

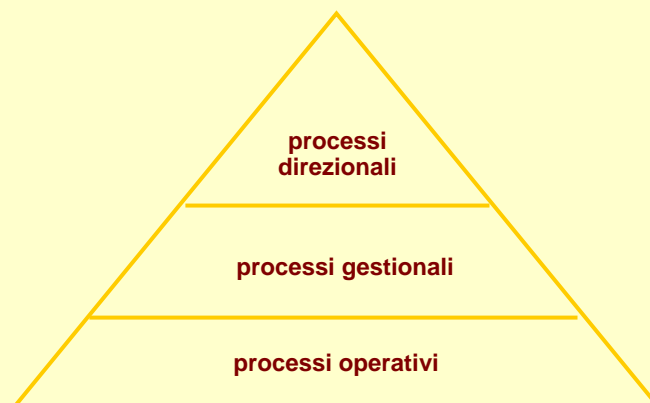
**Tecnologie per i sistemi
informativi**

**I data warehouse e la
loro progettazione**

Docente: Letizia Tanca
Politecnico di Milano
tanca@elet.polimi.it

1

Processi



2

Le query che vorremmo poter formulare ai livelli superiori

- Incassi registrati lo scorso anno per ciascuna regione e ciascuna categoria di prodotto
- Correlazione tra l'andamento dei titoli azionari dei produttori di computer e i profitti trimestrali degli ultimi 5 anni
- Quali sono gli ordini che massimizzano gli incassi?
- Quale di due nuove terapie risulterà in una diminuzione della durata media dei ricoveri?

3

Il problema

- Una promessa della tecnologia relazionale: "flexible data access":
 - Uno strumento per l'utente finale in cui tutte le query siano ugualmente formulabili
- Ma la tecnologia relazionale non ha mantenuto questa promessa:
 - Complessità e rigidità delle applicazioni
 - Enfasi su OLTP
- Le conseguenze:
 - Masse di dati per la gestione operativa
 - Scarso utilizzo dei dati per la gestione strategica

4

OLTP

Tradizionale elaborazione di transazioni, che realizzano i processi operativi dell'azienda-ente

- Operazioni spesso predefinite e relativamente semplici
- Ogni operazione coinvolge "pochi" dati
- Dati di dettaglio, aggiornati
- Le proprietà "acide" (atomicità, correttezza, isolamento, durabilità) delle transazioni sono essenziali

Le dimensioni delle basi di dati sono dell'ordine dei gbyte

La principale metrica di prestazione e' il throughput delle transazioni

5

OLAP

Elaborazione di operazioni per il supporto alle decisioni

- Operazioni complesse e casuali
- Ogni operazione può coinvolgere molti dati
- Dati aggregati, storici, anche non attualissimi
- Le proprietà "acide" non sono rilevanti, perché le operazioni sono di sola lettura

Le **dimensioni** del warehouse raggiungono facilmente i **terabyte**

Le prestazioni considerate sono il **throughput delle interrogazioni e il loro tempo di risposta**

6

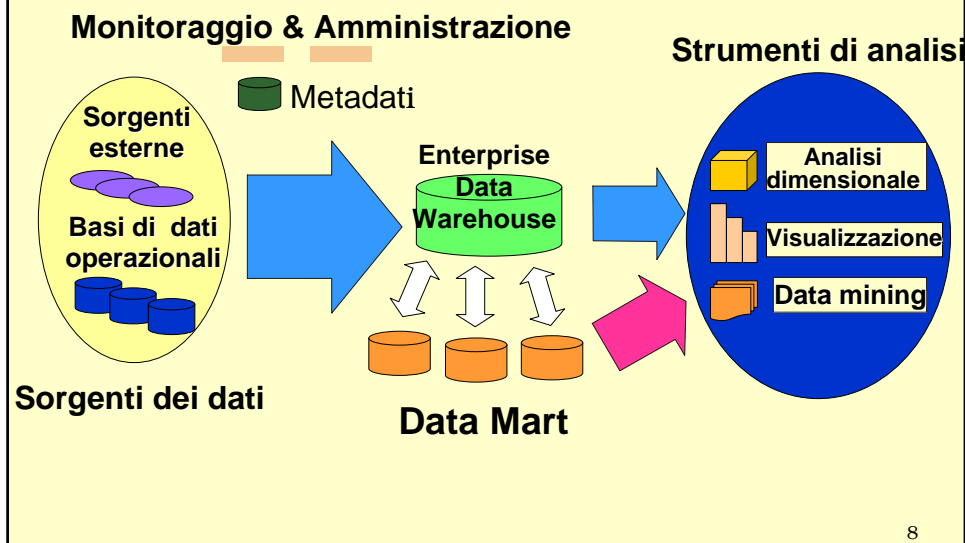
Data warehouse

Una base di dati

- utilizzata principalmente per il supporto alle decisioni direzionali
- integrata — aziendale e non dipartimentale
- orientata ai dati — non alle applicazioni
- Dati storici — con un ampio orizzonte temporale, e indicazione (di solito) di elementi di tempo
- non volatile — i dati sono caricati e acceduti fuori linea
- mantenuta separatamente dalle basi di dati operazionali

