

Il ciclo di vita del Data Warehouse

Prof. Stefano Rizzi

Perché?

- Molte organizzazioni mancano della necessaria esperienza e capacità per affrontare con successo le sfide implicite nei progetti di data warehousing
- Uno dei fattori che maggiormente minaccia la riuscita dei progetti è la mancata adozione di un **approccio metodologico**, che minimizza i rischi di insuccesso essendo basato su un'analisi costruttiva degli errori commessi



Fattori di rischio

- ✓ Rischi legati alla gestione del progetto
 - ✓ Rischi legati alle tecnologie
 - ✓ Rischi legati ai dati e alla progettazione
 - ✓ Rischi legati all'organizzazione
- Il rischio di ottenere un risultato insoddisfacente nei progetti di data warehousing è particolarmente alto a causa delle elevatissime aspettative degli utenti
 - Nella cultura aziendale contemporanea è infatti diffusissima la credenza che attribuisce al data warehousing il ruolo di panacea
 - In realtà una larga parte della responsabilità della riuscita del progetto ricade sulla qualità dei dati sorgente e sulla lungimiranza, disponibilità e dinamismo del personale dell'azienda

3



Approccio top-down

- Analizza i bisogni globali dell'intera azienda e pianifica lo sviluppo del DW per poi progettare e realizzarlo nella sua interezza
 - 👍 Promette ottimi risultati poiché si basa su una visione globale dell'obiettivo e garantisce in linea di principio di produrre un DW consistente e ben integrato
 - 👎 Il preventivo di costi onerosi a fronte di lunghi tempi di realizzazione scoraggia la direzione dall'intraprendere il progetto
 - 👎 Affrontare contemporaneamente l'analisi e la riconciliazione di tutte le sorgenti di interesse è estremamente complesso
 - 👎 Riuscire a prevedere a priori nel dettaglio le esigenze delle diverse aree aziendali impegnate è pressoché impossibile, e il processo di analisi rischia di subire una paralisi
 - 👎 Il fatto di non prevedere la consegna a breve termine di un prototipo non permette agli utenti di verificare l'utilità del progetto e ne fa scemare l'interesse e la fiducia

4

Approccio bottom-up

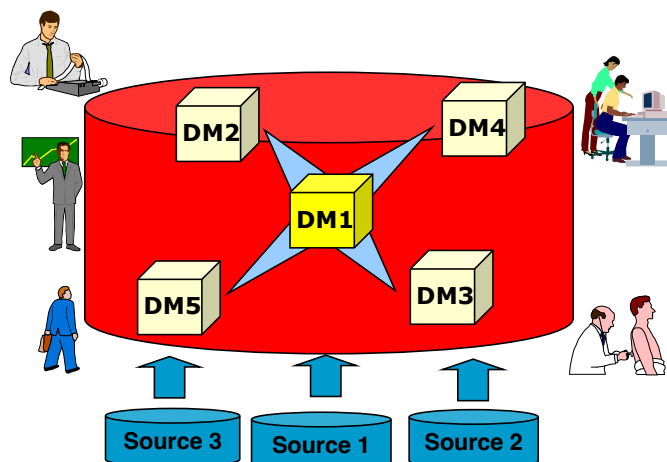
- Il DW viene costruito in modo incrementale, assemblando iterativamente più data mart, ciascuno dei quali incentrato su un insieme di fatti collegati a uno specifico settore aziendale e di interesse per una certa categoria di utenti

- 👍 Determina risultati concreti in tempi brevi
- 👍 Non richiede elevati investimenti finanziari
- 👍 Permette di studiare solo le problematiche relative al data mart in oggetto
- 👍 Fornisce alla dirigenza aziendale un riscontro immediato sull'effettiva utilità del sistema in via di realizzazione
- 👍 Mantiene costantemente elevata l'attenzione sul progetto
- 👍 Determina una visione parziale del dominio di interesse

5

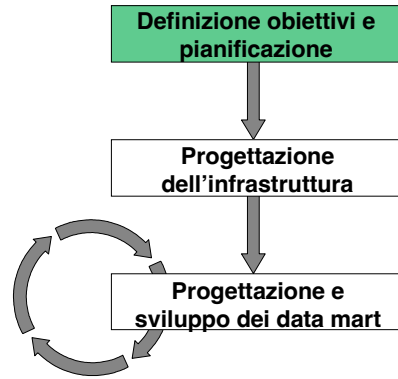
Il primo data mart da prototipare...

- ✓ deve essere quello che gioca il ruolo più strategico per l'azienda
- ✓ deve ricoprire un ruolo centrale e di riferimento per l'intero DW
- ✓ si deve appoggiare su fonti dati già disponibili e consistenti



6

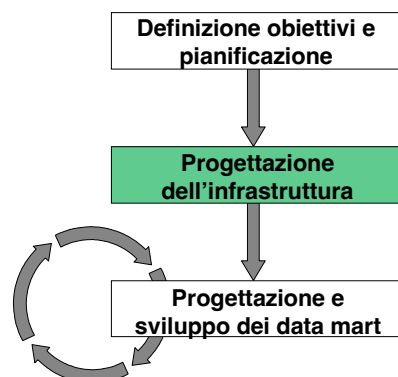
Il ciclo di sviluppo



- individuazione degli obiettivi e dei confini del sistema
- stima delle dimensioni
- scelta dell'approccio per la costruzione
- valutazione dei costi e del valore aggiunto
- analisi dei rischi e delle aspettative
- studio delle competenze del gruppo di lavoro

7

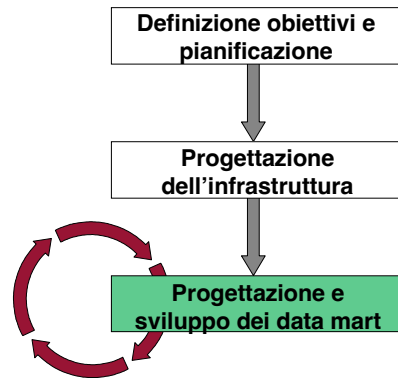
Il ciclo di sviluppo



Si analizzano e si comparano le possibili soluzioni architetture valutando le tecnologie e gli strumenti disponibili, al fine di realizzare un progetto di massima dell'intero sistema.

8

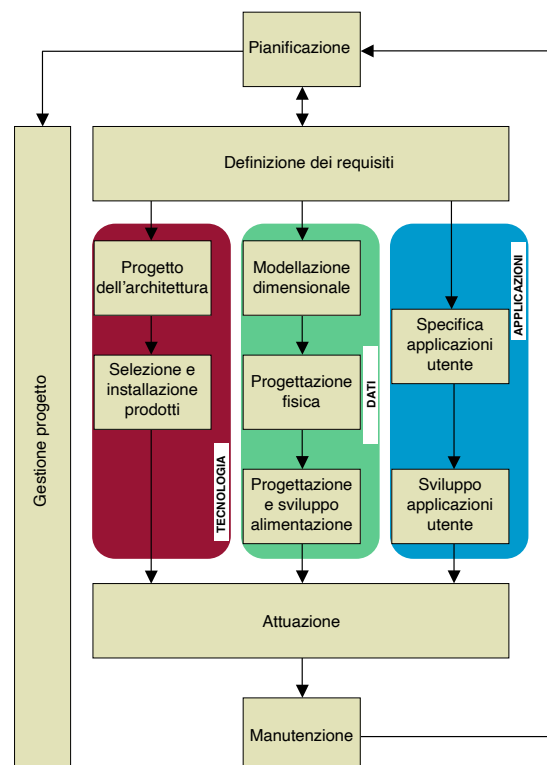
Il ciclo di sviluppo



Ciascuna iterazione comporta la creazione di un nuovo data mart e di nuove applicazioni, che vengono via via integrate nel sistema di data warehousing.

