# Course Review

Wan Fang

Southern University of Science and Technology

# Design and Learning with Data

| Ask & Prepare | → | Process & analyze | → | Share/Act |
|---|---|---|---|---|

Question 🔍 Data  →  Insights & Visualzations
For analysts themselves (mostly)

→

✓ Infographics
✓ Dashboards
✓ Data videos
✓ Product/services
For other people

Data Discovery

Design Storytelling Communication

Data Literacy

Data Thinking

Context

Data Quality/ Cleaning

Descriptive Statistics/ EDA

Inferential statistics

A/B Testing

Visual Encoding Design

Interaction

Uncertainty

Tableau

# Data Storytelling

Data Literacy

- The ability to effectively communicate insights from a dataset using narratives and visualizations.
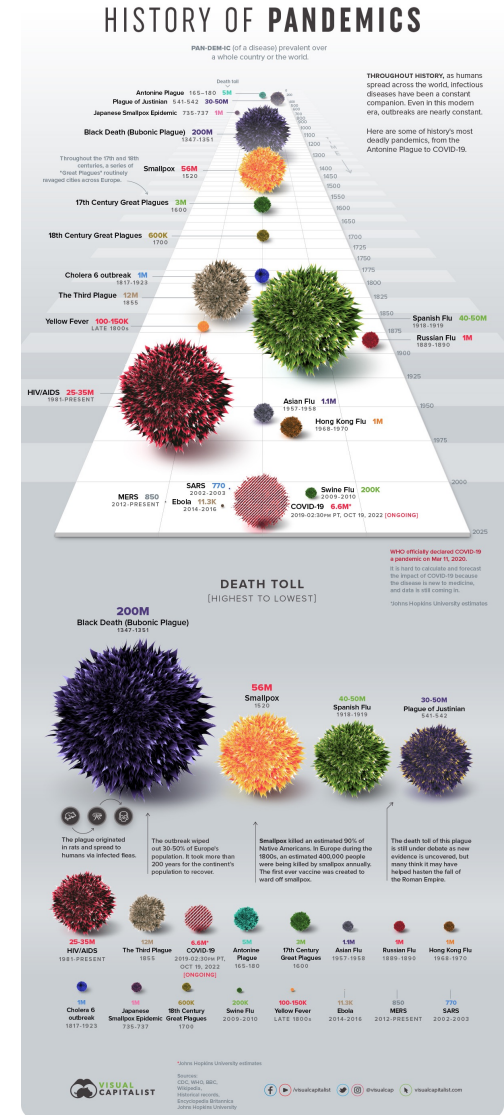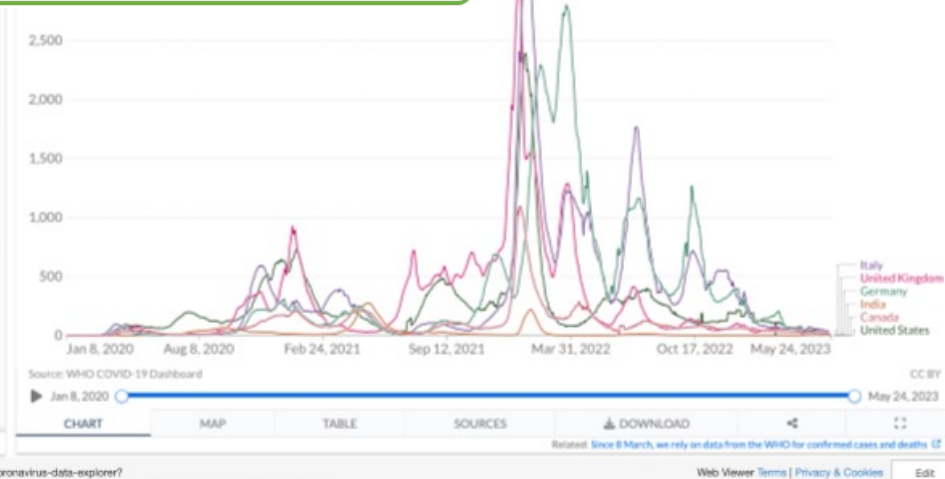
# Data Visualization

Data Literacy
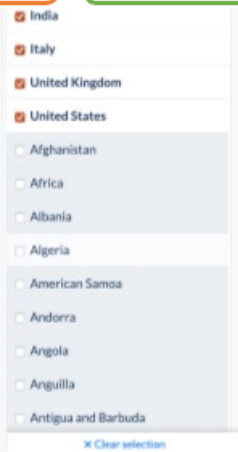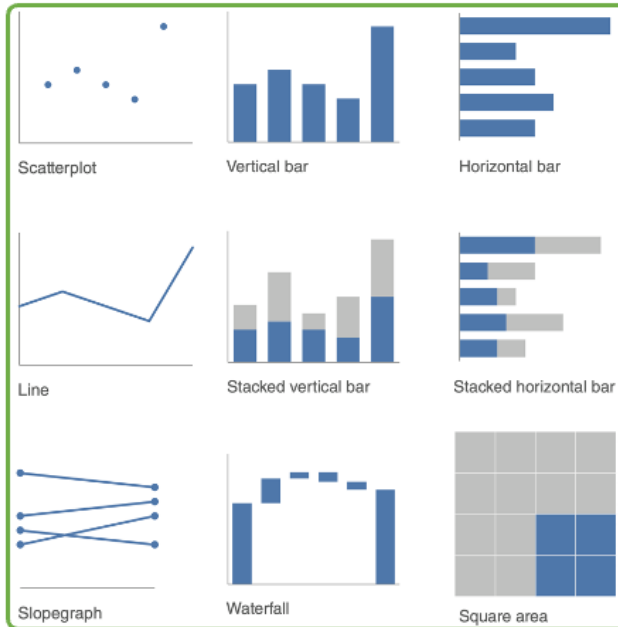
**Data Literacy**

# Dimensional Visualization of Data

## 1D: Nominal



## 1D Ordinal

- When you are interested in a single column containing ordinal values (i.e., counting and ranking are allowed)
  - E.g., **# of cylinders** column of the car dataset



