

The background of the slide is a complex, abstract composition. It features a network graph with numerous nodes and edges, rendered in shades of red, orange, and green. The nodes are small circles, and the edges are thin lines connecting them. The overall color palette is muted, with a lot of grey and white space. There are also some faint, larger-scale patterns and textures, including a grid of small plus signs in the top left and bottom left corners. The text is centered in a large, bold, black font.

# **Pattern Discovery in Data Mining: Course Overview**

# What Is Pattern Discovery?

---

- ❑ Considering massive shopping transaction data, pattern discovery may help answer the following questions:
  - ❑ What groups of items are frequently bought together?
  - ❑ If a person buys diapers at night, what is the probability of this person buying beer as well?
  - ❑ If a customer buys an iPhone 5 or iPhone 7, what other electronic products will the customer be most likely to buy in the next 3 months?

# The Value of Pattern Discovery

---

- ❑ What is the value of pattern discovery?
  - ❑ Pattern discovery helps you find hidden and inherent data patterns in massive data
  - ❑ Pattern mining will play a unique and critical role in mining massive data!
- ❑ What roles does pattern discovery play in the Data Mining Specialization?
  - ❑ You will learn scalable methods to find patterns (e.g., the set of data items strongly correlated to each other) from massive data
  - ❑ You will learn how to mine a large variety of patterns
  - ❑ You will also learn how to evaluate the value of patterns
  - ❑ Pattern discovery will help classification, clustering and other data mining tasks

# Broad Applications of Pattern Discovery

---

- ❑ Predicting shopping transaction data:
  - ❑ For a customer who buys products A and B, what is the likelihood of the customer buying product C?
- ❑ Predicting webpage click streams:
  - ❑ Now, which webpage is most likely to be clicked next?
- ❑ Mining software bugs: Where is the likely bug in this program?
- ❑ Identifying objects or sub-structures in images, videos, and social media
- ❑ Finding quality phrases, entities, and attributes in massive text
- ❑ Finding repeating DNA and protein sequences in genomes
- ❑ Finding “hidden” communities in a massive social network



