# GETTING TO KNOW DATA (PART 1)

SUPAPORN ERJONGMANEE
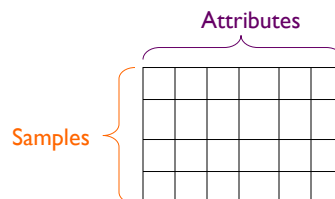
DEPARTMENT OF COMPUTER ENGINEERING
FACULTY OF ENGINEERING
KASETSART UNIVERSITY

1

---

2

## TYPES OF DATA SETS

- Record

    - Relational records

    - Data matrix, e.g., numerical matrix, crosstabs

    - Document data: text documents: term-frequency vector

    - Transaction data

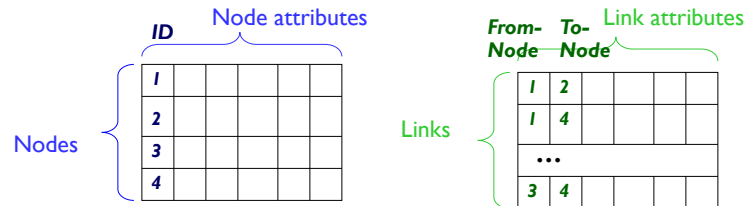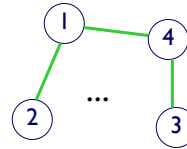| | team | coach | play | ball | score | game | win | lost | timeout | season |
|---|---|---|---|---|---|---|---|---|---|---|
| Document 1 | 3 | 0 | 5 | 0 | 2 | 6 | 0 | 2 | 0 | 2 |
| Document 2 | 0 | 7 | 0 | 2 | 1 | 0 | 0 | 3 | 0 | 0 |
| Document 3 | 0 | 1 | 0 | 0 | 1 | 2 | 2 | 0 | 3 | 0 |

Attributes

Samples

*Source:*
*J. Han, M. Kamber and J. Pei, "Chapter 2 Know Your Data" in Data Mining: Concepts and Techniques, Morgan Kaufmann, July 2011.*

2

# TYPES OF DATA SETS

- Graph and network
  - World Wide Web
  - Social or information networks
  - Molecular Structures



*Source:*
*J. Han, M. Kamber and J. Pei, "Chapter 2 Know Your Data" in Data Mining: Concepts and Techniques, Morgan Kaufmann, July 2011.*

3

# TYPES OF DATA SETS (CONT.)

- Ordered
  - Video data: sequence of images
  - Temporal data: time-series
  - Sequential Data: transaction sequences
  - Genetic sequence data
- Others
  - Spatial data: maps
  - Image data
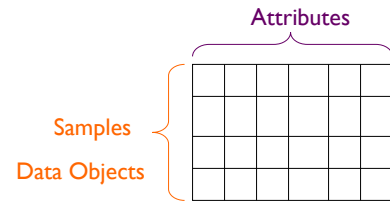  - Video data



*Image source:*
*https://mesquiteproject.wikispaces.com/file/view/DNAMatrix.gif/518627818/DNAMatrix.gif*

*Source (edited):*
*J. Han, M. Kamber and J. Pei, "Chapter 2 Know Your Data" in Data Mining: Concepts and Techniques, Morgan Kaufmann, July 2011.*

4

## Data Objects

- Database columns -> attributes.

- Database rows -> data objects
    - Data sets are made up of data objects.
    - Also called *samples , examples, instances, data points, objects, tuples.*
    - A **data object** represents an entity.

- Examples:
    - sales database: customers, store items, sales
    - medical database: patients, treatments
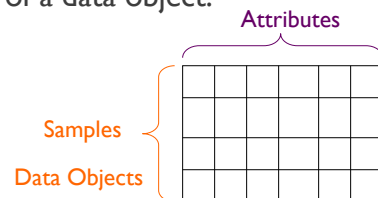    - university database: students, professors, courses

Attributes

Samples

Data Objects

*Source (edited):*
*J. Han, M. Kamber and J. Pei, "Chapter 2 Know Your Data" in Data Mining: Concepts and Techniques, Morgan Kaufmann, July 2011.*

---

## Attributes

- **Attribute (**or **dimensions, features, variables**):
    - a data field, representing a characteristic or feature of a data object.
    - *E.g., customer_ID, name, address*

- Types:
    - Qualitative data
    - Quantitative data

Attributes

Samples

Data Objects

*Source (edited):*
*J. Han, M. Kamber and J. Pei, "Chapter 2 Know Your Data" in Data Mining: Concepts and Techniques, Morgan Kaufmann, July 2011.*

## Data

- A set of values

- Type of data:

  1. Qualitative: characteristic or description data

     - Example: color, gender, country

  2. Quantitative: numerical data

     - Example: height, weight, temperature, area, scores

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 7

Department of Computer Engineering
Kasetsart University

7

## Qualitative Data

- Also call categorical data

- Characteristic or description data

- Immeasurable

- Intervals between values may not be the same

- Can be separated further into

  - Nominal data, Ordinal Data

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 8

Department of Computer Engineering
Kasetsart University

8

## Nominal Data

- Data separated in classes
  - Classes do not always relate to one another
  - Cannot really sort classes (not in order)
- Example:
  - Gender: male, female
  - Regions: America, Asia, Europe
  - Directions: North, East, West, South