## 1D: Quantitative



## 2D: Nominal x Nominal



## 2D Nominal x Quantitative

- If you are interested in how one nominal and one quantitative columns
  - E.g., origin and horsepower columns of the car dataset



## 2D: Quantitative x Quantitative
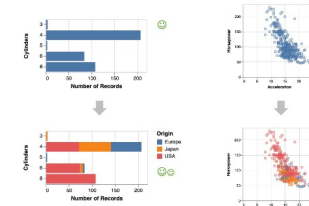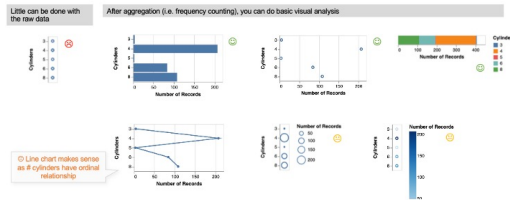


## 3D ANY

- Each visualization can accommodate 1-2 extra columns with color or size encodings. Why not explore higher-dimensions?
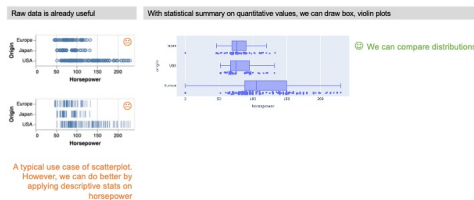


## Higher Dimension

- Single charts usually cannot accommodate larger than 5 dimensions.
  - However, we can use **composite charts**.
  - For example, we have used scatterplot matrix in the previous tutorial.

# The Importance of Context

Context

## 6 basic problem types

- Making predictions

- Categorizing things

- Spotting something unusual

- Identifying themes

- Discovering connections

- Finding patterns

## Craft effective questions

- SMART methodology
  - Specific — does the question address the problem? Does it have a context?
  - Measurable — does it give the answer that can be measured?
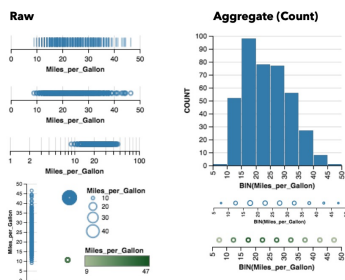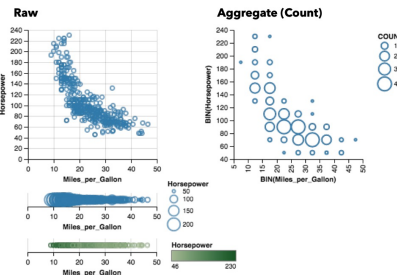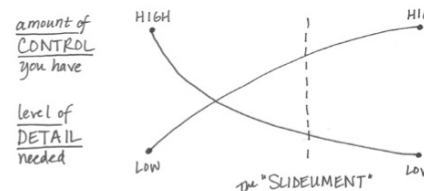  - Action-oriented — will the info that we get help us devise an action plan?
  - Relevant — is it about a particular problem we are trying to solve?
  - Time-bounded — are the answers relevant to the specific time being studied?

## Who, What, and How

- ***To whom are you communicating?***
  - It is important to have a good understanding of who your audience is and how they perceive you. This can help you to **identify common ground** that will help you ensure they hear your message.

- ***What do you want your audience to know or do?***
  - You should be clear how you want your audience to act and take into account how you will communicate to them and the overall tone that you want to set for your communication.

- ***How can you use data to help make your point?***
  - It's **only after** you can concisely answer these first two questions that you're ready to move forward with the third.



LIVE PRESENTATION . . . . . . . . WRITTEN DOC OR EMAIL

You | audience

audience

amount of CONTROL you have

level of DETAIL needed

HIGH | HIGH

LOW | LOW

The "SLIDEUMENT"

### Prompting action

Here are some action words to help act as thought starters as you determine what you are asking of your audience:

accept | agree | begin | believe | change | collaborate | commence | create | defend | desire | differentiate | do | empathize | empower | encourage | engage | establish | examine | facilitate | familiarize | form | implement | include | influence | invest | invigorate | know | learn | like | persuade | plan | promote | pursue | recommend | receive | remember | report | respond | secure | support | simplify | start | try | understand | validate

### Ignore the nonsupporting data?

You might assume that showing only the data that backs up your point and ignoring the rest will make for a stronger case. I do not recommend this. Beyond being misleading by painting a one-sided story, this is very risky. A discerning audience will poke holes in a story that doesn't hold up or data that shows one aspect but ignores the rest. The right amount of context and supporting and opposing data will vary depending on the situation, the level of trust you have with your audience, and other factors.

# Data X

Data Thinking

- **Data science vs Data analytics**
- *Data ecosystem / Data life cycle / Data Privacy & Ethics*
- *Data Integrity / Data & Analytics Skills*

**Descriptive analytics**
- looks at data to examine, understand, and describe something that's already happened.
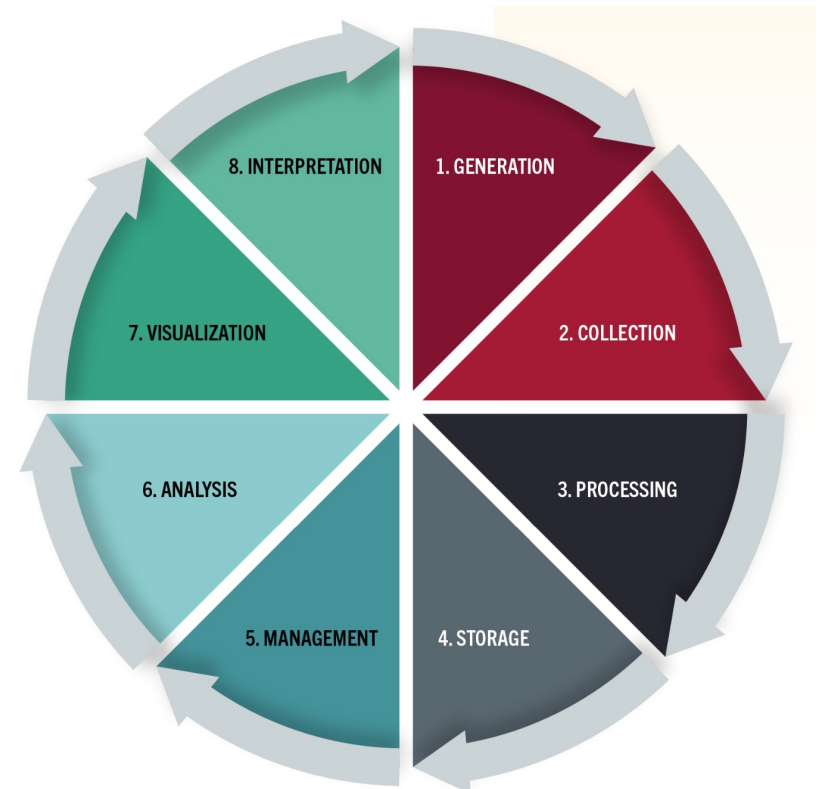
**Diagnostic analytics**
- goes deeper than descriptive analytics by seeking to understand the "why" behind what happened.

