

# Introduzione ai Data Warehouse

## Fondamenti, Architettura e Applicazioni

Un viaggio nel mondo dei data warehouse: dalla teoria alla pratica, scoprirete come trasformare dati operazionali in insight strategici per il business.

**Unità IA.4** - Analisi Descrittive usando Data Warehouse

Corso: Data Analyst e AI per Programmatori

Lezione 2

## Cos'è un Data Warehouse?

Un Data Warehouse è un sistema centralizzato che aggrega e memorizza grandi volumi di dati storici e operazionali provenienti da diverse fonti. La sua finalità principale è supportare l'analisi, il reporting e la Business Intelligence, fornendo una visione consolidata e coerente per decisioni strategiche.

A differenza dei database transazionali, ottimizzati per le operazioni quotidiane (OLTP), il Data Warehouse è progettato per interrogazioni complesse e l'analisi dei dati (OLAP), garantendo rapidità nell'estrazione di informazioni significative.

---

## Esempio Concreto: Retail Company

Immagina una grande catena di negozi al dettaglio con centinaia di punti vendita e un e-commerce. Ogni giorno, generano enormi quantità di dati: vendite per prodotto, per negozio, dati dei clienti, inventario, promozioni attive, ecc.

Senza un Data Warehouse, analizzare queste informazioni sarebbe quasi impossibile. Il Data Warehouse raccoglie questi dati da tutti i sistemi (POS, CRM, ERP, e-commerce), li pulisce, li integra e li organizza in un formato che permette ai manager di rispondere a domande come:

Quali sono i prodotti più venduti in ogni regione nell'ultimo trimestre?

Qual è l'impatto di una specifica promozione sulle vendite online rispetto ai negozi fisici?

Quali sono i segmenti di clienti più profittevoli e quali prodotti acquistano più frequentemente?



# Perché Serve un Data Warehouse?

I database operazionali (OLTP) sono ottimizzati per gestire le transazioni quotidiane, ma presentano limiti significativi quando si tratta di condurre analisi complesse e supportare decisioni strategiche. Il data warehouse nasce per risolvere queste sfide, offrendo un ambiente dedicato all'analisi dei dati.

## Problemi dei Database Operazionali



### Performance

Query analitiche complesse rallentano il sistema operativo, impattando negativamente le transazioni quotidiane e la produttività aziendale.

### Dati Frammentati

Informazioni sparse in sistemi diversi (CRM, ERP, e-commerce) rendono difficile ottenere una vista unificata del business.

### Mancanza di Storico

I dati vecchi vengono eliminati o archiviati, impedendo analisi di trend a lungo termine e previsioni accurate.

### Inconsistenza

Definizioni diverse degli stessi concetti in sistemi diversi complicano l'integrazione e generano confusione nelle analisi.

## Vantaggi del Data Warehouse



### Separazione Operazioni/Analisi

Sistema separato: le query analitiche non impattano le performance operative, garantendo continuità al business.

### Vista Unificata

Tutti i dati integrati in un unico repository con formati e definizioni coerenti, eliminando ambiguità e duplicazioni.

### Storico Completo

Conserva anni di dati storici per analisi di trend, pattern temporali e supporto alle decisioni basate su evidenze storiche.

### Ottimizzazione per Query

Strutture dati ottimizzate (indici, aggregazioni pre-calcolate, partizionamento) garantiscono risposte rapide anche su grandi volumi.



# Data Warehouse in Azione: Un Esempio Concreto

Per comprendere il valore di un Data Warehouse, confrontiamo due scenari in un settore competitivo come la vendita di elettronica. Vediamo come la gestione dei dati possa fare la differenza tra il successo e il declino.



## Negozio A: "ElettroStop" (Database Tradizionale)

ElettroStop si affida a database operazionali per gestire le transazioni quotidiane. Le informazioni su vendite, magazzino e feedback dei clienti sono frammentate e non integrate. Non esiste un archivio storico dettagliato a lungo termine.

