# Knowledge Discovery & Data Mining — Similarity and Distance Measures — Instructor: Yong Zhuang

Yong Zhuang

yong.zhuang@gvsu.edu

#### **Recall: data cleaning and Integration**

- Data Preprocessing: An Overview
  - Data Quality
  - Major Tasks in Data Preprocessing
- Data Cleaning
  - Missing Values,
  - Noise(Denoising): Binning, Regression, Low-pass filter 0
  - Outliers,
  - Data Cleaning as a Process Ο
- Data Integration
  - Ο
  - Handling Redundancy: At the tuple level; Between attributes Ο

Schema integration, Entity identification problem, Detect and resolve data value conflicts





#### **Recall: data transformation**

- Data Transformation
  - Transformation functions
  - Data normalization
    - Min-max
    - Z-score
    - Decimal scaling
  - Data discretization: Binning, Clustering analysis, Histogram analysis 0



#### **Recall: data reduction**

- Data compression
  - Discrete wavelet transform (DWT)
- Sampling
  - Sampling without replacement
  - Sampling with replacement
  - Cluster or Stratified Sampling



#### Outline

- Similarity and distance measures
  - Proximity Measures for
    - Nominal Attributes
    - Binary Attributes
    - Numeric Attributes
    - Ordinal attributes
    - Mixed types
  - Cosine Similarity
  - Entropy & Cross Entropy
  - Kullback-Leibler divergence



#### Similarity and distance measures

#### • Similarity

- Numerical measure of how alike two data objects are
- Value is higher when objects are more alike Ο
- Often falls in the range [0,1]
- Dissimilarity (e.g., distance)
  - Numerical measure of how different two data objects are
  - Lower when objects are more alike Ο
  - Minimum dissimilarity is often 0 Ο
  - Upper limit varies
- Proximity refers to a similarity or dissimilarity



## **Dissimilarity matrix**

#### Data matrix (or object-by-attribute stru

- This structure stores the n data point dimensions (n objects ×p attributes)
- Two-mode

#### **Dissimilarity matrix (or object-by-object structure)**:

- A triangular matrix
- d(i, j) is the measured dissimilarity or "difference" between objects i and j
- d(i, j) >= 0, close to 0 when objects i and j are highly similar or "near" each other, and becomes larger the more they differ.
- d(i, j) = d(j, i).
- **One-mode**

	ture):	
ts	with p	С







#### Simple matching

objects i and j can be computed based on the ratio of mismatches

d(i, j)

sim(i, j) = 1 - d(i, j)Similarity: 

**Encoding:** creating a new binary attribute for each of the M states of a nominal attribute.



Dissimilarity: m: # of matches, p: total # of variables, then dissimilarity between two

$$= \frac{p - m}{p}$$

$$j) = \frac{m}{p}$$



**Example.** Suppose that we have the sample data of following table, so the dissimilarity matrix is





$$d(i, j) = \frac{p - m}{p}$$
Object Identifier Test-1 (nominal  
1 code A  
2 code B  
3 code C  
4 code A







**Example.** Suppose that we have the sample data of following table, so the dissimilarity matrix is

$$\begin{bmatrix} 0 \\ d(2,1) & 0 \\ d(3,1) & d(3,2) & 0 \\ d(4,1) & d(4,2) & d(4,3) \end{bmatrix}$$

Only one nominal attribute, so p = 1



$$d(i, j) = \frac{p - m}{p}$$
Object Identifier Test-1 (nominal  
1 code A  
2 code B  
3 code C  
4 code A







**Example.** Suppose that we have the sample data of following table, so the dissimilarity matrix is



$$d(i, j) = \frac{p - m}{p}$$
Object Identifier Test-1 (nominal  
1 code A  
2 code B  
3 code C  
4 code A









If all binary attributes are thought of as having the same weight, we have the 2×2 contingency table, where q is the number of attributes that equal 1 for both objects i and j, r is the number of attributes that equal 1 for object i but equal 0 for object j, s is the number of attributes that equal 0 for object i but equal 1 for object j, and t is the number of attributes that equal 0 for both objects i and j. The total number of attributes is p, where p = q + r + s + t.



