



University of Colorado
Boulder

CSCI 4502/5502

Data Mining

Fall 2020
Lecture 02 (Aug 27)

Review

- ◆ Introduction: Instructor, class
- ◆ Administrative information
- ◆ Course overview
- ◆ Policies
- ◆ Chapter I: Introduction to Data Mining





University of Colorado
Boulder

Chapter I

Introduction to Data Mining

Into the Digital Era

- ◆ People's daily lives

- ◆ 4 billion Internet users
 - ◆ 500 million tweets/day



- ◆ Scientific discovery

- ◆ LHC: 15 PB/year; LSST: 20 TB/night



- ◆ IDC Digital Universe Report

- ◆ 0.8ZB (2009) => 35ZB (2020)
 - ◆ 4.4ZB (2013) => 44ZB (2020)



University of Colorado
Boulder

Fall 2020 Data Mining

Why Data Mining?

- ◆ Data explosion: KB, MB, GB, TB, PB, EB, ZB, ...
- ◆ data creation, transmission, storage, sharing, processing
- ◆ We are drowning in data, but starving for knowledge!
- ◆ Need automated analysis of massive data



What Is Data Mining?

- ◆ Data mining (knowledge discovery from data)
- ◆ extraction of interesting patterns or knowledge from huge amounts of data
- ◆ **interesting:** valid, previously unknown, potentially useful, ultimately understandable by human
- ◆ **huge amounts of data:** scalability, efficiency



DM Application Areas

◆ Science

- ◆ astrophysics, bioinformatics, drug discovery, sustainable energy, oceanography, seismology, ...

◆ Business

- ◆ market analysis, fraud detection, target marketing, churn prediction, product recommendation, ...

◆ Web

- ◆ search engines, advertising, online social networks, trending, ...

