Data management and visualization

Iniziato	lunedì, 20 febbraio 2023, 16:24
Stato	Completato
Terminato	lunedì, 20 febbraio 2023, 16:24
Tempo impiegato	8 secondi
Valutazione	0,00 su un massimo di 31,00 (0 %)



Are the values encoded in a uniformly proportional way?

Risposta non data

Punteggio max.: 4,00

```
Administrations(<u>VaxID</u>, <u>PatientID</u>, <u>TimeID</u>, <u>ClinicID</u>,
NumVaccinations, TotalCost
WelfareCoveredCost)
Vaccine(<u>VaxID</u>, BatchNumber,VaccineName, AssociatedPathology)
Patient(<u>PatientID</u>, AgeRange, AutoimmuneDisease, HeartDisease,
NeurologicalPathologies)
Clinic (<u>ClinicID</u>, ClinicName, Province, Region)
Time(<u>TimeID</u>, Month, 3m, 4m, 6m, year)
```

The data warehouse collects the total number and cost of vaccinations administered to patients over time. The clinic name is supposed to be unique for each clinic.

Separately for each clinic in the Piedmont region and for each month, select:

- A. the average cost per vaccination
- B. the ratio of total monthly cost with respect to the total cost of the year
- C. the moving average of the monthly number of vaccinations over the last 5 months

Write the requested SQL query.

Risposta non data

Punteggio max.:

0,25



Question

Which one of the following questions represents the purpose of this visualization?

- (a) What is the trend of the Chicago population over several years?
- (b) What is the average number of murders in Chicago in the last decade?
- (c) What is the year with the highest number of murders in Chicago?
- (d) What is the number of murders per year in Chicago?
- (e) What is the correlation between the number of murders and the population of Chicago over several years?

Risposta errata.

La risposta corretta è: What is the correlation between the number of murders and the population of Chicago over several years?

Risposta non data

Punteggio max.:

0,50

Design a data warehouse to analyse the purchasing behaviour of consumer electronic devices manufactured by a tech company, in terms of average price per product, and the total number of products sold.

- The tech company manufactures and sells different electronic **devices**
 - All devices have a model and a series. Each model belongs to one series only. Each series collects different models. Each series belongs to one of the following three lines: "basic", "mid-range", and "premium".
 - Device models can be divided into the following types: "smart-phone", "smart-watch", or "accessory".
 - For each model, the following features must be recorded:
 - camera (present/absent),
 - 5G (present/absent),
 - warranty period (1, 2, or 3 years)
- The **products** sold by the tech company are manufactured in and shipped to different places around the world.
 - Both manufacturing places and shipping destinations must be analysed in terms of commercial area, continent and country.
 - A commercial area includes a group of countries, possibly belonging to different continents. Each country belongs to one commercial area only.
- The **consumer** buying the products is characterized by
 - Her/his age, in one of the following ranges: 18-30, 31-45, 45-65, or 65+ years
 - Loyalty program enrolment: yes/no
- The analysis must be performed for each **date**, day of the week, holiday (yes/no), month, 4-month period, semester, and year.

Select, among the following dimensions, those that meet the requirements described in the problem specification (at most one answer is correct).









Risposta errata.

La risposta corretta è: None of the others is correct.

Risposta non data

Punteggio max.:

0,50



Visual Clarity

Are the data in the graph clearly identifiable and understandable (properly described)?

```
Domanda 6
```

Risposta non data

Punteggio max.: 3,00

The following document structure represents a document of a movie. Each document collects some information about the movie directors and the movie release in each country.

```
{
_id: ObjectId(),
title: <string>,
abstract:<string>,
keywords: [ <string>, <string>, ... ],
category: [ <string>, <string>, ... ],
directors: [{
     _id: ObjectId(),
       name: <string>,
       surname: <string>,
       nationality: <string>}
...
],
releases: [
       {date: <date>,
        country: <string>
       },
        ...
        ],
review_score: <float>
}
```

Considering only the movies released in Italy and belonging to the category "comedy", separately for each nationality of the director, compute the average review score and the total number of movies.

Sort the results in descending order according to the average review score.

```
db.collection_name.aggregate([
{$match: {
    "category": "comedy",
    "releases.country": "Italy"
}
,
{$unwind: '$directors},
{$group: {
    '_id': '$directors.nationality,
    'avg_score': {'$avg': '$review_score'},
    'tot': { '$sum: 1}
}
},
{$sort: { avg_score: -1} }
])
```

Risposta non data

Punteggio max.: 2,00

The following document structure represents a document of a movie. Each document collects some information about the movie directors and the movie release in each country.