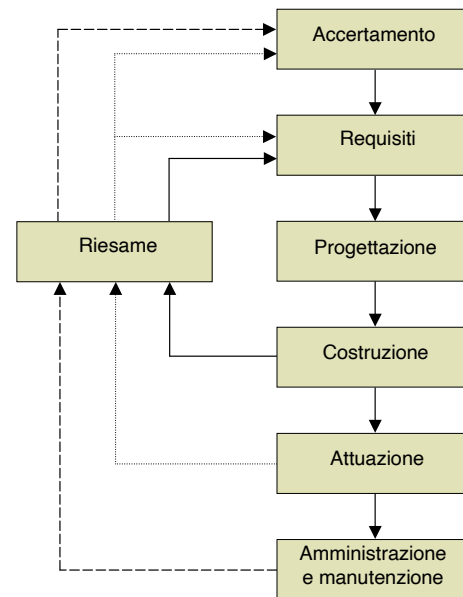
9

Il "Business Dimensional Lifecycle" (Kimball)



10

La "Rapid Warehousing Methodology" (SAS)

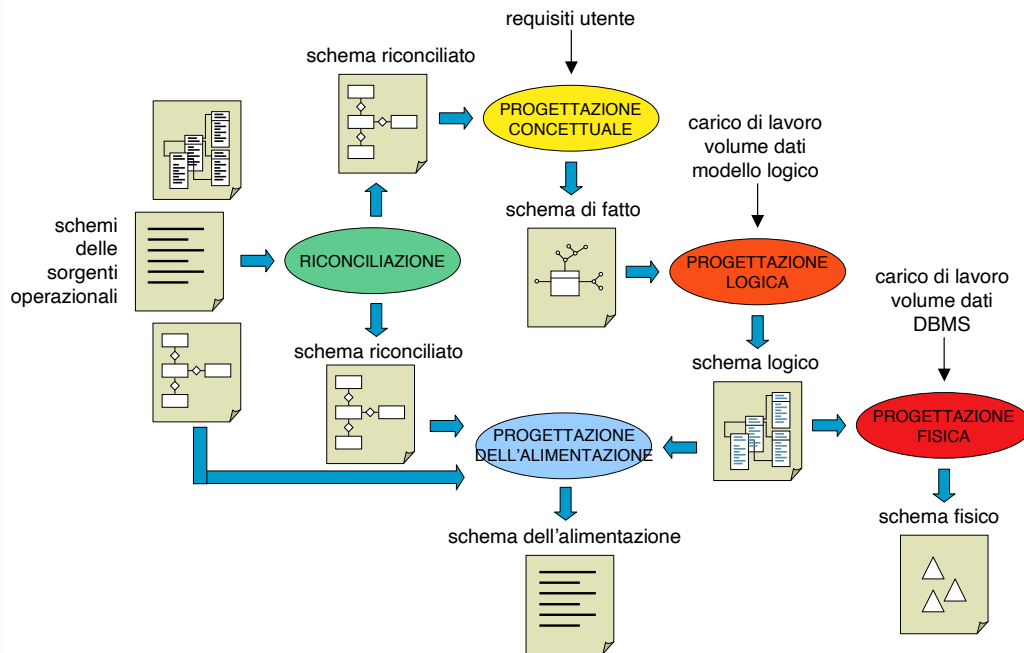


11

La progettazione del data mart



La progettazione del data mart

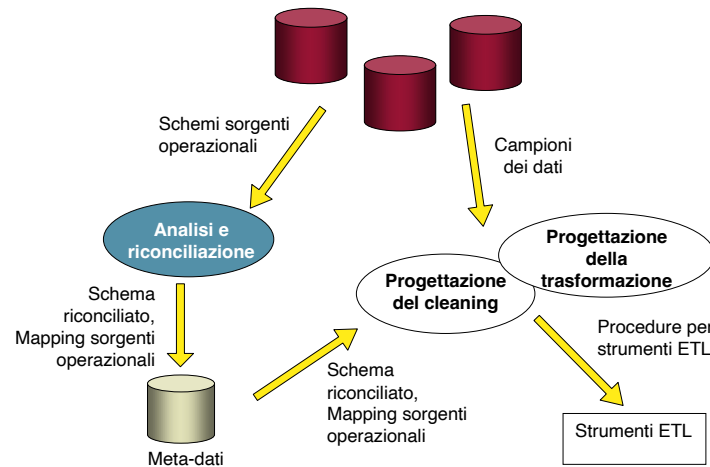


13

Analisi e riconciliazione delle sorgenti operazionali

Prof. Stefano Rizzi

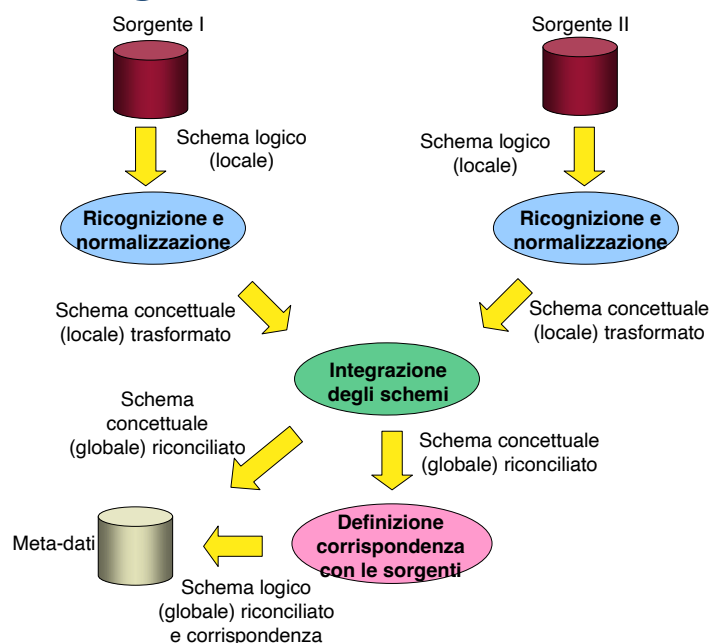
Progettazione del livello riconciliato



- ✓ La fase di integrazione è incentrata sulla componente intensionale delle sorgenti operazionali, ossia riguarda la consistenza degli schemi che le descrivono
- ✓ Pulizia e trasformazione dei dati operano a livello estensionale, ossia coinvolgono direttamente i dati veri e propri

15

Analisi e riconciliazione delle sorgenti operazionali



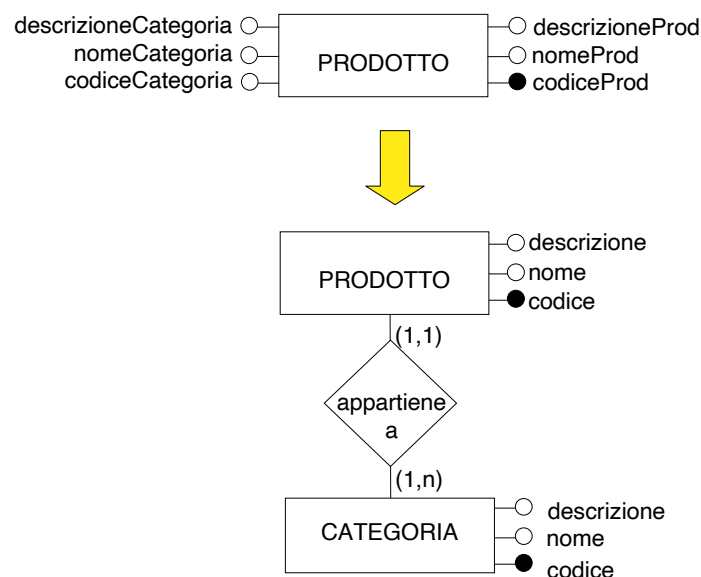
16

Ricognizione e normalizzazione

- Il progettista, confrontandosi con gli esperti del dominio applicativo, acquisisce un'approfondita conoscenza delle sorgenti operazionali attraverso:
 - ✓ **ricognizione**, che consiste in un esame approfondito degli schemi locali mirato alla piena comprensione del dominio applicativo;
 - ✓ **normalizzazione**, il cui obiettivo è correggere gli schemi locali al fine di modellare in modo più accurato il dominio applicativo
- Ricognizione e normalizzazione devono essere svolte anche qualora sia presente una sola sorgente dati; qualora esistano più sorgenti, l'operazione dovrà essere ripetuta per ogni singolo schema

17

Ricognizione e normalizzazione



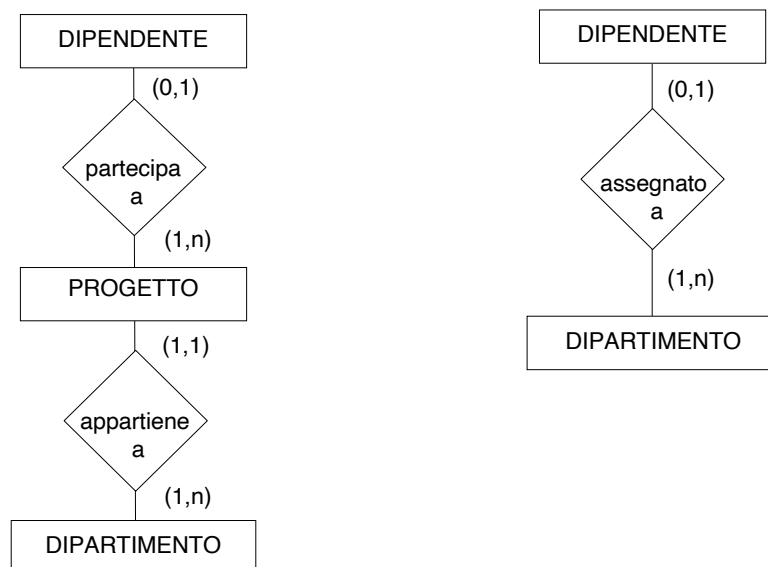
18

Integrazione

- L'integrazione di un insieme di sorgenti dati eterogenee (basi di dati relazionali, file dati, sorgenti legacy) consiste nell'individuazione delle corrispondenze tra i concetti rappresentati negli schemi locali e nella risoluzione dei conflitti evidenziati, finalizzate alla creazione di un unico schema globale i cui elementi possano essere correlati con i corrispondenti elementi degli schemi locali (*mapping*)
- La fase di integrazione non si deve limitare a evidenziare le differenze di rappresentazione dei concetti comuni a più schemi locali, ma deve anche identificare l'insieme di concetti distinti e memorizzati in schemi differenti che sono correlati attraverso proprietà semantiche (*proprietà interschema*)
- Per poter ragionare sui concetti espressi negli schemi delle diverse sorgenti dati è necessario utilizzare **un unico formalismo** in modo da fissare i costrutti utilizzabili e la potenza espressiva

