

# Elements of Machine Learning and Data Science

Part I: Data Science — Exam Notes (Living Document)

Emir Pisirici

January 31, 2026

Exam likelihood: High (overall Data Science part)

This document is structured to match the lecture topics exactly and is designed for adding **exam-style notes**, **common traps**, and **visual summaries**.

## Contents

<b>1</b>	<b>Introduction to Data Science</b>	<b>3</b>
1.1	Introduction . . . . .	3
1.2	Tabular Data . . . . .	3
1.3	Data Science Process . . . . .	3
1.3.1	ETL vs ELT (Definitions + Differences) . . . . .	3
1.3.2	CRISP-DM . . . . .	3
1.3.3	PDCA . . . . .	4
1.3.4	DMAIC . . . . .	4
1.4	Data Types . . . . .	4
1.5	Descriptive Statistics . . . . .	4
1.6	Basic Visualizations . . . . .	5
1.7	Feature Transformations . . . . .	6
1.8	“How to lie with statistics” . . . . .	6
<b>2</b>	<b>Decision Trees</b>	<b>7</b>
2.1	Introduction to Decision Trees . . . . .	7
2.2	Entropy and Information Gain . . . . .	7
2.3	ID3 Algorithm . . . . .	9
2.4	Quantifying Information Gain . . . . .	12
2.5	Pruning . . . . .	13
2.6	Continuous Data (Threshold splits) . . . . .	14
2.7	Ensembles (Bagging/Random Forest/Boosting) . . . . .	17
<b>3</b>	<b>Clustering</b>	<b>18</b>
3.1	Introduction to Unsupervised Learning . . . . .	18
3.2	Introduction to Clustering . . . . .	18
3.3	Similarity and Dissimilarity . . . . .	19
3.4	K-means and K-medoids . . . . .	20
3.5	Agglomerative Clustering . . . . .	21
3.6	DBSCAN . . . . .	23
3.7	Closing . . . . .	25

<b>4</b>	<b>Frequent Itemsets</b>	<b>26</b>
4.1	Introduction . . . . .	26
4.2	Properties of Frequent Itemsets . . . . .	27
4.3	Apriori Algorithm . . . . .	28
4.4	FP-Growth Algorithm . . . . .	30
<b>5</b>	<b>Association Rules</b>	<b>32</b>
5.1	Introduction . . . . .	32
5.2	Generating Association Rules . . . . .	33
5.3	Evaluation (support, confidence, lift, conviction) . . . . .	34
5.4	Applications . . . . .	37
5.5	Simpson's Paradox . . . . .	37
<b>6</b>	<b>Time Series</b>	<b>39</b>
6.1	Temporal Data . . . . .	39
6.2	Introduction to Time Series . . . . .	39
6.3	Analysis . . . . .	41
6.4	Forecasting . . . . .	42

# 1 Introduction to Data Science

## 1.1 Introduction

## 1.2 Tabular Data

## 1.3 Data Science Process

Exam likelihood: High

Framework questions are easy to grade and strongly test “big picture” understanding.

Examiner favorite (what they love to ask)

Typical asks: **ETL vs ELT**, **CRISP-DM phases**, and mapping a scenario to the correct phase. Also: where data leakage/bias lives (data understanding + evaluation).

### 1.3.1 ETL vs ELT (Definitions + Differences)

Cheat sheet / must-memorize

**ETL:** Extract → Transform → Load (transform before target).

**ELT:** Extract → Load → Transform (transform inside target platform).

**Key contrast:** where transformations happen; governance vs flexibility; raw history availability.

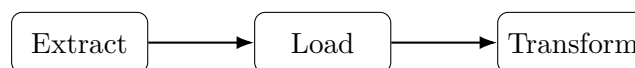
Common pitfall

People confuse “ELT = no cleaning”. Wrong. It means cleaning happens *after loading*, often in warehouse/lakehouse layers (staging → curated).

Visual



**ETL**



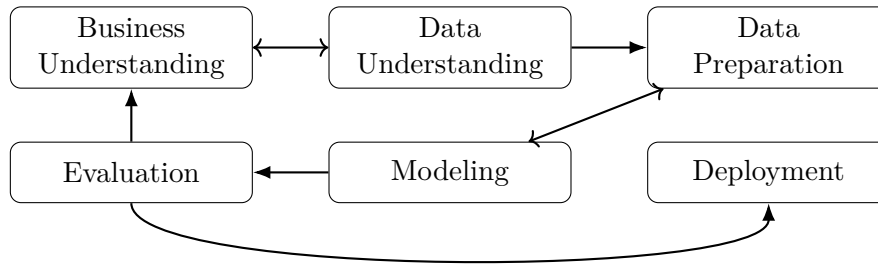
**ELT**

### 1.3.2 CRISP-DM

Cheat sheet / must-memorize

**CRISP-DM:** Business Understanding → Data Understanding → Data Preparation → Modeling → Evaluation → Deployment (iterative loops).

## Visual



### 1.3.3 PDCA

Cheat sheet / must-memorize

**PDCA:** Plan → Do → Check → Act (continuous improvement loop).

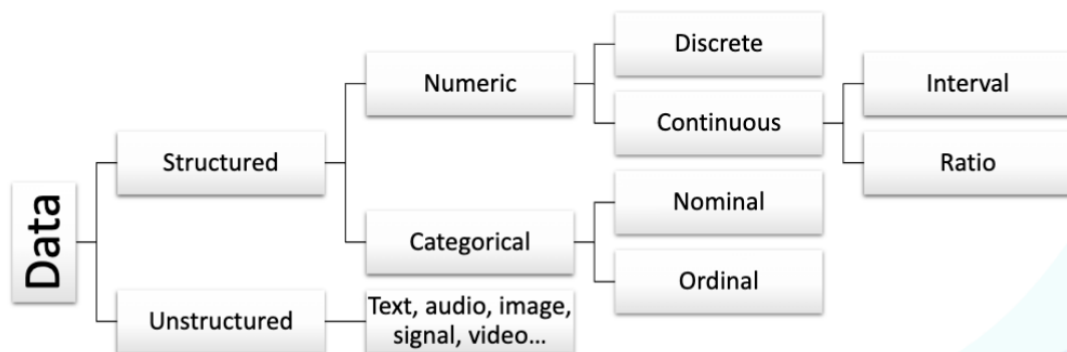
### 1.3.4 DMAIC

Cheat sheet / must-memorize

**DMAIC:** Define → Measure → Analyze → Improve → Control. Often used for process/quality improvement + monitoring and part of the Six Sigma methodology.

## 1.4 Data Types

## Visual



Advantages: clarifies variable handling and measurement choices.

Limitations: real data can mix types or sit between categories.

## 1.5 Descriptive Statistics

Exam likelihood: High

Frequent: compute variance/STD/covariance/correlation by hand; read a correlation matrix. Answer: apply sample formulas, compute values, and interpret sign/magnitude; matrix is symmetric with 1s on the diagonal.

Examiner favorite (what they love to ask)

Explain why covariance depends on units, and why correlation is normalized in  $[-1, 1]$ . Answer: covariance scales with units; correlation divides by SDs so it is unitless and bounded.

Why (motivation): Quantify spread and association between variables.  
 What (definition): Variance/STD measure spread; covariance/correlation measure linear association.  
 How (procedure/usage): Compute formulas, then interpret sign/magnitude and check the correlation matrix.

#### Cheat sheet / must-memorize

- **Variance (sample):**  $s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$
- **Std dev:**  $s = \sqrt{s^2}$
- **Covariance (sample):**  $\text{cov}(X, Y) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$
- **Correlation:**  $r = \frac{\text{cov}(X, Y)}{s_X s_Y}$  (unitless, -1 to 1)
- **Correlation matrix:** table of pairwise correlations; symmetric with 1s on the diagonal.

#### Common pitfall

Correlation  $\neq$  causation; a strong correlation can be driven by a confounder or Simpson's paradox.

Advantages: quick, compact summaries of spread and association.  
 Limitations: can hide distribution shape and outliers.

#### Visual

1	$r_{12}$	$r_{13}$
$r_{21}$	1	$r_{23}$
$r_{31}$	$r_{32}$	1

Correlation matrix

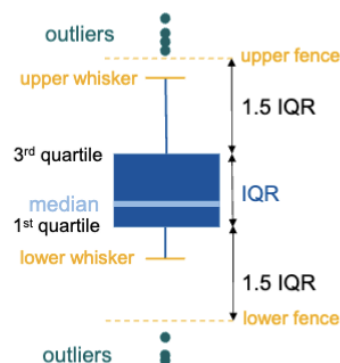
**Key takeaways:** Know formulas + interpretations; correlation matrix is symmetric with 1s on the diagonal.

## 1.6 Basic Visualizations

#### Visual

### Box Plot

- **Median** value (middle), depicted by bar
- **IQR** – Interquartile Range (covers 50% of middle instances), depicted by box
- **Upper fence** – 3<sup>rd</sup> quartile + 1.5 IQR  
**Upper whisker** – maximal value below upper fence
- **Lower fence** – 1<sup>st</sup> quartile - 1.5 IQR  
**Lower whisker** – minimal value above lower fence
- **Outliers** – drawn separately



Advantages: makes distribution and outliers visible at a glance.  
 Limitations: can hide multimodality or sample size differences.

## 1.7 Feature Transformations

Exam likelihood: High

Typical: pick the right transform (scale, log, encode) and explain why. Answer: choose encoding by category type and scaling/log for skew; justify the choice.

Examiner favorite (what they love to ask)

Identify data leakage in preprocessing; name the correct order for train/test transformations. Answer: fit transforms on training data only, then apply to validation/test.

Why (motivation): Turn raw categorical/continuous variables into model-ready features.

What (definition): Encoding or discretizing features without changing the target meaning.

How (procedure/usage): Choose encoding by category type; choose binning by distribution.

Cheat sheet / must-memorize

- **One-hot encoding:** create a 0/1 column per category (nominal).
- **Binary encoding:** represent categories as binary digits (compact one-hot).
- **Ordinal encoding:** map ordered categories to ranks (only if order is real).
- **Binning:** convert continuous to categories.
- **Equal-width binning:** fixed interval sizes across the range.
- **Equal-frequency binning:** same number of samples per bin.

Common pitfall

Fitting transforms on the full dataset (leakage). Always fit on training data, then apply to validation/test.

Advantages: improves model performance and stability.

Limitations: can reduce interpretability or introduce leakage if misapplied.

**Key takeaways:** Use one-hot for nominal, ordinal for ordered labels, and binning for simplification.

## 1.8 “How to lie with statistics”

## 2 Decision Trees

### 2.1 Introduction to Decision Trees

Exam likelihood: High

Intro questions often ask you to explain how trees split data and what leaves represent. Answer: describe root/internal/leaf roles and how splits partition the feature space.

Examiner favorite (what they love to ask)

Draw a small tree from a toy dataset or explain interpretability vs overfitting. Answer: build simple if-then splits and note that deeper trees overfit without pruning.

Why (motivation): Learn a function from labeled training instances to make predictions.

What (definition): A tree partitions the feature space by sequential if-then splits; leaves output a class or value.

How (procedure/usage): Choose splits to improve class purity or reduce prediction error.

Cheat sheet / must-memorize

- **Goal:** learn a function  $f(X)$  from labeled data to predict labels/values.
- **Tree structure:** root node, internal (non-leaf) nodes, leaf nodes.
- **Split rule:** one feature + threshold; paths are if-then rules.
- **Leaf meaning:** prediction (class/value) for that region of the space.

