

UNIVERSITÀ DEGLI STUDI DI NAPOLI FEDERICO II

DIPARTIMENTO DI INGEGNERIA ELETTRICA E TECNOLOGIE DELL'INFORMAZIONE



CORSO DI LAUREA MAGISTRALE IN INFORMATICA
A.A. 2017-18

Progetto Basi di Dati II Modulo B

Autori

Bizzarri Flavio
N97000281

Cuomo Daniele
N97000270

Sommario

Si vuole realizzare un data warehouse destinato all'analisi di dati riguardanti misurazioni effettuate su veicoli da parte dell'Istituto Motori di Napoli. Il DW realizzato è di tipo ROLAP (Relational On-Line Analytical Processing), implementato con PostgreSQL, il quale fornisce sia tutte caratteristiche necessarie al data warehousing sia un'estensione spaziale.

Indice

1	Progettazione	1
1.1	Interrogazioni	1
1.2	Diagrammi	2
1.3	ETL	5
1.4	Viste Materializzate	5
2	Analisi	8
2.1	ETL	8
2.1.1	Trasformazione in CSV	8
2.1.2	Import in tabella temporanea	9
2.1.3	Import nello schema	9
2.1.4	Aggiornamento degli indici	10
2.1.5	Aggiornamento delle viste materializzate	10
2.2	Analisi Performance Query	11
2.2.1	Query I	11
2.2.2	Query C	11
2.2.3	Query E	12
2.2.4	Query M	14
2.2.5	Query T	14

Capitolo 1

Progettazione

In questa capitolo sono esposte le scelte progettuali usate come linee guida per l'implementazione di una base di dati orientata al data warehousing.

1.1 Interrogazioni

Nella tabella 1.1 sono riportate le query che si intende sottoporre al sistema. Ogni query ha ad essa associata un numero identificativo ed una descrizione dell'analisi che si intende effettuare.

Query	Obiettivo
I Impatto ambientale medio in corse da 5 km (CO ₂ e NO _x , massa)	Pensata per analizzare varianti della stessa interrogazione su diversi livelli di granularità
C Consumo medio per intervalli di velocità prefissati, su tutto il dataset	Utile all'implementazione e l'analisi di viste materializzate che raggruppano i dati secondo delle fasce di velocità. Le fasce scelte, espresse in km/h, sono le seguenti: 0-50, 50-90, 90-130
E Efficienza dell'auto per intervalli di RPM (rotazioni per minuto)	Altra interessante interrogazione creata allo scopo di sfruttare le viste materializzate. Il dataset fornisce tutti i parametri necessari al calcolo dell'efficienza, o rendimento istantaneo
M Per ogni test, media di NO _x , CO ₂ , Potenza e Velocità	Quest'interrogazione serve a valutare le prestazioni ottenute dall'esecuzione su di un partizionamento verticale con le colonne sparse tra più tabelle
T Media e deviazione standard delle temperature	Quest'interrogazione serve a valutare le prestazioni ottenute dall'esecuzione con le colonne concentrate su di un'unica tabella restituita da un partizionamento verticale

Le query sopra riportate hanno guidato lo sviluppo del sistema in ogni sua fase e hanno permesso di sperimentare differenti tecniche di implementazione di un sistema ROLAP. Una definizione formale delle stesse in linguaggio SQL sarà presentata più avanti nel corso di questo documento.

1.2 Diagrammi

Di seguito è riportato il diagramma UML che rappresenta lo schema dei fatti implementato secondo il modello relazionale ROLAP.