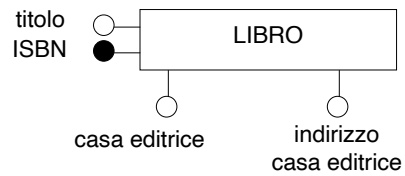
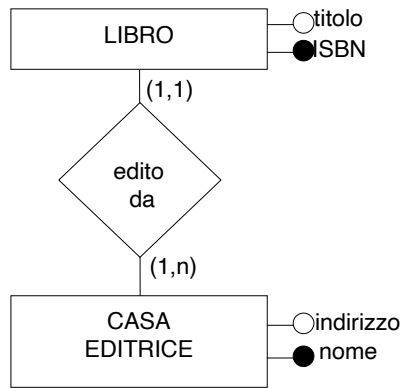
19

Problemi: diversa prospettiva



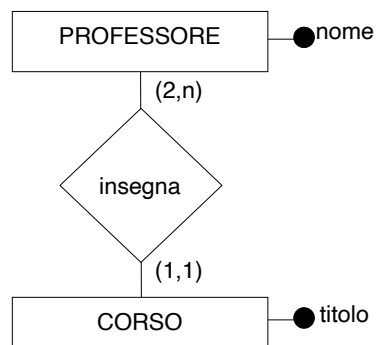
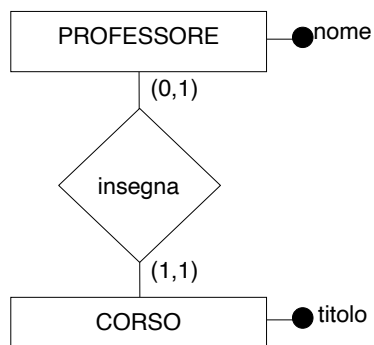
20

Problemi: costrutti equivalenti



21

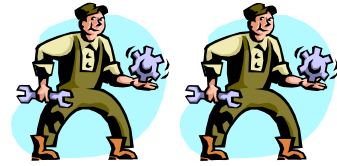
Problemi: incompatibilità



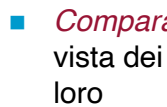
22

Relazioni tra concetti comuni

- **Identità:** vengono utilizzati gli stessi costrutti, il concetto è modellato dallo stesso punto di vista e non vengono commessi errori di specifica



- **Equivalenza:** sono stati utilizzati costrutti diversi (ma equivalenti) e non sussistono errori di specifica o diversità di percezione



- **Comparabilità:** i costrutti utilizzati e i punti di vista dei progettisti non sono in contrasto tra loro



- **Incompatibilità:** gli schemi sono in contrasto a causa dell'incoerenza nelle specifiche

23

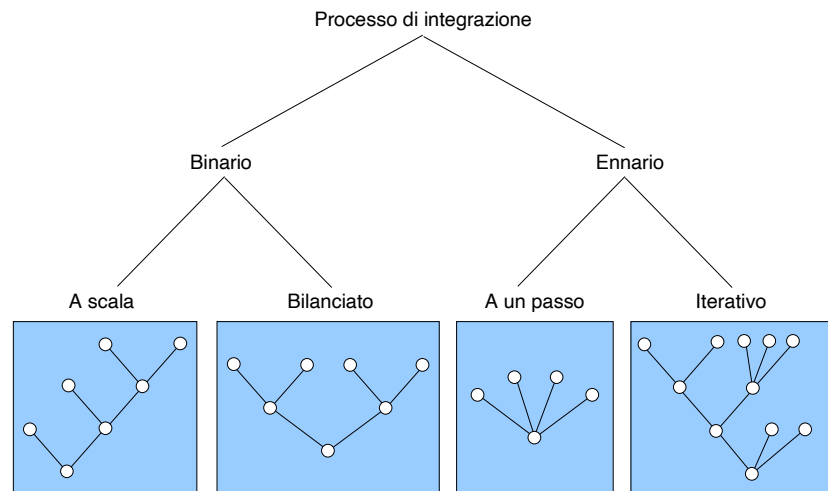
Fasi dell'integrazione

1. **Preintegrazione**
2. **Comparazione degli schemi**
3. **Allineamento degli schemi**
4. **Fusione e ristrutturazione degli schemi**

24

1. Preintegrazione

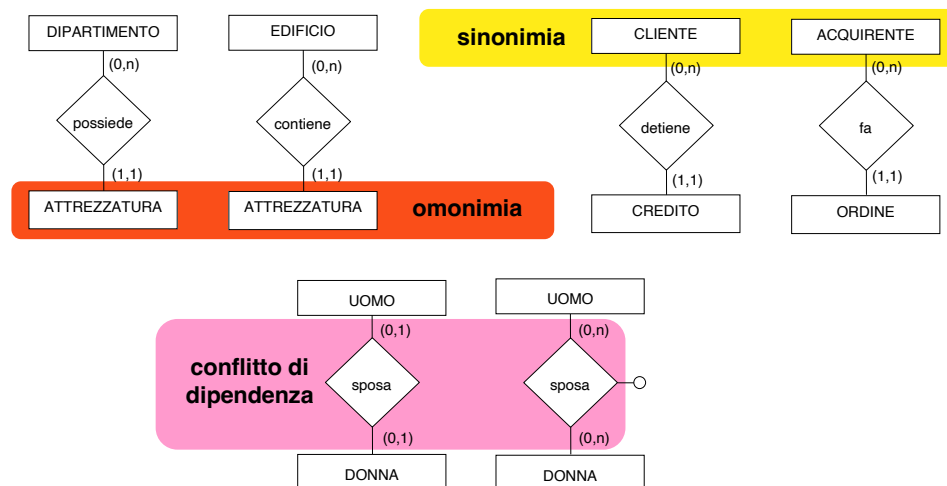
- Viene definita la strategia di integrazione



25

2. Comparazione degli schemi

- Un'analisi comparativa dei diversi schemi che mira a identificare le correlazioni e i conflitti tra i concetti in essi espressi



26



3. Allineamento degli schemi

- Scopo di questa fase è la risoluzione dei conflitti evidenziatisi al passo precedente, che si ottiene applicando primitive di trasformazione agli schemi sorgenti o allo schema riconciliato temporaneamente definito
 - ✓ Tipiche primitive di trasformazione riguardano il cambio dei nomi e dei tipi degli attributi, la modifica delle dipendenze funzionali e dei vincoli esistenti sugli schemi
 - ✓ Non sempre i conflitti possono essere risolti, poiché derivano da inconsistenze di base del sistema informativo; in questo caso la soluzione deve essere discussa con gli utenti che dovranno fornire indicazioni su qual è la più fedele interpretazione del mondo reale
 - ✓ In caso di incertezza si preferiscono le trasformazioni che avvantaggiano gli schemi ritenuti centrali nella struttura del data mart

27



4. Fusione degli schemi

- Gli schemi allineati vengono fusi a formare un unico schema riconciliato; l'approccio più diffuso è quello di sovrapporre i concetti comuni a cui saranno collegati tutti i rimanenti concetti provenienti dagli schemi locali.
- Dopo questa operazione si renderanno necessarie ulteriori trasformazioni mirate a migliorare la struttura dello schema riconciliato rispetto a:
 - ✓ Completezza
 - ✓ Minimalità
 - ✓ Leggibilità

28



Definizione delle corrispondenze

// DB1 Magazzino

ORDINI2001(chiaveO, chiaveC, data ordine, impiegato)

CLIENTE(chiaveC, nome, indirizzo, città, regione, stato)

.....

// DB2 Amministrazione

CLIENTE(chiaveC, partitalva, nome, telefono, fatturato)

FATTURE(chiaveF, data, chiaveC, importo, iva)

STORICO_ORDINI2000(chiaveO, chiaveC, data ordine, impiegato)

.....