{ _id: ObjectId(), title: <string>, abstract:<string>, keywords: [<string>, <string>, ...], category: [<string>, <string>, ...], directors: [{ _id: ObjectId(), name: <string>, surname: <string>, nationality: <string>} . . .], releases: [{date: <date>, country: <string> }, ...], review_score: <float> }

Write a MongoDB query to find all the movies in the category "Fantasy" having George Lucas among the directors and with a review score higher than 4. Show only the title and the keywords.

```
db.shows.find(
  { categories: "Fantasy",
    directors: { $elemMatch: { name: "George", surname: "Lucas" } },
    review_score: {$gt: 4} },
    {title: 1, keywords: 1, _id:0}
)
```

Risposta non data

Punteggio max.:

1,50

Which one of the following elements should **not** be removed from a visualization according to the guidelines for visual utility?

- (a) Background shades
- (b) Separated legend
- (c) Heavy grid lines
- (d) Decorative colors
- (e) Logos and other pictures

Risposta errata.

La risposta corretta è: Separated legend

Domanda 9

Risposta non data

Punteggio max.: 4,00 Administrations(<u>VaxID</u>, <u>PatientID</u>, <u>TimeID</u>, <u>ClinicID</u>, NumVaccinations, TotalCost WelfareCoveredCost) Vaccine(<u>VaxID</u>, BatchNumber,VaccineName, AssociatedPathology) Patient(<u>PatientID</u>, AgeRange, AutoimmuneDisease, HeartDisease, NeurologicalPathologies) Clinic (<u>ClinicID</u>, ClinicName, Province, Region) Time(TimeID, Month, 3m, 4m, 6m, year)

The data warehouse collects the total number and cost of vaccinations administered to patients over time. The clinic name is supposed to be unique for each clinic.

Considering the year 2022, separately for each **vaccine name**, **associated pathology**, and **month**:

A. the total number of vaccinations

B. the percentage of monthly vaccinations with respect to all the monthly vaccines with the same associated pathology

C. assign an increasing rank to each vaccine for each month, according to its total number of monthly vaccinations

Write the requested SQL query.

```
SELECT VaccineName, month,
    SUM(NumVaccinations) as A,
    100 * SUM(NumVax) / SUM(SUM(NumVax)) OVER (PARTITION BY
month, AssociatedPathology) as B
    RANK() OVER (PARTITION BY month, AssociatedPathology ORDER BY
SUM(NumVaccinations) DESC) as C
FROM Administrations A, Vaccine V, Time T
WHERE A.VaxID=V.VaxID and T.TimeId=A.TimeId and year=2022
GROUP BY VaccineName, AssociatedPathology, month
```

Risposta non data

Punteggio max.:

0,50

Design a data warehouse to analyse the purchasing behaviour of consumer electronic devices manufactured by a tech company, in terms of average price per product, and the total number of products sold.

- The tech company manufactures and sells different electronic **devices**
 - All devices have a model and a series. Each model belongs to one series only. Each series collects different models. Each series belongs to one of the following three lines: "basic", "mid-range", and "premium".
 - Device models can be divided into the following types: "smart-phone", "smart-watch", or "accessory".
 - For each model, the following features must be recorded:
 - camera (present/absent),
 - 5G (present/absent),
 - warranty period (1, 2, or 3 years)
- The **products** sold by the tech company are manufactured in and shipped to different places around the world.
 - Both manufacturing places and shipping destinations must be analysed in terms of commercial area, continent and country.
 - A commercial area includes a group of countries, possibly belonging to different continents. Each country belongs to one commercial area only.
- The **consumer** buying the products is characterized by
 - Her/his age, in one of the following ranges: 18-30, 31-45, 45-65, or 65+ years
 - Loyalty program enrolment: yes/no
- The analysis must be performed for each date, day of the week, holiday (yes/no), month, 4-month period, semester, and year.

Select, among the following dimensions, those that meet the requirements described in the problem specification (at most one answer is correct).

(a) holiday (Y/N) 6M year Devices sold date month 4M dayOfTheWeek holiday (Y/N)







Risposta non data

Punteggio max.:

1,50

Design a data warehouse to analyse the purchasing behaviour of consumer electronic devices manufactured by a tech company, in terms of average price per product, and the total number of products sold.

