# Data Discovery
# Lecture 7

Wan Fang

Southern University of Science and Technology

# Agenda

- Descriptive Statistics
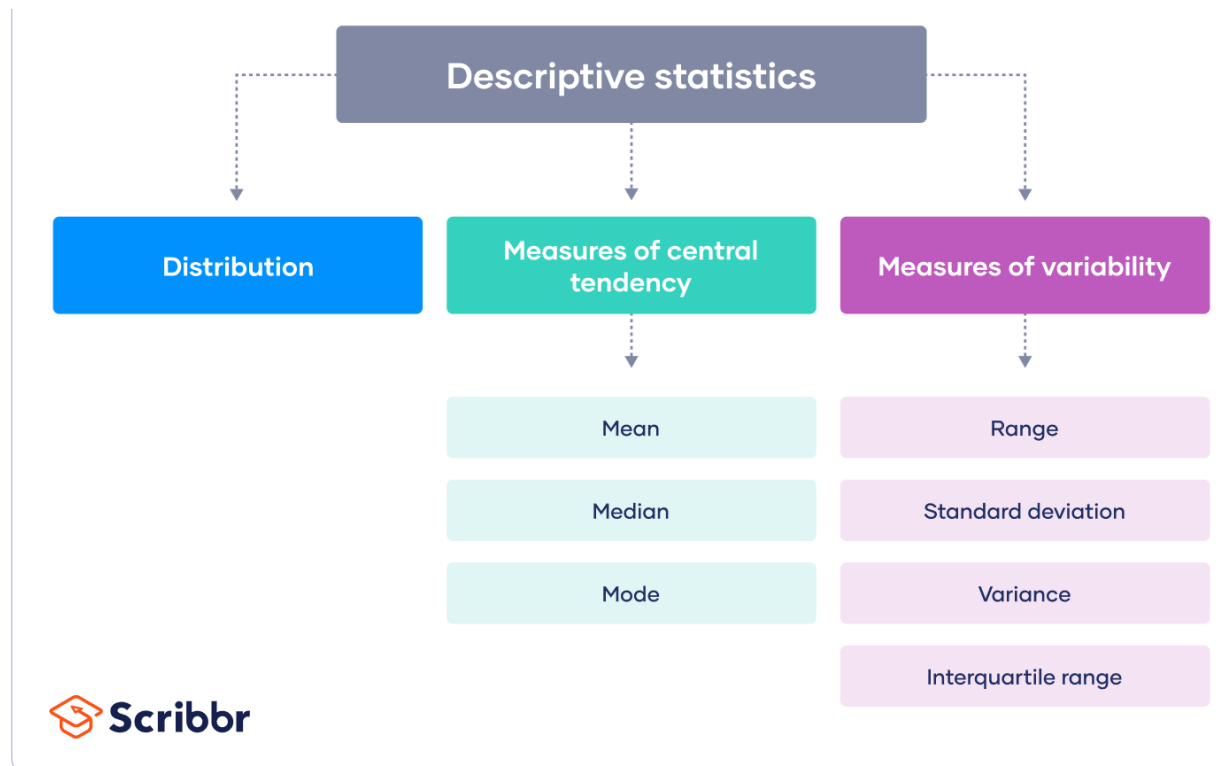
- Types of Exploratory Data Analysis

# Descriptive Statistics

Wan Fang

Southern University of Science and Technology

# Types of Descriptive Statistics

- The **distribution** concerns the frequency of each value.

- The **central tendency** concerns the averages of the values.

- The **variability** or dispersion concerns how spread out the values are.

# Frequency Distribution

- A data set is made up of a distribution of values, or scores.
  - In tables or graphs, you can summarize ***the frequency of every possible value of a variable in numbers or percentages***.

---

### Research example

You want to study the popularity of different leisure activities by gender.

You distribute a survey and ask participants how many times they did each of the following in the past year:

- Go to a library
- Watch a movie at a theater
- Visit a national park

Your data set is the collection of responses to the survey.

Now you can use descriptive statistics to find out the overall frequency of each activity (distribution), the averages for each activity (central tendency), and the spread of responses for each activity (variability).

---

# Frequency Distribution

- A data set is made up of a distribution of values, or scores.
  - In tables or graphs, you can summarize ___the frequency of every possible value of a variable in numbers or percentages___.

| Simple frequency distribution table | Grouped frequency distribution table |
|---|---|

For the variable of gender, you list all possible answers on the left hand column. You count the number or percentage of responses for each answer and display it on the right hand column.

| Gender | Number |
|---|---|
| Male | 182 |
| Female | 235 |
| Other | 27 |

From this table, you can see that more women than men or people with another gender identity took part in the study.

# Frequency Distribution

- A data set is made up of a distribution of values, or scores.
  - In tables or graphs, you can summarize ***the frequency of every possible value of a variable in numbers or percentages***.

| Simple frequency distribution table | Grouped frequency distribution table |
| --- | --- |

In a grouped frequency distribution, you can group numerical response values and add up the number of responses for each group. You can also convert each of these numbers to percentages.

| Library visits in the past year | Percent |
| --- | --- |
| 0–4 | 6% |
| 5–8 | 20% |
| 9–12 | 42% |
| 13–16 | 24% |
| 17+ | 8% |

From this table, you can see that most people visited the library between 5 and 16 times in the past year.