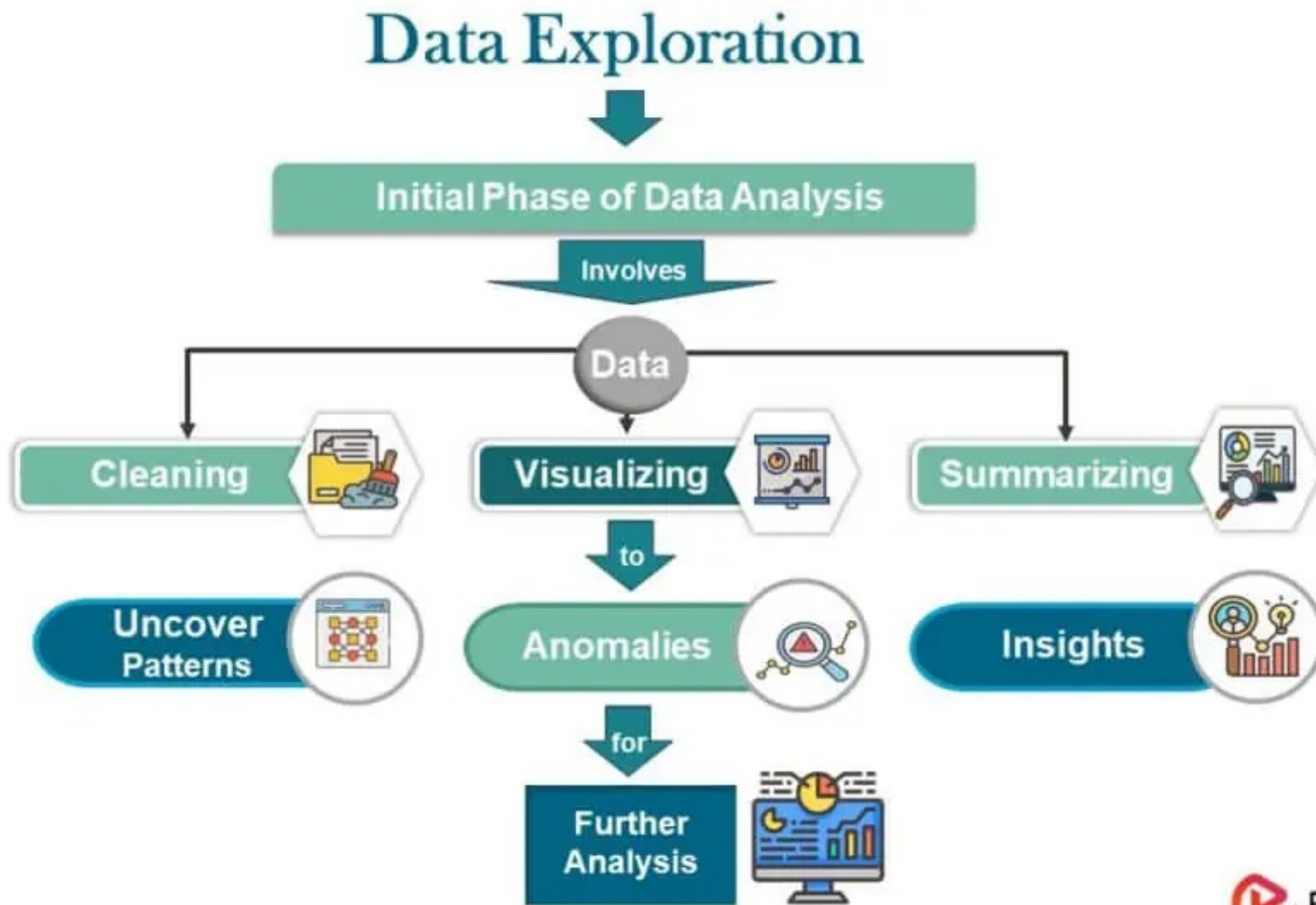**Predictive analytics**
- relies on historical data, past trends, and assumptions to answer questions about what will happen in the future.

**Prescriptive analytics**
- identifies specific actions an individual or organization should take to reach future targets or goals.

