# Project : Data cleaning

**Name :** AMARJEET KUMAR

**Email:** amarjeetbilla63@gmail.com

**Batch:** DA - BATCH ^6

**Phone:** 9871434547

---

## ◆ STEP 0: Inspect Raw Data

```sql
SELECT *
FROM customer_orders
LIMIT 10;
```



---

## ◆ STEP 1: Clean `first_name` (Spaces + Case)

```sql
SELECT
  first_name,
  TRIM(first_name) AS step1_trimmed,
  UPPER(TRIM(first_name)) AS cleaned_first_name
```

```
FROM customer_orders;
```



## ◆ STEP 2: Clean `last_name`

```
SELECT
  last_name,
  UPPER(TRIM(last_name)) AS cleaned_last_name
FROM customer_orders;
```

## ◆ STEP 3: Create `full_name` (CONCAT)

```sql
SELECT
  CONCAT(
    UPPER(TRIM(first_name)),
    ' ',
    UPPER(TRIM(last_name))
  ) AS full_name
```

```
FROM customer_orders;
```



---

## ◆ STEP 4: Clean `email` (Standardization)

```sql
SELECT
  email,
  LOWER(email) AS cleaned_email
FROM customer_orders;
```

## ◆ STEP 5: Clean `mobile_number` (Extract last 10 digits)

```sql
SELECT
  mobile_number,
  SUBSTR(mobile_number, LENGTH(mobile_number) - 9, 10) AS cleaned_mobile
FROM customer_orders;
```

```
35
36    -- ○ STEP 5: Clean mobile_number (Extract last 10 digits)
37 •  select
38        mobile_number,
39        substr(mobile_number,length(mobile_number) -9,10) as cleaned_mobile_number
40    from customer_orders;
41
42
```

Result Grid | Filter Rows: | Export: | Wrap Cell Content: A

| mobile_number | cleaned_mobile_number |
|---|---|
| 919110053353 | 9110053353 |
| 0091-9749621470 | 9749621470 |
| 0091-9664130526 | 9664130526 |
| 919654049436 | 9654049436 |
| 919940992571 | 9940992571 |
| +91-9811514914 | 9811514914 |
| +91-9466825638 | 9466825638 |
| 0091-9997612044 | 9997612044 |

Result 10 ×

---

## ◆ STEP 6: Extract Year from `order_id`

```sql
SELECT
  order_id,
  SUBSTR(order_id, 5, 4) AS order_year
FROM customer_orders;
```

```
41
42    -- o STEP 6: Extract Year from order_id
43 •  SELECT order_id,
44       SUBSTR(order_id, 5, 4) AS ordered_year
45    FROM customer_orders;
46
47
48
```

| order_id | ordered_year |
| --- | --- |
| ORD-2022-0001 | 2022 |
| ORD-2021-0002 | 2021 |
| ORD-2024-0003 | 2024 |
| ORD-2022-0004 | 2022 |
| ORD-2022-0005 | 2022 |
| ORD-2023-0006 | 2023 |
| ORD-2022-0007 | 2022 |
| ORD-2022-0008 | 2022 |

Result 11 ×

---

### ◆ STEP 7: Round order_amount

```sql
SELECT
  order_amount,
  ROUND(order_amount, 2) AS cleaned_order_amount
FROM customer_orders;
```

```
46
47        -- ◦ STEP 7: Round order_amount
48 •    SELECT
49          order_amount,
50          ROUND(order_amount, 2) AS cleaned_order_amount
51       FROM customer_orders;
52
53
```

| Result Grid | Filter Rows: | Export: | Wrap Cell Content: |

| order_amount | cleaned_order_amount |
|---|---|
| ▶ 3662.377 | 3662.38 |
| 10669.923 | 10669.92 |
| 23175.878 | 23175.88 |
| 38255.816 | 38255.82 |
| 91497.485 | 91497.48 |
| 59842.699 | 59842.7 |
| 50308.318 | 50308.32 |
| 8336.326 | 8336.33 |

