

人工智能

智能之根：数据

陈川

中山大学 计算机学院

2024年



中山大學
SUN YAT-SEN UNIVERSITY

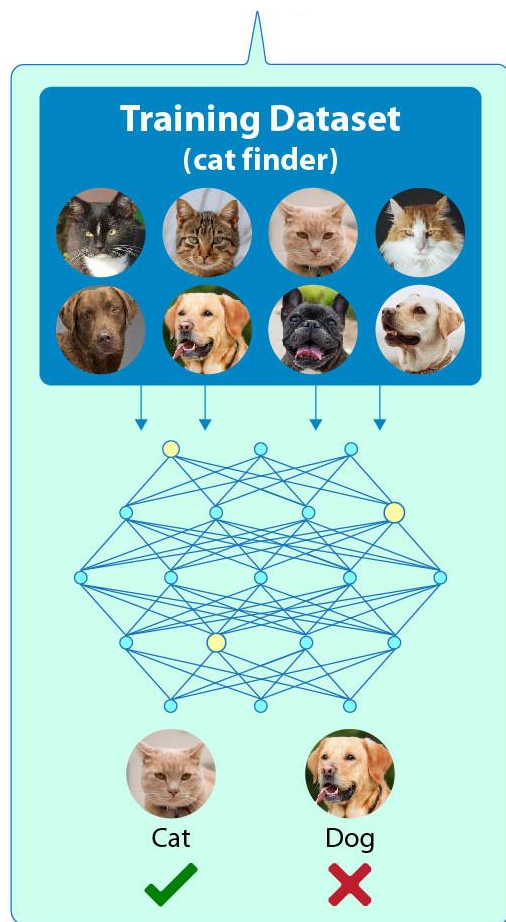
从数据、知识到智能



数据

知识

智能

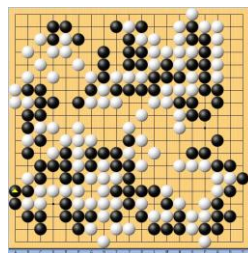


图像识别

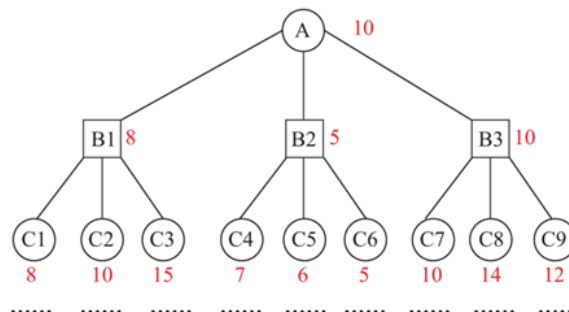
图片

神经网络
Neural Network

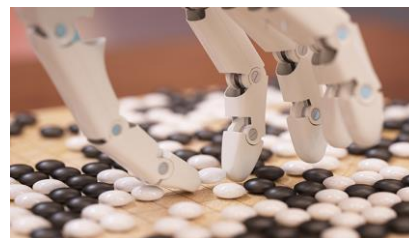
识别



棋局



博弈树
Game Tree



棋布策略

人机博弈

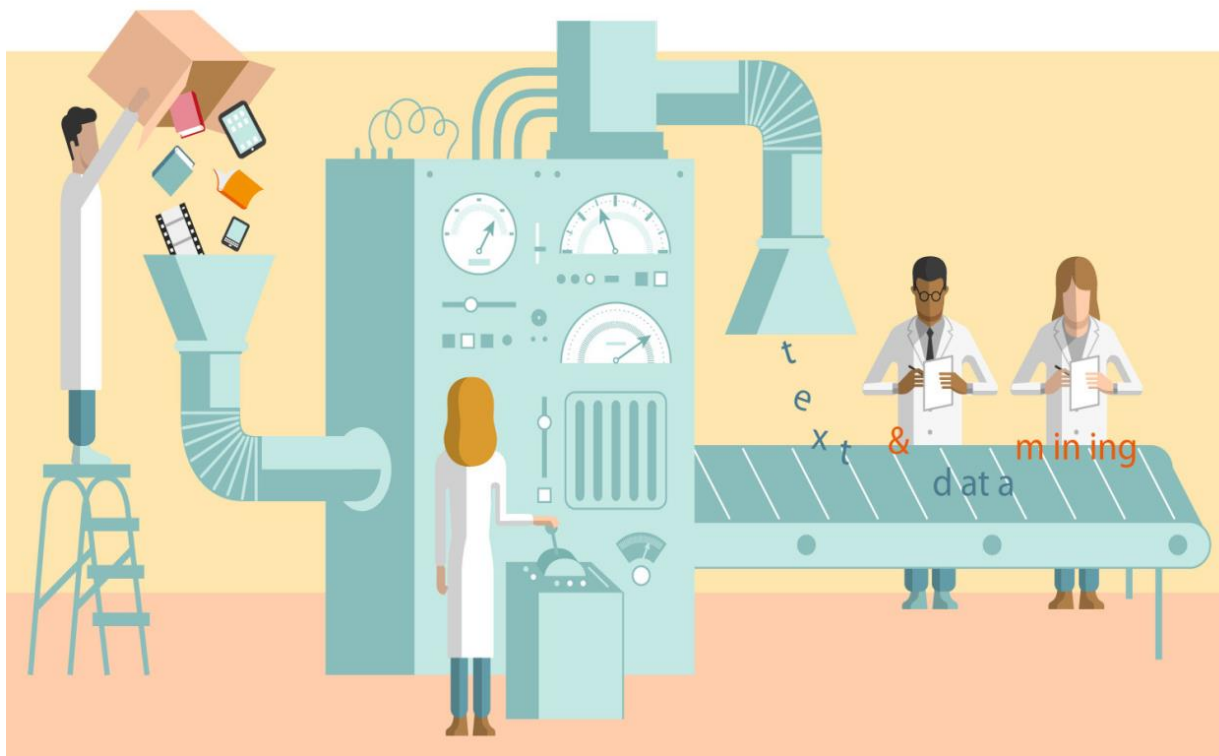
人工智能的本质到底是什么？

人工智能的本质 从数据中挖掘知识实现机器智能化 的过程

从数据、知识到智能



数据 → 知识 → 智能



过去有些教材里，这个过程也叫作“样例学习”：从样例（数据）的学习中习得样例（数据）中蕴含的知识、技能的过程，是智能体的一种基本学习方式

数据认识: Example



You receive an email from a medical researcher concerning a project that you are eager to work on.

Hi,

I've attached the data file.

Each line contains the information for a single patient and consists of five fields.

We want to predict the last field using the other fields.

Thanks and see you in a couple of days.

数据认识: Example



The first few rows of the file are as follows:

012	232	33.5	0	10.7
020	121	16.9	2	210.1
027	165	24.0	0	427.6
⋮				

Nothing looks strange. You put your doubts aside and start the analysis.

Two days later you arrive for the meeting, and before the meeting, you strike up a conversation with a statistician who is working on the project.

数据认识: Example



Statistician: So, you got the data for all the patients?

You: Yes. I haven't had much time for analysis, but I do have a few interesting results.

Statistician: Amazing. There were so many data issues with this set of patients that I couldn't do much.

You : Oh? I didn't hear about any possible problems.

Statistician: But surely you heard about what happened to field 4? It's supposed to be measured on a scale from 1 to 10, with 0 indicating a missing value, but because of a data entry error, all 10's were changed into 0's.

You : Interesting. Were there any other problems?

Statistician: Yes, fields 2 and 3 are basically the same, but I assume that you probably noticed that.

You : Yes, but these fields were only weak predictors of field 5.

012	232	33.5	0	10.7
020	121	16.9	2	210.1
027	165	24.0	0	427.6
:				

数据认识: Example



Statistician: Anyway, given all those problems, I'm surprised you were able to accomplish anything.

You : True, but my results are really quite good. Field 1 is a very strong predictor of field 5. I'm surprised that this wasn't noticed before.

Statistician: What? Field 1 is just an identification number.

You : Nonetheless, my results speak for themselves.

Statistician: Oh, no! I just remembered. We assigned ID numbers after we sorted the records based on field 5. There is a strong connection, but it's meaningless. Sorry.

Lesson: Get to know your data!

012	232	33.5	0	10.7
020	121	16.9	2	210.1
027	165	24.0	0	427.6

- 数据的概念
- 数据的类型
- 数据的质量
- 数据预处理
- 数据特征工程

What is Data?



Collection of data objects and their attributes

An attribute is a property or characteristic of an object

- Examples: eye color of a person, temperature, etc.
- **Attribute** 属性 is also known as **variable** 变量, **field** 域, **characteristic** 特性, or **feature** 特征

A collection of attributes describe an object

- **Object** is also known as **record** 记录, **point** 数据点, **case** 案例, **sample**, **entity**, or **instance**

Attributes
属性

objects
样本

<i>Tid</i>	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes



Attribute values are numbers or symbols assigned to an attribute

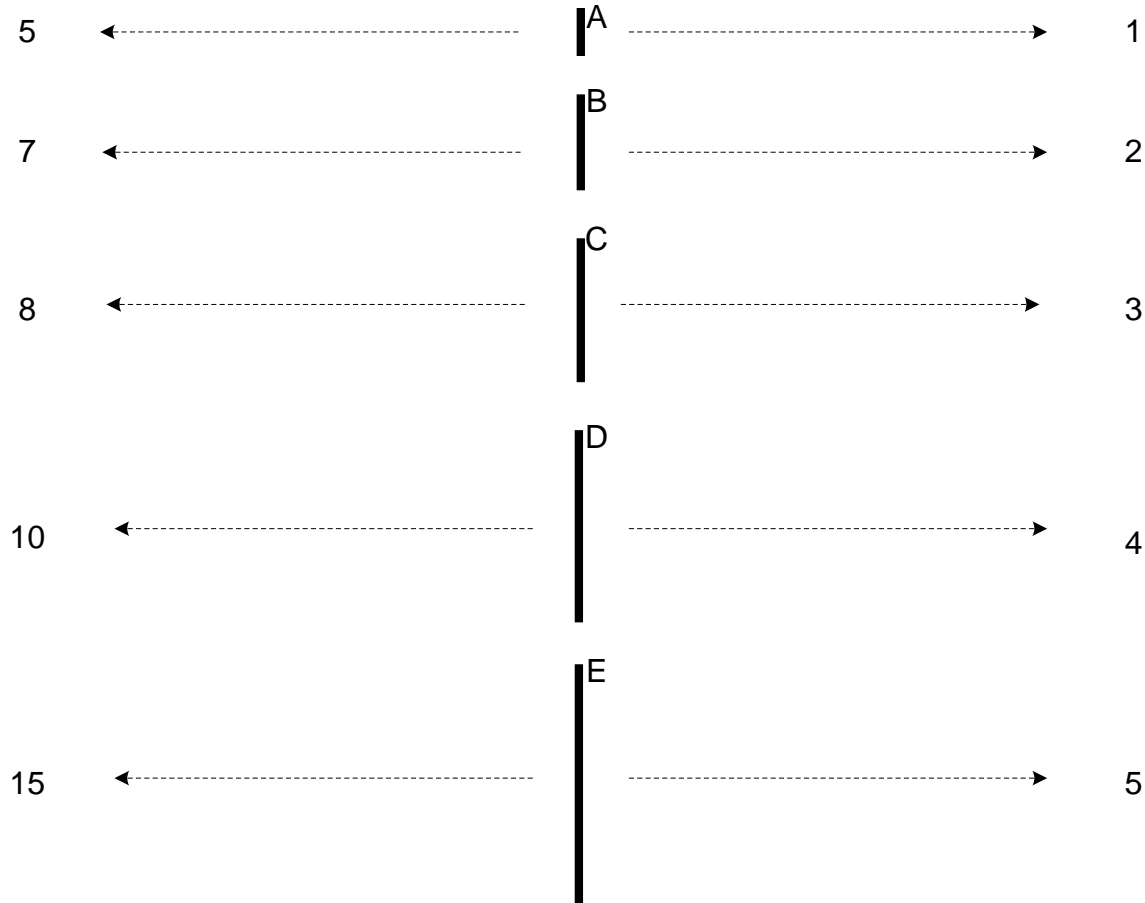
Distinction between attributes and attribute values

- Same attribute can be mapped to different attribute values
 - ◆ Example: height can be measured in feet or meters
- Different attributes can be mapped to the same set of values
 - ◆ Example: Attribute values for ID and age are integers
 - ◆ But properties of attribute values can be different
 - ID has no limit but age has a maximum and minimum value

Measurement of Length



注意：The way you measure an attribute is somewhat may not match the attributes properties. 属性值不一定反映出属性的性质



Types of Attributes 数据的属性分类



There are different types of attributes

- **Nominal** 类别属性

Examples: ID numbers, eye color, zip codes

- **Ordinal** 顺序属性

Examples: rankings (e.g., taste of potato chips on a scale from 1-10), grades, height in {tall, medium, short}

- **Interval** 区间属性

Examples: calendar dates, temperatures in Celsius or Fahrenheit.

- **Ratio** 比率属性

Examples: temperature in Kelvin, length, time, counts

Properties of Attribute Values



The type of an attribute depends on which of the following properties it possesses:

- Distinctness 区分性: $= \neq$
- Order 顺序性: $< >$
- Addition 可加性: $+ -$
- Multiplication 可比性: $* /$
- Nominal attribute 类别属性: distinctness
- Ordinal attribute 顺序属性: distinctness & order
- Interval attribute 区间属性: distinctness, order & addition
- Ratio attribute 比率属性: all 4 properties

Example of Attribute Values



Attribute Type	Description	Examples	Operations
Nominal	The values of a nominal attribute are just different names, i.e., <u>nominal attributes provide only enough information to distinguish one object from another.</u> ($=$, \neq)	zip codes, employee ID numbers, eye color, sex: { <i>male</i> , <i>female</i> }	mode, entropy, contingency correlation, χ^2 test
Ordinal	The values of an ordinal attribute provide enough information to order objects. ($<$, $>$)	hardness of minerals, { <i>good</i> , <i>better</i> , <i>best</i> }, grades, street numbers	median, percentiles, rank correlation, run tests, sign tests
Interval	For interval attributes, the differences between values are meaningful, i.e., a unit of measurement exists. ($+$, $-$)	calendar dates, temperature in Celsius or Fahrenheit	mean, standard deviation, Pearson's correlation, t and F tests
Ratio	For ratio variables, both differences and ratios are meaningful. ($*$, $/$)	temperature in Kelvin, monetary quantities, counts, age, mass, length, electrical current	geometric mean, harmonic mean, percent variation

Transformation of Attribute Values



Attribute Level	Transformation	Comments
Nominal	Any one-to-one mapping	If all employee ID numbers were reassigned, would it make any difference?
Ordinal	An order preserving change of values, i.e., $new_value = f(old_value)$ where f is a monotonic 单调 function.	An attribute encompassing the notion of good, better best can be represented equally well by the values {1, 2, 3} or by {0.5, 1, 10}.
Interval	$new_value = a * old_value + b$ where a and b are constants 华氏度=摄氏度 \times 1.8 + 32	Thus, the Fahrenheit and Celsius temperature scales differ in terms of where their zero value is and the size of a unit (degree).
Ratio	$new_value = a * old_value$	Length can be measured in meters or feet.