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 9

Department of Computer Engineering
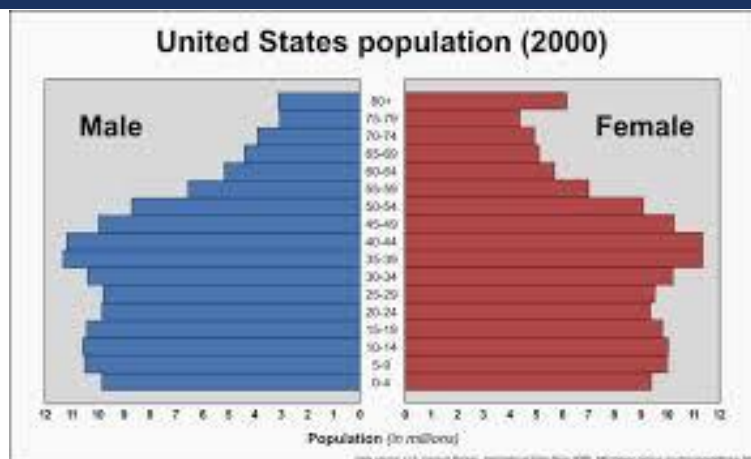Kasetsart University

9

## Nominal Data Example



*Image source: https://encrypted-tbn1.gstatic.com/images?q=tbn:ANd9GcSOwQW4R12eGZJo71pTF-dqPJb7gVY8fSqMevQFNWw_3izPp_gi*

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 10

Department of Computer Engineering
Kasetsart University

10

## Ordinal Data

- Data with ranks
  - Ranks are not actually numerical values (but some can be converted to numbers)
  - Can be sorted
- Immeasurable
- Intervals between values may not be the same
- Example:
  - Size: small, medium, large
  - Satisfaction degree: best, good, poor, worst

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 11

Department of Computer Engineering
Kasetsart University

11

## Ordinal Data Example



Image source: https://www2.barnsley.gov.uk/media/2624867/Customer%20Satisfaction%20graph,%20page%20content%20for%20detail.jpg

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 12

Department of Computer Engineering
Kasetsart University

12

## Quantitative Data

- Measurable

- Intervals between values are the same

- Can be separated further into

  - Interval data, Ratio Data

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 13

Department of Computer Engineering
Kasetsart University

13

## Interval Data

- <u>Ordered</u> numerical values measured in <u>interval</u> with <u>loose zero</u> point

  - Mostly used differences (addition/subtraction) to compare.

  - Cannot be directly compared in ratio (division/multiplication)

- Example:

  - Temperature

  - Times

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 14

Department of Computer Engineering
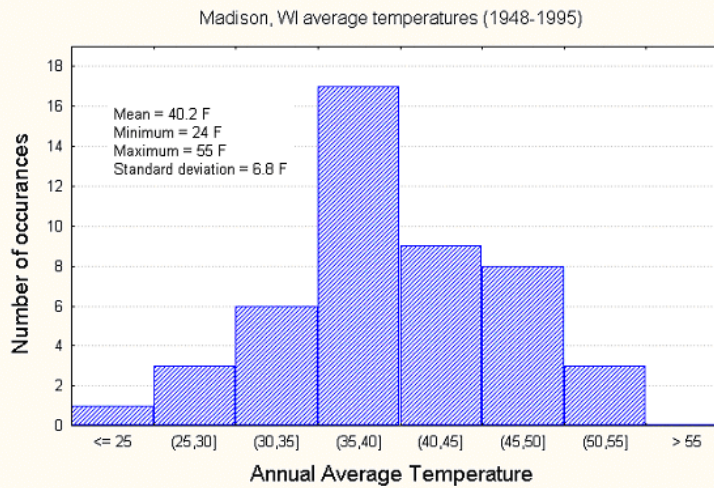Kasetsart University

14

## Interval Data Example



Madison, WI average temperatures (1948-1995)

Mean = 40.2 F
Minimum = 24 F
Maximum = 55 F
Standard deviation = 6.8 F

*Image source: http://itg1.meteor.wisc.edu/wxwise/AckermanKnox/chap14/madisonhist.gif*

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 15

Department of Computer Engineering
Kasetsart University

15

## Ratio Data

- Measurable

- Intervals between values are the same.