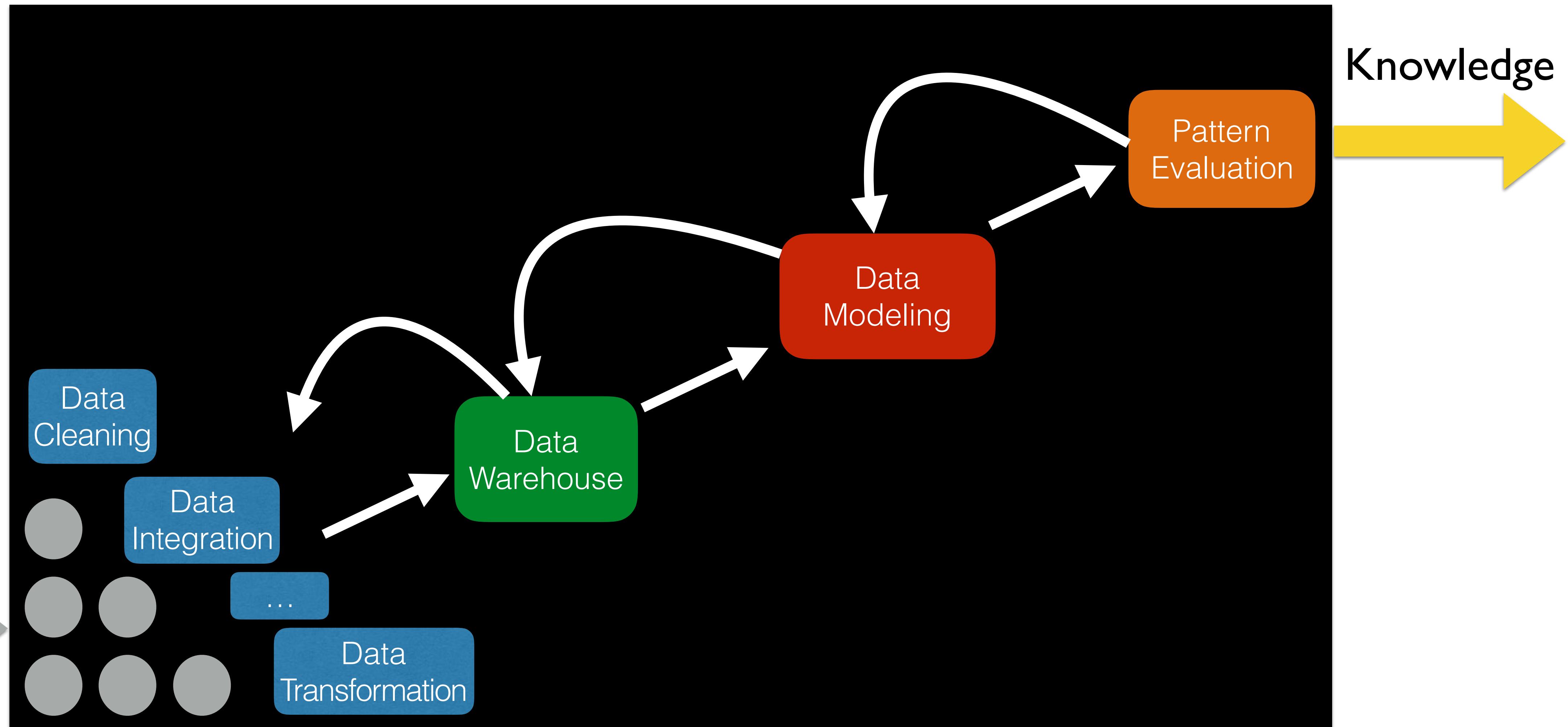
◆ Government

- ◆ surveillance, crime detection, transportation, development, ...

◆ And a lot more!



Data Mining Pipeline



Data Mining: Various Views

- ◆ **Data view**

- ◆ types of data to be mined

- ◆ **Knowledge view**

- ◆ types of knowledge to be discovered

- ◆ **Method view**

- ◆ types of techniques utilized

- ◆ **Application view**

- ◆ types of applications adapted



Data View (I)

- ◆ The 3Vs, 4Vs, and 5Vs

Volume

Variety

Veracity

Velocity

Value



Data View (2)

- ◆ Database-oriented
 - ◆ relational, transactional
 - ◆ data warehouse, NoSQL
- ◆ Sequence, stream, temporal, time-series data
 - ◆ trend analysis, anomaly
- ◆ Spatial, spatial-temporal data
- ◆ Text, multimedia, Web data
 - ◆ topic detection, similarity, popularity, sentiment
- ◆ Graph, social networks data
 - ◆ substructures, shared interests, influencers, information diffusion



Knowledge View

- ◆ Concept/class description
- ◆ Frequent patterns, associations, correlations
- ◆ Classification and prediction
- ◆ Cluster analysis
- ◆ Outlier analysis
- ◆ Evolution analysis



Concept/Class Description

- ◆ Data characterization (summarization)
- ◆ customers who spend \$1000 a year
- ◆ age 40-50, employed, good credit ratings
- ◆ Data discrimination (contrast)
 - ◆ frequent vs. infrequent customers: e.g., age, education, employed
 - ◆ dry vs. wet regions: e.g., precipitation, humidity, temperature



Frequent Patterns

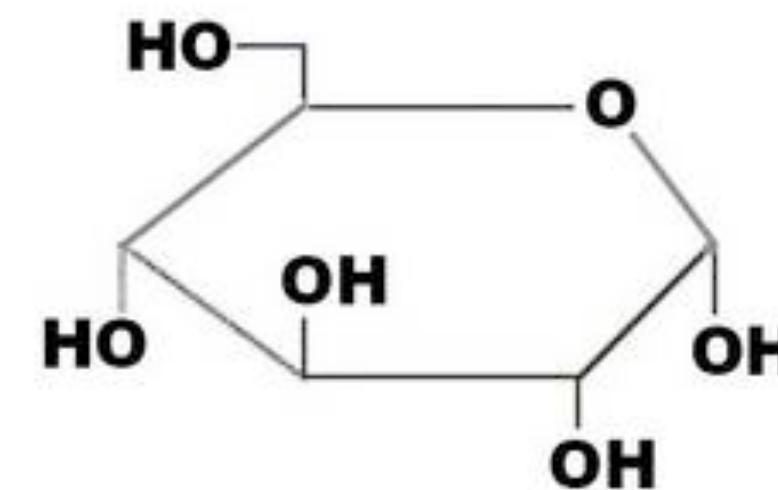
- ◆ Frequent itemsets

- ◆ e.g., (milk, bread, egg),
(beer, diaper)