Quando le vendite iniziano a calare, il team marketing vede solo i numeri aggregati. Non riesce a capire il perché: è un calo stagionale? La concorrenza? Non ha gli strumenti per analizzare i trend emergenti o le preferenze specifiche dei clienti.

La reazione è lenta e spesso basata su intuizioni: si fanno promozioni generiche, si ordinano prodotti "di successo" del passato, ma il declino continua perché le vere cause restano nascoste.

## Negozio B: "TechZone" (Con Data Warehouse)

TechZone ha implementato un Data Warehouse che raccoglie e integra dati da tutte le fonti: vendite al dettaglio, e-commerce, resi, interazioni con il servizio clienti e persino menzioni sui social media. I dati storici sono conservati e strutturati per l'analisi.

Appena le vendite mostrano un leggero calo, il Data Warehouse rivela rapidamente un pattern: un aumento delle richieste per nuovi modelli di smartphone e smart device non ancora disponibili nel loro inventario, ma che i clienti cercano attivamente.

Grazie a questa intuizione basata sui dati, TechZone reagisce in pochi giorni: ordina i nuovi prodotti richiesti, lancia una campagna di pre-ordine e riorganizza l'esposizione. Le vendite riprendono, evitando perdite significative e rafforzando la fedeltà dei clienti.



# OLTP vs OLAP: Due Mondi Diversi

Comprendere la differenza tra sistemi OLTP e OLAP è fondamentale per progettare architetture dati efficaci. Questi due paradigmi servono scopi complementari ma profondamente diversi nell'ecosistema informativo aziendale.

## OLTP Online Transaction Processing

Sistemi progettati per gestire **transazioni operative quotidiane**: inserimenti, aggiornamenti e cancellazioni veloci di singoli record. Ottimizzati per velocità e affidabilità nelle operazioni di routine come elaborazione ordini, registrazioni clienti e aggiornamenti inventario.

## OLAP Online Analytical Processing

Sistemi progettati per supportare **analisi complesse e decisioni strategiche**: query aggregate su grandi volumi di dati storici. Ottimizzati per flessibilità analitica, permettono di esplorare dati da diverse prospettive e generare insight di business.

Parameters	OLTP	OLAP
Purpose	It is a system for processing large volumes of real-time transactional data.	It is a system for the multidimensional analysis of consolidated business data.
Usage	It is used for adding, deleting, or updating databases to keep the data up-to-date.	It is used to make business decisions through queries and complex analyses of large amounts of data.
Focus	The system is more focused on transactional data maintenance and less on data analysis.	The system is focused on data analysis and not on maintaining day-to-day transactions.
Data Source	OLTP sources data from traditional database management systems.	OLAP has multiple data sources, which include real-time and historical databases, including OLTP.
Data Type	The data consists of a large number of short transactions.	The system processes large volumes of data from multiple sources.
Processing Time	Very low processing time at the scale of a few milliseconds.	Depending on the query, processing time is not as fast as OLTP systems and may range from a few seconds to hours.
Query	Related to adding, deleting, and updating data.	Related to data analysis.
Availability	OLTP systems are available round-the-clock and updated frequently to maintain data integrity.	OLAP systems don't need to be updated so frequently since their functions are analytic in nature.
Normalization	Data tables are normalized.	Data tables are not normalized.
Backup	Requires constant backup and recovery.	Can be backed up less frequently.
User volume	Supports large user volume simultaneously.	Accommodates multiple users but doesn't have a large user volume like OLTP.
Operations	Allows both read and write operations.	Usually supports read-only operations.
Process	Processes day-to-day data quickly.	Processes analytical queries consistently and at a fast pace.



# OLTP vs OLAP: Esempi Reali

Per illustrare chiaramente la differenza tra OLTP e OLAP, esaminiamo come questi due tipi di sistemi gestiscono i dati in contesti aziendali comuni, ognuno con il suo scopo e le sue ottimizzazioni specifiche.



## OLTP: Transazioni Operative Quotidiane

I sistemi OLTP sono ottimizzati per elaborare rapidamente un gran numero di transazioni individuali, garantendo integrità e disponibilità costante.

