Big Data Analysis (MA60306)

Bibhas Adhikari

Spring 2022-23, IIT Kharagpur

Lecture 3 January 6, 2022

3

・ 何 ト ・ ヨ ト ・ ヨ ト

Describing Data Terminologies

▷ Population, sample, observations - The units on which we measure data - such as persons, cars, animals are called observations. The collection of all units is called population which is denoted by Ω . A selection of observations $\omega_1, \ldots, \omega_n \in \Omega$ is called a sample.

Terminologies

- ▷ Population, sample, observations The units on which we measure data such as persons, cars, animals are called observations. The collection of all units is called population which is denoted by Ω . A selection of observations $\omega_1, \ldots, \omega_n \in \Omega$ is called a sample.
- ▷ Variable (random/statistical) qualitative and quantitative variables
 - Discrete variables
 - ▷ Continuous variables (measured rather than counted)

∃ ► < ∃ ►</p>

Terminologies

- ▷ Population, sample, observations The units on which we measure data such as persons, cars, animals are called observations. The collection of all units is called population which is denoted by Ω . A selection of observations $\omega_1, \ldots, \omega_n \in \Omega$ is called a sample.
- ▷ Variable (random/statistical) qualitative and quantitative variables
 - Discrete variables
 - ▷ Continuous variables (measured rather than counted)
- ▷ Scales different variables contain different amounts of information

Terminologies

- \triangleright Population, sample, observations The units on which we measure data such as persons, cars, animals are called observations. The collection of all units is called population which is denoted by Ω . A selection of observations $\omega_1, \ldots, \omega_n \in \Omega$ is called a sample.
- \triangleright Variable (random/statistical) qualitative and quantitative variables
 - Discrete variables
 - ▷ Continuous variables (measured rather than counted)
- ▷ Scales different variables contain different amounts of information
 - ▷ Nominal scale values can not be ordered

∃ ► < ∃ ►</p>

Terminologies

- ▷ Population, sample, observations The units on which we measure data such as persons, cars, animals are called observations. The collection of all units is called population which is denoted by Ω . A selection of observations $\omega_1, \ldots, \omega_n \in \Omega$ is called a sample.
- \triangleright Variable (random/statistical) qualitative and quantitative variables
 - Discrete variables
 - ▷ Continuous variables (measured rather than counted)
- ▷ Scales different variables contain different amounts of information
 - ▷ Nominal scale values can not be ordered
 - Ordinal scale can be ordered but cannot be interpreted/compared in a meaningful/numerical way

A B A A B A

Terminologies

- ▷ Population, sample, observations The units on which we measure data such as persons, cars, animals are called observations. The collection of all units is called population which is denoted by Ω . A selection of observations $\omega_1, \ldots, \omega_n \in \Omega$ is called a sample.
- $\triangleright~$ Variable (random/statistical) qualitative and quantitative variables
 - Discrete variables
 - ▷ Continuous variables (measured rather than counted)
- ▷ Scales different variables contain different amounts of information
 - ▷ Nominal scale values can not be ordered
 - Ordinal scale can be ordered but cannot be interpreted/compared in a meaningful/numerical way
 - Continuous scale can be ordered but can be interpreted/compared in a meaningful/numerical way

A B b A B b

Terminologies

- ▷ Population, sample, observations The units on which we measure data such as persons, cars, animals are called observations. The collection of all units is called population which is denoted by Ω . A selection of observations $\omega_1, \ldots, \omega_n \in \Omega$ is called a sample.
- ▷ Variable (random/statistical) qualitative and quantitative variables
 - Discrete variables
 - ▷ Continuous variables (measured rather than counted)
- ▷ Scales different variables contain different amounts of information
 - ▷ Nominal scale values can not be ordered
 - Ordinal scale can be ordered but cannot be interpreted/compared in a meaningful/numerical way
 - Continuous scale can be ordered but can be interpreted/compared in a meaningful/numerical way
- Grouped Data (categorical variable) instead of the original value, one may only know the category or group the value belongs to

Summary of variable classifications



э

<ロト <問ト < 目と < 目と

 \triangleright smallest unit of storage, a 0 or 1

(日) (四) (日) (日) (日)

2

- $\,\triangleright\,$ smallest unit of storage, a 0 or 1
- ▷ anything that can store two states now "transistors", used to be vacuum tubes
- \triangleright with *n* bits, can store 2^n patterns so one byte can store 256 patterns

- \triangleright smallest unit of storage, a 0 or 1
- ▷ anything that can store two states now "transistors", used to be vacuum tubes
- \triangleright with *n* bits, can store 2^n patterns so one byte can store 256 patterns
- \triangleright eight bits = one byte

- \triangleright smallest unit of storage, a 0 or 1
- ▷ anything that can store two states now "transistors", used to be vacuum tubes
- \triangleright with *n* bits, can store 2^n patterns so one byte can store 256 patterns
- \triangleright eight bits = one byte