Discrete and Continuous Attributes



Discrete 离散 Attribute

- Has only a finite or countably infinite set of values
- Examples: zip codes, counts, or the set of words in a collection of documents
- Often represented as integer variables.
- Note: binary attributes are a special case of discrete attributes

Continuous 连续 Attribute

- Has real numbers as attribute values
- Examples: temperature, height, or weight.
- Practically, real values can only be measured and represented using a finite number of digits.
- Continuous attributes are typically represented as floating-point variables.

- 数据的概念
- **数据的类型**
- 数据的质量
- 数据预处理
- 数据特征工程

Types of data sets 数据的类型



Record

- Data Matrix
- Document Data
- Transaction Data

Graph

- World Wide Web
- Molecular Structures

Ordered

- Spatial Data
- Temporal Data
- Sequential Data
- Genetic Sequence Data

Record Data 记录型数据



Data that consists of a collection of records, each of which consists of a fixed set of attributes

<i>Tid</i>	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

记录型数据的形式: **Data Matrix**



If data objects have the same fixed set of numeric attributes, then the data objects can be thought of as points in a multi-dimensional space, where each dimension represents a distinct attribute

Such data set can be represented by an m by n matrix, where there are m rows, one for each object, and n columns, one for each attribute

Projection of x Load	Projection of y load	Distance	Load	Thickness
10.23	5.27	15.22	2.7	1.2
12.65	6.25	16.22	2.2	1.1

记录型数据: **Document Data**

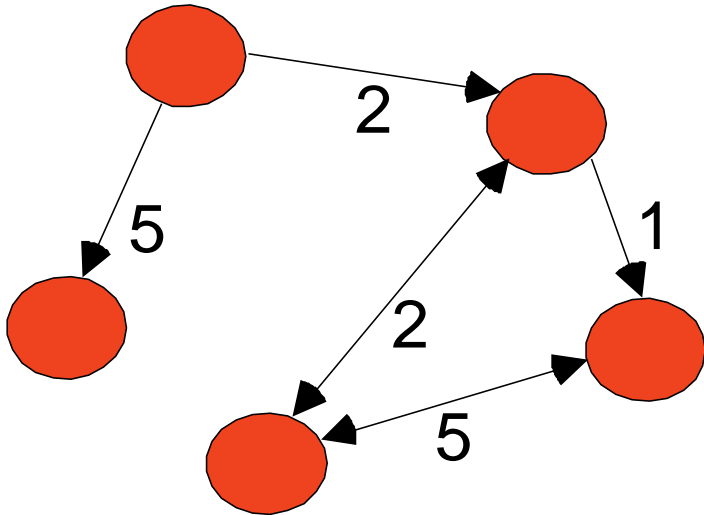


Each document becomes a **'term' vector**,

- each term is a component (attribute) of the vector,
- the value of each component is the number of times the corresponding term occurs in the document.

	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

Examples: Generic graph and HTML Links

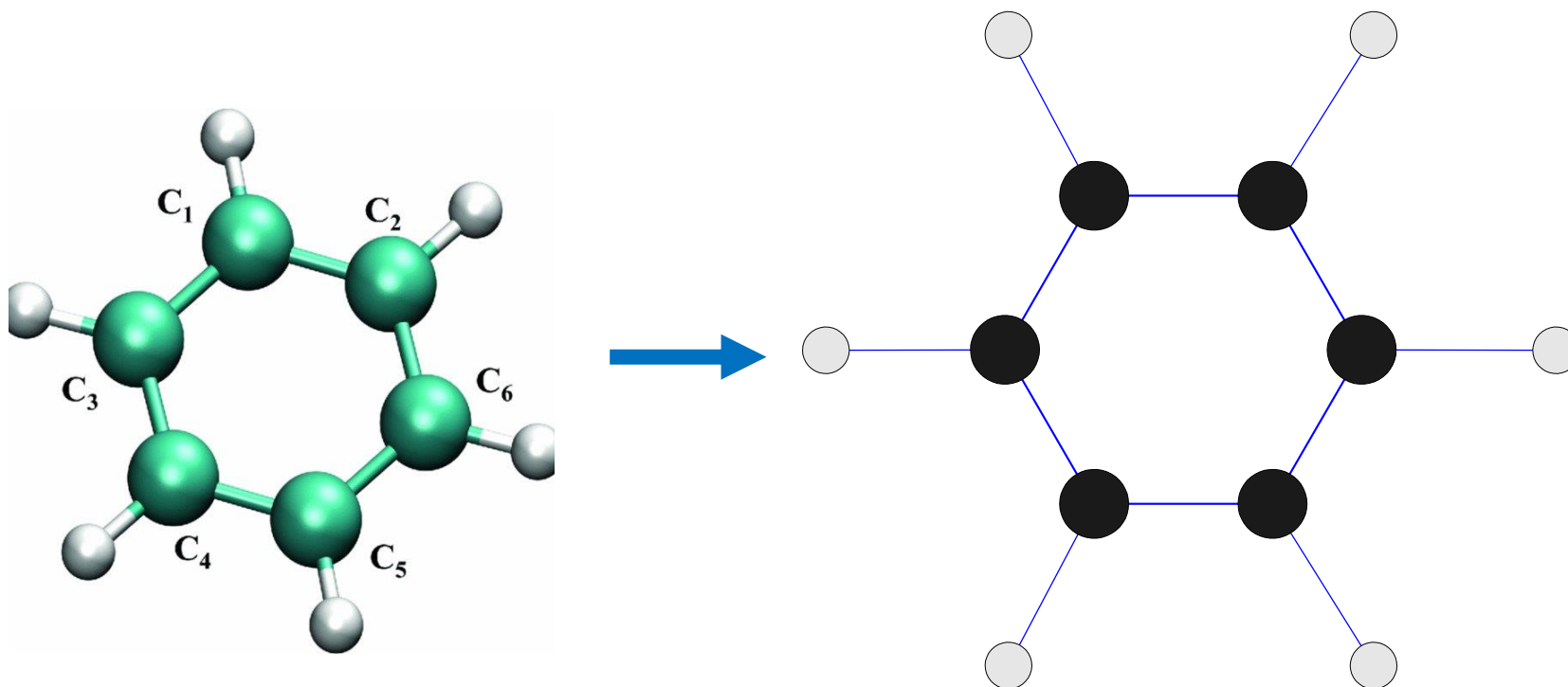


```
<a href="papers/papers.html#bbbb">  
Data Mining </a>  
<li>  
<a href="papers/papers.html#aaaa">  
Graph Partitioning </a>  
<li>  
<a href="papers/papers.html#aaaa">  
Parallel Solution of Sparse Linear System of Equations </a>  
<li>  
<a href="papers/papers.html#ffff">  
N-Body Computation and Dense Linear System Solvers
```

图数据: **Chemical Data**



Benzene Molecule (苯分子) : C_6H_6



图数据：社交网络



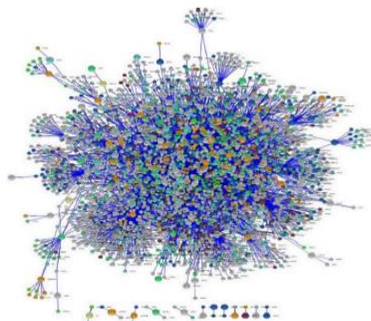
社交网络



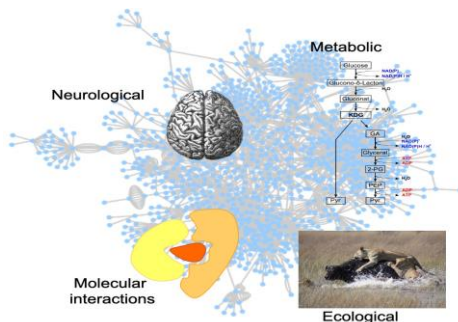
社会行为分析



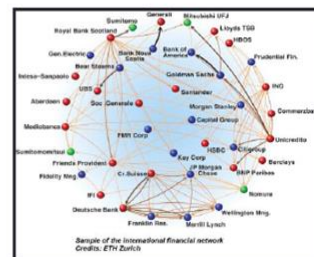
生物网络



多组学分析



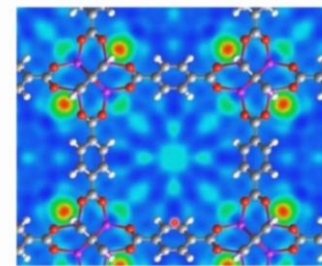
金融网络



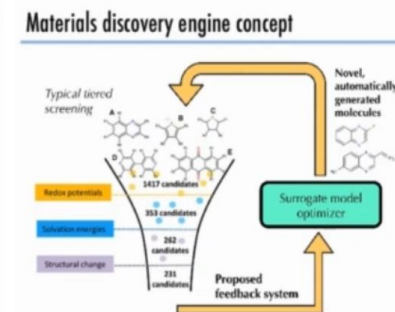
金融风险评估



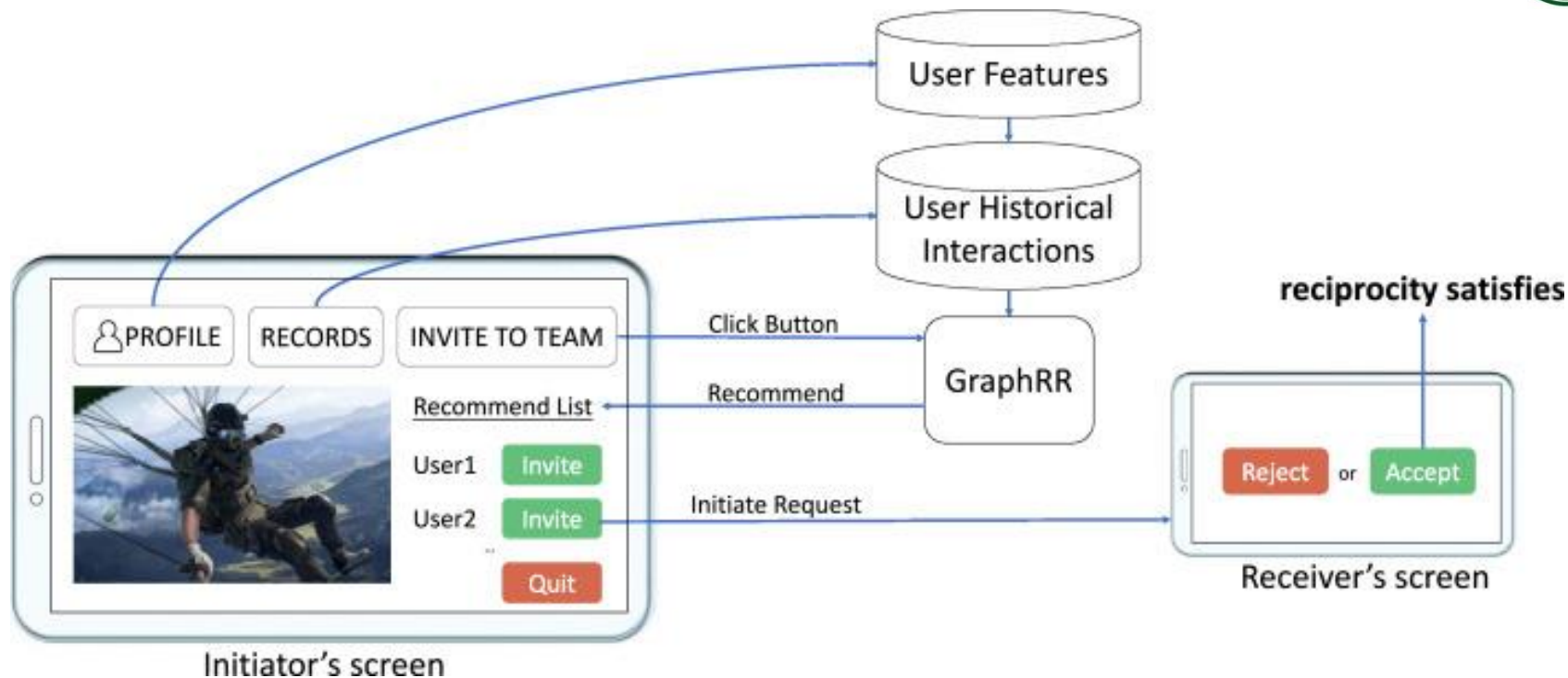
分子网络



新型药物发现



图数据：社交网络



GraphRR: A multiplex Graph based Reciprocal friend Recommender system with applications on online gaming service

Ordered Data: 有序数据



Sequential Data

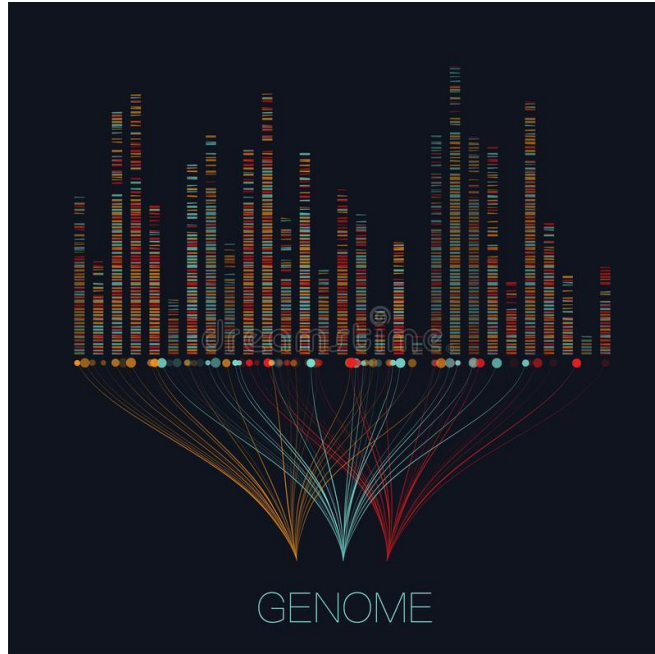
Time	Customer	Items Purchased
t1	C1	A, B
t2	C3	A, C
t2	C1	C, D
t3	C2	A, D
t4	C2	E
t5	C1	A, E

Customer	Time and Items Purchased
C1	(t1: A,B) (t2:C,D) (t5:A,E)
C2	(t3: A, D) (t4: E)
C3	(t2: A, C)

