

# EE6435 Lecture 2

Assignment Project Exam Help

**1. Basic data properties**

<https://powcoder.com>

**2. Data exploration techniques**

Add WeChat powcoder

**3. Introduction to classification problems**

# Outline

- Data properties
  - Attributes and Objects
  - Types of Data
  - Data Quality
- Basic data exploration techniques
- Introduction to classification problems
  - Decision tree

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Most of the slides are from “Introduction to data mining”

# 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 is also known as variable, field, characteristic, dimension, or feature
- A collection of attributes describe an **object**
  - Object is also known as record, point, case, sample, entity, or instance

**Attributes**

Tid	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

**Objects**

# A More Complete View of Data

- Data may have parts

Assignment Project Exam Help

- The different parts of the data may have relationships

<https://powcoder.com>

- More generally, data may have structure

Add WeChat powcoder

- Data can be incomplete

- We will discuss these in more details later

# Attribute Values

- **Attribute values** are numbers or symbols assigned to an attribute for a particular object

Assignment Project Exam Help

- 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

# Types of Attributes

- There are different types of attributes
  - **Nominal (categorical, no order)**
    - Examples: ID numbers, eye color, zip codes
  - **Ordinal**
    - Examples: rankings (e.g., taste of potato chips on a scale from 1-10), grades, height {tall, medium, short}
  - **Interval**
    - Examples: calendar dates, temperatures in Celsius or Fahrenheit.
  - **Ratio**
    - Examples: length, counts

# Properties of Attribute Values

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

- Distinctness:  $= \neq$

- Order:

- Differences are meaningful :

- Ratios are meaningful  $* /$

Assignment Project Exam Help

+ -

<https://powcoder.com>

Add WeChat powcoder

- Nominal attribute: distinctness
- Ordinal attribute: distinctness & order
- Interval attribute: distinctness, order & meaningful differences
- Ratio attribute: all 4 properties/operations

		Attribute Type	Description	Examples	Operations
Categorical	Qualitative	Nominal	Nominal attribute values only distinguish. (=, ≠)	zip codes, employee ID numbers, eye color, sex: { <i>male</i> , <i>female</i> }	mode, entropy, contingency correlation, $\chi^2$ test
		Ordinal	Ordinal attribute values also 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
Numeric	Quantitative	Interval	For interval attributes, differences between values are meaningful. (+, -)	calendar dates, temperature in Celsius or Fahrenheit	mean, standard deviation, Pearson's correlation, <i>t</i> and <i>F</i> tests
		Ratio	For ratio variables, both differences and ratios are meaningful. (*, /)	temperature in Kelvin, monetary quantities, counts, age, mass, length, current	geometric mean, harmonic mean, percent variation

This categorization of attributes is due to S. S. Stevens



# Discrete and Continuous Attributes

- Discrete 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

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

- Continuous 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.

# In-class exercise 1

- You need to design a student record database for students at CityU.
- Question: what attributes do you want to choose for each student record? What are the possible attribute values for each attribute? Design the attributes to cover nominal, ordinal, interval, and ratio.

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

# Types of data sets

- Record
  - Data Matrix
  - Document Data
  - Transaction Data
- Graph
  - World Wide Web
  - Molecular Structures
- Ordered
  - Spatial Data
  - Temporal Data
  - Sequential Data
  - Genetic Sequence Data

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

# Important Characteristics of Data

- Dimensionality (number of attributes)
  - High dimensional data brings a number of challenges
- Sparsity [Assignment Project Exam Help](https://powcoder.com)
  - Only present <https://powcoder.com>
- Resolution [Add WeChat powcoder](https://powcoder.com)
  - Patterns depend on the scale
- Size
  - Type of analysis may depend on size of data