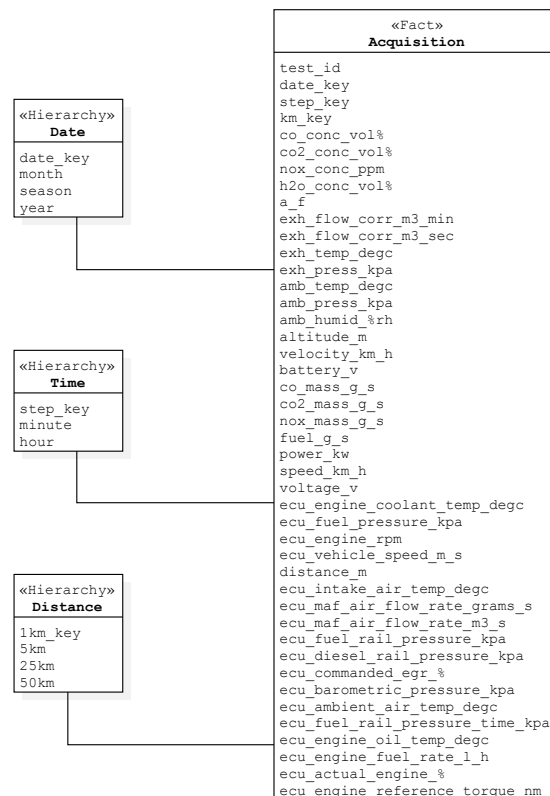


Figura 1.1: Diagramma UML schema dei fatti

È stata inoltre progettata e implementata una variante dello schema proposto che sfrutta la tecnica del partizionamento verticale ovvero la possibilità di dividere la tabella dei fatti in più tabelle, ognuna delle quali rappresenta una particolare sfaccettatura del *fatto*. Il partizionamento viene solitamente adoperato per agevolare quelle interrogazioni che riguardano solo una particolare area di interesse. Per poter applicare questa tecnica è necessario aggiungere una chiave tra le tabelle al fine di poterle ricongiungere per analizzare il dato nella sua completezza.