- Can be computed using

  - Difference (addition/subtraction)

  - Ratio (multiplication/division)

- Example:

  - Weight, Length, Revenue

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 16

Department of Computer Engineering
Kasetsart University

16

## Ratio Data Example

**Heathrow Temperature Forecast**
Generated at: 19 Jan 12:00 UTC
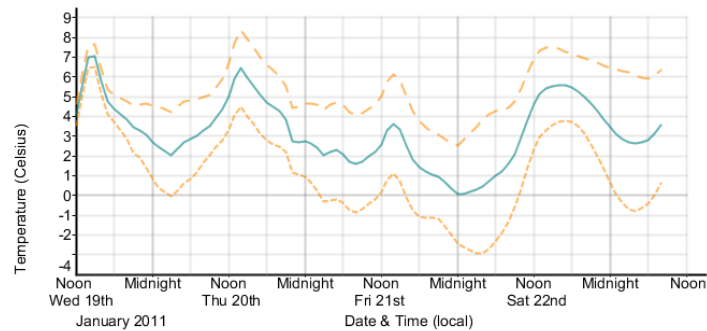Best Forecast    5% Confidence    95% Confidence

*Image source: http://www.metraweather.com/~metracom/sites/default/files/Hourly_Forecast_Temperature_Heathrow.png*

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 17

Department of Computer Engineering
Kasetsart University

17

## Comparison: Type of Data

■ Comparison

|  | Nominal | Ordinal | Interval | Ratio |
|---|---|---|---|---|
| Can order values |  | ✓ | ✓ | ✓ |
| Can compute differences of values |  |  | ✓ | ✓ |
| Can add or subtract values |  |  | ✓ | ✓ |
| Can divide or multiple values |  |  |  | ✓ |
| Has fixed zero points |  |  |  | ✓ |

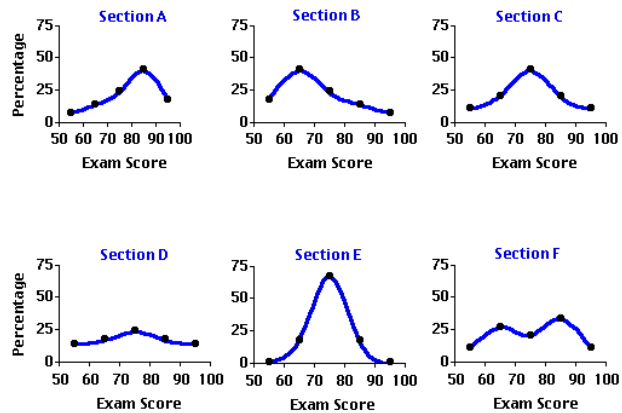*Image source: https://www.mymarketresearchmethods.com/types-of-data-nominal-ordinal-interval-ratio/*

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 18

Department of Computer Engineering
Kasetsart University

18

## Type of Data (2)

1. Discrete Data
   - Countable values (positive, zero, negative)
   - Can be either numerical or categorical data
   - Can be finite or infinite sequences
2. Continuous Data
   - Specific value in ranges
   - Can be finite or infinite ranges
   - Ranges can be joint or disjoint.

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 19

Department of Computer Engineering
Kasetsart University

19

## Questions to Ask about Data

- What are my data ?

- What are the attributes of data?

  - For each attribute, what is its type?

    - Quantitative, Nominal, Ordinal, Interval

- Data type affects computation

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 20

Department of Computer Engineering
Kasetsart University

20

## Measure of Central Tendency

- Mode
- Median
- Mean



Section A, Section B, Section C, Section D, Section E, Section F — Percentage vs Exam Score

*Image source: http://vassarstats.net/textbook/f0203x.gif*

Supaporn Erjongmanee
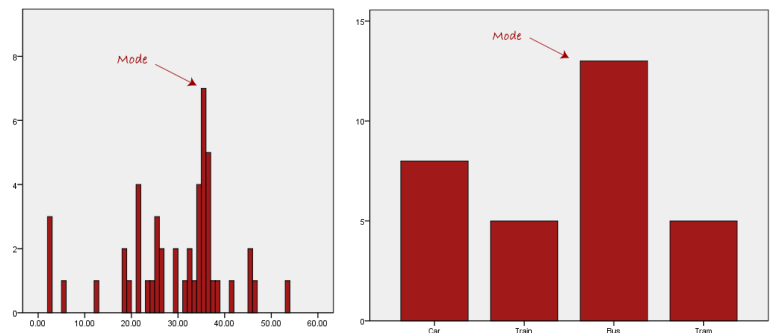fengspe@ku.ac.th

Getting to Know Data
Slide 21

Department of Computer Engineering
Kasetsart University

21

## Mode

- Most frequent value
- Easily spot as a peak in histogram
- Good for nominal data



Mode

Mode — Car, Train, Bus, Tram

*Source: [2]*

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 22

Department of Computer Engineering
Kasetsart University

22

## Mode (cont.)

- Be careful when used with continuous data
  - Difficult to specify detailed value (e.g., 65.3)
- Avoid if mode is not with the majority