7

Architettura per il data warehousing



8

DW e Data Mart

- Una DW spesso integra diversi Data Mart
- Gli utenti normalmente si rivolgono a un particolare Data Mart
- I Data Mart condividono dati tra di loro
- Ciascun Data Mart è responsabile di un particolare aspetto della realtà aziendale, **non necessariamente** di un reparto.
- Es: data mart della rilevazione e rendicontazione di tutte le attività di un ospedale

9

MODELLI DEI DATI PER OLAP

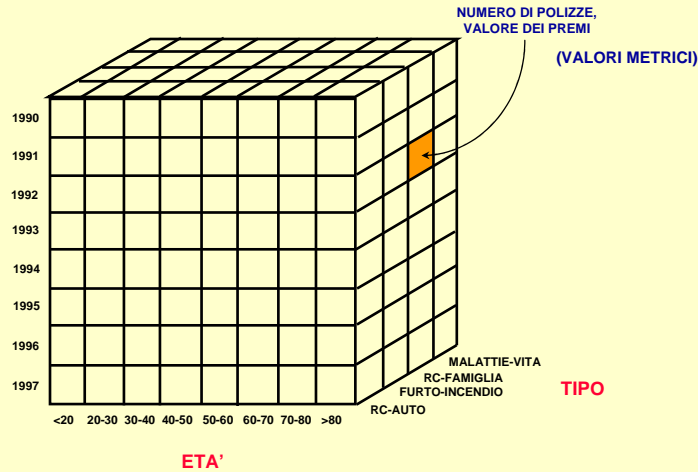
- devono supportare **analisi e calcoli sofisticati** su diverse dimensioni e gerarchie
- il modello logico dei dati piu' adatto e' una struttura multidimensionale - il **data cube**
- le **dimensioni** del cubo sono costituite dagli attributi secondo i quali si vogliono fare le ricerche (**chiavi**)
- ogni dimensione puo' "contenere" a sua volta una **gerarchia**
 - DATA {GIORNO - MESE – TRIMESTRE - ANNO}
 - PRODOTTO {NOME - TIPO - CATEGORIA}
(LAND ROVER - FUORISTRADA - AUTOVEICOLI)
- le **celle** del cubo contengono i valori **metrici** relativi ai valori dimensionali

10

MODELLI LOGICI DEI DATI PER OLAP

ESEMPIO PER UNA COMPAGNIA DI ASSICURAZIONI

DIMENSIONI:
ANNO, ETA', TIPO



11

Rappresentazione multidimensionale

Concetti rilevanti:

- **fatto** — un concetto sul quale centrare l'analisi - modella un evento che accade nell'azienda
- **misura** — una proprietà atomica di un fatto da analizzare - ne descrive un aspetto quantitativo
- **dimensione** — descrive una prospettiva lungo la quale effettuare l'analisi

12

Dimensioni e gerarchie di livelli

Ciascuna dimensione è organizzata in una **gerarchia** che rappresenta i possibili **livelli di aggregazione** per i dati



13

Progettazione di Data Warehouse

- La progettazione di una data warehouse è diversa dalla progettazione di una base di dati operativa
 - i dati da memorizzare hanno caratteristiche diverse
 - vincolata dalle basi di dati esistenti
 - guidata da criteri progettuali diversi
- Enfasi sulla generalizzazione e chiarezza concettuale
 - poche entità
 - ampia copertura

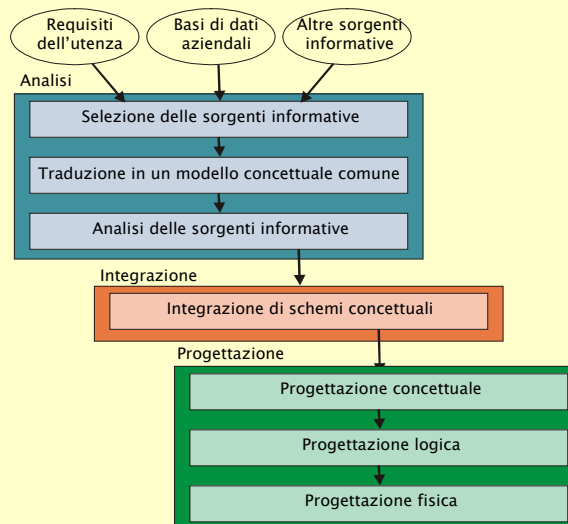
14

Progettazione di data warehouse

- Limitata frammentazione (denormalizzazione): l'analisi ha bisogno di una visione sintetica dei dati
- Nella progettazione sarà importante mettere in evidenza gli aspetti comuni
- Bisogna associare a ciascuna entità un significato che corrisponda all'intuizione dell'utente
- Fondamentale il ruolo dei metadati nel guidare all'uso del sistema
- Attività principali
 - analisi — delle sorgenti informative esistenti
 - integrazione
 - progettazione — concettuale, logica e fisica