8. INTERPRETATION
1. GENERATION
2. COLLECTION
3. PROCESSING
4. STORAGE
5. MANAGEMENT
6. ANALYSIS
7. VISUALIZATION

**Data Quality/ Cleaning**

# Big Data Quality Assessment Framework

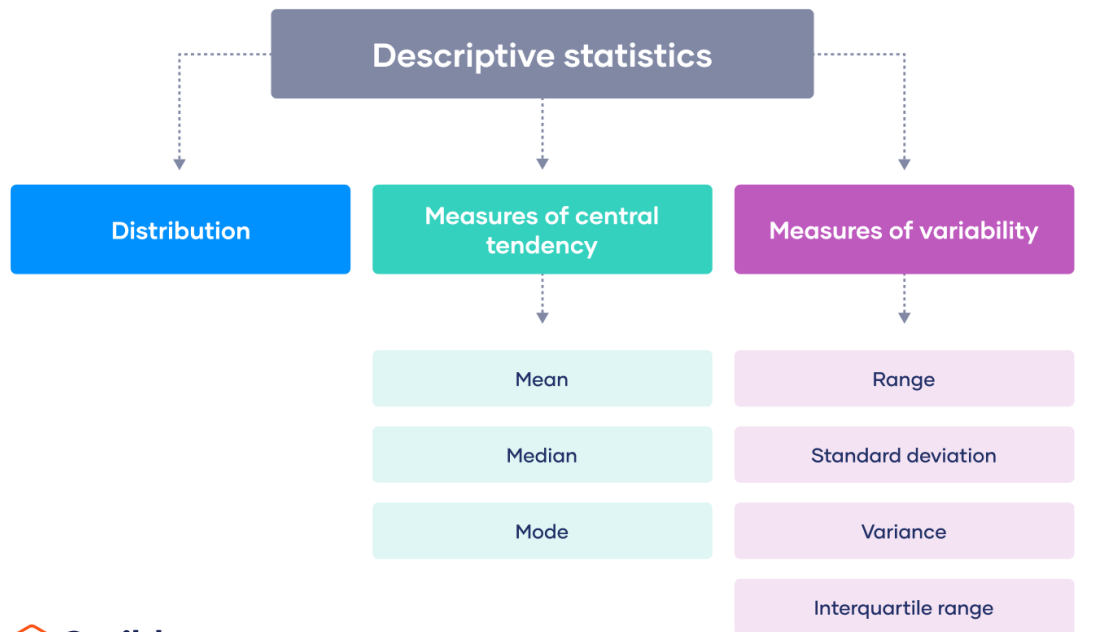| Dimensions | Elements | Indicators |
|---|---|---|
| 1) Availability | 1) Accessibility | ■ Whether a data access interface is provided<br>■ Data can be easily made public or easy to purchase |
| | 2) Timeliness | ■ Within a given time, whether the data arrive on time<br>■ Whether data are regularly updated<br>■ Whether the time interval from data collection and processing to release meets requirements |
| 2) Usability | 1) Credibility | ■ Data come from specialized organizations of a country, field, or industry<br>■ Experts or specialists regularly audit and check the correctness of the data content<br>■ Data exist in the range of known or acceptable values |
| 3) Reliability | 1) Accuracy | ■ Data provided are accurate<br>■ Data representation (or value) well reflects the true state of the source information<br>■ Information (data) representation will not cause ambiguity |
| | 2) Consistency | ■ After data have been processed, their concepts, value domains, and formats still match as before processing<br>■ During a certain time, data remain consistent and verifiable<br>■ Data and the data from other data sources are consistent or verifiable |
| | 3) Integrity | ■ Data format is clear and meets the criteria<br>■ Data are consistent with structural integrity<br>■ Data are consistent with content integrity |
| | 4) Completeness | ■ Whether the deficiency of a component will impact use of the data for data with multi-components<br>■ Whether the deficiency of a component will impact data accuracy and integrity |
| 4) Relevance | 1) Fitness | ■ The data collected do not completely match the theme, but they expound one aspect<br>■ Most datasets retrieved are within the retrieval theme users need<br>■ Information theme provides matches with users' retrieval theme |
| 5) Presentation Quality | 1) Readability | ■ Data (content, format, etc.) are clear and understandable<br>■ It is easy to judge that the data provided meet needs<br>■ Data description, classification, and coding content satisfy specification and are easy to understand |

**garbage in, garbage out**
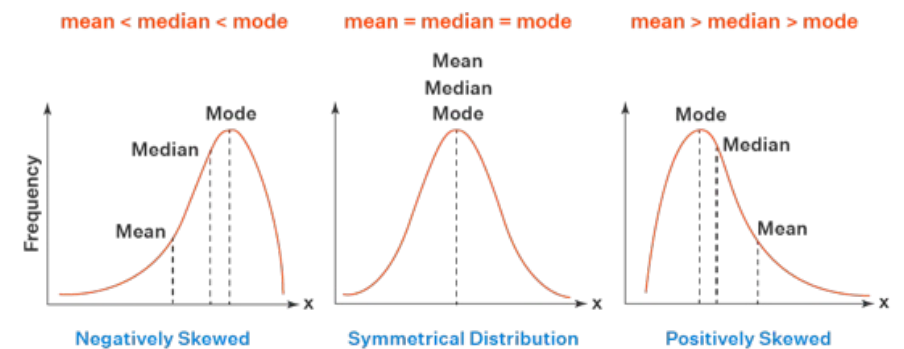
[garbage in, garbage out] 🔊

DEFINITION

used to express the idea that in computing and other fields, incorrect or poor-quality input will produce faulty output.