## ◆ STEP 8: Round rating

```sql
SELECT
  rating,
  ROUND(rating, 1) AS cleaned_rating
FROM customer_orders;
```

```
52
53        -- ○ STEP 8: Round rating
54  •   SELECT
55        rating,
56        ROUND(rating, 1) AS cleaned_rating
57      FROM customer_orders;
58
59
```

| rating | cleaned_rating |
|--------|----------------|
| 1.882  | 1.9            |
| 4.892  | 4.9            |
| 1.651  | 1.7            |
| 4.287  | 4.3            |
| 2.598  | 2.6            |
| 3.883  | 3.9            |
| 2.496  | 2.5            |
| 1.661  | 1.7            |

Result 13 ×

---

### ◆ STEP 9: Standardize city

```
SELECT
  city,
  UPPER(city) AS cleaned_city
FROM customer_orders;
```

```
 56        ROUND(rating, 1) AS cleaned_rating
 57    FROM customer_orders;
 58
 59    --  ◊ STEP 9: Standardize city
 60 •  select
 61        city , upper(city) as cleaned_city
 62    from customer_orders;
 63
```

| city | cleaned_city |
|---|---|
| PUNE | PUNE |
| delhi | DELHI |
| bangalore | BANGALORE |
| MUMBAI | MUMBAI |
| hyderabad | HYDERABAD |
| CHENNAI | CHENNAI |
| MUMBAI | MUMBAI |
| hyderabad | HYDERABAD |

### ◆ STEP 10: Delivery Time Calculation (DATEDIFF)

```sql
SELECT
  order_date,
  delivery_date,
  DATEDIFF(delivery_date, order_date) AS delivery_days
FROM customer_orders;
```

```
63
64     -- ○ STEP 10: Delivery Time Calculation (DATEDIFF)
65 •   select
66         delivery_date,
67         order_date,
68         datediff(delivery_date,order_date) as delivery_days_remains
69     from customer_orders;
70
```

| delivery_date | order_date | delivery_days_remains |
|---|---|---|
| 2022-07-20 | 2022-07-18 | 2 |
| 2021-01-21 | 2021-01-14 | 7 |
| 2024-02-08 | 2024-02-05 | 3 |
| 2022-07-07 | 2022-07-01 | 6 |
| 2022-07-04 | 2022-07-03 | 1 |
| 2023-11-14 | 2023-11-10 | 4 |
| 2022-11-29 | 2022-11-28 | 1 |
| 2022-10-30 | 2022-10-26 | 4 |

Result 17 ✕

### ◆ STEP 11: Customer Tenure Calculation

```sql
SELECT
  signup_date,
  DATEDIFF(NOW(), signup_date) AS days_with_company
FROM customer_orders;
```

```
71    -- ○ STEP 11: Customer Tenure Calculation
72 •  select
73        signup_date ,
74        datediff(now(),signup_date) as days_with_componey
75    from customer_orders;
76
77
78
```

| signup_date | days_with_componey |
|---|---|
| 2021-11-22 | 1496 |
| 2018-10-02 | 2643 |
| 2021-09-29 | 1550 |
| 2020-07-22 | 1984 |
| 2021-11-19 | 1499 |
| 2023-07-30 | 881 |
| 2022-01-22 | 1435 |
| 2022-09-27 | 1187 |

◆ **STEP 12: CASE WHEN – Order Value Category**

```sql
SELECT
  order_amount,
  CASE
    WHEN order_amount >= 50000 THEN 'High Value'
    WHEN order_amount >= 20000 THEN 'Medium Value'
    ELSE 'Low Value'
  END AS order_category
FROM customer_orders;
```

```
77      -- STEP 12: CASE WHEN - Order Value Category
78 •    select
79        CONCAT( UPPER(TRIM(first_name)),' ',UPPER(TRIM(last_name))) AS full_name , order_amount,
80 ⊖      case
81          when order_amount >=50000 then 'ambani mere ____ pe '
82          when order_amount >=30000 then 'bss itni aukat hai teri ,so you  arre ambani ,hehehe'
83          when order_amount >=20000 then 'teri aukat nhi hai order_dene ki toh apni ____ mra'
84          else 'bhai tu ekk km kr , tel__chatai ka dhndha start kr de '
85        end as aukat_ka_mirror
86      from customer_orders;
o7
```