The encoding paradigm: Here

Multi-byte units

unit	abbreviation	total bytes	nearest decimal equivalent
kilobyte	KB	1,024 ¹	1000 ¹
megabyte	MB	$1,024^2$	1000 ²
gigabyte	GB	1,024 ³	1000 ³
terabyte	ТВ	1,024 ⁴	1000 ⁴
petabyte	PB	1,024 ⁵	1000 ⁵
exabyte	EB	1,024 ⁶	1000 ⁶
zettabyte	ZB	1,024 ⁷	1000 ⁷
yottabyte	YB	1,024 ⁸	1000 ⁸

æ

A D N A B N A B N A B N

Multi-byte units

unit	abbreviation	total bytes	nearest decimal equivalent
kilobyte	KB	$1,024^{1}$	1000 ¹
megabyte	MB	$1,024^2$	1000 ²
gigabyte	GB	1,024 ³	1000 ³
terabyte	ТВ	$1,024^{4}$	1000 ⁴
petabyte	PB	1,024 ⁵	1000 ⁵
exabyte	EB	1,024 ⁶	1000 ⁶
zettabyte	ZB	1,024 ⁷	1000 ⁷
yottabyte	YB	1,024 ⁸	1000 ⁸

Observation this is why 1GB is greater than 1 billion bytes

3

Data Collection - Primary and secondary data - gather data on a sample and draw conclusions about the population of interest

Survey - collecting data by asking questions

Data Collection - Primary and secondary data - gather data on a sample and draw conclusions about the population of interest

- ▷ Survey collecting data by asking questions
- Experiment generated by the researcher with full control over one or many variables of interest

Data Collection - Primary and secondary data - gather data on a sample and draw conclusions about the population of interest

- ▷ Survey collecting data by asking questions
- Experiment generated by the researcher with full control over one or many variables of interest
- ▷ Observational without a survey or conducting an experiment

Creating a data set - the data is stored in a data matrix (= data set) with p columns and n rows

3

(日) (四) (日) (日) (日)

Creating a data set - the data is stored in a data matrix (= data set) with p columns and n rows

ω	Variable 1	Variable 2		<i>Variable</i> p	
1	Γ × ₁₁	×12		x_{1p}	٦
2	x ₂₁	x ₂₂		× _{2p}	
÷	:	:	÷	÷	
n	$ L x_{n1} $	× _{n2}		Х _{пр}	l

3

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 >

Creating a data set - the data is stored in a data matrix (= data set) with p columns and n rows

ω	Variable 1	Variable 2		Variable	р
1	Γ × ₁₁	x ₁₂		x_{1p}	٦
2	x ₂₁	x ₂₂		x _{2p}	
÷	:	:	÷	:	
n	$ L x_{n1} $	× _{n2}		× _{np}	

A statistical convention: Variables are are aften called *features*- the columns

3

(日) (四) (日) (日) (日)

Dataset manipulation

- A dataset is a "rectangular" formatted table of data in which all the values of the same variable must be in a single column
- ▷ A dataset is not:
 - $\triangledown\,$ the results of tabulating a dataset
 - $\triangledown\,$ any set of summary statistics on a dataset
 - $\triangledown\,$ a series of relational tables
- Many of the datasets we use have been artificially reshaped in order to fulfill this criterion of rectangularity

The first column is usually an identifier or index column, where each unit/row is given a unique name or ID

¹P. Lafaye de Micheaux, R. Drouilhet, and B. Liquet. The R Software: Fundamentals of Programming and Statistical Analysis. Springer, New York, 2014.

- The first column is usually an identifier or index column, where each unit/row is given a unique name or ID
- Certain columns (features) can correspond to the design of the experiment, specifying, for example, to which experimental group the unit belongs

¹P. Lafaye de Micheaux, R. Drouilhet, and B. Liquet. The R Software: Fundamentals of Programming and Statistical Analysis. Springer, New York, 2014.

- ▷ The first column is usually an identifier or index column, where each unit/row is given a unique name or ID
- Certain columns (features) can correspond to the design of the experiment, specifying, for example, to which experimental group the unit belongs
- Other columns represent the observed measurements of the experiment. Usually, these measurements exhibit variability; that is, they would change if the experiment were to be repeated

¹P. Lafaye de Micheaux, R. Drouilhet, and B. Liquet. The R Software: Fundamentals of Programming and Statistical Analysis. Springer, New York, 2014.

- The first column is usually an identifier or index column, where each unit/row is given a unique name or ID
- Certain columns (features) can correspond to the design of the experiment, specifying, for example, to which experimental group the unit belongs
- Other columns represent the observed measurements of the experiment. Usually, these measurements exhibit variability; that is, they would change if the experiment were to be repeated

A well-known repository of data sets is the Machine Learning Repository maintained by the University of California at Irvine (UCI), found at HERE

¹P. Lafaye de Micheaux, R. Drouilhet, and B. Liquet. The R Software: Fundamentals of Programming and Statistical Analysis. Springer, New York, 2014.