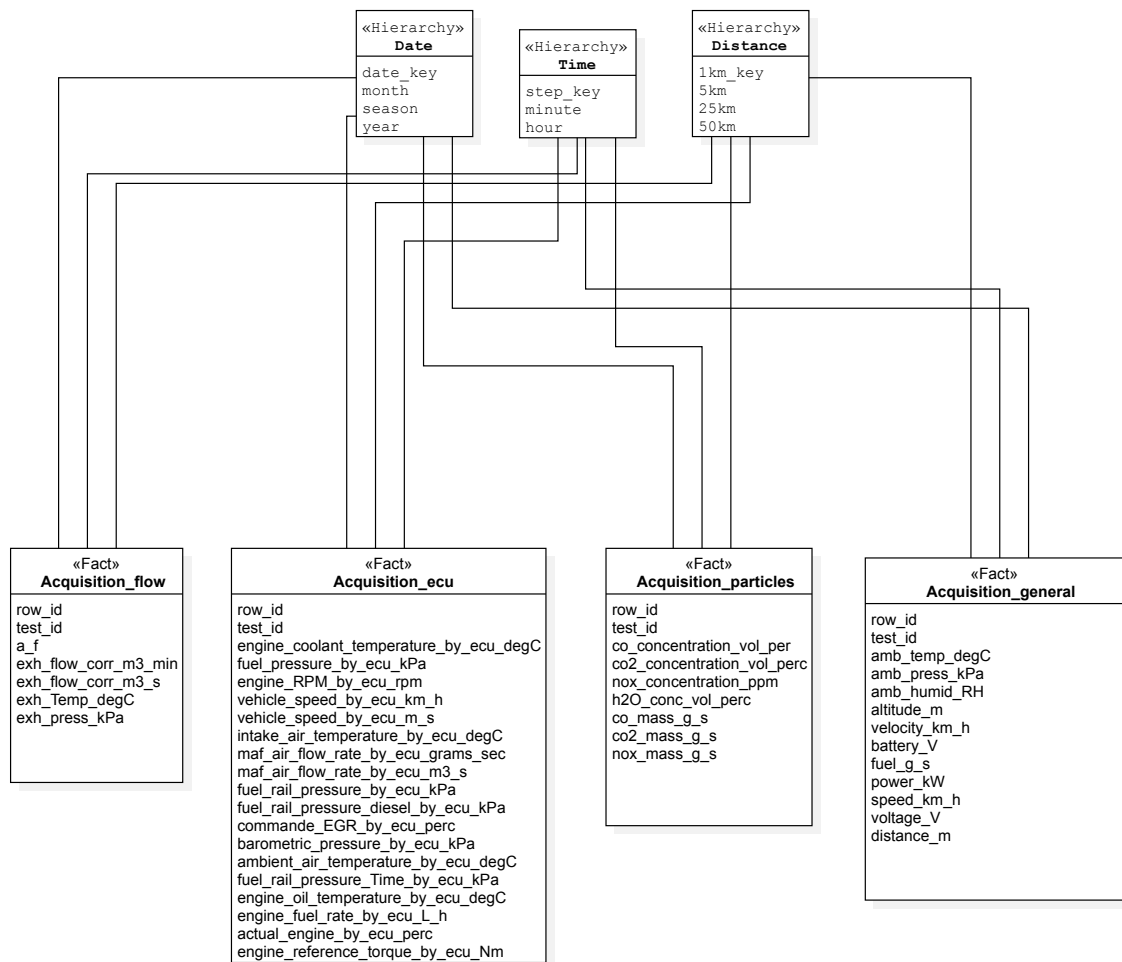


Figura 1.2: Diagramma UML schema dei fatti partizionato verticalmente

1.3 ETL

Al fine di avere un numero rilevante di dati per l'analisi dei tempi è stato implementato un meccanismo di duplicazione dei file. La procedura ETL sviluppata si compone di differenti passaggi descritti nella seguente tabella.

#	Passaggio	Descrizione	Implementazione
1	Trasformazione dei file XLSX in formato CSV	Eliminazione delle righe contenenti dati inconsistenti	Java
		Calcolo dei campi formula	
		Creazione di una data fittizia	
		Creazione di un test_id	
		Salvataggio in formato CSV	
2	Importazione dei dati in tabella provvisoria	Viene sfruttata il comando copy del DBMS per importare i dati in maniera efficiente	PostgreSQL
3	Shutdown degli indici		PostgreSQL
4	Aggiornamento dello schema	Aggiornamento delle dimensioni Aggiornamento tabella dei fatti	PostgreSQL
5	Riattivazione degli indici		PostgreSQL
6	Aggiornamento viste		PostgreSQL

Per righe inconsistenti si intendono tutte quelle righe ove per valori indicanti volumi e/o concentrazioni vi sono valori negativi; inoltre viene dedotta la distanza percorsa ove mancante sfruttando la distanza percorsa e le velocità istantanee rilevate negli istanti precedenti.

1.4 Viste Materializzate

Sono state create le seguenti viste materializzate al fine di implementare efficacemente le query 2 e 3.

Il primo frammento dichiara una vista contenente per ogni corsa (*test_id*) il carburante consumato suddiviso per fasce di velocità.

```
CREATE MATERIALIZED VIEW IF NOT EXISTS fuel_compare_speed AS
SELECT t1.id, t1.fuel_litres as consumo_litri_meno_di_50km_h, t2.fuel_litres
       as consumo_litri_meno_di_90km_h, t3.fuel_litres as
       consumo_litri_meno_di_130km_h
FROM
(SELECT test_id as id, round(SUM((fuel_g_s)/1000)::numeric/0.8,2) AS
       fuel_litres, round(AVG(velocity_kmh)::numeric,1) as average_speed_kmh
FROM acquisition_fact
```



```
where speed_km_h <=50
GROUP BY test_id) t1
JOIN (
SELECT test_id as id, round(SUM((fuel_g_s)/1000)::numeric/0.8,2) AS
    fuel_litres, round(AVG(velocity_km_h)::numeric,1) as average_speed_km_h
FROM acquisition_fact
WHERE speed_km_h >=50 and speed_km_h<=90
GROUP BY test_id
) t2
ON t1.id=t2.id
JOIN (
SELECT test_id as id , round(SUM((fuel_g_s)/1000)::numeric/0.8,2) AS
    fuel_litres, round(AVG(velocity_km_h)::numeric,1) as average_speed_km_h
FROM acquisition_fact
WHERE speed_km_h >=90
GROUP BY test_id
)t3
ON t1.id=t3.id;
```

