# INTRODUCTION TO DATA MINING

BEIYU LIN

# COURSE INFORMATION
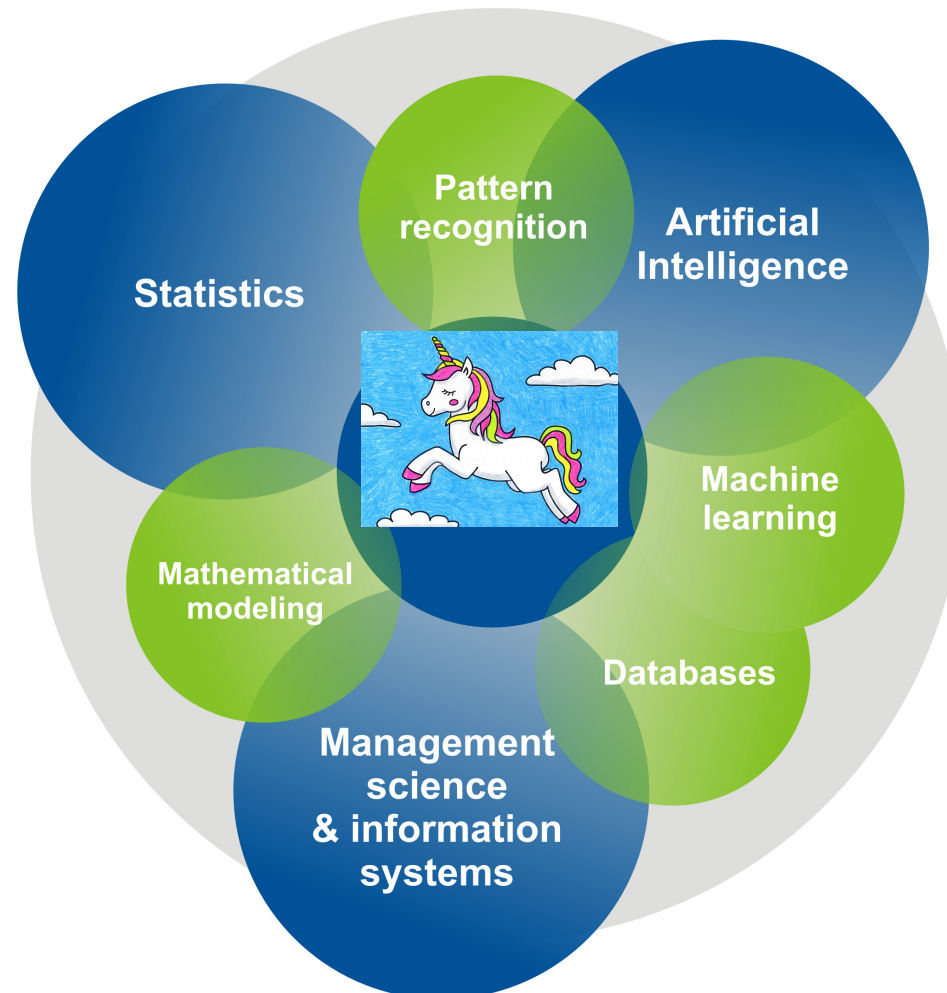
Outline: https://beiyulincs.github.io/teach/fall_21/dm.html

