

Practical Exam: Grocery Store Sales

FoodYum is a grocery store chain that is based in the United States.

Food Yum sells items such as produce, meat, dairy, baked goods, snacks, and other household food staples.

As food costs rise, FoodYum wants to make sure it keeps stocking products in all categories that cover a range of prices to ensure they have stock for a broad range of customers.

Data

The data is available in the table products.

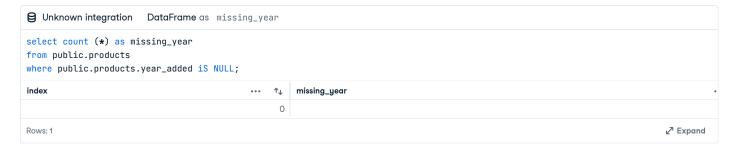
The dataset contains records of customers for their last full year of the loyalty program.

Column Name	Criteria
product_id	Nominal. The unique identifier of the product. Missing values are not possible due to the database structure.
product_type	Nominal. The product category type of the product, one of 5 values (Produce, Meat, Dairy, Bakery, Snacks). Missing values should be replaced with "Unknown".
brand	Nominal. The brand of the product. One of 7 possible values. Missing values should be replaced with "Unknown".
weight	Continuous. The weight of the product in grams. This can be any positive value, rounded to 2 decimal places. Missing values should be replaced with the overall median weight.
price	Continuous. The price the product is sold at, in US dollars. This can be any positive value, rounded to 2 decimal places. Missing values should be replaced with the overall median price.
average_units_sold	Discrete. The average number of units sold each month. This can be any positive integer value. Missing values should be replaced with 0.
year_added	Nominal. The year the product was first added to FoodYum stock. Missing values should be replaced with 2022.
stock_location	Nominal. The location that stock originates. This can be one of four warehouse locations, A, B, C or D Missing values should be replaced with "Unknown".

Task 1

Last year (2022) there was a bug in the product system. For some products that were added in that year, the year_added value was not set in the data. As the year the product was added may have an impact on the price of the product, this is important information to have.

Write a query to determine how many products have the year_added value missing. Your output should be a single column, missing_year, with a single row giving the number of missing values.



Task 2

Given what you know about the year added data, you need to make sure all of the data is clean before you start your analysis. The table below shows what the data should look like.

Write a query to ensure the product data matches the description provided. Do not update the original table.

Column Name	Criteria
product_id	Nominal. The unique identifier of the product. Missing values are not possible due to the database structure.
product_type	Nominal. The product category type of the product, one of 5 values (Produce, Meat, Dairy, Bakery, Snacks). Missing values should be replaced with "Unknown".
brand	Nominal. The brand of the product. One of 7 possible values. Missing values should be replaced with "Unknown".
weight	Continuous. The weight of the product in grams. This can be any positive value, rounded to 2 decimal places. Missing values should be replaced with the overall median weight.
price	Continuous. The price the product is sold at, in US dollars. This can be any positive value, rounded to 2 decimal places. Missing values should be replaced with the overall median price.
average_units_sold	Discrete. The average number of units sold each month. This can be any positive integer value. Missing values should be replaced with 0.
year_added	Nominal. The year the product was first added to FoodYum stock. Missing values should be replaced with last year (2022).
stock_location	Nominal. The location that stock originates. This can be one of four warehouse locations, A, B, C or D Missing values should be replaced with "Unknown".

