

Data; Information; Knowledge; Wisdom

Raw; Meaning; Context; Applied

Table; bar chart; pie chart; stem and leaf plot

Stem	Leaf
6	0 1 1 4 4 4 4 6 6 8 8 8 9

Mode: The value appears most often

Range: Min to max

pth: p% are smaller

Interquartile range: First quartile=25%

LSF: $b = \frac{\sum_{i=1}^n (x_i - \bar{X})(y_i - \bar{Y})}{\sum_{i=1}^n (x_i - \bar{X})^2}$ covariance

Least square fit $\sum_{i=1}^n (x_i - \bar{X})^2$ variance

$a = \bar{Y} - b\bar{X}$ Fit a line $Y = a + bX$ such that it minimizes the error S

Correlation Coefficient: >0 正 <0 负 ≈0 不相关

$r = \frac{\sum_{i=1}^n (x_i - \bar{X})(y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{X})^2 \sum_{i=1}^n (y_i - \bar{Y})^2}}$

Var(X) = $E((X - \mu)^2) = E(X^2) - (E(X))^2 = \sigma^2$

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

Variance of a sample: $\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{X})^2$ Standard deviation of a sample: $\sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{X})^2}$

Binomial Probability Model 二项分布 P53

成功的次数

$P(X=x) = C_n^x P^x (1-P)^{n-x} \quad X \sim bin(n, p)$

mean: $\mu = n \cdot p$

n 次独立重复的伯努利实验

standard deviation: $\sigma = \sqrt{np(1-p)}$

Poisson Probability Model 泊松分布 P56

 $P(X=x) = e^{-\lambda} \frac{\lambda^x}{x!}$

$X \sim Po(\lambda)$

入: 单位时间内事件的平均发生次数

$\mu = \lambda$ 平均值(也等于期望)

 $\sigma = \sqrt{\lambda}$ 标准差

单位时间内随机事件发生的次数

个数: 泊松分布函数

正态分布: 标准正态分布

Expectation and Variance of Continuous Random Variables 定义和离散的不

期望值: $E(X) = \mu = \int_{-\infty}^{\infty} xf(x)dx$ 方差: $Var(X) = \sigma^2 = \int_{-\infty}^{\infty} (x - \mu)^2 f(x)dx$

Normal Probability Model / Gaussian Distribution 正态分布

$y = f_x(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad X \sim N(\mu, \sigma^2)$

Standard Normal Distribution 标准正态分布 P65

$y = f_x(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} \quad X \sim N(0, 1)$

如果 X 服从 $N(\mu, \sigma^2)$, 那么 $Z = \frac{X-\mu}{\sigma}$ 服从 $N(0, 1)$

Confidence Intervals 置信区间 P71

95% confidence intervals is $(-1.96, 1.96)$ 计算真实值: $M \pm 2 \frac{\sigma}{\sqrt{n}}$ → 样本数Uniform Distribution 匀分布 $X \sim U(a, b)$

$y = f_x(x) = \begin{cases} \frac{1}{b-a}, & a \leq x \leq b \\ 0, & \text{otherwise} \end{cases}$

x: 均匀分布期望即均值

Mean: 代入 $E(X)$ 通式 $= \frac{b+a}{2}$

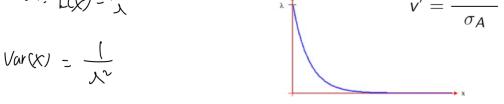
Variance: 代入 $= \frac{(b-a)^2}{12}$

Exponential Distribution 指数分布 $X \sim Exp(\lambda)$

$y = f_x(x) = \begin{cases} \lambda e^{-\lambda x}, & x > 0 \\ 0, & x \leq 0 \end{cases}$

Mean: $E(X) = \frac{1}{\lambda}$

Var(X) = $\frac{1}{\lambda^2}$



贝叶斯公式

与两个因素有关, 暗 A 是 x 状态, B-1-J 为先验条件下, X 发生的概率

Bayes Theorem

In many situations, you will know one conditional distribution $P(x|y)$ and $P(y)$ but you are really interested in the other conditional distribution $P(y|x)$.

$$P(y|x) = \frac{P(x,y)}{P(x)} = \frac{P(x|y)P(y)}{P(x)}$$

$$P(y/x) \cdot P(x) = P(x|y) \cdot P(y) \Rightarrow P(x|y) = \frac{P(y/x) \cdot P(x)}{P(y)}$$