- The tech company manufactures and sells different electronic devices
 - All devices have a model and a series. Each model belongs to one series only. Each series collects different models. Each series belongs to one of the following three lines: "basic", "mid-range", and "premium".
 - Device models can be divided into the following types: "smart-phone", "smart-watch", or "accessory".
 - For each model, the following features must be recorded:
 - camera (present/absent),
 - 5G (present/absent),
 - warranty period (1, 2, or 3 years)
- The **products** sold by the tech company are manufactured in and shipped to different places around the world.
 - Both manufacturing places and shipping destinations must be analysed in terms of commercial area, continent and country.
 - A commercial area includes a group of countries, possibly belonging to different continents. Each country belongs to one commercial area only.
- The **consumer** buying the products is characterized by
 - Her/his age, in one of the following ranges: 18-30, 31-45, 45-65, or 65+ years
 - Loyalty program enrolment: yes/no
- The analysis must be performed for each date, day of the week, holiday (yes/no), month, 4-month period, semester, and year.

Select all and only the required measures of the fact table in the conceptual schema design among the following (multiple-choice question). Hint: do consider the dimensions defined by the previous answers.

Scegli una o più alternative:

- (a) Total income (euros)
- (b) Average number of consumers per region (count)
- (c) Total number of products (count)
- (d) Average income per factory (euros)
- (e) Total number of orders (count)
- (f) Total warranty period (time)
- (g) Total number of consumers (count)
- (h) Average number of factories per country (count)
- (i) Average number of products per type (count)

- (j) Total number of available devices (count)
- (k) Maximum number of orders (count)
- (I) Average warranty period per product (time)
- (m) Total number of factories (count)

Risposta errata.

La risposta corretta è: Total income (euros), Total number of products (count)

Domanda 12

Risposta non data

Non valutata

This is a blank question to be used as your personal notepad during the exam.

Anything written here will NOT be evaluated.



```
Fact(AutID, BuyerID, TimeID, JID, TotalIncome,
TotalArtworksSold)
Junk(JID, PriceRange, Collection)
Author(AutID, ArtisticMovement, Alive, Nation, Region)
Buyer(PatientID, Type, ReliabilityIndex, Nation, Region)
Time(TimeID, Month, 4m, 6m, year)
```

PriceRange and Collection as 2 separate tables are also acceptable

Risposta non data

Punteggio max.:

1,50

The schema versioning pattern has the advantage of:

- (a) none of the answers is correct
- (b) controlling the schema migration
 - (c) avoiding join operations
- (d) reducing the required indexes during migration

Risposta errata.

La risposta corretta è: controlling the schema migration



Risposta non data

Punteggio max.:

0,25



Design data

Design the visualization based on the following data structure.

MURDERS	Scegli	~
POPULATION	Scegli	*
YEAR	Scegli	~

Risposta errata.

La risposta corretta è: MURDERS – Measure, POPULATION – Measure, YEAR – Dimension

Risposta non data

Punteggio max.:

0,75



Visual Utility

All the elements in the graph convey useful information?



```
) (c)
```

```
db.videos.aggregate([
    {$match: { releaseDate : {$gte: new Date ("2022-01-01"), $lte: new Date ("20
    22-12-31") } } },
    {$unwind: "$tags"},
         {$group: {
                    _id: {
    category: "$tags"
        },
               n: {$sum: 1}
                    }
               },
         {$match: { n: { $gte: 1000 } } }
    ])
(d) none of the other answers is correct
(e)
    db.videos.aggregate([
    {$match: { releaseDate : {$gte: new Date ("2022-01-01"), $lte: new Date ("20
    22-12-31") } } },
    {$unwind: "$tags"},
         {$group: {
                    _id: {
    category: "$tags"
        },
               n: {$sum: 1}
                    }
               }
         }
    ])
(f)
    db.videos.aggregate([
          {$group: {
                    _id: {
    category: "$tags"
        },
               n: {$sum: 1}
                    }
               },
    {$match: { releaseDate : {$gte: new Date ("2022-01-01"), $lte: new Date ("20
    22-12-31")}, n: { $gte: 1000 } } }
    ])
```

```
Risposta errata.
```

La risposta corretta è:

```
db.videos.aggregate([
{$match: { releaseDate : {$gte: new Date ("2022-01-01"), $lte: new Date ("2022-1
2-31") } },
{$unwind: "$tags"},
{$group: {
    __id: {
    category: "$tags"
    },
    n: {$sum: 1}
        }
    },
    {$match: { n: { $gte: 1000 } } }
])
```

Risposta non data

Punteggio max.:

1,25



Design schema & Sketch

Fill in the required schema elements; formulas can be used if required. Then describe in words the design proposal.

Risposta non data

Punteggio max.:

1,25



Data

Is the data quality appropriate? Select true answers only.