# Record Data

- 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

# 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

Assignment Project Exam Help

- 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

Add WeChat powcoder

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

# 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

# Transaction Data

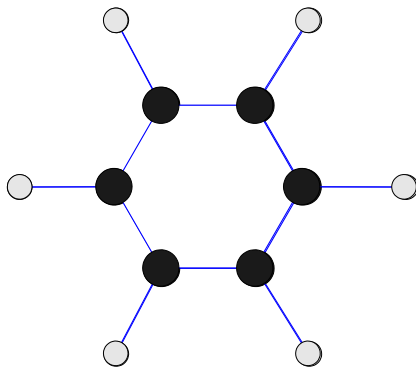
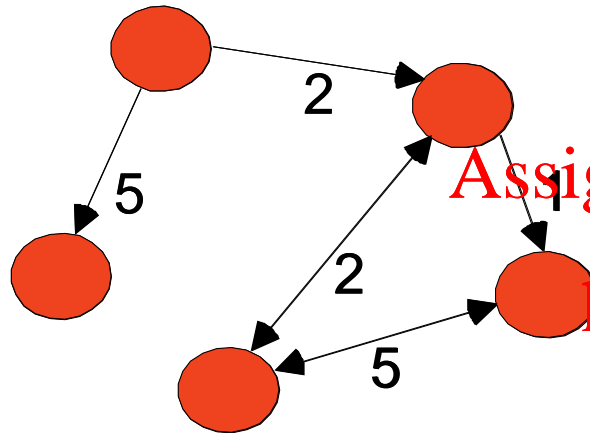
- A special type of record data, where
  - Each record (transaction) involves a set of items.
  - For example, consider a grocery store. The set of products purchased by a customer during one shopping trip constitute a transaction, while the individual products that were purchased are the items.

<i>TID</i>	<i>Items</i>
1	Bread, Coke, Milk
2	Beer, Bread
3	Beer, Coke, Diaper, Milk
4	Beer, Bread, Diaper, Milk
5	Coke, Diaper, Milk



# Graph Data

- Examples: Generic graph, a molecule, and webpages



Benzene Molecule: C<sub>6</sub>H<sub>6</sub>

## Useful Links:

- [Bibliography](#)
- [Other Useful Websites](#)
  - [ACM SIGKDD](#)
  - [KDDuggets](#)
  - [The Data Mine](#)

## Knowledge Discovery and Data Mining Bibliography

(Gets updated frequently, so visit often!)

- [Books](#)
- [General Data Mining](#)

## Book References in Data Mining and Knowledge Discovery

Usama Fayyad, Gregory Piatetsky-Shapiro, Padhraic Smyth, and Ramasamy uthurasamy, "Advances in Knowledge Discovery and Data Mining", AAAI Press/the MIT Press, 1996.

J. Ross Quinlan, "C4.5: Programs for Machine Learning", Morgan Kaufmann Publishers, 1993.  
Michael Berry and Gordon Linoff, "Data Mining Techniques (For Marketing, Sales, and Customer Support)", John Wiley & Sons, 1997.

## General Data Mining

Usama Fayyad, "Mining Databases: Towards Algorithms for Knowledge Discovery", Bulletin of the IEEE Computer Society Technical Committee on data Engineering, vol. 21, no. 1, March 1998.

Christopher Matheus, Philip Chan, and Gregory Piatetsky-Shapiro, "Systems for knowledge Discovery in databases", IEEE Transactions on Knowledge and Data Engineering, 5(6):903-913, December 1993.

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

# Ordered Data

- Sequences of transactions

**Items/Events**

Assignment Project Exam Help

(A B) (D) (C E)

(B D) (C) (E)

(C D) (B) (A E)

**An element of  
the sequence**

# Ordered Data

- Genomic sequence data

**GGTTCCGCCTTCAGCCCCGCGCC**  
**CGCAGGGCCCCGCCCGCGCCGTC**  
**GAGAAGGGCCCGCCTGGCGGGCG**  
**GGGGGAGGCGGGGCCGCCCGAGC**  
**CCAACCGAGTCCGACCAGGTGCC**  
**CCCTCTGCTCGGCCTAGACCTGA**  
**GCTCATTAGGCGGCAGCGGACAG**  
**GCCAAGTAGAACACGCGAAGCGC**  
**TGGGCTGCCTGCTGCGACCAGGG**

# Ordered Data

- Spatio-Temporal Data

**Average Monthly  
Temperature of  
land and ocean**

Jan