# Descriptive Statistics

Descriptive statistics

**Statistical graphics**

# Univariate/Multivariate Graphical
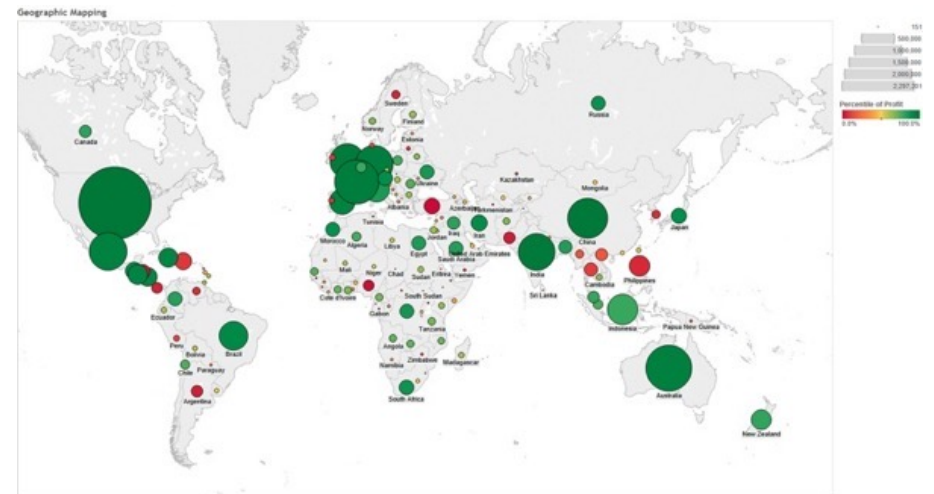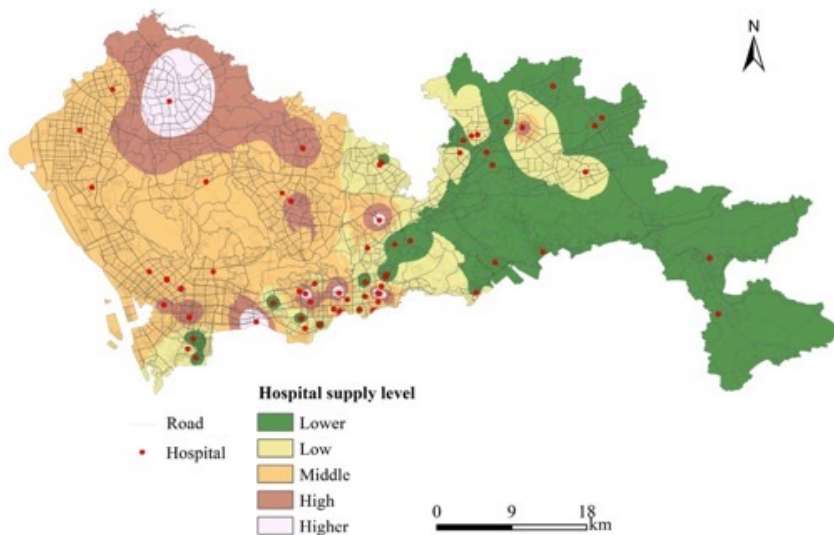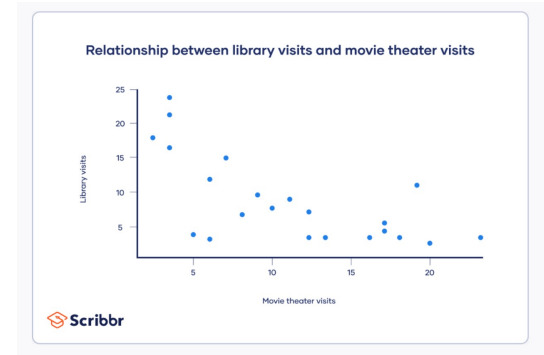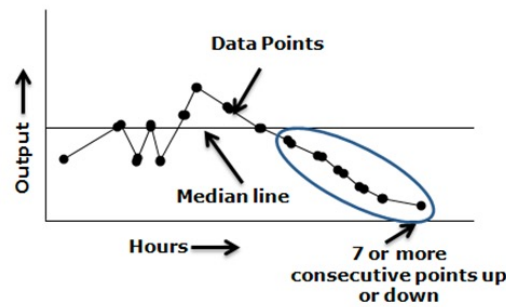


Figure 4. Spatial distribution of road network and hospital supply level in Shenzhen.

# With Python

**Descriptive Statistics/ EDA**

# Inferential Statistics

Inferential statistics

Step 1: Write your hypotheses and plan your research design

Step 2: Collect data from a sample

Step 3: Summarize your data with descriptive statistics

Step 4: Test hypotheses or make estimates with inferential statistics

Step 5: Interpret your results

Example: Statistical hypotheses to test an effect

- *Null hypothesis*: A 5-minute meditation exercise will have no effect on math test scores in teenagers.
- *Alternative hypothesis*: A 5-minute meditation exercise will improve math test scores in teenagers.

Example: Statistical hypotheses to test a correlation

- *Null hypothesis*: Parental income and GPA have no relationship with each other in college students.
- *Alternative hypothesis*: Parental income and GPA are positively correlated in college students.

| Null hypothesis is ... | True | False |
|---|---|---|
| Rejected | Type I error / False positive / Probability = α | Correct decision / True positive / Probability = 1 − β |
| Not rejected | Correct decision / True negative / Probability = 1 − α | Type II error / False negative / Probability = β |

Scribbr



Null hypothesis (H₀) distribution — Alternative hypothesis (H₁) distribution

1 − α   1 − β

β   α

Type II error rate   Type I error rate

**A/B Testing**



- **A/B testing is a special kind of quantitative experiments**
  - where "A" representing the old control, and "B" representing a single / multiple experimental changes
  - where the goal is to maximize single/multiple measurements of user behavior
  - A/B testing is most useful when combined with qualitative methods (e.g. observation, survey, interview)

- **A/B testing has multiple goals as listed below**
  - Finding a winner among many design variations
  - Testing hypotheses of cause (design) and effect (user behavior)
  - Gaining generalizable understanding of the task and the user population
  - Developing a reusable platform for continuously running a series of A/B tests

- **A/B testing consists of five big steps**
  1. Choose a single (or couple) quantifiable measure(s) of user behavior
  2. Construct hypotheses
  3. Create design variations
  4. Run the experiment
  5. Analyze the results

**A/B Testing**

**1. Equally allocate users to design variations (if possible)**

Don't assign majority of users to specific conditions. It will slow down the experiment by increasing the minimum sample sizes to get statistical significance.

**2. Get sufficient samples to get statistical significance of at least 95% (p<0.05)**

300-400 samples per variation is usually considered enough. If you cannot reach the 95% (p=0.05) with 300-400 samples * # variations, it is very unlikely that the variations have strong impact on the measured goal.

**3. Run the experiment for multiple behavioral cycles**

Lots of user behaviors tend to have weekly patterns (e.g. on Monday, people tend to be quite busy). To minimize the risk of weekly biases, run the experiment for at least one week (preferably two weeks or a month). However, it depends on what the task is and who the users are.

**4. Pay attention to external factors**

If it is Christmas or Valentine's day, an online marketing campaign might be highly affected by the external factor. While it's impossible to get a perfectly clean environment, consider external factors in your interpretation.