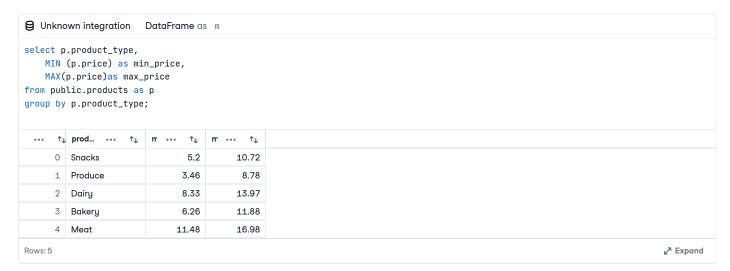
```
Unknown integration DataFrame as clean_data
-- Write your query for task 2 in this cell
WITH ProductStats AS (
    SELECT
        PERCENTILE_CONT(0.5) WITHIN GROUP (ORDER BY CAST(REGEXP_REPLACE(CAST(public.products.weight AS TEXT), '[^0-9.]', '', 'g')
AS NUMERIC)) AS median_weight,
        PERCENTILE_CONT(0.5) WITHIN GROUP (ORDER BY CAST(REGEXP_REPLACE(CAST(public.products.price AS TEXT), '[^0-9.]', '', 'g') AS
NUMERIC)) AS median_price
    FROM
        public.products
)
SELECT
    p.product_id,
    COALESCE(p.product_type, 'Unknown') AS product_type,
    COALESCE(p.brand, 'Unknown') AS brand,
    COALESCE(CAST(REGEXP_REPLACE(CAST(p.weight AS TEXT), '[^0-9.]', '', 'g') AS NUMERIC), (SELECT median_weight FROM ProductStats))
    COALESCE(CAST(REGEXP_REPLACE(CAST(p.price AS TEXT), '[^0-9.]', '', 'g') AS NUMERIC), (SELECT median_price FROM ProductStats))
AS price,
    COALESCE(p.average_units_sold, 0) AS average_units_sold,
    COALESCE(p.year_added, 2022) AS year_added,
    UPPER(TRIM(COALESCE(p.stock_location, 'Unknown'))) AS stock_location
FROM
    public.products AS p
CROSS JOIN
    ProductStats;
      ↑↓ p... •••
                                       brand ...
                                                         ... ↑↓
                                                                       ↑ average_units_...
                                                                                                                stock_lo... ···
                  \uparrow_{\bot}
                       \uparrow_{\perp}
                                                                                                      ••• ↑↓
      0
                   1
                      Bakeru
                                       TopBrand
                                                         602.61
                                                                      11
                                                                                              15
                                                                                                          2022
                                                                                                                С
                   2
                                                         478.26
                                                                    8.08
                                                                                                                C
      1
                      Produce
                                       SilverLake
                                                                                              22
                                                                                                          2022
                                                                    6.16
      2
                   3
                      Produce
                                       TastyTreat
                                                         532.38
                                                                                              21
                                                                                                          2018
                                                                                                                B
      3
                   4
                      Bakery
                                       StandardYums
                                                         453.43
                                                                    7.26
                                                                                              21
                                                                                                          2021
                                                                                                                D
      4
                   5
                      Produce
                                       GoldTree
                                                         588.63
                                                                    7.88
                                                                                              21
                                                                                                          2020 A
      5
                   6
                      Meat
                                       TopBrand
                                                         612.06
                                                                    16.2
                                                                                              24
                                                                                                          2017
                   7
                                       GoldTree
                                                         320.49
                                                                    8.01
      6
                      Produce
                                                                                              21
                                                                                                          2019
      7
                   8
                      Meat
                                       SilverLake
                                                         535.19
                                                                   15.77
                                                                                              28
                                                                                                          2021
                                                                                                                Α
      8
                   9
                                       StandardYums
                                                         375.07
                                                                   11.57
                                                                                              30
                      Meat
                                                                                                          2020 A
      9
                                                         506.34
                                                                                                                С
                  10
                      Meat
                                       TastyTreat
                                                                   13.94
                                                                                              27
                                                                                                          2018
     10
                      Dairy
                                       StandardYums
                                                         345.07
                                                                    9.26
                                                                                              26
                                                                                                          2020
                                                                                                                В
                  11
     11
                  12
                      Bakeru
                                       StandardYums
                                                         345.58
                                                                    6.87
                                                                                              21
                                                                                                          2022
                                                                                                                D
     12
                  13
                      Snacks
                                       SmoothTaste
                                                         512.54
                                                                    8.65
                                                                                              19
                                                                                                          2016
                                                                                                                Α
     13
                  14
                      Meat
                                       StandardYums
                                                         395.76
                                                                   11.92
                                                                                              30
                                                                                                          2019 A
     14
                  15
                       Produce
                                       SilverLake
                                                         324.92
                                                                    7.94
                                                                                              23
                                                                                                          2021
                                                                                                                D
     15
                  16
                      Dairy
                                       SmoothTaste
                                                         446.76
                                                                   10.79
                                                                                              23
                                                                                                          2017
Rows: 1,700

∠ Expand
```

Task 3

To find out how the range varies for each product type, your manager has asked you to determine the minimum and maximum values for each product tupe.

Write a query to return the product_type, min_price and max_price columns.



Task 4

The team want to look in more detail at meat and dairy products where the average units sold was greater than ten.

Write a query to return the product_id, price and average_units_sold of the rows of interest to the team.