# Major Reference Readings for the Course

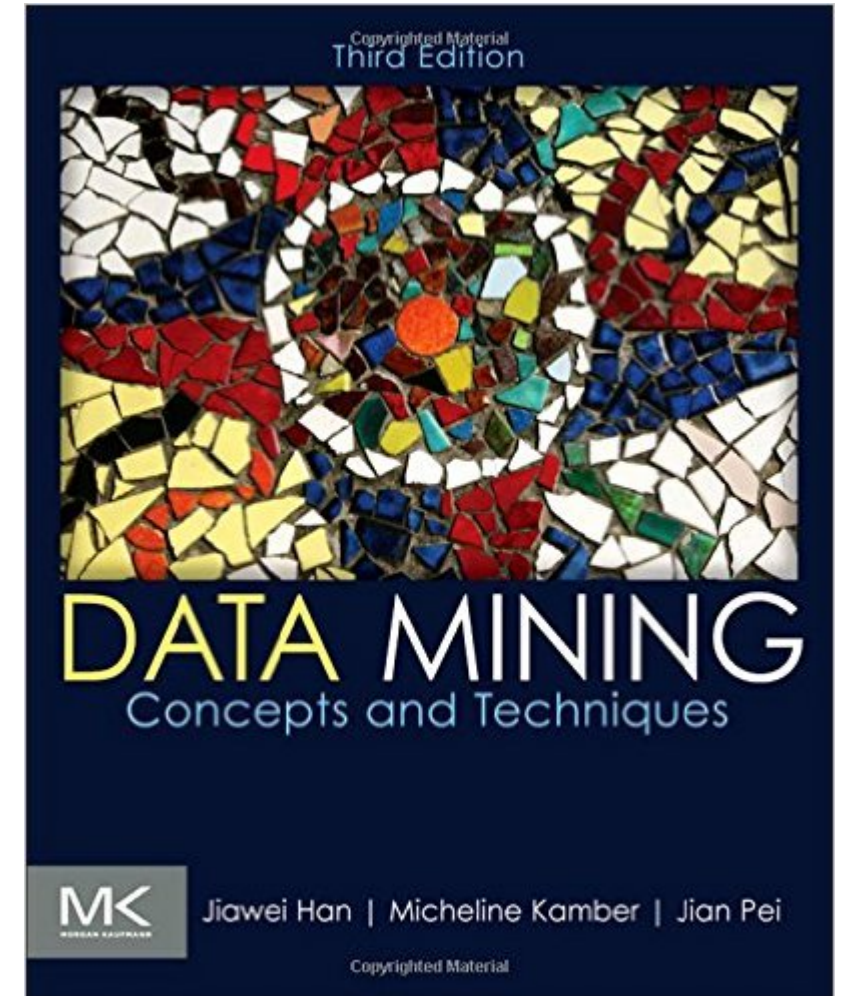
## □ Textbook

- Han, J., Kamber, M., & Pei, J. (2011). *Data Mining: Concepts and Techniques (3<sup>rd</sup> ed)*. Morgan Kaufmann

## □ Chapters most related to the course

- Chapter 1: Introduction
- Chapter 6: Mining Frequent Patterns, Associations, and Correlations: Basic Concepts and Methods
- Chapter 7: Advanced Pattern Mining

- Other references will be listed at the end of each lecture video



The background of the slide is a complex, abstract composition. It features a network of thin, light-colored lines forming a web-like structure. Overlaid on this are various data visualizations: a grid of small, colorful dots (green, blue, orange) and a series of horizontal bars with a color gradient from orange to red. The overall aesthetic is technical and data-driven.

# **What Is Pattern Discovery? Why Is It Important?**

# What Is Pattern Discovery?

---

## ❑ What are patterns?

- ❑ **Patterns**: A set of items, subsequences, or substructures that occur frequently together (or strongly correlated) in a data set
- ❑ Patterns represent **intrinsic** and **important properties** of datasets

## ❑ **Pattern discovery**: Uncovering patterns from massive data sets

## ❑ Motivation examples:

- ❑ What products were often purchased together?
- ❑ What are the subsequent purchases after buying an iPad?
- ❑ What code segments likely contain copy-and-paste bugs?
- ❑ What word sequences likely form phrases in this corpus?

# Pattern Discovery: Why Is It Important?

---

- ❑ Finding **inherent regularities** in a data set
- ❑ **Foundation** for many essential data mining tasks
  - ❑ Association, correlation, and causality analysis
  - ❑ Mining sequential, structural (e.g., sub-graph) patterns
  - ❑ Pattern analysis in spatiotemporal, multimedia, time-series, and stream data
  - ❑ Classification: Discriminative pattern-based analysis
  - ❑ Cluster analysis: Pattern-based subspace clustering
- ❑ Broad applications
  - ❑ Market basket analysis, cross-marketing, catalog design, sale campaign analysis, Web log analysis, biological sequence analysis



The background features a complex, abstract design. It includes a grid of small plus signs, a network of red lines connecting green dots, and a large, faint, stylized letter 'A' in the center. The overall color palette is muted, with shades of brown, beige, and light blue.