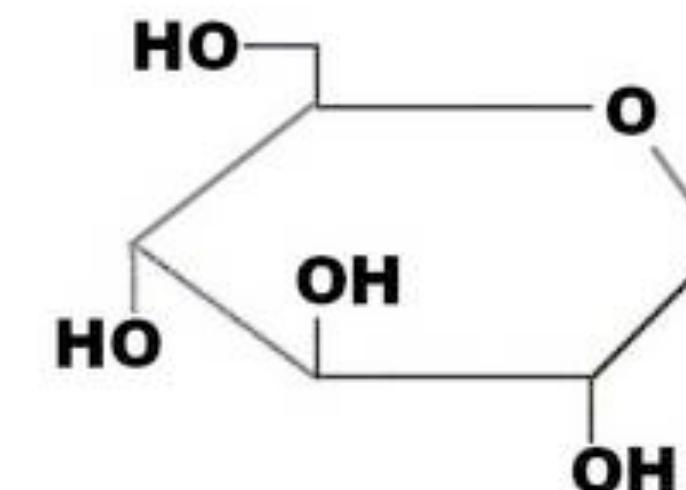
- ◆ Frequent sequences

- ◆ e.g., <printer, paper>,
<dinner, movie>

- ◆ Frequent structures



GLUCOSE



1,5-ANHYDROGLUCITOL

[http://www.endotext.org/diabetes/
diabetes12/figures/figure12.jpg](http://www.endotext.org/diabetes/diabetes12/figures/figure12.jpg)



Associations

- ◆ Association analysis
 - ◆ buys (X, milk) => buys (X, bread)
 - ◆ [support = 0.5%, confidence = 75%]
- ◆ Minimum support (or confidence) threshold