Let A_1, A_2, \dots, A_n be a set of mutually exclusive events that together form the sample space S.Let B be any event from the same sample space, such that $P(B) > 0$. Then

$$P(A_k|B) = \frac{P(B|A_k)P(A_k)}{\sum_{i=1}^n P(B|A_i)P(A_i)} = P(B)$$

Example

For a magazine, the probability that the reader is male given that the reader is at least 35 years old is 0.3. The probability that a reader is male, given that the reader is under 35, is 0.85. If 75% of the reader are under 35, what is the probability that a randomly chosen reader is

- a) Male
- b) Female
- c) Under 35 and it is given the reader is a female

• Solution:

- (a) Let A_1 be the event of the reader being at least 35 years old, A_2 the event of the reader being under 35 years old, M be the event of the reader is being a male, and F be the event of the reader is being a female.

$P(A_2) = 0.75, P(A_1) = 1 - 0.75 = 0.25$

$P(M|A_1) = 0.3, P(M|A_2) = 0.65$

$P(F|A_1) = 0.7, P(F|A_2) = 0.35$

$P(M) = P(A_1, M) + P(A_2, M)$ 同时满足两个条件

$= P(A_1)P(M|A_1) + P(A_2)P(M|A_2) = 0.25 \times 0.3 + 0.75 \times 0.65 = 0.5625$

- (b) $P(F) = 1 - P(M) = 1 - 0.5625 = 0.4375$

- (c) $P(A_2|F) = \frac{P(F|A_2)P(A_2)}{P(F|A_1)P(A_1) + P(F|A_2)P(A_2)} = \frac{0.35 \times 0.75}{0.7 \times 0.25 + 0.35 \times 0.75} = 0.6$

Data Preprocessing

Data Cleaning; Transformation; Integration

Normalization; Missing Data Imputation;

Noise Identification; Feature Selection;

Instance Selection; Discretization;

Feature Extraction; Instance Generation

Categorical: Nominal & Ordinal

Numerical: Discrete & Continuous

Finding Redundant Attributes:

 χ^2 Correlation Test quantifies the correlation among two nominal attributes contain c and r different values each:

$$\chi^2 = \sum_{i=1}^c \sum_{j=1}^r \frac{(o_{ij} - e_{ij})^2}{e_{ij}}$$

where o_{ij} is the frequency of (A_i, B_j) and:

$$e_{ij} = \frac{\text{count}(A = a_i) \times \text{count}(B = b_j)}{m}$$

The larger the χ^2 value, the more likely the variables are related

卡方相关检验: 适用于 nominal attributes

Chi-Square Correlation Test: An Example

(C)

	Play Chess	Not Play Chess	Sum (row)
Like science fiction	250	200	450
Not like science fiction	50	1000	1050
Sum	300	1200	1500

- Consider "Play Chess" and "Like science fiction", the expected value = $\frac{300+450}{1500} = 90$

- The expected values for the above:

$$\chi^2 = \frac{\text{实际}-\text{期望/期望}}{90}^2 + \frac{(250-90)^2}{210} + \frac{(200-360)^2}{360} + \frac{(1000-840)^2}{840} = 507.93$$

- It shows that "Like science fiction" and "Play Chess" are correlated.

Finding Redundant Attributes

Pearson's product moment coefficient 适用于 numerical attributes

$$r_{A,B} = \frac{\sum_{i=1}^m (a_i - \bar{A})(b_i - \bar{B})}{\sqrt{m} \sigma_A \sigma_B}$$

covariance 适用于两个变量同时变化

$$\text{Cov}(A, B) = E((A - \bar{A})(B - \bar{B})) = \frac{1}{m} \sum_{i=1}^m (a_i - \bar{A})(b_i - \bar{B})$$

$$r_{A,B} = \frac{\text{Cov}(A, B)}{\sigma_A \sigma_B}$$

Binning: Sort the data, partition into bins

Equal-depth bins: Smooth by bin means:

- Bin 1: 4, 8, 9, 15
- Bin 2: 21, 21, 24, 25
- Bin 3: 26, 28, 29, 34

- Bin 1: 9, 9, 9, 9
- Bin 2: 23, 23, 23, 23
- Bin 3: 29, 29, 29, 29

Similarity

Manhattan distance: Euclidean distance: Cosine similarity:

