
- More concretely: find parameter values that maximise the probability of the data

Example: Bernoulli distribution

- Sampling from: $X \sim \text{Be}(p)$
- Data are 0's and 1's
- Then pmf is

$$f(x | p) = p^x(1-p)^{1-x}, \quad x = 0, 1, \quad 0 \leq p \leq 1$$

- Observe values x_1, \dots, x_n of X_1, \dots, X_n (iid)
- The probability of the data (the random sample) is

$$\Pr(X_1 = x_1, \dots, X_n = x_n | p) = \prod_{i=1}^n f(x_i | p) = \prod_{i=1}^n p^{x_i}(1-p)^{1-x_i}$$

- Regard the sample x_1, \dots, x_n as known (since we have observed it) and regard the probability of the data as a function of p .
- When written this way, this is called the likelihood of p :

$$\begin{aligned} L(p) &= L(p | x_1, \dots, x_n) \\ &= \Pr(X_1 = x_1, \dots, X_n = x_n | p) \\ &= p^{\sum x_i} (1-p)^{n-\sum x_i} \end{aligned}$$

① pdf \Rightarrow likelihood function

② $\ln L(p)$

③ $\partial \ln L(p) / \partial p = 0$

④ \Rightarrow answer

- Want to find the value of p that maximizes this likelihood.

- It often helps to find the value of θ that maximizes the log of the likelihood rather than the likelihood

- This is called the log-likelihood

$$\ln L(p) = \ln p^{\sum x_i} + \ln(1-p)^{n-\sum x_i}$$

- The final answer (the maximising value of p) is the same, since the log of non-negative numbers is a one-to-one function whose inverse is the exponential, so any value θ that maximises the log-likelihood also maximises the likelihood.

- Putting $x = \sum_{i=1}^n x_i$ so that x is the number of 1's in the sample,

$$\ln L(p) = x \ln p + (n-x) \ln(1-p)$$

- Find the maximum of this log-likelihood with respect to p by differentiating and equating to zero,

$$\frac{\partial \ln L(p)}{\partial p} = x \frac{1}{p} + (n-x) \frac{-1}{1-p} = 0$$