Getting to Know Data
Slide 23

Department of Computer Engineering
Kasetsart University

23

---

## Median

- Middle value in the sorted data
- Can be used with outliers or skewed data

Let n = data size

median = $(\frac{n+1}{2})^{th}$ value

2, 2, 5, 7, 8, 10, 12

**Median = 7**

0, 2, 5, 6, 7, 8, 8, 8, 9, 10

Median = $\frac{7+8}{2}$ = $\frac{15}{2}$ = 7.5

Getting to Know Data
Slide 24

Department of Computer Engineering
Kasetsart University

24

## Mean

- Most commonly used to measure of center
- Can be used for both discrete and continuous data

$$\text{Mean} = \frac{\sum_{i=1}^{n} x_i}{n}$$

- Every value takes part in calculation
- Often stand for <u>typical value</u>
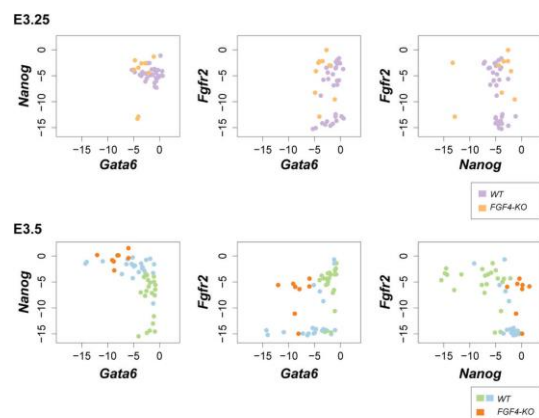  - Located at center
- Minimize error in predicting other values

*Source: [2]*

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 25

Department of Computer Engineering
Kasetsart University

## Additional Measurement of Central Tendency

- Harmonic mean
  - Generally use for average rate

$$\text{Harmonic mean} = \frac{n}{\sum_{i=1}^{n} \frac{1}{x_i}}$$

- Geometric mean
  - Generally, use for average compound growth rate

$$\text{Geometric mean} = \sqrt[n]{a_0 a_1 a_2 \ldots a_{n-1}}$$

$a_0 = 1 + r_0$

$a_1 = 1 + r_1$

...

$a_{n-1} = 1 + r_{n-1}$

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 26

Department of Computer Engineering
Kasetsart University

## Time to pick suitable measure of central tendency

- Recommending….

| Type of Data | Measure of Central Tendency |
|---|---|
| Nominal, Categorical | Mode |
| Ordinal | Median |
| Interval & Ratio (not skewed) | Mean |
| Interval & Ratio (skewed) | Median |

*Source: [2]*

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 27

Department of Computer Engineering
Kasetsart University

27

## Measure of Variability

- Range
- Variance
- Standard deviation
- Coefficient of variation
- Mean absolute deviation
- Inter-quartile range



*Image source: http://www.nature.com/ncb/journal/v16/n1/images/ncb2881-sf6.jpg*

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 31

Department of Computer Engineering
Kasetsart University

31

## Range

- Simplest form of variability measurements
- Beware of outliers

$$\text{Range} = \text{Max} - \text{Min}$$

12, 25, 27, 29, 36, 38, 40, 43, 50, 54, 62

Range = 62 - 12 = 50

*Image source: http://www.regentsprep.org/regents/math/algtrig/ats1/Range.gif*

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 32

Department of Computer Engineering
Kasetsart University

32

---
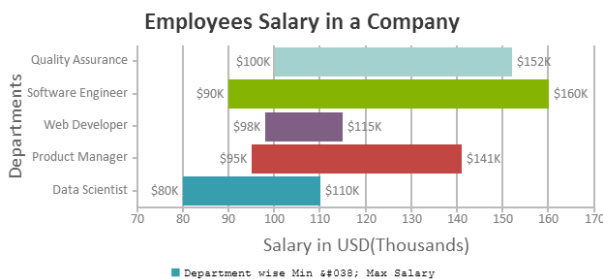
## Range Plot

- Compare minimum and maximum for multiple items



**Employees Salary in a Company**

**Min/Max Air Temperature Raleigh-Durham Airport December 1996**

*Image source: http://canvasjs.com/wp-content/uploads/images/gallery/javascript-range-column-range-bar-charts/employee-salary.jpg*
*https://www.ncsu.edu/labwrite/res/gh/rangebar-vert.gif*

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 33
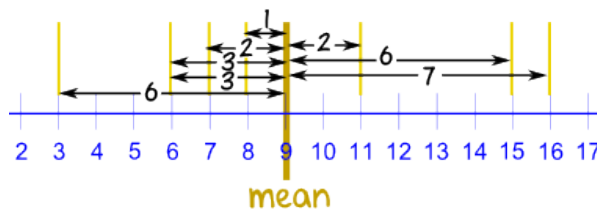
Department of Computer Engineering
Kasetsart University

33

## Variance & Standard Deviation

- Average *difference (squared distance)* between all data and the mean
- Fit for
  - Continuous data
  - Quantitative data, not categorical data
- Avoid if data are skewed or have outliers
- Unit of variance is squared
- Standard deviation = $\sqrt{variance}$

Population variance

$$\sigma^2 = \frac{\sum_{i=1}^{n}(x_i - \mu)^2}{n}$$

Sample variance

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

Sample variance (divided by n -1) is unbiased estimate of population variance

34

---

## Standard Deviation for Normal Data

- How does standard deviation tell us about spread of normal data?