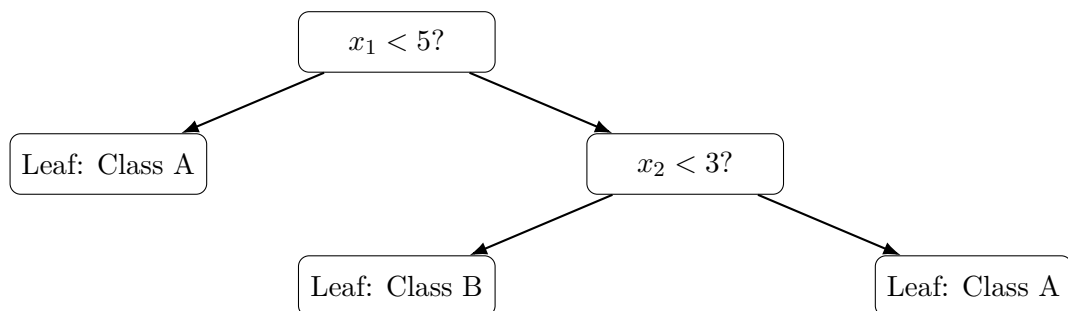
Common pitfall

Overly deep trees memorize training data; control with max depth, min samples, or pruning.

Advantages: interpretable rules and fast inference.

Limitations: high variance and prone to overfitting without pruning.

Visual



**Key takeaways:** Trees learn  $f$  from labeled data using splits; nodes/leaf roles are core.

### 2.2 Entropy and Information Gain

Exam likelihood: Very High

Almost guaranteed: compute entropy and information gain for candidate splits. Answer: compute parent entropy, weighted child entropies, then IG.

### Examiner favorite (what they love to ask)

Given a small labeled dataset, compare splits and pick the one with highest information gain.  
Answer: calculate IG for each split and select the largest.

Why (motivation): Choose splits that make child nodes as pure as possible.

What (definition): Entropy measures impurity; information gain is impurity reduction.

How (procedure/usage): Compute parent entropy, child entropies, then  $IG = \text{parent} - \text{weighted children}$ .

### Cheat sheet / must-memorize

- **Entropy:**  $H(S) = -\sum_c p_c \log_2 p_c$  (define  $0 \log 0 = 0$ ).
- **Weighted child entropy:**  $\sum_k \frac{|S_k|}{|S|} H(S_k)$ .
- **Information gain:**  $IG(S, \text{split}) = H(S) - \sum_k \frac{|S_k|}{|S|} H(S_k)$ .
- **Goal:** choose split with highest IG (most impurity reduction).

### Common pitfall

Forgetting to weight child entropies by subset size; using raw entropy sums gives wrong IG.

Advantages: principled split criterion that increases purity.

Limitations: IG is biased toward many-valued attributes.





### Examiner favorite (what they love to ask)

Explain stopping conditions and why ID3 prefers high information gain. Answer: stop if pure/no attributes; IG maximizes impurity reduction.

Why (motivation): Build a decision tree that best separates labeled data.

What (definition): ID3 is a greedy, top-down tree induction algorithm using information gain.

How (procedure/usage): Compute IG for each attribute, split on the best, and recurse.

### Cheat sheet / must-memorize

- **Input:** labeled dataset  $S$  with categorical attributes (ID3 original).
- **Step 1:** if all labels same  $\rightarrow$  make a leaf.
- **Step 2:** if no attributes left  $\rightarrow$  leaf with majority class.
- **Step 3:** choose attribute with highest IG.
- **Step 4:** split  $S$  by attribute values and recurse.
- **Output:** a decision tree; leaves store class label.

### Common pitfall

ID3 favors attributes with many values; without corrections (e.g., gain ratio) it can overfit.

Advantages: simple, fast, and easy to explain.

Limitations: greedy (not optimal), biased to many-valued attributes.

## ID3 Algorithm

### ID3 algorithm:

1. **if** all the instances in  $X$  have the same classification
  - (a) **return** a decision tree with one leaf node with consensus value as a label
2. **else if** there are no features left
  - (a) **return** a decision tree with one leaf node with majority value as a label
3. **else if** the dataset is empty
  - (a) **return** a decision tree with one leaf node with majority parent value as a label
4. **else**
  - (a) pick a feature that maximizes information gain
  - (b) once a feature is picked along a path from the root, it cannot be used again
  - (c) create subproblems based on the selected feature

three  
stopping  
criteria

recursively  
constructing  
the tree

### ID3 Algorithm

#### Example

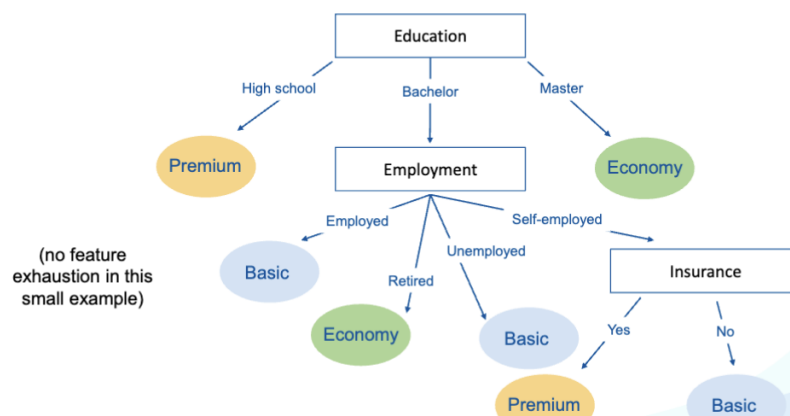
$$H(\text{Customer}) = 1.5567$$

ID	Insurance	Education	Employment	Customer
1	Yes	Bachelor	Employed	Basic
2	Yes	High school	Unemployed	Premium
3	Yes	Bachelor	Self-employed	Premium
4	No	Bachelor	Self-employed	Basic
5	No	Master	Employed	Economy
6	Yes	Bachelor	Retired	Economy
7	Yes	High school	Employed	Premium

Split by feature	Possible Values	Instances	Entropy	Overall Entropy	Information Gain
Insurance	No	4, 5	1	1.265	$1.5567 - 1.265 = 0.2917$
	Yes	1, 2, 3, 6, 7	1.3710		
Education	High school	2, 7	0	0.8571	$1.5567 - 0.8571 = 0.6996$
	Master	5	0		
	Bachelor	1, 3, 4, 6	1.5		
Employment	Employed	1, 5, 7	1.5850	0.9650	$1.5567 - 0.9650 = 0.5917$
	Unemployed	2	0		
	Self-employed	3, 4	1		
	Retired	6	0		

### ID3 Algorithm

#### Example



**Key takeaways:** ID3 is greedy; compute IG, split, recurse, stop with pure/majority leaves.

## 2.4 Quantifying Information Gain

### Exam likelihood: Very High

Often compute IG for a specific split and compare candidate attributes. Answer: show parent entropy, weighted children, and IG (optionally gain ratio).

### Examiner favorite (what they love to ask)

Show all intermediate steps: parent entropy, each child entropy, weighted sum, IG. Answer: compute each step explicitly and report the final IG (and GR if asked).

Why (motivation): Convert “best split” into a concrete, comparable number.

What (definition):  $IG = \text{parent entropy} - \text{weighted child entropies}$ ; Split Info measures how evenly the split divides data; Gain Ratio normalizes IG.

How (procedure/usage): Compute IG, then divide by split info to get gain ratio.

### Cheat sheet / must-memorize

- **Step 1:** compute parent entropy  $H(S)$ .
- **Step 2:** split by attribute values.
- **Step 3:** compute each child entropy  $H(S_k)$ .
- **Step 4:** compute weighted sum  $\sum_k \frac{|S_k|}{|S|} H(S_k)$ .
- **Step 5:**  $IG = H(S) - \sum_k \frac{|S_k|}{|S|} H(S_k)$ .
- **Split info (lecture:  $H(d)$ ):** entropy of split proportions (how evenly data is partitioned).  
 $H(d) = SI = - \sum_k \frac{|S_k|}{|S|} \log_2 \frac{|S_k|}{|S|}$ .
- **Gain ratio:**  $GR = \frac{IG}{H(d)}$  (same as  $IG/SI$ ; penalizes many-valued attributes).

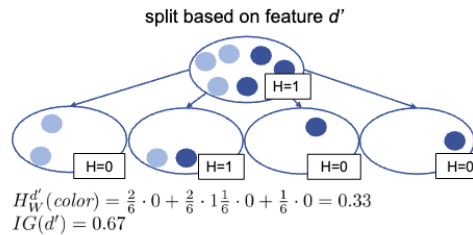
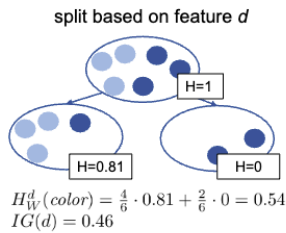
### Common pitfall

Information gain is biased toward attributes with many values; use gain ratio to correct. Also: weight by subset size and keep log base consistent.

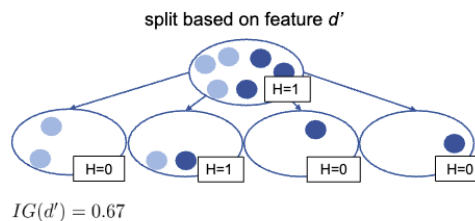
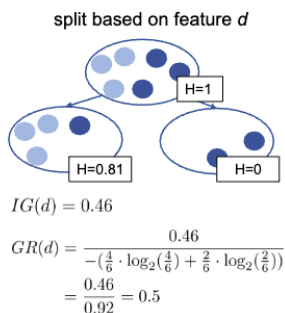
Advantages: makes split comparisons explicit and quantitative.

Limitations: IG can favor attributes with many values without normalization.

### Information Gain Ratio - Example



### Information Gain Ratio - Example



Feature  $d$  splits the 6 instances into one partition of size 4 and one partition of size 2

$$GR(d) = \frac{IG(d)}{H(d)} = \frac{H(t) - H_W^d(t)}{-\sum_{k=1}^K (P(d=k) \cdot \log_2(P(d=k)))}$$

**Key takeaways:** IG is a weighted impurity reduction; higher is better.

## 2.5 Pruning

Exam likelihood: High

Often: explain why pruning reduces overfitting and name pre- vs post-pruning. Answer: pruning removes low-value branches to reduce variance; name pre/post pruning.

Examiner favorite (what they love to ask)

Given a tree, identify which branches to prune using validation error or complexity. Answer: prune branches that do not improve validation performance or reduce cost-complexity.

Why (motivation): Reduce overfitting by simplifying a deep tree.

What (definition): Pruning removes splits/branches that do not improve generalization.

How (procedure/usage): Stop early (pre-pruning) or cut back after training (post-pruning).

Cheat sheet / must-memorize

- **Pre-pruning:** stop splitting early using rules like: max depth, min samples per node, min impurity decrease, min samples per leaf.
- **Post-pruning:** grow full tree, then prune using validation error or cost-complexity.
- **Goal:** simpler tree with similar or better validation performance.

### Common pitfall

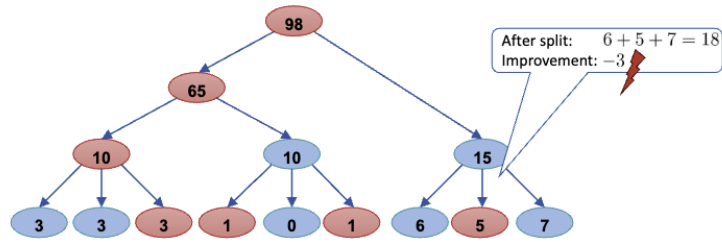
Pruning too aggressively can underfit; always tune on validation data, not test data.

Advantages: improves generalization by reducing variance.

Limitations: too much pruning increases bias.

### Visual

#### Post-pruning



- Decision tree learned on a **training set**
- Numbers indicate misclassifications based on a **validation set**

**Key takeaways:** Pruning trades depth for generalization; use validation to choose.

## 2.6 Continuous Data (Threshold splits)

### Exam likelihood: High

Common: show how a continuous feature is split using a threshold and compute the best split score. Answer: sort values, test midpoints, and choose the split with highest IG (or variance reduction).

### Examiner favorite (what they love to ask)

Given sorted values, test candidate thresholds and pick the one with maximum information gain. Answer: evaluate IG at valid midpoints and choose the maximum.

Why (motivation): Many real-world features are continuous; trees need a simple yes/no split.

What (definition): A threshold split picks a value  $t$  so that  $x_j \leq t$  goes left and  $x_j > t$  goes right.