- **Banca:** Un cliente preleva denaro da un ATM, un bonifico viene registrato o un pagamento con carta di credito viene autorizzato in tempo reale.
- **E-commerce:** Un utente aggiunge un prodotto al carrello, completa un ordine o aggiorna il proprio indirizzo di spedizione.
- **Ospedale:** Un medico registra l'anamnesi di un paziente, un infermiere aggiorna la somministrazione di un farmaco o si programma un appuntamento.

## OLAP: Analisi e Decisioni Strategiche

I sistemi OLAP sono progettati per l'analisi complessa di grandi volumi di dati storici, permettendo di estrarre insight per la pianificazione strategica.

- **Banca:** L'ufficio marketing analizza i trend di spesa dei clienti per identificare nuovi segmenti di mercato o valutare l'efficacia delle promozioni.
- **E-commerce:** Il team di gestione analizza le vendite per categoria di prodotto nell'ultimo anno per prevedere la domanda o ottimizzare l'inventario.
- **Ospedale:** L'amministrazione analizza i tassi di successo di trattamenti specifici, i costi per paziente o l'efficienza dei reparti.

# OLTP vs OLAP: Confronto Dettagliato

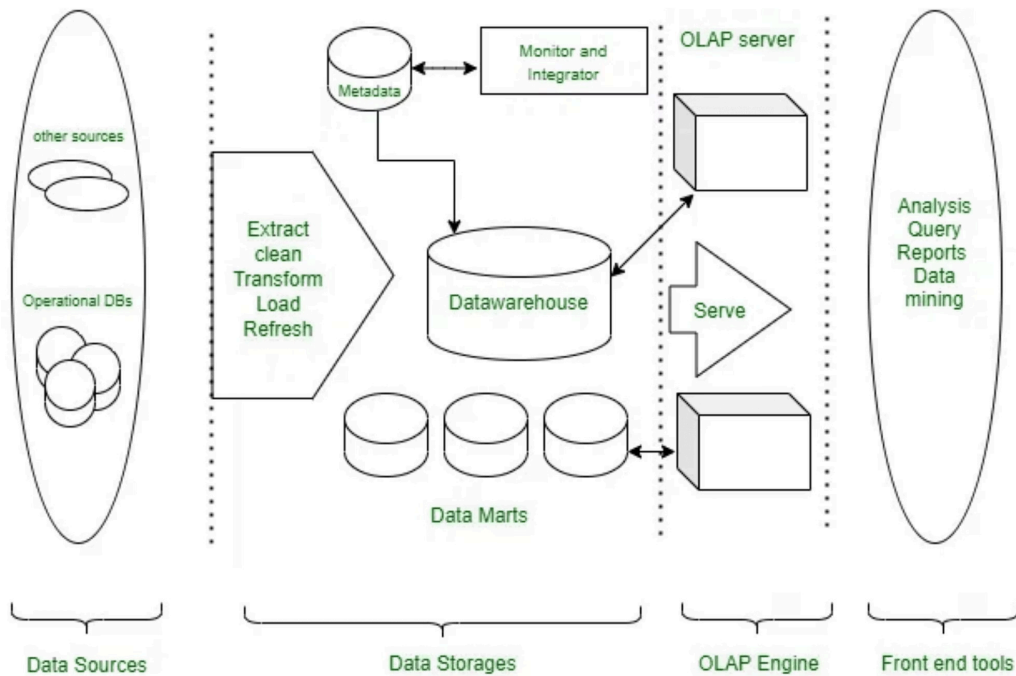
Questa tabella sintetizza le differenze chiave tra i due paradigmi, evidenziando come ciascuno sia ottimizzato per esigenze specifiche e complementari nell'architettura dati aziendale.

