

## Data minig

### Limiti di OLAP

Data l'elevata quantità e la complessità delle relazioni, le informazioni non sono completamente identificabili.

Si predispone una situazione in cui i dati rimangono inutilizzati o sotto-utilizzati (molti dati ma poche informazioni).

Gli strumenti OLAP non sono più sufficienti perchè operano per supportare processi decisionali, quindi sviluppano percorsi di analisi da ipotesi che è l'utente a formulare.

### Fasi del processo di mining

Il data mining è spesso definito KDD (Knowledge Discovery in Databases) e viene diviso nelle seguenti fasi elementari:

- **Pulizia dei dati:** vengono eliminate le incorrettezze.
- **Integrazione dei dati:** uniformare i dati.
- **Selezione dei dati.**
- **Trasformazione dei dati:** riorganizzazione dei dati.
- **Data mining:** il vero e proprio processo di analisi.
- **Valutazione dei pattern:** l'insieme delle condizioni viene ridotto a quelle interessanti.
- **Presentazione delle conoscenze.**

Le prime fasi coincidono con il popolamento del DWH, può essere considerato un'evoluzione delle indagini OLAP.

### Da OLAP a OLAM

OLAM (On Line Analytical Mining)

Partendo dai DWH abbiamo dati ben strutturati, puliti e completi. Ciononostante il processo di mining non può essere interamente automatico, infatti i pattern rilevati potrebbero essere troppi e non interessanti, il processo di data mining deve quindi essere interattivo con gli utenti che specificano la direzione in cui indagare.

Un processo interattivo permette di affinare le ricerche.

### Architettura dei sistemi di Data Mining

L'architettura dei sistemi di data mining si appoggia ai seguenti componenti:

- **Data warehouse**
- **Base di conoscenza (knowledge base):** l'insieme di regole e conoscenze note, verranno utilizzate per guidare le ricerche.
- **Motore di data mining (data mining engine):** l'insieme delle funzioni di analisi dei dati.

- **Valutazione delle condizioni (pattern evaluation):** i moduli che fanno focalizzare la ricerca sulle condizioni interessanti.
- **Sistema di presentazione:** l'interfaccia con la quale l'utente fa le ricerche.

#### I 4 principi di analisi



#### Statistiche elementari e analisi relative

##### Generalizzazione

Deve fornire una visione ad alto livello tramite l'accorpamento di concetti e riassumendo caratteristiche di base.

Il principio di base è che gli elementi che un utente può analizzare devono essere un numero limitato.

Un diffuso tipo di generalizzazione è l'aggregazione dei sistemi OLAP, i sistemi di data mining amplificano il potenziale mettendo a disposizione anche delle metodologie di induzione.

##### 1. Caratterizzazione

Serve a comprendere le caratteristiche di una classe, che siano queste (caratteristiche) di tendenza o di dispersione.

Viene spesso rappresentata con tabelle, grafici e boxplot.

##### 2. Discriminazione

Con questa modalità invece, le caratteristiche di una classe vengono messe a confronto con quelle di un'altra classe ad essa paragonabile.

Viene quindi eseguito un confronto diretto sulle tabelle o sui grafici.

## Analisi associative

Meccanismi che permettono di identificare situazioni che ne implicano altre con un'elevata frequenza.

Devono essere individuati pattern che rappresentano implicazioni logiche come  $A \rightarrow B$ .

La significatività di un'associazione viene definita con due parametri:

1. **Confianza:** misura la certezza di un pattern, è definita come  $P(A|B)$ .
2. **Supporto:** la frequenza con cui il pattern è stato verificato nel DB, è definito come la percentuale degli elementi che verifica la regola.