# Measures of Central Tendency

- Estimate the center, or average, of a data set.
  - The **mean**, median and mode are 3 ways of finding the average.

| **Mean** | **Median** | **Mode** |

The **mean**, or *M*, is the most commonly used method for finding the average.

To find the mean, simply add up all response values and divide the sum by the total number of responses. The total number of responses or observations is called *N*.

**Mean number of library visits**

| | |
|---|---|
| **Data set** | 15, 3, 12, 0, 24, 3 |
| **Sum of all values** | 15 + 3 + 12 + 0 + 24 + 3 = 57 |
| **Total number of responses** | $N = 6$ |
| **Mean** | Divide the sum of values by $N$ to find $M$: 57/6 = **9.5** |

# Measures of Central Tendency

- Estimate the center, or average, of a data set.
  - The mean, **median** and mode are 3 ways of finding the average.

**Mean     Median     Mode**

The **median** is the value that's exactly in the middle of a data set.

To find the median, order each response value from the smallest to the biggest. Then, the median is the number in the middle. If there are two numbers in the middle, find their mean.

**Median number of library visits**

| | |
|---|---|
| **Ordered data set** | 0, 3, 3, 12, 15, 24 |
| **Middle numbers** | 3, 12 |
| **Median** | Find the mean of the two middle numbers: (3 + 12)/2 = **7.5** |

# Measures of Central Tendency

- Estimate the center, or average, of a data set.
  - The mean, median and **mode** are 3 ways of finding the average.

**Mean**  **Median**  **Mode**

The **mode** is the simply the most popular or most frequent response value. A data set can have no mode, one mode, or more than one mode.

To find the mode, order your data set from lowest to highest and find the response that occurs most frequently.
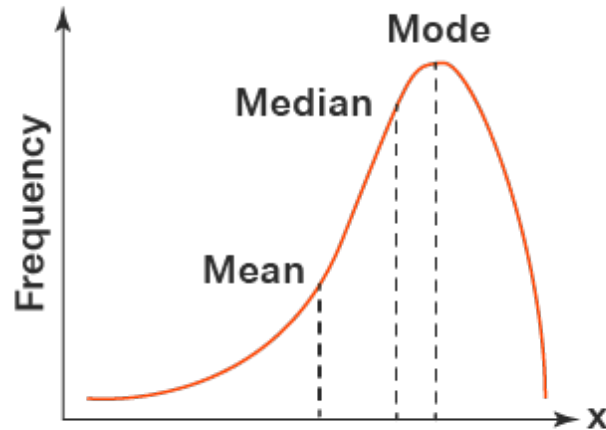
**Mode number of library visits**

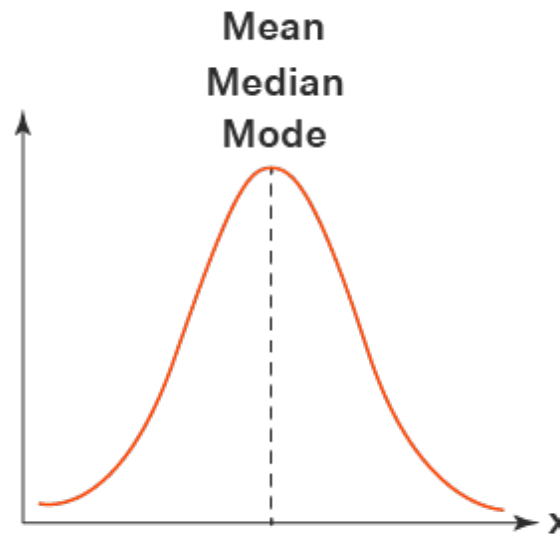| | |
|---|---|
| **Ordered data set** | 0, 3, 3, 12, 15, 24 |
| **Mode** | Find the most frequently occurring response: **3** |

# Measures of Central Tendency

- **Mode**: the most popular response or value in the data set.

- **Median**: the value in the exact middle of the data set when ordered from low to high.

- **Mean**: the sum of all values divided by the number of values.
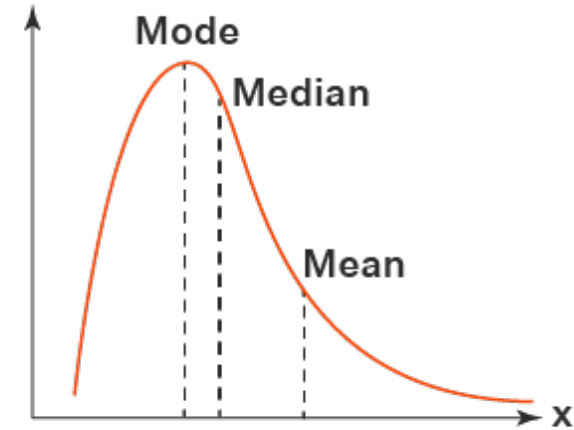
mean < median < mode    mean = median = mode    mean > median > mode

**Negatively Skewed**    **Symmetrical Distribution**    **Positively Skewed**

# Measures of Central Tendency



2020年各城市薪资水平（元）

# Measures of Variability

- Give you a sense of how spread out the response values are.
  - The range, standard deviation and variance each reflect different aspects of spread.