contingency table



Symmetric binary attributes: symmetric

#### **Asymmetric binary attributes:**

- the two states are not equally important,
- that of two 0s (a negative match).
- asymmetric binary dissimilarity

$$d(i, j) = \frac{r+s}{q+r+s}$$

- asymmetric binary similarity: Ο
  - is called the Jaccard coefficient

$$sim(i, j) = \frac{q}{q+r+s} = 1 - d$$

binary dissimilarity 
$$d(i, j) = \frac{r+s}{q+r+s}$$

#### the agreement of two 1s (a positive match) is then considered more significant than **Object** j







**Example.** Suppose that we have the sample data of following table, so the distance between each pair of the three patients—Jack, Mary, and Jim—is

Name	Gender	Fever	Cough	Test-1	Test-2	Test-3	Test-4
Jack	Μ	Y	N	Р	N	N	N
Mary	F	Y	N	Р	N	Р	N
Jim	Μ	Y	Р	N	N	N	N

- d(Jack, Jim) =
- d(Jack, Mary) =
- d(Jim, Mary) =



$$d(i, j) = \frac{r+q}{q+r}$$

Object j

		1	0	S (1
Object i	1	q	r	q
	0	S	t	S
	Sum(col.)	q+s	<i>r+t</i>	p

#### contingency table





**Example.** Suppose that we have the sample data of following table, so the distance between each pair of the three patients—Jack, Mary, and Jim—is

Name	Gender	Fever	Cough	Test-1	Test-2	Test-3	Test-4
Jack	Μ	Y	Ν	Ρ	N	N	N
Mary	F	Y	Ν	Р	N	Р	N
Jim	Μ	Y	Р	Ν	N	N	N

$$d(Jack, Jim) = \frac{1+1}{1+1+1} = 0.67$$
$$d(Jack, Mary) = \frac{0+1}{2+0+1} = 0.33$$
$$d(Jim, Mary) = \frac{1+2}{1+1+2} = 0.75$$

$$d(i, j) = \frac{r+q}{q+r}$$

Object j

		1	0	S (1
Object i	1	q	r	q
	0	S	t	S
	Sum(col.)	q+s	<i>r+t</i>	p

#### contingency table





Distance measures are commonly used for computing the dissimilarity of objects described by numeric attributes.

- **Euclidean distance:** The most popular distance measure
  - Ο by p numeric attributes.
  - The Euclidean distance between objects i and j is defined as

$$d(i, j) = \sqrt{(x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 + \dots + (x_{ip} - x_{jp})^2}$$

Manhattan (or city block) distance:  $d(i, j) = |x_{i1} - x_{j1}| + |x_{i2}|$ 

Let i = (xi1, xi2, ..., xip) and j = (xj1, xj2, ..., xjp) be two objects described

$$|x_{j2}| + \cdots + |x_{ip} - x_{jp}|$$





Both the Euclidean and the Manhattan distance satisfy the following mathematical properties:

- **Nonnegativity:**  $d(i, j) \ge 0$ : Distance is a nonnegative number.
- **Identity of indiscernibles:** d(i, i) = 0: The distance of an object to itself is 0.
- **Symmetry:** d(i, j) = d(j, i): Distance is a symmetric function.
- **Triangle inequality:**  $d(i, j) \le d(i, k) + d(k, j)$ : Going directly from object i to object j in space is no more than making a detour over any other object k.

A measure that satisfies these conditions is known as **metric**.



17

distances. It is defined as

$$d(i, j) = \sqrt[h]{|x_{i1} - x_{j1}|^h} + |x_{i2} - x_{j2}|^h + \dots + |x_{ip} - x_{jp}|^h$$

- where h is a real number such that  $h \ge 1$ Ο Manhattan distance when h = 1 (L1 norm) Euclidean distance when h = 2 (L2 norm)
- generalization of the Minkowski distance for  $h \rightarrow \infty$

$$d(i, j) = \lim_{h \to \infty} \left( \sum_{f=1}^{p} |x_{if} - x_{jf}|^h \right)^{\frac{1}{h}} = \max_{f}^{p} |x_{if} - x_{jf}|$$

Minkowski distance: is a generalization of the Euclidean and Manhattan

Supremum distance (Lmax, L∞ norm, and the Chebyshev distance): a



18

**Example.** Let x1 = (1, 2) and x2 = (3, 5) represent two objects as shown



**Euclidean distance:** 

 $d(i, j) = \sqrt{(x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 + \dots + (x_{ip} - x_{jp})^2}$ 