CREATE VIEW CLIENTE AS

SELECT CL1.chiaveC, CL1.nome, CL1.indirizzo, CL1.città, CL1.regione,
CL1.stato, CL2.partitalva, CL2.telefono, CL2.fatturato

FROM DB1.CLIENTE AS CL1, DB2.CLIENTE AS CL2

WHERE CL1.chiaveC = CL2.chiaveC;

CREATE VIEW ORDINI AS

SELECT * FROM DB1.ORDINI2001

UNION

SELECT * FROM DB2.STORICO_ORDINI2000;

29

Analisi dei requisiti



Prof. Stefano Rizzi

Obiettivi

- La fase di analisi dei requisiti ha l'obiettivo di raccogliere le esigenze di utilizzo del data mart espresse dai suoi utenti finali
- Essa ha un'importanza strategica poiché influenza le decisioni da prendere riguardo:
 - ✓ lo schema concettuale dei dati
 - ✓ il progetto dell'alimentazione
 - ✓ le specifiche delle applicazioni per l'analisi dei dati
 - ✓ l'architettura del sistema
 - ✓ il piano di avviamento e formazione
 - ✓ le linee guida per la manutenzione e l'evoluzione del sistema.

31

Fonti

- La “fonte” principale da cui attingere i requisiti sono i futuri utenti del data mart (*business users*)
 - ✓ La differenza nel linguaggio usato da progettisti e utenti, e la percezione spesso distorta che questi ultimi hanno del processo di warehousing, rendono il dialogo difficile e a volte infruttuoso
- Per gli aspetti più tecnici, saranno gli amministratori del sistema informativo e/o i responsabili del CED a fungere da riferimento per il progettista
 - ✓ In questo caso, i requisiti che dovranno essere catturati riguardano principalmente vincoli di varia natura imposti sul sistema di data warehousing



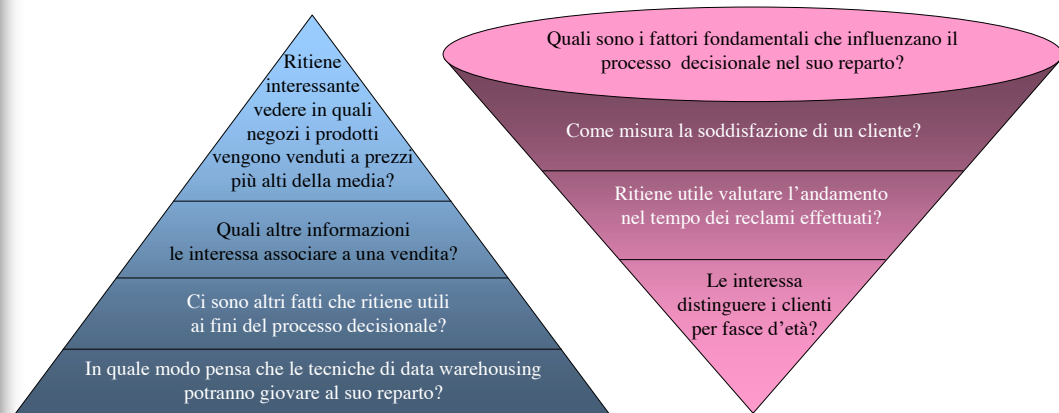
32

Le interviste

- **A piramide.** Approccio induttivo: l'intervistatore parte da domande molto dettagliate per poi ampliare l'argomento dell'intervista mediante domande aperte che richiedono risposte più generali.
 - ✓ Questo tipo di intervista permette di superare la riluttanza di un intervistato scettico poiché inizialmente non richiede un forte coinvolgimento da parte dell'intervistato.
- **A imbuto.** Approccio deduttivo: l'intervistatore parte da domande molto generali per poi restringere l'argomento dell'intervista a temi specifici
 - ✓ Questo approccio è utile nel caso in cui l'intervistato sia emozionato o eccessivamente deferente, poiché il fatto che le domande di carattere generale (normalmente in forma aperta) non prevedano una risposta "sbagliata" allevia la tensione dell'intervistato.

33

Le interviste



34



Le domande

Ruolo	Domande chiave
Dirigente	Quali sono gli obiettivi aziendali? Come misuri il successo della tua azienda? Quali sono oggi i principali problemi dell'azienda? In che modo ti aspetti che una maggiore disponibilità di informazioni possa migliorare la situazione aziendale?
Direttore di reparto	Quali sono gli obiettivi del tuo reparto? Come misuri il successo del tuo reparto? Descrivi i soggetti coinvolti nel tuo settore di interesse. Ci sono colli di bottiglia nell'accesso ai dati? Che analisi di routine esegui? Che tipi di analisi ti piacerebbe poter eseguire? A che livello di dettaglio occorre vedere le informazioni? Quanta informazione storica è necessaria?
Amministratore del sistema informativo	Illustra le caratteristiche delle principali fonti dati disponibili. Che strumenti vengono usati per analizzare i dati? Come vengono gestite le richieste di analisi ad hoc? Quali sono i principali problemi di qualità dei dati?

35



I fatti

- I **fatti** sono i concetti su cui gli utenti finali del data mart baseranno il processo decisionale; ogni fatto descrive una categoria di eventi che si verificano in azienda
 - ✓ Fissare le dimensioni di un fatto è importante poiché significa determinarne la **granularità**, ovvero il più fine livello di dettaglio a cui i dati saranno rappresentati. La scelta della granularità di un fatto nasce da un delicato compromesso tra due esigenze contrapposte: quella di raggiungere un'elevata flessibilità d'utilizzo e quella di conseguire buone prestazioni
 - ✓ Per ogni fatto occorre definire l'**intervallo di storicizzazione**, ovvero l'arco temporale che gli eventi memorizzati dovranno coprire

36

I fatti

	<i>Data mart</i>	<i>Fatti</i>
commerciale/ manfatturiero	approvvigionamenti	acquisti, inventario di magazzino, distribuzione
	produzione	confezionamento, inventario, consegna, manifattura
	gestione domanda	vendite, fatturazione, ordini, spedizioni, reclami
	marketing	promozioni, fidelizzazione, campagne pubblicitarie
finanziario	bancario	conti correnti, bonifici, prestiti ipotecari, mutui
	investimenti	acquisto titoli, transazioni di borsa
	servizi	carte di credito, domiciliazioni bollette
sanitario	scheda di ricovero	ricoveri, dimissioni, interventi chirurgici, diagnosi
	pronto soccorso	accessi, esami, dimissioni
	medicina di base	scelte, revoche, prescrizioni
trasporti	merci	domanda, offerta, trasporti
	passengeri	domanda, offerta, trasporti
	manutenzione	interventi
telecomunicazioni	traffico	traffico in rete, chiamate
	CRM	fidelizzazione, reclami, servizi
	gestione domanda	biglietteria, noleggi auto, soggiorni
turismo	CRM	frequent-flyers, reclami
	logistica	trasporti, scorte, movimentazione
gestionale	risorse umane	assunzioni, dimissioni, promozioni, incentivi
	budgeting	budget commerciale, budget di marketing
	infrastrutture	acquisti, opere

37

Glossario dei requisiti

<i>Fatto</i>	<i>Possibili dimensioni</i>	<i>Possibili misure</i>	<i>Storicità</i>
inventario di magazzino	prodotto, data, magazzino	quantità in magazzino	1 anno
vendite	prodotto, data, negozio	quantità venduta, importo, sconto	5 anni
linee d'ordine	prodotto, data, fornitore	quantità ordinata, importo, sconto	3 anni

38



Il carico di lavoro preliminare

- Il riconoscimento di fatti, dimensioni e misure è strettamente collegato all'identificazione di un *carico di lavoro preliminare*.
 - ✓ Oltre che dall'interazione diretta con l'utente, indicazioni al riguardo potranno essere ricavate da un esame della reportistica correntemente in uso in azienda.
 - ✓ In questa fase il carico di lavoro può essere espresso in linguaggio naturale; esso sarà comunque utile per valutare la granularità dei fatti e le misure di interesse, nonché per iniziare ad affrontare il problema dell'aggregazione

39



Il carico di lavoro preliminare

<i>Fatto</i>	<i>Interrogazione</i>
inventario di magazzino	Quantità media di ciascun prodotto presente mensilmente in tutti i magazzini. Prodotti per i quali è stata esaurita la scorta contemporaneamente in tutti i magazzini in almeno un'occasione durante la settimana passata. Andamento giornaliero delle scorte complessive per ciascun tipo di prodotto.
vendite	Quantità totali di ciascun tipo di prodotto vendute durante l'ultimo mese. Incasso totale giornaliero di ciascun negozio. Per un dato negozio, incassi relativi alle diverse categorie di prodotti durante un certo giorno. Riepilogo annuale degli incassi per regione relativamente a un dato prodotto.
linee d'ordine	Quantità totale ordinata annualmente presso un certo fornitore. Importo giornaliero ordinato nell'ultimo mese per un certo tipo di prodotto. Sconto massimo applicato da ciascun fornitore durante l'ultimo anno per ciascuna categoria di prodotto.