- This gives $p = x/n$

- Therefore, the maximum likelihood estimator is $\hat{p} = X/n = \bar{X}$

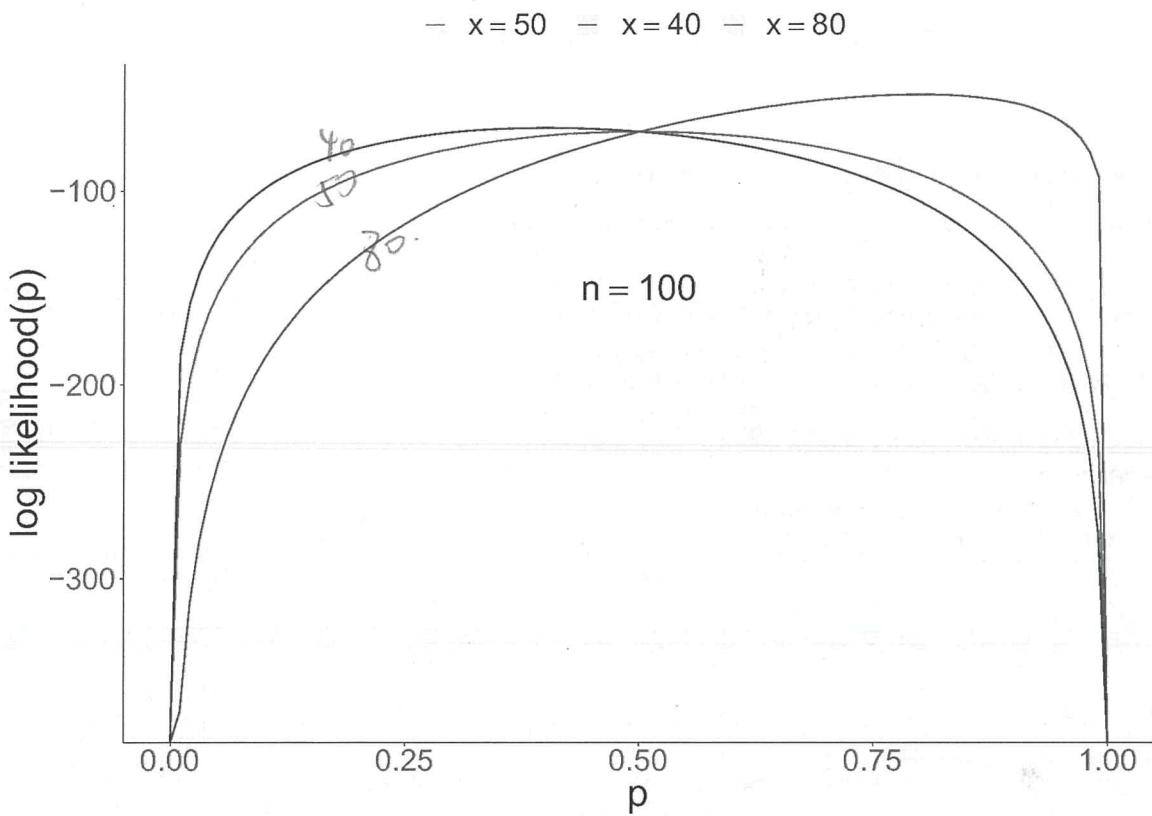


Figure 1: Log-likelihoods for Bernoulli trials with parameter p

Maximum likelihood: general procedure

- Random sample (iid): X_1, \dots, X_n
- Likelihood function with m parameters $\theta_1, \dots, \theta_m$ and data x_1, \dots, x_n is:

$$L(\theta_1, \dots, \theta_m) = \prod_{i=1}^n f(x_i | \theta_1, \dots, \theta_m)$$

- $\frac{1}{\lambda^n}$

- If X is discrete, for f use the pmf
- If X is continuous, for f use the pdf
- The maximum likelihood estimates (MLEs) or the maximum likelihood estimators (MLEs) $\hat{\theta}_1, \dots, \hat{\theta}_m$ are values that maximize $L(\theta_1, \dots, \theta_m)$.
- Note: same abbreviation and notation for both the estimators (random variable) and the estimates (realised values).
- Often (but not always) useful to take logs and then differentiate and equate derivatives to zero to find MLE's.
- Sometimes this is too hard, but we can maximise numerically. No closed-form expression in this case.

Example: Exponential distribution

Sampling (iid) from: $X \sim \text{Exp}(\lambda)$

$$f(x | \lambda) = \frac{1}{\lambda} e^{-x/\lambda}, \quad x > 0, \quad 0 < \lambda < \infty$$

$$L(\lambda) = \frac{1}{\lambda^n} \exp\left(-\frac{\sum_{i=1}^n x_i}{\lambda}\right)$$

$$\ln L(\lambda) = -n \ln(\lambda) - \frac{1}{\lambda} \sum_{i=1}^n x_i$$

$$\frac{\partial \ln L(\lambda)}{\partial \lambda} = -n \frac{1}{\lambda} + \sum x_i \frac{1}{\lambda^2} = 0$$

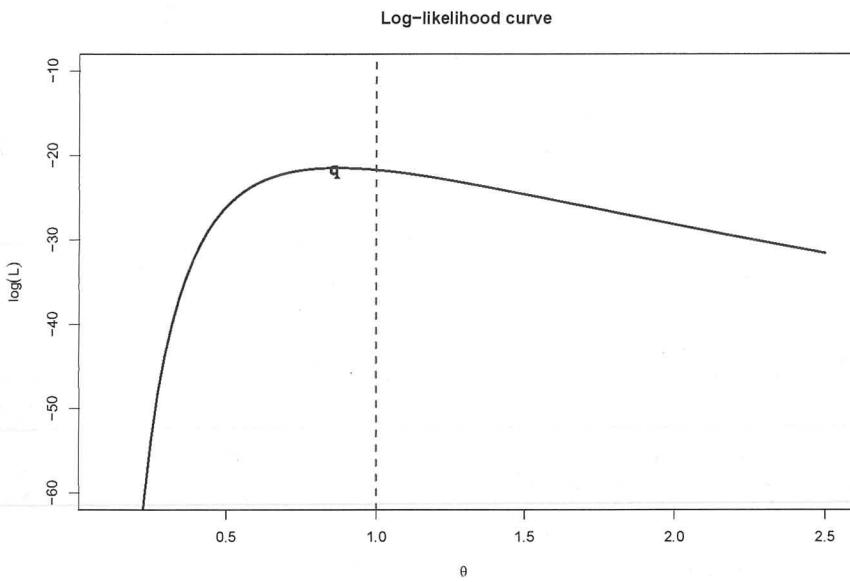
$$\frac{n}{\lambda} = \frac{\sum x_i}{\lambda^2} \Rightarrow \lambda = \frac{\sum x_i}{n}$$

$$\frac{\partial \ln L(\lambda)}{\partial \lambda} = -\frac{n}{\lambda} + \frac{\sum x_i}{\lambda^2} = 0$$