- Manhattan distance:
- $d(i, j) = |x_{i1} x_{j1}| + |x_{i2} x_{j2}| + \dots + |x_{ip} x_{jp}|$ 
  - **Supremum distance:**

$$\int_{\infty} \left( \sum_{f=1}^{p} |x_{if} - x_{jf}|^h \right)^{\frac{1}{h}} = \max_{f} |x_{if} - x_{ff}|^{\frac{1}{h}} = \max_{f} |x_{if} - x_{ff}|^{\frac{1}{h}}$$









#### **Example.** Let x1 = (1, 2) and x2 = (3, 5) represent two objects as shown



Yong Zhuang

 $x_2 = (3, 5)$ Euclidean distance  $=(2^2+3^2)^{1/2}=3.61$ Manhattan distance = 2 + 3 = 5

> Supremum distance = 5 - 2 = 3





- The values of an ordinal attribute have a meaningful order or ranking about them. e.g. drink size = {small, medium, large} Suppose that f is an ordinal attribute and has Mf ordered states. Let 1, ..., Mf represent ranking of these ordered states. The dissimilarity of f can be calculated by:
- 1.Normalize the rank *rif* of the object i and attribute f by
- 2.Compute the dissimilarity using distance methods

$$z_{if} = \frac{r_{if} - 1}{M_f - 1}$$



**Example.** Suppose that we have the sam distance, the dissimilarity matrix is?



Yong Zhuang

$$\frac{z_{if}}{M_f} = \frac{r_{if} - 1}{M_f - 1}$$

<b>Object Identifier</b>	Test-2 (ordinal
1	excellent
2	fair
3	good
4	excellent





**Example.** Suppose that we have the sam distance, the dissimilarity matrix is?

- **Z**1f **=**
- **Z**2f **=**
- **Z**3f **=**
- **Z**4f **=**





$$\frac{z_{if}}{M_f} = \frac{r_{if} - 1}{M_f - 1}$$

<b>Object Identifier</b>	Test-2 (ordinal
1	excellent
2	fair
3	good
4	excellent





**Example.** Suppose that we have the sam distance, the dissimilarity matrix is?

- $M_f = 3$ , [fair, good, excellent] = [1,2,3]
- $Z_{1f} = 1$
- $Z_{2f} = 0$
- **Z**3f **=** 0.5
- Z4f = 1



$$\frac{z_{if}}{M_f} = \frac{r_{if} - 1}{M_f - 1}$$

<b>Object Identifier</b>	Test-2 (ordinal
1	excellent
2	fair
3	good
4	excellent





**Example.** Suppose that we have the sam distance, the dissimilarity matrix is?

 $M_f = 3$ , [fair, good, excellent] = [1,2,3]

 $Z_{1f} = 1$ 

 $Z_{2f} = 0$ 

 $Z_{3f} = 0.5$ 

Z4f = 1

 $\begin{bmatrix} 0 \\ 1.0 & 0 \\ 0.5 & 0.5 \\ 0 & 1.0 \end{bmatrix}$ 



$$z_{if} = \frac{r_{if} - 1}{M_f - 1}$$

		<b>Object Identifier</b>	Test-2 (ordinal
		1	excellent
0		2	fair
		3	good
0.5	U_	4	excellent





A database may contain all attribute types

Nominal, symmetric binary, asymmetric binary, numeric, ordinal between objects i and j is defined as

where the indicator  $\delta_{ii}^{(f)} = 0$  if

- 1.  $x_{if}$  or  $x_{if}$  is missing (i.e., there is no measurement of attribute f for object i or object j ),
- 2.  $x_{if} = x_{if} = 0$  and attribute f is asymmetric binary;

3. otherwise, 
$$\delta_{ij}^{(f)} = 1$$
.

- Suppose that the data set contains p attributes of mixed types. The dissimilarity d(i, j)

$$d(i, j) = \frac{\sum_{f=1}^{p} \delta_{ij}^{(f)} d_{ij}^{(f)}}{\sum_{f=1}^{p} \delta_{ij}^{(f)}}$$





The contribution of attribute f to the dissimilarity between i and j (i.e.,  $d_{ii}$ ) is computed dependent on its type:

- If f is numeric:  $d_{ij}^{(f)} = \frac{|x_{if} x_{jf}|}{\max_f \min_f}$ , where  $\max_f$  and  $\min_f$  are the maximum and minimum values of attribute f, respectively;
- If f is nominal or binary:  $d_{ij}^{(f)} = 0$  if  $x_{if} = x_{jf}$ ; otherwise,  $d_{ij}^{(f)} = 1$ ; and

• If f is ordinal: compute the ranks  $r_{if}$  and  $z_{if} = \frac{r_{if}-1}{M_f-1}$ , and treat  $z_{if}$  as numeric.





