

Python (part 2)



KHOA CÔNG NGHỆ THÔNG TIN
TRƯỜNG ĐẠI HỌC KHOA HỌC TỰ NHIÊN

fit@hcmus

Why Jupyter notebook

- Today
 - ▣ More about data science process
 - ▣ Demo: using Python to do a data science process

Data science process

Data science process

- ☐ Ask a meaningful question
- ☐ Collect data
- ☐ Explore data
- ☐ Preprocess data
- ☐ Analyze data
 - the answer
- ☐ Communicate results / make decision

What is the purpose of exploring data?

- ☐ To **understand more about data**

- ☐ From that, we can:

 - ☐ Identify problems in data

If there is a problem, we might need to:

 - Preprocess it (we might need to preprocess it right away in order to continue to explore data)
 - Or even go back to collecting data

 - ☐ Refine original questions (if there are) or pose questions which can be answered with this data

What info about (tabular) data do we need to explore?

- ☐ **How many** rows and how many columns?
- ☐ What is the meaning of each row? Are there **rows having different meaning from the majority**?
- ☐ Are there **duplicated rows**?
- ☐ What is the meaning of each column?
- ☐ What is the current data type of each column? Are there columns having **inappropriate data types**?
- ☐ With each numerical column, how are values distributed?
 - ☐ What is the percentage of **missing values**?
 - ☐ Min? max? Are they **abnormal**?
- ☐ With each categorical column, how are values distributed?
 - ☐ What is the percentage of **missing values**?
 - ☐ How many different values? Show a few
 - ☐ Are they **abnormal**?

What info about (tabular) data do we need to explore?

- ☐ Previous slide shows basic info about data we should explore
- ☐ In complex cases, we may want to explore additional info about data
- ☐ For example, if we want to know more about column distribution, we can compute additional info using descriptive statistics

What is the purpose of preprocessing data?

- During data exploration, when we see a problem about data, we may need to do preprocessing operations to **fix the problem in order to continue to explore data**
- After identifying a specific question (we often reach this point after exploring data and understanding more about data) and the corresponding specific analysis method, we may need to do additional preprocessing operations with the goal: **preparing data which are ready for applying the specific analysis method**

Demo: using Python to do a data science process

Example from previous lecture

- **Ask a question:** what is the current state of understanding Python of students?
- **Collect data:** let students do quiz about Python in moodle, results can be downloaded as a csv file

Data exploration

- After understanding more about data, we will come back to “**ask a question**” step to refine original questions, or pose new questions

Refined question

- The question “what is the current state of understanding Python of students?” can be refined into 2 more specific questions:
- 1. How are values of the “Grade/10.00” column distributed?
 - ▣ *When exploring data, we just knew missing percentage, min, max; here we want to know more ...*
- 2. According to the criterion of being answered correctly by most students, which quiz is in first place, which quiz second, ...?

Answer the question

- ☐ **Preprocess data** (optional) + **analyze data** to answer question 1
- ☐ **Question 1:** How are values of the “Grade/10.00” column distributed?

From data of this column, how to answer this question?

- ☐ Option 1: Look at full data and feel ...
- ☐ Option 2: Summarize data using **descriptive statistics** :

Descriptive statistics

- ☐ We will use orange color to denote descriptive statistics for categorical data, and normal color for numerical data
- ☐ Summarize data with center: mean, median, **mode**
- ☐ Summarize data with center + range: mean & standard deviation, lower quartile & median & upper quartile
- ☐ Summarize data with full distribution: histogram, **bar plot**

Median

- Median = 50th percentile
- pth percentile ($0 \leq p \leq 100$) of a list of **values** is a value which tells us: there are about p% of values in the list $<$ this **value**
- Example: “75th percentile of the quiz scores = 8” means there are about 75% of students having quiz scores $<$ 8

Median

- There are different ways to compute pth percentile, here is one way:
- 1. Sort values in the list in ascending order
- 2. Find the location corresponding to p% of values in the list:
$$\text{location} = p/100 \times \text{the number of values in the list}$$

If it is not an integer, round it up
- 3. pth percentile = value at this location in the sorted list

Median

☐ Given this list: 1, 5, 3, 4, 2

Median = 50th percentile = 3

☐ Given this list: 1, 5, 3, 4, 2, 6

Median = 50th percentile = 3

Median vs mean

- ☐ Given this list: 1, 2, 3, 4, 5
 - ☐ Median = 3
 - ☐ Mean = 3
- ☐ Given this list: 1, 2, 3, 4, **500**
 - ☐ Median = 3
 - ☐ Mean = **102**
- ☐ Median is less affected by outlier (a value lying far away from the range of most values) than mean!

