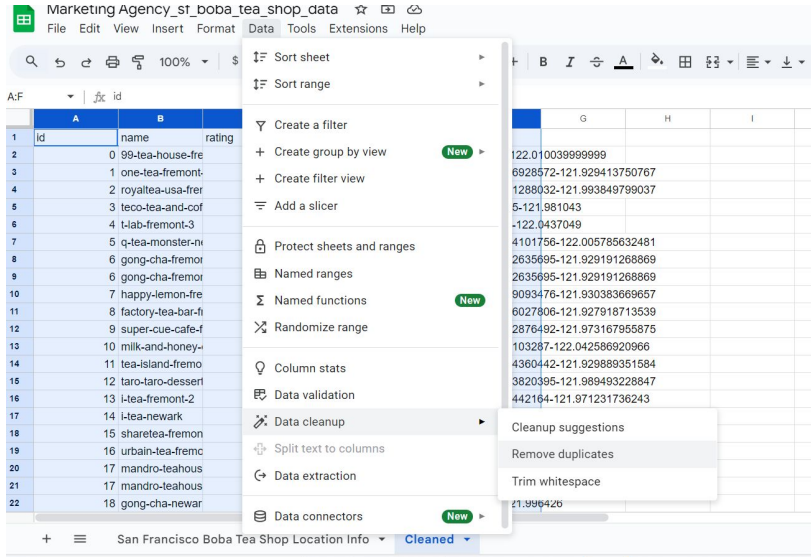# Analytics Portfolio

Jessica Gutstein

# Dirty to Clean Data - Google Sheets Project
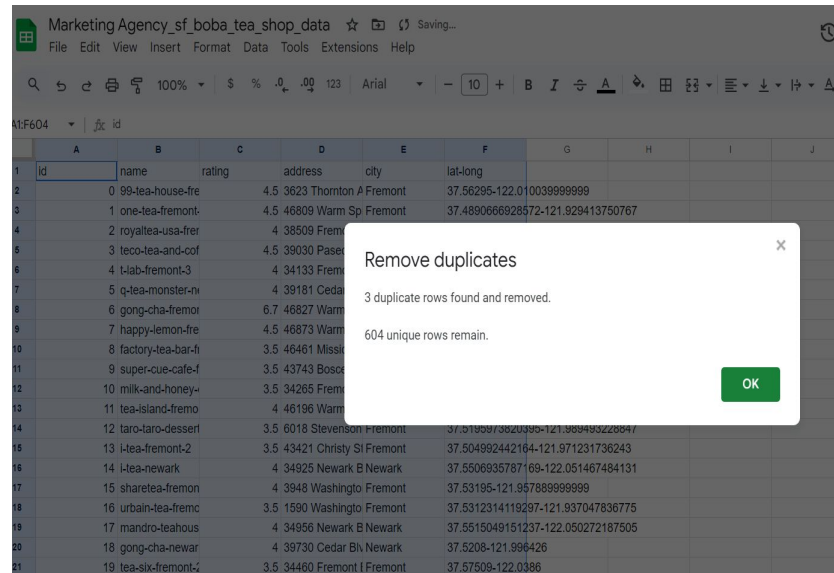
Skills: identify dirty elements using **Google Sheets** in a dataset, remove duplicate data, and use the `COUNTIF` and `SPLIT` functions to help clean data.

**Dataset: Marketing Agency**

**https://docs.google.com/spreadsheets/d/1p5TiRjV-criLJTI7A0tBdMF0tEUiBIdjY3xL0KOGGHg/edit?usp=sharing**

# Functions - COUNTIF



- Yelp Rates in Column C can only be between (0-5).
- Using COUNTIF to determine if any values are over 5.
- 9 entries that have a rating greater than 5.
- Solution:
1. Research Boba locations on yelp to find the accurate rankings, in this dataset example, the incorrect entries will be replaced with 5.

# Sorting/Filtering



Solution: Replaced all 9 incorrect data entries with 5. To validate correction, the COUNTIF function will show 0 entries >5.