Scegli una o più alternative:

- (a) Precision is not appropriate as murders shoud be represented with more decimal digits.
- (b) Precision is appropriate for the task as the population is reported in millions with a decimal digit.
- (c) The visualization is very clear because the labels are associated with the bars.
- (d) Data is not complete because several years are missing.
- (e) Data is not credible because no source is reported.
- (f) Data about murders is credible because the source is the Chicago Police Department (CPD).
- (g) Data is consistent because some decades are compared with single years.
- (h) The number of murders and the population are accurate data points because they are represented as percentages.
- (i) Data is reasonably updated for the task because the last data point is from 2021.
- (j) Data is complete because starting from 2010 all years are reported.

Risposta errata.

La risposta corretta è: Data is not complete because several years are missing., Data is reasonably updated for the task because the last data point is from 2021., Precision is appropriate for the task as the population is reported in millions with a decimal digit.,

Data about murders is credible because the source is the Chicago Police Department (CPD).

Risposta non data

Punteggio max.:

0,50

Design a data warehouse to analyse the purchasing behaviour of consumer electronic devices manufactured by a tech company, in terms of average price per product, and the total number of products sold.

- The tech company manufactures and sells different electronic **devices**
 - All devices have a model and a series. Each model belongs to one series only. Each series collects different models. Each series belongs to one of the following three lines: "basic", "mid-range", and "premium".
 - Device models can be divided into the following types: "smart-phone", "smart-watch", or "accessory".
 - For each model, the following features must be recorded:
 - camera (present/absent),
 - 5G (present/absent),
 - warranty period (1, 2, or 3 years)
- The **products** sold by the tech company are manufactured in and shipped to different places around the world.
 - Both manufacturing places and shipping destinations must be analysed in terms of commercial area, continent and country.
 - A commercial area includes a group of countries, possibly belonging to different continents. Each country belongs to one commercial area only.
- The **consumer** buying the products is characterized by
 - Her/his age, in one of the following ranges: 18-30, 31-45, 45-65, or 65+ years
 - Loyalty program enrolment: yes/no
- The analysis must be performed for each **date**, day of the week, holiday (yes/no), month, 4-month period, semester, and year.

Select, among the following dimensions, those that meet the requirements described in the problem specification (at most one answer is correct).









Risposta errata.

La risposta corretta è:



Risposta non data

Punteggio max.: 4,00 Design a MongoDB database to collect restaurant menus and dish recipes according to the following requirements.

The dish recepits are characterized by a name, a description, the level of complexity (from 1 to 5), and the chef who proposed the dish. For the chef, the name, country of origin, website address (if available), and email address (if available) are known.

The ingredients used in each recipe must be recorded. Each ingredient consists of the name, the quantity, and (if necessary) the unit of measure. For example, a cake may use 500 ml of milk as an ingredient. In that case, the ingredient name is "milk", the quantity is 500, and the unit of measure is "ml". Each ingredient can be present only once in a recipe.

The menus proposed by different restaturants must be recorded. Each menu is characterized by the name of its restaurant, and the list of dishes.

For the restaurant, its location is known. The restaurant location is characterized by the street, the street number, the zip code, the city, the country and the geographical coordinates (latitude and longitude).

For each dish, its price must be recorded (i.e., the currency and the value).

Given a menu, the database must be designed to efficiently provide all the dish names and their chef names without additional lookups.

Write a sample document for each collection of the database.

Explicitly indicate the design patterns used.

Dish { _id: ObjectId(), name: <string>, description: <string>, complexity_level: <int or float>, chef: { name: <string>, country: <string>, website: <url>, email: <string> }, ingredients: { // otherwise a list with sub-dictionaries, i.e., attribute pattern "milk": {value: 500, uom: "ml"}, "eggs": {value: 2}, }, }

Menu

```
{
_id: ObjectId(),
name: <string>,
restaurant: {
     "street": <string>,
     "street_number": <string>,
     "zip_code": <number>,
     "city": <string>,
     "country": <string>,
     // coordinates only, without Point, are acceptable [x,y] as legacy format
     "loc": {"type":"Point",
          "coordinates":[: <number>, <number>]
}
dishes: [
{
 _id: Objectid(),
 dish_name: <string>,
 chef_name: <string>,
 price: { val: <number>, currency: <string> } // otherwise price_val and price_currency
 }, ...
]
}
```

Pattern used:

- Extended reference pattern in menu collection for dishes name and chef
- Polymorphic pattern for chef information
- Polymorphic or Attribute pattern for ingredients