40



Altri requisiti

- **Vincoli di progettazione logica e fisica** (spazio disponibile)
- **Progetto dell'alimentazione** (periodicità dell'alimentazione)
- **Architettura del sistema di data warehousing** (tipo di architettura da implementare, numero dei livelli, presenza di data mart dipendenti o indipendenti, materializzazione del livello riconciliato)
- **Applicazioni per l'analisi dei dati** (disamina delle tipologie di interrogazioni e dei rapporti analitici normalmente richiesti)
- **Piano di avviamento**
- **Piano di formazione**

41

Progettazione concettuale



Prof. Stefano Rizzi



Quale formalismo?

- Mentre è universalmente riconosciuto che un DW si appoggia sul modello multidimensionale, non c'è accordo sulla metodologia di progetto concettuale.
- Il modello Entity/Relationship è molto diffuso nelle imprese come formalismo per la documentazione dei sistemi informativi relazionali, ma *non può essere usato per modellare il DW*.

43



Il Dimensional Fact Model

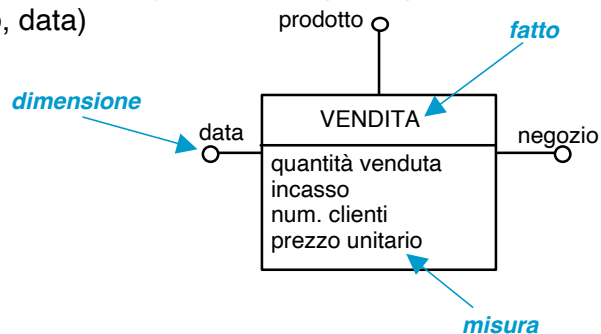
- Il DFM è un modello concettuale grafico per data mart, pensato per:
 - ✓ supportare efficacemente il progetto concettuale;
 - ✓ creare un ambiente su cui formulare in modo intuitivo le interrogazioni dell'utente;
 - ✓ permettere il dialogo tra progettista e utente finale per raffinare le specifiche dei requisiti;
 - ✓ creare una piattaforma stabile da cui partire per il progetto logico (*indipendentemente dal modello logico target*);
 - ✓ restituire una documentazione a posteriori espressiva e non ambigua.
- La rappresentazione concettuale generata dal DFM consiste in un insieme di **scemi di fatto**. Gli elementi di base modellati dagli scemi di fatto sono i fatti, le misure, le dimensioni e le gerarchie

44

II DFM: costrutti di base

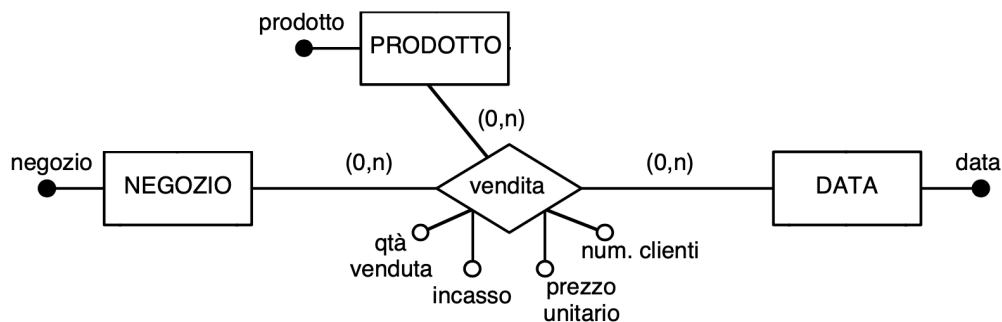
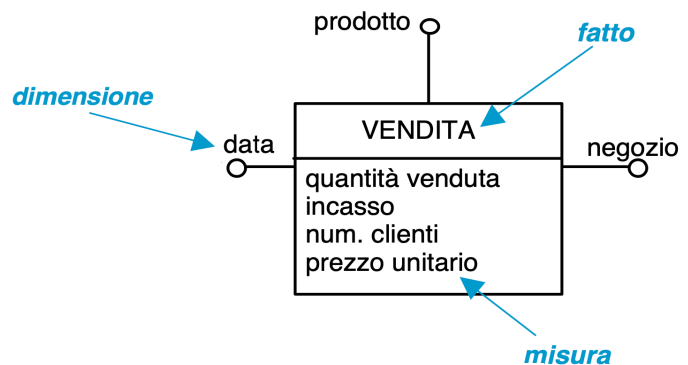
- Un **fatto** è un concetto di interesse per il processo decisionale; tipicamente modella un insieme di eventi che accadono nell'impresa (ad esempio: vendite, spedizioni, acquisti, ...). È essenziale che un fatto abbia aspetti dinamici, ovvero evolva nel tempo
- Una **misura** è una proprietà numerica di un fatto e ne descrive un aspetto quantitativo di interesse per l'analisi (ad esempio, ogni vendita è misurata dal suo incasso)
- Una **dimensione** è una proprietà con dominio finito di un fatto e ne descrive una coordinata di analisi (dimensioni tipiche per il fatto vendite sono prodotto, negozio, data)

Un fatto esprime una
associazione
multi-a-molti
tra le dimensioni



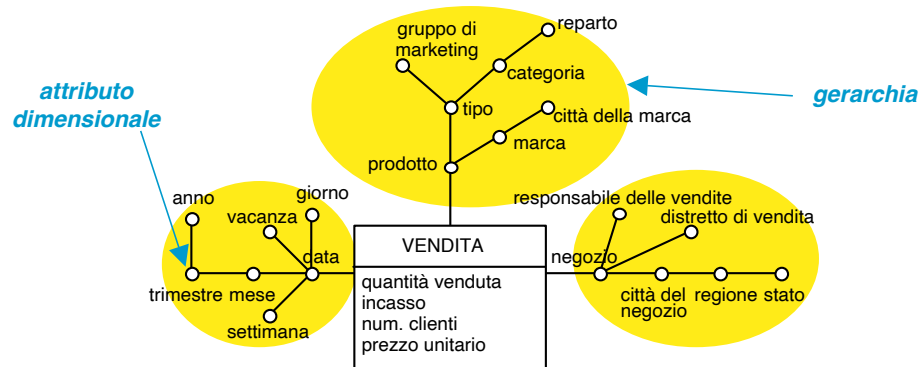
45

Un fatto esprime una
associazione
multi-a-molti
tra le dimensioni



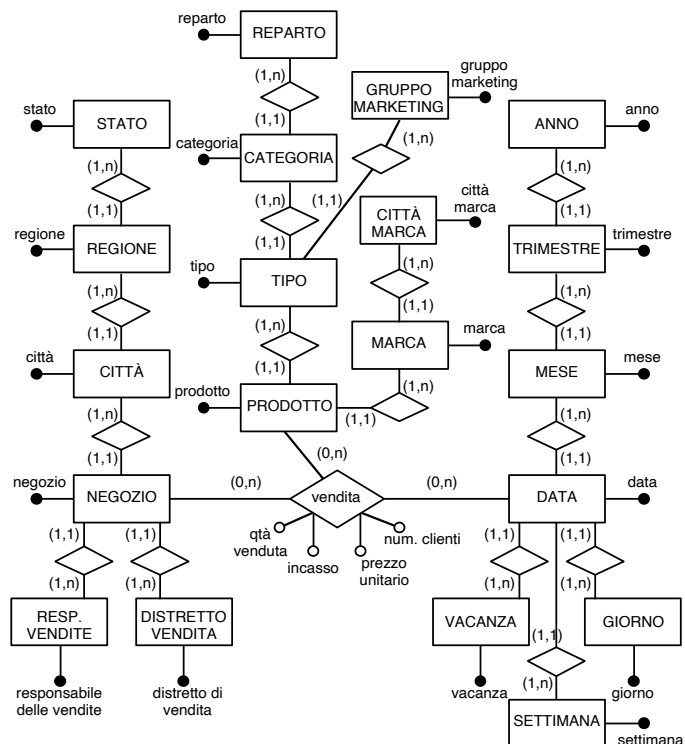
II DFM: costrutti di base

- Con il termine generale **attributi dimensionali** si intendono le dimensioni e gli eventuali altri attributi, sempre a valori discreti, che le descrivono (per esempio, un prodotto è descritto dal suo tipo, dalla categoria cui appartiene, dalla sua marca, dal reparto in cui è venduto)
- Una **gerarchia** è un albero direzionato i cui nodi sono attributi dimensionali e i cui archi modellano associazioni multi-a-uno tra coppie di attributi dimensionali. Essa racchiude una dimensione, posta alla radice dell'albero, e tutti gli attributi dimensionali che la descrivono



46

II DFM: corrispondenza con l'E/R



47



“Naming conventions”

- Tutti gli attributi dimensionali in ciascuno schema di fatto devono avere nomi diversi
- Eventuali nomi uguali devono essere differenziati qualificandoli con il nome di un attributo dimensionale che li precede nella gerarchia
 - ✓ Ad esempio, *warehouse city* è la città in cui si trova un magazzino, mentre *store city* è la città in cui si trova un negozio
- I nomi degli attributi non dovrebbero riferirsi esplicitamente al fatto a cui appartengono
 - ✓ Ad esempio, si evitino *shipped product* e *shipment date*
- Attributi con lo stesso significato in schemi diversi devono avere lo stesso nome