# **Basic Concepts: Frequent Patterns and Association Rules**

# Basic Concepts: Frequent Itemsets (Patterns)

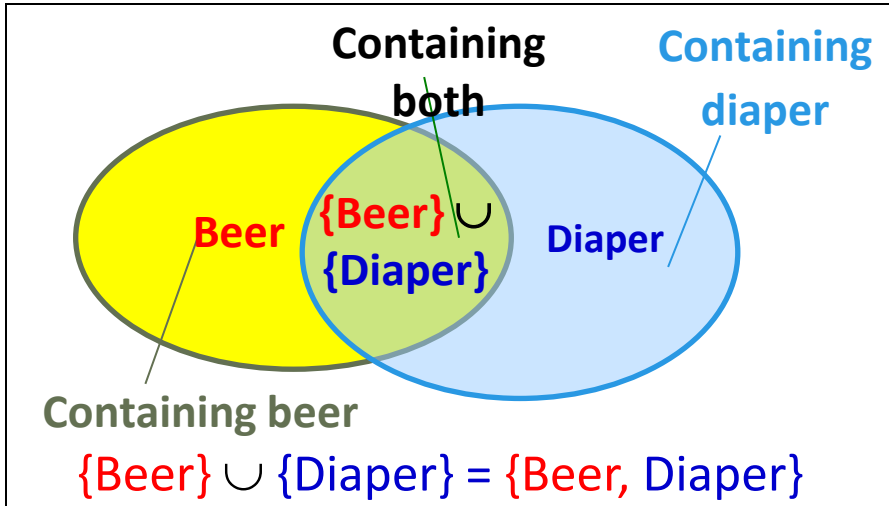
- ❑ **Itemset**: A set of one or more items
- ❑ **k-itemset**:  $X = \{x_1, \dots, x_k\}$
- ❑ **(absolute) support (count)** of  $X$ :  
Frequency or the number of occurrences of an itemset  $X$
- ❑ **(relative) support**,  $s$ : The fraction of transactions that contains  $X$  (i.e., the **probability** that a transaction contains  $X$ )
- ❑ An itemset  $X$  is **frequent** if the support of  $X$  is no less than a *minsup* threshold (denoted as  $\sigma$ )

Tid	Items bought
10	Beer, Nuts, Diaper
20	Beer, Coffee, Diaper
30	Beer, Diaper, Eggs
40	Nuts, Eggs, Milk
50	Nuts, Coffee, Diaper, Eggs, Milk

- ❑ Let *minsup* = 50%
- ❑ Freq. 1-itemsets:
  - ❑ Beer: 3 (60%); Nuts: 3 (60%)
  - ❑ Diaper: 4 (80%); Eggs: 3 (60%)
- ❑ Freq. 2-itemsets:
  - ❑ {Beer, Diaper}: 3 (60%)

# From Frequent Itemsets to Association Rules

Tid	Items bought
10	Beer, Nuts, Diaper
20	Beer, Coffee, Diaper
30	Beer, Diaper, Eggs
40	Nuts, Eggs, Milk
50	Nuts, Coffee, Diaper, Eggs, Milk



Note: Itemset:  $X \cup Y$ , a subtle notation!

- Association rules:  $X \rightarrow Y (s, c)$ 
  - **Support**,  $s$ : The probability that a transaction contains  $X \cup Y$
  - **Confidence**,  $c$ : The conditional probability that a transaction containing  $X$  also contains  $Y$
  - $c = \text{sup}(X \cup Y) / \text{sup}(X)$
- **Association rule mining**: Find **all** of the rules,  $X \rightarrow Y$ , with minimum support and confidence
- Frequent itemsets: Let  $\text{minsup} = 50\%$ 
  - Freq. 1-itemsets: Beer: 3, Nuts: 3, Diaper: 4, Eggs: 3
  - Freq. 2-itemsets:  $\{\text{Beer, Diaper}\}$ : 3
- Association rules: Let  $\text{minconf} = 50\%$ 
  - $\text{Beer} \rightarrow \text{Diaper}$  (60%, 100%)
  - $\text{Diaper} \rightarrow \text{Beer}$  (60%, 75%) (Q: Are these all rules?)





# **+ + Compressed Representation: Closed Patterns and Max- Patterns**



# Challenge: There Are Too Many Frequent Patterns!

- A long pattern contains a combinatorial number of sub-patterns
- How many frequent itemsets does the following TDB<sub>1</sub> contain?

