# Lecture 1: 24 January, 2024

Madhavan Mukund

`https://www.cmi.ac.in/~madhavan`

Data Mining and Machine Learning

January–May 2022

# What is this course about?

**Data Mining**

- Identify "hidden" patterns in data

- Also data collection, cleaning, uniformization, storage
  - Won't emphasize these aspects

# What is this course about?

## Data Mining

- Identify "hidden" patterns in data

- Also data collection, cleaning, uniformization, storage
  - Won't emphasize these aspects

## Machine Learning

- "Learn" mathematical models of processes from data

- Supervised learning — learn from experience

- Unsupervised learning — search for structure

# Supervised Learning

**Extrapolate from historical data**

- Predict board exam scores from model exams

- Should this loan application be granted?

- Do these symptoms indicate CoViD-19?

# Supervised Learning

Extrapolate from historical data

- Predict board exam scores from model exams

- Should this loan application be granted?

- Do these symptoms indicate CoViD-19?

"Manually" labelled historical data is available

- Past exam scores: model exams and board exam

- Customer profiles: age, income, ..., repayment/default status

- Patient health records, diagnosis

# Supervised Learning

## Extrapolate from historical data

- Predict board exam scores from model exams

- Should this loan application be granted?

- Do these symptoms indicate CoViD-19?

## "Manually" labelled historical data is available

- Past exam scores: model exams and board exam

- Customer profiles: age, income, ..., repayment/default status

- Patient health records, diagnosis

Historical data → model to predict outcome

# Supervised learning . . .

What are we trying to predict?

Numerical values

- Board exam scores

- House price (valuation for insurance)

- Net worth of a person (for loan eligibility)

# Supervised learning . . .

What are we trying to predict?

Numerical values

- Board exam scores

- House price (valuation for insurance)

- Net worth of a person (for loan eligibility)

Categories

- Email: is this message junk?

- Insurance claim: pay out, or check for fraud? } Binary

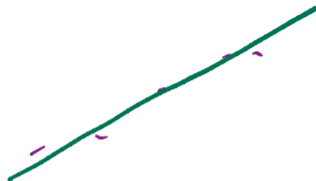- Credit card approval: reject, normal, premium — Multi category

# Supervised learning . . .

How do we predict?

- Build a mathematical model
  - Different types of models
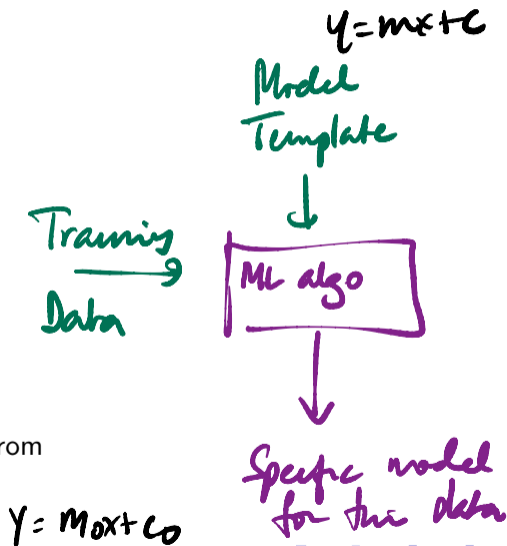  - Parameters to be tuned

How do we predict?

- Build a mathematical model
  - Different types of models
  - Parameters to be tuned

- Fit parameters based on input data
  - Different historical data produces different models
  - e.g., each user's junk mail filter fits their individual preferences

## How do we predict?

- Build a mathematical model
  - Different types of models
  - Parameters to be tuned

- Fit parameters based on input data
  - Different historical data produces different models
  - e.g., each user's junk mail filter fits their individual preferences

- Study different models, how they are built from historical data

$$y = mx + c$$

Model Template

Training Data → ML algo

Specific model for this data

$$y = m_0 x + c_0$$
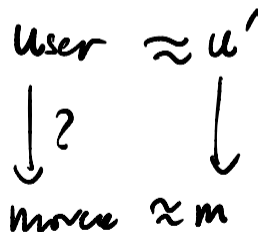
# Unsupervised learning