有序数据： Genomic Sequence Data



Genomic sequence data



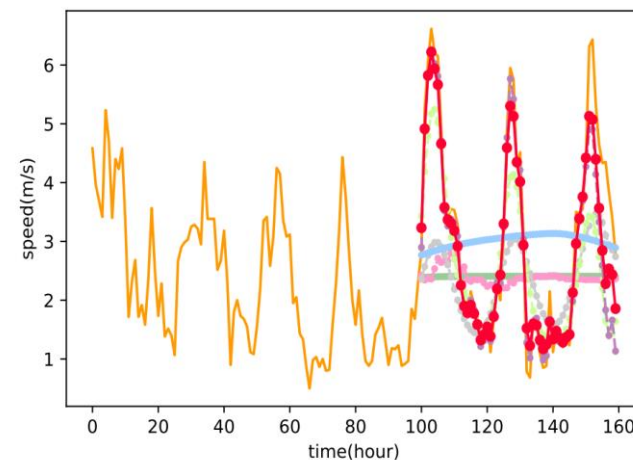
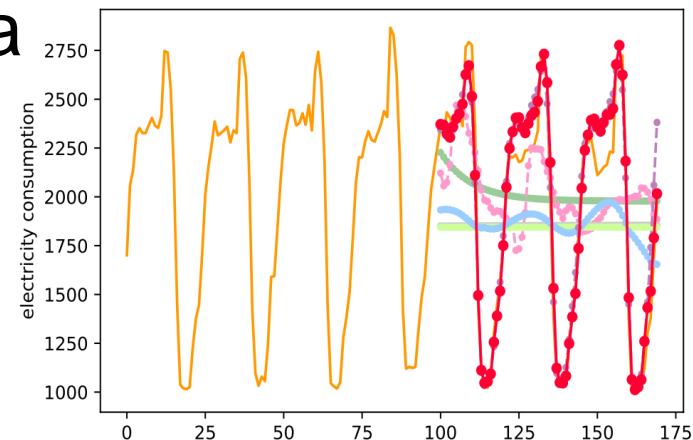
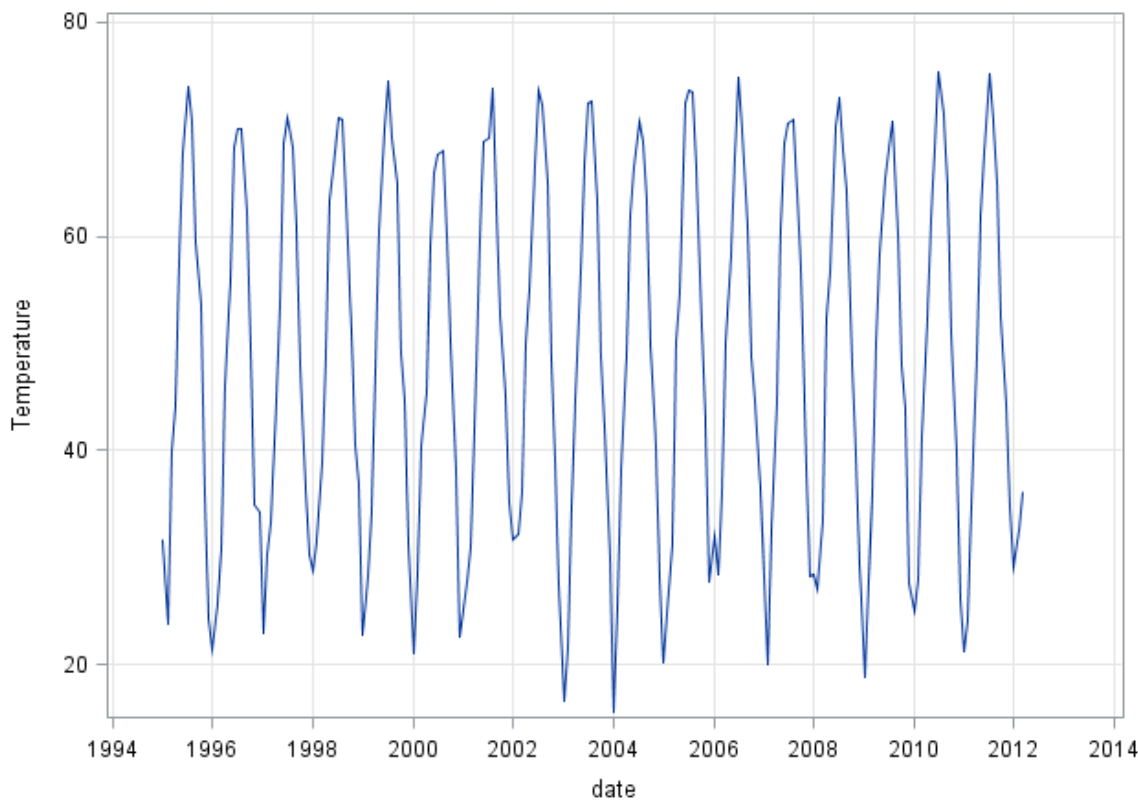
```
ATCTCTTGGCTCCAGCATCGATGAAGAACGCA
TCATTTAGAGGAAGTAAAAGTCGTAACAAGGT
GAACTGTCAAAACTTTTAACAACGGATCTCTT
TGTTGCTTCGGCGGGCGCCCGCAAGGGTGCCCG
GGCCTGCCGTGGCAGATCCCCAACGCCGGGCC
TCTCTTGGCTCCAGCATCGATGAAGAACGCAG
CAGCATCGATGAAGAACGCAGCGAAACGCGAT
CGATACTTCTGAGTGTTCTTAGCGAACTGTCA
CGGATCTCTTGGCTCCAGCATCGATGAAGAAC
ACAACGGATCTCTTGGCTCCAGCATCGATGAA
CGGATCTCTTGGCTCCAGCATCGATGAAGAAC
GATGAAGAACGCAGCGAAACGCGATATGTAAT
```


有序数据: Time Series Data



Special type of sequential data
Temporal autocorrelation

Series Values for Temperature

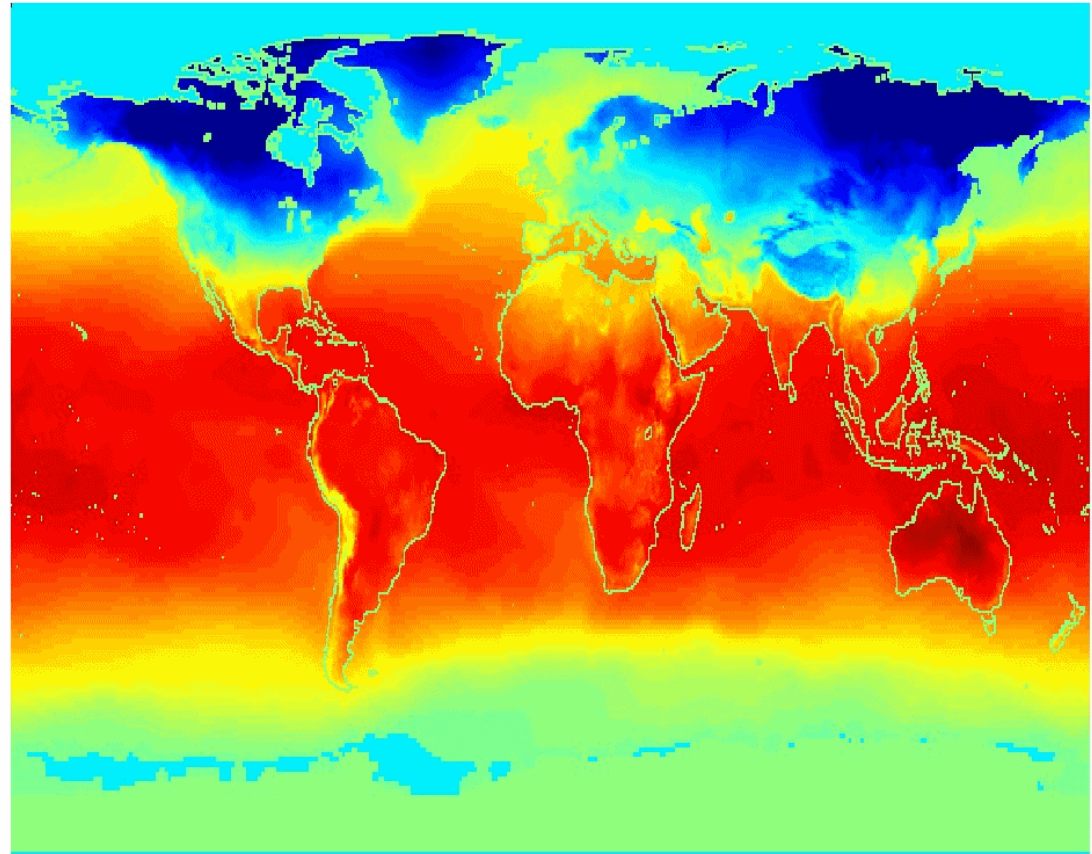


有序数据: **Spatio-Temporal Data**



Jan

**Average Monthly
Temperature of
land and ocean**



有序数据: Spatio-Temporal Data



Examples of sequence data

Speech recognition



“The quick brown fox jumped
over the lazy dog.”

Music generation



Sentiment classification

“There is nothing to like
in this movie.”



DNA sequence analysis

AGCCCCTGTGAGGAACTAG



AG**CCCCTGTGAGGAACT**AG

Machine translation

Voulez-vous chanter avec
moi?



Do you want to sing with
me?

Video activity recognition



Running

Name entity recognition

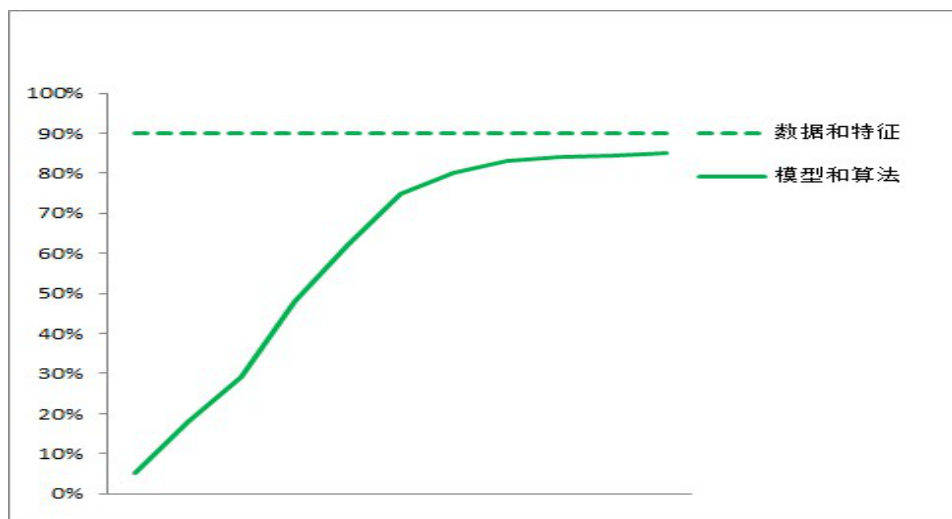
Yesterday, Harry Potter
met Hermione Granger.



Yesterday, **Harry Potter**
met **Hermione Granger**.

- 数据的概念
- 数据的类型
- **数据的质量**
- 数据预处理
- 数据特征工程

Why data quality matters?

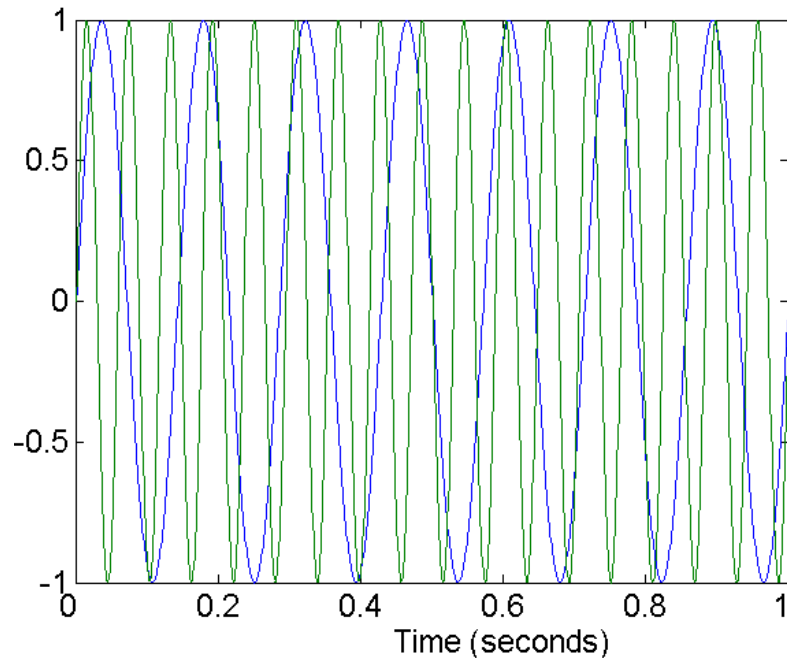


Examples of data quality problems:

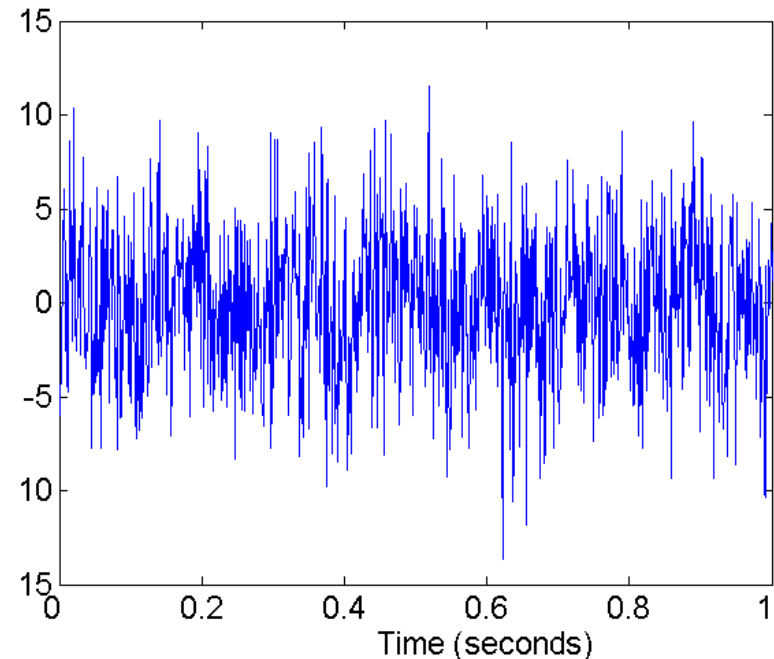
- Noise and outliers
- Missing values
- Duplicate data

Noise refers to modification of original values

- Examples: distortion of a person's voice when talking on a poor phone and “snow” on television screen