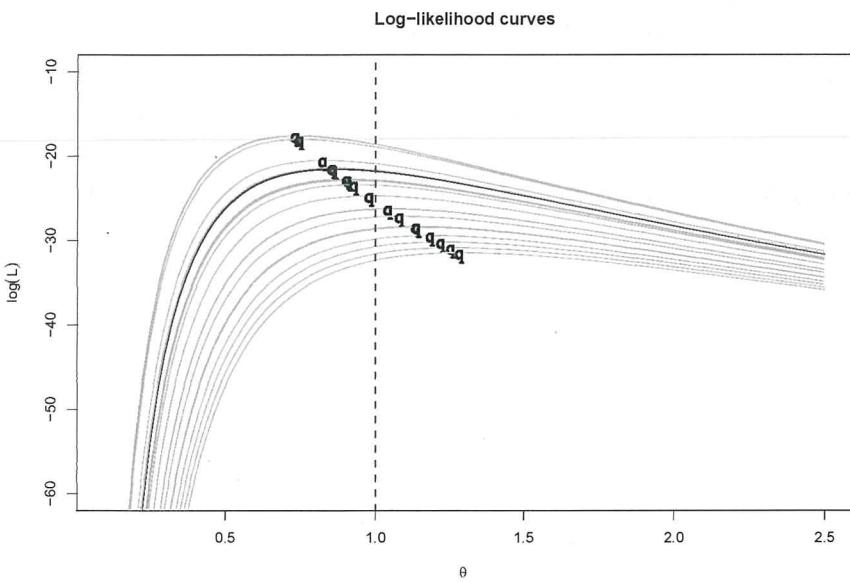
This gives: $\hat{\lambda} = \bar{X}$

Example: Exponential distribution (simulated)

```
> x <- rexp(25) # simulate 25 observations from Exp(1)
> x
[1] 0.009669867 3.842141708 0.394267770 0.098725403
[5] 0.386704987 0.024086824 0.274132718 0.872771164
[9] 0.950139285 0.022927997 1.538592014 0.837613769
[13] 0.634363088 0.494441270 1.789416017 0.503498224
[17] 0.000482703 1.617899321 0.336797648 0.312564298
[21] 0.702562098 0.265119483 3.825238461 0.238687987
[25] 1.752657238
> mean(x) # maximum likelihood estimate
[1] 0.8690201
```



What if we repeat the sampling process several times?



Example: Geometric distribution

Sampling (iid) from: $X \sim \text{Geom}(p)$

$$L(p) = \prod_{i=1}^n p(1-p)^{x_i-1} = p^n(1-p)^{\sum x_i - n}, \quad 0 \leq p \leq 1$$

$$\frac{\partial \ln L(p)}{\partial p} = \frac{n}{p} - \frac{\sum_{i=1}^n x_i - n}{1-p} = 0$$

This gives: $\hat{p} = 1/\bar{X}$

Example: Normal distribution

Sampling (iid) from: $X \sim N(\theta_1, \theta_2)$

two parameters,

$$L(\theta_1, \theta_2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\theta_2}} \exp\left[-\frac{(x_i - \theta_1)^2}{2\theta_2}\right]$$

$$\ln L(\theta_1, \theta_2) = -\frac{n}{2} \ln(2\pi\theta_2) - \frac{1}{2\theta_2} \sum_{i=1}^n (x_i - \theta_1)^2$$

Take partial derivatives with respect to θ_1 and θ_2 .

$$\frac{\partial \ln L(\theta_1, \theta_2)}{\partial \theta_1} = \frac{1}{\theta_2} \sum_{i=1}^n (x_i - \theta_1)$$

$$\frac{\partial \ln L(\theta_1, \theta_2)}{\partial \theta_2} = -\frac{n}{2\theta_2} + \frac{1}{2\theta_2^2} \sum_{i=1}^n (x_i - \theta_1)^2$$