Caratteristica	OLTP (Transazionale)	OLAP (Analitico)
Scopo	Operazioni quotidiane e gestione transazioni	Analisi e decisioni strategiche di business
Tipo di Query	Semplici, predefinite (SELECT, INSERT, UPDATE)	Complesse, ad-hoc con aggregazioni e JOIN multipli
Operazioni	INSERT, UPDATE, DELETE frequenti	SELECT con JOIN, GROUP BY, funzioni analitiche
Tempo Risposta	Millisecondi (requisito critico)	Secondi o minuti (accettabile per analisi)
Utenti	Molti (100-1000+ concorrenti)	Pochi (10-100, principalmente analisti)
Volume Dati	GB, dati correnti	TB o PB, dati storici multi-anno
Schema	Normalizzato (3NF) per evitare ridondanza	Denormalizzato (Star/Snowflake) per performance
Aggiornamenti	Frequenti e in tempo reale	Periodici (batch notturni o settimanali)



# Architettura Three-Tier del Data Warehouse

L'architettura a tre livelli separa logicamente le responsabilità del sistema, garantendo scalabilità, manutenibilità e performance ottimali. Ogni tier ha un ruolo specifico nel flusso dei dati dalle sorgenti agli utenti finali.



1

## Bottom Tier

### Data Sources

Il livello delle **sorgenti dati eterogenee** comprende tutti i sistemi che generano o contengono informazioni aziendali: database operazionali OLTP (PostgreSQL, MySQL, Oracle), file flat in vari formati (CSV, JSON, XML), API esterne di servizi cloud, log di sistema e applicazioni, sensori IoT, feed da social media e piattaforme web.

2

## Middle Tier

### Data Warehouse Core

Il **cuore del sistema** include il database del data warehouse con schema dimensionale ottimizzato (Star o Snowflake), server ETL per l'orchestrazione dei processi di trasformazione, metadata repository che documenta strutture e lineage dei dati, aggregazioni pre-calcolate per accelerare query frequenti, e indici ottimizzati per pattern di accesso analitici.

3

## Top Tier

### Presentation / Analytics

Il livello di **presentazione e analisi** fornisce strumenti per consumare i dati: BI tools come Metabase, Tableau, Power BI per visualizzazioni interattive, dashboard operative e strategic, report standardizzati e ad-hoc, OLAP tools per analisi multidimensionali (drill-down, slice-and-dice), API REST per accesso programmatico e integrazione con applicazioni personalizzate.



# Architettura Concreta di un Data Warehouse

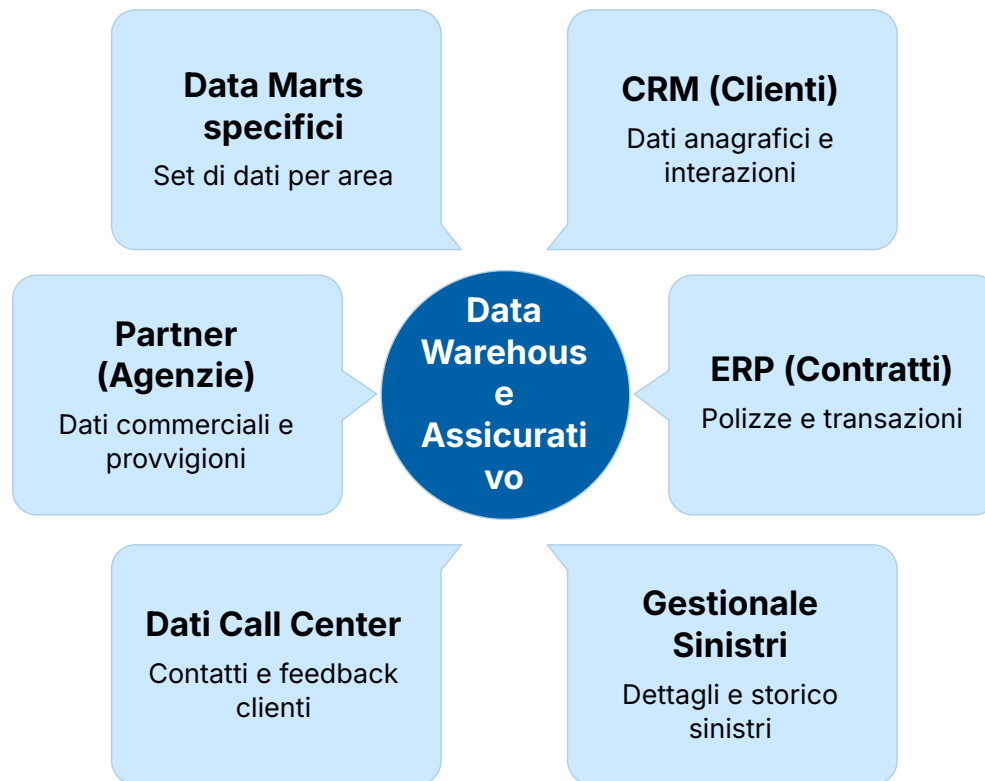
Visualizziamo il percorso dei dati all'interno di un Data Warehouse, dal loro stato grezzo nei sistemi operazionali fino alla loro forma raffinata e pronta per l'analisi. Questo flusso mostra come i dati vengono trasformati, aggregati e presentati per generare valore.