| Result Grid | 🔢 | 🔁 Filter Rows: | | Export: 🖫 | Wrap Cell Content: 🔳 |

| | full_name | order_amount | aukat_ka_mirror |
|---|---|---|---|
| ▶ | ANITA SHARMA | 3662.377 | bhai tu ekk km kr , tel__chatai ka dhndha start k... |
| | KIRAN GUPTA | 10669.923 | bhai tu ekk km kr , tel__chatai ka dhndha start k... |
| | VIKAS PATEL | 23175.878 | teri aukat nhi hai order_dene ki toh apni ____ mra |
| | POOJA SINGH | 38255.816 | bss itni aukat hai teri ,so you  arre ambani ,heh... |
| | POOJA VERMA | 91497.485 | ambani mere ____ pe |

◆ **STEP 13: CASE  WHEN – Customer Type**

```
SELECT
  signup_date,
  CASE
    WHEN DATEDIFF(NOW(), signup_date) <= 30 THEN 'New'
    WHEN DATEDIFF(NOW(), signup_date) <= 180 THEN 'Regular'
    ELSE 'Loyal'
  END AS customer_type
FROM customer_orders;
```

```
87
88      -- ∘ STEP 13: CASE WHEN – Customer Type
89 •   SELECT signup_date,
90  ⊖    CASE
91           WHEN DATEDIFF(NOW(), signup_date) <= 30 THEN 'New'
92           WHEN DATEDIFF(NOW(), signup_date) <= 180 THEN 'Regular'
93           ELSE 'Loyal'
94         END AS customer_type
95      FROM customer_orders;
```

| Result Grid | Filter Rows: | Export: | Wrap Cell Content: |

| signup_date | customer_type |
|---|---|
| 2021-11-22 | Loyal |
| 2018-10-02 | Loyal |
| 2021-09-29 | Loyal |
| 2020-07-22 | Loyal |
| 2021-11-19 | Loyal |
| 2023-07-30 | Loyal |
| 2022-01-22 | Loyal |

Result 20 ✕

## ◆ STEP 14: FINAL CLEANED VIEW (Industry Practice)

```sql
CREATE VIEW customer_orders_cleaned AS
SELECT
  customer_id,
  UPPER(TRIM(first_name)) AS first_name,
  UPPER(TRIM(last_name)) AS last_name,
  CONCAT(
    UPPER(TRIM(first_name)), ' ',
    UPPER(TRIM(last_name))
  ) AS full_name,
  LOWER(email) AS email,
  SUBSTR(mobile_number, LENGTH(mobile_number) - 9, 10) AS mobile_number,
```

```sql
  order_id,
  SUBSTR(order_id, 5, 4) AS order_year,
  order_date,
  delivery_date,
  DATEDIFF(delivery_date, order_date) AS delivery_days,
  ROUND(order_amount, 2) AS order_amount,
  UPPER(city) AS city,
  signup_date,
  DATEDIFF(NOW(), signup_date) AS customer_tenure_days,
  CASE
    WHEN order_amount >= 50000 THEN 'High Value'
    WHEN order_amount >= 20000 THEN 'Medium Value'
    ELSE 'Low Value'
  END AS order_category,
  ROUND(rating, 1) AS rating
FROM customer_orders;
```

```sql
103     --  STEP 14: FINAL CLEANED VIEW (Industry Practice)
104 •   create view  vw_for_cleaned_customer_orders as
105     select
106         customer_id,
107     upper(trim(first_name)) as clean_fname,
108         UPPER(TRIM(last_name)) AS cleaned_last_name,
109     CONCAT( UPPER(TRIM(first_name)),' ',UPPER(TRIM(last_name))) AS full_name ,
110         lower(email) as clear_email ,
111     substr(mobile_number,length(mobile_number) -9,10) as cleaned_mobile_number ,
112          order_id,
113     SUBSTR(order_id, 5, 4) AS ordered_year,
114         ROUND(order_amount, 2) AS cleaned_order_amount,
115     delivery_date,
116         upper(city) as cleaned_city ,
117     order_date,
118         datediff(delivery_date,order_date) as delivery_days_remains,
119     signup_date ,
120         datediff(now(),signup_date) as days_with_componey,
121     case
122         when order_amount >=50000 then 'ambani mere _____ pe '
123         when order_amount >=30000 then 'bss itni aukat hai teri ,so you  arre ambani ,hehehe'
124         when order_amount >=20000 then 'teri aukat nhi hai order_dene ki toh apni _____ mra'
125         else 'bhai tu ekk km kr , tel__chatai ka dhndha start kr de '
126     end as order_category
127     from customer_orders;
```