- **Range**
  - The range gives you an idea of how far apart the most extreme response scores are.
  - To find the range, simply subtract the lowest value from the highest value.

**Range of visits to the library in the past year**

**Ordered data set:**  0, 3, 3, 12, 15, 24

**Range:** 24 − 0 = **24**

# Measures of Variability

- **Standard deviation**
  - The standard deviation (*s* or *SD*) is the average amount of variability in your dataset.
    - It tells you, on average, how far each score lies from the mean.
    - The larger the standard deviation, the more variable the data set is.

*Six steps for finding the standard deviation:*
1. *List each score and find their mean.*
2. *Subtract the mean from each score to get the deviation from the mean.*
3. *Square each of these deviations.*
4. *Add up all of the squared deviations.*
5. ***Divide the sum of the squared deviations by N – 1.***
6. ***Find the square root of the number you found.***

Standard deviations of visits to the library in the past year

In the table below, you complete **Steps 1 through 4**.

| Raw data | Deviation from mean | Squared deviation |
|---|---|---|
| 15 | 15 − 9.5 = 5.5 | 30.25 |
| 3 | 3 − 9.5 = -6.5 | 42.25 |
| 12 | 12 − 9.5 = 2.5 | 6.25 |
| 0 | 0 − 9.5 = -9.5 | 90.25 |
| 24 | 24 − 9.5 = 14.5 | 210.25 |
| 3 | 3 − 9.5 = -6.5 | 42.25 |
| *M* = 9.5 | Sum = 0 | Sum of squares = 421.5 |

**Step 5:** 421.5/5 = 84.3

**Step 6:** √84.3 = 9.18

From learning that **s = 9.18**, you can say that on average, each score deviates from the mean by 9.18 points.

# Variance

- The average of squared deviations from the mean.
  - Variance reflects the degree of spread in the data set.
  - The more spread the data, the larger the variance is in relation to the mean.
  - To find the variance, simply square the standard deviation.
  - The symbol for variance is $s^2$.

Variance of visits to the library in the past year

**Data set:** 15, 3, 12, 0, 24, 3

$s$ = 9.18

$s^2$ = **84.3**

# Univariate Descriptive Statistics

- Focus on only **one variable** at a time.
  - It's important to examine data from each variable separately using multiple measures of distribution, central tendency and spread.
  - Programs like SPSS and Excel can be used to easily calculate these.

  - *If you were to only consider the mean as a measure of central tendency, your impression of the "middle" of the data set can be skewed by outliers, unlike the median or mode.*

  - *Likewise, while the range is sensitive to outliers, you should also consider the standard deviation and variance to get easily comparable measures of spread.*

| Visits to the library | |
| --- | --- |
| N | 6 |
| Mean | 9.5 |
| Median | 7.5 |
| Mode | 3 |
| Standard deviation | 9.18 |
| Variance | 84.3 |
| Range | 24 |

# Types of exploratory data analysis

## Univariate non-graphical

- This is simplest form of data analysis, where the data being analyzed consists of just one variable. The main purpose of univariate analysis is to describe the data and find patterns that exist within it.

## Multivariate nongraphical

- Multivariate data arises from more than one variable. Multivariate non-graphical EDA techniques generally show the relationship between two or more variables of the data through cross-tabulation or statistics.

## Univariate graphical

- Graphical methods provide a full picture of the data.

## Multivariate graphical

- Multivariate data uses graphics to display relationships between two or more sets of data.

# Univariate non-graphical

- **Frequency for categorical data**

| Statistic/College | H&SS | MCS | SCS | other | Total |
|---|---|---|---|---|---|
| Count | 5 | 6 | 4 | 5 | 20 |
| Proportion | 0.25 | 0.30 | 0.20 | 0.25 | 1.00 |
| Percent | 25% | 30% | 20% | 25% | 100% |