◆ Support

- ◆ chance of A and B appearing together

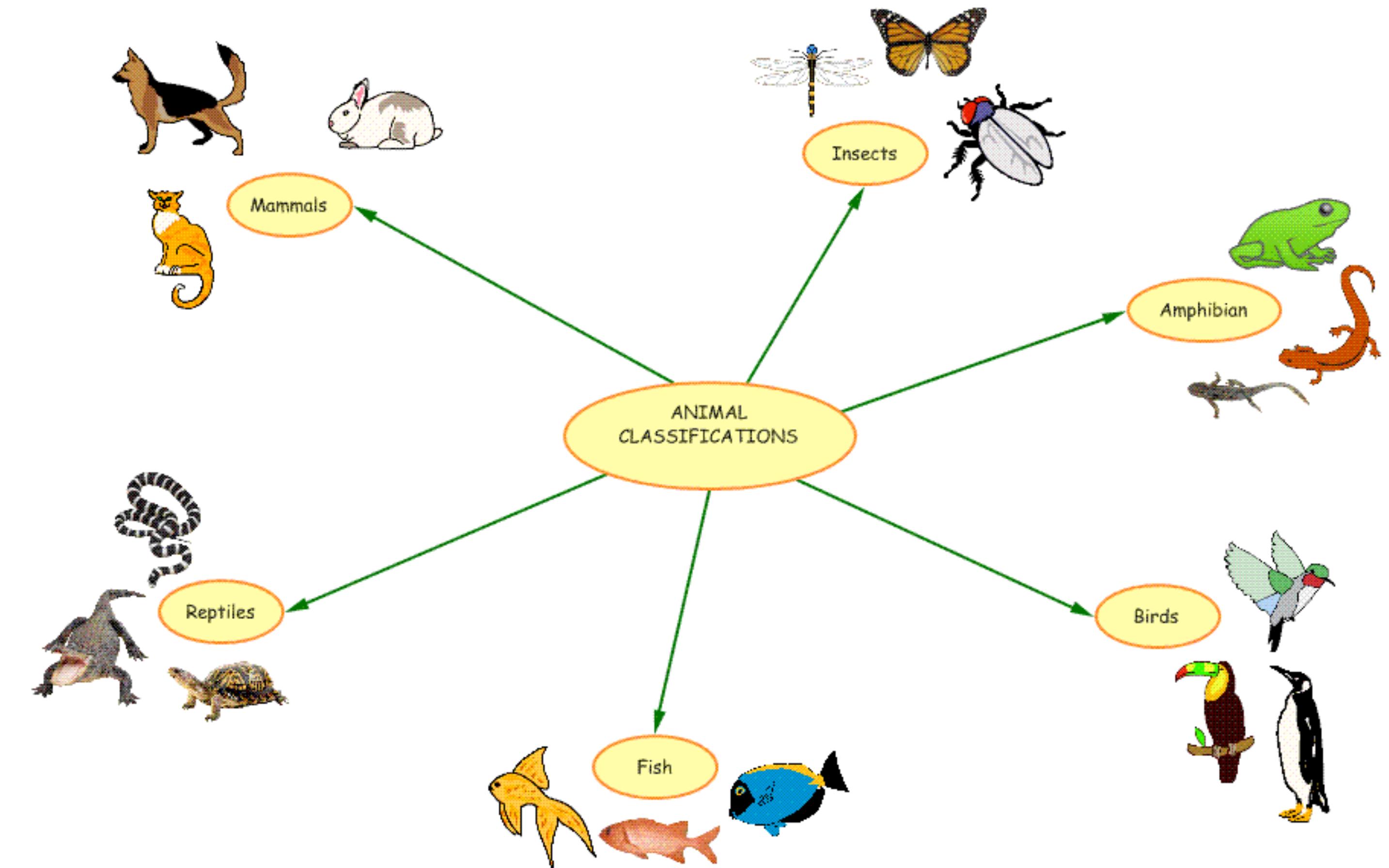
◆ Confidence

- ◆ if A appears, chance of B appears



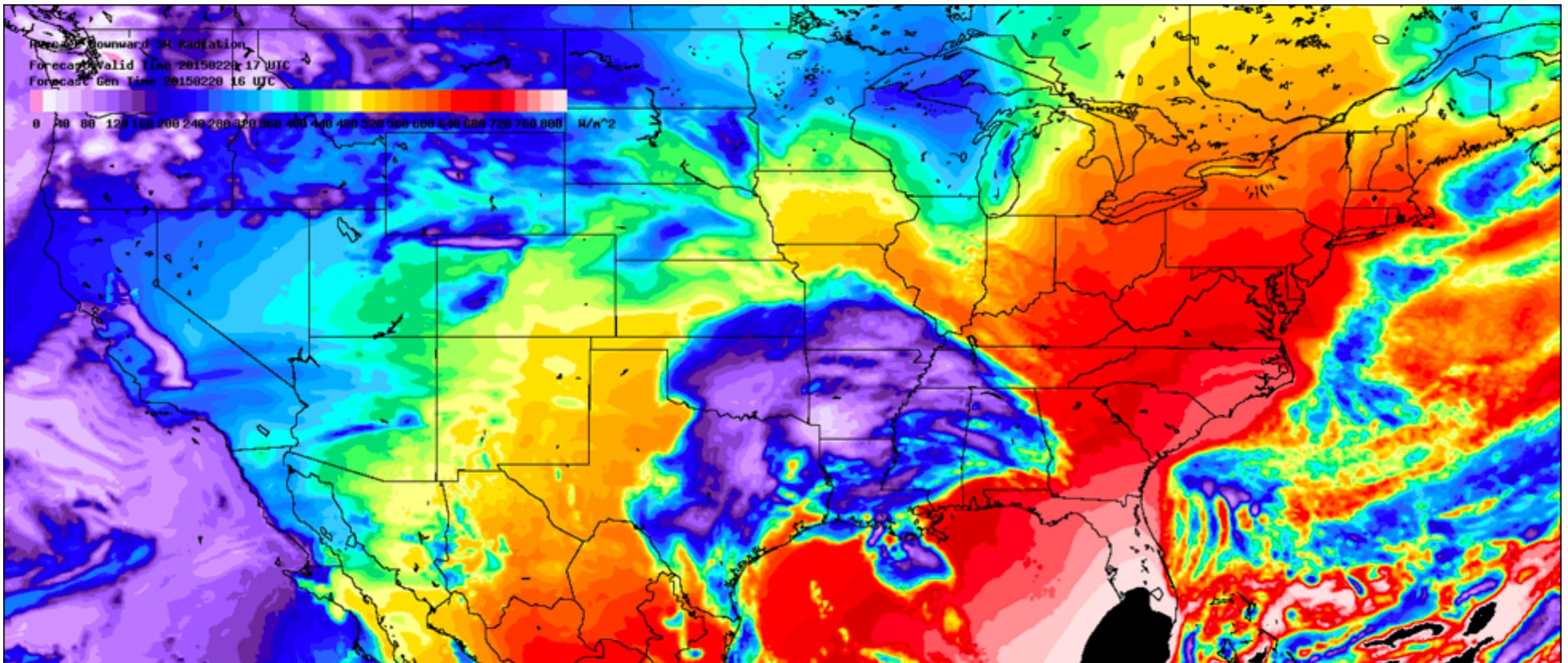
Classification

- ◆ Finding a model that describes and distinguishes data classes or concepts
- ◆ Training data
- ◆ IF-THEN rules, decision tree, neural network



Prediction

- ◆ Numerical prediction: continuous-valued instead of class labels
- ◆ E.g., weather, stock price, traffic



https://nar.ucar.edu/sites/default/files/labs/ral/2017/2.3.weather_prediction_3.png