Come osservato nel diagramma, il volume dei dati diminuisce progressivamente ad ogni fase, mentre la loro qualità e il loro valore informativo aumentano. Le frecce rappresentano questo processo: **grezze e voluminose (rosse)** dalle sorgenti, **raffinate durante la trasformazione (gialle)**, fino a diventare **concentrate e pronte per l'analisi (verdi)** negli strumenti di Business Intelligence.

# Architettura di un Data Warehouse per una Compagnia Assicurativa

Un Data Warehouse per una compagnia assicurativa integra dati provenienti da sistemi eterogenei, offrendo una visione unificata e storica essenziale per l'analisi strategica, la gestione del rischio e l'ottimizzazione delle operazioni.



Questo schema evidenzia come i dati grezzi da diverse fonti operative vengano consolidati, puliti e trasformati nel Data Warehouse, per poi essere distribuiti in Data Marts specializzati e utilizzati da strumenti di Business Intelligence per supportare decisioni informate e strategiche in ogni settore dell'azienda.

# Il Processo ETL

## Extract - Transform - Load: Il Cuore del Data Warehouse

L'ETL è il processo critico che alimenta il data warehouse, trasformando dati grezzi e frammentati in informazioni strutturate e affidabili. Ogni fase ha responsabilità specifiche e tecniche dedicate per garantire qualità e integrità dei dati.



### Extract

#### Estrazione

- Connessione a sorgenti multiple (database, API, file)
- Lettura dati operazionali (OLTP)
- Estrazione incrementale o completa
- Gestione formati diversi (CSV, JSON, XML)

L'estrazione può essere **full** (completa) o **incremental** (solo cambiamenti), schedulata periodicamente o triggered da eventi.



### Transform

#### Trasformazione

- Pulizia dati (duplicati, valori mancanti, outlier)
- Standardizzazione formati (date, valute, unità)
- Integrazione da sorgenti diverse
- Aggregazioni e calcoli derivati
- Validazione qualità dati

Questa fase è la più complessa e richiede logica di business per risolvere inconsistenze e arricchire i dati.



### Load

#### Caricamento

- Caricamento nel Data Warehouse
- Aggiornamento Fact Tables
- Gestione Dimension Tables (SCD)
- Creazione indici e aggregazioni
- Verifica integrità referenziale

Il caricamento può essere **full refresh** o **incremental**, con strategie diverse per gestire cambiamenti storici.



# ETL in Azione: L'Acquisizione di una Compagnia

In caso di acquisizione, l'integrazione rapida dei dati è cruciale. Per superare la complessità e la lunghezza dei tempi di fusione dei sistemi operativi, si ricorre a un processo ETL notturno che consolida le informazioni della nuova entità, rendendole disponibili per l'analisi strategica fin dal giorno seguente.



## Estrazione Quotidiana

Ogni sera, un job ETL automatizzato si connette ai sistemi operazionali (OLTP) della compagnia acquisita. Vengono identificati ed estratti solo i dati nuovi o modificati dal giorno precedente, minimizzando il carico sui sistemi sorgente e ottimizzando le performance.