How (procedure/usage): Sort values, test candidate thresholds (midpoints between distinct neighbors), compute split score, choose best, recurse.

Variance in a node/leaf (regression trees): For a node with targets  $y_1, \dots, y_n$  and mean  $\bar{y}$ , the node variance (impurity) is the average squared deviation from the mean; a good split minimizes the weighted variance of the child nodes (equivalently, maximizes variance reduction).

**Mini example (regression):** Data  $(x, y) = (1, 2), (2, 2), (3, 3), (4, 10), (5, 11), (6, 12)$ .

Parent variance  $\approx 19.22$ . Try two thresholds:

- $t = 3.5$ : left  $y = \{2, 2, 3\}$  var  $\approx 0.22$ , right  $y = \{10, 11, 12\}$  var  $\approx 0.67$ . Weighted variance  $= (3/6) \cdot 0.22 + (3/6) \cdot 0.67 \approx 0.45$ .
- $t = 1.5$ : left  $y = \{2\}$  var  $= 0$ , right  $y = \{2, 3, 10, 11, 12\}$  var  $\approx 17.84$ . Weighted variance  $\approx (5/6) \cdot 17.84 = 14.87$ .

So  $t = 3.5$  is preferred because it gives much lower weighted variance.

### Cheat sheet / must-memorize

- **Split rule:**  $x_j \leq t$  vs  $x_j > t$  (binary split).
- **Candidates:** midpoints between sorted unique values; often only where class label changes.
- **Score (classification):** maximize information gain (IG).
- **Variance in a node:**  $\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2$  (SSE/variance impurity).
- **Score (regression):** minimize weighted child variance / MSE (maximize variance reduction).
- **Efficiency:** sort once ( $O(n \log n)$ ), then sweep thresholds.

### Common pitfall

Using thresholds that create empty or nearly empty children (e.g., duplicates without midpoints); always use distinct values and check child sizes.

Advantages: enables trees to handle continuous features.

Limitations: many candidate splits; sensitive to noise/outliers.

## Continuous Descriptive Features - Example

ID	Insurance	Income	Employment	Customer
2	Yes	0	Unemployed	Premium
3	Yes	1000	Self-employed	Premium
4	No	2000	Self-employed	Basic
7	Yes	3000	Employed	Premium
1	Yes	3500	Employed	Basic
5	No	5000	Employed	Economy
6	Yes	5100	Retired	Economy

Four candidate thresholds

Thresholds: middle values of continuous feature in between changed target features

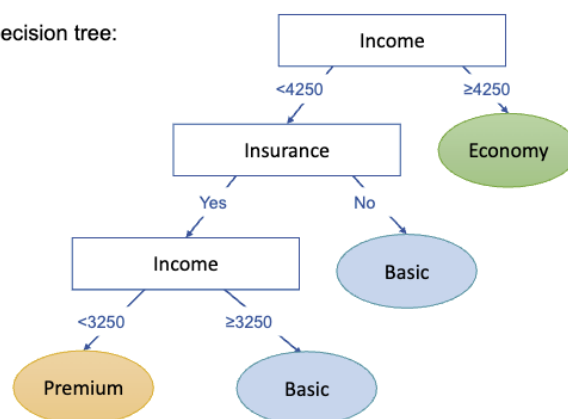
## Continuous Descriptive Features - Example

Threshold	Instances	Partition Entropy	Overall Entropy	Information Gain
$\geq 1500$	2, 3 1, 4, 5, 6, 7	0 1.5219	1.0871	0.1981
$\geq 2500$	2, 3, 4 1, 5, 6, 7	0.9183 1.5	1.2507	0.306
$\geq 3250$	2, 3, 4, 7 1, 5, 6	0.8113 0.9183	0.8572	0.6995
$\geq 4250$	1, 2, 3, 4, 7 5, 6	0.9710 0	0.6935	0.8631

Compute as usual

## Continuous Descriptive Features - Example

Resulting decision tree:



The same continuous feature can now be used multiple times!

**Key takeaways:** Continuous features are split by thresholds; evaluate candidate midpoints and pick the best impurity reduction.



## 2.7 Ensembles (Bagging/Random Forest/Boosting)

Exam likelihood: High

Short theory questions about why ensembles help and how they are built are common. Answer: averaging many trees reduces variance and stabilizes predictions.

Examiner favorite (what they love to ask)

Explain bias–variance intuition: averaging many trees reduces variance and improves stability. Answer: averaging cancels noise, lowering variance and improving generalization.

Why (motivation): Single trees are high-variance; ensembles stabilize predictions and improve accuracy.

What (definition): An ensemble combines many base models (often trees) into one predictor.

How (procedure/usage): Train multiple trees (e.g., on resampled data) and aggregate their outputs by vote or average.

Cheat sheet / must-memorize

- **Goal:** reduce variance and improve generalization.
- **Combine:** predictions are averaged (regression) or majority-voted (classification).
- **Intuition:** many weak/unstable trees  $\rightarrow$  one strong, stable predictor.

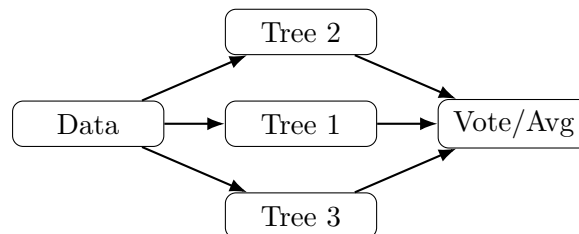
Common pitfall

Thinking more trees always fix bias; ensembles mainly reduce variance, not poor features or labels.

Advantages: strong accuracy and stability.

Limitations: reduced interpretability and higher compute.

Visual



**Key takeaways:** Ensembles combine many trees to reduce variance and stabilize predictions.

## 3 Clustering

### 3.1 Introduction to Unsupervised Learning

Exam likelihood: High

Often asked: define unsupervised learning and contrast it with supervised learning. Answer: unsupervised has no labels; it discovers structure like clusters.

Examiner favorite (what they love to ask)

Explain what “no labels” means and give a concrete task like clustering or dimensionality reduction. Answer: only  $X$  is given; tasks include clustering and dimensionality reduction.

Why (motivation): Labels are expensive or unavailable; we still want structure in the data.

What (definition): Unsupervised learning finds patterns, groups, or low-dimensional structure without target labels.

How (procedure/usage): Choose a goal (grouping, structure, anomaly detection), define a similarity measure, then apply an algorithm.

Cheat sheet / must-memorize

- **No labels:** only  $X$  (features), no  $y$ .
- **Main tasks:** clustering, dimensionality reduction, anomaly detection.
- **Key idea:** “similar” points should end up close or in the same group.

Common pitfall

Interpreting clusters as ground truth classes without validation; clusters depend on distance metrics and scaling.

Advantages: finds structure without labels.

Limitations: results can be subjective and metric-dependent.

**Key takeaways:** Unsupervised learning discovers structure without labels; clustering groups similar points.

### 3.2 Introduction to Clustering

Exam likelihood: High

Expect definitions and a short compare/contrast of clustering goals or algorithms. Answer: clustering groups similar points; results depend on metric and method.

Examiner favorite (what they love to ask)

Explain what a cluster is and why distance/similarity choice matters. Answer: clusters maximize within-group similarity; metric choice changes the grouping.

Why (motivation): We want to group similar observations when no labels exist.

What (definition): Clustering partitions data into groups so points in the same cluster are more similar to each other than to points in other clusters.

How (procedure/usage): Choose a similarity/distance metric, pick a clustering method, then evaluate/interpret the groups.

#### Cheat sheet / must-memorize

- **Goal:** high within-cluster similarity, low between-cluster similarity.
- **Distance matters:** Euclidean, Manhattan, cosine, etc., change the result.
- **Outputs:** cluster labels (hard) or membership scores (soft).

#### Common pitfall

Clustering is not “ground truth”; different scaling or distance metrics can change clusters drastically.

Advantages: reveals hidden group structure.

Limitations: sensitive to metric, scale, and algorithm choice.

**Key takeaways:** Clustering groups similar points; metric choice and scaling drive the result.

### 3.3 Similarity and Dissimilarity

#### Exam likelihood: High

Common: compute a distance/similarity and explain when to use each metric. Answer: apply the correct formula and justify the metric for the data type.

#### Examiner favorite (what they love to ask)

Compare Euclidean vs Manhattan vs cosine; note the role of scaling/normalization. Answer: Euclidean (L2) vs Manhattan (L1) for continuous data, cosine for direction; scale features first.

Why (motivation): Clustering depends on “closeness”; wrong metric gives wrong groups.

What (definition): Similarity measures how alike points are; dissimilarity (distance) measures how far apart they are.

How (procedure/usage): Choose a metric that matches data type and scale; normalize features when needed.

**Goal statement:** Maximize similarity within the same group and maximize dissimilarity between different groups.

#### Cheat sheet / must-memorize

- **Jaccard (nominal/binary):**  $J(A, B) = \frac{|A \cap B|}{|A \cup B|}$ ; distance  $d_J = 1 - J$ .
- **Minkowski (continuous):**  $d_p(x, y) = (\sum_i |x_i - y_i|^p)^{1/p}$ .
- **Manhattan:**  $d_1 = \sum_i |x_i - y_i|$  (Minkowski with  $p = 1$ ).
- **Euclidean:**  $d_2 = \sqrt{\sum_i (x_i - y_i)^2}$  (Minkowski with  $p = 2$ ).
- **Chebyshev:**  $d_\infty = \max_i |x_i - y_i|$  (Minkowski as  $p \rightarrow \infty$ ).
- **Scaling:** standardize/normalize when units differ.

#### Common pitfall

Using Euclidean on unscaled features (one large-unit feature dominates).

Advantages: lets you tailor similarity to data type.

Limitations: wrong metric or scaling yields misleading clusters.

**Which to use (rule of thumb):**

- **Jaccard:** nominal/binary attributes or set-like data (presence/absence).

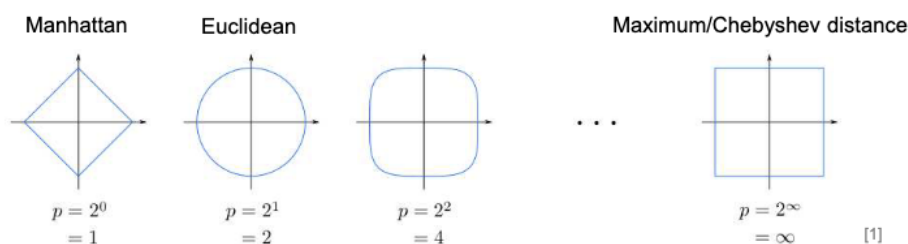
- **Manhattan:** high-dimensional data or when you want axis-aligned distance.
- **Euclidean:** continuous features with comparable scale and spherical clusters.
- **Chebyshev:** when the maximum coordinate difference is what matters (strict tolerance).

## Visual

### Continuous Features – Minkowski Distance (Metric)

Generalization of Manhattan and Euclidean distance to any natural dimension  $p \geq 1$  (also called  $L^p$  norm)

**Note:** only Manhattan, Euclidean and Chebyshev distances are relevant for the exam.



### Minkowski Distance – Relevant Examples

$$d(\mathbf{x}_i, \mathbf{x}_j) = \sqrt[p]{|x_{i1} - x_{j1}|^p + |x_{i2} - x_{j2}|^p + \dots + |x_{iD} - x_{jD}|^p}$$

**Manhattan distance  $p = 1$**

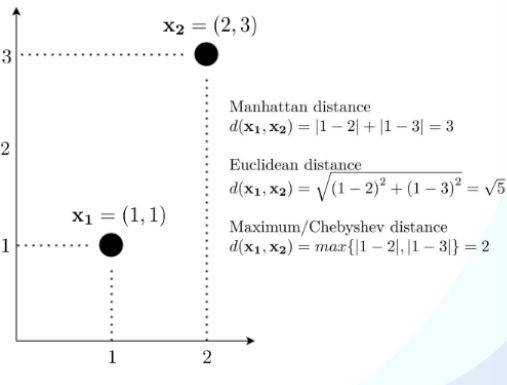
$$d(\mathbf{x}_i, \mathbf{x}_j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots + |x_{iD} - x_{jD}|$$

**Euclidean distance  $p = 2$**

$$d(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{(x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 + \dots + (x_{iD} - x_{jD})^2}$$

**Maximum/Chebyshev distance  $p \rightarrow \infty$**

$$d(\mathbf{x}_i, \mathbf{x}_j) = \lim_{p \rightarrow \infty} \left( \sum_{d=1}^D |x_{id} - x_{jd}|^p \right)^{\frac{1}{p}} = \max_{d \in \{1, \dots, D\}} |x_{id} - x_{jd}|$$



**Key takeaways:** Metric choice and scaling control the notion of similarity and the resulting clusters.

## 3.4 K-means and K-medoids

Exam likelihood: High

Often: list K-means steps and compare K-means vs K-medoids. Answer: K-means = init/assign/update/repeat; K-medoids uses medoids and is more robust.