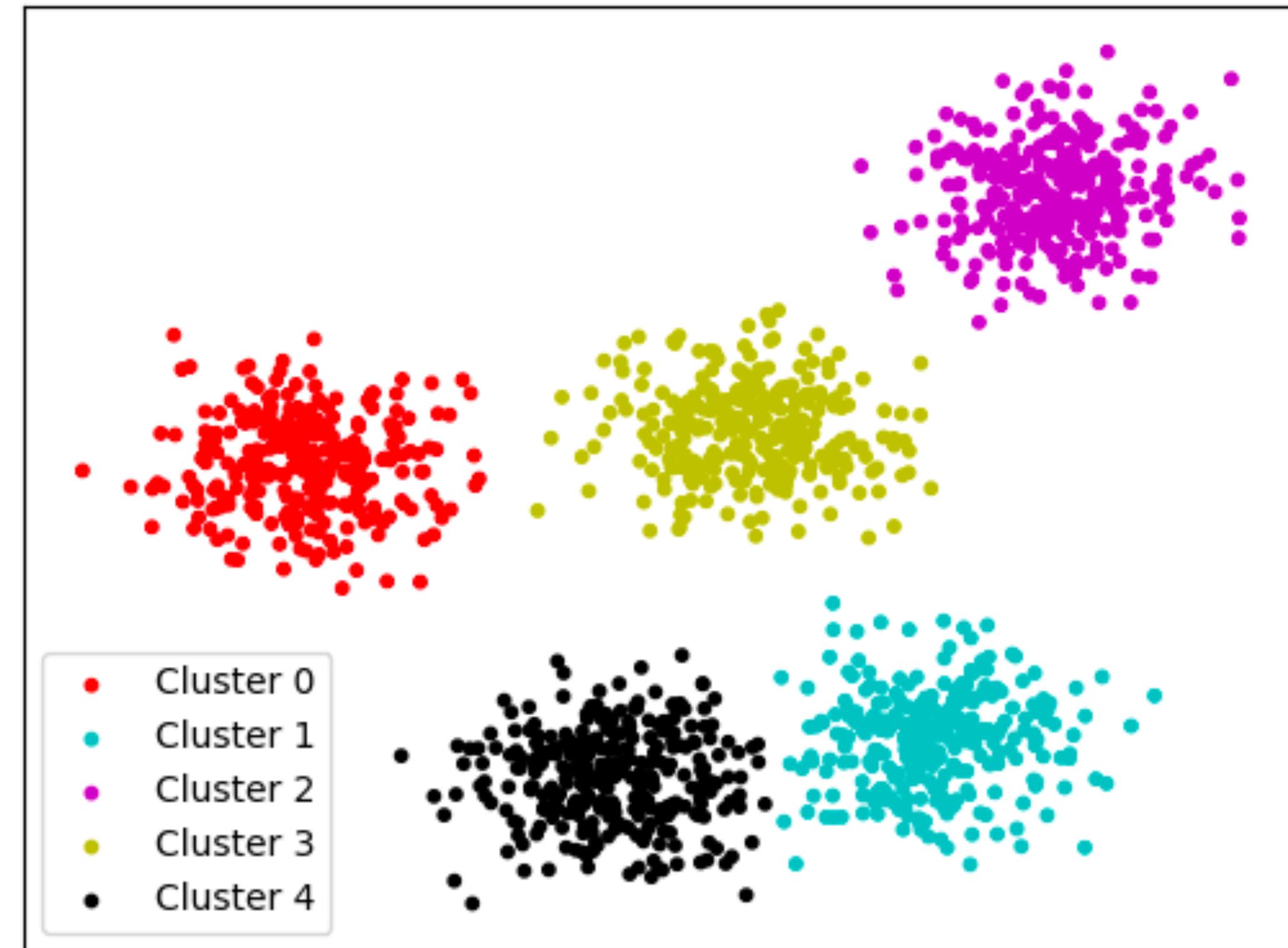


University of Colorado
Boulder

Fall 2020 Data Mining

Cluster Analysis

- ◆ Class labels unknown
- ◆ Intraclass similarity
 - ◆ maximize, closeness
- ◆ Interclass similarity
 - ◆ minimize, separation
- ◆ Hierarchical