## ◆ STEP 15: Validate Cleaned Data

```sql
SELECT *
FROM customer_orders_cleaned
LIMIT 10;
```



```sql
--  STEP 15: Validate Cleaned Data

select *
from vw_for_cleaned_customer_orders
limit 20;
```

| cleaned_mobile_number | order_id | ordered_year | cleaned_order_amount | delivery_date | cleaned_city | order_date | delivery_days_remains | signup_date | days_with_componey | order_category |
|---|---|---|---|---|---|---|---|---|---|---|
| 9110053353 | ORD-2022-0001 | 2022 | 3662.38 | 2022-07-20 | PUNE | 2022-07-18 | 2 | 2021-11-22 | 1496 | bhai tu ekk km kr , tel_chatai ka dhndha sta |
| 9749621470 | ORD-2021-0002 | 2021 | 10669.92 | 2021-01-21 | DELHI | 2021-01-14 | 7 | 2018-10-02 | 2643 | bhai tu ekk km kr , tel_chatai ka dhndha sta |
| 9664130526 | ORD-2024-0003 | 2024 | 23175.88 | 2024-02-08 | BANGALORE | 2024-02-05 | 3 | 2021-09-29 | 1550 | teri aukat nhi hai order_dene ki toh apni __ |
| 9654049436 | ORD-2022-0004 | 2022 | 38255.82 | 2022-07-07 | MUMBAI | 2022-07-01 | 6 | 2020-07-22 | 1984 | bss itni aukat hai teri ,so you arre ambani ,k |
| 9940992571 | ORD-2022-0005 | 2022 | 91497.48 | 2022-07-04 | HYDERABAD | 2022-07-03 | 1 | 2021-11-19 | 1499 | ambani mere ____ pe |
| 9811514914 | ORD-2023-0006 | 2023 | 59842.7 | 2023-11-14 | CHENNAI | 2023-11-10 | 4 | 2023-07-30 | 881 | ambani mere ____ pe |
| 9466825638 | ORD-2022-0007 | 2022 | 50308.32 | 2022-11-29 | MUMBAI | 2022-11-28 | 1 | 2022-01-22 | 1435 | ambani mere ____ pe |
| 9997612044 | ORD-2022-0008 | 2022 | 8336.33 | 2022-10-30 | HYDERABAD | 2022-10-26 | 4 | 2022-09-27 | 1187 | bhai tu ekk km kr , tel_chatai ka dhndha sta |
| 9215984311 | ORD-2022-0009 | 2022 | 44090.96 | 2022-03-15 | MUMBAI | 2022-03-10 | 5 | 2020-01-17 | 2171 | bss itni aukat hai teri ,so you arre ambani ,k |
| 9354794895 | ORD-2021-0010 | 2021 | 72473.58 | 2021-05-20 | CHENNAI | 2021-05-19 | 1 | 2018-12-11 | 2573 | ambani mere ____ pe |
| 9701642483 | ORD-2023-0011 | 2023 | 86236.42 | 2023-08-22 | DELHI | 2023-08-15 | 7 | 2021-07-20 | 1621 | ambani mere ____ pe |
| 9897677788 | ORD-2021-0012 | 2021 | 40403.51 | 2021-04-20 | PUNE | 2021-04-14 | 6 | 2019-09-29 | 2281 | bss itni aukat hai teri ,so you arre ambani ,k |