35

## Coefficient of Variation

■ Ratio between standard deviation and mean

Population

$$\text{Coefficient of variation} = \frac{\sigma}{\mu}$$

$$cv_{student} = \frac{s}{\bar{x}} = \frac{6.22}{174.54} = 0.0356$$

Sample

$$\text{Coefficient of variation} = \frac{s}{\bar{x}}$$

**vs.**

$$cv_{elephant} = 15.6$$

Elephants have more variability in height than student

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 36

Department of Computer Engineering
Kasetsart University

36

## Mean Absolution Deviation

■ Average *absolute distance* between all data and the mean

$$\text{Mean absolution deviation} = \frac{\sum_{i=1}^{n} |x_i - \bar{x}|}{n}$$



*Image source: http://www.mathsisfun.com/data/images/mean-deviation.gif*

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 37

Department of Computer Engineering
Kasetsart University

37

# Quartile

- Values that split data into 4 equal partitions



100%

75%

50%

25%

Min

Max

*Quartile 1*

*Quartile 3*

*Quartile 2*

*Median*

If it is regular number, call percentile

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 38

ngineering
Kasetsart University

38

---

# Interquartile Range (IQR)

- Good to use for data with outlier or skewed

- Not consider all data
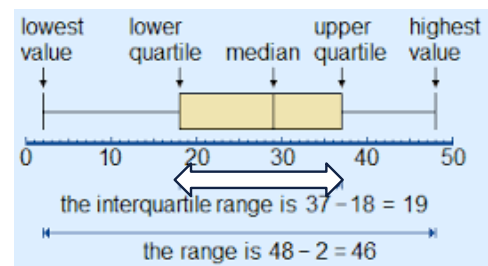
Q1 = lower quartile = 25% percentile
Q3 = upper quartile = 75% percentile
IQR = Q3 - Q1



lowest value    lower quartile    median    upper quartile    highest value

0    10    20    30    40    50

the interquartile range is 37 − 18 = 19

the range is 48 − 2 = 46

*Image source: http://www.dr-aart.nl/Statistiek_bestanden/boxplot1EN.png*

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 39

Department of Computer Engineering
Kasetsart University

39

## SD vs. IQR

- Comparison between normal distribution and box plots



*Image source: https://upload.wikimedia.org/wikipedia/commons/thumb/1/1a/Boxplot_vs_PDF.svg/250px-Boxplot_vs_PDF.svg.png*

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 40

Department of Computer Engineering
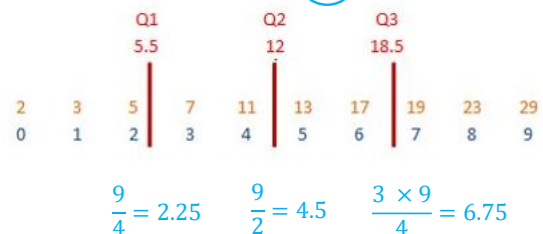Kasetsart University

40

---

## Quartile in Excel

- Quartile.exc vs. Quartile.inc



*Source: http://datapigtechnologies.com/blog/index.php/why-excel-has-multiple-quartile-functions-and-how-to-replicate-the-quartiles-from-r-and-other-statistical-packages/*

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 41

Department of Computer Engineering
Kasetsart University

41

## Which Measurement of Variability to Use?

- ■ Range ➡ Easiest to use. Not suitable for data with outliers

- ■ Sample variance

- ■ Sample standard deviation

  Most commonly-used

- ■ Inter quartile range ➡ Good for data with outliers

- ■ Coefficient of variation ➡ Tell more story about the data:
  how std is compared to the mean
  No unit
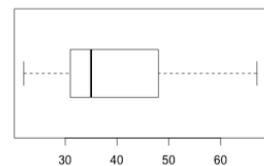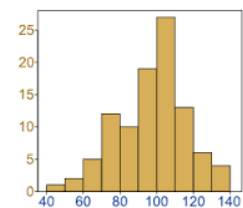  Sensitive when mean $\rightarrow 0$
  Not suitable for multiple replicates of data

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 42

Department of Computer Engineering
Kasetsart University

42

---

## OUTLIERS

- ■ Out-of-the-norm data

- ■ Threshold is needed to cut outliers



*Source: [2]*

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 43

Department of Computer Engineering
Kasetsart University

43

## Outliers

- Outliers can be determined from IQR or SD



If data value > Q3 + 1.5IQR  or
   data value < Q1 - 1.5IQR,
we consider such value to be outlier.

If data value > mean + 3SD or
   data value < mean -3SD,
we consider such value to be outlier.

*Image source: https://upload.wikimedia.org/wikipedia/commons/thumb/1/1a/Boxplot_vs_PDF.svg/250px-Boxplot_vs_PDF.svg.png*

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 44

Department of Computer Engineering
Kasetsart University

44

---

## Are descriptive statistics enough?

- <u>Descriptive statistics are not answer to everything</u>

- Be careful of outlier and skewed data

- Always GRAPH your data

   - Histogram

   - Boxplot



*Image sources: https://www.mathsisfun.com/data/images/histogram.gif*
*http://www.johnquarto.com/wp-content/uploads/2013/09/Boxplot-PartyPeopleAll.png*