Mode

□ Mode = the most frequent value

Meand & standard deviation

Standard deviation (SD) of a list of values is a value ≥ 0 which tells us the average deviation of values in the list from mean

Example: “mean of quiz scores = 7 and SD = 0” means all students have quiz scores = 7; large SD means students have quiz scores spread widely around 7

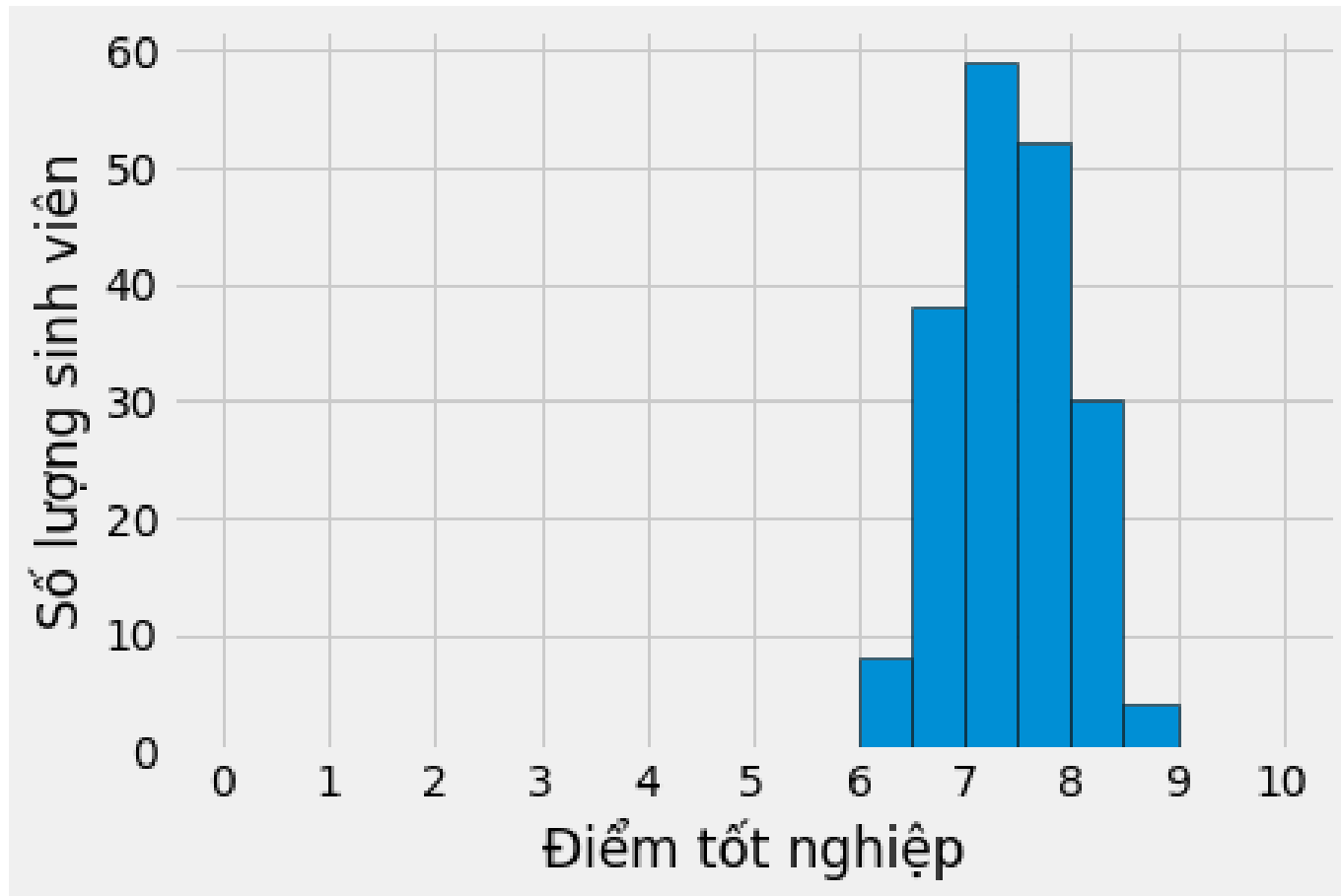
Chebychev discovered: With *any* list of values, the range mean $\pm z$ SD collects *at least* $(1 - \frac{1}{z^2}) \times 100\%$ of values in the list