**Example.** compute a dissimilarity matrix for the objects in following table

<b>Object Identifier</b>	Test-1 (nominal)	Test-2 (ordinal)	Test-3 (numeric)
1	code A	excellent	45
2	code B	fair	22
3	code C	good	64
4	code A	excellent	28

- If f is numeric:  $d_{ij}^{(f)} = \frac{|x_{if} x_{jf}|}{\max_f \min_f}$ , where  $\max_f$  and  $\min_f$  are the maximum and minimum values of attribute f, respectively;
- If f is nominal or binary:  $d_{ii}^{(f)} = 0$  if  $x_{if} = 0$
- If f is ordinal: compute the ranks  $r_{if}$  and z

= 
$$x_{jf}$$
; otherwise,  $d_{ij}^{(f)} = 1$ ; and  
 $z_{if} = \frac{r_{if} - 1}{M_f - 1}$ , and treat  $z_{if}$  as numeric.





#### Example. compute a dissimilarity matrix for the objects in following table

<b>Object Identifier</b>	Test-1 (nominal)	Test-2 (ordinal)	Test-3(numeric)
1	code A	excellent	45
2	code B	fair	22
3	code C	good	64
4	code A	excellent	28



Yong Zhuang



#### **Example.** compute a dissimilarity matrix for the objects in following table

	<b>Object Identifier</b>	Test-1 (nominal)	Test-2 (ordinal)	Test-3(numeric)
	1	code A	excellent	45
	2	code B	fair	22
	3	code C	good	64
	4	code A	excellent	28
$d_{ij}^{(1)} =$	$\begin{bmatrix} 0 & & & \\ 1 & 0 & & \\ 1 & 1 & 0 & \\ 0 & 1 & 1 & 0 \end{bmatrix}$	$d_{ij}^{(2)} = \begin{bmatrix} 0 \\ 1.0 & 0 \\ 0.5 & 0.5 \\ 0 & 1.0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0.5 \end{bmatrix}$	



#### Example. compute a dissimilarity matrix for the objects in following table

	<b>Object Identifier</b>	Test-1 (nominal)	Test-2 (ordinal)	Test-3(numeric)		
	1	code A	excellent	45		
	2	code B	fair	22		
	3	code C	good	64		
	4	code A	excellent	28		
$d_{ij}^{(1)} =$	$\begin{bmatrix} 0 & & & \\ 1 & 0 & & \\ 1 & 1 & 0 & \\ 0 & 1 & 1 & 0 \end{bmatrix}$	$d_{ij}^{(2)} = \begin{bmatrix} 0 \\ 1.0 & 0 \\ 0.5 & 0.5 \\ 0 & 1.0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0.5 \end{bmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$	$= \begin{bmatrix} 0 \\ 0.55 & 0 \\ 0.45 & 1.00 \\ 0.40 & 0.14 \end{bmatrix} ($		





#### **Example.** compute a dissimilarity matrix for the objects in following table

	<b>Object Identifier</b>	Test-1 (nominal)	Test-2 (ordinal)	Test-3(numeric)			
	1	code A	excellent	45 22 64			
	2	code B	fair				
	3	code C	good				
	4	code A	excellent	28			
$d_{ij}^{(1)} =$	$\begin{bmatrix} 0 & & & \\ 1 & 0 & & \\ 1 & 1 & 0 & \\ 0 & 1 & 1 & 0 \end{bmatrix}$	$d_{ij}^{(2)} = \begin{bmatrix} 0 \\ 1.0 & 0 \\ 0.5 & 0.5 \\ 0 & 1.0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0.5 \end{bmatrix} d_{ij}^{(3)}$	$= \begin{bmatrix} 0 \\ 0.55 & 0 \\ 0.45 & 1.00 \\ 0.40 & 0.14 \end{bmatrix}$			
$\delta_{ij}^{(f)} = 1,  d(3, 1) = \frac{1(1) + 1(0.50) + 1(0.45)}{3} = 0.65$							

0 0.86





#### Example. compute a dissimilarity matrix for the objects in following table

<b>Object Identifier</b>	Test-1 (nominal)	Test-2 (ordinal)	Test-3(numeric)
1	code A	excellent	45
2	code B	fair	22
3	code C	good	64
4	code A	excellent	28

$$d(i, j) = \begin{bmatrix} 0 \\ 0.85 \\ 0.65 \\ 0.13 \end{bmatrix} 0.$$



#### 0 .83 0 .71 0.79 0



two vectors are pointing in roughly the same direction.