# Functions - SPLIT

# Formulas



- Longitude values should be negative so that they are accurate coordinates for mapping.
- `=G2*-1`

# SQL - Common Queries

# Identifies variable within a specific row

# Updating a variable within a row

# SQL-Filtering



```
1  SELECT
2  *
3  -- name
4  FROM sql-portfolio-444521.customer_data.customer_address
5  --WHERE customer_id = 2645
6  --WHERE state = 'FL'
7  WHERE state IN ('FL', 'CA')
```

# Cleaning String Variables - removing duplicates

# Functions - LENGTH / DISTINCT / SUBSTR

```
1  SELECT
2  LENGTH(country) AS letters_in_country
3  FROM sql-portfolio-444521.customer_data.customer_address
```

Press Alt+F1 f

Query results

SAVE RESULTS ▾

JOB INFORMATION    RESULTS    CHART    JSON

| w | letters_in_country |
|---|---|
| 10 | 2 |
| 11 | 3 |
| 12 | 2 |
| 13 | 2 |

Notes

```
SELECT
  country
FROM
  sql-portfolio-444521.customer_data.customer_address
WHERE
  LENGTH(country) > 2
```

Press Alt+F1 f

Query results

SAVE RESULTS ▾

JOB INFORMATION    RESULTS    CHART    JSON

| | country |
|---|---|
| 1 | USA |
| 2 | USA |

```
SELECT
  DISTINCT customer_id
FROM
  sql-portfolio-444521.customer_data.customer_address
WHERE
  SUBSTR(country,1,2) = 'US'
```

Press

Query results

SAVE RESULTS ▾

JOB INFORMATION    RESULTS    CHART    JSON

| | customer_id |
|---|---|
| 1 | 2463 |
| 2 | 5306 |
| 3 | 9886 |
| 4 | 3821 |

Results per page:    50 ▾    1 – 25 of 25

# TRIM()

OH is not greater than 2 characters, so using the TRIM function will remove any spaces

```sql
SELECT
  state
FROM
  sql-portfolio-444521.customer_data.customer_address
WHERE
  LENGTH(state) > 2
```

**Query results**

| | JOB INFORMATION | RESULTS | CHART | JSON |
|---|---|---|---|---|
| ow | state ▾ | | | |
| 1 | OH | | | |

```sql
SELECT
  DISTINCT customer_id
FROM
  sql-portfolio-444521.customer_data.customer_address
WHERE
  TRIM(state) = 'OH'
```

**Query results**

| | JOB INFORMATION | RESULTS | CHART | JSO |
|---|---|---|---|---|
| w | customer_id ▾ | | | |
| 1 | 1928 | | | |
| 2 | 7571 | | | |
| 3 | 1980 | | | |
| 4 | 4687 | | | |

# Dirty to Clean Data - SQL Project

## Dataset: Automobiles

https://docs.google.com/spreadsheets/d/1uc563psJrCU6Gf0ftaWLXKXAF7czdONsKnaCgC5bz9c/edit?usp=sharing

Contains historical sales data, including details such as car features and prices. You can use this data to find the top 10 most popular cars and trims

Goal of project: clean data to avoid presenting inaccurate inventory, which could result in the company to lose money on their car investment

Automobile - UCI Machine Learning Repository - verifies that the fuel_type column should contain only two unique string values 'diesel' and 'gas'.

See screenshot to verify fuel_type has the two unique strings.

# MAX/MIN Functions

```
1  SELECT
2    MIN(length) AS min_length,
3    MAX(length) AS max_length
4  FROM sql-portfolio-444521.cars.car_info
```

Within the automobiles dataset, I wanted to know the (max/min) within the length column. When the query returned the values, I asked to make two new columns for the entries to fall into.

## Query results

| | JOB INFORMATION | RESULTS | CH |
|---|---|---|---|

| w | min_length ▼ | max_length ▼ | |
|---|---|---|---|
| 1 | 141.1 | 208.1 | |

# SQL - Missing Data

ALWAYS check data for null or missing values. Confirm replace with team or data owner and make correction.



```
   Missing_data        ▶ RUN    📩 SAVE QUERY  ▾    ⋮    ✅ This query will process 28.71
1  SELECT
2   *
3  FROM
4    sql-portfolio-444521.cars.car_info
5
6  WHERE
7    num_of_doors IS NULL;
```

Press Alt+F1 for Accessib

**Query results**  📩 SAVE RESULTS ▾   📊 OPEN IN ▾

| | JOB INFORMATION | RESULTS | CHART | JSON | EXECUTION DETAILS |

| Row | make ▾ | fuel_type ▾ | num_of_doors ▾ |
| --- | --- | --- | --- |
| 1 | dodge | gas | null |
| 2 | mazda | diesel | null |

```
1  UPDATE
2    sql-portfolio-444521.cars.car_info
3  SET
4    num_of_doors = "four"
5  WHERE
6    make = "dodge"
7    AND fuel_type = "gas"
8    AND body_style = "sedan";
```

**Query results**  📩 SAVE

| JOB INFORMATION | RESULTS | EXECUTION DETAILS |

ℹ  This statement modified 3 rows in car_info.