- Supervised learning builds models to reconstruct "known" patterns given by historical data

- Unsupervised learning tries to identify patterns without guidance

# Unsupervised learning

- Supervised learning builds models to reconstruct "known" patterns given by historical data

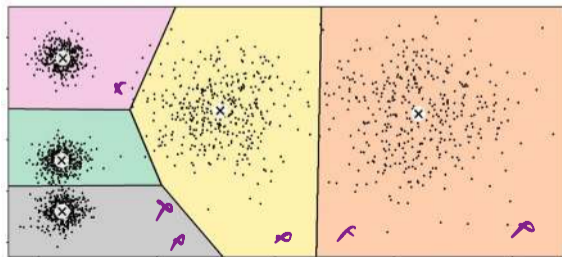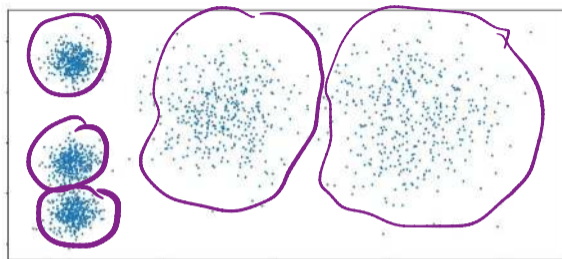- Unsupervised learning tries to identify patterns without guidance

## Customer segmentation

- Different types of newspaper readers

- Age vs product profile of retail shop customers
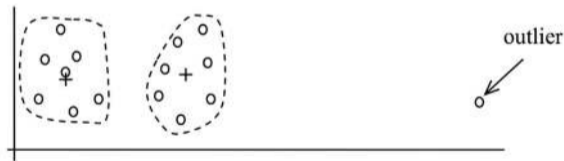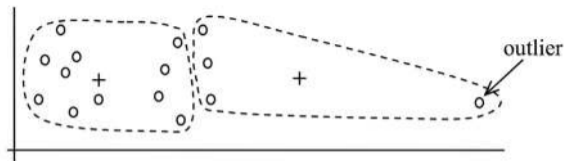
- Viewer recommendations on video platform

$$u \approx u'$$

$$m' \approx m$$

$$\text{User} \approx u'$$

$$\text{Movie} \approx m$$

# Clustering

- Organize data into "similar" groups — clusters

- Define a similarity measure, or distance function

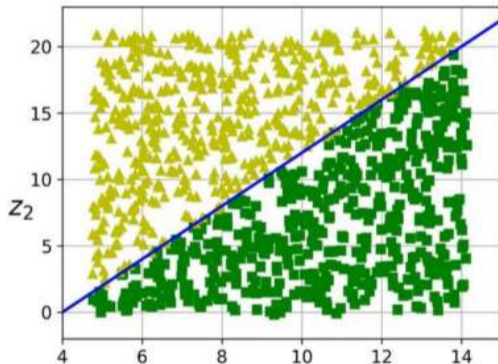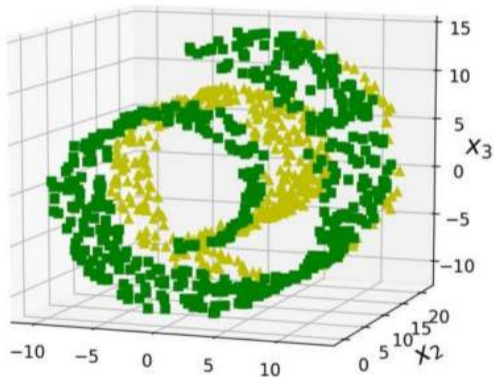- Clusters are groups of data items that are "close together"

# Outliers

- Outliers are anomalous values
  - Net worth of Bill Gates, Mukesh Ambani

- Outliers distort clustering and other analysis

- How can we identify outliers?