Examiner favorite (what they love to ask)

Run 1 iteration of K-means by hand (assign & update) or explain why K-medoids is more robust to outliers. Answer: compute distances, assign clusters, recompute means; medoids resist outliers.

Why (motivation): Need a simple, scalable way to cluster continuous data.

What (definition): K-means uses centroids (means) to minimize within-cluster SSE; K-medoids uses

actual data points as centers and minimizes total within-cluster distance.

How (procedure/usage): Initialize  $k$  centers, assign each point to nearest center, update centers, repeat until convergence.

**Comparison:** K-means is faster but sensitive to outliers and scaling; K-medoids is more robust but slower.

**Difference:** K-means centers are means (can be non-data points); K-medoids centers are actual data points.

#### Cheat sheet / must-memorize

- **K-means objective:** minimize  $\sum_{k=1}^K \sum_{x_i \in C_k} \|x_i - \mu_k\|^2$ .
- **K-means steps:** (1) choose  $K$  and initialize  $\mu_k$ ; (2) assign each point to nearest  $\mu_k$ ; (3) update  $\mu_k$  to the mean of its cluster; (4) repeat 2–3 until assignments stop changing.
- **K-means optimization:** alternating minimization (assignment step + centroid update); decreases SSE each iteration.
- **K-medoids center:** medoid  $m_k$  is a data point minimizing total distance in its cluster.
- **K-medoids steps:** (1) choose  $K$  medoids  $m_k$ ; (2) assign each point to nearest medoid; (3) for each cluster, swap medoid candidate to reduce total distance; (4) repeat 2–3 until no improvement.
- **K-medoids optimization:** minimize  $\sum_k \sum_{x_i \in C_k} d(x_i, m_k)$  via swap heuristics (e.g., PAM); monotonic improvement, not guaranteed global optimum.
- **Robustness:** K-medoids less sensitive to outliers than K-means.

#### Common pitfall

K-means assumes roughly spherical clusters, is sensitive to scaling/initialization, and handles outliers poorly. K-medoids is more robust but slower on large datasets.

Advantages: K-means is fast; K-medoids is robust.

Limitations: K-means sensitive to outliers; K-medoids is slower.

**Key takeaways:** K-means is fast but sensitive; K-medoids is more robust but costlier.

### 3.5 Agglomerative Clustering

#### Exam likelihood: High

Common: list agglomerative steps and compute one merge with a given linkage. Answer: start with singletons, merge closest by linkage, update distances, repeat.

#### Examiner favorite (what they love to ask)

Given a small distance matrix, perform one or two agglomerative merges and state the linkage used. Answer: compute linkage distances and merge the smallest pair step by step.

Why (motivation): We want a hierarchy of clusters and do not want to pre-specify  $K$ .

What (definition): Agglomerative hierarchical clustering starts with each point as its own cluster and repeatedly merges the closest clusters.

How (procedure/usage): Choose a distance metric + linkage rule, then merge until one cluster remains (cut the dendrogram to pick  $K$ ).

### Cheat sheet / must-memorize

- **Algorithm (simplified):** start with  $n$  singleton clusters; compute inter-cluster distances; merge closest pair; update distances; repeat.
- **Single linkage:**  $d(A, B) = \min_{x \in A, y \in B} d(x, y)$ .
- **Complete linkage:**  $d(A, B) = \max_{x \in A, y \in B} d(x, y)$ .
- **Average linkage:**  $d(A, B) = \frac{1}{|A||B|} \sum_{x \in A} \sum_{y \in B} d(x, y)$ .
- **Centroid linkage:**  $d(A, B) = \|\mu_A - \mu_B\|$ .
- **Ward:** merge that yields the smallest increase in within-cluster SSE.

### Common pitfall

Single linkage can “chain” clusters; results depend heavily on scaling and the chosen linkage.

Advantages: produces a full hierarchy; no need to choose  $K$  upfront.

Limitations:  $O(n^2)$  memory/time; early merges are irreversible.

### Linkage Measures

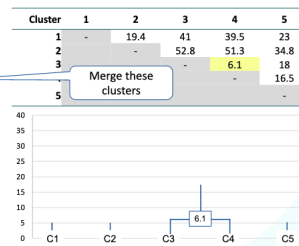
- The distance between clusters is otherwise known as **linkage measure**
- Four widely used linkage measures:



### Simplistic Agglomerative Clustering – Example

Cluster	Country	Meat & Fish	Plants	Eggs & Milk
1	Albania	11.7	50.1	9.4
2	Bulgaria	15	65.7	9.9
3+4	Austria Belgium	25 27.3	37.2 38.4	24.2 21.6
5	Czechoslovakia	23.1	44.4	15.3

- Compare countries based on sources of protein
- Use mean linkage between clusters (compare centroids)
- Use Manhattan distance between instances



### Simplistic Agglomerative Clustering – Example

Cluster	Country	Meat & Fish	plants	Eggs & Milk
1	Albania	11.7	50.1	9.4
2	Bulgaria	15	65.7	9.9
3+4	Austria Belgium	26.15 27.3	37.8 38.4	22.9 21.6
5	Czechoslovakia	23.1	44.4	15.3

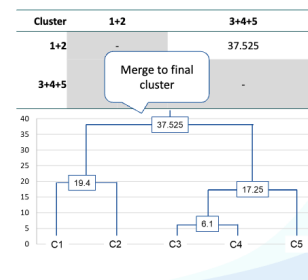
- Compare countries based on sources of protein
- Use mean linkage between clusters (compare centroids)
- Use Manhattan distance between instances



### Simplistic Agglomerative Clustering – Example

Cluster	Country	Meat & Fish	Plants	Eggs & Milk
1+2	Albania Bulgaria	13.35	57.9	9.65
3+4+5	Austria Belgium Czechoslovakia	24.625	41.1	19.1

- Compare countries based on sources of protein
- Use mean linkage between clusters (compare centroids)
- Use Manhattan distance between instances



**Key takeaways:** Agglomerative builds a dendrogram using a linkage rule; choose the cut for  $K$ .

## 3.6 DBSCAN

Exam likelihood: High

Often: define  $\varepsilon$  and MinPts, classify core/border/noise points. Answer: core has  $\geq$  MinPts in  $\varepsilon$ -neighborhood; border within  $\varepsilon$  of a core; noise otherwise.

### Examiner favorite (what they love to ask)

Given a small plot, label core/border/noise and explain why DBSCAN finds non-spherical clusters.  
Answer: use density reachability; clusters can follow arbitrary shapes.

Why (motivation): K-based methods struggle with arbitrary shapes and noise.

What (definition): DBSCAN groups dense regions using  $\varepsilon$ -neighborhoods and MinPts; points in sparse regions become noise.

How (procedure/usage): Pick  $\varepsilon$  and MinPts, mark core points, expand clusters via density reachability, label remaining as noise/border.

### Cheat sheet / must-memorize

- **Parameters:**  $\varepsilon$  (radius), MinPts (minimum neighbors).
- **Core point:** has at least MinPts points in its  $\varepsilon$ -neighborhood.
- **Border point:** within  $\varepsilon$  of a core but not itself core.
- **Noise:** not reachable from any core.
- **Density-reachable/connected:** chains of core points link a cluster.

### Common pitfall

Bad  $\varepsilon$ /MinPts choices: too small  $\rightarrow$  many noise points; too large  $\rightarrow$  merged clusters.

**Advantages:** no need to choose  $K$ ; finds arbitrary-shaped clusters; handles noise well.

**Limitations:** struggles with varying density; parameter sensitive; distance metrics degrade in high dimensions.

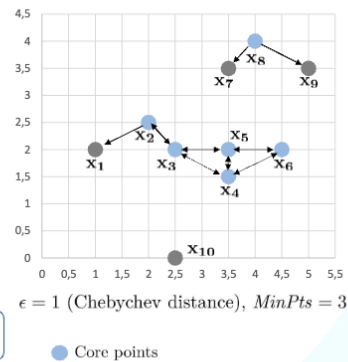


## Density-Based (DBSCAN)

## DBSCAN

- Two parameters:
  - Neighborhood size  $\epsilon$ : every point within distance  $\epsilon$  of an instance  $x_i$  is in the neighborhood
  - Density Threshold *MinPts*: a neighborhood with at least *MinPts* instances is considered dense
- Instance  $x_i$  is a **core point** if its neighborhood is dense (at least *MinPts*, including  $x_i$ , are within distance  $\epsilon$ )
- An instance  $x_j$  is **directly reachable** from  $x_i$  if
  - $x_i$  is a **core point**
  - $x_j$  is within distance  $\epsilon$  from  $x_i$

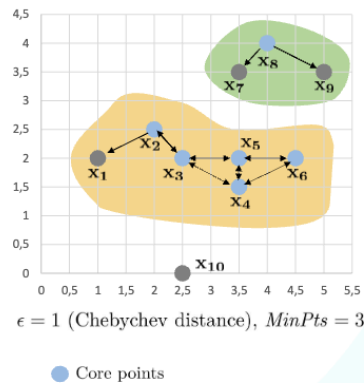
Reachability is **transitive**  
but is **not symmetric**



## DBSCAN - Algorithm

- Until** no unvisited core point is left **repeat**:
  - Pick any core point  $x_i$  that is not part of a cluster
  - Create a new cluster containing  $x_i$
  - Add all points reachable from  $x_i$  to the cluster
- Return the clusters
- All remaining points are marked as outliers

$C_1 = \{x_2, x_1, x_3, x_4, x_5, x_6\}$   
 $C_2 = \{x_8, x_7, x_9\}$   
 Outliers =  $\{x_{10}\}$



**Key takeaways:** DBSCAN clusters by density and naturally flags noise; choose  $\epsilon$  and *MinPts* carefully.

## 3.7 Closing

## 4 Frequent Itemsets

### 4.1 Introduction

Exam likelihood: High

Often: define itemset, support, and why frequent itemsets matter. Answer: itemset is a set of items; support is fraction of transactions; frequent if above minsup.

Examiner favorite (what they love to ask)

Compute support for a small transaction set and identify frequent itemsets above a threshold. Answer: count transactions containing the itemset, divide by total, compare to minsup.

Why (motivation): Discover co-occurrence patterns in transaction data (e.g., market baskets).

What (definition): An itemset is a set of items; a frequent itemset appears in at least a minimum fraction of transactions.

How (procedure/usage): Count itemset support across transactions and keep those meeting a minimum support threshold.

Cheat sheet / must-memorize

- **Transaction DB:** each transaction is a set of items.
- **Itemset:** a subset of items (e.g.,  $\{A, B\}$ ).
- **Support:**  $\text{supp}(X) = \frac{\text{\#transactions containing } X}{\text{\#transactions}}$ .
- **Frequent:**  $\text{supp}(X) \geq \text{min\_support}$ .
- **Apriori idea:** if an itemset is frequent, all its subsets are frequent.

Common pitfall

Mixing up support count (absolute) vs support (fraction); always state which one you use.

Advantages: highlights common co-occurrence patterns.

Limitations: can generate many itemsets; sensitive to min support.