48

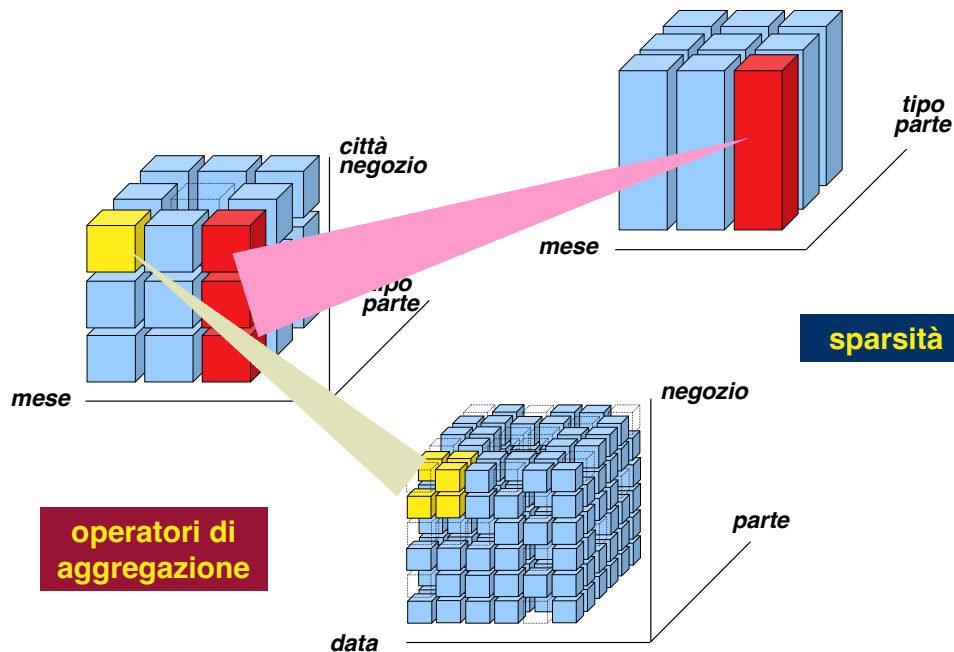


Eventi e aggregazione

- Un *evento primario* è una particolare occorrenza di un fatto, individuata da una ennupla costituita da un valore per ciascuna dimensione. A ciascun evento primario è associato un valore per ciascuna misura
 - ✓ Con riferimento alle vendite, un possibile evento primario registra per esempio che, il 10/10/2001, nel negozio NonSoloPappa sono state vendute 10 confezioni di detersivo Brillo per un incasso complessivo pari a 25 euro
- Dato un insieme di attributi dimensionali (*pattern*), ciascuna ennupla di loro valori individua un *evento secondario* che aggrega tutti gli eventi primari corrispondenti. A ciascun evento secondario è associato un valore per ciascuna misura, che riassume in sé tutti i valori della stessa misura negli eventi primari corrispondenti
 - ✓ Pertanto, le gerarchie definiscono il modo in cui gli eventi primari possono essere aggregati e selezionati significativamente per il processo decisionale; mentre la dimensione in cui una gerarchia ha radice ne definisce la granularità più fine di aggregazione, agli altri attributi dimensionali corrispondono granularità via via crescenti

49

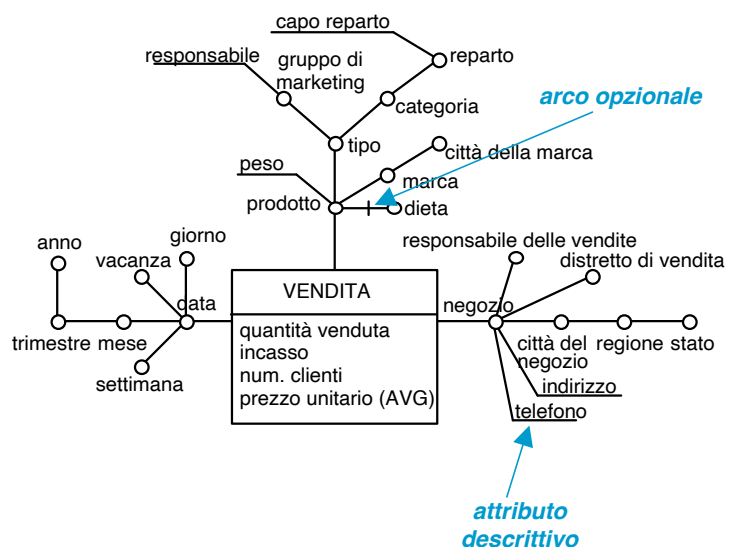
Eventi e aggregazione



50

Il DFM: costrutti avanzati

- Un **attributo descrittivo** contiene informazioni aggiuntive su un attributo dimensionale di una gerarchia, a cui è connesso da una associazione -a-uno. Non viene usato per l'aggregazione poiché ha valori continui e/o poiché deriva da un'associazione uno-a-uno
- Alcuni archi dello schema di fatto possono essere **opzionali**

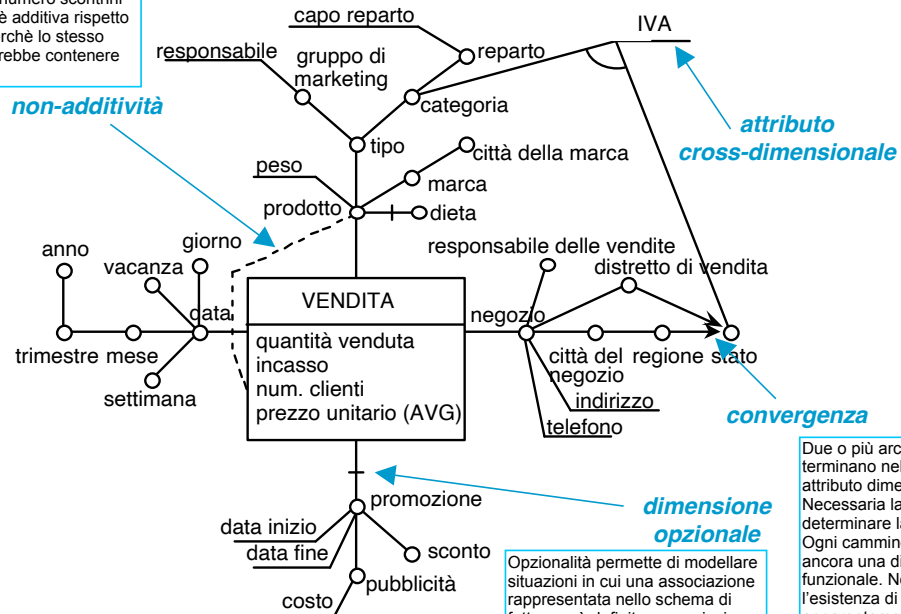


51

II DFM: costrutti avanzati

Num.Clienti (numero scontrini emessi) non è additiva rispetto a prodotto perché lo stesso scontrino potrebbe contenere più prodotti.

non-addittività



Addittività: una misura è detta additiva su una dimensione se i suoi valori possono essere aggregati lungo la corrispondente gerarchia tramite l'operatore somma, altrimenti è non-addittiva.

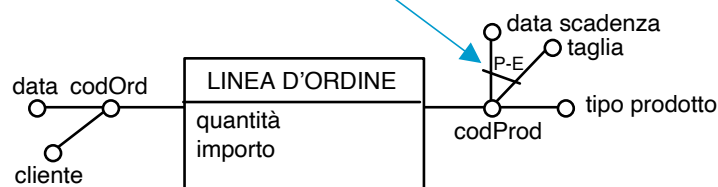
Opzionalità permette di modellare situazioni in cui una associazione rappresentata nello schema di fatto non è definita per un insieme di eventi.
Dieta assume valore solo per i prodotti di tipo alimentare.
Promozione opzionale significa che esisteranno alcuni eventi primari identificati solo da tre dimensioni

Due o più archi che terminano nello stesso attributo dimensionale. Necessaria la freccia per determinare la direzione. Ogni cammino rappresenta ancora una dipendenza funzionale. Non sempre l'esistenza di due attributi apparentemente uguali (città) determina una convergenza.

52

II DFM: costrutti avanzati

copertura di arco opzionale



La copertura è totale se a ciascun valore dell'attributo "a" è sempre abbinato un valore di almeno uno dei figli; se invece esistono valori di "a" per i quali tutti i figli sono indefiniti, allora la copertura è detta parziale.