L'estratto a seguire serve a calcolare una vista materializzata riportante l'efficienza dell'auto nelle diverse corse. I risultati sono stati distribuiti su 3 fasce di rotazioni per minuto (RPM).

```
CREATE MATERIALIZED VIEW IF NOT EXISTS efficiency_compare_rpm AS
SELECT t1.id, round((t1.rendimento*100)::numeric, 2) as
    efficiency_perc_max2000rpm, round((t2.rendimento*100)::numeric,2) as
    efficiency_perc_max3000rpm, round((t3.rendimento*100)::numeric,2) as
    efficiency_perc_over3000rpm
FROM
(SELECT test_id as id, avg(power_kw)/ (avg(fuel_g_s)/1000 * 458000) as
    rendimento
FROM acquisition_fact
WHERE engine_RPM_by_ecu_rpm <2000
GROUP BY test_id) t1
JOIN (
SELECT test_id as id, avg(power_kw)/ (avg(fuel_g_s)/1000 * 458000) as
    rendimento
FROM acquisition_fact
WHERE engine_RPM_by_ecu_rpm >=2000 and engine_RPM_by_ecu_rpm <=3000
GROUP BY test_id
) t2
```

```
ON t1.id=t2.id
JOIN (
SELECT test_id as id, avg(power_kw)/( avg(fuel_g_s)/1000 * 458000) as
    rendimento
FROM acquisition_fact
WHERE engine_RPM_by_ecu_rpm >=3000
GROUP BY test_id
)t3
ON t1.id=t3.id;
```

Capitolo 2

Analisi

In questo capitolo sono riportati i risultati e i tempi ottenuti per ogni fase: dalle procedure ETL all'esecuzione delle query.

2.1 ETL

2.1.1 Trasformazione in CSV

In questa sezione copriremo un'analisi dei tempi ottenuti nella fase che permette, a partire dai dati originali, di ottenere un file in formato CSV contenente i dati ripuliti da inconsistenze e in un formato adatto ad essere importato nel Datawarehouse. In particolare nella fase di pulizia, a partire dal file XSLX, si estraggono unicamente i record dove non appaiono valori negativi per quantità intrinsecamente positive (concentrazione, volume, ecc.); per i record ove manca il valore "relative" e/o la distanza percorsa si procede ad un calcolo a partire dall'ultima rilevazione valida estratta. Inoltre tutte le righe che presentano una chiara assenza di dati (70% delle colonne) vengono scartate. Infine, durante la fase di trasformazione, per ogni file viene generata una data e un test_id e per separare la parte intera da quella decimale si sostituisce la virgola con il punto.

I dati ottenuti si riferiscono ad una media aritmetica ottenuta testando 200 file (~1 milione di righe) generati a partire dal file originario fornito dall'Istituto Motori di Napoli.

Righe XSLX	Righe CSV	Righe perse	Peso file XSLX	Peso file CSV
5764	5706	1%	2.138KB	2.283KB

Fase di pulizia	Fase di trasformazione	Tempo totale
3.04s	0,05s	3,05s

Questa fase mette in luce la buona qualità dei dati forniti dall'istituto: ci si aspetta in media di perdere pochissimi dati a causa di inconsistenze. Per quanto riguarda invece

le dimensioni si nota come il formato Comma-separated values sia leggermente meno efficiente nella compressione rispetto al formato proprietario di Microsoft® ma al tempo stesso permetta una più veloce elaborazione.

2.1.2 Import in tabella temporanea

Al fine di importare i file CSV nella tabella dei fatti ci si appoggia ad una tabella temporanea al fine di agevolare le operazioni successive. Questa tabella viene troncata alla fine della procedura. L'import sfrutta la funzione *COPY* messa a disposizione dal DBMS PostgreSQL.

Numero file	Numero righe	Tempo
1	5706	214ms
50	285.300	8,428s
100	570.600	16.522s

Tabella 2.1: I dati si riferiscono ad una media di 5 esecuzioni

I risultati mostrano la bontà della funzione *COPY* che sfruttando un inserimento batch di 1000 righe per volta riesce ad abbattere in modo consistente i tempi.

2.1.3 Import nello schema

In questa fase i dati, precedentemente inseriti in una tabella temporanea, vengono travasati nello schema proposto dopo aver disabilitato tutti gli indici. Le prove effettuate prendono in considerazione diverse dimensioni della tabella dei fatti oltre al caso in cui lo schema risulti partizionato.