Two Sine Waves

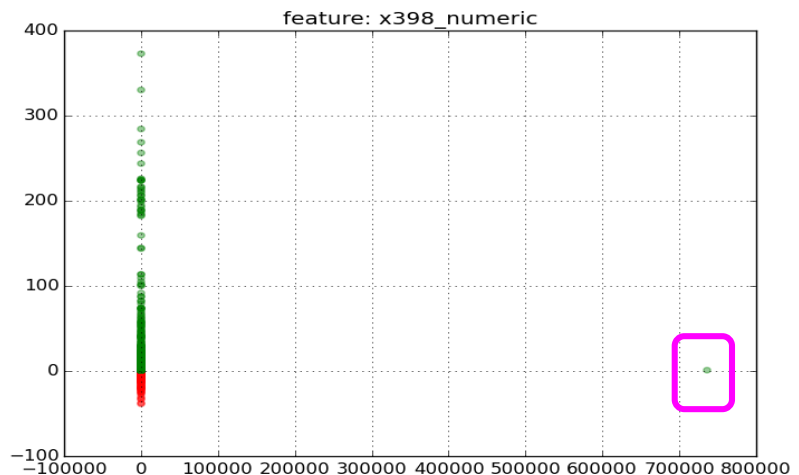
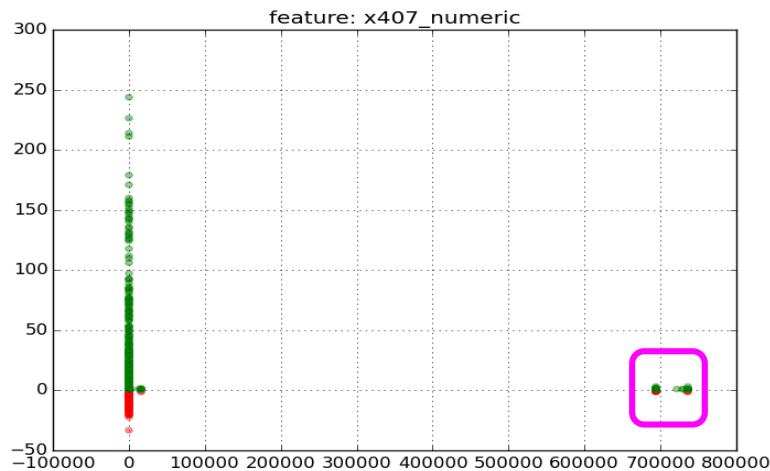
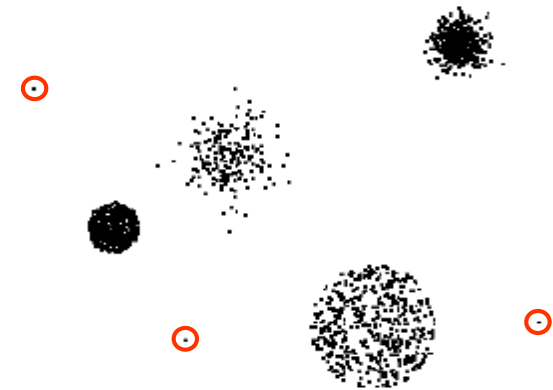


Two Sine Waves + Noise

Outliers



Outliers are data objects with characteristics that are considerably different than most of the other data objects in the data set



Missing Values



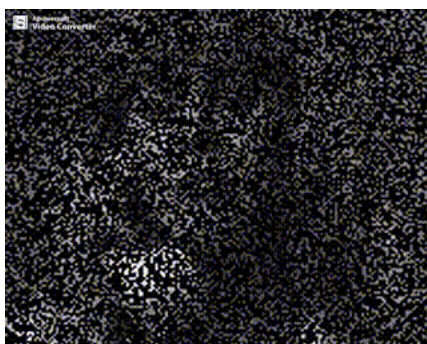
Reasons for missing values

- Information is not collected
(e.g., people decline to give their age and weight)
- Attributes may not be applicable to all cases
(e.g., annual income is not applicable to children)

图像/视频复原

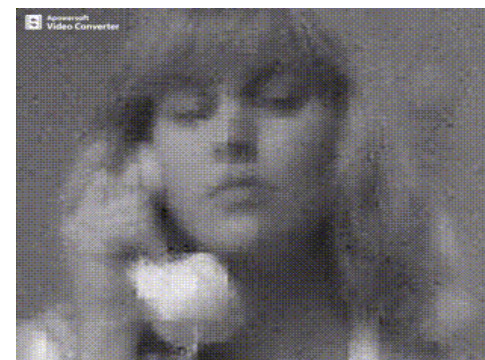


Observed (PSNR 2.03 dB)



3-channel matrix

Blue					
Green					
	255	134	93	22	
Red					
	255	134	202	22	
255	231	42	22	4	30
123	94	83	2	92	124
34	44	187	92	14	142
34	76	232	124	14	
67	83	194	202		



Missing Values



Observed (PSNR 2.03 dB)

3-channel matrix

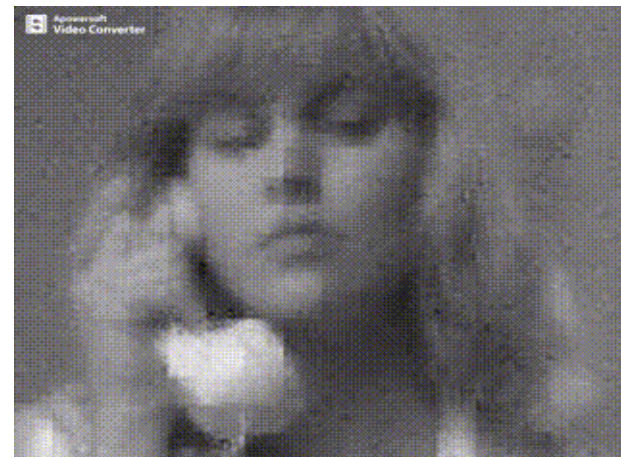
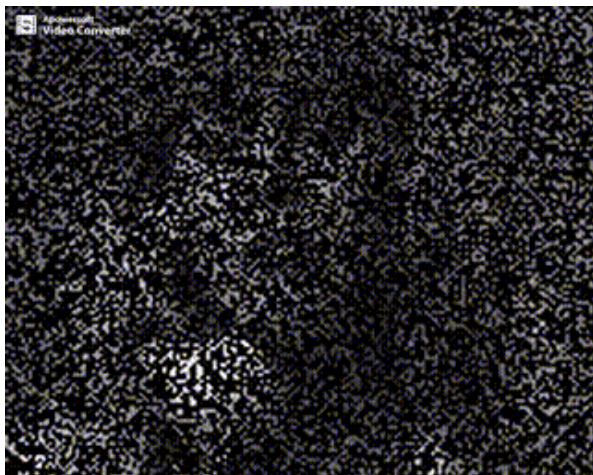
Blue					
Green					
	255	134	93	22	
Red					
	255	134	202	22	
255	231	42	22	4	30
123	94	83	2	92	124
34	44	187	92	14	142
34	76	232	124	14	
67	83	194	202		

AI模型

图像/视频复原



图像/视频数据的矩阵存储



Duplicate Data



Data set may include data objects that are duplicates, or almost duplicates of one another

- Major issue when merging data from heterogeneous sources

Examples:

- Same person with multiple email addresses

- 数据的概念
- 数据的类型
- 数据的质量
- **数据预处理**
- 数据特征工程



Aggregation

Sampling

Dimensionality Reduction

Discretization

Attribute Transformation

Feature creation

Feature subset selection

Combining two or more attributes (or objects) into a single attribute (or object)

Purpose

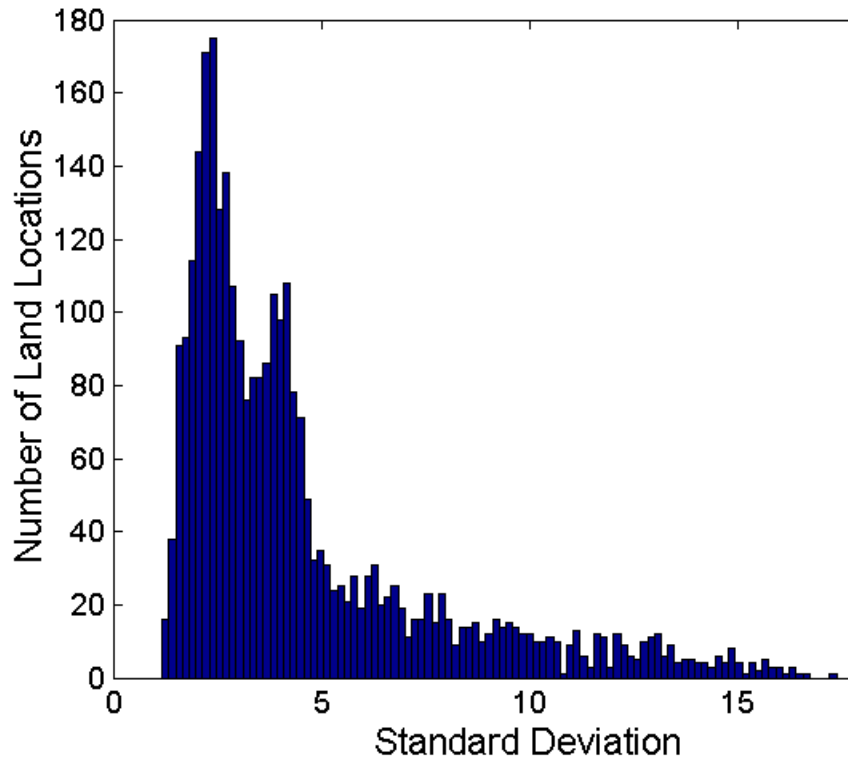
- Data reduction
 - ◆ Reduce the number of attributes or objects
- Change of scale
 - ◆ Cities aggregated into regions, states, countries, etc
- More “stable” data
 - ◆ Aggregated data tends to have less variability

Aggregation

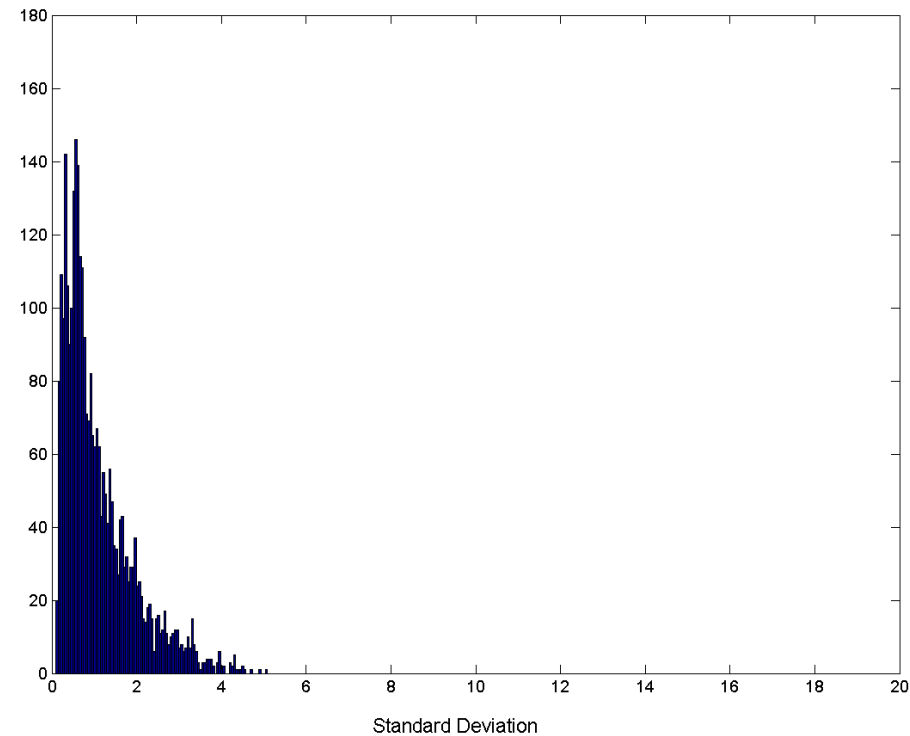


Variation of Precipitation in Australia

1982-1993的降水量，国土按经纬度分成3030个网格



Standard Deviation of Average
Monthly Precipitation



Standard Deviation of Average
Yearly Precipitation



Aggregation

Sampling

Dimensionality Reduction

Discretization

Attribute Transformation

Feature creation

Feature subset selection

Sampling is the main technique employed for data selection.

- It is often used for both the preliminary investigation of the data and the final data analysis. 数据初步调研与最终分析

Statisticians sample because **obtaining** the entire set of data of interest is too expensive or time consuming.

Sampling is used in data mining because **processing** the entire set of data of interest is too expensive or time consuming.

Sampling ...



The key principle for effective sampling is the following:

- Using a sample will work almost as well as using the entire data sets, if the sample is representative
- A sample is representative if it has approximately the same property (of interest) as the original set of data



Types of Sampling



Simple Random Sampling

- There is an equal probability of selecting any particular item

Sampling without replacement (无放回抽样)

- As each item is selected, it is removed from the population

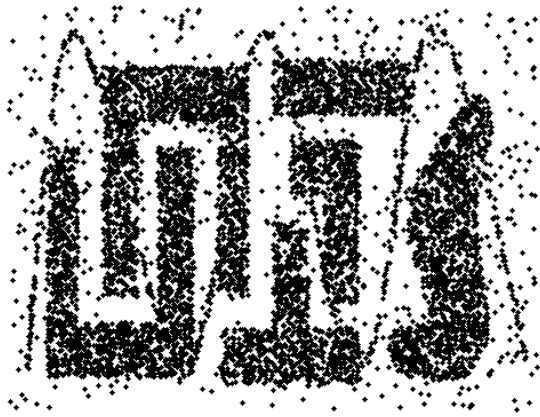
Sampling with replacement (有放回抽样)

- Objects are not removed from the population as they are selected for the sample.
 - ◆ In sampling with replacement, the same object can be picked up more than once (每个对象被选中的概率保持不变)

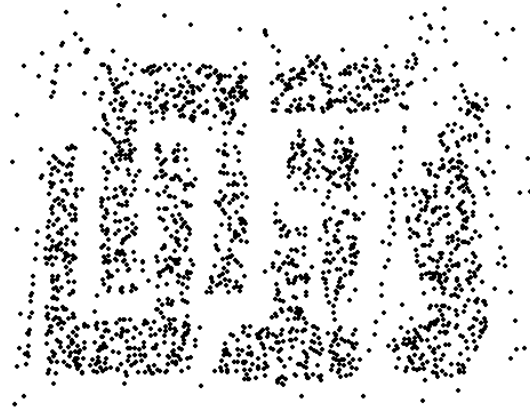
Stratified sampling (分层抽样)

- Split the data into several partitions; then draw random samples from each partition

Sample Size



8000 points



2000 Points



500 Points

Progressive Sampling



Start with a small sample

Increase the sample size

Need to evaluate the sample to judge if it is large enough

Marginal effect (边际效应)