- Often used to measure document similarity in text analysis.
- A document can be represented by thousands of attributes, each recording the frequency of a particular word (such as keywords) or phrase in the document. Thus each document is an object represented by what is called a **term-frequency vector**.

similarity function, we have

sim(x, y)

where ||x|| is the Euclidean norm of vector x ||y|| is the Euclidean norm of vector y

**Cosine similarity:** measures the similarity between two vectors of an inner product space. It is measured by the cosine of the angle between two vectors and determines whether

Let x and y be two term-frequency vectors for comparison. Using the cosine measure as a

$$= \frac{x \cdot y}{||x||||y||}$$
  
 =  $(x_1, x_2, \dots, x_p)$ , defined as  $\sqrt{x_1^2 + x_2^2 + \dots + x_p^2}$ 





table, That is, x = (5, 0, 3, 0, 2, 0, 0, 2, 0, 0) and y = (3, 0, 2, 0, 1, 1, 0, 1, 0, 1).

How similar are x and y?

#### **Document vector or term-frequency vector.**

Document	Team	Coach	Hockey	Baseball	Soccer	Penalty	Score	Win	Loss	Sea
1	5	0	3	0	2	0	0	2	0	0
2	3	0	2	0	1	1	0	1	0	1
3	0	7	0	2	1	0	0	3	0	0
4	0	1	0	0	1	2	2	0	3	0

**Example.** Suppose that x and y are the first two term-frequency vectors in the following  $sim(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x} \cdot \mathbf{y}}{||\mathbf{x}||||\mathbf{y}||}$ 





table, That is, x = (5, 0, 3, 0, 2, 0, 0, 2, 0, 0) and y = (3, 0, 2, 0, 1, 1, 0, 1, 0, 1).

How similar are x and y?

 $x \cdot y = 5 \times 3 + 0 \times 0 + 3 \times 2 + 0 \times 0 + 2 \times 1 + 0 \times 1 + 0 \times 0 + 2 \times 1$  $+0 \times 0 + 0 \times 1 = 25$  $||\mathbf{x}|| = \sqrt{5^2 + 0^2 + 3^2 + 0^2 + 2^2 + 0^2 + 0^2 + 2^2 + 0^2 + 0^2 + 0^2 = 6.48$  $||\mathbf{v}|| = \sqrt{3^2 + 0^2 + 2^2 + 0^2 + 1^2 + 1^2 + 0^2 + 1^2 + 0^2 + 1^2 = 4.12}$ 

**Example.** Suppose that x and y are the first two term-frequency vectors in the following  $sim(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x} \cdot \mathbf{y}}{||\mathbf{x}||||\mathbf{y}||}$ 





table, That is, x = (5, 0, 3, 0, 2, 0, 0, 2, 0, 0) and y = (3, 0, 2, 0, 1, 1, 0, 1, 0, 1).

How similar are x and y?

 $+0 \times 0 + 0 \times 1 = 25$  $||\mathbf{x}|| = \sqrt{5^2 + 0^2 + 3^2 + 0^2 + 2^2 + 0^2 + 0^2 + 2^2 + 0^2 + 0^2 + 0^2 = 6.48$  $||\mathbf{v}|| = \sqrt{3^2 + 0^2 + 2^2 + 0^2 + 1^2 + 1^2 + 0^2 + 1^2 + 0^2 + 1^2 = 4.12$ 

SIII

Yong Zhuang

**Example.** Suppose that x and y are the first two term-frequency vectors in the following  $sim(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x} \cdot \mathbf{y}}{||\mathbf{x}||||\mathbf{y}||}$  $x \cdot y = 5 \times 3 + 0 \times 0 + 3 \times 2 + 0 \times 0 + 2 \times 1 + 0 \times 1 + 0 \times 0 + 2 \times 1$ 

$$(x, y) = 0.94$$





# **Information Theory**

#### **Claude Shannon**



Massachusetts Institute of Technology (MS, PhD)

Claude Elwood Shannon was an American mathematician, electrical engineer, computer scientist and cryptographer known as the "father of information theory" and as the "father of the Information Age".... and was one of the **founding fathers** of artificial intelligence.