- Mean ± 1 SD collects at least 0% of values (hmm ...)
- Mean ± 2 SD collects at least 75% of values
- Mean ± 3 SD collects at least 88.9% of values

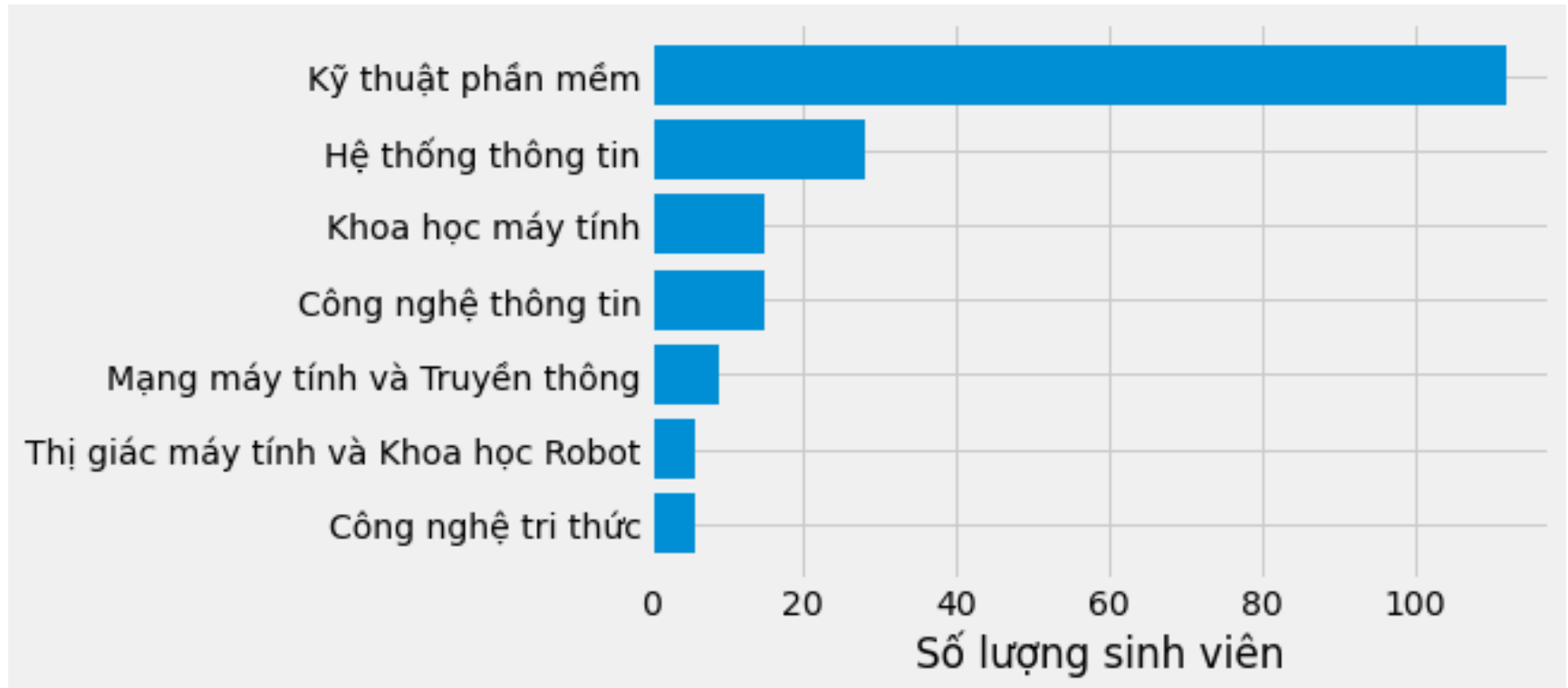
Lower quartile & median & upper quartile

- ☐ Lower quartile = 25th percentile
- ☐ Median = 50th percentile
- ☐ Upper quartile = 75th percentile

Histogram



Histogram



Reference