Aggregation

Sampling

Dimensionality Reduction

Discretization

Attribute Transformation

Feature creation

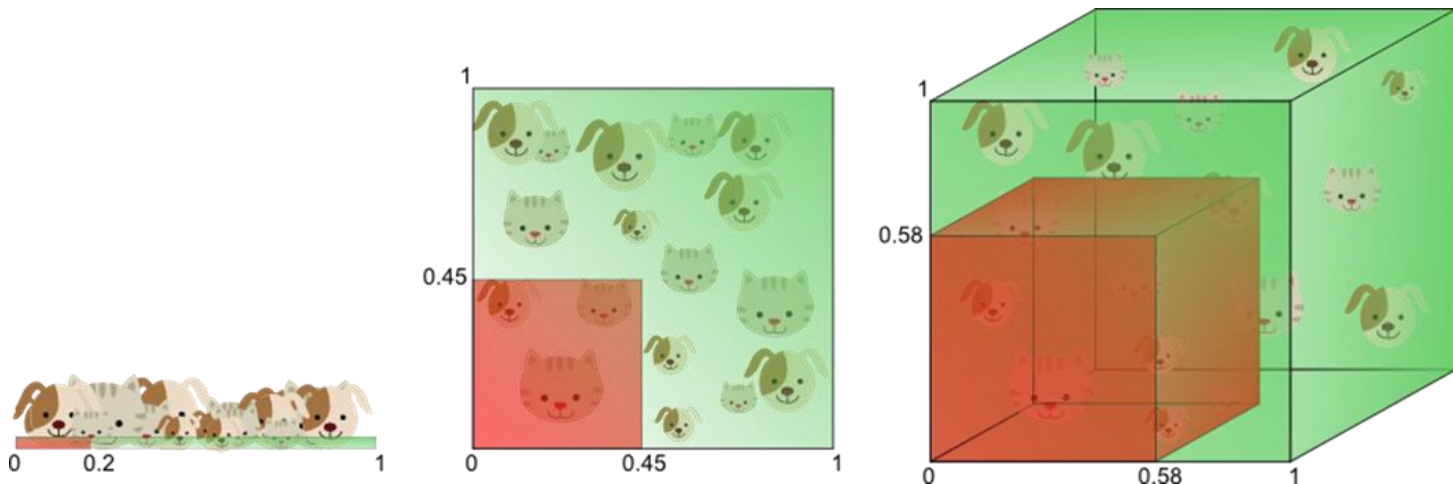
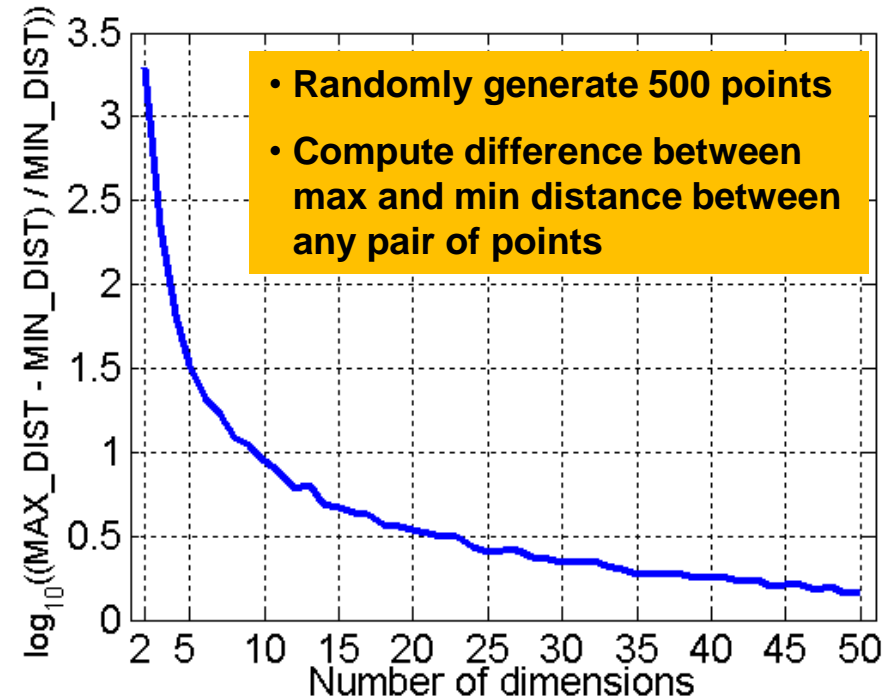
Feature subset selection

Curse of Dimensionality



When dimensionality increases, data becomes increasingly sparse in the space that it occupies

Definitions of density and distance between points, which is critical for clustering and outlier detection, become less meaningful



Dimensionality Reduction

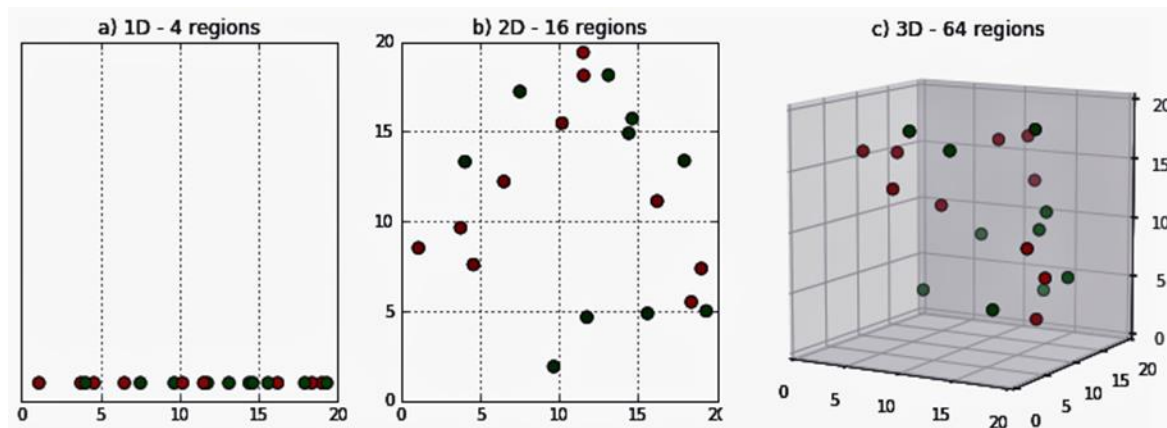


Purpose:

- Avoid curse of dimensionality
- Reduce amount of time and memory required by data mining algorithms
- Allow data to be more easily visualized
- May help to eliminate irrelevant features or reduce noise

Techniques

- Principle Component Analysis (PCA)
- Singular Value Decomposition (SVD)
- Others: supervised and non-linear techniques



Dimensionality Reduction: PCA

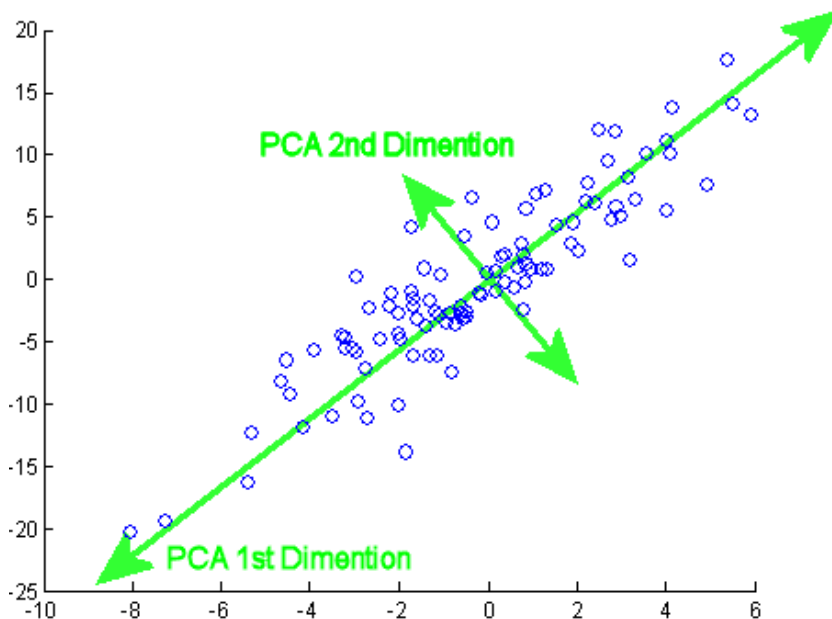


变量之间是有一定的相关关系的

当两个变量之间有一定相关关系时，可以解释为这两个变量的信息有一定的重叠

主成分分析是对于原先提出的所有变量，将重复的变量（关系紧密的变量）删去多余，建立尽可能少的新变量，使得这些新变量是两两不相关的

这些新变量在反映信息方面尽可能保持原有的信息





Aggregation

Sampling

Dimensionality Reduction

Discretization

Attribute Transformation

Feature creation

Feature subset selection

Discretization



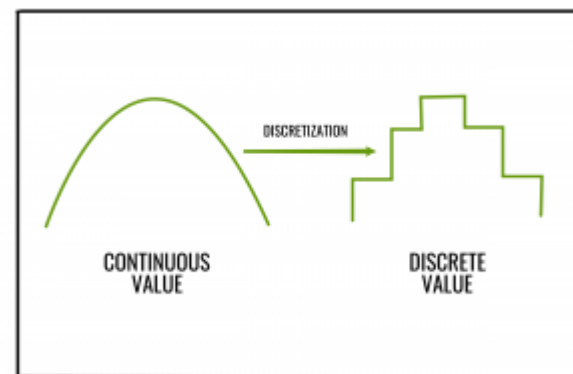
Discretization: Transform a continuous attribute to categorical attribute 连续数据的离散化

The best discretization depends on the algorithm being used

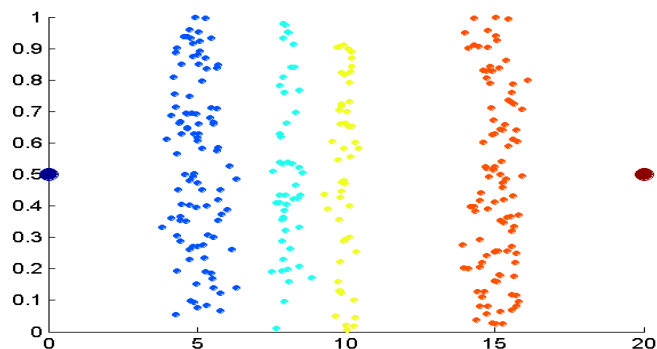
- How many categories?
- How to map the continuous attributes to these categories?
- How many split points to choose and where to place them?

Solutions

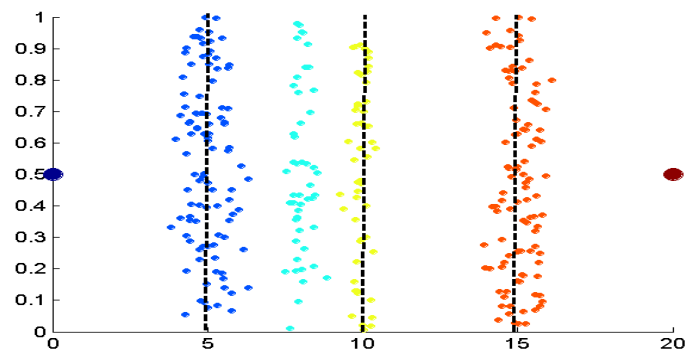
- Unsupervised discretization
- Supervised discretization



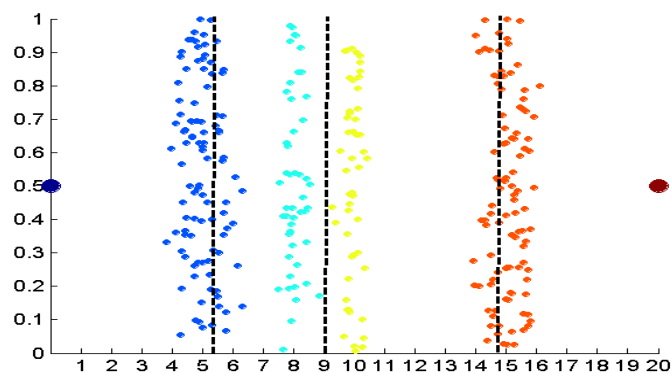
Discretization Without Labels



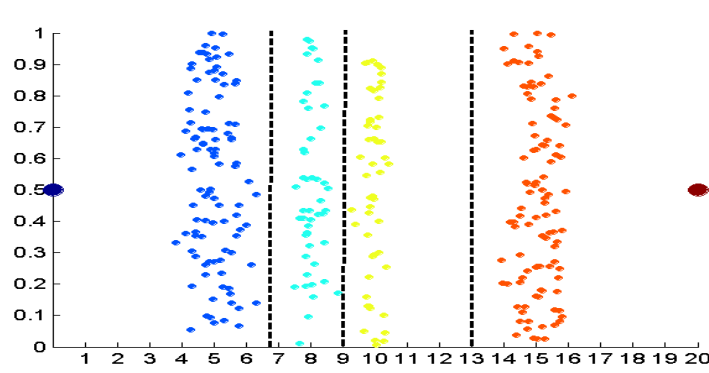
Data



Equal interval width



Equal frequency



K-means

Discretization with Labels: Entropy

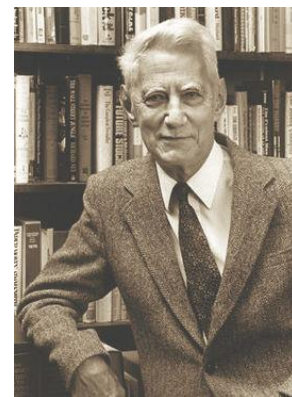


Entropy (熵)

- 熵的概念是由德国物理学家克劳修斯于1865年所提出。熵最初是被用在热力学方面的
- 香农1948年的一篇论文 《A Mathematical Theory of Communication》 提出了**信息熵**的概念，解决了对信息的量化度量问题，并且以后信息论也被作为一门单独的学科

要搞清楚一件非常不确定的事，就需要了解大量的信息。相反，如果我们对某件事已经有了较多的了解，我们不需要太多的信息就能把它搞清楚。

对于任意一个随机变量 X ，熵定义如下：“变量的不确定性越大，熵也就越大，把它搞清楚所需要的信息量也就越大。”



Entropy (熵)



- 世界杯谁是冠军?
- 世界杯赛后问一个知道结果的观众“哪支球队是冠军”? 他不愿意直接告诉我, 而要我猜, 并且我每猜一次, 他要收一元钱才肯告诉我是否猜对了, 那么我需要付给他多少钱才能知道谁是冠军呢?
- 我可以把球队编上号, 从1到32, 然后提问: “冠军的球队在1-16号中吗?” 假如他告诉我猜对了, 我会接着问: “冠军在1-8号中吗?” 假如他告诉我猜错了, 我自然知道冠军队在9-16中。这样最多只需要五次, 我就能知道哪支球队是冠军
- 谁是世界杯冠军这条消息的信息量值五块钱

热情好客的中大学子
正准备迎接远道而来的法国朋友



Entropy (熵)



不需要猜五次就能猜出谁是冠军，巴西、法国、阿根廷这样的球队得冠军的可能性比美国、越南等队大的多。

第一次猜测时不需要把 32 个球队等分成两个组，而可以把少数几个最可能的球队分成一组，把其它队分成另一组。然后我们猜冠军球队是否在那几只热门队中。

重复这样的过程，根据夺冠概率对剩下的候选球队分组，直到找到冠军队。也许三次或四次就猜出结果。

当每个球队夺冠的可能性（概率）不等时，“谁世界杯冠军”的信息量比五比特少。香农指出，它的准确信息量应该是

— “谁是世界杯冠军”的信息量：

$$= - (p_1 \log p_1 + p_2 \log p_2 + \dots + p_{32} \log p_{32}),$$

— p_1, \dots, p_{32} 是 32 个球队各自夺冠的概率

课外阅读：《数学之美》第六章“信息的度量与作用”

Supervised Discretization



基于熵的离散化方法

— 最大化区间的纯度

$$e_i = - \sum_{j=1}^k p_{ij} \log_2 p_{ij}$$

类的个数

第 i 个区间中类 j 的概率 (值的比例).