Dimensione tab. fatti	Numero righe importate	Partizionamento	Tempo
0	285.300	No	3,925s
0	285.300	Si	7,753s
570.600	285.300	No	3,730s
570.600	285.300	Si	7,404s
1.141.200	285.300	No	4,535s
1.141.200	285.300	Si	8,572s

Tabella 2.2: I dati si riferiscono ad una media di 5 esecuzioni

Dalla tabella emerge come il partizionamento comporti quasi un raddoppio del tempo necessario all'inserimento: ciò è dovuto al dover spalmare un singolo record della tabella temporanea su 4 differenti tabelle. Questo slow-down è risolvibile pensando a 4 inserimenti in parallelo: infatti le 4 tabelle, anche se logicamente collegate, durante l'inserimento non necessitano di condividere alcuna informazione.

2.1.4 Aggiornamento degli indici

In questa fase vengono riattivati gli indici delle chiavi tra tabella dei fatti e dimensioni. Questi indici sono assolutamente necessari per velocizzare le query ma possono essere disabilitati durante la fase di update dello schema.

Dimensione tab. fatti	Partizionamento	Tempo
285.300	No	2,530s
285.300	Si	11,534s
570.600	No	6,194s
570.600	Si	24,715s
1.141.200	No	16,903s
1.141.200	Si	67,469s
1.426.500	No	24,141s
1.426.500	Si	98,218s

Tabella 2.3: I dati si riferiscono ad una media di 5 esecuzioni

Dai dati sopra mostrati emerge ancora una volta come l'introduzione di un partizionamento verticale comporti un notevole rallentamento. In particolare la forbice tra i tempi registrati aumenta all'aumentare della dimensione dei fatti. Ciò induce a pensare attentamente all'introduzione di un partizionamento in fase di progettazione valutando il rapporto costo/benefici.

2.1.5 Aggiornamento delle viste materializzate

A valle dell'inserimento dei nuovi record si rende necessario l'aggiornamento delle viste materializzate utilizzate per velocizzare le query 3 e 4.

Dimensione tab. fatti	Efficiency_compare_rpm_mv	Fuel_compare_speed_mv
285.300	0,307s	0,238s
570.600	0,983s	0,898s
1.141.200	2,464s	2,173s
1.426.500	3,004s	2,863s

Tabella 2.4: I dati si riferiscono ad una media di 5 esecuzioni

I tempi osservati mostrano come un raddoppio del numero di record comporti un più che raddoppio del tempo necessario all'aggiornamento delle viste. Notevole è infatti l'incremento di tempo per materializzare la vista con 1.141.200 record: 2.5 volte quello necessario per materializzare la vista con la metà dei record.

2.2 Analisi Performance Query

In questa sezione verranno analizzate le query implementate nelle differenti implementazioni e presentati gli snippet di codice SQL

2.2.1 Query I

Impatto ambientale medio in corse da 5 km

A seguire sono proposte due diverse formulazioni della stessa interrogazione. Il primo estratto chiede al sistema di recuperare l'appartenenza ad una certa tratta (di 5 km) dalla tabella *distance_hierarchy*. Nella seconda forma la tratta viene calcolata, eseguendo quindi un'operazione di roll-up implicito.

Roll-up esplicito

```
SELECT F.test_id, H.km_5, AVG(F.co2_mass_g_s), AVG(F.nox_mass_g_s)
FROM acquisition_fact F INNER JOIN distance_hierarchy H
    ON (F.km_key = H.km_key)
GROUP BY F.test_id, H.km_5
```

Roll-up implicito

```
SELECT F.test_id, FLOOR(F.km_key/5), AVG(F.co2_mass_g_s), AVG(F.nox_mass_g_s)
FROM acquisition_fact F
GROUP BY F.test_id, FLOOR(F.km_key/5)
```

Esito

L'estrazione dei tempi di esecuzione delle due query può fornire una misura di utilità di una dimensione. In questo caso è interessante notare come un *fatto* possa essere ricondotto alla sua dimensione di appartenenza attraverso un calcolo richiedente 2 operazioni elementari.

2.2.2 Query C

Consumo medio per intervalli di velocità prefissati

L'utilizzo della vista materializzata permette di abbattere i tempi di esecuzione della query fino a un trecentesimo: in questo modo il costo di aggiornamento della vista viene immediatamente ammortizzato.