15

Progettazione di data warehouse



16

Dati in ingresso

Le informazioni in ingresso necessarie alla progettazione di un data warehouse

- **requisiti** — le esigenze aziendali di analisi
- **descrizione delle basi di dati operazionali** — con una documentazione sufficiente per la loro comprensione
- **descrizione di altre sorgenti informative** — l'analisi richiede spesso la correlazione con dati non di proprietà dell'azienda ma comunque da essa accessibili — ad esempio, dati ISTAT o sull'andamento dei concorrenti

17

COSTRUZIONE DEL WAREHOUSE

- **i dati provengono da sorgenti diverse, probabilmente "sporche"**
 - **sistemi legacy non documentati**
 - **sistemi di produzione senza check di integrità' interni**
 - **sorgenti esterne di dubbia qualità'**
 - scarsa qualità del data entry
 - assenza di dati in alcuni campi

e' indispensabile restituire la qualità' ai dati per potervi basare decisioni affidabili

18

COSTRUZIONE DEL WAREHOUSE

- **strumenti per la **qualita'** dei dati**
 - per la **migrazione**
trasformano e riformattano i dati dalle diverse fonti
 - per la **pulizia (scrubbing)**
usano la conoscenza del dominio per pulire e omogeneizzare
<jerry l. jonson, 16 clarke st., altuna, pa> =
<gerry l. johnson, 16 clark street, altoona, penn> ???
 - per il **controllo (auditing)**
scoprono regole e relazioni tra i dati e ne verificano il rispetto
- **strumenti per il **caricamento** dei dati**
verificano violazioni di integrita' referenziale; ordinano, aggregano, costruiscono dati derivati; costruiscono indici e altri percorsi di accesso

19

Analisi delle sorgenti informative esistenti

- **Selezione delle sorgenti informative**
 - analisi preliminare del patrimonio informativo aziendale — analisi di qualità delle singole sorgenti
 - correlazione del patrimonio informativo con i requisiti
 - identificazione di **priorità tra schemi**
- **Traduzione in un modello concettuale di riferimento**
 - attività preliminare alla correlazione e all'integrazione di schemi — che si svolge meglio con riferimento a schemi concettuali
- **Analisi delle sorgenti informative**
 - identificazione di massima dei fatti, o concetti su cui basare l'analisi, delle loro misure (proprietà atomiche), e delle dimensioni (concetti su cui aggregare le misure)

20

Integrazione di sorgenti informative

- *L'integrazione di sorgenti informative* è l'attività di fusione dei dati rappresentati in più sorgenti in un'unica *base di dati globale* (fisica o virtuale) che rappresenta l'intero patrimonio informativo a disposizione
- Lo scopo principale dell'integrazione è l'*identificazione* di tutte le porzioni delle diverse sorgenti informative che si riferiscono a uno stesso aspetto della realtà di interesse, per *unificare* la loro rappresentazione
- L'approccio è orientato alla *identificazione, analisi e risoluzione di conflitti* — terminologici, strutturali, di codifica

21

Progettazione del data warehouse

- L'integrazione delle sorgenti informative ha prodotto una descrizione globale del patrimonio informativo aziendale
- Questo è però solo il risultato dell'integrazione di dati operazionali — non descrive tutti i dati di interesse per il DW
- Progettazione del data warehouse
 - **concettuale** — completare la rappresentazione dei concetti dimensionali necessari per l'analisi — ad esempio, dati storici e geografici
 - **logica** — identificare il miglior compromesso tra la necessità di aggregare i dati e quella di normalizzarli
 - **fisica** — individuare la distribuzione dei dati e le relative strutture di accesso

22

Progettazione del DW

- Introduzione di elementi dimensionali nella base di dati integrata
- Attività
 - identificazione di fatti, misure e dimensioni
 - ristrutturazione dello schema concettuale
 - rappresentazione di fatti mediante entità
 - individuazione di nuove dimensioni
 - raffinamento dei livelli di ogni dimensione
 - derivazione di un grafo dimensionale
 - progettazione logica e fisica