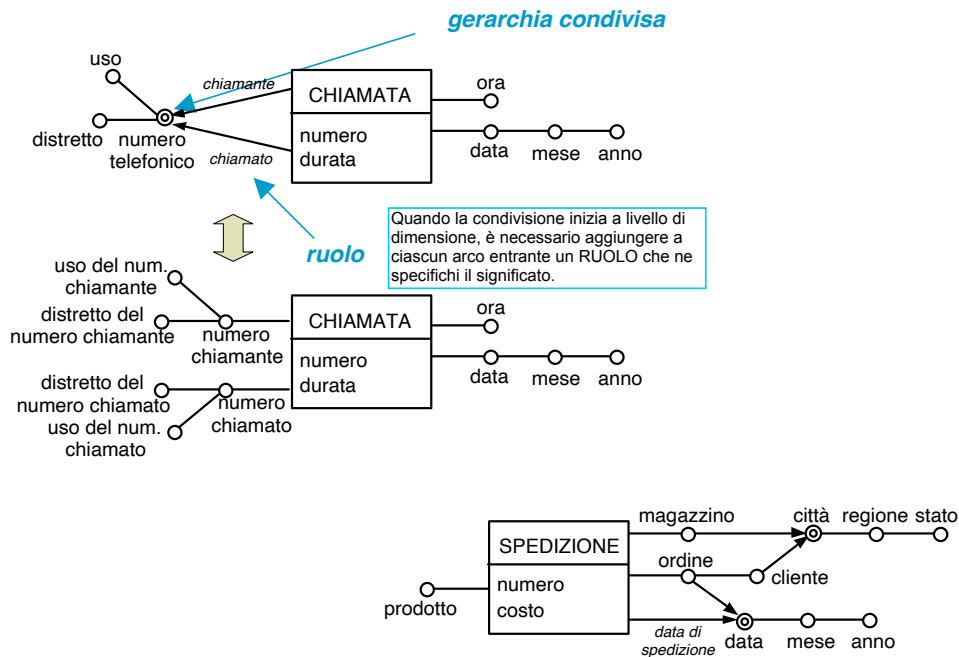
La copertura è esclusiva se in corrispondenza di ciascun valore di "a" si ha al massimo un valore per uno dei figli; se invece esistono valori di "a" abbinati a valori di due o più figli, la copertura è detta sovrapposta.

Quattro tipo di copertura: T-E, T-S, P-E e P-S.

Esempio: supponiamo che i prodotti possano essere di tre tipi: alimentari, abbigliamento e casalinghi. La data di scadenza e la taglia sono definiti solo per alimentari e abbigliamento rispettivamente, quindi la copertura è parziale ed esclusiva.

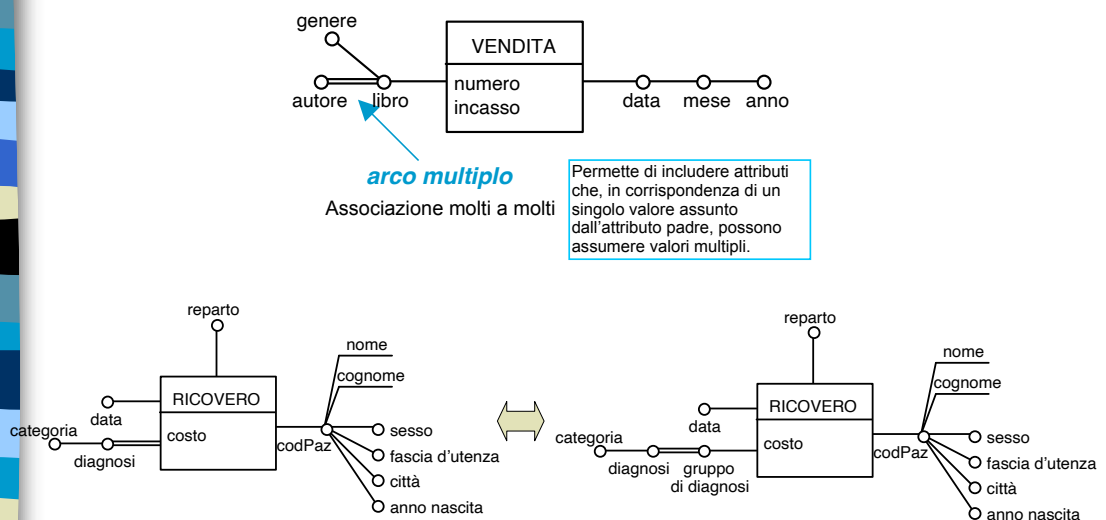
53

II DFM: costrutti avanzati



54

II DFM: costrutti avanzati



55

Additività

- L'aggregazione richiede di definire un operatore adatto per comporre i valori delle misure che caratterizzano gli eventi primari in valori da abbinare a ciascun evento secondario
- Da questo punto di vista è possibile distinguere tre categorie di misure:
 - ✓ **Misure di flusso**: si riferiscono a un periodo, al cui termine vengono valutate in modo cumulativo (il numero di prodotti venduti in un giorno, l'incasso mensile, il numero di nati in un anno)
 - ✓ **Misure di livello**: vengono valutate in particolari istanti di tempo (il numero di prodotti in inventario, il numero di abitanti di una città)
 - ✓ **Misure unitarie**: vengono valutate in particolari istanti di tempo, ma sono espresse in termini relativi (il prezzo unitario di un prodotto, la percentuale di sconto, il cambio di una valuta)

Misura di flusso: Quantità Venduta (di prodotti) è additiva -> La quantità venduta in un mese è la somma delle quantità vendute nei singoli giorni del mese

Misura di livello: Numero Prodotti (in inventario) NON è additiva, ma può essere aggregata con altri operatori. Esempio: numero minimo/massimo di prodotti

Misure di flusso

Misure di livello

Misure unitarie

Misura unitaria: Prezzo Unitario (di prodotti) NON è additiva, ma può essere aggregata con altri operatori. Esempio: media dei prezzi unitari di più prodotti

Gerarchie temporali

Gerarchie non temp orali

SUM, AVG, MIN, MAX

SUM, AVG, MIN, MAX

AVG, MIN, MAX

SUM, AVG, MIN, MAX

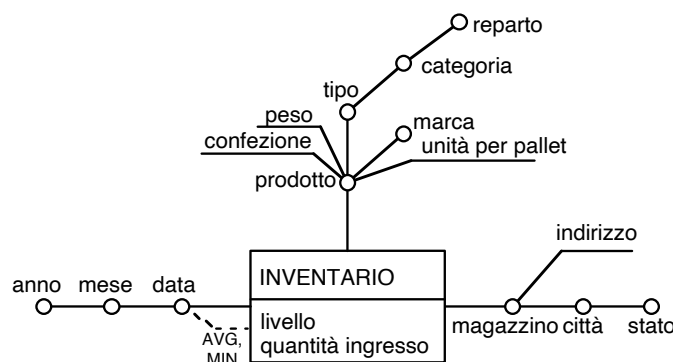
AVG, MIN, MAX

AVG, MIN, MAX

56

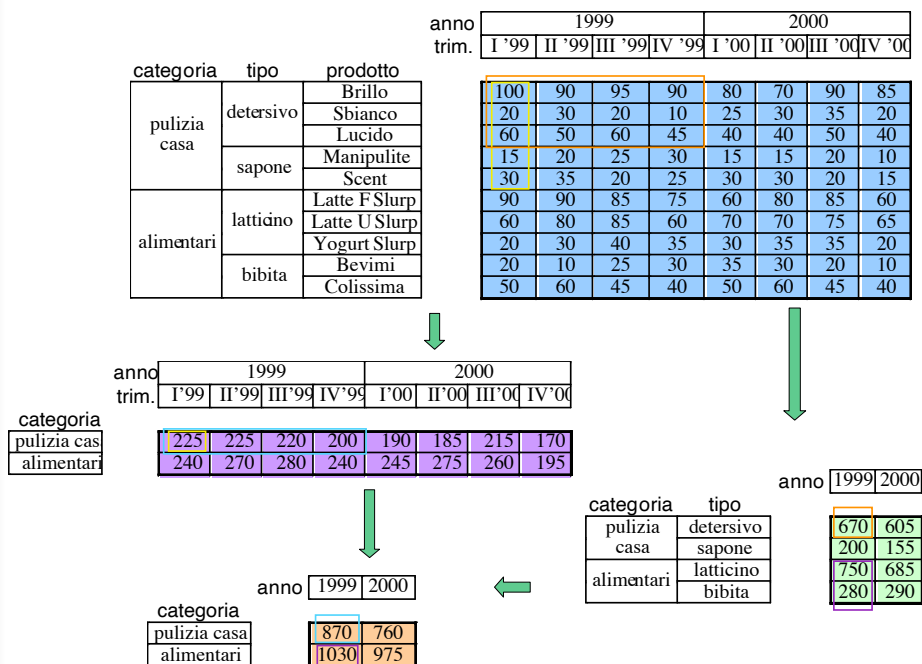
Additività

- Una misura è detta **additiva** su una dimensione se i suoi valori possono essere aggregati lungo la corrispondente gerarchia tramite l'operatore di somma, altrimenti è detta **non-additiva**. Una misura non-additiva è **non-aggregabile** se nessun operatore di aggregazione può essere usato su di essa