Dimensione tab. fatti	Senza vista	Con vista
570.600	0,905s	0,020
1.141.200	2,286s	0,026s
1.426.500	2,925s	0.028s

Tabella 2.5: Test effettuati su schema non partizionato

```
SELECT *
FROM fuel_compare_speed_mv;
```

```
SELECT t1.id, t1.fuel_litres as consumo_litri_meno_di_50km_h, t2.fuel_litres
      as consumo_litri_meno_di_90km_h, t3.fuel_litres as
      consumo_litri_meno_di_130km_h
FROM
(SELECT test_id as id, round(SUM((fuel_g_s)/1000)::numeric/0.8,2) AS
      fuel_litres, round(AVG(velocity_kmh)::numeric,1) as average_speed_kmh
FROM acquisition_fact
where speed_kmh <=50
GROUP BY test_id) t1
JOIN (
SELECT test_id as id, round(SUM((fuel_g_s)/1000)::numeric/0.8,2) AS
      fuel_litres, round(AVG(velocity_kmh)::numeric,1) as average_speed_kmh
FROM acquisition_fact
where speed_kmh >=50 and speed_kmh<=90
GROUP BY test_id
) t2
ON t1.id=t2.id
JOIN (
SELECT test_id as id , round(SUM((fuel_g_s)/1000)::numeric/0.8,2) AS
      fuel_litres, round(AVG(velocity_kmh)::numeric,1) as average_speed_kmh
FROM acquisition_fact
where speed_kmh >=90
GROUP BY test_id
)t3
ON t1.id=t3.id;
```

2.2.3 Query E

Efficienza dell'auto per intervalli di RPM

Si conferma quanto osservato nel caso della *query 2*: utilizzare una vista materializzata permette di abbattere i tempi di esecuzione della query e il costo del refresh risulta ampiamente ammortizzato.

Dimensione tab. fatti	Senza vista	Con vista
570.600	0,937s	0,018s
1.141.200	2,182s	0,028s
1.426.500	2,892s	0,030s

Tabella 2.6: Test effettuati su schema non partizionato

```
select *
from efficiency_compare_rpm_mv;
```

```
SELECT t1.id, round((t1.rendimento*100)::numeric, 2) as
    efficiency_perc_max2000rpm, round((t2.rendimento*100)::numeric,2) as
    efficiency_perc_max3000rpm, round((t3.rendimento*100)::numeric,2) as
    efficiency_perc_over3000rpm
FROM
    (SELECT test_id as id, avg(power_kw)/ (avg(fuel_g_s)/1000 * 458000) as
        rendimento
    FROM acquisition_fact
    where engine_RPM_by_ecu_rpm <2000
    GROUP BY test_id) t1
JOIN (
    SELECT test_id as id, avg(power_kw)/ (avg(fuel_g_s)/1000 * 458000) as
        rendimento
    FROM acquisition_fact
    where engine_RPM_by_ecu_rpm >=2000 and engine_RPM_by_ecu_rpm <=3000
    GROUP BY test_id
    ) t2
ON t1.id=t2.id
JOIN (
    SELECT test_id as id, avg(power_kw)/( avg(fuel_g_s)/1000 * 458000) as
        rendimento
    FROM acquisition_fact
    where engine_RPM_by_ecu_rpm >=3000
    GROUP BY test_id
```



```
)t3  
ON t1.id=t3.id;
```

2.2.4 Query M

Per ogni test, media di NOx, CO2, Potenza e Velocità

TODO. Qui la query dovrebbe essere più veloce senza il partizionamento.

Dimensione tab. fatti	Partizionamento	Tempo
570.600	No	6,194s
570.600	Si	24,715s
1.141.200	No	16,903s
1.141.200	Si	67,469s
1.426.500	No	24,141s
1.426.500	Si	98,218s

NO Partizionamento

SI partizionamento

2.2.5 Query T

Media e deviazione standard delle temperature

TODO. Qui la query dovrebbe essere più veloce con il partizionamento.

Dimensione tab. fatti	Partizionamento	Tempo
570.600	No	6,194s
570.600	Si	24,715s
1.141.200	No	16,903s
1.141.200	Si	67,469s
1.426.500	No	24,141s
1.426.500	Si	98,218s

NO Partizionamento

SI partizionamento