23

Strategie di progetto

- **Approccio top-down:**
 - Interessante perché garantisce la coerenza interna del progetto
 - Spesso fallimentare perché impresa lunga e ardua che scoraggia l'utenza
 - L'analisi contemporanea di tutte le fonti informative è compito molto complesso
 - La previsione contemporanea delle esigenze informative di tutti gli utenti è difficile e rischia di paralizzare il processo
- **Approccio bottom-up:**
 - Costruzione incrementale assemblando più data mart
 - Si abbina a tecniche di prototipazione veloce
 - Incoraggia l'utenza che vede velocemente il prodotto operativo

24

Approccio bottom up

Si sceglie di progettare per primo il data mart più strategico per l'azienda, che **guida** l'integrazione successiva

Ciclo di vita progetto bottom up

- Definizione degli obiettivi e pianificazione
- Progettazione dell'infrastruttura tecnologica
- Progettazione dei data mart e loro progressiva integrazione

25

La progettazione dei data mart

- **Partire dagli schemi dei dati operazionali esistenti:**
 - In questo modo non si rischia di promettere all'utente delle funzionalità che è impossibile fornire
- **Si riesce a tradurre lo schema concettuale dei dati operazionali in quello del data mart in modo quasi algoritmico**

26

Progettazione di data warehouse



27

Raccolta e analisi dei requisiti utente

Scopo fondamentale di questa fase è la scelta dei FATTI del data mart

- **Tecniche:**
 - **A risposte aperte:**
 - **Pro:** Più coinvolgenti, insegnano al progettista il linguaggio del dominio applicativo...
 - **Contro:** possono comportare risposte lunghe e dispersive
 - **A risposte chiuse:**
 - **Pro:** riducono i tempi di intervista, permettono il confronto tra risposte...
 - **Contro:** possono risultare noiose per l'intervistato
 - **A risposte probatorie**
 - **Pro:** permettono di capire il livello di competenza dell'intervistato...
 - **Contro:** possono comportare ansia nell'intervistato che si sente indagato

28

Esempi di domande

- **Per un dirigente:**
 - Quali sono gli obiettivi aziendali?
 - Come misuri il successo dell'azienda?
- **Per un capo reparto:**
 - Quali sono gli obiettivi del tuo reparto?
 - Quali sono i colli di bottiglia nell'accesso ai dati?
 - Quali analisi eseguite di solito sui dati? Quali vorreste eseguire?
- **Per l'amministratore del sistema informativo:**
 - Quali sono le caratteristiche delle fonti di dati disponibili?
 - Che strumenti vengono attualmente usati per analizzare i dati?
 - Quali analisi vi vengono richieste di solito sui dati?

29

Il glossario dei requisiti utente

FATTO	POSSIBILI DIMENSIONI	POSSIBILI MISURE	STORICITA'
<i>inventario di magazzino</i>	<i>prodotto, data, magazzino</i>	<i>quantità in magazzino</i>	<i>1 anno</i>
<i>vendite</i>	<i>prodotto, data, negozio</i>	<i>quantità venduta, importo, sconto</i>	<i>5 anni</i>
<i>linee d'ordine</i>	<i>prodotto, data, fornitore</i>	<i>quantità ordinata, importo, sconto</i>	<i>3 anni</i>

30

Il carico di lavoro preliminare

FATTO	POSSIBILI INTERROGAZIONI AGGREGATE
<i>inventario di magazzino</i>	Scorte medie di ciascun prodotto presenti mensilmente Prodotti per i quali è stata esaurita la scorta almeno una volta in tutti i magazzini durante l'ultima settimana
<i>vendite</i>	Quantità totali di ciascun tipo di prodotto vendute ogni mese per regione Incasso totale giornaliero di ciascun negozio Per ciascun negozio, riepilogo delle vendite di ciascuna categoria di prodotto per ogni giorno
<i>linee d'ordine</i>	Quantità totale ordinata annualmente presso un dato fornitore Importo giornaliero ordinato nell'ultimo mese per ciascuna categoria di prodotto

31

Ci serve a capire, ad esempio, che:

- **Per il magazzino:**
 - Serve sapere, per ogni prodotto, il tipo
 - La gerarchia temporale deve includere almeno il mese e la settimana
 - Occorre poter calcolare le medie delle quantità
- **Per le vendite:**
 - Serve sapere, per ogni prodotto, il tipo e la categoria
 - La gerarchia della dimensione negozio deve includere l'attributo regione
- **Per le linee d'ordine:**
 - Serve sapere, per ogni prodotto, il tipo e la categoria
 - Durante l'aggregazione la misura importo viene sommata sulle gerarchie dei prodotti e dei fornitori