# SQL - Correcting Misspelling

'Tow' is now corrected to 'two'. The DISTINCT function only picks up 'two'.



```
1  SELECT
2    DISTINCT num_of_cylinders
3    FROM sql-portfolio-444521.cars.car_info;
```

Query results

| Row | num_of_cylinders ▾ |
|---|---|
| 3 | three |
| 4 | two |
| 5 | tow |
| 6 | twelve |

Results per page:



```
1  UPDATE
2    sql-portfolio-444521.cars.car_info
3  SET
4    num_of_cylinders = "two"
5  WHERE
6    num_of_cylinders = "tow";
```

Query results                    SAVE RESULTS

JOB INFORMATION    RESULTS    EXECUTION DETAILS    EXEC

ⓘ    This statement modified 1 row in car_info.



```
1  SELECT
2    DISTINCT num_of_cylinders
3    FROM
4      sql-portfolio-444521.cars.car_info
5
```

Query results

| JOB INFORMATION | RESULTS | C |
|---|---|---|

| w | num_of_cylinders ▾ |
|---|---|
| 2 | eight |
| 3 | three |
| 4 | two |
| 5 | twelve |

Results per

# MAX/MIN ranges (removing/deleting <> outliers)

```
1  SELECT
2    MAX (compression_ratio) AS max_compression_ratio,
3    MIN (compression_ratio) AS min_compression_ratio
4  FROM
5    sql-portfolio-444521.cars.car_info
6  --WHERE
7    --compression_ratio <> 70;
```

## Query results

| | JOB INFORMATION | RESULTS | CHART | JS( |
|---|---|---|---|---|
| | max_compression_ra | min_compression_ra | | |
| 1 | 70.0 | 7.0 | | |

SAVE RESULTS ▾

```
1  SELECT
2    MAX (compression_ratio) AS max_compression_ratio,
3    MIN (compression_ratio) AS min_compression_ratio
4  FROM
5    sql-portfolio-444521.cars.car_info
6  WHERE
7    compression_ratio <> 70;
```

## Query results

| | JOB INFORMATION | RESULTS | CHART | JS( |
|---|---|---|---|---|
| | max_compression_ra | min_compression_ra | | |
| 1 | 23.0 | 7.0 | | |

SAVE RESULTS ▾

```
1  SELECT
2    COUNT(*) AS num_rows_to_delete
3  FROM
4    sql-portfolio-444521.cars.car_info
5  WHERE
6    compression_ratio = 70;
7
```

## Query results

| | JOB INFORMATION | RESULTS |
|---|---|---|
| ow | num_rows_to_delete | |
| 1 | 1 | |

Untitled query   ▶ RUN   SAVE ▾   DOWNLOAD   ⁝ ✓

```
1   SELECT
2     COUNT(*) AS num_rows_to_delete
3   FROM
4     sql-portfolio-444521.cars.car_info
5   WHERE
6     compression_ratio = 70;
7   DELETE
8     sql-portfolio-444521.cars.car_info
9   WHERE
10    compression_ratio = 70;
11
```
Press Alt+F1 for Accessibility Optio

Query results        SAVE RESULTS ▾   OPEN IN ▾   ↕

JOB INFORMATION   **RESULTS**   EXECUTION DETAILS   EXECUTION GRAPH

ⓘ   This statement removed 1 row from car_info.   **GO TO TABLE**

# SQL - Multiple Functions

# GROUP BY



- Applying the GROUPBY to the non-aggregated function for MAX price to align the make with MAX/AVG prices.

# ORDER BY

```sql
1  SELECT
2    make,
3    num_of_doors,
4    avg(price) AS avg_price,
5  FROM
6    sql-portfolio-444521.cars.car_info
7  GROUP BY make, num_of_doors
8  ORDER BY avg_price DESC;
```

Press Alt+F1 for Acc

## Query results

SAVE RESULTS ▾      OPEN I

| | JOB INFORMATION | RESULTS | CHART | JSON | EXECUTION DETAIL |
|---|---|---|---|---|---|

| Row | make ▾ | num_of_doors ▾ | avg_price ▾ | |
|---|---|---|---|---|
| 1 | mercedes-benz | two | 36210.66666666... | |
| 2 | jaguar | two | 36000.0 | |
| 3 | jaguar | four | 33900.0 | |
| 4 | mercedes-benz | four | 32108.8 | |

Results per page:    50 ▾    1 – 37 of 37    |<    <

Job history

# Advanced Function - CAST()

- String -> Float conversion.
- Verify change through SCHEMA.
- We needed to determine the customer purchase price from greatest to least, but the numbers were being ordered by strings not numerical values. By changing to FLOAT, we can accurately order the values.

| | Field name | Type |
|---|---|---|
| ☐ | date | DATETIME |
| ☐ | transaction_id | INTEGER |
| ☐ | customer_id | INTEGER |
| ☐ | product | STRING |
| ☐ | product_code | STRING |
| ☐ | product_color | STRING |
| ☐ | product_price | FLOAT |
| ☐ | purchase_size | INTEGER |
| ☐ | purchase_price | STRING |
| ☐ | revenue | FLOAT |