#### **<u>The Mathematical Theory of Communication (1948)</u></u>**

#### Yong Zhuang





# **Information Theory & Entropy**

**Goal:** is to reliable and efficiently transmit a message from a sender to a recipient. In digital age, message are composed of bits. Bit = 0 or 1, when we communicate a message, we want as much useful information as possible to get through.

#### What is Entropy?

- Entropy measures the **uncertainty** in a probability distribution.

The entropy of a random variable X with a probability mass function p(x) is defined by

$$H(X) = -\sum_{x} p(x) \log_2 p(x).$$
 (1.1)

We use logarithms to base 2. The entropy will then be measured in bits. The entropy is a measure of the average uncertainty in the random variable. It is the number of bits on average required to describe the random variable.

• It quantifies the average amount of information you gain from observing an outcome.



#### Entropy

# likely. Each outcome has a probability of 1/8. so the entropy is \_\_\_\_\_?





**Example:** Consider rolling a **fair eight-sided die** where each face (1-8) is equally

$$H(X) = -\sum_{x} p(x) \log_2 p(x)$$





#### Entropy

# likely. Each outcome has a probability of 1/8. so the entropy is ?

H(X) = -8 >

H(X) =

H(X) =

This means, on average, you get 3 bits of information per roll.

**Example:** Consider rolling a **fair eight-sided die** where each face (1-8) is equally

$$\times \left(\frac{1}{8}\log_2\left(\frac{1}{8}\right)\right) \qquad H(X) = -\sum_x p(x)\log_2 p(x) + \log_2 p(x$$

H(X) = 3 bits





#### **Example:** Now consider a **biased eight-sided die** where:

- P(1)=P(2)=0.35
- P(3)=P(4)=0.1
- P(5)=P(6)=0.04
- P(7)=P(8)=0.01

so the entropy is \_\_\_\_ ?







$$H(X) = -\sum_{x} p(x) \log_2 p(x)$$





#### **Example:** Now consider a **biased eight-sided die** where:

- P(1)=P(2)=0.35
- P(3)=P(4)=0.1
- P(5)=P(6)=0.04
- P(7)=P(8)=0.01

so the entropy is ?

 $H(X) = -\left(0.35 \log_2(0.35) + 0.35 \log_2(0.35) + 0.1 \log_2(0.1) + 0.1 \log_2(0.1) + 0.04 \log_2(0.04) + 0.04 \log_2(0.04) + 0.01 \log_2(0.01) + 0.01 \log_2(0.01)\right)$ 

- $0.35 \log_2(0.35) \approx -0.530$
- $0.1 \log_2(0.1) \approx -0.332$ H(X)
- $0.04 \log_2(0.04) \approx -0.185$
- $0.01 \log_2(0.01) \approx -0.066$

This means, on average, you get 2.226 bits of information per roll.

Yong Zhuang



$$H(X) = -\sum_{x} p(x) \log_2 p(x)$$

$$= -(2 imes - 0.530 + 2 imes - 0.332 + 2 imes - 0.185 + 2 imes - 0.066)$$

H(X) = 1.06 + 0.664 + 0.37 + 0.132 = 2.226 bits









This means, on average, you get 2.226 bits of information per roll.

Yong Zhuang





#### Example: Now consider a biased eight-sided die where:

- $P(1)=P(2)=\frac{0.35}{0.01}$
- P(3)=P(4)=0.1
- P(5)=P(6)=0.04
- P(7)=P(8)=0.01-0.35
- H(x) = 2.226

Yong Zhuang







# **Cross-Entropy and Message Encoding**

#### What is Cross-Entropy?

from a distribution p (true) when using a code based on distribution q(predicted).

$$\begin{split} q &= \left\{\frac{1}{2^2} = 0.25, \frac{1}{2^2} = 0.25, \frac{1}{2^3} = 0.125, \frac{1}{2^3} = 0.125, \frac{1}{2^4} = 0.0625, \frac{1}{2^4} = 0.0625, \frac{1}{2^5} = 0.03125, \frac{1}{2^5}$$