Pearson correlation: Euclidean distance: Cosine similarity:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Pearson correlation: Euclidean distance: Cosine similarity:

Euclidean distance: Cosine similarity:

Cosine similarity: Pearson correlation:

Rule-based Approaches

Early Approaches

- Using the observed value from the most recently updated source
- Taking the average, maximum, or minimum for numerical values
- Majority voting

Advantage of these approaches

- The result is generally easier to debug and to understand

	S1	S2	S3		S1	S2	S3
Jagadish	UM	ATT	UM	Jagadish	UM	ATT	UM
Dewitt	MSR	MSR	UW	Dewitt	MSR	MSR	UW
Bernstein	MSR	MSR	MSR	Bernstein	MSR	MSR	MSR
Carey	UCI	ATT	BEA	Carey	UCI	ATT	BEA
Franklin	UCB	UCB	UMD	Franklin	UCB	UCB	UMD

Naïve Voting

• Supports difference of opinion, allows conflict resolution

• Works well for independent sources that have similar accuracy

• When sources have different accuracies

- Need to give more weight to votes by knowledgeable sources

• When sources copy from other sources

- Need to reduce the weight of votes by copiers

• Problem: the wisdom of minority

Iterative Approach

• Recent studies combine iteratively estimating the quality of the source with truth discovery

• The basic principle is:

- Sources that provide true information are more reliable, and

- it is more likely that the information provided by reliable sources is true

• Iteratively until converges:

- Truth computation step
- Source weight estimation step

• Truth computation:

- The truth is inferred through weighted voting.

$$\text{The votes of a value } v: \text{vote}(v) = \left(\sum_{s \in S} \frac{w_s}{|S|} \right)^g$$

• S_p - a set of sources providing the value v

• $|V_s|$ - the number of values by source s

• g - the growth of belief, e.g. 1.2

• Source weight estimation: $w_s = \left(\sum_{v \in V_s} \text{vote}(v) \right)^{\frac{1}{g}}$

Optimization-based Approaches MSLDS001 Fall 2022

• Similar to those iterative-based methods

• Models the truth discovery as an optimization problem

• Infer the reliability of the source and reliable information and to update truths and reliability weights from sources, iteratively to convergence

$$\arg \min_{\{w_s\}, \{v_o^s\}} \sum_{o \in O} w_s \cdot d(v_o^s, v_o^*)$$

Probability-based Approaches

• Uses probabilistic models to jointly calculate the reliability of the source and the correctness of the values

For example:

- The likelihood:

$$\prod_{s \in S} p(w_s | \beta) \prod_{o \in O} \left(\prod_{s \in S} p(v_o^s | \alpha) \prod_{s \in S} p(v_o^s | v_o^*, w_s) \right)$$

Big Data Privacy Issues

• User's privacy may be breached under the following circumstances:

- Personal information when combined with external datasets may lead to the inference of new facts about users.
- Personal information is collected and used to add value to business.
- Sensitive data are stored and processed in a location not secured properly and data leakage may occur during the storage and processing phases.

Four Types of Attributes in Data

• Identifier

- Can be used to uniquely identify a person

• For example: name, driving license number, mobile number

• Quasi-identifier

- Cannot uniquely identify a person by the attribute

• May be used to re-identify a person when linked with some external dataset

• For example: age, gender, postcode

• Sensitive attribute

- Attributes that a person may want to conceal

• For example: salary, disease

• Non-sensitive attribute

- Privacy of a person will not be violated when the attribute is disclosed

Privacy Protection Methods

Anonymization

• Anonymization (De-identification) of personal records

• The subject identity of the data records are removed, concealed or hidden

• Can be performed by either the users or data providers

• Users may perform anonymization before doing analytics on the big datasets

• Providers may perform anonymization before storing the datasets which then be used, sold or share

Re-identification Attacks 可能会面临重识别攻击

• Even when identifying details are removed, individuals can still be re-identified

• An intentional act to identify individuals and revealing personal details

• Correlations among datasets lead to a unique fingerprint of a single individual

Protection against Re-identification Attacks

To ensure no individual's record is unique in a given dataset

K-Anonymity

• No individual's record in the dataset released is distinguishable from at least $k-1$ other records

For accomplishing k-anonymity:

• Suppression

- All or some values are replaced by “*”

• Generalization

- Individual values of attributes are replaced by a broader category or range

• E.g.