- The first column is usually an identifier or index column, where each unit/row is given a unique name or ID
- Certain columns (features) can correspond to the design of the experiment, specifying, for example, to which experimental group the unit belongs
- Other columns represent the observed measurements of the experiment. Usually, these measurements exhibit variability; that is, they would change if the experiment were to be repeated

A well-known repository of data sets is the Machine Learning Repository maintained by the University of California at Irvine (UCI), found at HERE

- ⊳ Iris.data
- ▷ abalone.data

\triangleright nutrition_elderly.xls¹

¹P. Lafaye de Micheaux, R. Drouilhet, and B. Liquet. The R Software: Fundamentals of Programming and Statistical Analysis. Springer, New York, 2014.

Tidy data - Hadley Wickham's three rules for "tidy" datasets ▷ Each variable must have its own column

3

(日) (四) (日) (日) (日)

Tidy data - Hadley Wickham's three rules for "tidy" datasets

- Each variable must have its own column
- Each observation must have its own row

э

▲ 東 ▶ | ▲ 更 ▶

Tidy data - Hadley Wickham's three rules for "tidy" datasets

- Each variable must have its own column
- Each observation must have its own row
- Each value must have its own cell

э

4 15 16 16 15

Tidy data - Hadley Wickham's three rules for "tidy" datasets

- Each variable must have its own column
- $\triangleright\,$ Each observation must have its own row
- ▷ Each value must have its own cell

Read Tidy data

3

Tidy data - Hadley Wickham's three rules for "tidy" datasets

- Each variable must have its own column
- Each observation must have its own row
- Each value must have its own cell

Read Tidy data



Question Change the xls file into CSV file

Let $x = [x_1, \ldots, x_n]^T$ be a column vector with *n* numbers.

イロト 不得 トイヨト イヨト 二日

Let $x = [x_1, ..., x_n]^T$ be a column vector with *n* numbers. \triangleright sample mean:

$$\overline{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$$

イロト 不得 トイヨト イヨト 二日

Let $x = [x_1, \ldots, x_n]^T$ be a column vector with *n* numbers.

▷ sample mean:

$$\overline{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$$

▷ sample quantile: The *p*-sample quantile (0 of*x* $is a value <math>x_p$ such that at least a fraction *p* of the data is less than or equal to x_p and at least a fraction 1 - p of the data is greater than or equal to x_p

Let $x = [x_1, \ldots, x_n]^T$ be a column vector with *n* numbers.

⊳ sample mean:

$$\overline{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$$

- ▷ sample quantile: The *p*-sample quantile (0 of*x* $is a value <math>x_p$ such that at least a fraction *p* of the data is less than or equal to x_p and at least a fraction 1 p of the data is greater than or equal to x_p
- ▷ The *sample median* is the sample 0.5-quantile

<日

<</p>

Let $x = [x_1, \ldots, x_n]^T$ be a column vector with *n* numbers.

⊳ sample mean:

$$\overline{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$$

- ▷ sample quantile: The *p*-sample quantile (0 of*x* $is a value <math>x_p$ such that at least a fraction *p* of the data is less than or equal to x_p and at least a fraction 1 p of the data is greater than or equal to x_p
- ▷ The *sample median* is the sample 0.5-quantile
- \triangleright The p-sample quantile is also called the $100 \times p$ percentile. The 25, 50, and 75 sample percentiles are called the *first, second, and third quartiles* of the data.

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

dispersion (spread) of the data is measured by

▷ sample range

$$\max_{i} x_{i} - \min_{i} x_{i}$$

▷ sample variance

$$s^{2} = \frac{1}{n-1} \sum_{i=1}^{n} (x_{i} - \overline{x})^{2}$$

3

・ 何 ト ・ ヨ ト ・ ヨ ト

dispersion (spread) of the data is measured by

▷ sample range

 $\max_i x_i - \min_i x_i$

▷ sample variance

$$s^{2} = \frac{1}{n-1} \sum_{i=1}^{n} (x_{i} - \overline{x})^{2}$$

 \triangleright sample standard deviation $s = \sqrt{s^2}$

3

A B < A B </p>

dispersion (spread) of the data is measured by

▷ sample range

$$\max_i x_i - \min_i x_i$$

▷ sample variance

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \overline{x})^2$$

 \triangleright sample standard deviation $s = \sqrt{s^2}$

Question Calculate the above measures for a feature of the nutrition data

E 6 4 E 6

Visualizing data

- Plotting Qualitative Variables Barplot
- Plotting Quantitative Variables -
 - (a) Boxplot a graphical representation of the five-number summary of the data consisting of the minimum, maximum, and the first, second, and third quartiles
 - (b) Histogram a graphical representation of the distribution of a quantitative feature
 - (c) Empirical Cumulative Distribution Function denoted by F_n , is a step function which jumps an amount k/n at observation values, where k is the number of tied observations at that value. For observations x_1, \ldots, x_n ,

$$F_n(x) = \frac{\text{number of } x_i \leq x}{n} = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{x_i \leq x\}$$

where $\mathbbm{1}$ denotes the indicator function