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 45

Department of Computer Engineering
Kasetsart University

45

## Basic Data Visualization

- Histogram
- Boxplot
- Scatter plot

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 46

Department of Computer Engineering
Kasetsart University

46

## Histogram

- Specific bar graph representing <u>distribution of data</u>
- x-axis: bins of data values
- y-axis: <u>frequency</u> of data values
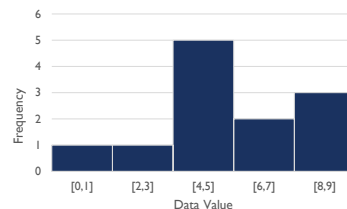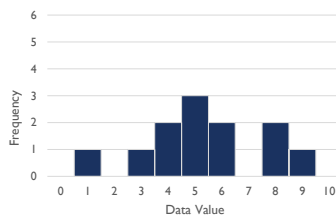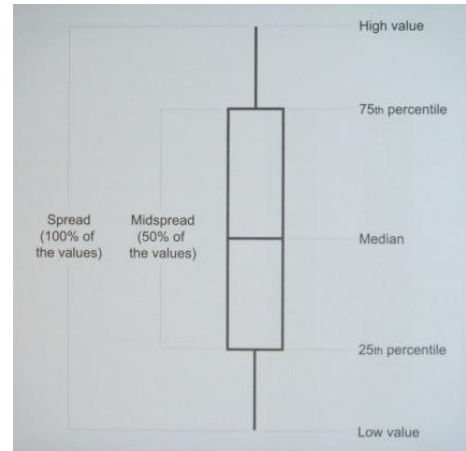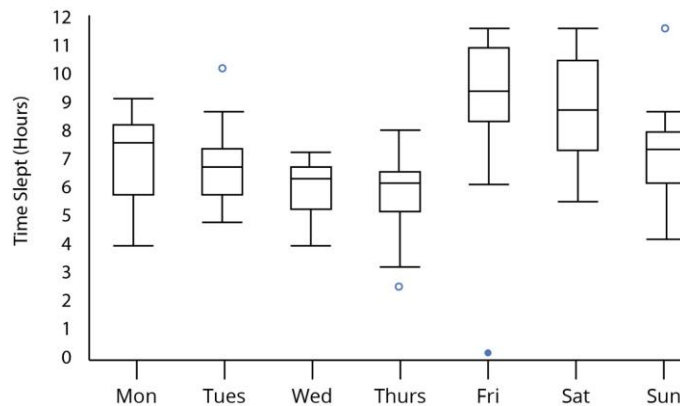- Example:

8 5 4 9 6 5
6 3 8 1 4 5





*Image source: https://openclipart.org*

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 47

Department of Computer Engineering
Kasetsart University

47

# Box Plots

- Also call box-and-whisker plot
- Use statistical values to plot distribution of data



*Image source: Figures 10.34 [1]*

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 48

Department of Computer Engineering
Kasetsart University

48

---

# Box Plots (cont.)



Outliers

75th percentile

Median

25th percentile

Outliers

- More detailed box plot

Definition of "outliers" must be given.

*Image source: Figure 10.38, [1]*

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 49

Department of Computer Engineering
Kasetsart University

49

## Boxplot (cont.)

- Compare multiple data sets

- Example:

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 50

Department of Computer Engineering
Kasetsart University

50

## Scatter Plot

- To visualize relationship between multiple variables

- To measure relationship, we use <u>correlation</u>

- Type of relation

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 51

Department of Computer Engineering
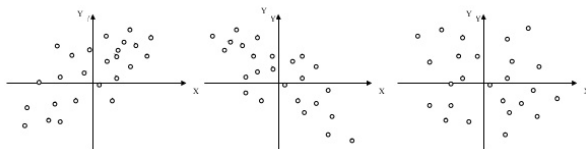Kasetsart University

51

## Correlation

$$\rho_{X,Y} = \frac{Cov(X,Y)}{\sigma_X \sigma_Y} = \frac{\sum_x \sum_y (x - \mu_x)(y - \mu_y)\, p(x,y)}{\sigma_X \sigma_Y}$$

- Range of $\rho_{X,Y}$:  $-1 \leq \rho_{X,Y} \leq 1$



Positive correlation      Negative correlation      No correlation

We use scatter plot to visualize correlation

*Image source: http://www.slideshare.net/AhmedShahid/t-tests-anovas-and-regression*

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 52

Department of Computer Engineering
Kasetsart University

52

---

## Correlation

$$\rho_{X,Y} = \frac{Cov(X,Y)}{\sigma_X \sigma_Y} = \frac{\sum_x \sum_y (x - \mu_x)(y - \mu_y)\, p(x,y)}{\sigma_X \sigma_Y}$$

Correlation does not imply causation.

- Range of $\rho_{X,Y}$:  $-1 \leq \rho_{X,Y} \leq 1$



Positive correlation      Negative correlation      No correlation

Correlation tells how two values track each other.

If X increases, how about Y?