- **Central Tendency**
  - The three generally estimated are mean, median, and mode.
- **Range**
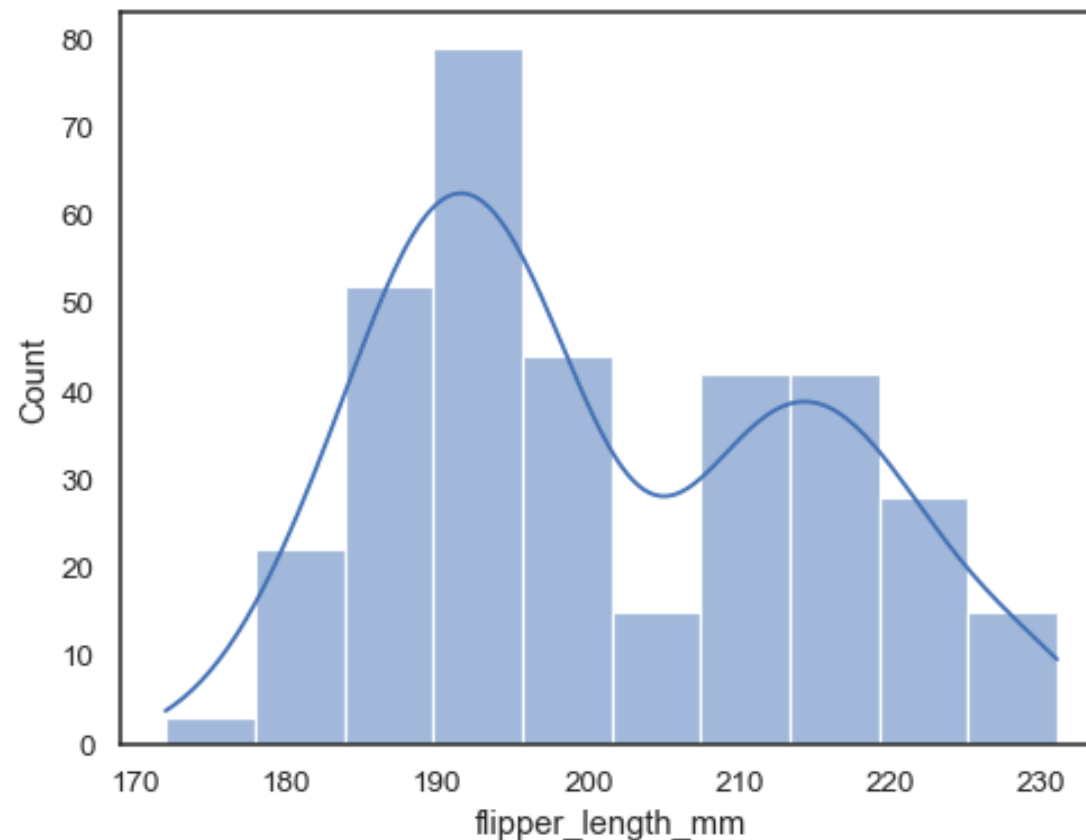  - The range is the difference between the maximum and minimum value in the data.
- **Variance and Standard Deviation**
  - indicates the spread of all data points in a data set.
- **Skewness, Outliers**

# Univariate Graphical

- **Histograms**
  - A bar plot in which each bar represents the frequency (count) or proportion (count/total count) of cases for a range of values.

# Univariate Graphical

- **Box plots**
  - graphically depict the five-number summary of minimum, first quartile, median, third quartile, and maximum.

# Multivariate nongraphical

- **Cross-tabulation**
  - The basic bivariate non-graphical EDA technique

| Subject ID | Age Group | Sex |
|---|---|---|
| GW | young | F |
| JA | middle | F |
| TJ | young | M |
| JMA | young | M |
| JMO | middle | F |
| JQA | old | F |
| AJ | old | F |
| MVB | young | M |
| WHH | old | F |
| JT | young | F |
| JKP | middle | M |

| Age Group / Sex | Female | Male | Total |
|---|---|---|---|
| young | 2 | 3 | 5 |
| middle | 2 | 1 | 3 |
| old | 3 | 0 | 3 |
| Total | 7 | 4 | 11 |

Table 4.2: Cross-tabulation of Sample Data

Table 4.1: Sample Data for Cross-tabulation

# Multivariate nongraphical

- **Correlation coefficient**
  - The correlation between two random variables is a number that runs from -1 through 0 to +1 and indicates a strong inverse relationship, no relationship, and a strong direct relationship, respectively.