□ TDB<sub>1</sub>:      T<sub>1</sub>: {a<sub>1</sub>, ..., a<sub>50</sub>}; T<sub>2</sub>: {a<sub>1</sub>, ..., a<sub>100</sub>}

□ Assuming (absolute) *minsup* = 1

□ Let's have a try

1-itemsets: {a<sub>1</sub>}: 2, {a<sub>2</sub>}: 2, ..., {a<sub>50</sub>}: 2, {a<sub>51</sub>}: 1, ..., {a<sub>100</sub>}: 1,

2-itemsets: {a<sub>1</sub>, a<sub>2</sub>}: 2, ..., {a<sub>1</sub>, a<sub>50</sub>}: 2, {a<sub>1</sub>, a<sub>51</sub>}: 1 ..., ..., {a<sub>99</sub>, a<sub>100</sub>}: 1,

..., ..., ..., ...

99-itemsets: {a<sub>1</sub>, a<sub>2</sub>, ..., a<sub>99</sub>}: 1, ..., {a<sub>2</sub>, a<sub>3</sub>, ..., a<sub>100</sub>}: 1

100-itemset: {a<sub>1</sub>, a<sub>2</sub>, ..., a<sub>100</sub>}: 1

□ In total:  $\binom{100}{1} + \binom{100}{2} + \dots + \binom{100}{100} = 2^{100} - 1$  sub-patterns!

A too huge set for  
any computer to  
compute or store!



# Expressing Patterns in Compressed Form: Closed Patterns

---

- ❑ How to handle such a challenge?
- ❑ Solution 1: **Closed patterns**: A pattern (itemset)  $X$  is **closed** if  $X$  is *frequent*, and there exists *no super-pattern*  $Y \supset X$ , with the same support as  $X$ 
  - ❑ Let Transaction DB  $TDB_1$ :  $T_1: \{a_1, \dots, a_{50}\}$ ;  $T_2: \{a_1, \dots, a_{100}\}$
  - ❑ Suppose  $minsup = 1$ . How many closed patterns does  $TDB_1$  contain?
    - ❑ Two:  $P_1: \{\{a_1, \dots, a_{50}\}: 2\}$ ;  $P_2: \{\{a_1, \dots, a_{100}\}: 1\}$
- ❑ **Closed pattern** is a **lossless compression** of frequent patterns
  - ❑ Reduces the # of patterns but does not lose the support information!
  - ❑ You will still be able to say:  $\{\{a_2, \dots, a_{40}\}: 2\}$ ,  $\{\{a_5, a_{51}\}: 1\}$

# Expressing Patterns in Compressed Form: Max-Patterns

---

- ❑ Solution 2: **Max-patterns**: A pattern  $X$  is a **max-pattern** if  $X$  is frequent and there exists no frequent super-pattern  $Y \supset X$
- ❑ Difference from close-patterns?
  - ❑ Do not care the real support of the sub-patterns of a max-pattern
  - ❑ Let Transaction DB  $TDB_1$ :  $T_1: \{a_1, \dots, a_{50}\}$ ;  $T_2: \{a_1, \dots, a_{100}\}$
  - ❑ Suppose  $minsup = 1$ . How many max-patterns does  $TDB_1$  contain?
    - ❑ One:  $P: \{\{a_1, \dots, a_{100}\}: 1\}$
- ❑ **Max-pattern** is a **lossy compression**!
  - ❑ We only know  $\{a_1, \dots, a_{40}\}$  is frequent
  - ❑ But we do not know the real support of  $\{a_1, \dots, a_{40}\}$ , ..., any more!
- ❑ Thus in many applications, mining close-patterns is more desirable than mining max-patterns

# Recommended Readings

---

- ❑ R. Agrawal, T. Imielinski, and A. Swami, “Mining association rules between sets of items in large databases”, in Proc. of SIGMOD'93
- ❑ R. J. Bayardo, “Efficiently mining long patterns from databases”, in Proc. of SIGMOD'98
- ❑ N. Pasquier, Y. Bastide, R. Taouil, and L. Lakhal, “Discovering frequent closed itemsets for association rules”, in Proc. of ICDT'99
- ❑ J. Han, H. Cheng, D. Xin, and X. Yan, “Frequent Pattern Mining: Current Status and Future Directions”, Data Mining and Knowledge Discovery, 15(1): 55-86, 2007