They may be hidden factor

*Image source: http://www.slideshare.net/AhmedShahid/t-tests-anovas-and-regress...*

Supaporn Erjongmanee
fengspe@ku.ac.th
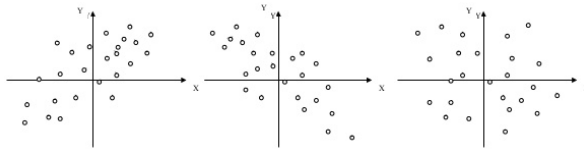
Getting to Know Data
Slide 53

Department of Computer Engineering
Kasetsart University

53

## Sample Correlation

$$\hat{\rho}_{X,Y} = \frac{1}{n-1} \sum_{i=1}^{n} \left( \frac{(x - \bar{x})(y - \bar{y})}{s_X s_y} \right)$$

- Range of $\hat{p}_{X,Y}$:  $-1 \leq \hat{p}_{X,Y} \leq 1$



Positive correlation     Negative correlation     No correlation

*Image source: http://www.slideshare.net/AhmedShahid/t-tests-anovas-and-regression*

Supaporn Erjongmanee
fengspe@ku.ac.th
Getting to Know Data
Slide 54
Department of Computer Engineering
Kasetsart University

54

---

## Correlation

- Strength



Strong
[0.6,0.8)

Weak
[0.2,0.4)

None
[0.0,0.2)

Very strong
[0.8,1.0]

Moderate
[0.4,0.6)

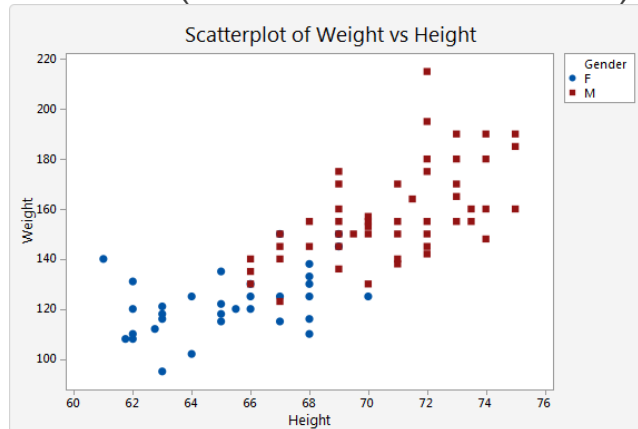Source: http://www.statstutor.ac.uk/resources/uploaded/pearsons.pdf
Image source: https://support.minitab.com/en-us/minitab-express/1/help-and-how-to/graphs/scatterplot/create-the-graph/choose-a-scatterplot/

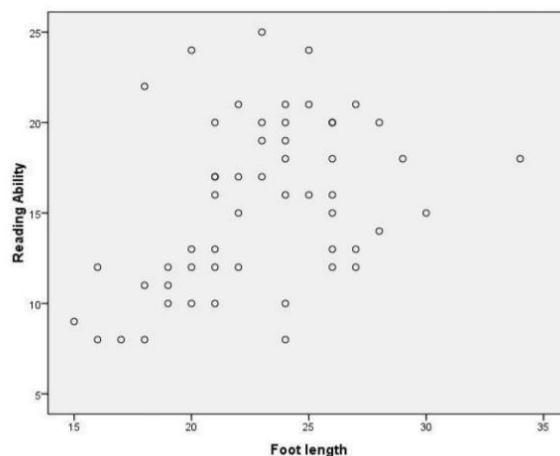Supaporn Erjongmanee
fengspe@ku.ac.th
Getting to Know Data
Slide 55
Department of Computer Engineering
Kasetsart University

55

## Correlation

- Example: 3 variables (share same X, Y variables)



Scatterplot of Weight vs Height

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 56

Department of Computer Engineering
Kasetsart University

56

## Correlation

- Example 3:

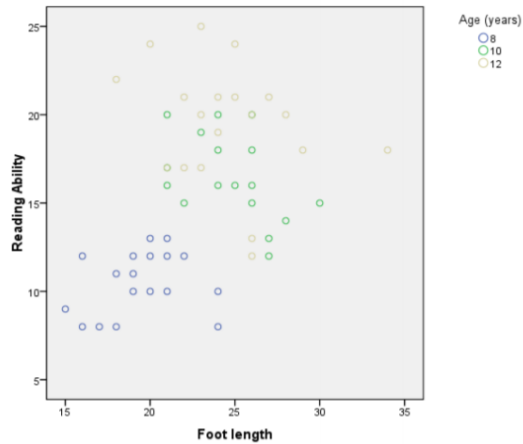Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 58

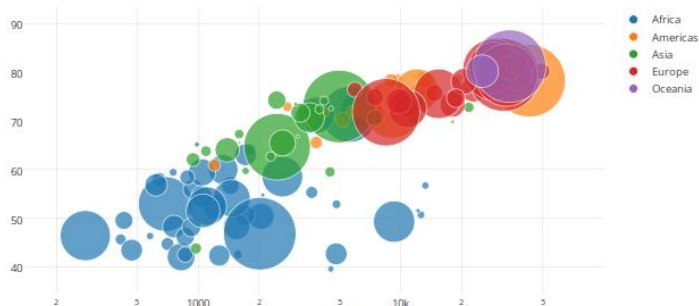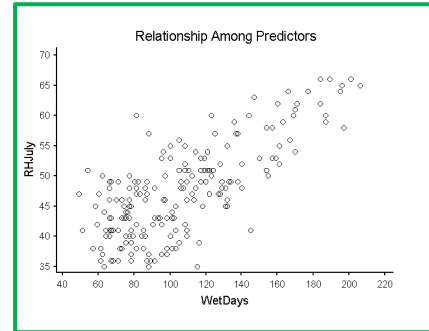Department of Computer Engineering
Kasetsart University