- A/B testing is not silver bullet. It has a lot of limitations and pitfalls.

| When A/B testing is not worth | What's the problem? How can we make it worth again? |
|---|---|
| **You don't have meaningful traffic (i.e. # users)** | • Without meaningful traffic you won't be able to tell anything about statistical significance.<br>• Reduce # variations<br>• Wait until you have a big enough population<br>• Construct a hypothesis that you can indirectly test on a similar platform (e.g. running A/B test of your streaming platform on YouTube users), crowdsourcing platform (e.g. Amazon Mechanical Turk), or as a lab experiment |
| **You don't have enough resources for running A/B test.** | • You will regret, "We should've finished our design first."<br>• Wait until you have enough time<br>• Consider using existing tools for A/B testing. They are not mature yet, but still useful for marketing. |
| **You don't have an informed hypothesis.** | • Gather more information and perform EDA<br>• Treat A/B test like real science. Learn from lectures, case studies, and tutorials. Ask psychologists and data scientists who have conducted randomized controlled experiments what hypothesis you can test for the given situation |

# Visual Encoding Design

| Example | Encoding | Ordered | Useful values | Quantitative | Ordinal | Categorical | Relational |
|---|---|---|---|---|---|---|---|
| | position, placement | yes | infinite | Good | Good | Good | Good |
| 1, 2, 3; A, B, C | text labels | optional alpha or num | infinite | Good | Good | Good | Good |
| | length | yes | many | Good | Good | | |
| | size, area | yes | many | Good | Good | | |
| | angle | yes | medium | Good | Good | | |
| | pattern density | yes | few | Good | Good | | |
| | weight, boldness | yes | few | | Good | | |
| | saturation, brightness | yes | few | | Good | | |
| | color | no | few (<20) | | | Good | |
| | shape, icon | no | medium | | | Good | |
| | pattern texture | no | medium | | | Good | |
| | enclosure, connection | no | infinite | | | Good | Good |
| | line pattern | no | few | | | | Good |
| | line endings | no | few | | | | Good |
| | line weight | yes | few | | Good | | |

Interaction

# Ben Shneiderman's information-seeking mantra



**1. Overview first**　　　2. Zoom and Filter　　　3. Detail-on-demand

The initial view of
provides the data

1. Overview first　　　**2. Zoom and Filter**　　　3. Detail-on-demand

Sorted by Category (A->Z)　　　　　Sorted by Ratio (small->large)

1. Overview first　　　2. Zoom and Filter　　　**3. Detail-on-demand**

**India**
continent=Asia
gdp per capita=2,452.21
life expectancy=64.698
population=1110396331

**Mouse-over Tooltip**
Most EDA tools allow viewers to see
raw data of individual items as tool-tip

**Data table filtered by user selection**
Most EDA tools offer data tables that shows raw data of
currently selected items
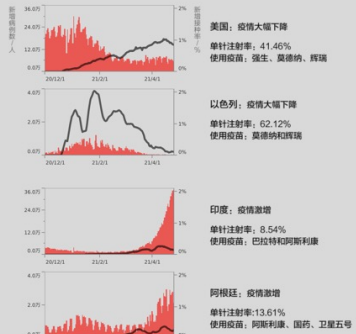
Tableau

18

# 生命之桥

## 新冠疫苗分配的富与贫

新冠疫苗是过河之路，帮助我们越过疫情的激流；新冠疫苗是连接之桥，将每个国家、每个生命彼此联结。

我们用艺术可视化的方式，具有交互性的可视化海报，以横向阅读的叙事方式，展现了新冠疫苗分配不公平的情况，希望引起人们的关注，让更多人参与到公平疫苗分配之中。
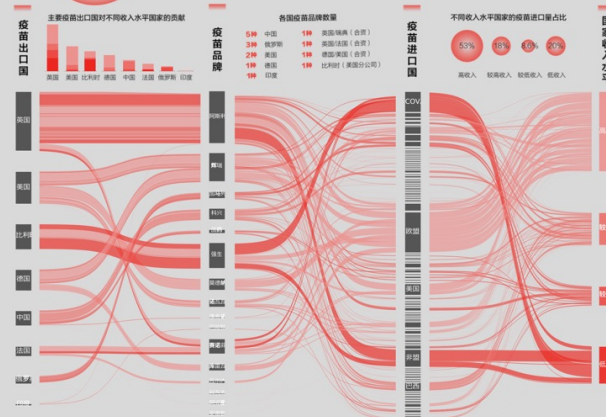
请向右翻阅，横向阅读 →

## 新冠疫苗的接种对疫情的控制有重要作用

我们分别选取了疫苗接种率较高的美国和以色列，以及较低的印度和阿根廷。通过展现新增病例和疫苗接种率的关系，体现疫苗对于疫情控制的重要性。

**美国：** 疫情大幅下降
单针注射率：41.46%
使用疫苗：强生、莫德纳、辉瑞

**以色列：** 疫情大幅下降
单针注射率：62.12%
使用疫苗：莫德纳和辉瑞

**印度：** 疫情激增
单针注射率：8.54%
使用疫苗：巴拉特和阿斯利康

**阿根廷：** 疫情激增
单针注射率：13.61%
使用疫苗：阿斯利康、卫星五号

数据来源：https://www.kaggle.com/gpreda/covid-world-vaccination-progress_截止2021.4.25
https://ourworldindata.org/covid-vaccinations_截止2021.4.25

## 新冠疫苗流向

我们选取了8个疫苗主要出口国，分析得出疫苗主要流向高收入国家，而低收入国家由于covax计划也获得较多疫苗。

疫苗出口国 | 主要疫苗出口国家对不同收入水平国家的贡献

各国疫苗品牌数量
5种 中国 | 1种 俄罗斯
3种 美国 | 1种 美国/法国（合资）
2种 美国 | 1种 美国/德国（合资）
1种 德国 | 1种 比利时/美国（合资）
1种 印度

疫苗品牌

疫苗进口国

不同收入水平国家的疫苗进口量占比
8.5% | 18% | 36% | 37%
高收入 | 较高收入 | 较低收入 | 低收入

国家收入水平

↑ 点击各个节点，可筛选桑基图

数据来源：https://public.tableau.com/profile/duke.global.health.innovation.center#/vizhome/COVID-19VaccineAdvanceMarketCommitmentsbyCountry_16131542122100_截止2021.4.9