类别1 = X,

$p_{11}=3/4,$ **X X**
 $p_{21}=1/4,$ **X 0**

类别2 = 0

0 0 $p_{12}=1/4$
0 X $p_{22}=3/4$

熵：区间纯度的度量

- 只包含一个类：熵为0
- 包含所有类，并且每类出现的概率相等：熵最大

划分连续属性的简单方法：

- 将初始值切分成两部分，让结果区间产生最小的熵
- 然后选取熵最大的区间，重复该过程
- 直到区间数量达到用户指定个数



Aggregation

Sampling

Dimensionality Reduction

Discretization

Attribute Transformation

Feature creation

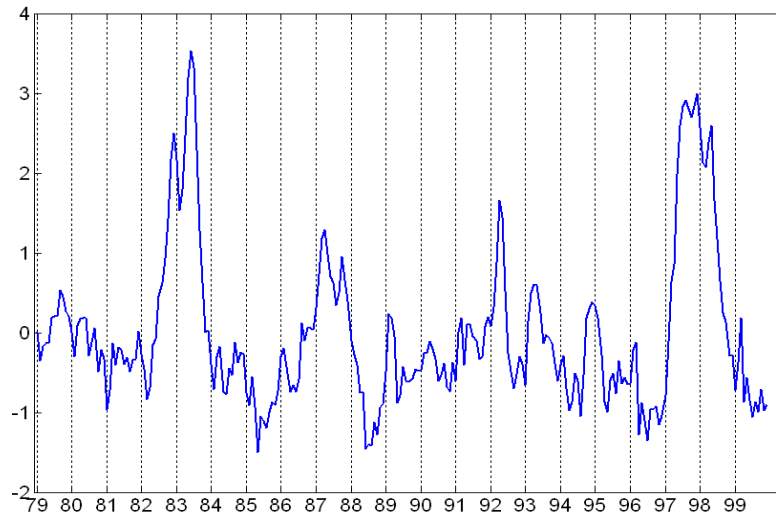
Feature subset selection

Attribute Transformation



A function that maps the entire set of values of a given attribute to a new set of replacement values such that each old value can be identified with one of the new values

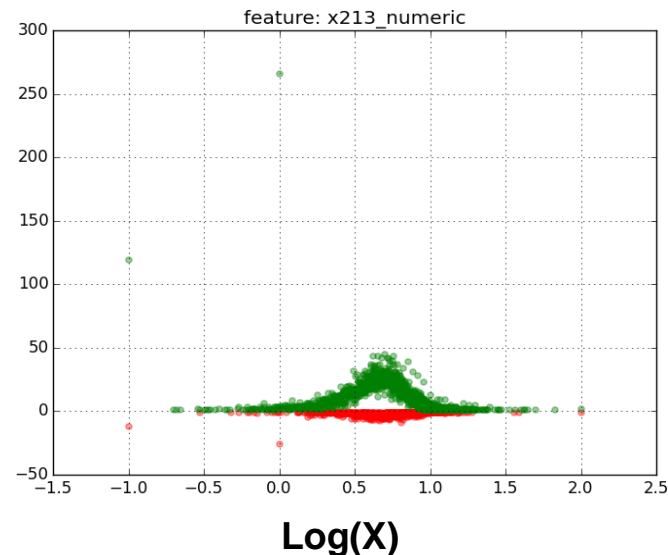
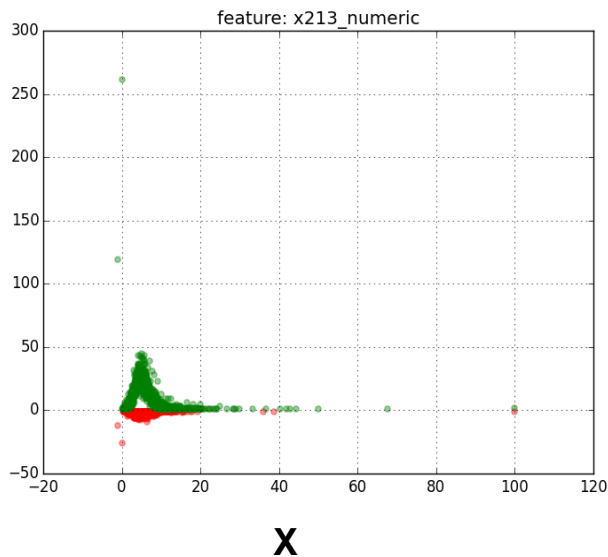
- Simple functions: x^k , $\log(x)$, e^x , $|x|$
- Standardization and Normalization



Attribute Transformation



Standardization	$(x - \text{mean}) / \text{sd}$	
Max-Min	$(x - \text{min}) / (\text{max} - \text{min})$	0 - 1
Sigmoid	$1 / (1 + \exp(-x))$	0 - 1
Tanh	$(\exp(x) - \exp(-x)) / (\exp(x) + \exp(-x))$	-1 - 1
Log	$\log(x)$	
Nesting	$\log(\text{Normalization/Max-Min/Sigmoid/Tanh})$	





Aggregation

Sampling

Dimensionality Reduction

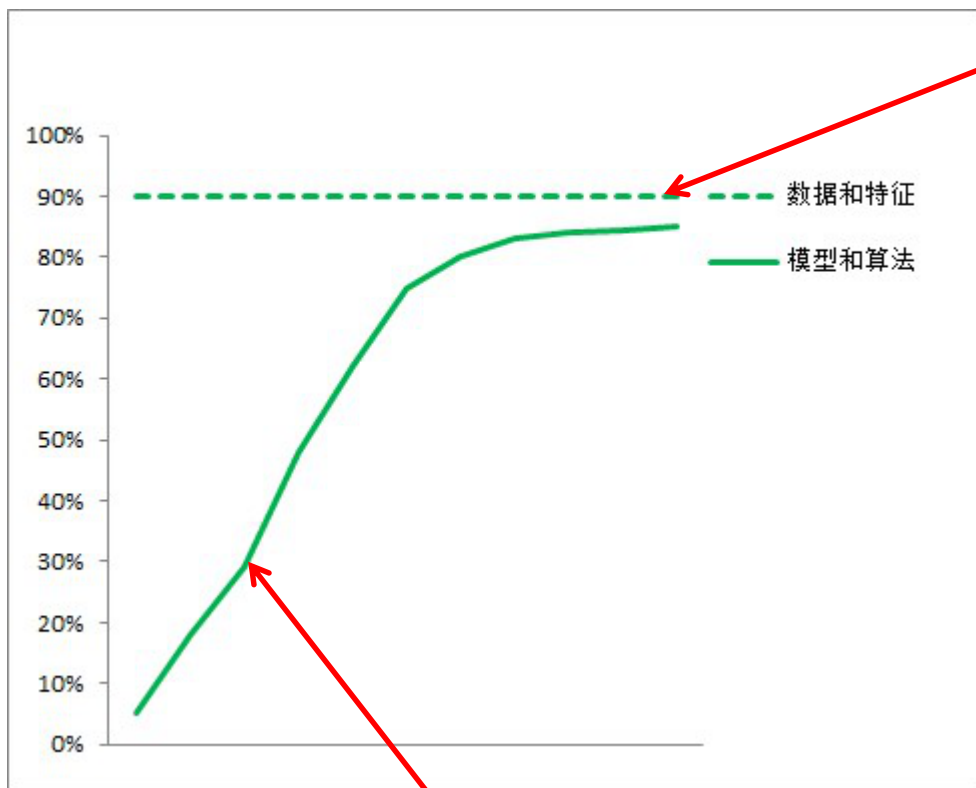
Discretization

Attribute Transformation

Feature creation

Feature subset selection

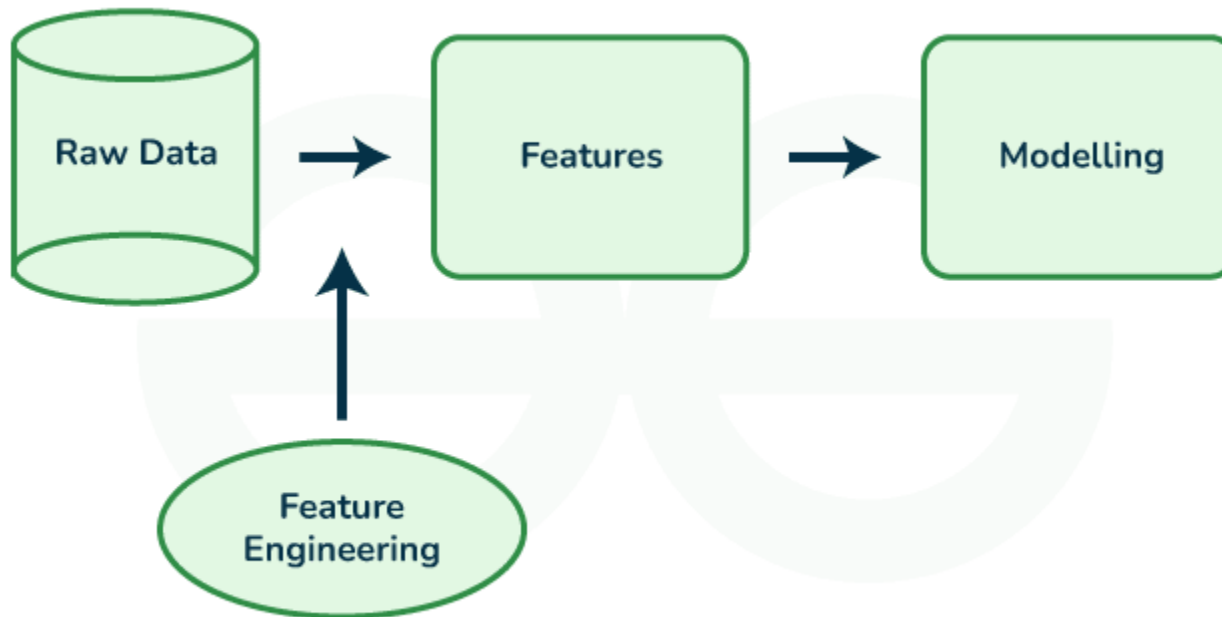
Feature Creation



数据和特征决定了数据挖掘的上限

模型和算法只是帮助我们逼近这个上限

Feature Creation



Feature Creation: New Feature



由**原始数据**创建新的特征，从而更有效地捕捉原始数据中的重要信息

常用方法

- 特征提取 (Feature Extraction)
- 空间映射 (Mapping Data to New Space)
- 特征构造 (Feature Construction)

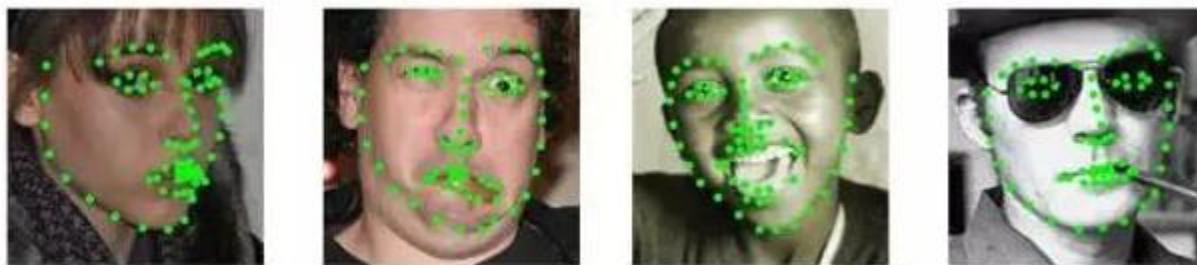
Feature Creation



特征提取（ Feature Extraction ）：由原始数据创建新的特征

例子：对图片是否包含人脸进行二分类

- 原始数据是像素
- 提取人脸相关的边缘特征、区域特征等

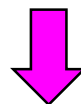
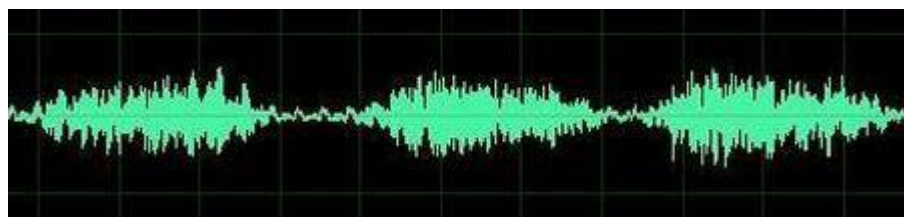