Syllabus: https://beiyulincs.github.io/teach/fall_21/syllabus_cs_458.pdf

# OUTLINE

- **What is Data Mining?**

- **Why Data Learning is important?**

- **Data Learning and its Applications**

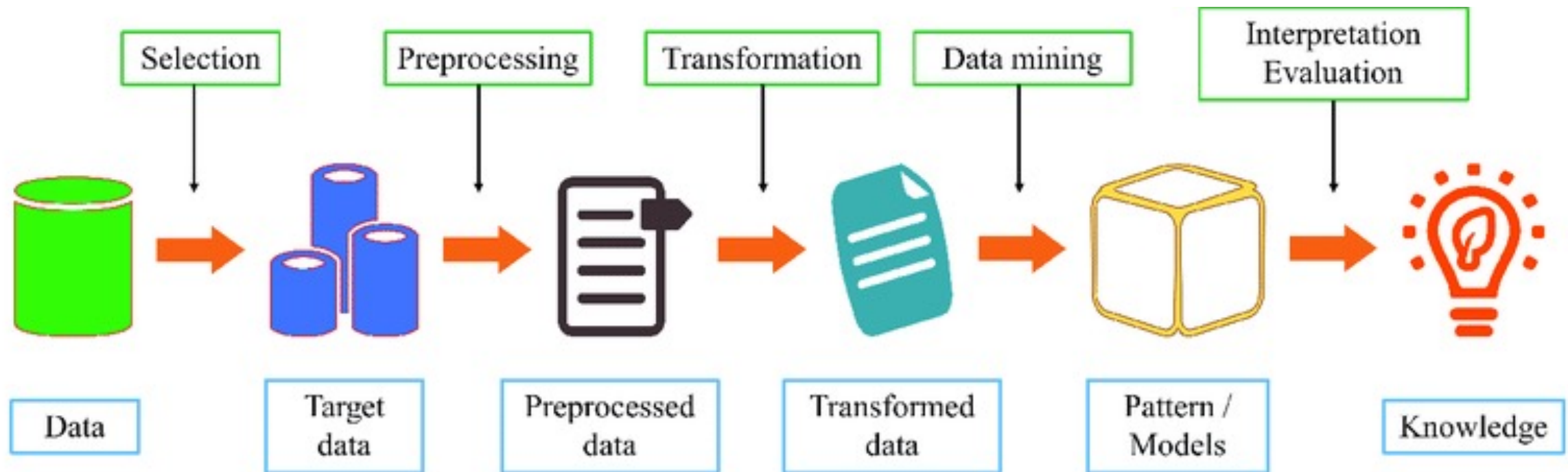- **Real Life Examples**

# DATA MINING

# WHAT IS DATA MINING



The process of discovering meaningful **new**
- correlations
- Patterns
- Trends

By learning from large amounts of stored data
**Via**
- Pattern recognition
- Statistical
- Mathematical
- Machine learning methods

# DATA MINING PROCESS

# WHAT IS DATA MINING

- Data mining (knowledge discovery from data (KDD))
  - – Extraction of interesting (non-trivial, implicit, previously unknown and potentially useful) patterns or knowledge from huge amount of data
  - – Data mining: a misnomer (outliner)?

- • Alternative names
  - – Knowledge discovery in data,
  - – knowledge extraction,
  - – data/pattern analysis,
  - – data archeology,
  - – data dredging, information harvesting, business intelligence, etc.

# DATA MINING IN REAL LIFE



Image Classification

# DATA MINING IN REAL LIFE

Document Categorization

Speech Recognition

Protein Classification

Fraud Detection

Playing Games

Spam Detection

# WHY DATA MINING

"The world is one big data problem."
(by Andrew McAfee, co-director of the MIT Initiative)

"Data is the new science. Big Data holds the answers." (Pat Gelsinger, CEO, VMWare)

transforming raw data into useful knowledge



- understanding, integrated, actionable

**wisdom**

*given insight, becomes*

- contextual, synthesized, learning

**knowledge**

*given meaning becomes*

- useful, organised, structured

**information**

*given context, becomes*
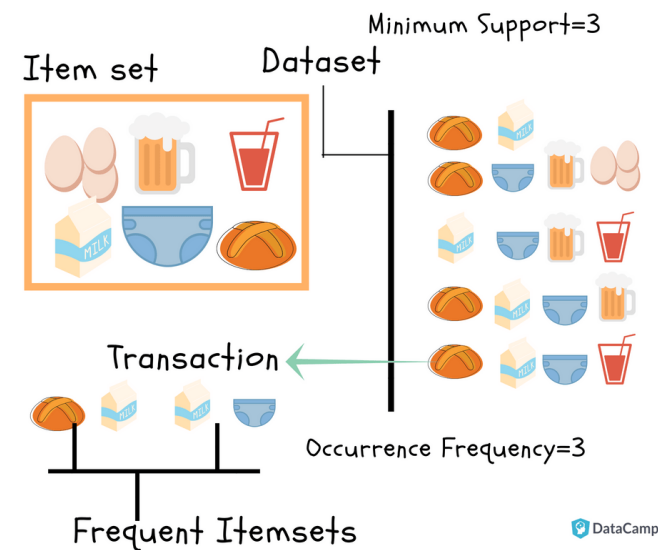
- signals, know-nothing

**data**

# WHY DATA MINING

- Data analysis and decision support
  - – Market analysis and management
    - • Target marketing, customer relationship management (CRM), market basket analysis, cross selling, market segmentation
  - – Risk analysis and management
    - • Forecasting, customer retention, improved underwriting, quality control, competitive analysis
  - – Fraud detection and detection of unusual patterns (outliers)