$$r_{xy} = \frac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^{n}(x_i - \bar{x})^2}\sqrt{\sum_{i=1}^{n}(y_i - \bar{y})^2}}$$
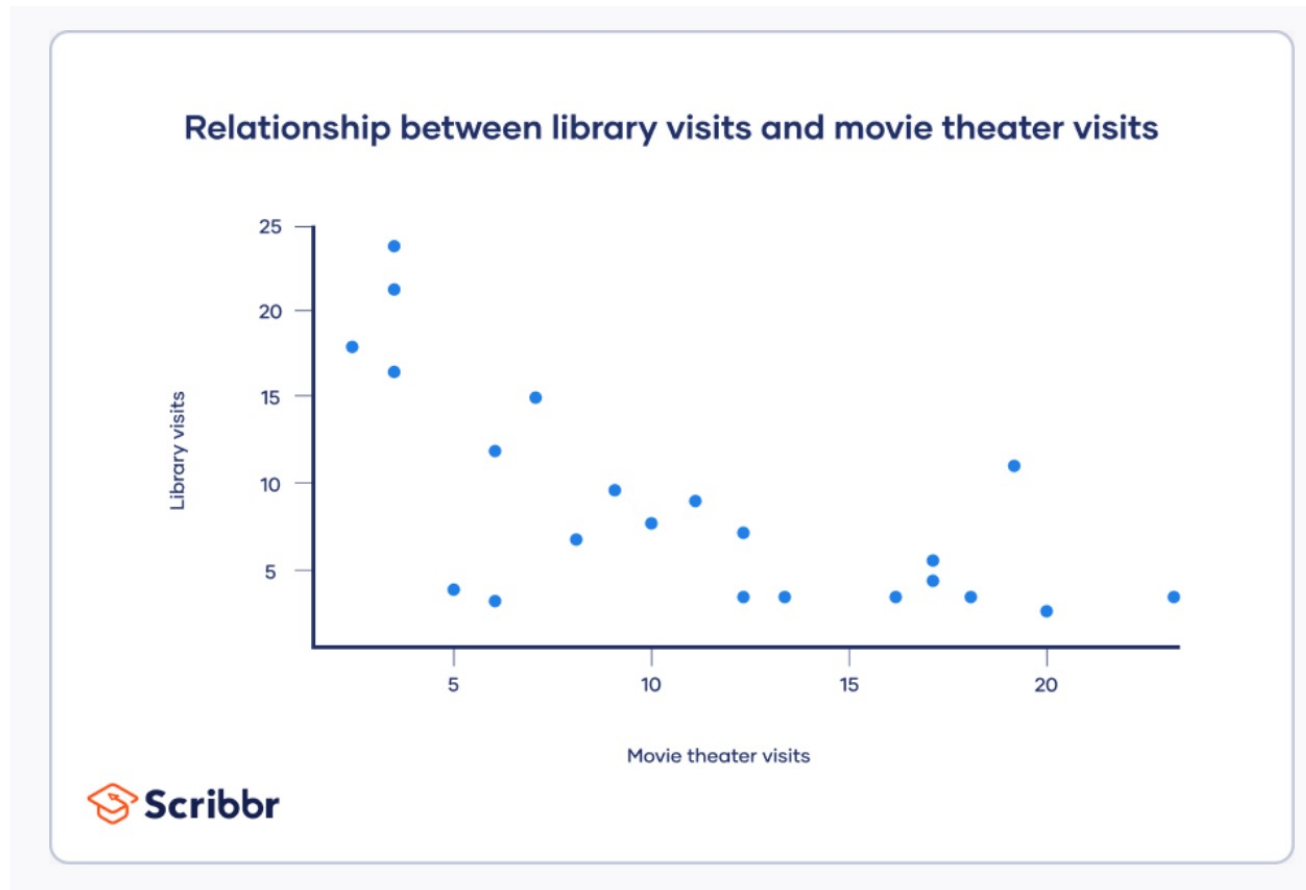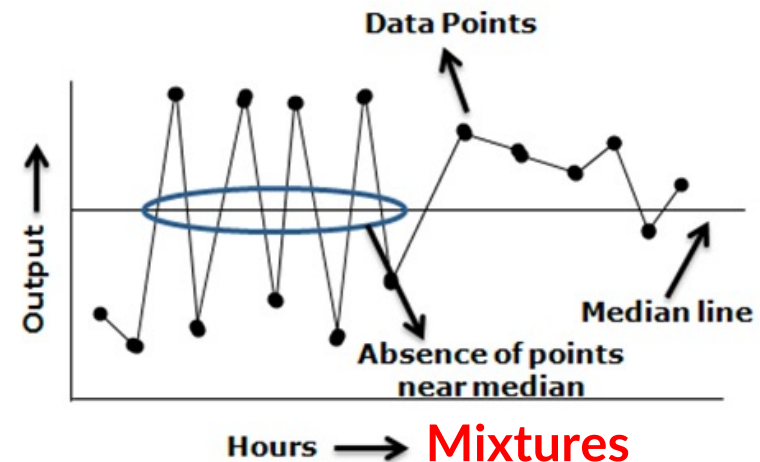
# Multivariate nongraphical