## Caricamento nel Data Warehouse

I dati, ora puliti e armonizzati, vengono caricati nelle rispettive tabelle del Data Warehouse principale. Questo processo aggiorna le Fact Tables (es. polizze emesse, sinistri) e le Dimension Tables (es. clienti, prodotti), rendendo le nuove informazioni parte integrante del patrimonio dati aziendale.



## Trasformazione e Armonizzazione

I dati estratti subiscono un'intensa fase di trasformazione. Questo include la pulizia da errori, la standardizzazione dei formati (es. date, valute), la risoluzione di duplicati e l'applicazione di regole di business per allineare le informazioni al modello dati del Data Warehouse centrale.



## Validazione e Disponibilità

Dopo il caricamento, vengono eseguiti controlli automatici di validazione per assicurare l'integrità e la coerenza dei dati. Le informazioni della compagnia acquisita sono ora integrate e pronte per essere analizzate tramite gli strumenti di Business Intelligence, fornendo una visione unificata e aggiornata.

Questo approccio permette una visione consolidata e quasi in tempo reale delle operazioni post-acquisizione, supportando decisioni strategiche tempestive.

# Data Mart, Data Lake, Data Warehouse

Nell'ecosistema dei dati aziendali coesistono diverse soluzioni, ciascuna progettata per esigenze specifiche. Comprendere le differenze è essenziale per scegliere l'architettura giusta per ogni scenario di business.



## Data Warehouse

Repository **centralizzato** per dati **strutturati e puliti**, ottimizzato per analisi enterprise-wide e business intelligence. Utilizza schema dimensionale (Star/Snowflake), supporta query SQL complesse, e garantisce alta qualità dei dati attraverso processi ETL rigorosi.



## Data Mart

**Sottoinsieme del DW** focalizzato su un'area specifica del business o dipartimento (Marketing, Sales, Finance). Più agile e veloce da implementare, con dati pre-aggregati per esigenze specifiche. Può essere *dependent* (deriva dal DW) o *independent* (alimentato direttamente dalle sorgenti).



## Data Lake

Storage per **grandi volumi di dati grezzi** in **qualsiasi formato** (strutturati, semi-strutturati, non strutturati). Ideale per big data, machine learning e data science. Schema-on-read: la struttura viene applicata solo al momento della lettura. Tecnologie comuni: Hadoop, AWS S3, Azure Data Lake.

Caratteristica	Data Warehouse	Data Mart	Data Lake
Scopo	Analisi enterprise	Analisi dipartimentale	Storage big data, ML
Dati	Strutturati, puliti	Strutturati, puliti	Qualsiasi formato (raw)
Schema	Schema-on-write	Schema-on-write	Schema-on-read
Utenti	Analisti, manager	Team specifico	Data scientist, ML engineer
Costi	Alto (storage + ETL)	Medio	Basso (storage commodity)

# Processo di Implementazione di un Data Warehouse

L'implementazione di un data warehouse è un progetto complesso che richiede pianificazione metodica e coinvolgimento di stakeholder tecnici e business. Seguire un approccio strutturato garantisce successo e ROI.

01

## Analisi Requisiti

Intervistare stakeholder chiave, identificare KPI e metriche critiche, definire domande di business (es. "Quali prodotti generano più margine?"), identificare sorgenti dati disponibili e loro qualità, stimare volumi e frequenza di aggiornamento.

03

## Sviluppo

Creare database DW con DDL ottimizzati, implementare tabelle con indici e partizioni, sviluppare processi ETL (Python, SQL, Apache Airflow), creare viste e aggregazioni per query comuni, implementare logging e error handling.

05

## Deployment

Caricamento dati storici (backfill), configurazione scheduling ETL (cron, Airflow), training utenti finali su strumenti BI, documentazione completa (architettura, data dictionary, user guide), setup monitoraggio e alerting.

02

## Progettazione

Progettare schema dimensionale (Star o Snowflake), definire Fact Tables (metriche) e Dimension Tables (contesto), progettare processo ETL con mapping sorgente-destinazione, pianificare aggregazioni e viste materializzate, definire strategia per Slowly Changing Dimensions.