- Other Applications
  - – Text mining (news group, email, documents) and Web mining
  - – Stream data mining
  - – DNA and bio-data analysis
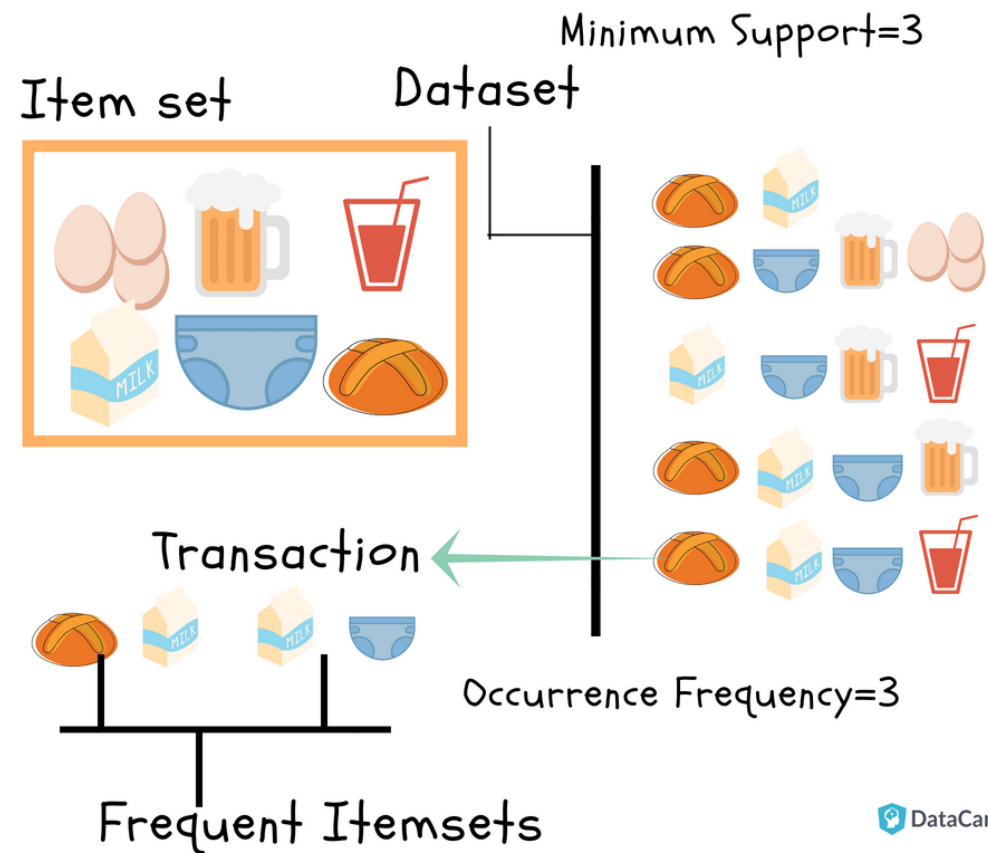
REAL TIME STREAM BIG DATA PROCES

# DATA MINING ALGORITHMS

- Association rule

- Categorization (supervised learning)

- Clustering (unsupervised learning)

- Mining Internet of Things (IoT) data

# ASSOCIATION RULE



Item set

Dataset

Minimum Support=3

Transaction

Occurrence Frequency=3

Frequent Itemsets

DataCamp

# SUPERVISED CLASSIFICATION

Decide whom credit card application should be approved.





Goal: use a person's information seen so far to produce good prediction rule for future applications.

# MODELS – SUPERVISED LEARNING

- **Learn a classification model** from the data

- Use the model to classify future loan applications into

  - Yes (approved) and

  - No (not approved)

- What is the class for following case/instance?

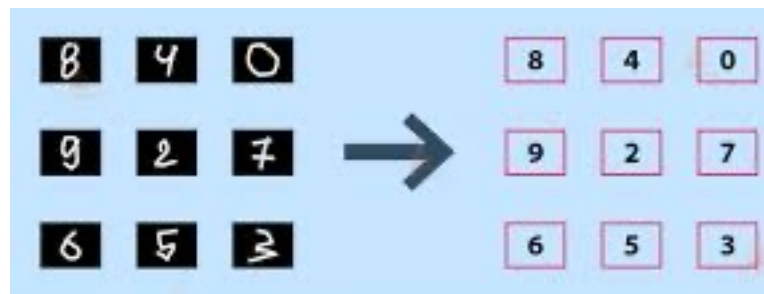| Age | Has_Job | Own_house | Credit-Rating | Class |
|-----|---------|-----------|---------------|-------|
| young | false | false | good | ? |