## 疫苗接种

据世界卫生组织统计，截至4月13日

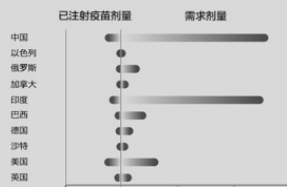全球共有15个国家的接种剂量数超1千万剂，共占全球总接种量的84%

与此同时，将近130个国家，超过25亿人口没有接种

人均疫苗接种剂量
总接种剂量

高收入
较高收入
较低收入
低收入

↑ 悬停于"国家收入水平"节点，可高亮地图

© 2023 Mapbox © OpenStreetMap

数据来源：https://www.kaggle.com/gpreda/covid-world-vaccination-progress_截止2021.4.25
https://ourworldindata.org/covid-vaccinations_截止2021.4.25

---

## 囤积新冠疫苗

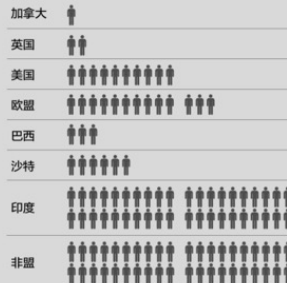报告显示，疫苗接种率达到人口的90%，可以达到群体免疫，而当今世界疫苗接种情况仍不容乐观，已注射剂量远远达不到需求剂量。

本图表中需求剂量
=人口 * 2（剂量）* 90%

部分高收入国家需求量小，却在大量囤积疫苗；与此同时没有一个中低收入国家拥有足够疫苗。

### 疫苗已注射剂量远小于需求剂量

已注射疫苗剂量 | 需求剂量

中国
以色列
俄罗斯
加拿大
印度
巴西
德国
沙特
美国
英国

数据来源：https://www.covid19-vaccine.live_截止2021.5.14

### 不同收入水平国家的订购剂量

高收入
较高收入
较低收入
低收入
COVAX

### 疫苗需求剂量对比

| 国家 | 剂量 |
|---|---|
| 加拿大 | |
| 英国 | |
| 美国 | |
| 欧盟 | |
| 巴西 | |
| 沙特 | |
| 印度 | |
| 非盟 | |

### 人均疫苗占有剂量对比（基于订购剂量计算）

| | 剂量 | 比例 |
|---|---|---|
| | 10.0 | 1人/10剂 |
| | 7.6 | 1人/7.6剂 |
| | 3.7 | 1人/3.7剂 |
| | 6.4 | 1人/6.4剂 |
| | 2.2 | 1人/2.2剂 |
| | 0.21 | 4.5人/1剂 |
| | 0.28 | 3.5人/1剂 |
| | 0.20 | 5人/1剂 |

## COVAX：与低收入国家分享新冠疫苗的国际行动

世卫组织强力批判在严酷的全球疫情形势下囤积疫苗的"疫苗民族主义"行为，为促进新冠疫苗的公平性，提出COVAX这一解决方案。

| 参与国家 | 参与的低收入国家 |
|---|---|
| 190 | 92 |

| 累计送达国家 | 累计运送疫苗剂量 |
|---|---|
| 121 | 54,007,670 |

**疫苗**
46亿 达到群体免疫所需要购买的疫苗量
20亿 2021年预计要的疫苗剂量
5400万 已筹集到的疫苗量

**资金**
68亿美元 2021年疫苗研发与分配等所需的资金投入
24亿美元 已筹集到的资金

### 紧急使用清单的疫苗品牌对比

综合来看，中国国药较适合于公平分配计划。

| | 疫苗技术 | 接种剂量 | 有效性 | 价格 | 严重副作用 | 储存条件 |
|---|---|---|---|---|---|---|
| 国药 | 灭活 | 2 | 79% | $44/剂 | 无 | 2℃~8℃ |
| 莫德纳 | mRNA | 2 | 94% | $25~$37/剂 | 无 | -25℃~-15℃（6个月）2℃~8℃（30天） |
| 强生 | 腺病毒载体 | 1 | 66% | $10/剂 | 无 | 2℃~8℃（3个月） |
| 阿斯利康/牛津 | 腺病毒载体 | 2 | 70% | $3~$4/剂 | 血栓 | 2℃~8℃（6个月） |
| 辉瑞 | mRNA | 2 | 95% | $20/剂 | 面瘫/严重过敏 | -80℃~-60℃（6个月）2℃~8℃（5天） |

数据来源：http://www.360doc.com/content/21/0322/09/3843034_968229393.shtml
https://news.un.org/zh/story/2021/05/1083742