32

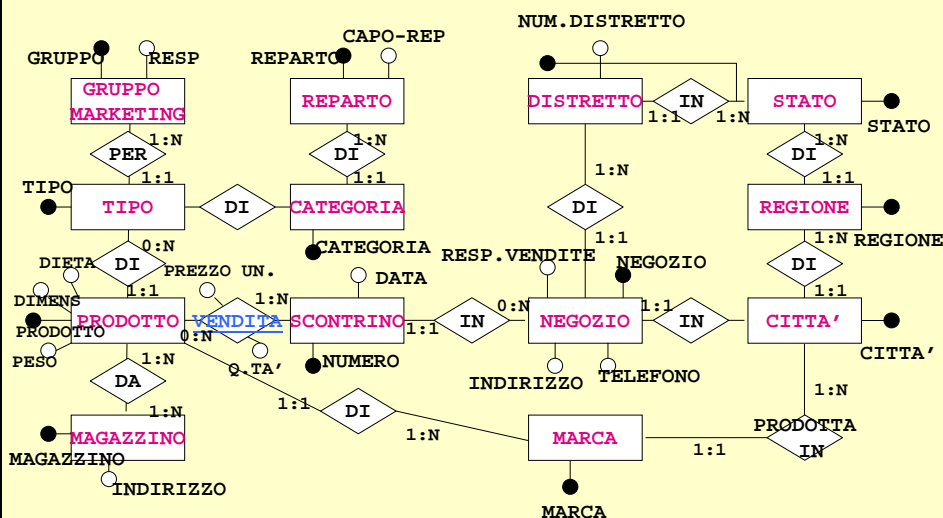
Progettazione di schemi concettuali di DW a partire da schemi ER

1) Definizione dei fatti

- I **fatti** sono concetti di interesse per il processo decisionale, e corrispondono a eventi che accadono nel mondo aziendale
- Le entità che rappresentano dati aggiornati frequentemente (es. VENDITE) sono **buoni candidati** per essere dei fatti - es. NEGOZIO, CITTÀ non lo sono
- Per ogni fatto:
 - Costruzione dell'**albero degli attributi**

33

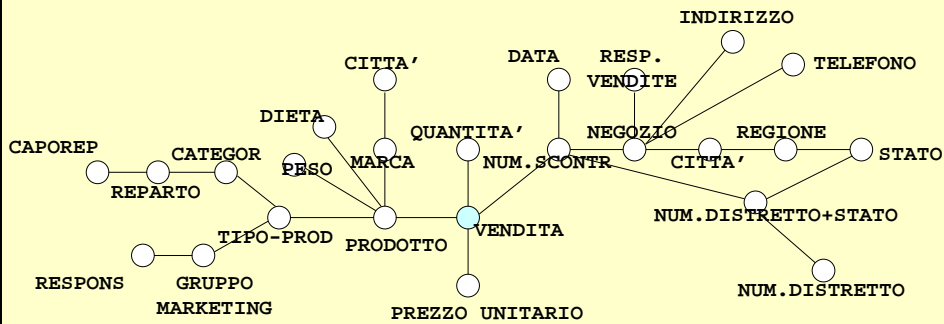
BASE DI DATI OPERAZIONALE



34

L'albero degli attributi

Fatto: VENDITE



VENDITA = PRODOTTO + NUM. SCONTRINO

35

Costruzione dell'albero degli attributi

Fatto: VENDITE

- ogni vertice corrisponde a un attributo dello schema sorgente
- La radice corrisponde all'identificatore del fatto
- Per ogni vertice v dell'albero, l'attributo corrispondente determina funzionalmente gli attributi corrispondenti ai discendenti di v
- Le dipendenze funzionali si riconoscono, sono:
 - Identificatori
 - Associazioni "a-uno"

36

Si costruisce ricorsivamente attraversando il diagramma ER

- Si parte dall'entità corrispondente al fatto scelto (*radice*)
- Per ogni entità E attraversata, si crea nell'albero un vertice *v* corrispondente all'identificatore di E
- Si aggiunge un vertice figlio per ogni attributo di E (inclusi gli attributi che compongono l'id)
- Per ogni relazione "a-uno" uscente da E verso una entità G, si aggiungono anche tutti i suoi attributi come figli di *v*
- Si ripete il procedimento per G