### Frequent Itemsets – Notation

- $\mathcal{I} = \{I_1, I_2, \dots, I_D\}$  is the **set of all possible items**
- $\mathcal{A} \subseteq \mathcal{I}$  is an **itemset**
- A **transaction**  $\mathcal{T}$  is a **non-empty itemset**
- A **dataset**  $\mathcal{X}$  is a collection of transactions
- Technically  $\mathcal{X} \in \mathbb{M}(\mathbb{P}(\mathcal{I}))$  such that  $\emptyset \notin \mathcal{X}$  ( $\mathbb{M}$  is the multiset and  $\mathbb{P}$  is the powerset operator)

### Frequent Itemsets – Example

ID	Che	Bre	Chi	Mil	...	Pas
1	2 (true)	0 (false)	0 (false)	3 (true)	0 (false)	2 (true)
2	0 (false)	0 (false)	1 (true)	1 (true)	0 (false)	0 (false)
3	2 (true)	1 (true)	0 (false)	0 (false)	0 (false)	0 (false)
4	0 (false)	1 (true)	0 (false)	0 (false)	0 (false)	0 (false)

$\mathcal{X} \in \mathbb{M}(\mathbb{P}(\mathcal{I}))$

- Set of all items  $\mathcal{I} = \{Che, Bre, Chi, Mil, \dots, Pas\}$
- Transaction  $\mathcal{T}_1 = \{Che, Mil, Pas\} \subseteq \mathcal{I}$
- Dataset with four transactions  $\mathcal{X} = \{\{Che, Mil, Pas\}, \{Chi, Mil\}, \{Che, Bre\}, \{Bre\}\}$
- Dataset with ten transactions  $\mathcal{X} = \{\{Che, Mil, Pas\}^4, \{Chi, Mil\}^3, \{Che, Bre\}^2, \{Bre\}^1\}$

### Frequent Itemsets – Support

$$\text{support}(\mathcal{A}) = \frac{|\{T \in \mathcal{X} \mid \mathcal{A} \subseteq T\}|}{|\mathcal{X}|}$$

(relative)

Fraction of transactions  $\mathcal{T}$  in dataset  $\mathcal{X}$  that cover the itemset  $\mathcal{A}$

$$\text{support\_count}(\mathcal{A}) = |\{T \in \mathcal{X} \mid \mathcal{A} \subseteq T\}|$$

(absolute, also called frequency or count)

- Minimum **support threshold**:
  - **min\_sup**: lower bound for **support**( $\mathcal{A}$ )
  - **min\_sup\_count**: lower bound for **support\_count**( $\mathcal{A}$ )
- An itemset is **frequent** if its support is higher than **min\_sup** (or **min\_sup\_count**)
- Frequent itemsets are used to find **association rules**

### Support – Example

ID	Che	Bre	Chi	Mil	...	Pas
1	2 (true)	0 (false)	0 (false)	3 (true)	0 (false)	2 (true)
2	0 (false)	0 (false)	1 (true)	1 (true)	0 (false)	0 (false)
3	2 (true)	1 (true)	0 (false)	0 (false)	0 (false)	0 (false)
4	1 (true)	1 (true)	0 (false)	1 (true)	0 (false)	0 (false)

Dataset  $\mathcal{X} = \{\{Che, Mil, Pas\}, \{Chi, Mil\}, \{Che, Bre\}, \{Che, Bre, Mil\}\}$

Itemset  $\mathcal{A} = \{Che, Mil\} \subseteq \mathcal{I}$

$\text{support\_count}(\mathcal{A}) = |\{T \in \mathcal{X} \mid \mathcal{A} \subseteq T\}| = |\{T_1, T_4\}| = 2$

$\text{support}(\mathcal{A}) = \frac{|\{T \in \mathcal{X} \mid \mathcal{A} \subseteq T\}|}{|\mathcal{X}|} = \frac{|\{T_1, T_4\}|}{4} = \frac{2}{4}$

$\mathcal{A}$  is frequent if  $\text{min\_sup} \leq 0.5$

Itemset  $\mathcal{B} = \{Mil\} \subseteq \mathcal{I}$

$\text{support\_count}(\mathcal{B}) = |\{T \in \mathcal{X} \mid \mathcal{B} \subseteq T\}| = |\{T_1, T_2, T_4\}| = 3$

$\text{support}(\mathcal{B}) = \frac{|\{T \in \mathcal{X} \mid \mathcal{B} \subseteq T\}|}{|\mathcal{X}|} = \frac{|\{T_1, T_2, T_4\}|}{4} = \frac{3}{4}$

$\mathcal{B}$  is frequent if  $\text{min\_sup} \leq 0.75$

**Key takeaways:** Frequent itemsets reveal common co-occurrences; support is the key measure.

## 4.2 Properties of Frequent Itemsets

Exam likelihood: High

Often: define closed vs maximal itemsets and compare them with frequent itemsets. Answer: closed has no superset with same support; maximal has no frequent superset.

Examiner favorite (what they love to ask)

Given a small set of frequent itemsets, identify which are closed and which are maximal. Answer: maximal = no frequent superset; closed = no superset with equal support.

Why (motivation): Frequent itemsets can be many; closed/maximal summarize them compactly. What (definition): Closed itemsets keep full support information; maximal itemsets keep only the largest frequent ones.

How (procedure/usage): Mine frequent itemsets first, then filter for closed or maximal conditions.

### Cheat sheet / must-memorize

- **Frequent:**  $\text{supp}(X) \geq \text{min\_support}$ .
- **Closed:** no proper superset of  $X$  has the *same* support.
- **Maximal:** no proper superset of  $X$  is frequent.
- **Relationship:** Maximal  $\subset$  Closed  $\subset$  Frequent.
- **Info loss:** Maximal loses support counts; closed preserves supports for all frequent itemsets.

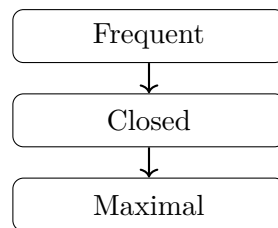
### Common pitfall

Confusing “closed” with “maximal”: closed allows frequent supersets as long as support drops; maximal allows no frequent supersets.

Advantages: closed/maximal reduce output size.

Limitations: maximal loses support info; closed can still be large.

### Visual



**Key takeaways:** Closed keeps support info; maximal is a compact summary but loses support counts.

## 4.3 Apriori Algorithm

### Exam likelihood: High

Common: outline Apriori steps and generate  $C_k$  and  $L_k$  for a tiny dataset. Answer: build  $L_1$ , generate/prune  $C_k$ , count supports, repeat to get  $L_k$ .

### Examiner favorite (what they love to ask)

Show candidate generation and pruning using the Apriori property. Answer: join  $L_{k-1}$  to form  $C_k$ , prune any candidate with an infrequent subset.

Why (motivation): Brute-force itemset search is exponential; Apriori prunes early.

What (definition): Apriori uses the downward-closure property: all subsets of a frequent itemset are frequent.

How (procedure/usage): Iteratively generate candidates and prune using frequent subsets.

### Cheat sheet / must-memorize

- **Step 1:** find frequent 1-itemsets  $L_1$ .
- **Step 2:** generate candidate  $k$ -itemsets  $C_k$  from  $L_{k-1}$  (join step).
- **Step 3:** prune any candidate in  $C_k$  with an infrequent  $(k-1)$ -subset.
- **Step 4:** scan DB to count supports; keep  $L_k$  (frequent candidates).
- **Stop:** when  $L_k$  is empty.

## Common pitfall

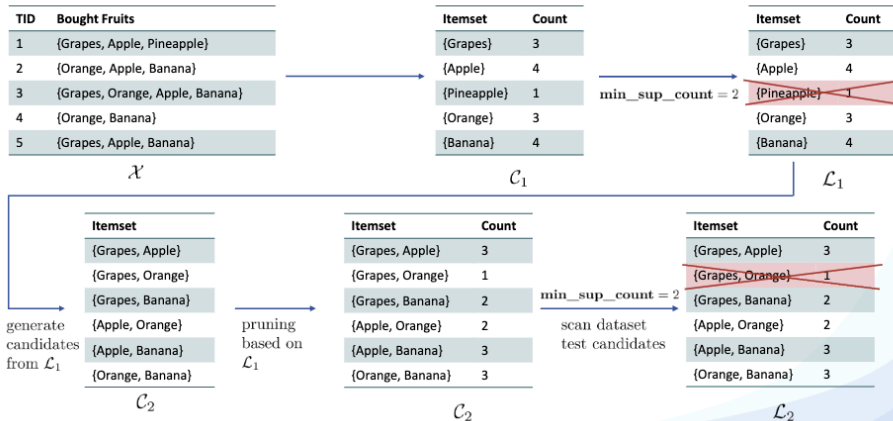
Forgetting the prune step (using the Apriori property) leads to too many candidates.

**Advantages:** strong pruning via downward-closure; easy to implement and explain.

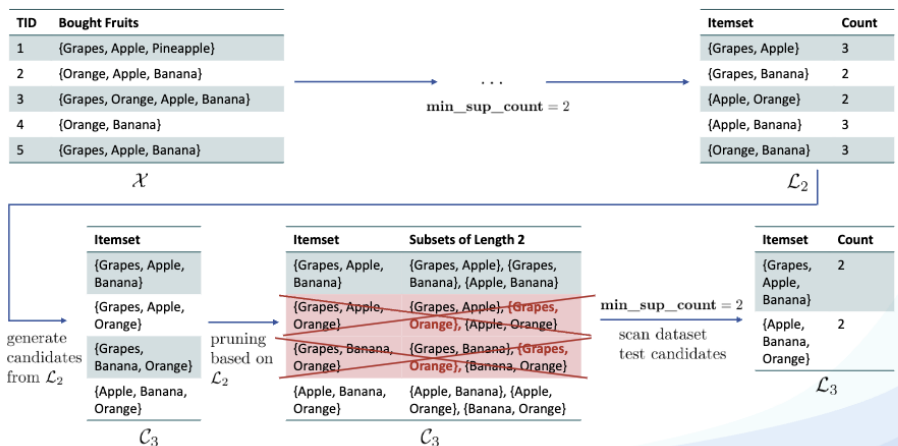
**Limitations:** challenging candidate generation; each candidate must be tested against the whole dataset.

## Visual

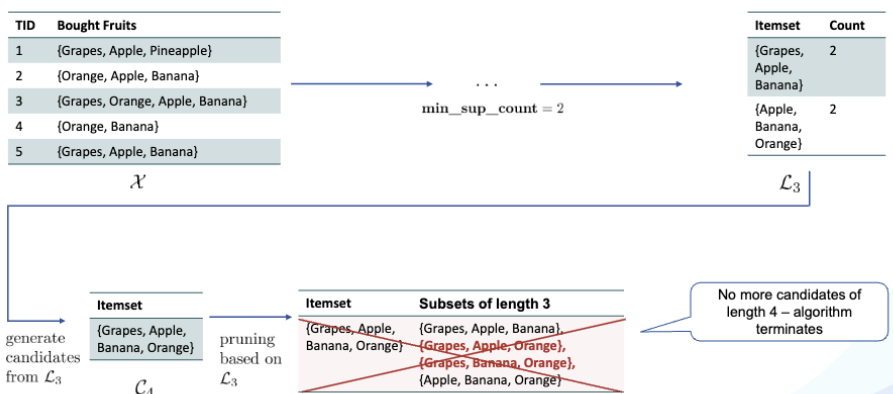
### Example



### Example



### Example



**Key takeaways:** Apriori alternates candidate generation and pruning using downward-closure.

## 4.4 FP-Growth Algorithm

Exam likelihood: High

Often: explain why FP-Growth is faster than Apriori and outline its steps. Answer: no candidate generation; build FP-tree in two passes and mine with DFS.

Examiner favorite (what they love to ask)

State the two DB passes and describe the FP-tree + conditional pattern base idea. Answer: pass 1 counts and orders items; pass 2 builds FP-tree; mine conditional bases/trees by DFS.