- **Correlation coefficient**



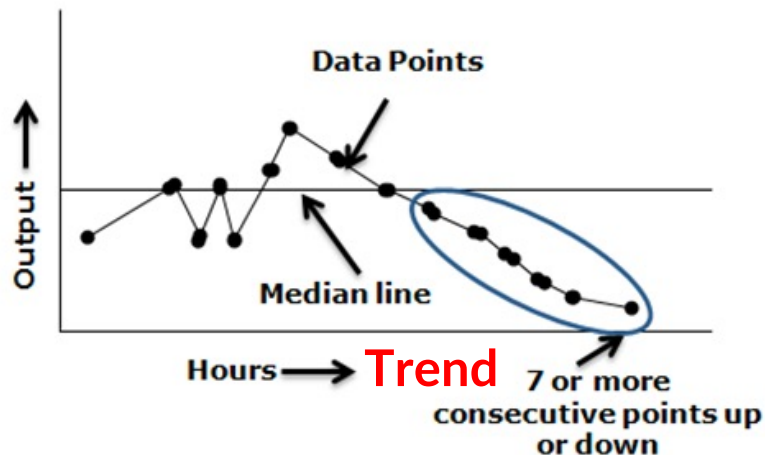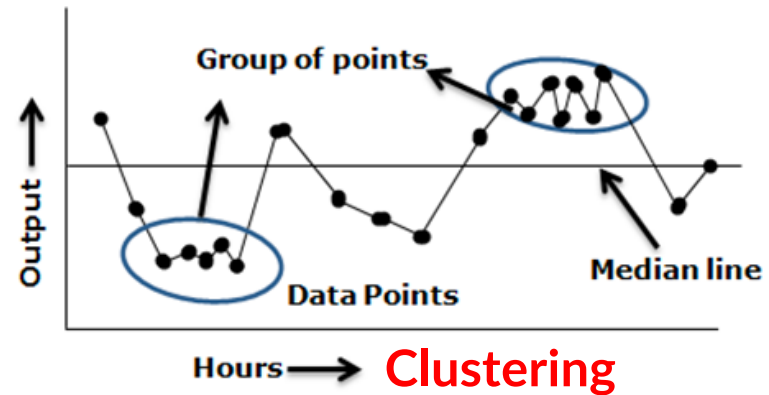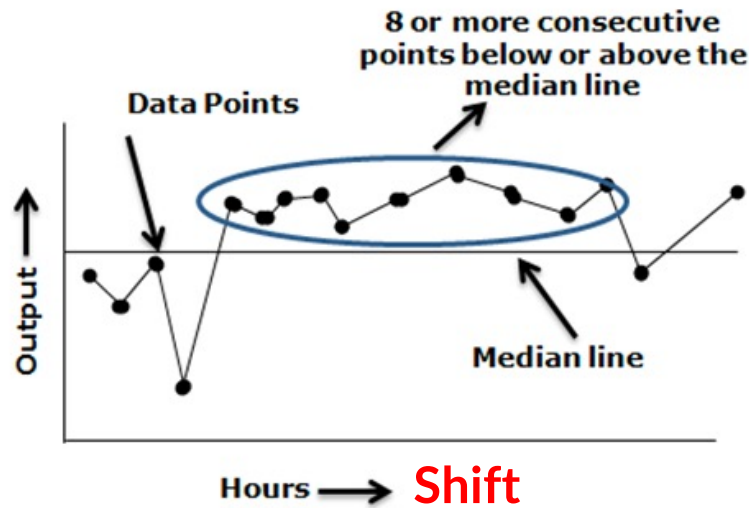https://en.wikipedia.org/wiki/Pearson_correlation_coefficient

# Multivariate Graphical

- **Scatter plot**, plot data points on a horizontal and a vertical axis to show how much one variable is affected by another.
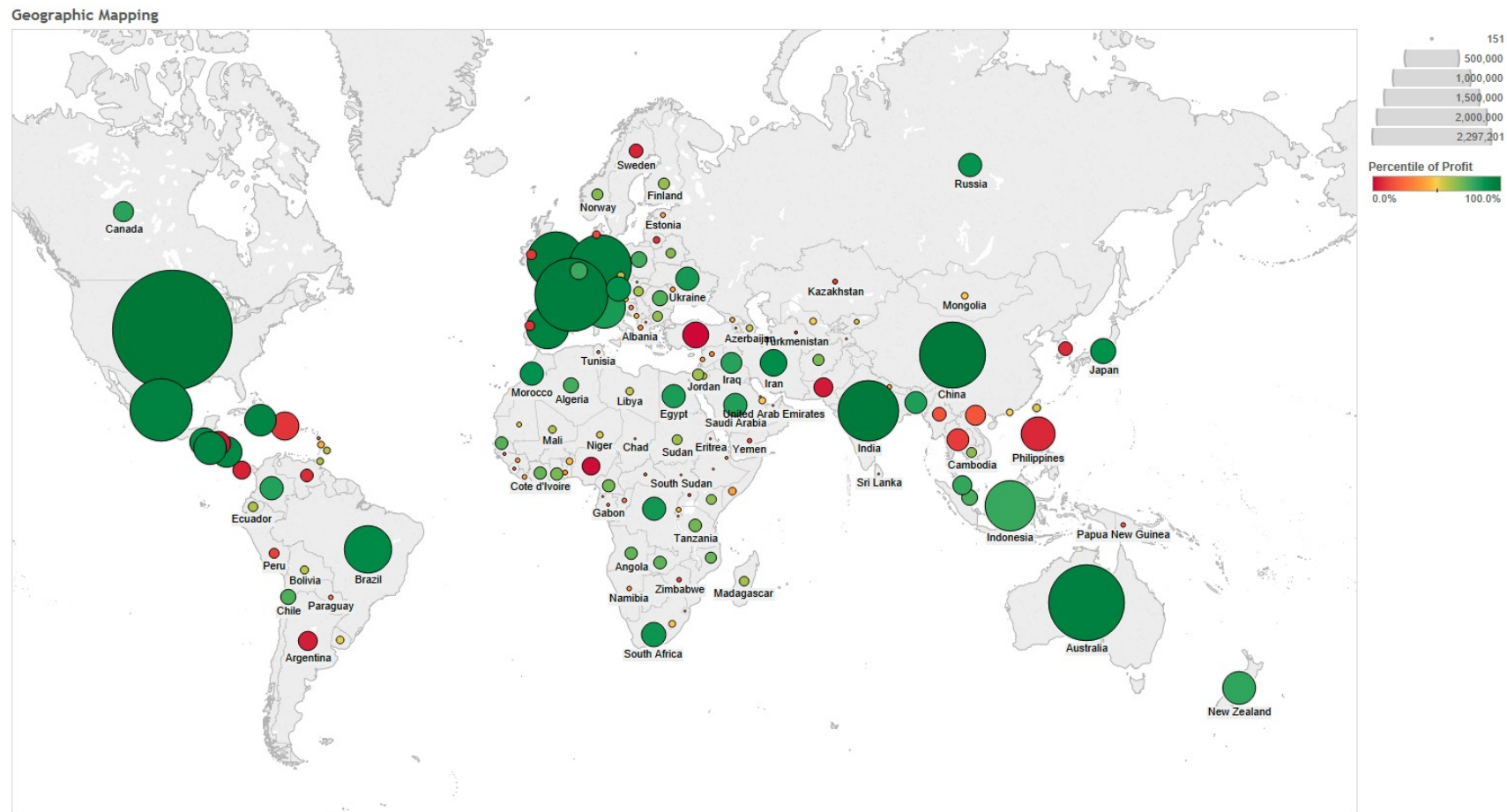
Relationship between library visits and movie theater visits

