

# Machine Learning

## Association Rules – Examples

Claudio Sartori

DISI

Department of Computer Science and Engineering – University of Bologna, Italy

[claudio.sartori@unibo.it](mailto:claudio.sartori@unibo.it)

# Overview

- **Healthcare**: clinical decision support
- **Cybersecurity**: intrusion detection
- **Bioinformatics**: gene expression patterns
- **Public health**: outbreak and risk factor analysis
- **Manufacturing**: fault diagnosis and maintenance
- **Education**: learning analytics
- **Smart cities**: traffic and incident analysis

1	Healthcare	3
2	Cybersecurity	9
3	Bioinformatics	16
4	Public Health	20
5	Manufacturing	24
6	Education	28
7	Smart Cities	32

# Healthcare: Clinical Decision Support

## Context

- Hospitals and research institutes use association rules to find co-occurring diagnoses, treatments, symptoms, and lab results, such as hypertension  $\Rightarrow$  diabetes
- Supports improved diagnosis, risk prediction, and treatment protocols.

## Data structures

- Transactions: individual patient visits.
- Items: diagnoses, medications, lab abnormalities, symptoms, procedures.

# Healthcare - Data & Preprocessing

## Data acquisition

- Electronic Health Records (EHRs).
- ICD diagnosis codes.
- Medication lists.
- Lab results.

## Preprocessing tasks

- Remove personally identifiable information.
- Standardize ICD codes (for example, ICD-10).
- Discretize continuous lab values into categorical bins (for example, Glucose\_High).
- Filter noisy or rare events.

# Healthcare - Postprocessing & Outcomes

## Postprocessing

- Remove medically implausible or coincidental patterns.
- Rank rules by confidence, lift, and leverage.
- Validate rules with clinician review.

## Actionable outcomes

- Identification of patient groups at elevated risk.
- Improved triage and screening guidelines.
- Medication interaction warnings.
- Data-driven refinement of clinical pathways.

# Healthcare: Typical Questions

- Which diagnoses frequently appear together?  
 $\{\text{Hypertension}\} \Rightarrow \{\text{Chronic\_Kidney\_Disease}\}$
- What medication combinations tend to follow certain diagnoses?  
 $\{\text{Type\_II\_Diabetes}, \text{ Obesity}\} \Rightarrow \{\text{Metformin}\}$
- Which symptoms strongly predict a future diagnosis?  
 $\{\text{Night\_Sweats}, \text{ Weight\_Loss}\} \Rightarrow \{\text{Tuberculosis}\}$
- Are there unexpected adverse drug combinations?  
 $\{\text{Drug\_A}, \text{ Drug\_B}\} \Rightarrow \{\text{Abnormal\_Liver\_Enzymes}\}$

# Healthcare - Patient Record Transactions

## Transaction T101 (Patient Visit)

-----  
Diagnosis: [Hypertension]  
Diagnosis: [Obesity]  
Medication: [Metformin]  
Symptom: [Fatigue]  
Lab Result: [Glucose\_High]

-----

T101 = {Hypertension, Obesity, Metformin,  
Fatigue, Glucose\_High}

## Transaction T102

-----  
Diagnosis: [Asthma]  
Medication: [Inhaler]  
Symptom: [Wheezing]  
Environmental: [Pollen\_High]

-----

T102 = {Asthma, Inhaler, Wheezing,  
Pollen\_High}

1	Healthcare	3
2	Cybersecurity	9
3	Bioinformatics	16
4	Public Health	20
5	Manufacturing	24
6	Education	28
7	Smart Cities	32

# Cybersecurity: Intrusion Detection

## Context

- Security operations centers use association rules to detect **patterns of suspicious behavior**.
- Focus on co-occurring events across logs and emerging attack sequences.

## Data structures

- Transactions: sessions, sequences of log events, or daily logs.
- Items: event types, IP categories, resources accessed, anomalies.

# Cybersecurity - Data & Preprocessing

## Data acquisition

- SIEM log events.
- Firewall logs.
- Authentication logs.
- System event feeds.

## Preprocessing

- Normalize timestamps.
- Map raw events to categorical labels (for example, port\_scan\_detected).
- Remove duplicate or irrelevant logs.
- Aggregate logs into session windows (for example, per user per hour).

# Meaning of SIEM Log Events

## Definition

- SIEM = Security Information and Event Management.
- SIEM log events are security-relevant records aggregated, normalized, and analyzed by a SIEM platform.

## Sources

- OS logs (Windows Event Logs, syslog)
- Network devices (routers, firewalls)
- Security tools (IDS, IPS, antivirus, EDR)
- Authentication systems (AD, LDAP)
- Cloud logs (AWS CloudTrail, Azure)

## Examples

- Failed logins
- Privilege escalation attempts
- Firewall rule violations
- Malware detection events
- Port scans, unusual traffic

## Purpose

- Detect threats via correlation
- Identify anomalies
- Support forensic investigations
- Enable compliance reporting

# Cybersecurity - Postprocessing & Outcomes

## Postprocessing

- Filter rules with lift  $> 1.5$  to reduce false positives.
- Cross-check with known MITRE ATT&CK techniques.
- Validate rules with security experts.