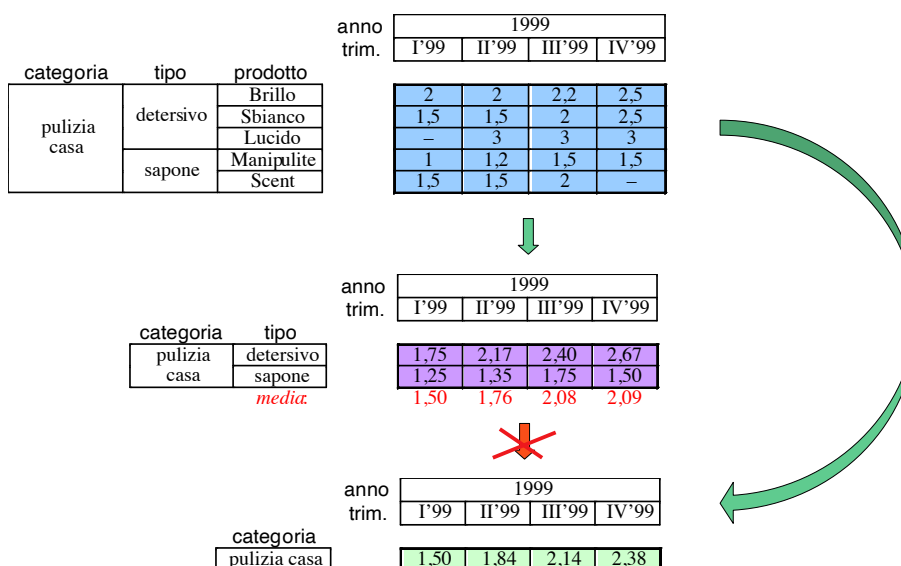
57

Misure additive



58

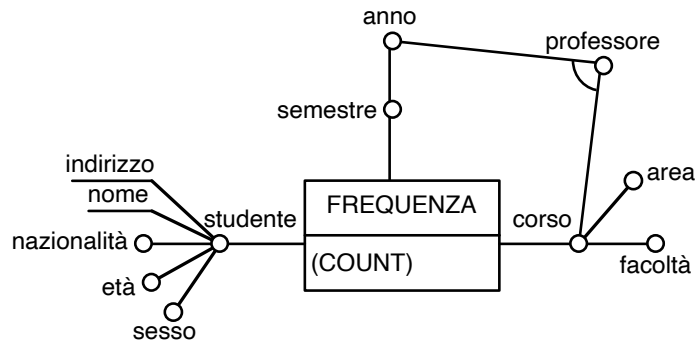
Misure non-additive



59


Schemi di fatto vuoti

- Uno schema di fatto si dice **vuoto** se non ha misure
 - ✓ In questo caso, il fatto registra solo il verificarsi di un evento



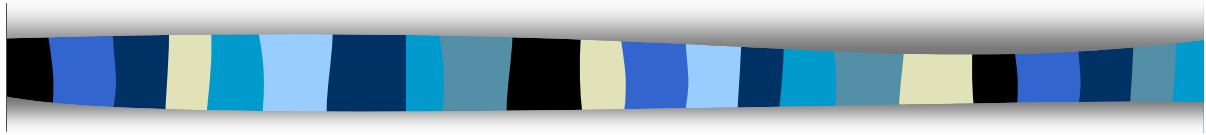
60

Progettazione concettuale: approcci

- Basata sui requisiti
 - ✓ Il progettista deve essere in grado di enucleare, dalle interviste condotte presso l'utente, un'indicazione precisa circa i fatti da rappresentare, le misure che li descrivono e le gerarchie attraverso cui aggregarli utilmente. Il problema del collegamento tra lo schema concettuale così determinato e le sorgenti operazionali viene affrontato in un secondo tempo
- Basata sulle sorgenti 
 - ✓ È possibile definire lo schema concettuale in funzione della struttura delle sorgenti, evitando il complesso compito di stabilire il legame con esse a posteriori. Inoltre, è possibile derivare uno schema concettuale prototipale dagli schemi operazionali in modo pressoché automatico

61

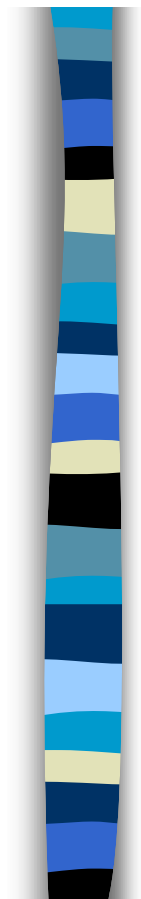
Carico di lavoro e volume dati



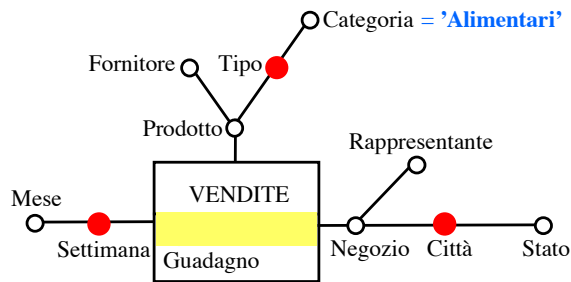
Prof. Stefano Rizzi

Il carico di lavoro

- Il carico di lavoro di un sistema OLAP è per sua natura estemporaneo
- È necessario identificare in fase di progettazione un carico di lavoro di riferimento
 - ✓ Reportistica standard
 - ✓ Colloqui con gli utenti
- Le interrogazioni OLAP sono facilmente caratterizzabili
 - ✓ Pattern di aggregazione
 - ✓ Misure richieste
 - ✓ Clausole di selezione



Il carico di lavoro



VENDITE(Negozio.Città, Settimana, Prodotto.Tipo;
Prodotto.Categoria='Alimentari').Quantità

*Totale della quantità venduta per i diversi tipi di prodotto, in ogni settimana e città
ma solo per i prodotti alimentari*

87

Dinamicità del carico di lavoro

- Il carico di lavoro preliminare non è di per sé sufficiente a ottimizzare le prestazioni del sistema
 - ✓ L'interesse degli utenti cambia nel tempo
 - ✓ Il numero di interrogazioni aumenta al crescere della confidenza degli utenti con il sistema
- Per ottimizzare la struttura logica del data mart è necessaria una fase di tuning attuabile solo dopo che il sistema è stato messo in funzione
- Il carico di lavoro reale può essere desunto dal log delle interrogazioni sottoposte al sistema

88



Il volume dati

- Consiste nelle informazioni necessarie a determinare/stimare la dimensione del data mart.
 - ✓ Numero di valori distinti degli attributi nelle gerarchie
 - ✓ Lunghezza degli attributi
 - ✓ Numero di eventi di ogni fatto
- Deve essere calcolato considerando la quantità di dati necessari a coprire l'intervallo temporale deciso per il data mart.

89



Il volume dati

- È utilizzato sia durante la progettazione logica sia durante la progettazione fisica per determinare:
 - ✓ la dimensione delle tabelle
 - ✓ la dimensione degli indici
 - ✓ i costi di accesso
- La bontà delle stime è spesso compromessa a causa del problema della sparsità.

90

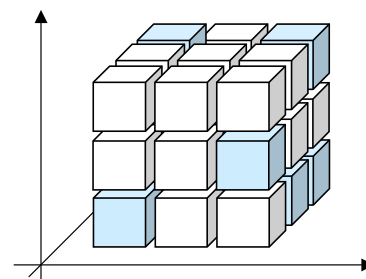
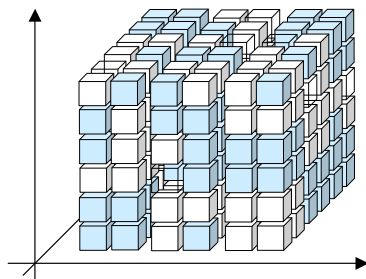
Il problema della sparsità

- Nel modello multidimensionale, a un insieme di coordinate corrisponde un possibile evento anche se questo non è realmente avvenuto
- Normalmente il numero di eventi accaduti è di gran lunga inferiore a quelli possibili
- Tenere traccia degli eventi non accaduti comporta uno spreco di risorse e riduce le prestazioni del sistema
 - ✓ ROLAP: memorizza solo gli eventi accaduti
 - ✓ MOLAP: richiede tecniche complesse per ridurre al minimo lo spazio necessario a tenere traccia degli eventi non accaduti

91

Il problema della sparsità

- La sparsità dei dati inficia le stime sulla cardinalità dei dati aggregati



- La sparsità si riduce all'aumentare del livello di aggregazione dei dati

92