Set both of these to zero and solve. This gives: $\hat{\theta}_1 = \bar{x}$ and $\hat{\theta}_2 = n^{-1} \sum_{i=1}^n (x_i - \bar{x})^2$. The maximum likelihood estimators are therefore:

$$\hat{\theta}_1 = \bar{X}, \quad \hat{\theta}_2 = \underbrace{\frac{1}{n} \sum_{i=1}^n}_{\sim} (X_i - \bar{X})^2 = \underbrace{\frac{n-1}{n}}_{\sim} S^2$$

Note: $\hat{\theta}_2$ is biased.

Stress and cancer: VEGFC

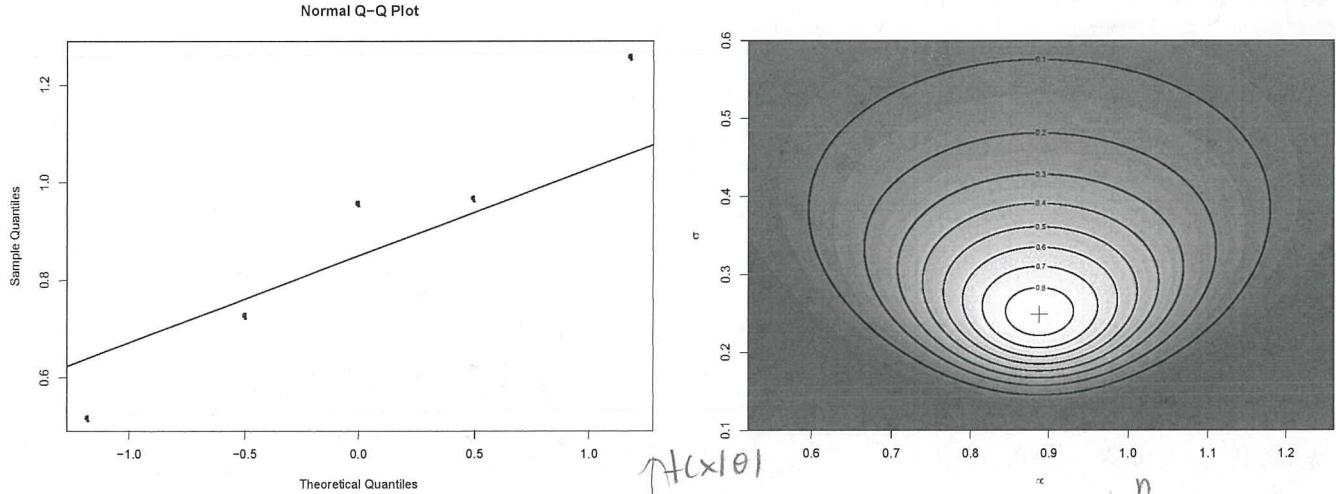
```

> x <- c(0.97, 0.52, 0.73, 0.96, 1.26)
> n <- length(x)
> mean(x) # MLE for population mean
[1] 0.888

> sd(x) * sqrt((n - 1) / n) # MLE for the pop. st. dev.
[1] 0.2492709

> qqnorm(x) # Draw a QQ plot
> qqline(x) # Fit line to QQ plot

```



Challenge problem (boundary problem)

Take a random sample of size n from the shifted exponential distribution, with pdf:

$$f(x | \theta) = e^{-(x-\theta)} \quad (x > \theta)$$

Equivalently:

$$X_i \sim \theta + \text{Exp}(1)$$

Derive the MLE for θ . Is it biased? Can you create an unbiased estimator from it?

Invariance property

Suppose we know $\hat{\theta}$ but are actually interested in $\phi = g(\theta)$ rather than θ itself. Can we estimate ϕ ?

Yes! It is simply $\hat{\phi} = g(\hat{\theta})$.

This is known as the *invariance property of the MLE*. In other words, transformations don't affect the value of the MLE.

Consequence: MLEs are usually biased since expectations are **not** invariant under transformations.

Is the MLE a good estimator?

Some useful results:

- Asymptotically unbiased
- Asymptotically optimal variance ('efficient')
- Asymptotically normally distributed

The proofs of these rely on the CLT. More details of the mathematical theory will be covered towards the end of the semester.