SCHEMA | DETAILS | PREVI

```
SELECT
  CAST(purchase_price AS FLOAT64)
FROM
  sql-portfolio-444521.customer_data2.customer_purchase
ORDER BY
  CAST(purchase_price AS FLOAT64) DESC
```

# Filtering Date Periods (Correcting String -> Date)



```
SELECT
  date,
  purchase_price
FROM
  sql-portfolio-444521.customer_data2.customer_purchase
WHERE
  date BETWEEN '2020-12-01' AND '2020-12-31'
```

Press Alt+F1 f

uery results

SAVE RESULTS ▾

| | JOB INFORMATION | RESULTS | CHART | JSON |
|---|---|---|---|---|

| | date ▾ | purchase_price ▾ |
|---|---|---|
| 1 | 2020-12-12T00:00:00 | 13.99 |
| 2 | 2020-12-28T00:00:00 | 27.98 |
| 3 | 2020-12-28T00:00:00 | 160.965 |
| 4 | 2020-12-30T00:00:00 | 269.55 |

---

date_filtering_periods   ▶ RUN

```
1  SELECT
2    CAST(date AS date) AS date_only,
3    purchase_price
4  FROM
5    sql-portfolio-444521.customer_data2.custome
6  WHERE
7    date BETWEEN '2020-12-01' AND '2020-12-31'
```

Query results

SAVE RES

| | JOB INFORMATION | RESULTS | CHART |
|---|---|---|---|

| | date_only ▾ | purchase_price ▾ |
|---|---|---|
| 1 | 2020-12-12 | 13.99 |
| 2 | 2020-12-28 | 27.98 |
| 3 | 2020-12-28 | 160.965 |
| 4 | 2020-12-30 | 269.55 |

# CONCAT()

| product_code | product_color |
|---|---|
| SKU83503 | brass |
| SKU83503 | brass |
| SKU83503 | white |

```sql
1  SELECT
2    CONCAT(product_code, product_color) AS new_product_code
3  FROM
4    sql-portfolio-444521.customer_data2.customer_purchase
5  WHERE
6    product = "couch"
```

## Query results                                    ☐ SAVE RESULT

| JOB INFORMATION | RESULTS | CHART | JSON |
|---|---|---|---|

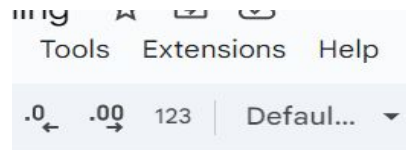| w | new_product_code ▾ |
|---|---|
| 1 | SKU31871grey |
| 2 | SKU31871grey |
| 3 | SKU31871grey |
| 4 | SKU31871blue |

Results per page:  50 ▾    1 –

# Advanced Function - COALESCE ()

```
SELECT
  COALESCE (product, product_code) AS product_info
FROM
  sql-portfolio-444521.customer_data2.customer_purchase
```

- This function replaces NULL's with another identifier in the column you are trying to calculate total. This would be a placeholder for another field that shows as NULL.

# Verifying Data is Cleaned - Excel/Google Sheets Using Pivot Tables/ Find & Replace

# Verifying Data is Cleaning Using SQL CASE()

```sql
1    SELECT
2      CASE
3        WHEN Attrition = false THEN 'Active'
4        WHEN Attrition = true THEN 'Non-Active'
5      END AS Attrition,
6      CASE
7        WHEN Age BETWEEN 18 AND 19 THEN '10s'
8        WHEN Age BETWEEN 20 AND 29 THEN '20s'
9        WHEN Age BETWEEN 30 AND 39 THEN '30s'
10       WHEN Age BETWEEN 40 AND 49 THEN '40s'
11       WHEN Age BETWEEN 50 AND 59 THEN '50s'
12       WHEN Age BETWEEN 60 AND 69 THEN '60s'
13     END AS Age,
14     MAX(MonthlyIncome) AS max_monthlyincome,
15     MIN(MonthlyIncome) AS min_monthlyincome,
16     ROUND(AVG(MonthlyIncome),2) AS average_monthlyincome,
17     APPROX_QUANTILES(MonthlyIncome,100)[OFFSET(25)] AS `25th_percentile`,
18     APPROX_QUANTILES(MonthlyIncome,100)[OFFSET(75)] AS `75th_percentile`
19   FROM
20     sql-portfolio-444521.HR_People_Attrition_data.attrition_data
21   GROUP BY Attrition, Age
22   ORDER BY Age DESC
```