Yong Zhuang

• Cross-Entropy measures the average number of bits required to transmit outcomes

$$H(p,q) = -\sum_i p(x_i) \log q(x_i)$$

Where:

- $p(x_i)$  is the probability of the true distribution for event  $x_i$ ,
- $q(x_i)$  is the probability of the predicted distribution for event  $x_i$ ,





# **Cross-Entropy and Message Encoding**

#### What is Cross-Entropy?

from a distribution p (true) when using a code based on distribution q(predicted).

$$H(p,q) = -\sum_i p(x_i) \log i$$

- It also measures the difference between the true probability distribution p and the true one.
- If p = q, the cross-entropy will be equal to the entropy of p. If p and q differ, the KL Divergence.

• Cross-Entropy measures the average number of bits required to transmit outcomes

 $\log q(x_i) = H(p) + D_{KL}(p||q)$ 

predicted distribution q, quantifying how well the predicted distribution approximates the

cross-entropy will be greater than the entropy of p, the difference between p and q is the



47

# Kulback-Leibler divergence

used in the data mining literature to measure the difference between two probability distributions over the same variable x.

- closely related to relative entropy, information divergence, and information for discrimination
- is a **nonsymmetric** measure of the difference between two probability distributions p (x) and q(x)
- the KL divergence of q(x) from p(x), denoted DKL(p(x)||q(x)), is a measure of the information loss when q(x) is used to approximate p(x).



Kullback-Leibler divergence(the KL divergence): a measure that has been popularly



# Kullback-Leibler divergence

Let p(x) and q(x) be two probability distributions of a discrete random variable x. That is, both p(x) and q(x) sum up to 1, and p(x) > 0 and q(x) > 0 for any x in X. DKL(p(x)||q(x)) is defined as

$$D_{KL}(p(x)||q(x)) = \sum_{x \in X} p(x) \ln \frac{p(x)}{q(x)}$$

Typically p(x) represents the "true" distribution of data. The measure q(x) typically represents a theory, model, description, or approximation of p(x).

- it is not a distance measure, because it is not a metric measure.
- (x) to p(x).
- DKL $(p(\mathbf{x})||q(\mathbf{x}))$  is a nonnegative measure. DKL $(p(\mathbf{x})||q(\mathbf{x})) \ge 0$ , • DKL $(p(\mathbf{x})||q(\mathbf{x})) = 0$  if and only if  $p(\mathbf{x}) = q(\mathbf{x})$

It is not symmetric: the KL from p(x) to q(x) is generally not the same as the KL from q





## Kullback-Leibler divergence

1/5, c : 1/5) and Q: (a : 5/9, b : 3/9, d : 1/9). Compute the KL divergence DKL(P ||Q)







# **Example.** Suppose there are two sample distributions P and Q as follows: P : (a : 3/5, b :



## Kulback-Leibler divergence

**Example.** Suppose there are two sample distributions P and Q as follows: P : (a : 3/5, b : 1/5, c : 1/5) and Q: (a : 5/9, b : 3/9, d : 1/9). Compute the KL divergence DKL(P ||Q)



No sample d in P, and no sample c in Q?

#### **Avoiding the Zero-Probability Problem**







# Kullback-Leibler divergence: smoothing

1/5, c : 1/5) and Q: (a : 5/9, b : 3/9, d : 1/9). Compute the KL divergence DKL(P ||Q)

- Introduce a small constant e = 0.001,
- **smoothing:** the missing symbols can be added to each distribution accordingly, with the small probability e.
  - P': (a : 3/5 e/3, b : 1/5 e/3, c : 1/5 e/3, d : e)
  - Q': (a: 5/9 e/3, b: 3/9 e/3, c: e, d: 1/9 e/3)

**Example.** Suppose there are two sample distributions P and Q as follows: P : (a : 3/5, b :

DKL(P',Q') can be calculated.



#### Summary

- Similarity and distance measures
  - Proximity Measures for
    - Nominal Attributes
    - Binary Attributes
    - Numeric Attributes
    - Ordinal attributes
    - Mixed types
  - Cosine Similarity
  - Entropy & Cross Entropy
  - Kullback-Leibler divergence