https://miro.medium.com/max/1052/1*RsF6MMkuv0eECd_6m0-otw.png



Outlier Analysis

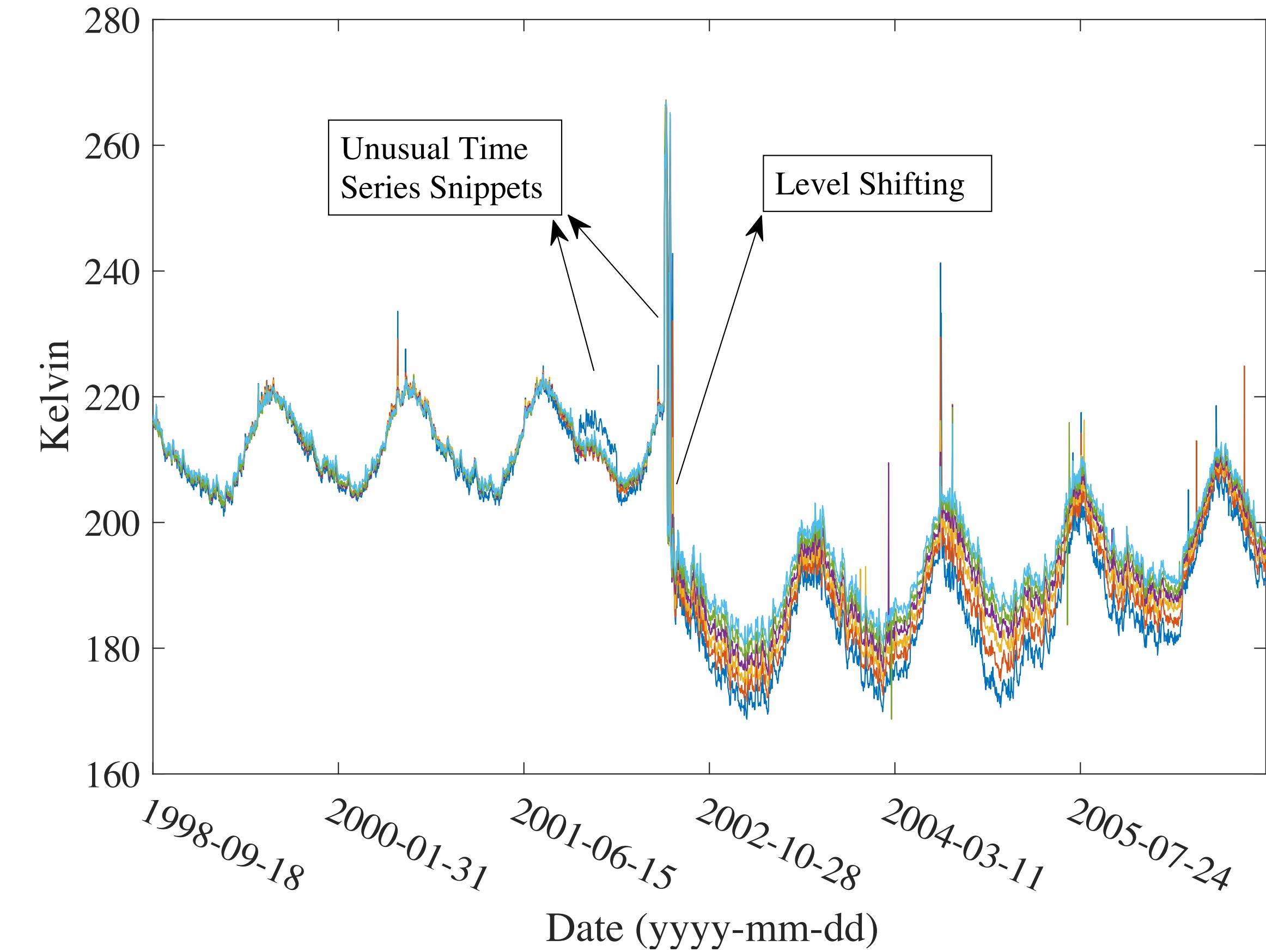
- ◆ Outliers