Why (motivation): Apriori's candidate explosion and repeated full scans are costly.

What (definition): FP-Growth mines frequent itemsets by building an FP-tree and doing DFS on conditional pattern bases (no candidate generation).

How (procedure/usage): Make two passes, build the FP-tree, then recursively mine conditional FP-trees (depth-first).

Cheat sheet / must-memorize

- **Pass 1:** count item supports; keep frequent items and order them.
- **Pass 2:** build FP-tree by inserting ordered transactions.
- **DFS mining:** for each item, build its conditional pattern base and conditional FP-tree; recurse.
- **Output:** frequent itemsets without candidate generation.

Common pitfall

Not ordering items by global frequency before building the FP-tree (tree becomes large and inefficient).

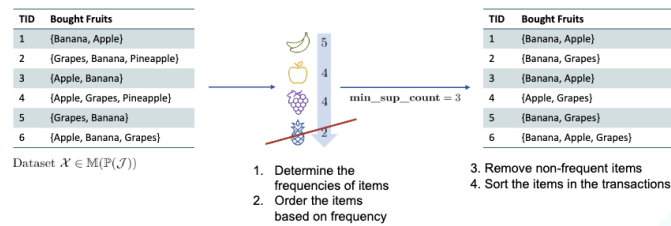
**Advantages:** avoids candidate generation; only two DB passes; efficient for dense datasets.

**Limitations:** FP-tree can be large in very sparse data; more complex to implement.

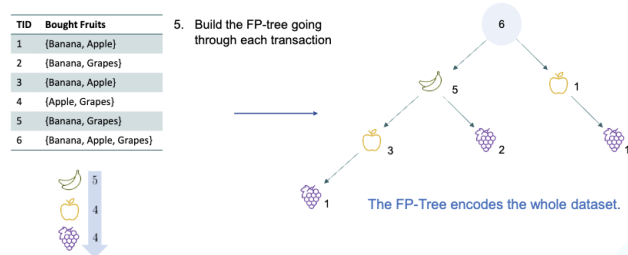
## FP-Growth Steps

1. Determine the frequency of each item (first pass through the dataset)
  2. Sort  $\mathcal{I} = \{I_1, \dots, I_D\}$  based on their frequencies ( $I_1$  is most frequent,  $I_D$  is the least frequent)
  3. Remove the non-frequent items from all itemsets (keep the remainder of the transactions)
  4. The remaining items in each transactions are ordered by frequency (same as above)
  5. This can be used to build a so-called prefix tree (second pass through the dataset)
6. The resulting FP-tree contains all information needed to find the frequent itemsets of any length (no need to traverse the dataset again)

### Constructing FP-Tree – Example



### Constructing FP-Tree – Example

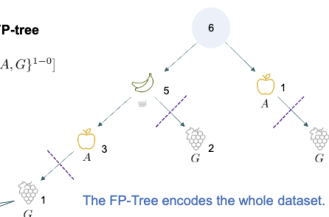


### FP-Tree – Cannot Cut Naïvely

We can read the transactions from the FP-tree

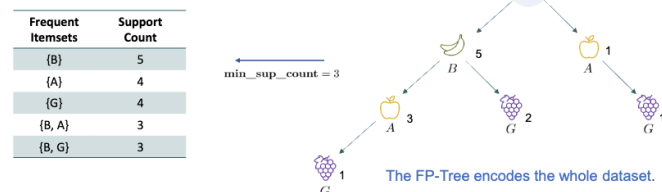
$$\begin{aligned} \mathcal{X} &= [\{B, A, G\}^{1-0}, \{B, A\}^{3-1}, \{B, G\}^{2-0}, \{A, G\}^{1-0}] \\ &= [\{B, A, G\}, \{B, A\}^2, \{B, G\}^2, \{A, G\}] \end{aligned}$$

Even though G exists only once in this subtree, we can't cut it naively, because its support may be higher than the threshold ( $4 \geq 3$ )



### FP-Tree – Frequent Itemsets

Next: Mining the FP-tree to obtain frequent itemsets



**Key takeaways:** FP-Growth uses an FP-tree and DFS to mine patterns with only two passes.

## 5 Association Rules

### 5.1 Introduction

#### Exam likelihood: High

Often: define an association rule and its parts (antecedent/consequent) and relate rules to frequent itemsets. Answer: rule  $X \rightarrow Y$  with  $X \cap Y = \emptyset$ ; rules are generated from frequent itemsets.

#### Examiner favorite (what they love to ask)

Given a tiny basket dataset, propose a reasonable rule and interpret it in words. Answer: e.g.,  $\{\text{bread}\} \rightarrow \{\text{butter}\}$  means baskets with bread often also contain butter.

Why (motivation): We want simple, actionable “if-then” patterns that summarize co-purchase behavior.

What (definition): An association rule is an implication  $X \rightarrow Y$  between disjoint itemsets;  $X$  is the antecedent,  $Y$  the consequent.

How (procedure/usage): Mine frequent itemsets first, then generate candidate rules and keep those that pass quality thresholds.

#### Cheat sheet / must-memorize

- **Rule form:**  $X \rightarrow Y$ , with  $X \cap Y = \emptyset$ .
- **Antecedent:** left side (condition); **Consequent:** right side (outcome).
- **From itemsets:** rules are derived from frequent itemsets of size  $\geq 2$ .
- **Interpretation:** “If  $X$  occurs in a transaction,  $Y$  tends to occur too.”

#### Common pitfall

Treating rules as causal (“ $X$  causes  $Y$ ”). Association rules are correlational and can be misleading without lift/base-rate checks.

**Advantages:** highly interpretable patterns; useful for merchandising, cross-sell, and recommendation hints.

**Limitations:** many trivial/spurious rules; sensitive to thresholds; does not imply causation.



QUESTION

## Association Rules - Notation

- $\mathcal{I} = \{I_1, I_2, \dots, I_D\}$  is the set of all possible items
- A transaction  $\mathcal{T} \in \mathbb{P}(\mathcal{I}) \setminus \{\emptyset\}$  is a non-empty itemset
- A dataset  $\mathcal{X} \in \mathbb{M}(\mathbb{P}(\mathcal{I}))$  (such that  $\emptyset \notin \mathcal{X}$ ) is a multiset of transactions (Here,  $\mathbb{M}$  is the multiset and  $\mathbb{P}$  is the powerset operator)
- $\mathcal{A} \Rightarrow \mathcal{B}$  with  $\mathcal{A} \subseteq \mathcal{I}, \mathcal{B} \subseteq \mathcal{I}$  and  $\mathcal{A} \cap \mathcal{B} = \emptyset$  is an association rule
- For example,  $\{\text{Cheese, Bread}\} \Rightarrow \{\text{Milk}\}$

Same as before

$$\mathcal{A} \Rightarrow \mathcal{B}$$

## Support and Confidence - Example



ID	Bought Items
1	{Bread, Cheese, Milk, Pasta}
2	{Bread, Cheese, Chips}
3	{Cheese, Pasta, Milk}
4	{Bread, Cheese, Milk}
5	{Bread, Pasta}

$$\begin{aligned} \text{support}(\{\text{Bread} \Rightarrow \{\text{Cheese, Milk}\}\}) &= \text{support}(\{\text{Bread, Cheese, Milk}\}) = \frac{2}{5} \\ &= \text{support}(\{\{\text{Cheese, Milk}\} \Rightarrow \{\text{Bread}\}\}) \\ &= \text{support}(\{\text{Bread, Cheese}\} \Rightarrow \{\text{Milk}\}) \end{aligned}$$

Symmetric: moving an item does not change the value

## Support and Confidence - Example



ID	Bought Items
1	{Bread, Cheese, Milk, Pasta}
2	{Bread, Cheese, Chips}
3	{Cheese, Pasta, Milk}
4	{Bread, Cheese, Milk}
5	{Bread, Pasta}

$$\begin{aligned} \text{conf}(\{\text{Bread} \Rightarrow \{\text{Cheese, Milk}\}\}) &= \frac{\text{support}(\{\text{Bread, Cheese, Milk}\})}{\text{support}(\{\text{Bread}\})} = \frac{2}{4} \\ \text{conf}(\{\{\text{Cheese, Milk}\} \Rightarrow \{\text{Bread}\}\}) &= \frac{\text{support}(\{\text{Bread, Cheese, Milk}\})}{\text{support}(\{\text{Cheese, Milk}\})} = \frac{2}{3} \end{aligned}$$

General rule:  
 $\text{conf}(\{A, B\} \Rightarrow \{C\}) \geq \text{conf}(\{A\} \Rightarrow \{B, C\})$ 

## Key takeaways:

- Association rules express “if-then” co-occurrence patterns.
- Rules come from frequent itemsets and need quality thresholds to be useful.

## 5.2 Generating Association Rules

Exam likelihood: High

Often: describe how to generate rules from a frequent itemset and how to prune them. Answer: split a frequent itemset  $I$  into  $X$  and  $Y = I \setminus X$ ; keep rules with confidence above  $\text{min\_conf}$  (and usually lift).

Examiner favorite (what they love to ask)

Given a frequent itemset  $I = \{A, B, C\}$ , list candidate rules and compute confidence for one. Answer: all non-empty proper splits, e.g.,  $AB \rightarrow C$ ,  $A \rightarrow BC$ , etc.

Why (motivation): Mining frequent itemsets is not enough; we need directional rules that are strong and useful.

What (definition): A rule  $X \rightarrow Y$  is generated from a frequent itemset  $I$  by choosing  $X \subset I$  and  $Y = I \setminus X$ .

How (procedure/usage): For each frequent itemset, enumerate non-empty splits, compute confidence (and lift), and prune by thresholds.

#### Cheat sheet / must-memorize

- **Candidates:** for itemset  $I$ , all rules  $X \rightarrow I \setminus X$  with  $\emptyset \subset X \subset I$ .
- **Confidence:**  $\text{conf}(X \rightarrow Y) = \frac{\text{supp}(X \cup Y)}{\text{supp}(X)}$ .
- **Pruning:** if a rule fails  $\text{min\_conf}$ , any rule with smaller antecedent may still pass (use careful pruning).
- **Typical flow:** mine frequent itemsets  $\rightarrow$  generate rules  $\rightarrow$  filter by  $\text{min\_conf}$  (and lift).

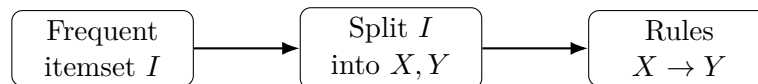
#### Common pitfall

Generating rules from infrequent itemsets or forgetting to enforce  $X \cap Y = \emptyset$ .

**Advantages:** systematic generation from frequent itemsets; simple to explain and compute.

**Limitations:** rule count can explode; high confidence can be misleading without lift.

#### Visual



#### Key takeaways:

- Rules come from splitting frequent itemsets into antecedent and consequent.
- Use confidence (and lift) thresholds to prune weak or misleading rules.

### 5.3 Evaluation (support, confidence, lift, conviction)

#### Exam likelihood: High

Often: compute support, confidence, and lift for a rule and interpret the result. Answer: use counts (or a confusion matrix) to compute  $P(X \cap Y)$ ,  $P(Y|X)$ , and  $P(Y|X)/P(Y)$ .

#### Examiner favorite (what they love to ask)

Interpret lift values: 1, much greater than 1, and much less than 1. Answer: 1 means independence;  $\gg 1$  strong positive association;  $\ll 1$  negative association.

Why (motivation): There are many rules; we need measures that balance coverage, reliability, and base rates.

What (definition): Support measures how often the rule occurs, confidence measures how often it is correct when it fires, and lift adjusts confidence by the base rate of  $Y$  (independence check).

How (procedure/usage): Compute metrics from counts (or confusion matrix), then filter by  $\text{min\_sup}/\text{minconf}$  and prefer high lift (and reasonable support).