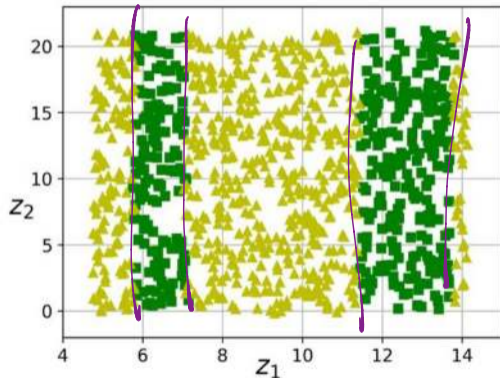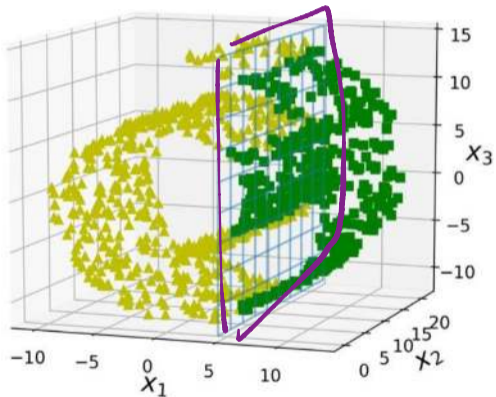# Preprocessing for supervised learning

Dimensionality reduction

# Preprocessing for supervised learning

Need not be a good idea — perils of working blind!

# Summary

Machine Learning

- Supervised learning
    - Build predictive models from historical data

- Unsupervised learning
    - Search for structure
    - Clustering, outlier detection, dimensionality reduction

# Summary

Machine Learning

- Supervised learning
  - Build predictive models from historical data

- Unsupervised learning
  - Search for structure
  - Clustering, outlier detection, dimensionality reduction

*If intelligence were a cake, unsupervised learning would be the cake, supervised learning would be the icing on the cake, . . .*

Yann Le Cun, ACM Turing Award 2018

# Market-Basket Analysis

- People who buy $X$ also tend to buy $Y$

- Rearrange products on display based on customer patterns

# Market-Basket Analysis

- People who buy *X* also tend to buy *Y*

- Rearrange products on display based on customer patterns
  - The diapers and beer legend
  - The true story, http://www.dssresources.com/newsletters/66.php

# Market-Basket Analysis

- People who buy $X$ also tend to buy $Y$

- Rearrange products on display based on customer patterns
    - The diapers and beer legend
    - The true story, `http://www.dssresources.com/newsletters/66.php`

- Applies in more abstract settings
    - Items are concepts, basket is a set of concepts in which a student does badly
        - Students with difficulties in concept $A$ also tend to do misunderstand concept $B$
    - Items are words, transactions are documents

- Set of items $I = \{i_1, i_2, \ldots, i_N\}$        *N large*

- A transaction is a set $t \subseteq I$ of items

- Set of transactions $T = \{t_1, t_2, \ldots, t_M\}$        *M large*

# Formal setting

- Set of items $I = \{i_1, i_2, \ldots, i_N\}$

- A transaction is a set $t \subseteq I$ of items

- Set of transactions $T = \{t_1, t_2, \ldots, t_M\}$

- Identify association rules $X \rightarrow Y$ — *Sets of items* *"itemsets"*
    - $X, Y \subseteq I, X \cap Y = \emptyset$
    - If $X \subseteq t_j$ then it is likely that $Y \subseteq t_j$

# Formal setting

- Set of items $I = \{i_1, i_2, \ldots, i_N\}$

- A transaction is a set $t \subseteq I$ of items

- Set of transactions $T = \{t_1, t_2, \ldots, t_M\}$

- Identify association rules $X \rightarrow Y$
    - $X, Y \subseteq I$, $X \cap Y = \emptyset$
    - If $X \subseteq t_j$ then it is likely that $Y \subseteq t_j$

- Two thresholds
    - How frequently does $X \subseteq t_j$ imply $Y \subseteq t_j$?
    - How significant is this pattern overall?