- ◆ do not comply with the general model

- ◆ Noise or exception

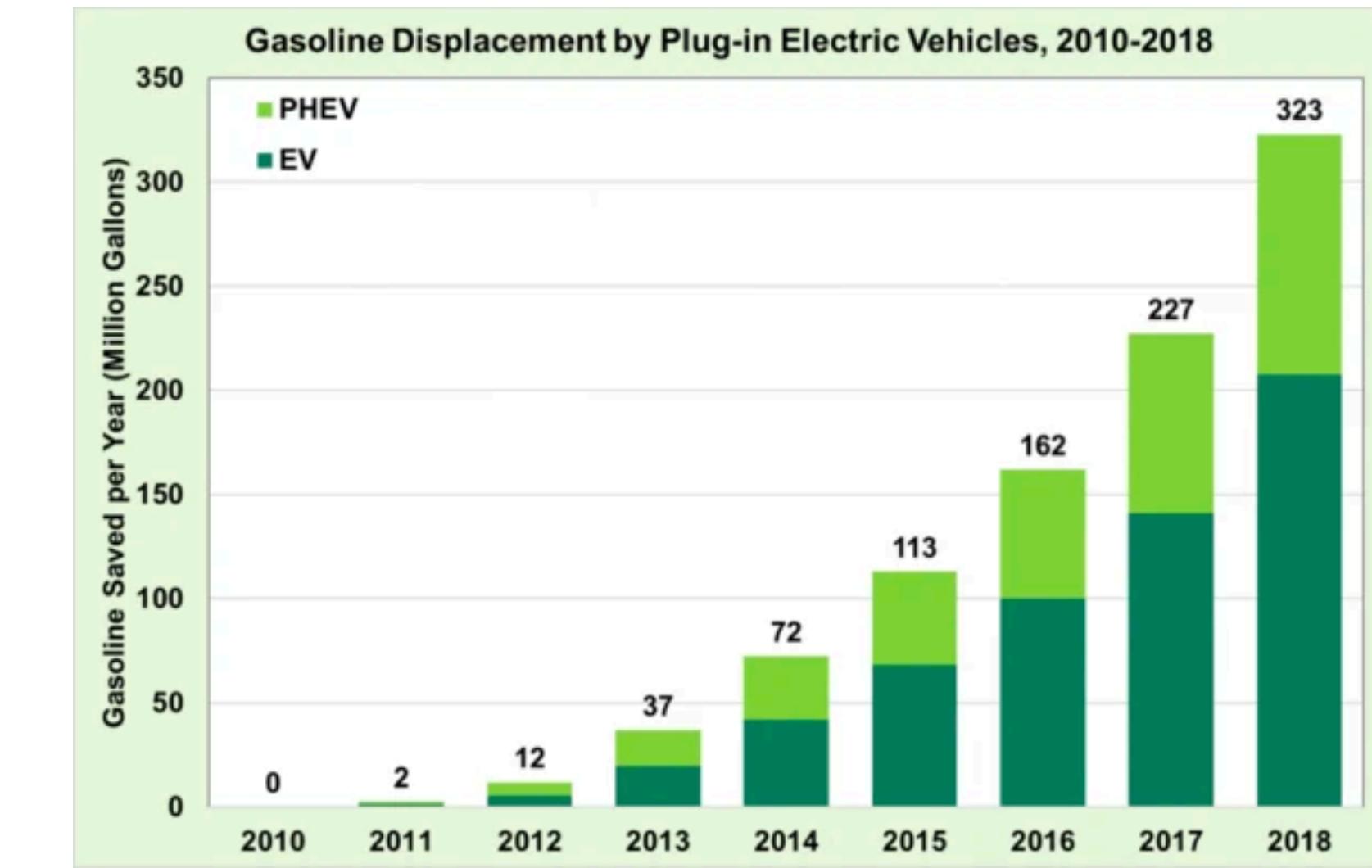
- ◆ Fraud detection, rare event analysis