# Multivariate Graphical

- **Run chart**, which is a line graph of data plotted over time.

# Multivariate Graphical

- **Bubble chart**, which is a data visualization that displays multiple circles (bubbles) in a two-dimensional plot.



https://community.tableau.com/s/idea/0874T000000HAgfQAG/detail

# Multivariate Graphical

- **Heat map**, which is a graphical representation of data where values are depicted by color.
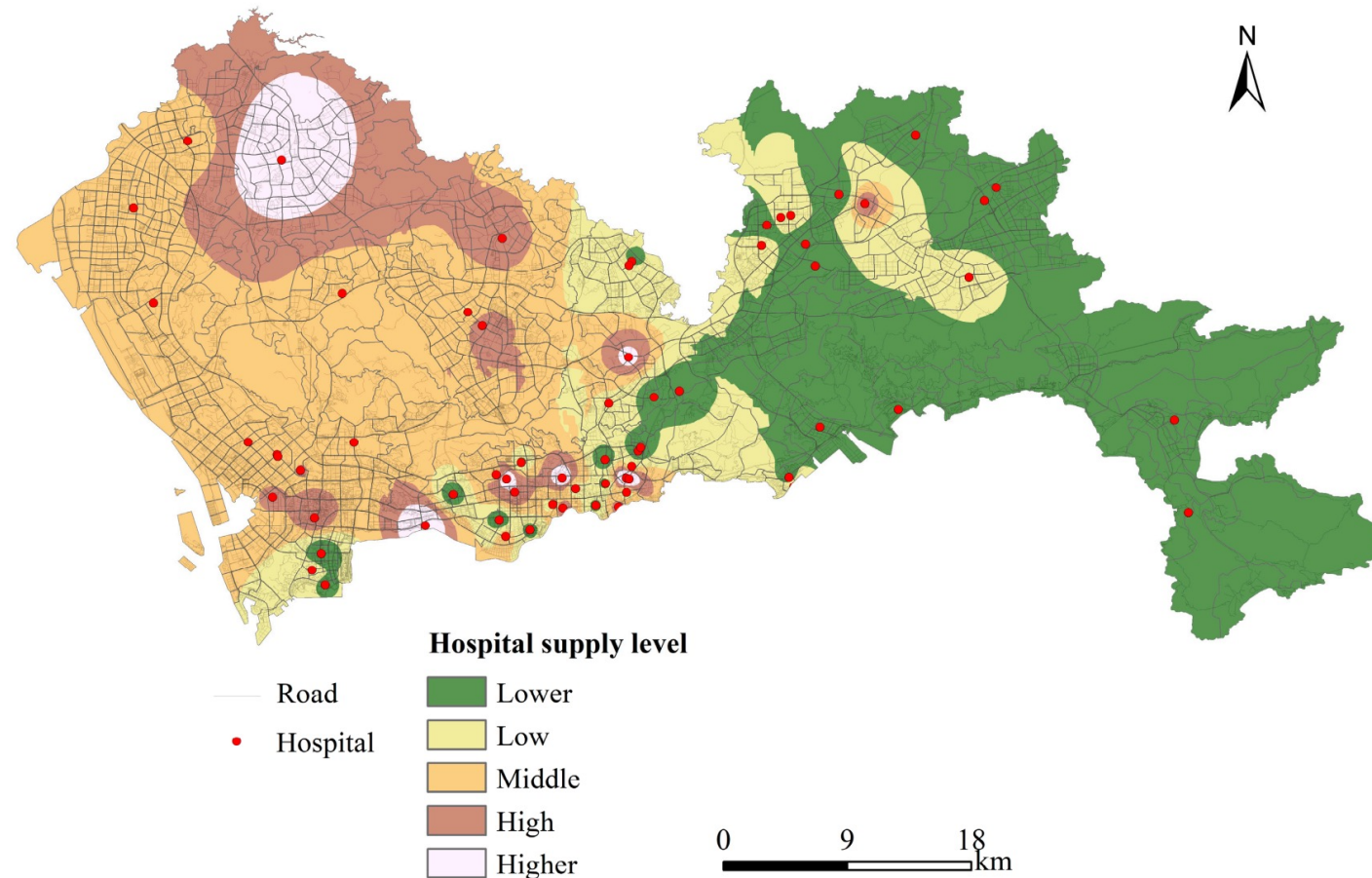


**Figure 4.** Spatial distribution of road network and hospital supply level in Shenzhen.

# Summary of EDA

- You should always perform appropriate EDA before further analysis of your data.