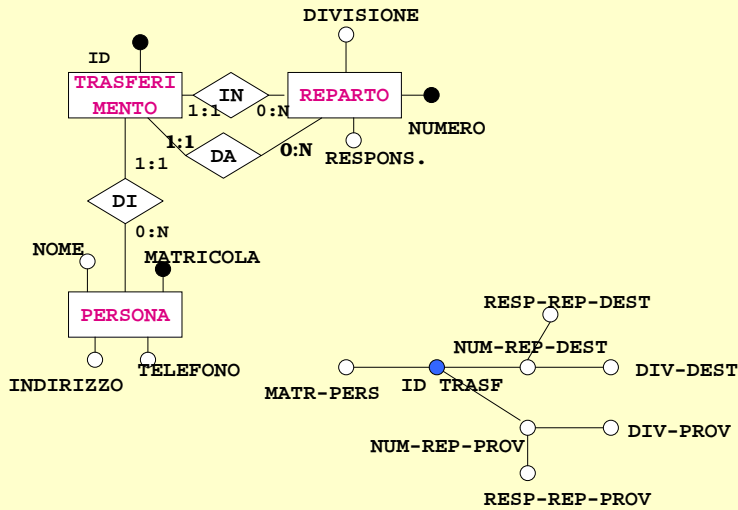
37

Casi particolari

- Relazioni cicliche (es. parte-sottoparte) → si spezzano dopo un certo num. di iterazioni
- Cicli nello schema → si spezzano legando all'entità + conveniente (event. duplicando)
- associazioni n-arie o molti-a molti si possono *reificare*, cioè trasformare in entità
- Attributi opzionali: segnati da un trattino
- Gerarchie: equiv. ad associazioni 0-1 (trattino di opzionalità)
- Attributo composto: un vertice con figli

38

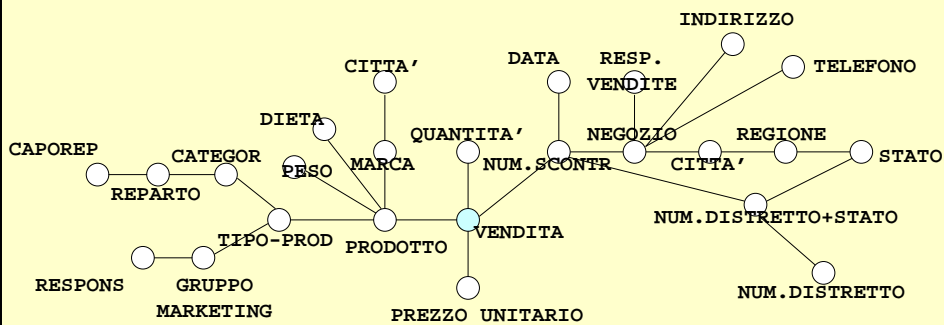
Cicli nello schema



39

L'albero degli attributi

Fatto: **VENDITE** →
relazione reificata



VENDITA = PRODOTTI + NUM. SCONTRINO

40

Potatura e innesto

- **Potatura** di un vertice v : si elimina v con tutti i discendenti
- **Innesto** di v con padre v' : si elimina v e si collegano i figli di v direttamente a v'
- Es. potatura di *num-scontrino* e innesto di *data* e *negozio* su *vendite*

41

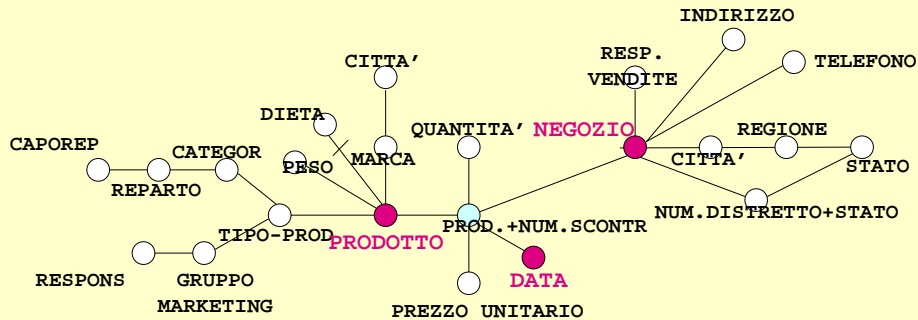
Progettazione di schemi concettuali di DW

2) Identificazione delle **dimensioni**

- Le dimensioni vanno scelte nell'albero degli attributi tra i vertici **figli della radice**
- Possono corrispondere ad **attributi discreti** o ad **intervalli di attributi continui**
- Il **tempo** è una dimensione sempre significativa
 - Se non appare come figlio della radice, vale la pena di ristrutturare l'albero (potare e innestare) per farlo comparire

42

Esempio DIMENSIONI



SCEGLIAMO COME DIMENSIONI
LA TERNA <DATA,NEGOZIO,PRODOTTO>

43

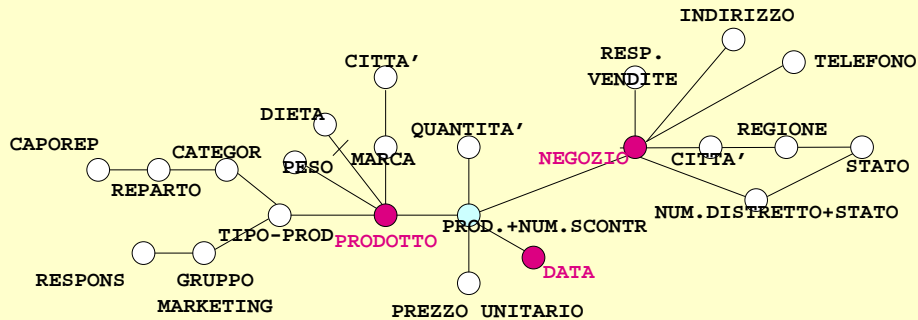
Progettazione di schemi concettuali di DW