## Actionable outcomes

- Early alerts for suspicious event combinations.
- Improved threat signatures.
- Automated risk scoring.
- Prioritization of machines for forensic review.

# Cybersecurity: Typical Questions

- Which event combinations precede a confirmed intrusion?  
 $\{\text{Multi\_Fail\_Login}, \text{Unusual\_Time\_Access}\} \Rightarrow \{\text{Unauthorized\_Access}\}$
- Which patterns of behavior distinguish normal from anomalous activity?
- Which attack chains co-occur across different machines?  
 $\{\text{Port\_Scan}, \text{SMB\_Exploit}\} \Rightarrow \{\text{Ransomware\_Deployment}\}$
- Which user or device profiles correlate with higher breach likelihood?

# Cybersecurity - Event Log Transactions

## Transaction S301 (User Session)

---

Event: [Failed\_Login]  
Event: [VPN\_Login]  
Event: [ privilege\_escalation ]  
Resource: [Admin\_Panel]  
Time: [Unusual\_Time]

---

S301 = {Failed\_Login, VPN\_Login,  
    Privilege\_Escalation,  
    Access\_Admin\_Panel,  
    Unusual\_Time}

## Transaction S302

---

Event: [Port\_Scan]  
Event: [SMB\_Exploit]  
Event: [File\_Encryption]

---

S302 = {Port\_Scan, SMB\_Exploit,  
    File\_Encryption}

1	Healthcare	3
2	Cybersecurity	9
3	<b>Bioinformatics</b>	<b>16</b>
4	Public Health	20
5	Manufacturing	24
6	Education	28
7	Smart Cities	32

# Bioinformatics: Gene Expression

## Context

- Association rules reveal relationships among gene expressions, protein interactions, and pathways.
- Applied to large genomic or transcriptomic datasets.

## Data structures

- Transactions: samples, experiments, expression profiles.
- Items: gene up/down states, protein interactions, pathway activations.

# Bioinformatics: Typical Questions

- Which sets of genes are **co-expressed** under certain conditions?  
 $\{\text{Gene\_A\_up}\} \Rightarrow \{\text{Gene\_B\_up}\}$
- Which expression patterns **predict disease phenotypes**?  
 $\{\text{Gene\_X\_up}, \text{Gene\_Y\_down}\} \Rightarrow \{\text{Tumor\_Aggressive}\}$
- What protein interaction chains commonly appear together?
- Which pathways are **co-activated** in specific cancers?

# Bioinformatics - Gene Expression Profiles

Transaction G12 (Tumor Sample)

-----  
Gene\_A: [Upregulated]  
Gene\_B: [Downregulated]  
Gene\_C: [Upregulated]  
Protein\_X: [Interaction\_Active]

-----

G12 = {Gene\_A\_up, Gene\_B\_down,  
Gene\_C\_up, Protein\_X\_interact}

Transaction G13

-----  
Gene\_D: [Upregulated]  
Pathway\_Y: [Activated]

-----

G13 = {Gene\_D\_up, Pathway\_Y\_active}

1	Healthcare	3
2	Cybersecurity	9
3	Bioinformatics	16
4	Public Health	20
5	Manufacturing	24
6	Education	28
7	Smart Cities	32

# Public Health: Outbreak and Risk Factors

## Context

- Public health agencies mine **co-occurring symptoms, conditions, and social determinants.**
- Goal: understand and anticipate disease outbreaks.

## Data structures

- Transactions: individual case reports or region-day aggregates.
- Items: symptoms, demographics, exposures, environmental conditions.

# Public Health: Typical Questions

- Which **symptom clusters** strongly indicate a specific disease?  
{Rash, Fever} ⇒ {Measles}
- Which environmental conditions co-occur with **disease spikes**?  
{High\_Humidity, Standing\_Water} ⇒ {Dengue\_Outbreak}
- Which **risk factors** tend to appear together in severe cases?  
{Smoking, Air\_Pollution} ⇒  
{Severe\_Respiratory\_Issues}
- Which combinations of travel history and symptoms predict imported cases?

# Public Health - Epidemiology Case Transactions

## Transaction C901 (Case Report)

Symptom: [Fever]

Symptom: [Rash]

Exposure: [Travel\_Region\_X]

Demographic: [Child]

C901 = {Fever, Rash, Travel\_Region\_X, Child}

## Transaction C902

Symptom: [Cough]

Environment: [Air\_Pollution\_High]

Behavior: [Smoking]

C902 = {Cough, Air\_Pollution\_High, Smoking}

1	Healthcare	3
2	Cybersecurity	9
3	Bioinformatics	16
4	Public Health	20
5	Manufacturing	24
6	Education	28
7	Smart Cities	32

# Manufacturing: Fault Diagnosis

## Context

- Factories apply association rules to uncover failure patterns.
- Supports predictive maintenance and reliability engineering.