常用的特征提取技术都是针对**具体领域**的
数据挖掘用于新领域时，需开发新的特征提取方法

Feature Creation



将数据进行空间映射，使用不同的视角挖掘数据

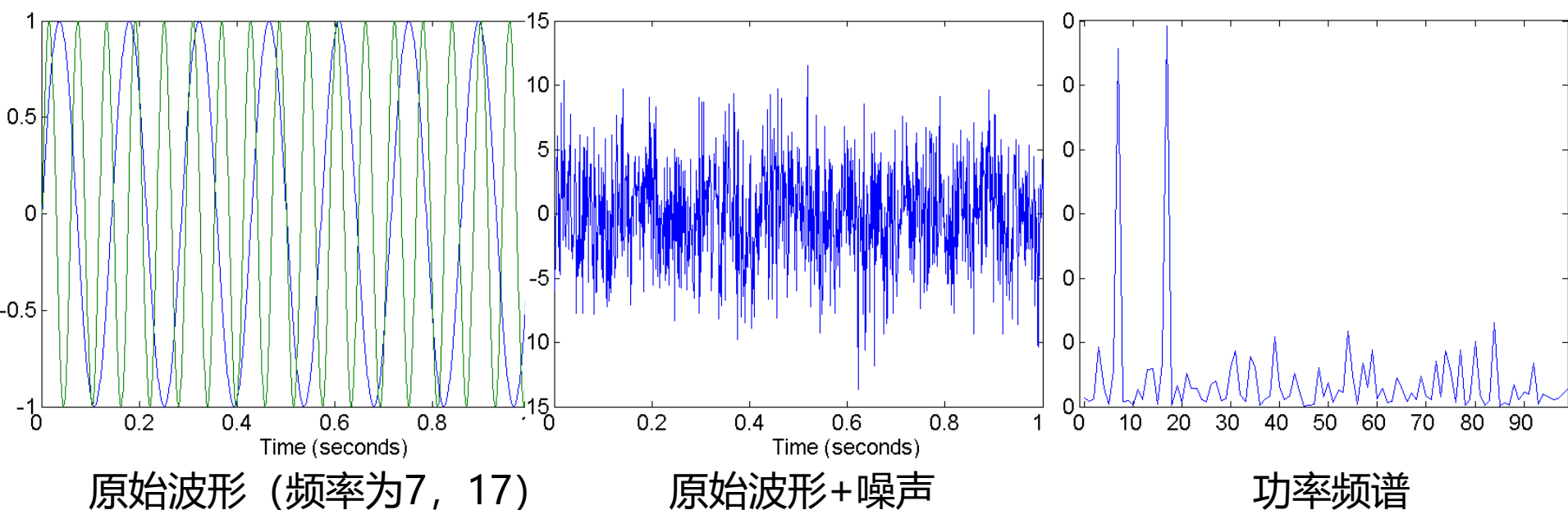


数据空间映射之后，可以更好地提取特征

Feature Creation



空间映射：傅里叶变换

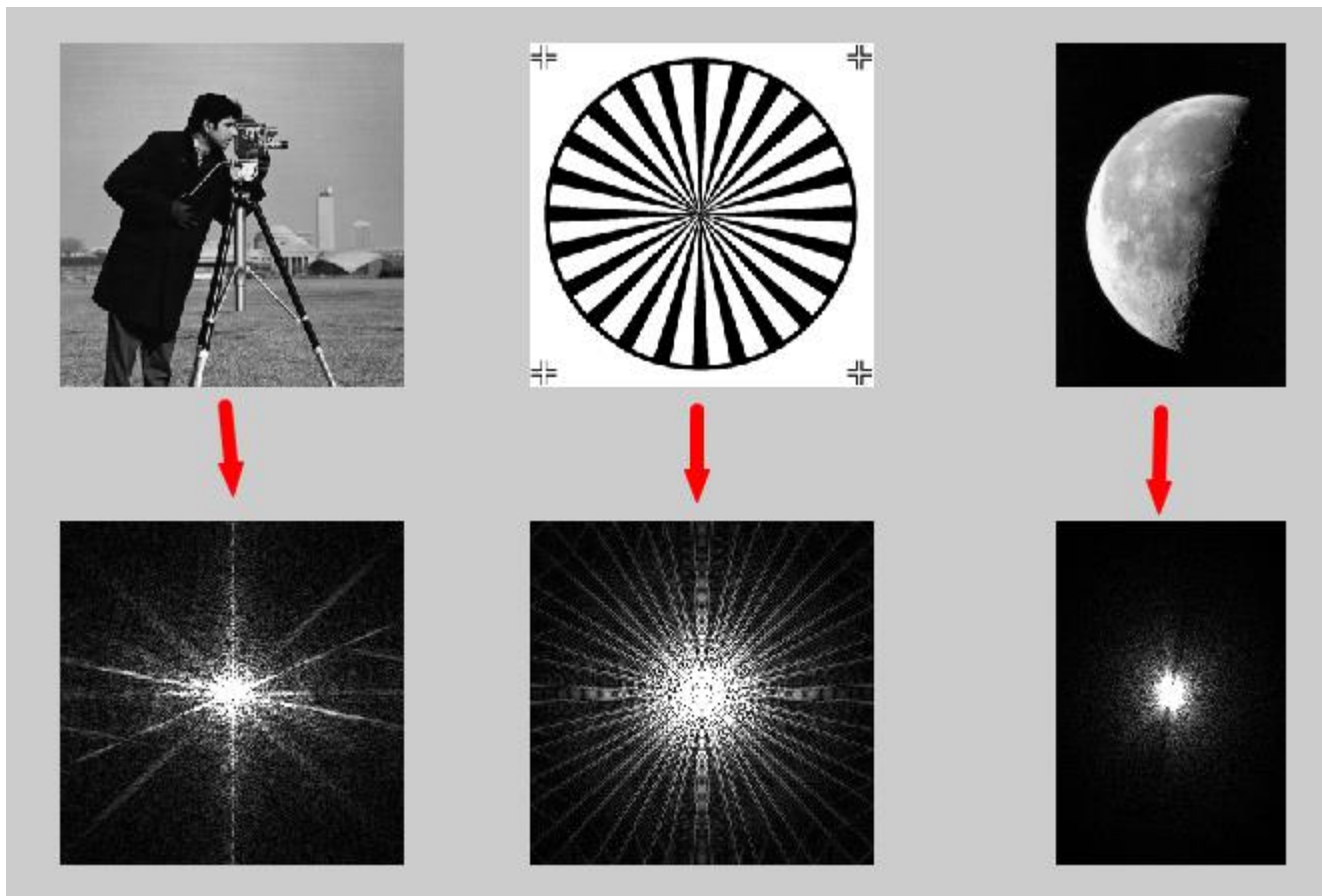


对时间序列实施傅立叶变换，转换成频率信息明显的表示

Feature Creation



空间映射：傅里叶变换



Feature Creation



空间映射：傅里叶变换

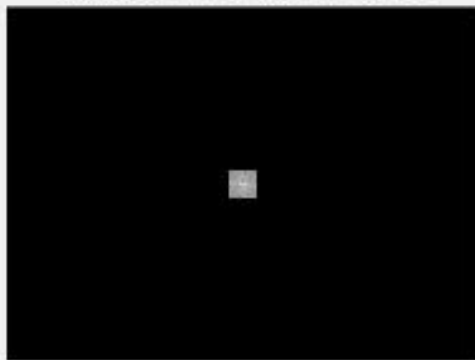
噪声图像



傅里叶变换后幅度图



去除外围幅值后幅度图



去噪后的图像



https://blog.csdn.net/qq_37691909

特征构造 (Feature Construction)：原始特征包含了必要信息，但是形式不适合，因此需由原特征构造新特征

例子：人工制品分类

- 使用不同材料制造：木材、陶土、铜、黄金等
- 希望根据制造材料对它们进行分类
- 原始特征：质量、体积
- 构造的新特征：密度 = 质量 / 体积

常用的方法：使用专家的意见构造特征

有些数据挖掘算法要求输入是二元属性形式

类别特征包括：

- 无序类别 (Categorical)
- 有序类别 (Ordinal)

Feature Creation: 无序类别



无序类别 (Categorical)

分类值	整数值	X_1	X_2	X_3
<i>blue</i>	0	0	0	0
<i>green</i>	1	0	0	1
<i>red</i>	2	0	1	0
<i>black</i>	3	0	1	1
<i>white</i>	4	1	0	0

独热编码(One hot Encoding): 把每个无序特征转化为一个数值向量

分类值	X_1	X_2	X_3	X_4	X_5
<i>blue</i>	1	0	0	0	0
<i>green</i>	0	1	0	0	0
<i>red</i>	0	0	1	0	0
<i>black</i>	0	0	0	1	0
<i>white</i>	0	0	0	0	1

Feature Creation: 有序类别



- 有序类别 (Ordinal)

Status		Vectorization
Bad	→	[1, 0, 0]
Normal		[0, 1, 0]
Good		[0, 0, 1]

- 向量表示方法 (Multi-hot Encoding) : 值之间有顺序的含义

当status特征向量输入模型时, 对于status这个类别特征模型会学习出 w_1, w_2, w_3 三个权重, 如果是good的话将会是 $w_1 w_2 w_3$ 的叠加, 如果取值为bad的话只有 w_1 , 从而体现出有序性

Status		Vectorization
Bad	→	[1, 0, 0]
Normal		[1, 1, 0]
Good		[1, 1, 1]

$$\sum w_i x_i$$

Feature Creation: 特征组合



基本特征仅仅是真实特征分布在低维空间的映射，不足以描述真实分布，加入组合特征是为了在更高维空间拟合真实分布，使得预测更准确

线性模型对于非线性关系缺乏准确刻画，特征组合正好可以加入非线性表达，增强模型的表达能力

例如有两个类别特征color和light，分别取值red, green, blue和on, off。两个特征可以分别转化为3维和2维的向量，对他们做笛卡尔乘积转化后可以组合出6维的向量

X	on	off
red		
green		
blue		

样本	color=red& light=on	color=red& light=off	color=green& light=on	color=green& light=off	color=blue& light=on	color=blue& light=off
1	1	0	0	0	0	0
2	0	1	0	0	0	0
3

可以通过笛卡尔乘积的方式来组合2个或多个特征

Feature Creation



连续特征 (continuous features)

Student ID	Age	Weight(kg)	Height(cm)
0	18	56	174
1	21	61	176
2	25	58	168

连续特征处理：

- MinMax Scale:
$$z_i = \frac{x_i - \min(x)}{\max(x) - \min(x)}$$
- Log transform: $z_i = \log(1 + x_i)$
- 描述统计特征: min max mean median mode std var ...

Feature Creation: Example



一学生属性如下，如何转化为能够输入模型的特征向量？

Student ID	age	weight	height	gender	status
0	18	65	178	M	Bad

连续特征age、weight、height 进行MinMax scale: $[0, 0.29, 1]$

无序特征gender进行One hot Encoding: $[1, 0]$

有序特征status向量化: $[1, 0, 0]$

最终表示学生0的特征向量为: $[0, 0.29, 1, 1, 0, 1, 0, 0]$

Feature Creation: Example



时间信息包含有丰富的数据意义

— 例如: 2017-10-01 16:38:43

Year	Month	Day of Month	Week	Holiday
2017	10	1	Sunday	Yes

Season	Hour Type	Day of Holiday	Hour of Day
Autumn	Afternoon	1	16/24

Feature Creation: Example



地理位置信息包含有丰富的数据意义

— 例如：广东省广州市番禺区大学城

Province	City	Area	Distribution	City-level
广东	广州	番禺区	东南	1

Longitude	Latitude	Area Type	Temperature type
113.23	23.16	学校	Hot

Feature Creation: 案例



国家电网提供了88436名用户，2014~2016年每天的用电数据

基于用户的用电数据，挖掘窃电用户行为特征，识别窃电用户，这是一个二分类问题



案例分析：客户用电异常行为分析

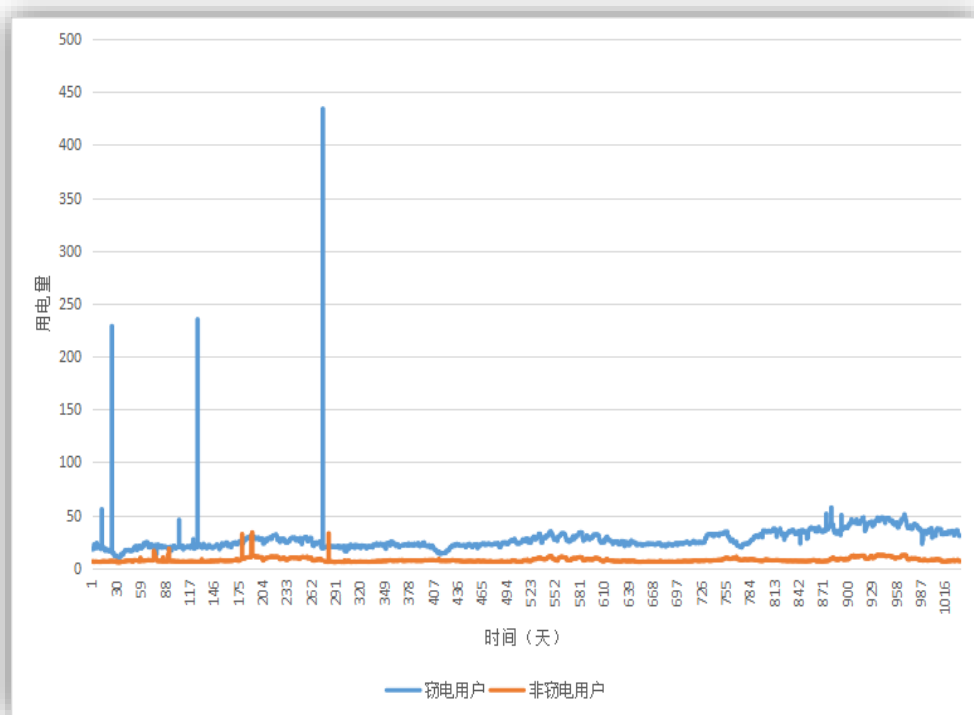


观察问题

- 窃电用户平均用电度数偏高
- 窃电用户用电量瞬时波幅偏高

特征工程

- 用户用电度数**最大值、均值、中位数**



案例分析：客户用电异常行为分析

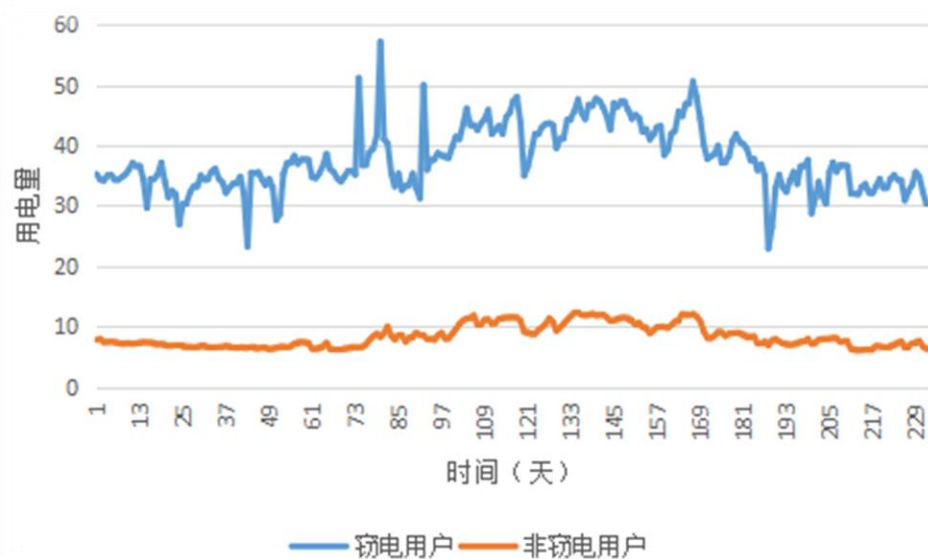


观察问题

- 窃电用户用电量的波动性比较大
- 非窃电用户用电的稳定性比较强

特征工程

- 用户用电度数**标准差、四分位数、异常值的个数**
- 稳定性衡量由前后等长一段时间的数据**相似度**计算



案例分析：客户用电异常行为分析



空间映射

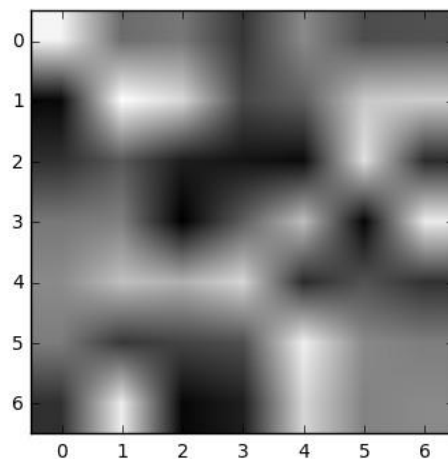
- 用户用电数据由一维向量转换成**二维矩阵**
- 二维矩阵转换成**灰度图**