The age of an individual can be replaced by an age group: 28 can be replaced by “<30”.

• K-anonymity may not be able to protect against target identification attacks

• L-diversity

• Enhanced version of K-anonymity

• Reduce the granularity of data representation

• For each group of records sharing a combination of quasi-identifier (key attributes), attribute values, there are at least L “well-represented” values for each confidential attribute

K-anonymity and L-diversity: An Example

4-anonymity

Zip Code	Non-Sensitive		Sensitive	
	Age	Nationality	Condition	Condition
1	130** < 30	*	Heart Disease	
2	130** < 30	*	Heart Disease	
3	130** < 30	*	Viral Infection	
4	130** < 30	*	Viral Infection	
5	148** ≥ 40	*	Cancer	
6	148** ≥ 40	*	Heart Disease	
7	148** ≥ 40	*	Viral Infection	
8	148** ≥ 40	*	Viral Infection	
9	130** 3+ ≥ 40	*	Cancer	
10	130** 3+ ≥ 40	*	Cancer	
11	130** 3+ ≥ 40	*	Cancer	
12	130** 3+ ≥ 40	*	Cancer	

3-diversity				
3 values				

Zip Code	Non-Sensitive		Sensitive	
	Age	Nationality	Condition	Condition
1	130** ≤ 40	*	Heart Disease	
2	130** ≤ 40	*	Viral Infection	
3	130** ≤ 40	*	Cancer	
4	130** ≤ 40	*	Cancer	
5	148** > 40	*	Cancer	
6	148** > 40	*	Heart Disease	
7	148** > 40	*	Viral Infection	
8	148** > 40	*	Viral Infection	
9	130** 3+ ≤ 40	*	Heart Disease	
10	130** 3+ ≤ 40	*	Viral Infection	
11	130** 3+ ≤ 40	*	Cancer	
12	130** 3+ ≤ 40	*	Cancer	

3 values

Other Privacy Protection Methods

Data aggregation

- Aggregating individual records within a report-based and summarized format before release
- Similar to K-anonymity
- Analysis at finer levels difficult
- Create problems of ecological inferences

Data suppression

- Not all data values are released
- Some are removed, withheld or disclosed
- May lead to inaccurate data mining and analysis

Data swapping

- Data values of selected records are swapped to hide the true owner of the records
- A high rate of swapping destroys relationships involving the swapped and unswapped variables

Data randomization

- Adding noise of randomly generated numerical values to data variables
- Distort the values of sensitive variables
- Difficult to deduce accurate matching
- Large variance noise distribution may introduce measurement errors and inaccurate regression coefficients

Data synthesis

- The values of sensitive variables are replaced with synthetic values generated by simulation

Comparison of Encryption Schemes^a

Encryption Scheme	Features	Limitations
Identity-based Encryption (IBE)	<ul style="list-style-type: none"> Access control is based on the identity of a user Complete access over all resources 	<ul style="list-style-type: none"> Time consuming in larger environment Granular access control is hard to implement Changing ciphertext receiver is not possible Data to be processed must be downloaded and decrypted
Attribute-based Encryption (ABE)	<ul style="list-style-type: none"> Access control is based on user's attribute Granular access control is possible More secure and flexible 	<ul style="list-style-type: none"> Computation overhead in handling different user categories Updating ciphertext receiver is not possible Data to be processed must be downloaded and decrypted
Proxy Re-Encryption	<ul style="list-style-type: none"> Can be deployed in IBE or ABE scheme settings Updating ciphertext receiver is possible 	<ul style="list-style-type: none"> Computational overhead Data to be processed must be downloaded and decrypted
Homomorphic Encryption	<ul style="list-style-type: none"> Computations are performed on the encrypted data Very secure 	Computational overhead is very high

MSLDS001 Fall 2022

29

Privacy Preserving Data Processing

Privacy Preserving Collaborative Data Mining

• Distributed collaborative approaches where organizations retain their own datasets and cooperate to learn the global data mining results