## Data structures

- Transactions: machine cycles, daily logs, fault incidents.
- Items: sensor anomalies, part replacements, error codes, vibration spikes.

# Manufacturing: Typical Questions

- Which **sensor readings** in combination precede a specific failure?  
 $\{\text{Temp\_High}, \text{Vibration\_High}\} \Rightarrow \{\text{Bearing\_Failure}\}$
- What **component co-failures** frequently occur together?  
 $\{\text{Pump\_Failure}\} \Rightarrow \{\text{Valve\_Replacement}\}$
- Which maintenance actions tend to **resolve related anomalies**?
- Which operational conditions correlate with **reduced lifespan**?

# Manufacturing - Machine Cycle Snapshot

## Transaction M203 (Machine Cycle)

[Sensor] Temp\_High  
[Sensor] Vib\_Spike  
[Error Code] E17  
[Maintenance] None

M203 = {Temp\_High, Vib\_Spike, Error\_E17}

## Transaction M204

[Sensor] Noise\_High  
[Repair] Bearing\_Replace

M204 = {Noise\_High, Bearing\_Replace}

1	Healthcare	3
2	Cybersecurity	9
3	Bioinformatics	16
4	Public Health	20
5	Manufacturing	24
6	Education	28
7	Smart Cities	32

# Education: Learning Analytics

## Context

- Universities analyze behavior patterns, resource usage, and outcomes using association rules.
- Aim: improve learning design and early-warning systems.

## Data structures

- Transactions: a student term, course session, or weekly activity.
- Items: actions (viewed lectures, submissions), skills, quiz scores.

# Education: Typical Questions

- Which behaviors predict **high performance**?  
 $\{ \text{Early\_Assignment\_Submission}, \text{Forum\_Participation} \} \Rightarrow \{\text{High\_Grade}\}$
- Which behaviors signal **dropout risk**?  
 $\{ \text{No\_Logins\_Week2}, \text{Missed\_Quiz\_1} \} \Rightarrow \{\text{Dropout}\}$
- Which **learning resources** are often used together?
- How do specific misconceptions co-occur across assessments?

# Education - Student Interaction Log

## Transaction A012 (Week 2 Behavior)

[Action] View\_Lecture\_3  
[Action] View\_Lecture\_4  
[Quiz] Quiz\_1\_Attempted  
[Performance] Low\_Score

A012 = {View\_Lecture\_3, View\_Lecture\_4,  
Quiz\_1\_Attempted, Low\_Score}

## Transaction A013

[Action] Forum\_Post  
[Action] Early\_Submission  
[Outcome] High\_Grade

A013 = {Forum\_Post, Early\_Submission, High\_Grade}

1	Healthcare	3
2	Cybersecurity	9
3	Bioinformatics	16
4	Public Health	20
5	Manufacturing	24
6	Education	28
7	Smart Cities	32

# Smart Cities: Traffic Patterns

## Context

- Urban planners use association rules on IoT and traffic data to find patterns leading to congestion or accidents.

## Data structures

- Transactions: road-segment time windows (for example, 5 minutes).
- Items: traffic volume, weather, accidents, low speed, congestion.

# Smart Transportation - Data & Preprocessing

## Data acquisition

- Roadside IoT sensors.
- Weather stations.
- GPS and speed detectors.
- Incident reports.

## Preprocessing

- Synchronize time windows across sensors.
- Convert continuous values to categories (for example, Speed\_Low, Density\_High).
- Handle missing sensor data.
- Aggregate into fixed windows (for example, 5-minute intervals).

# Pipeline 3: Smart Transportation - Postprocessing & Outcomes

## Postprocessing

- Rank rules by lift to detect strongest indicators.
- Validate with historical crash data.
- Cluster rules by road type (for example, highway vs urban).

## Actionable outcomes

- Dynamic speed-limit recommendations.
- Accident probability dashboards.
- Better placement of signage and sensors.
- Predictive alerts for drivers.

# Smart Cities: Typical Questions

- Which **conditions co-occur** immediately before accidents?  
 $\{Rain, Low\_Speed, High\_Density\} \Rightarrow \{Collision\}$
- What combinations of road conditions cause **predictable congestion**?  
 $\{Construction, Lane\_Closure\} \Rightarrow \{Severe\_Delay\}$
- How do weather patterns influence traffic flow?
- Which intersections exhibit **recurring joint anomalies**?

# Smart Cities - Road Segment Events

Transaction T550 (Segment: Highway-12)

---

Weather: [Rain]  
Traffic: [High\_Density]  
Speed: [Below\_30]  
Event: [Accident]

---

T550 = {Rain, High\_Density, Speed\_Low, Accident}

Transaction T551

---

Weather: [Clear]  
Traffic: [Moderate]  
Event: [No\_Accident]

---

T551 = {Clear, Moderate\_Traffic, No\_Accident}

# Wrap-up

- Association rules apply far beyond commerce and marketing.
- Common pattern: transactions of events + items as discrete attributes.
- Output: interpretable rules that support decision making and policy.