3) Identificazione delle **misure**

- Di solito le misure vanno scelte nell'albero degli attributi tra i vertici **figli della radice**
- Questo non è assolutamente necessario, nel caso non si verificasse **si può potare e innestare l'albero**
- Le misure di solito si ottengono aggregando **rispetto a tutte le dimensioni**

44

Esempio MISURE



Quantità venduta = SUM(VENDITA.quantità)

Incasso=SUM(VENDITA.quantità*VENDITA.prezzoUnitario)

prezzoUnitario=AVG(VENDITA.prezzoUnitario)

NumClienti=COUNT(*)

IN CORRISPONDENZA DELLA STESSA TERNA <DATA,NEGOZIO,PRODOTTI>

45

Tecnologie target per la progettazione logica

MOLAP (Multidimensional-**OLAP**) in
alternativa a **ROLAP** (Relational-**OLAP**)

- **MOLAP**: usa strutture interne non relazionali, prestazioni migliori
- **ROLAP**: in grado di gestire grandi quantità di dati mediante le classiche tecnologie relazionali

46

SERVER OLAP MULTIDIMENSIONALE (MOLAP)

- implementa **direttamente** il modello a cubo
 - strutture a **matrice multidimensionale**
- prestazioni **elevate e costanti** per l'elaborazione delle interrogazioni
 - metodi di accesso **specializzati**
 - aggregazione e compilazione eseguite in precedenza
- **limitata scalabilita'** a causa delle preelaborazioni
- richiede maggiori capacita' da parte dell'amministratore dei dati

47

SERVER OLAP RELAZIONALE (ROLAP)

- utilizza un **RDBMS standard** per realizzare la struttura multidimensionale, applicando operatori aggregati
- lo **schema** assume una configurazione a **stella**
- **Centro stella:** fatto
- **Punte stella:** dimensioni
- Vantaggi:
 - Si possono usare opportune interfacce di interrogazione
 - Si ottengono buone prestazioni
 - Si ha un'immediata definizione dello schema logico

48

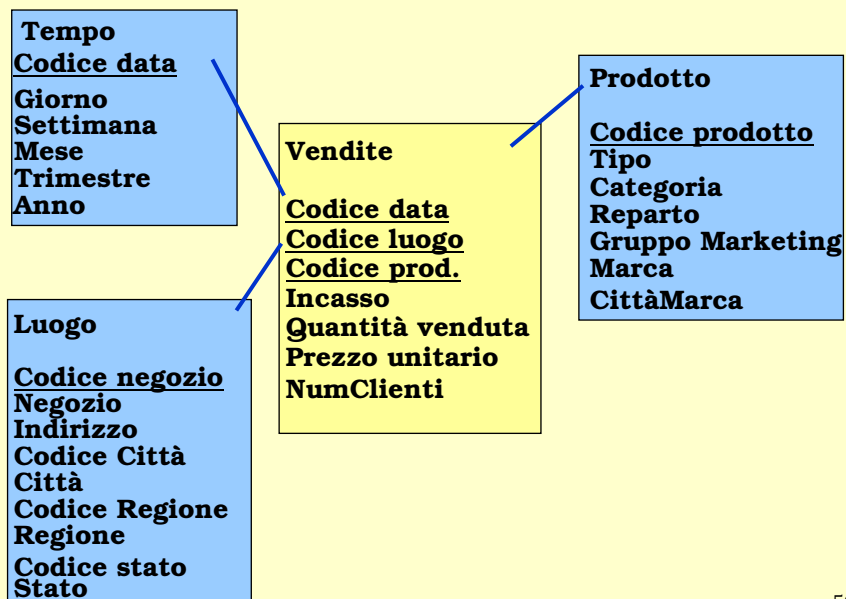
SERVER OLAP RELAZIONALE (ROLAP)

Lo **schema** assume una configurazione **a stella** o a **fiocco di neve**

- **tabella centrale dei fatti**
 - le tuple sono costituite dai **puntatori** (=chiavi esterne) alle tabelle di dimensione e dai **valori per le misure descritte**
 $f = (k_1, \dots, k_n, v_1, \dots, v_m)$
- **tabelle di dimensione**
 - contengono le tuple con gli **attributi relativi a quella dimensione** $d_1=(k_1, a_1, \dots, a_n)$
- **costellazione di fatti**
 - piu' tabelle dei fatti condividono tabelle di dimensione di uguale struttura

