

Top 10 Data Cleaning Steps



Remove Duplicates

Duplicate records can inflate metrics and distort analysis.

```
1 SELECT DISTINCT *  
2 FROM table_name;
```

```
1 WITH CTE AS (  
2   SELECT *, ROW_NUMBER() OVER (PARTITION BY id_column ORDER BY created_at DESC)  
   AS rn  
3   FROM table_name  
4 )  
5 DELETE FROM CTE WHERE rn > 1;
```

USE:

GROUP BY OR COUNT(*)



Handle NULL Values

NULLs can lead to incorrect aggregations or logic errors.

```
1 SELECT
2     COALESCE(column_name, 'Default Value') AS clean_column
3 FROM table_name;
```

```
1 UPDATE table_name
2 SET column_name = 'Default'
3 WHERE column_name IS NULL;
```



Trim Whitespace

Extra spaces cause mismatches in joins and comparisons.

```
1 UPDATE table_name  
2 SET column_name = TRIM(column_name);
```

USE

TRIM(), LTRIM() AND RTRIM() ..



Standardize Text Case

Inconsistent casing causes grouping and filtering issues.

```
1 UPDATE table_name  
2 SET column_name = LOWER(column_name);
```

```
1 UPDATE table_name  
2 SET column_name = UPPER(column_name);
```

**STANDARDIZE BEFORE
USING GROUP BY ON TEXT COLUMNS.**



Fix Invalid Date Formats

Invalid date formats break time-based analysis.

Convert string to proper date

```
1 SELECT CAST(column_name AS DATE)
2 FROM table_name;
```

**VALIDATE USING
ISDATE() IN SQL SERVER OR
USE REGEX IN ADVANCED DATABASES.**



Validate Data Types

Text stored in numeric/date fields can cause failures in queries.

Cast to correct Data type

```
1 SELECT CAST(column_name AS INT)
2 FROM table_name;
```



Remove Outliers

Outliers distort averages, trends, and modeling.

Remove top / bottom 1% salaries

```
1  SELECT *  
2  FROM table_name  
3  WHERE salary BETWEEN  
4     PERCENTILE_CONT(0.01) WITHIN GROUP (ORDER BY salary)  
5     AND  
6     PERCENTILE_CONT(0.99) WITHIN GROUP (ORDER BY salary);
```



Fix Inconsistent Categories

'NY', 'ny', and 'New York' should all mean the same thing.

Create a mapping table for large datasets.

```
1  UPDATE table_name
2  SET city = 'New York'
3  WHERE LOWER(city) IN ('ny', 'n.y.', 'new york');
```



Thanks for Reading!

Follow for more content on SQL