## NO ONE IS SAFE UNTIL EVERYONE IS SAFE.
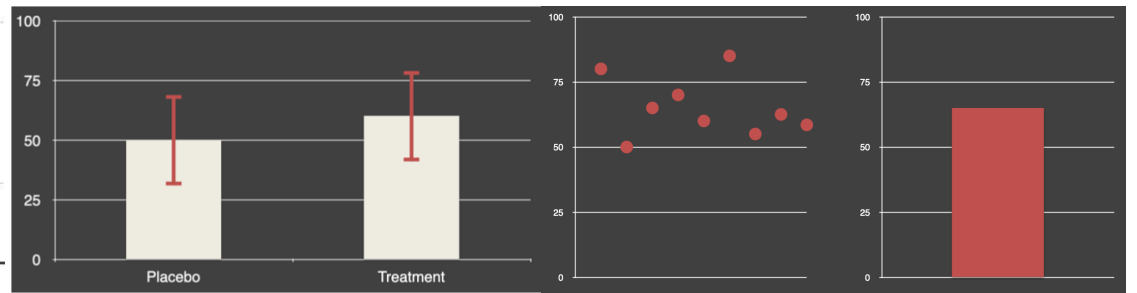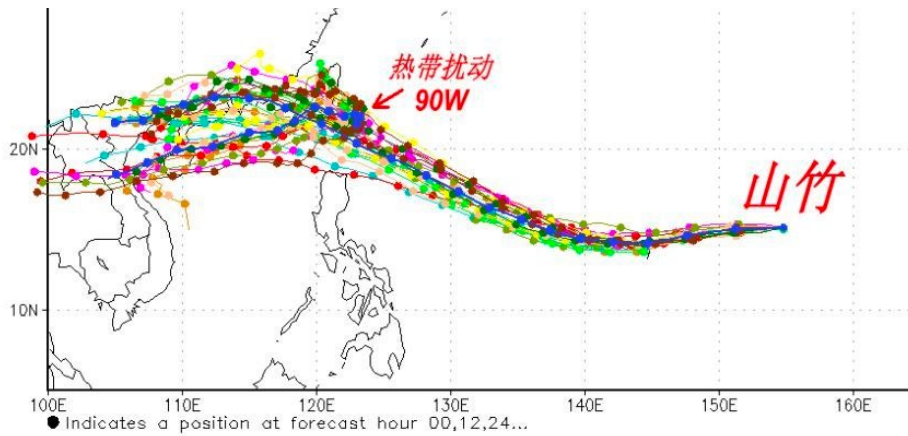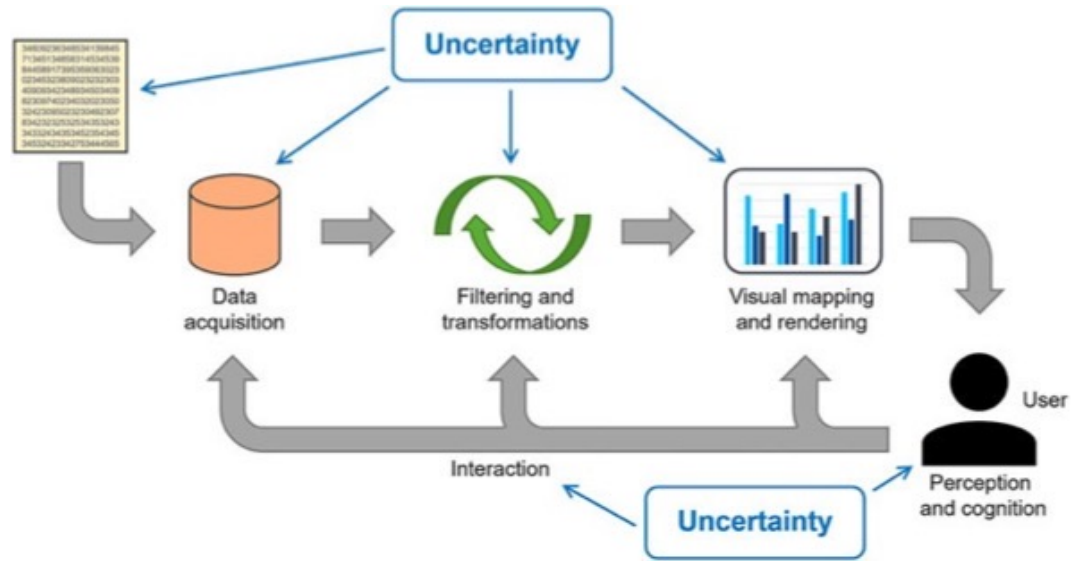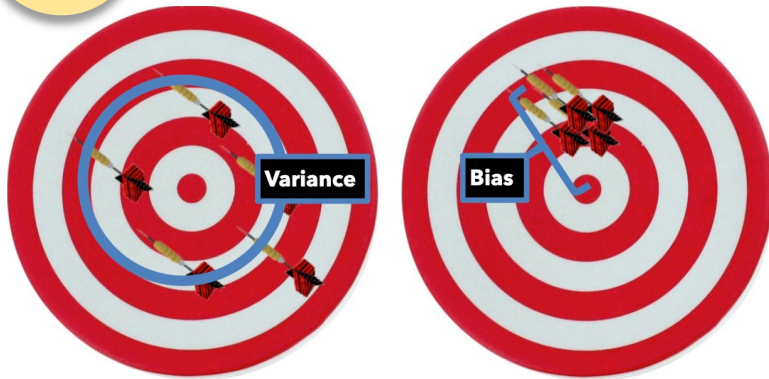
### 疫情无国界，新冠疫苗是生命之桥

面对百年疫情，新冠疫苗的研发与分配绝不是国与国、企业与企业之间的竞争，而是全人类与病毒的决战，捍卫人类生命福祉的社会理性应高于市场主导的经济逻辑。

人类是一个命运共同体，只有确保所有处于危险之中的人得到有效防护，人类才能彻底战胜病毒，真正实现整体安全。因为正如谭德塞所说，抗击疫情"必须从一个全球大家庭的角度出发"。此时此刻，我们能做的只有团结。

TONGJI UNIVERSITY COLLEGE OF DESIGN AND INNOVATION

Uncertainty

Variance

Bias

热带扰动 90W

山竹

20N

10N

100E   110E   120E   130E   140E   150E   160E
● Indicates a position at forecast hour 00,12,24...

Uncertainty

Data acquisition → Filtering and transformations → Visual mapping and rendering

Interaction

Uncertainty

User

Perception and cognition

100   75   50   25   0
Placebo   Treatment

1   2   3   4   5   6

# The Takeaways

In our life, study, work

Understand the intrinsic nature of data

Make use of it for our purpose

# Time for
# Course Evaluation

Your participation is important!

# Course Evaluation

**方法及步骤**

1.网页端：登录教务系统：https://tis.sustech.edu.cn/-业务办理-评教任务-2024春季学期学生评价任务。系统按课程类型设置评价任务（理论类、实验实践类、体育类、艺术类），如页面上有多个评价任务，请逐一进入并提交评价。

2.微信端：通过微信进入"南方科技大学"微信企业号--教学质量管理平台，在"我的任务-待评"中填写并提交本学期所选课程的所有听课评价。

# Course Project | Final Showcase

- Final Project Showcase on **Data Storytelling**
  - **Submit** first draft before Fri Jun 07 @ 12:00 | **Present** on Fri Jun 07 in-class
  - **Submit** final version before Sunday Jun 09 @ 23:30

  - <u>**~ 8 mins presentation**</u> regarding your data-driven story
  - <u>**~ 2 mins Q&A**</u>

# Thank you~

Wan Fang
Southern University of Science and Technology