58

## Correlation

■ Example 3 (cont.):

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 59

Department of Computer Engineering
Kasetsart University

59

---

## Correlation

■ Example 4: >=3 variables

■ Often use bubble scatter plot (4 variables) or scatter plot matrix

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 60

Department of Computer Engineering
Kasetsart University
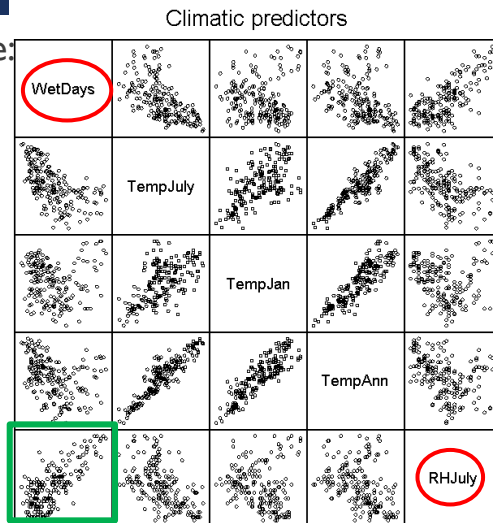
60

# SCATTERPLOT MATRIX

- Example:



Climatic predictors

*Image source:*
*https://www.pcord.com/nscatterplot.htm*
*https://www.pcord.com/nscatterplotmatrix.htm*

Supaporn Erjongmanee
fengspe@ku.ac.th

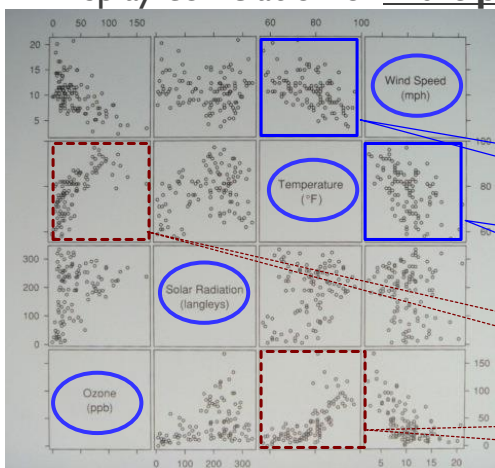Getting to Know Data
Slide 61

Department of Computer Engineering
Kasetsart University

61

---

# SCATTERPLOT MATRIX (CONT.)

- Display correlation of **multiple pairs** of variables in the same time.



4 variables:
- Wind Speed
- Temperature
- Solar Radiation
- Ozone

X = temperature, Y = wind speed

X = wind speed, Y = temperature

Same pair of variables

X = ozone, Y = temperature

X = temperature, Y = ozone

Same pair of variables

Only need half of scatterplot matrix

*Source: Figure 11.22, [1]*

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 62

Department of Computer Engineering
Kasetsart University

62

## Summary

- To first explore data, we can find

  - Outliers

  - Centrality: Mean, Median, Mode

  - Variability: Range, Variance, Standard Deviation, Coefficient of Variation, Mean Absolute Deviation, Interquartile Range

  - Correlation

  - Visualization: Histogram, Boxplot, Scatter Plot

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 63

Department of Computer Engineering
Kasetsart University

63

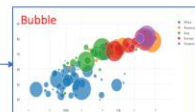## Summary (cont.)



*Image Source: https://www.netsolutions.com/insights/data-visualization-a-building-block-of-an-intelligent-enterprise/*

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 64

Department of Computer Engineering
Kasetsart University

64

## Reference

1. *http://www.socialresearchmethods.net/kb/sampprob.php*
2. *https://statistics.laerd.com/statistical-guides/measures-central-tendency-mean-mode-median.php*
3. *http://blog.minitab.com/blog/michelle-paret/using-the-mean-its-not-always-a-slam-dunk*
4. *https://statistics.laerd.com/statistical-guides/measures-of-spread-standard-deviation.php*
5. J.L. Devore and K.N.Berk, Modern Mathematical Statistics with Applications, Springer, 2012
6. *https://support.office.com/en-sg/article/Add-change-or-remove-a-trendline-in-a-chart-072d130b-c60c-4458-9391-3c6e4b5c5812*

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 65

Department of Computer Engineering
Kasetsart University

65

## References

6. Effectively Communicating Numbers: Selecting the Best Means and Manner of Display, Stephen Few, Principal, Perceptual Edge, 2005
7. A.L. Leon-Garcia, Probability and Random Processes for Electrical Engineering, Addison-Wesley, 1994.

Supaporn Erjongmanee
fengspe@ku.ac.th

Getting to Know Data
Slide 66

Department of Computer Engineering
Kasetsart University

66