### Cheat sheet / must-memorize

- **Support:**  $\text{supp}(X \rightarrow Y) = \text{supp}(X \cup Y) = P(X \cap Y)$ .
- **Confidence:**  $\text{conf}(X \rightarrow Y) = \frac{P(X \cap Y)}{P(X)} = P(Y|X)$ .
- **Lift:**  $\text{lift}(X \rightarrow Y) = \frac{P(Y|X)}{P(Y)} = \frac{P(X \cap Y)}{P(X)P(Y)}$ .
- **Conviction:**  $\text{conv}(X \rightarrow Y) = \frac{1-P(Y)}{1-\text{conf}}$  (directional strength).
- **Confusion matrix view:** predict  $Y$  when  $X$  occurs.
- **TP:**  $X$  and  $Y$ ; **FP:**  $X$  and not  $Y$ ; **FN:**  $Y$  and not  $X$ ; **TN:** neither.
- **Mapping:** confidence = precision =  $\frac{TP}{TP+FP}$ ; support =  $\frac{TP}{N}$ .

### Common pitfall

High confidence can be trivial if  $Y$  is very common; always check lift (and support).

**Advantages:** simple, interpretable rule quality metrics; lift corrects for base rates.

**Limitations:** can favor rare items or tiny supports; does not imply causation.

### Confusion matrix for association rules

Consider a set of transactions and an association rule  $\mathcal{A} \Rightarrow \mathcal{B}$

$\mathcal{A} \Rightarrow \mathcal{B}$	$\mathcal{B}$ is included	$\mathcal{B}$ is not included	
$\mathcal{A}$ is included	$\#AB$	$\#A\bar{B}$	$\#A$
$\mathcal{A}$ is not included	$\#\bar{A}B$	$\#\bar{A}\bar{B}$	$\#\bar{A}$
	$\#B$	$\#\bar{B}$	$\#ALL$

The lower the better

The higher the better

Not captured in any of the metrics

$$\text{support}(\mathcal{A} \Rightarrow \mathcal{B}) = \frac{\#AB}{\#ALL} \quad \text{conf}(\mathcal{A} \Rightarrow \mathcal{B}) = \frac{\#AB}{\#A}$$

### Support and Confidence Don't Tell The Full Story

Consider association rule  $\mathcal{A} \Rightarrow \mathcal{B}$

$\mathcal{A} \Rightarrow \mathcal{B}$	$\mathcal{B}$ is included	$\mathcal{B}$ is not included	
$\mathcal{A}$ is included	80	10	90
$\mathcal{A}$ is not included	0	10	10
	80	20	100

$$\text{support}(\mathcal{A} \Rightarrow \mathcal{B}) = \frac{80}{100}$$

$$\text{conf}(\mathcal{A} \Rightarrow \mathcal{B}) = \frac{80}{90}$$

Seems to be a good rule because if  $\mathcal{A}$  is not included,  $\mathcal{B}$  is also never included

### Support and Confidence Don't Tell The Full Story

Consider association rule  $\mathcal{A} \Rightarrow \mathcal{B}$

$\mathcal{A} \Rightarrow \mathcal{B}$	$\mathcal{B}$ is included	$\mathcal{B}$ is not included	
$\mathcal{A}$ is included	80	10	90
$\mathcal{A}$ is not included	10	0	10
	90	10	100

The distribution of counts in the second row does not influence support and confidence

$$\text{support}(\mathcal{A} \Rightarrow \mathcal{B}) = \frac{80}{100}$$

$$\text{conf}(\mathcal{A} \Rightarrow \mathcal{B}) = \frac{80}{90}$$

Same support and confidence, but seems to be a poor rule because if  $\mathcal{A}$  is not included,  $\mathcal{B}$  is always included

### Is the Rule Surprising?

Consider association rule  $\mathcal{A} \Rightarrow \mathcal{B}$

$\mathcal{A} \Rightarrow \mathcal{B}$	$\mathcal{B}$ is included	$\mathcal{B}$ is not included	
$\mathcal{A}$ is included	9	1	10
$\mathcal{A}$ is not included	81	9	90
	90	10	100

Different frequencies in the rows, but distribution is the same

$$\text{lift}(\mathcal{A} \Rightarrow \mathcal{B}) = \frac{\text{support}(\mathcal{A} \cup \mathcal{B})}{\text{support}(\mathcal{A}) \cdot \text{support}(\mathcal{B})} = \frac{P(\mathcal{A} \cup \mathcal{B})}{P(\mathcal{A}) \cdot P(\mathcal{B})} = \frac{\frac{\#AB}{\#ALL}}{\frac{\#A}{\#ALL} \cdot \frac{\#B}{\#ALL}}$$

$$\text{support}(\mathcal{A} \Rightarrow \mathcal{B}) = \frac{9}{100}$$

$$\text{conf}(\mathcal{A} \Rightarrow \mathcal{B}) = \frac{9}{10}$$

$$\text{lift}(\mathcal{A} \Rightarrow \mathcal{B}) = \frac{\frac{9}{100}}{\frac{10}{100} \cdot \frac{90}{100}} = 1$$

No surprise!

### Key takeaways:

- Lift = 1 means independence;  $\gg 1$  strong positive association;  $\ll 1$  negative association.
- Use support + confidence + lift together to judge rule quality.

## 5.4 Applications

## 5.5 Simpson's Paradox

### Exam likelihood: High

Often: define Simpson's paradox and explain how it can affect association rules or metrics.  
Answer: a trend in aggregated data reverses after stratifying by a confounder; subgroup analysis is required.

### Examiner favorite (what they love to ask)

Given a table with two subgroups, identify the paradox and name the confounder. Answer: the confounder changes the base rates; the combined trend is misleading.

Why (motivation): Aggregated metrics can hide or reverse real relationships, leading to bad rules or decisions.

What (definition): Simpson's paradox occurs when an association in pooled data reverses after conditioning on a confounding variable (a third variable).

How (procedure/usage): Always check key subgroups (stratify), compare within-group metrics, and look for confounders.

### Cheat sheet / must-memorize

- **Paradox signal:** pooled trend  $\neq$  subgroup trends (often reversed).
- **Cause:** confounding variable with different base rates across groups.
- **Fix:** stratify by the confounder; compare like with like.
- **Rule mining:** a high-lift rule may disappear or flip after stratification.

### Common pitfall

Reporting only aggregated confidence/lift and ignoring confounding variables that drive subgroup differences.

**Advantages:** forces careful interpretation; highlights need for subgroup analysis.

**Limitations:** requires enough data per subgroup; choosing stratification variables can be non-obvious.

### Simpson's Paradox - Example

Two classes: **old** and **young**

smoke $\Rightarrow$ cancer	has cancer	doesn't have cancer	
smokes	1 + 66	2 + 34	3 + 100
doesn't smoke	34 + 2	66 + 1	100 + 3
	35 + 68	68 + 35	103 + 103

**humans**  $\text{conf}(\text{smoke} \Rightarrow \text{cancer}) = \frac{67}{103} = 0.65 > \text{conf}(\text{not smoke} \Rightarrow \text{cancer}) = \frac{36}{103} = 0.35$

**old**  $\text{conf}(\text{smoke} \Rightarrow \text{cancer}) = \frac{1}{3} = 0.333 < \text{conf}(\text{not smoke} \Rightarrow \text{cancer}) = \frac{34}{100} = 0.34$

**young**  $\text{conf}(\text{smoke} \Rightarrow \text{cancer}) = \frac{66}{100} = 0.66 < \text{conf}(\text{not smoke} \Rightarrow \text{cancer}) = \frac{3}{3} = 0.666$

Smoking is healthy for **old** and **young** people, but not for all humans!



### Simpson's Paradox - Example

Two classes: **old** and **young**

smoke $\Rightarrow$ cancer	has cancer	doesn't have cancer	
smokes	1 + 66	2 + 34	3 + 100
doesn't smoke	34 + 2	66 + 1	100 + 3
	35 + 68	68 + 35	103 + 103

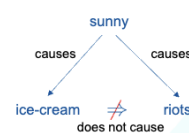
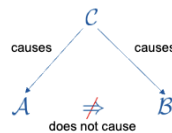
**humans**  $\text{lift}(\text{smoke} \Rightarrow \text{cancer}) = \frac{67}{103 \cdot \frac{67+36}{103}} = \frac{67 \cdot 206}{103 \cdot 103} = 1.301$  — Positively correlated

**old**  $\text{lift}(\text{smoke} \Rightarrow \text{cancer}) = \frac{1}{3 \cdot \frac{1+34}{100}} = \frac{1 \cdot 103}{3 \cdot 35} = 0.9809$  — Negatively correlated

**young**  $\text{lift}(\text{smoke} \Rightarrow \text{cancer}) = \frac{66}{103 \cdot \frac{66+3}{103}} = \frac{66 \cdot 103}{100 \cdot 68} = 0.9997$

### Confounding

- Simpson's paradox is related to **confounding**
- A (possibly hidden) confounding feature  $C$  (also called "lurking variable") may influence both  $A$  and  $B$ , and therefore "blur"  $A \Rightarrow B$



### Key takeaways:

- Aggregation can reverse associations; always check subgroups.
- Confounders drive the paradox; stratify before trusting a rule.

## 6 Time Series

### 6.1 Temporal Data

Exam likelihood: High

Often: define temporal data and list its key characteristics (order, granularity, seasonality).  
Answer: temporal data is indexed by time; ordering matters; it may show trend/seasonality and irregular sampling.

Examiner favorite (what they love to ask)

Explain why random train/test splits are wrong for time series. Answer: time order must be preserved to avoid leakage; use forward-chaining or hold-out of later periods.

Why (motivation): Time-indexed data appears in finance, sensors, logs, and business KPIs; ignoring time breaks analysis.

What (definition): Temporal data is a sequence of observations indexed by time; the order and spacing carry meaning.

How (procedure/usage): Inspect timestamps, granularity, missing/irregular intervals, and basic patterns (trend/seasonality) before modeling.

Cheat sheet / must-memorize

- **Time index:** timestamp or ordered time steps.
- **Granularity:** sampling rate (hourly/daily/etc.).
- **Patterns:** trend, seasonality, cycles, noise.
- **Irregularity:** gaps, unequal spacing, missing data.
- **Leakage risk:** never use future data to predict the past.

Common pitfall

Treating temporal data as i.i.d. and shuffling it; this destroys order and causes leakage.

**Advantages:** captures dynamics and change over time; supports forecasting and anomaly detection.

**Limitations:** non-stationarity and irregular sampling complicate modeling; needs careful preprocessing.

Visual

**Key takeaways:**

- Temporal data has ordered observations; order is part of the signal.
- Always respect time when splitting data and interpreting patterns.

### 6.2 Introduction to Time Series

Exam likelihood: High

Often: define a time series and its components (trend, seasonality, noise). Answer: a time series is an ordered sequence of observations; components include trend/seasonal/cycle/noise.

### Examiner favorite (what they love to ask)

Identify components from a small plot and state if the series is stationary. Answer: describe trend/seasonality; stationary means constant mean/variance over time.

Why (motivation): Time series models use structure in time to forecast and detect anomalies.

What (definition): A time series is a sequence  $\{x_t\}$  indexed by time; it can be decomposed into components.

How (procedure/usage): Visualize, decompose, and check stationarity before choosing models.

### Cheat sheet / must-memorize

- **Components:** trend, seasonality, cycle, noise.
- **Stationarity:** mean/variance (and autocovariance) stable over time.
- **Decomposition:** additive  $x_t = T_t + S_t + C_t + e_t$  or multiplicative  $x_t = T_t \cdot S_t \cdot C_t \cdot e_t$ .

### Common pitfall

Ignoring seasonality or trend and fitting a stationary model to a non-stationary series.

**Advantages:** explicit modeling of temporal structure; enables forecasting.

**Limitations:** non-stationarity and seasonality can break simple models.

### Visual

#### Sequences of Itemsets: Used as Input for Sequential Pattern Mining

Customer ID	Purchased Items	Time
1	A	15.12.22 12:25
1	A, B	15.12.22 12:45
2	B	15.12.22 13:01
3	C	30.12.22 18:01
3	A, C, D	11.01.23 17:25
4	B	31.12.22 17:32
...	...	...