- Perform whatever steps are necessary to become more familiar with your data,
    - check for obvious mistakes,
    - learn about variable distributions, and
    - learn about relationships between variables.

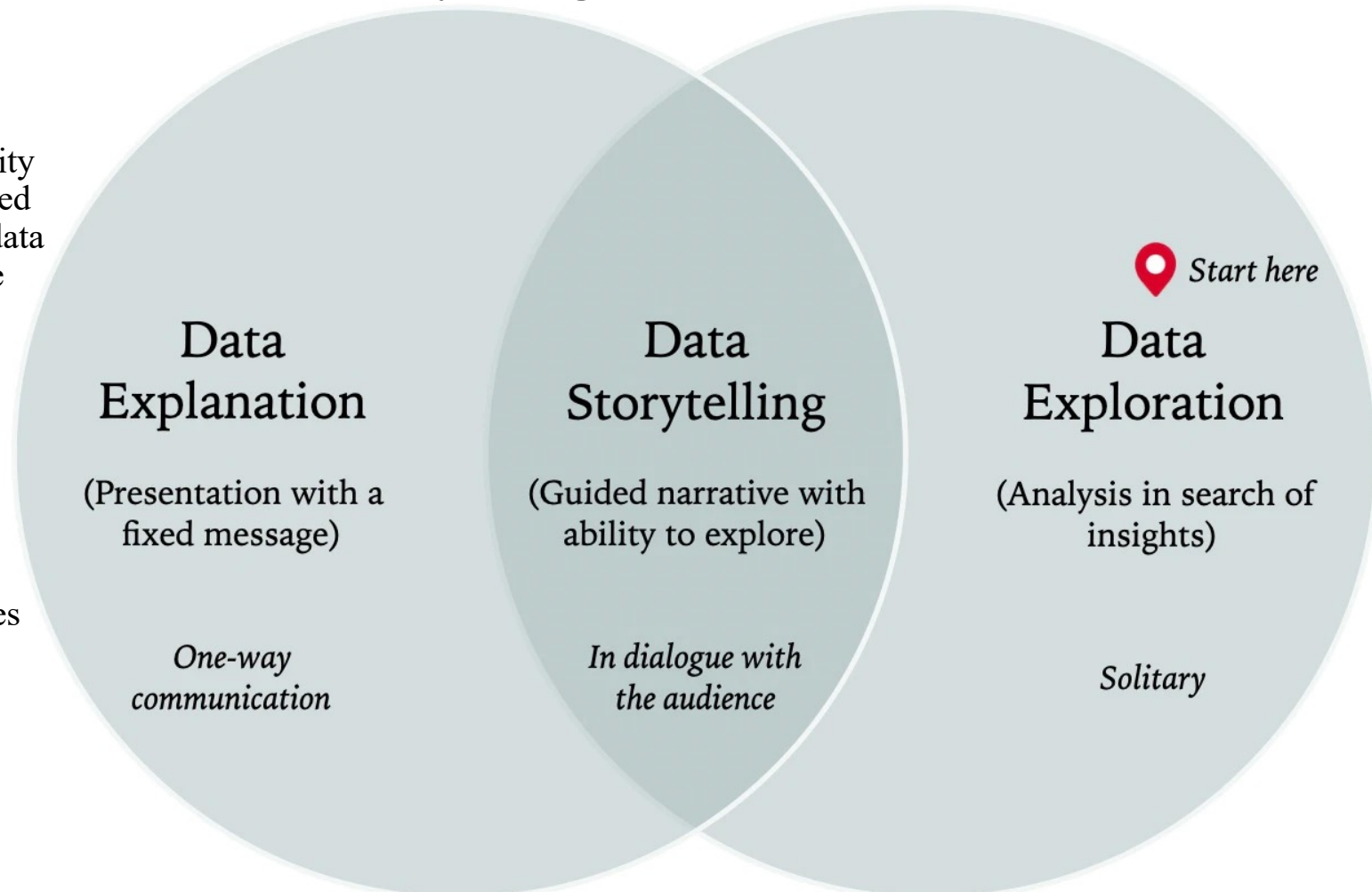- EDA is not an exact science – it is a very important art!

# Practice

- Top 5000 Albums of All Time - Spotify features
    - https://www.kaggle.com/datasets/lucascantu/top-5000-albums-of-all-time-spotify-features
    - https://www.kaggle.com/code/lucascantu/top-5000-spotify

# Can we combine exploratory and explanatory?

- Sure. There is a middle ground that combines data explanation and data exploration. We can call it **interactive data storytelling**.

- At this intersection, there is an opportunity to combine the guided narrative nature of data explanation with the ability to find new insights through exploration.

- Some examples of these three categories on the right.

Data
Explanation

(Presentation with a fixed message)

*One-way communication*

Data
Storytelling

(Guided narrative with ability to explore)

*In dialogue with the audience*

📍 *Start here*

Data
Exploration

(Analysis in search of insights)

*Solitary*

# Thank you~

Wan Fang
Southern University of Science and Technology