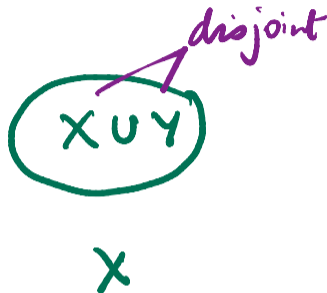- For $Z \subseteq I$, $Z.\text{count} = |\{t_j \mid Z \subseteq t_j\}|$

$$\leq M$$

$M$ total #

transaction

# Setting thresholds

- For $Z \subseteq I$, $Z.\text{count} = |\{t_j \mid Z \subseteq t_j\}|$

- How frequently does $X \subseteq t_j$ imply $Y \subseteq t_j$?

  - Fix a confidence level $\chi$

  - Want $\dfrac{(X \cup Y).\text{count}}{X.\text{count}} \geq \chi$

$\chi$

$\leq 1$

disjoint

$X \cup Y$

$\chi$

# Setting thresholds

- For $Z \subseteq I$, $Z.\text{count} = |\{t_j \mid Z \subseteq t_j\}|$

- How frequently does $X \subseteq t_j$ imply $Y \subseteq t_j$?
  - Fix a confidence level $\chi$
  - Want $\dfrac{(X \cup Y).count}{X.count} \geq \chi$

- How significant is this pattern overall?
  - Fix a support level $\sigma$
  - Want $\dfrac{(X \cup Y).count}{M} \geq \sigma$

$$X \to Y$$

- For $Z \subseteq I$, $Z.\text{count} = |\{t_j \mid Z \subseteq t_j\}|$

- How frequently does $X \subseteq t_j$ imply $Y \subseteq t_j$?
  - Fix a confidence level $\chi$
  - Want $\dfrac{(X \cup Y).count}{X.count} \geq \chi$

- How significant is this pattern overall?
  - Fix a support level $\sigma$
  - Want $\dfrac{(X \cup Y).count}{M} \geq \sigma$

- Given sets of items $I$ and transactions $T$, with confidence $\chi$ and support $\sigma$, find all valid association rules $X \to Y$

Fixed set of
Valid $X \to Y$

# Frequent itemsets

- $X \rightarrow Y$ is interesting only if $(X \cup Y).\text{count} \geq \sigma \cdot M$

- First identify all frequent itemsets
    - $Z \subseteq I$ such that $Z.\text{count} \geq \sigma \cdot M$

# Frequent itemsets

- $X \to Y$ is interesting only if $(X \cup Y).\text{count} \geq \sigma \cdot M$

- First identify all frequent itemsets
  - $Z \subseteq I$ such that $Z.\text{count} \geq \sigma \cdot M$

- Naïve strategy: maintain a counter for each $Z$

  - For each $t_j \in T$
    - For each $Z \subseteq t_j$
      - Increment the counter for $Z$

  - After scanning all transactions, keep $Z$ with $Z.\text{count} \geq \sigma \cdot M$

$|t| \leq \boxed{m} \quad 10 \sim 20$

$2^m$ subsets

Decompose $Z$ as

$X, Y$

$X \to Y$ is valid

(above confidence)

$M \cdot 2^m$

$Z \subseteq I \qquad |I| = N$

$2^N \sim 10^6$ potential subsets

| 1 | ℍℍ |
| 2 | 1 |
| 3 | ℍℍ |
| 4 | 1 |

# Frequent itemsets

- $X \rightarrow Y$ is interesting only if $(X \cup Y).count \geq \sigma \cdot M$

- First identify all frequent itemsets
  - $Z \subseteq I$ such that $Z.count \geq \sigma \cdot M$

- Naïve strategy: maintain a counter for each $Z$
  - For each $t_j \in T$
    For each $Z \subseteq t_j$
      Increment the counter for $Z$
  - After scanning all transactions, keep $Z$ with $Z.count \geq \sigma \cdot M$

- Need to maintain $2^{|I|}$ counters
  - Infeasible amount of memory
  - Can we do better?

$$X \quad\quad Y$$

Rolls Royce    leather seats

$N$ — cars + accessories on sale

$=$

$|I|$

$M = \{ t_1, t_2 \cdots \quad t_M \}$

For each $t \in T$

~~For each~~ $z \subseteq T$    $z \subseteq t$

Check if $z.count$
should be
incremented

For each $z \subseteq I$

For each $t \in T$

Does $t$
contain $z$ ?