Customer ID	Customer Sequence
1	$\langle\{A\}, \{A, B\}\rangle$
2	$\langle\{B\}\rangle$
3	$\langle\{C\}, \{A, C, D\}\rangle$
4	$\langle\{B\}\rangle$
...	...

Input is a **multiset** of sequences of itemsets

#### Frequent Sequence Patterns: Containment

Customer ID	Customer Sequence
1	$\langle\{beer, chips\}, \{wine, chips\}\rangle$
2	$\langle\{beer, chips\}, \{beer\}\rangle$
3	$\langle\{beer\}, \{beer\}, \{wine\}\rangle$
4	$\langle\{beer\}, \{chips\}, \{wine\}\rangle$
5	$\langle\{beer, chips, wine\}\rangle$

$\langle\{beer\}, \{wine\}\rangle$  is contained by 1,3,4  
 $\langle\{beer, wine\}\rangle$  is contained by 5  
 $\langle\{beer\}, \{beer\}\rangle$  is contained by 2,3  
 $\langle\{beer\}\rangle$  is contained by 1,2,3,4,5  
 $\langle\{beer\}, \{chips, wine\}\rangle$  is contained by 1

It is possible to apply the apriori principle to determine frequent sequence patterns!

### Key takeaways:

- Time series have components; visualize and decompose first.
- Stationarity matters for many models.



## 6.3 Analysis

### Exam likelihood: High

Often: explain ACF/PACF or identify lags and autocorrelation from a plot. Answer: ACF shows correlation with past lags; PACF shows direct lag effects.

### Examiner favorite (what they love to ask)

Given ACF/PACF plots, suggest a simple AR or MA order. Answer: cut-off in ACF suggests MA(q); cut-off in PACF suggests AR(p).

Why (motivation): We need to understand temporal dependence before choosing a forecasting model.

What (definition): Time-series analysis studies autocorrelation and lag structure; **ACF** (Autocorrelation Function) measures correlation of a series with its past lags. **White noise** is a series with zero mean, constant variance, and no autocorrelation at non-zero lags.

How (procedure/usage): Plot the series, check ACF, and inspect residuals for remaining structure.

### Cheat sheet / must-memorize

- **ACF**: correlation of  $x_t$  with  $x_{t-k}$  across lags.
- **Autocorrelation formula**:  $\rho_k = \frac{\sum_{t=k+1}^n (x_t - \bar{x})(x_{t-k} - \bar{x})}{\sum_{t=1}^n (x_t - \bar{x})^2}$ .
- **White noise**:  $\rho_k \approx 0$  for all  $k > 0$ .
- **Rule of thumb**: ACF cut-off suggests MA order.
- **Diagnostics**: residuals should look like white noise.

### Common pitfall

Choosing model orders from a single noisy plot; verify with diagnostics and hold-out tests.

**Advantages**: reveals dependence structure; guides model selection.

**Limitations**: plots can be noisy; patterns can be ambiguous.

## Correlation: Time Series Interpretation

datasets $\langle \dots, x_t, \dots \rangle$ $\langle \dots, y_t, \dots \rangle$ ( $T$ observations)	sample mean $\bar{x} = \frac{1}{T} \sum_t x_t$ $\bar{y} = \frac{1}{T} \sum_t y_t$	sample variance $\sigma_x^2 = \frac{1}{T-1} \sum_t (x_t - \bar{x})^2$ $\sigma_y^2 = \frac{1}{T-1} \sum_t (y_t - \bar{y})^2$
sample covariance $q_{xy} = \frac{1}{T-1} \sum_t (x_t - \bar{x})(y_t - \bar{y})$		
sample correlation coefficient $r_{xy} = \frac{q_{xy}}{\sigma_x \sigma_y} = \frac{\sum_t (x_t - \bar{x})(y_t - \bar{y})}{\sqrt{\sum_t (x_t - \bar{x})^2 \sum_t (y_t - \bar{y})^2}}$		

## Time Series Decomposition

- Various patterns coexist
- May be helpful to split a time series into several components
- Additive decomposition

$$y_t = S_t + T_t + R_t$$

trend-cycle component
seasonal component
remainder component

- Alternative: multiplicative decomposition  $y_t = S_t \times T_t \times R_t$

Autocorrelation: The Formula Given Lag  $k$ 

The linear relationship between lagged values of a time series.

The length of the time series

The coefficient between  $y_t$  and  $y_{t-k}$

Lag

Average of observation  $y$

Observation at time  $t$

$$r_k \approx \frac{\sum_{t=k+1}^T (y_t - \bar{y})(y_{t-k} - \bar{y})}{\sum_{t=1}^T (y_t - \bar{y})^2}$$

## Moving Average

Estimating the trend-cycle component  $T_t$ : **moving average**

- The estimate of the  $T_t$  is obtained by averaging values of the time series within  $k$  periods of  $t$ .
- A moving average of order  $m = 2k + 1$  is called **m-MA**.

$$\hat{T}_t = \frac{1}{m} \sum_{j=-k}^k y_{t+j}$$

$m = 2k + 1$

Annual electricity sales

Year	Sales (GWh)	5-MA
1989	2354.34	
1990	2379.71	
1991	2388.52	
1992	2468.99	
1993	2386.09	
1994	2494.47	
1995	2572.72	
1996	2762.72	
1997	2824.50	
1998	3000.20	
1999	3108.40	
2000	3197.50	
2001	3075.70	
2002	3180.60	
2003	3221.60	
2004	3178.30	
2005	3410.60	
2006	3527.68	
2007	3617.89	
2008	3555.00	

$k = 2$   
 $m = 5$

## Key takeaways:

- Use ACF/PACF to understand lag structure and pick model orders.
- Always validate with residual checks and time-based splits.

## 6.4 Forecasting

## Exam likelihood: High

Often: distinguish AR, MA, ARMA, and ARIMA; explain when differencing is needed. Answer: AR uses past values, MA uses past errors, ARMA combines both, ARIMA adds differencing for non-stationary series.

## Examiner favorite (what they love to ask)

Given a short series, explain how to split data and evaluate forecasts. Answer: keep chronological order; train on earlier data and test on later data; report MAE/MSE/RMSE or MAPE.

Why (motivation): We want models that capture temporal dependence and produce reliable forecasts.

What (definition): **AR** (Autoregressive) models regress on past values, **MA** (Moving Average) models regress on past errors, **ARMA** (Autoregressive Moving Average) combines AR+MA, and **ARIMA** (Autoregressive Integrated Moving Average) adds differencing to handle non-stationarity. How (procedure/usage): Check stationarity, difference if needed, choose AR/MA orders, fit, and evaluate on a time-ordered test set.

### Cheat sheet / must-memorize

- **AR( $p$ ):**  $x_t = c + \sum_{i=1}^p \phi_i x_{t-i} + \varepsilon_t$ .
- **MA( $q$ ):**  $x_t = c + \sum_{i=1}^q \theta_i \varepsilon_{t-i} + \varepsilon_t$ .
- **ARMA( $p, q$ ):** combine AR and MA terms.
- **ARIMA( $p, d, q$ ):** difference  $d$  times to make series stationary, then ARMA.
- **Stationarity:** constant mean/variance; use differencing to remove trend.
- **Evaluation:** MAE, MSE/RMSE, MAPE; time-ordered train/test split.
- **Granger causality:**  $X$  helps predict  $Y$  if past  $X$  improves forecast of  $Y$  beyond past  $Y$  alone.

### Common pitfall

Randomly shuffling time series for train/test or claiming “causality” from Granger tests; it is predictive, not causal in a real-world sense.

**Advantages:** interpretable baselines; fast to fit; good for short-term forecasting.

**Limitations:** assumes linear structure; sensitive to non-stationarity and outliers; limited for complex patterns.

### Autoregressive (AR) Models

- An **Autoregressive (AR)** model is a regression of the variable against itself.
- The variable of interest is forecasted using a **linear combination of past values of the variable**.

$$y_t = c + \phi_1 y_{t-1} + \phi_2 y_{t-2} + \dots + \phi_p y_{t-p} + \varepsilon_t$$

Annotations:  $\phi$  is "phi",  $\theta$  is "theta",  $\varepsilon$  is "epsilon".  
 coefficient:  $\phi_1, \phi_2, \dots, \phi_p$   
 past value of y at time t-1:  $y_{t-1}$   
 error term:  $\varepsilon_t$

- It is referred to as a **AR(p)** model, an **Autoregressive (AR)** model of order  $p$ .

### Moving Average (MA) Models

- A **Moving Average (MA)** model is a regression of the past errors.
- The variable of interest is forecasted using a linear combination of **past forecast errors**.

$$y_t = c + \varepsilon_t + \theta_1 \varepsilon_{t-1} + \theta_2 \varepsilon_{t-2} + \dots + \theta_q \varepsilon_{t-q}$$

Annotations:  $\phi$  is "phi",  $\theta$  is "theta",  $\varepsilon$  is "epsilon".  
 error terms:  $\varepsilon_t, \varepsilon_{t-1}, \dots, \varepsilon_{t-q}$

- Note that given  $c, \theta_1, \theta_2, \dots, \theta_q$ , it is possible to compute  $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_T$ .
- It is referred to as a **MA(q)** model, a **Moving Average (MA)** model of order  $q$ .

### Moving Average (MA) Models: Compute Errors

An example of an MA(1) Model  $y_t = c + \varepsilon_t + \theta_1 \varepsilon_{t-1}$

assumed to be given (e.g., by an Expectation-Maximization algorithm)

Time	Forecasted Values ( $\hat{y}_t$ )	Error at time t ( $\varepsilon_t$ )	Actual Values ( $y_t$ )
1	10	-1	9
2	$9.5 = 10 + 0.5(-1)$	2	11.5
3	$11 = 10 + 0.5(2)$	-1	10
4	$9.5 = 10 + 0.5(-1)$	2	10.5
5	$11 = 10 + 0.5(2)$	-1	10
6	$9 = 10 + 0.5(-1)$	0	9

$\varepsilon_t = y_t - (c + \theta_1 \varepsilon_{t-1})$  can be derived

### ARMA: Combine Autoregressive (AR) and Moving Average (MA) Models

$$y_t = c + \underbrace{\phi_1 y_{t-1} + \phi_2 y_{t-2} + \dots + \phi_p y_{t-p}}_{\text{Autoregressive (AR)}} + \underbrace{\theta_1 \varepsilon_{t-1} + \theta_2 \varepsilon_{t-2} + \dots + \theta_q \varepsilon_{t-q}}_{\text{Moving Average (MA)}} + \varepsilon_t$$

- Note that given  $c, \phi_1, \phi_2, \dots, \phi_p, \theta_1, \theta_2, \dots, \theta_q$ , it is possible to compute  $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_T$ .
- Use for example Expectation-Maximization (EM) algorithms

### ARIMA Models

$$y_t^I = c + \phi_1 y_{t-1}^I + \phi_2 y_{t-2}^I + \dots + \phi_p y_{t-p}^I + \theta_1 \varepsilon_{t-1} + \theta_2 \varepsilon_{t-2} + \dots + \theta_q \varepsilon_{t-q} + \varepsilon_t$$

- Combination of the elements present before
  - AR**: Autoregressive (lagged values)
  - I**: Integrated (differencing to make series stationary)
  - MA**: Moving Average (lagged errors)
- Hence **ARIMA** is **ARMA** with differencing
- Parameters:  $ARIMA(p, d, q)$ 
  - $p$  = order of the autoregressive part
  - $d$  = number of times for differencing
  - $q$  = order of the moving average part
- Different experimental strategies are possible to find suitable  $p, d$ , and  $q$ . (out of scope)
- Note that  $y_t$  can be reconstructed from  $y_t^I$  and past values.

Goal:

$$\text{Minimize } \frac{1}{T} \sum_{t=1}^T \varepsilon_t^2$$

or

$$\frac{1}{T} \sum_{t=1}^T |\varepsilon_t|$$

### Key takeaways:

- AR/MA/ARMA model dependence; ARIMA handles non-stationary series via differencing.
- Always evaluate with time-ordered splits and forecast error metrics.