49

Organizzazione a stella



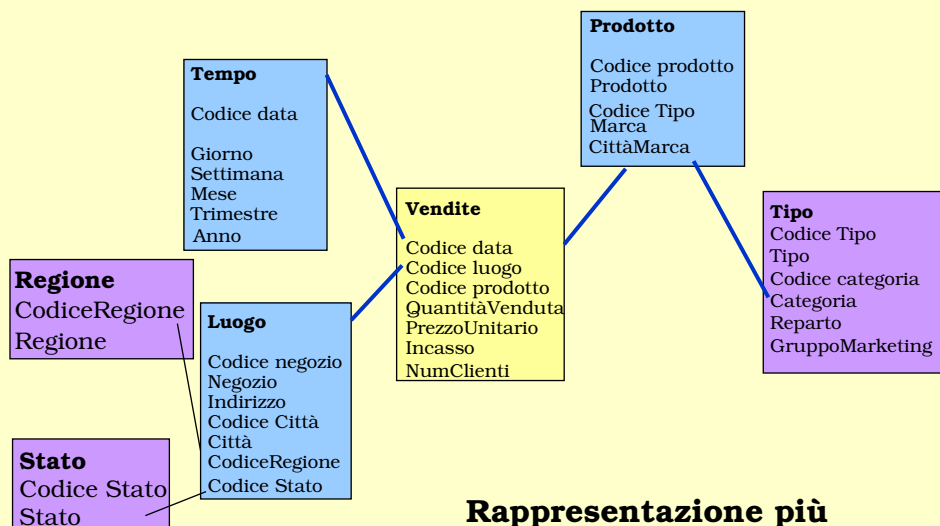
50

Modello a fiocco di neve (snow flake)

- Costituisce un'estensione del modello a stella
- Permette di evitare ridondanze eccessive nelle dimensioni
- Dalla tabella dei fatti si raggiungono tutte le tabelle delle dimensioni, sempre muovendosi lungo associazioni n:1

51

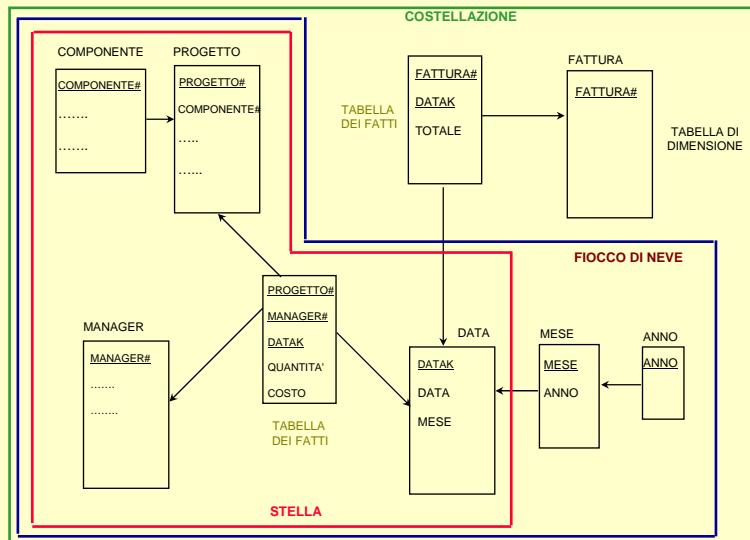
Organizzazione "snowflake"



Rappresentazione più "normalizzata" ma richiede più join

52

SCHEMI A STELLA



53

Strategia di implementazione

- Si realizza dapprima uno schema di massima e si individuano i dati comuni ai diversi data mart, definendo per essi uno schema condiviso (con l'approccio appena visto)
- Si realizzano poi uno o più Data Mart, con l'obiettivo di familiarizzare gli sviluppatori e gli utenti con la tecnologia OLAP
- Infine, si integrano i sistemi già sviluppati e si procede alla realizzazione della Enterprise Data Warehouse

54

Altri requisiti

Altri requisiti vengono raccolti per le fasi successive:

- **Per la progettazione logica e fisica: requisiti di spazio e tempo – ottimizzazione degli accessi tramite viste materializzate**
- **Per il progetto dell'alimentazione: periodicità dell'alimentazione, grado di "freschezza" richiesto**

55

Tecnologie per il progetto fisico

- **Indici join**
 - Precomputano il join tra la dimensione e la tabella dei fatti
- **Materializzazione di viste**
 - Vengono precalcolate le viste che possono essere utilizzate come base per rispondere alle query più frequenti

56