# MODELS – SUPERVISED LEARNING (IMAGE CLASSIFICATION)

Face Detection and Recognition



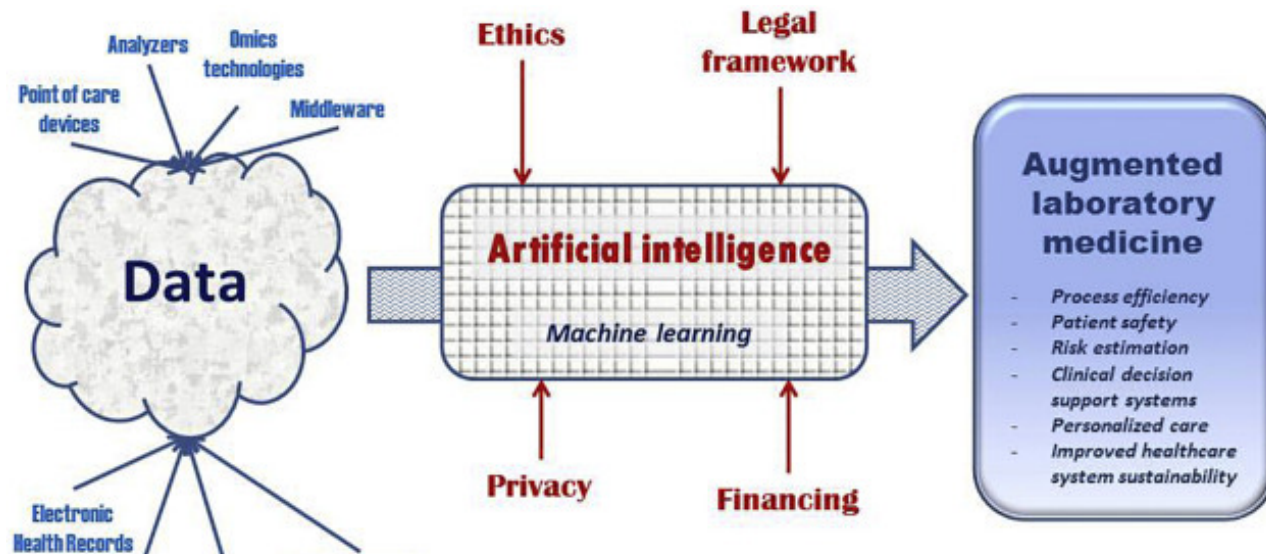Handwritten digit recognition (convert hand-written digits to characters 0..9)

# MODELS – SUPERVISED LEARNING (OTHER EXAMPLES)

## Weather Prediction



## Medicine



Computational Economics:
– predict if a stock will rise or fall
– predict if a user will click on an ad or not
• in order to decide which ad to show

# MODELS – SUPERVISED LEARNING (REGRESSION)
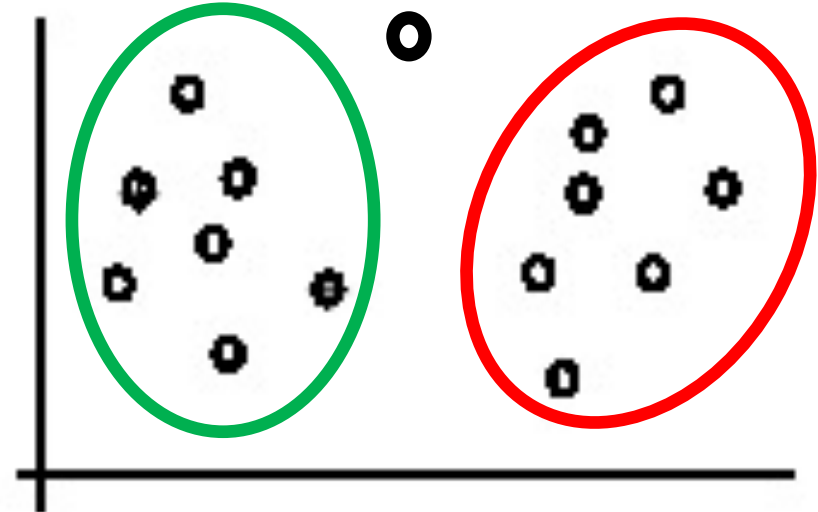
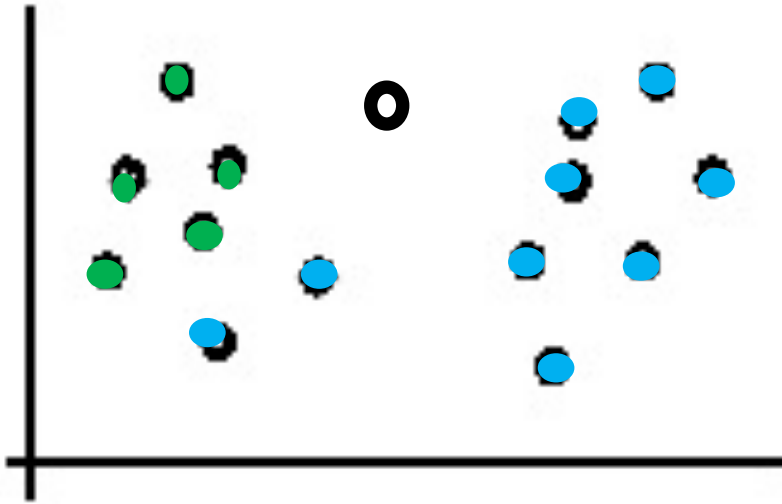Regression: Predicting a numeric value

Stock market



Weather prediction

Thank you!