Privacy Preserving Record Matching

Recent approaches include data transformation and mapping into vector spaces, and combination of secure multiparty computation and data sanitization (e.g., differential privacy and k-anonymity)

Privacy for Data Matching

Two-Party Data Matching Protocol

- Only records in the two databases that are similar are identified

- The identities of these records are revealed to both organisations

- Neither of the two parties must be able to learn anything else about the other party's confidential data

(1) Alice → Bob (2) Bob → Alice (3) Alice → Bob (4) Bob → Alice

Three-Party Data Matching Protocol

- All databases are sent to a trusted matching unit

- No database owner is able to learn anything about any other databases that are being matched

- No external adversary is able to learn anything about the source databases even if they get access to any data exchanged

- Only selected attributes of matched records are revealed to the research team

Exact Privacy-Preserving Matching Techniques

Techniques

- One-way hash-encoding function

- Only allows exact matching

- Approximate matching is not feasible

- Pre-process and standardize inputs for matching

- Three-party protocol: matching based only on hash-codes

- Two-party protocol: does not hide information from the database owners

- Add a secret key in three-party protocol

Approximate Privacy-Preserving Matching Techniques

- A two-party protocol for securely calculating edit distance

- The distance matrix is stored in a shared fashion between the two database owners and generated iteratively

- Based on homomorphic encryption approach

Dice coefficient = $2(3)/(3+4)$

= 6/7

• Bloom filters

- Five bits set to 1 are in common in both Bloom filters
- “peter” is with seven bits set
- “pete” is with 5 bits set
- The Dice coefficient is $2(5)/(7+5) = 10/12$

• Statistical analysis

• Statistical analysis is the science of collecting, exploring and presenting large amounts of data (a.k.a. dataset) to discover underlying patterns and trends¹.

• It is used extensively in science, from physics to social sciences.

• The following are the major tasks in statistical analysis:

- Describing and summarizing the data

- Identifying the relationship between variables

- Forecasting the outcomes

• Suppose we ask twenty students their weights and record them as: 65, 122, 131, 138, 142, 148, 151, 153, 155, 156, 157, 160, 162, 174, 178, 197, 201, 235, ...

Mean of a sample: $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = 155.85$

Minimum of a sample: 65

Maximum of a sample: 235

Range of a sample: $235 - 65 = 170$

Median of a sample: $M = \frac{x_{(n/2)}}{2}$

Mode of a sample: $M = \text{the value that appears most often}$

Range of a sample: $x_{(n)} - x_{(1)} = 235 - 65 = 170$

Interquartile range (IQR) of a sample: $IQR = x_{(75)} - x_{(25)} = 178.5 - 122.5 = 56$

First quartile: $x_{(25)} = 122.5$

Third quartile: $x_{(75)} = 178.5$

First quartile: $Q_1 = \frac{x_{(25)} + x_{(75)}}{2} = \frac{122.5 + 178.5}{2} = 150.5$

Third quartile: $Q_3 = \frac{x_{(75)} + x_{(100)}}{2} = \frac{178.5 + 235}{2} = 206.75$

• Variance of a sample: $s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$

Standard deviation of a sample: $s = \sqrt{s^2} = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$

Median of a sample: Middle number of the sorted list of x_1, x_2, \dots, x_n

If n is even, the median is the simple average of the middle two numbers

Mode of a sample: The value that appears most often in $\{x_1, x_2, \dots, x_n\}$

Range of a sample: Minimum to Maximum

pth percentile of a sample: The value so that roughly $p\%$ of the sample are smaller and $(100-p)\%$ of the sample are larger

Interquartile range (IQR) of a sample: Third quartile - First quartile

• First quartile: Median of the first half of the data

• Third quartile: Median of the second half of the data

Variance of a sample: $\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$

Standard deviation of a sample: $\sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$

Median of a sample: Middle number of the sorted list of x_1, x_2, \dots, x_n

If n is even, the median is the simple average of the middle two numbers

Mode of a sample: The value that appears most often in $\{x_1, x_2, \dots, x_n\}$

Range of a sample: Minimum to Maximum

pth percentile of a sample: The value so that roughly $p\%$ of the sample are smaller and $(100-p)\%$ of the sample are larger

Interquartile range (IQR) of a sample: Third quartile - First quartile

• First quartile: Median of the first half of the data

• Third quartile