观察问题

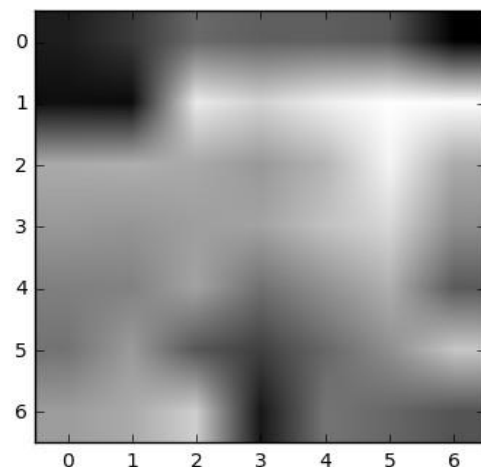
- 窃电用户的图形黑白相间
- 非窃电用户的图形相对空旷

特征工程

- 用卷积神经网络对图像进行自动化特征提取



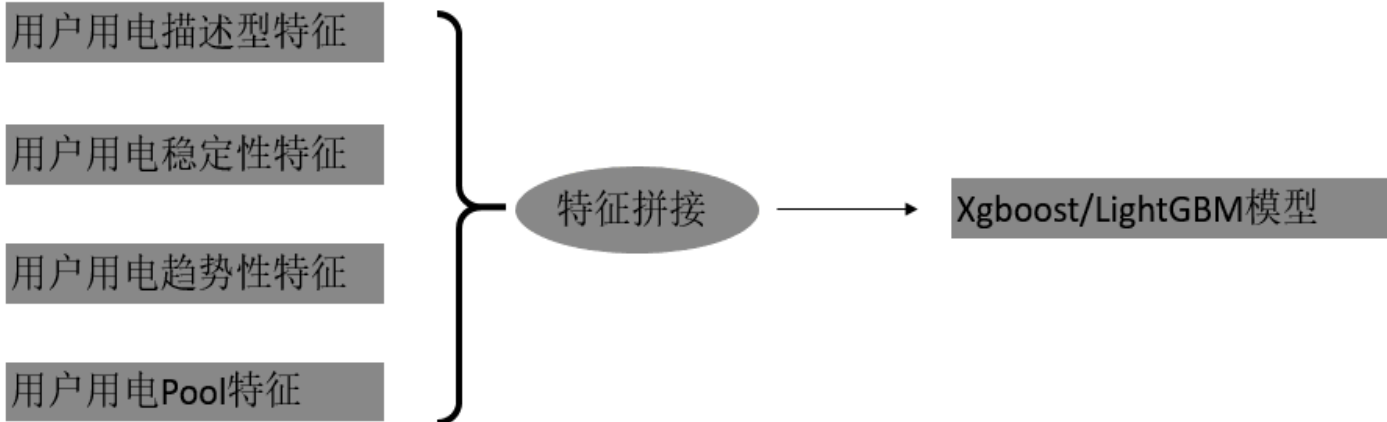
窃电用户



非窃电用户

将7周用电量数据转换成7*7灰度图矩阵

案例分析：客户用电异常行为分析



初赛

A榜			B榜		
排名	队伍名称	最高得分(B)	排名	队伍名称	最高得分(B)
1	TNT_000_	0.95459	1	我们又回来了-美林数据	0.95187
2	我们又回来了-美林数据	0.95187	2	Top	0.95164
3	Top	0.95164	3		

复赛

A榜			B榜		
排名	队伍名称	最高得分(B)	排名	队伍名称	最高得分(B)
1	我们又回来了-美林数据	0.94274	1	TNT_000_	0.92564
2	隐马尔可夫联盟	0.93373	2		
3	打酱油`拎壶冲	0.92871	3		
4	TNT_000_	0.92564	4		



Aggregation

Sampling

Dimensionality Reduction

Discretization

Attribute Transformation

Feature creation

Feature subset selection



Another way to reduce dimensionality of data
减少数据维度

Redundant features 剔除冗余特征

- duplicate much or all of the information contained in one or more other attributes
- Example: purchase price of a product and the amount of sales tax paid

Irrelevant features 剔除不相关特征

- contain no information that is useful for the data mining task at hand
- Example: students' ID is often irrelevant to the task of predicting students' GPA

Feature Subset Selection



Techniques:

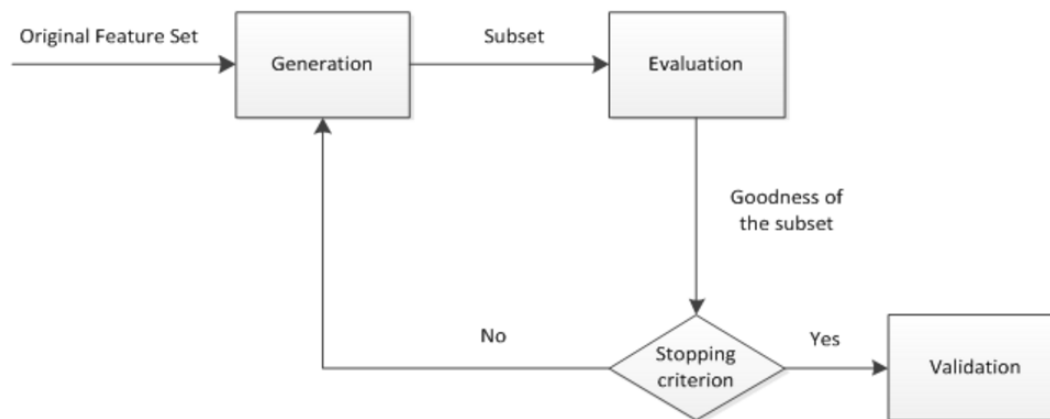
- Brute-force approaches (暴力):
 - ◆ Try all possible feature subsets as input to data mining algorithm
- Filter approaches (过滤):
 - ◆ Features are selected before data mining algorithm is run
- Wrapper approaches (包装):
 - ◆ Use the data mining algorithm as a black box to find best subset of attributes

Feature Subset Selection



Brute-force approaches 暴力筛选最佳特征集:

- (1) 产生过程
- (2) 评价函数
- (3) 停止准则
- (4) 验证过程



Feature Subset Selection



- Filter approaches:
- 思路：特征和目标变量之间的关联
 - 统计检验，如卡方检验、t检验
 - 相关系数，如皮尔森相关系数、
 - 互信息和最大信息系数（MIC）

	适用范围	是否标准化	计算复杂度	鲁棒性
Pearson	线性数据	是	低	低
spearman	线性、简单单调非线性数据	是	低	中等
Kendall	线性、简单单调非线性数据	是	低	中等
阈值相关	线性、非线性数据	是	高	高
最大相关系数	线性、非线性数据	是	高	中等
相位同步相关	时变序列	是	中等	中等
距离相关	线性、非线性数据	是	中等	高
核密度估计(KDE)	线性、非线性数据	否	高	高
k-最邻近距离(KNN)	线性、非线性数据	否	高	高
MIC	线性、非线性数据	是	低	高

Feature Subset Selection



医学数据集:

Leukemia 7129×72

Colon 2000×62

特征样例:

Gene \ sp.	Sample 1 (Cancer)	Sample 2 (Normal)	Sample k
Gene 1	29	19	16
Gene 2	5	17	40
.....
Gene n	13	8	2

Feature Subset Selection



特征选择结果:

Leukemia (SVM)

Number of genes	Train accuracy	Test accuracy
100	100	99.31
50	100	98.276
34	100	99.31
20	100	98.621
10	100	98.621
8	100	96.552
5	100	95.172
3	100	92.759
1	92.093	78.966

Feature Subset Selection



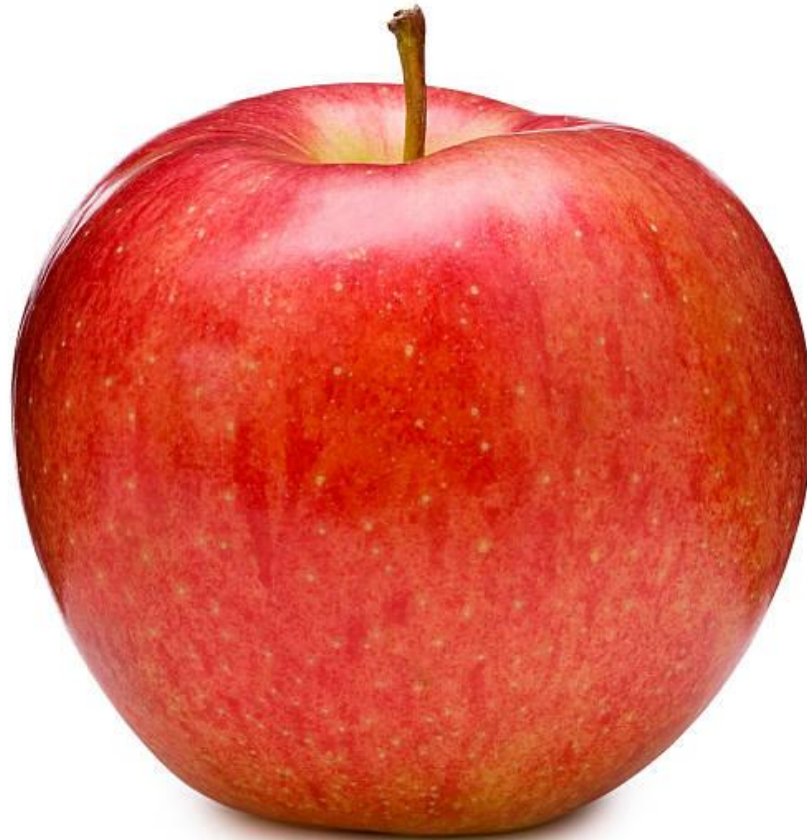
特征选择结果：

Colon (SVM)

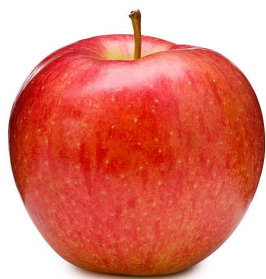
Number of genes	Train accuracy	Test accuracy
100	100	80.4
50	100	80.8
33	100	82
20	100	79.2
10	100	78.8
8	100	77.6
5	99.189	75.6
3	95.405	77.6
1	80	71.6

选择有效的特征能提高预测准确性！

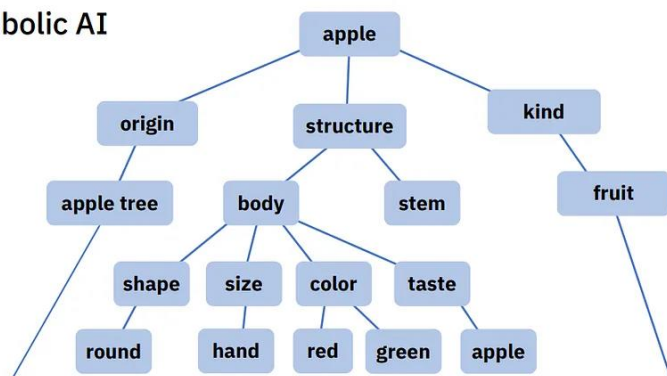
从数据到智能?



数据



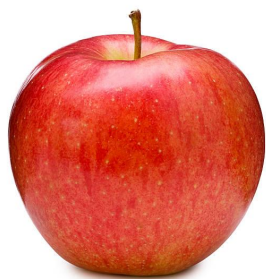
Symbolic AI



用符号系统和规则理解数据

符号主义（Symbolism）：

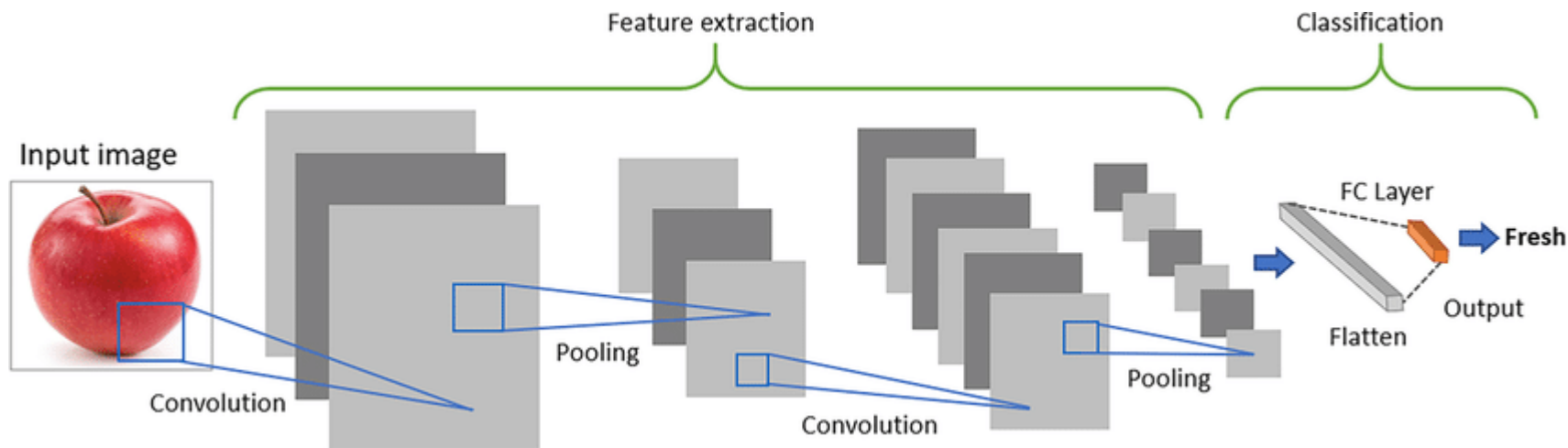
在符号主义的视角下，苹果会被描述为一组符号和属性的集合。比如，定义苹果为一个对象，具有颜色（红色、绿色）、形状（圆形）、质地（光滑）等属性，以及与其他对象（如“树”、“果汁”）的关系（“长在”、“可以制成”）。若要判断某个对象是否为苹果，符号主义的AI系统将检查这个对象是否符合它的符号逻辑规则集合，例如，如果一个对象是红色的、圆形的，并且可以从树上摘下来，那么系统可能会判断这个对象是苹果。



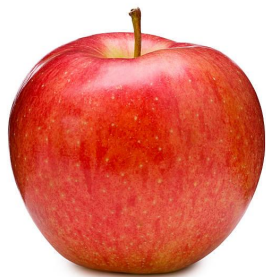
用网络和学习算法来理解数据

连接主义（Connectionism）：

连接主义的角度，苹果会通过神经网络来识别，这个网络通过大量的苹果图片训练而成。在这种情况下，苹果不是由明确的符号属性定义的，而是通过其外观的模式来识别。网络会分析苹果的图像数据，通过不同层次的特征检测（如边缘、颜色、形状等）来识别苹果。这种方法不需要显式定义“苹果”的属性，而是通过从实例中学习来间接理解和识别苹果。



从数据到智能



用环境和反馈来理解数据

行为主义（Behaviorism）：

在行为主义视角下，AI系统不会关注苹果的内部属性或图像模式，而是会学习苹果的行为关联，比如苹果与其它对象或环境的交互。例如，如果一个对象被拿起并被吃掉，而吃它的行为者表现出满足的反应，那么这个对象可能会被标记为“食物”。AI系统通过分析数据学习到，当人们看到圆形、红色的对象时，他们可能会将其拾起并食用，从这些行为上下文中推断出这可能是一个苹果。





Thanks