- ◆ E.g. credit fraud analysis



Trend and Evolution Analysis

- ◆ Trends, deviations
- ◆ Sequential pattern mining
 - ◆ e.g., traffic congestion
- ◆ Periodicity analysis
- ◆ E.g., music, applications, ...



<https://www.osti.gov/biblio/1506474-assessment-light-duty-plug-electric-vehicles-united-states>



Market Analysis/Management

- ◆ Data sources: credit card transactions, club cards, customer calls, ...
- ◆ What types of customers buy what products
- ◆ What factors attract new customers
- ◆ Target marketing, product recommendation, discount
- ◆ Fraud detection



Are the Patterns Interesting?

- ◆ Interesting pattern
 - ◆ valid on new/test data with some certainty
 - ◆ novel
 - ◆ potentially useful
 - ◆ ultimately understandable by humans
- ◆ Objective measures
 - ◆ e.g., support, confidence, false positive/negative, accuracy
- ◆ Subjective measures
 - ◆ Completeness, exclusiveness



Major Issues in Data Mining (I)

◆ Mining technology

- ◆ mining different knowledge from diverse data (maybe noisy or incomplete)
- ◆ pattern evaluation: interestingness
- ◆ efficiency, effectiveness, scalability
- ◆ parallel, distributed, incremental mining
- ◆ incorporation of background knowledge
- ◆ integration of discovered knowledge with existing knowledge



Major Issues in Data Mining (2)

- ◆ User interaction
 - ◆ data mining query languages, ad-hoc mining
 - ◆ expression and visualization of results
 - ◆ interactive mining at multiple granularities
- ◆ Applications and social impacts
 - ◆ domain-specific data mining
 - ◆ applications of data mining results
 - ◆ protect data security, integrity, privacy



Data Science Ethics

- ◆ Data ownership
- ◆ Privacy, anonymity
- ◆ Data and model validity
- ◆ Data and model bias (algorithmic fairness)
- ◆ Interpretation, application, societal consequence



Data Mining Resources

- ◆ ACM SIGKDD: <https://www.kdd.org/>
- ◆ Conferences
 - ◆ KDD: tutorials, research, applied data science, KDD Cup, sponsors
 - ◆ SDM, ICDM, WSDM, CIKM, ICDE, TheWebConference (formerly WWW), SIGIR, ICML, CVPR, NeurIPS (formerly NIPS), SIGMOD, VLDB, ...
- ◆ Journals
 - ◆ TKDE, TKDD, DMKD, TPAMI, ...



Summary

- ◆ Chapter I: Introduction to data mining
- ◆ Data mining: discover interesting patterns in huge amounts of data
- ◆ Data mining pipeline
- ◆ Different views: data, knowledge, method, application
- ◆ Measure of pattern interestingness
- ◆ Major issues in data mining