04

## Testing

Test qualità dati (completezza, accuratezza, consistenza), test performance query su volumi realistici, test processo ETL end-to-end con scenari edge case, validazione con utenti finali su query e report campione.

06

## Manutenzione

Monitoraggio continuo performance e qualità dati, ottimizzazione query lente, aggiunta nuove sorgenti dati e metriche, gestione crescita volume (archiving, partitioning), backup regolari e disaster recovery testing.



# Business Case del Corso

## Il Progetto che Svilupperemo Insieme

Nei prossimi moduli, costruiremo un data warehouse completo per **TechStore**, un e-commerce di elettronica in crescita. Questo case study realistico vi permetterà di applicare tutti i concetti teorici a un progetto pratico end-to-end.

### TechStore

E-commerce specializzato nella vendita di prodotti elettronici di alta qualità: laptop delle migliori marche, smartphone Android e iOS, tablet per lavoro e intrattenimento, accessori tech (cuffie, mouse, tastiere), e componenti hardware per appassionati di PC building. Opera online da 3 anni con crescita costante anno su anno.



100K+

Transazioni

Negli ultimi 2 anni

500+

Prodotti Attivi

Nel catalogo

20K+

Clienti Registrati

Base utenti

€5M+

Fatturato Annuale

Revenue 2023

### Obiettivi del Progetto Data Warehouse

#### Analisi Trend Vendite

Analizzare trend di vendita temporali (giornalieri, settimanali, mensili), performance prodotti per categoria e marca, seasonality e picchi di vendita.

#### Segmentazione Clienti

Segmentare clienti in base a comportamenti d'acquisto (RFM analysis), calcolare Customer Lifetime Value (CLV), identificare clienti ad alto valore.

#### Ottimizzazione Marketing

Valutare efficacia campagne marketing e ROI, analizzare conversion funnel e drop-off, ottimizzare pricing e promozioni.

#### Preparazione ML

Preparare dati puliti e strutturati per modelli predittivi: forecast vendite, recommendation systems, churn prediction.

# Riepilogo e Prossimi Passi

In questa prima lezione abbiamo gettato le fondamenta per comprendere il ruolo strategico dei data warehouse nell'architettura dati moderna. Ora siete pronti per approfondire la modellazione dimensionale.

## Concetti Chiave Appresi

1

### Data Warehouse

Repository orientato al soggetto, integrato, variante nel tempo e non volatile per supportare decisioni di business basate su dati storici.

2

### OLTP vs OLAP

Sistemi transazionali (operazioni quotidiane) vs sistemi analitici (decisioni strategiche) hanno scopi, architetture e caratteristiche di performance profondamente diverse.

3

### Architettura Three-Tier

Bottom Tier (sorgenti dati), Middle Tier (DW core + ETL), Top Tier (strumenti analitici e BI) garantiscono separazione di responsabilità e scalabilità.

4

### Processo ETL

Extract (estrazione da sorgenti), Transform (pulizia, standardizzazione, integrazione), Load (caricamento nel DW) è il cuore pulsante del data warehouse.

5

### Data Mart e Data Lake

Soluzioni complementari per esigenze dipartimentali (Data Mart) e big data non strutturati per ML e advanced analytics (Data Lake).

## Prossima Lezione



### Lezione 2 Modellazione Dimensionale

Approfondiremo le tecniche di data modeling per data warehouse:

- **Schema a Stella (Star Schema):** il modello più diffuso, ottimizzato per query performance
- **Fact Tables e Dimension Tables:** come strutturare metriche e contesto dimensionale
- **Schema a Fiocco di Neve (Snowflake):** normalizzazione delle dimensioni per ridurre ridondanza
- **Slowly Changing Dimensions (SCD):** gestire cambiamenti storici nelle dimensioni
- **Progettazione pratica per TechStore:** applicazione hands-on al nostro case study

### Esercitazione Pratica

Setup ambiente PostgreSQL locale e esplorazione del database operativo di TechStore. Familiarizzeremo con le tabelle sorgente prima di progettare il data warehouse.

