# SQL Data Cleaning & Exploratory Data Analysis (EDA)

## 1. Introduction

The objective of this project was to clean and analyze a dataset related to company layoffs. The dataset contained information about companies, industries, countries, and the number of employees laid off.

## 2. Data Cleaning Steps

The following steps were performed during data cleaning:

- - Removed duplicate rows
- - Standardized company names, industries, and country names
- - Handled missing/null values appropriately
- - Corrected date formats for consistency

## 3. Exploratory Data Analysis (EDA)

The following analyses were conducted:

- - Maximum layoffs observed in dataset
- - Companies with 100% layoffs
- - Top companies with most layoffs
- - Yearly layoff trends
- - Monthly layoff trends
- - Rolling total layoffs
- - Top companies by year

## 4. Key Insights & Findings

- Certain industries experienced significantly higher layoffs compared to others.
- Some companies had 100% of their workforce laid off.
- Layoffs showed clear yearly and monthly patterns, with spikes in specific periods.
- Rolling totals highlighted long-term layoff trends.

## 5. Conclusion

This project highlighted the importance of proper data cleaning before conducting analysis. The exploratory data analysis provided insights into how layoffs have impacted companies over time, revealing key trends and industry-specific effects.

## 6. Appendix - SQL Queries

The project included SQL queries for data cleaning and EDA. Example queries:

*SELECT company, SUM(total_laid_off) FROM layoffs GROUP BY company;*

*SELECT YEAR(date), SUM(total_laid_off) FROM layoffs GROUP BY YEAR(date);*

# SQL Data Cleaning & Exploratory Data Analysis (EDA) – Layoffs Dataset:

1. **Creating a Backup Table :**

```
1              -- DATA CLEANING
2 •   SELECT * FROM layoffs;
3
4 •   CREATE TABLE layoffs1 LIKE layoffs;
5 •                                          SELECT * FROM layoffs1;
6
7 •   INSERT INTO layoffs1
8     SELECT *
9     FROM layoffs;
10
```

2. **Identifying and Removing Duplicates :**

```
11 •   SELECT *,
12     ROW_NUMBER() OVER(PARTITION BY company, location, industry, total_laid_off, `date`, stage, country, funds_raised_millions) AS row_numbering
13     FROM layoffs1;
14
15 •   WITH duplicates AS
16   ⊝ (
17     SELECT *,
18     ROW_NUMBER() OVER(PARTITION BY company, location, industry, total_laid_off, `date`, stage, country, funds_raised_millions) AS row_numbering
19     FROM layoffs1
20     )
21     SELECT *
22     FROM duplicates
23     WHERE row_numbering > 1;
24
```

```sql
25 •    SELECT *
26      FROM layoffs1
27      WHERE company LIKE 'casper';
28
29 •    WITH duplicates AS
30    ⊖ (
31      SELECT *,
32      ROW_NUMBER() OVER(PARTITION BY company, location, industry, total_laid_off, `date`, stage, country, funds_raised_millions) AS row_numbering
33      FROM layoffs1
34      )
35      DELETE
36      FROM duplicates
37      WHERE row_numbering > 1;
38

39 • ⊖  CREATE TABLE `layoffs2` (
40          `company` text,
41          `location` text,
42          `industry` text,
43          `total_laid_off` int DEFAULT NULL,
44          `percentage_laid_off` text,
45          `date` text,
46          `stage` text,
47          `country` text,
48          `funds_raised_millions` int DEFAULT NULL,
49          `row_numbering` INT
50        ) ENGINE=InnoDB DEFAULT CHARSET=utf8mb4 COLLATE=utf8mb4_0900_ai_ci;
51
52 •                              SELECT * FROM layoffs2;
53 •    INSERT INTO layoffs2
54      SELECT *,
55      ROW_NUMBER() OVER(PARTITION BY company, location, industry, total_laid_off, `date`, stage, country, funds_raised_millions) AS row_numbering
56      FROM layoffs1;
57
58 •    SELECT *
59      FROM layoffs2
60      WHERE row_numbering > 1;
61
62 •    DELETE
63      FROM layoffs2
64      WHERE row_numbering > 1;
65
```

3. **Standardizing Text Data :**

```sql
66      -- STRANDANDIZING DATA
67 •    SELECT company, TRIM(company)
68      From layoffs2;
69
70 •    UPDATE layoffs2
71      SET company = TRIM(company);
72
73 •    SELECT DISTINCT industry
74      From layoffs2;
75
76 •    SELECT *
77      FROM layoffs2
78      WHERE industry LIKE 'CRYPTO%'
79      ;
80
81 •    UPDATE layoffs2
82      SET industry = 'Crypto'
83      WHERE industry LIKE 'Crypto%';
84
85 •    SELECT DISTINCT industry
86      FROM layoffs2;
```

```
85 •    SELECT DISTINCT industry
86      FROM layoffs2;
87
88 •    SELECT DISTINCT country
89      FROM layoffs2;
90
91 •    UPDATE layoffs2
92      SET country = 'United States'
93      WHERE country LIKE 'United States%';
```

## 4. Fixing Dates :

```
95 •   SELECT `date`,
96     STR_TO_DATE(`date`, '%m/%d/%Y')
97     FROM layoffs2;
98
99 •   UPDATE layoffs2
100    SET `date` = STR_TO_DATE(`date`, '%m/%d/%Y') ;
101
102 •  ALTER TABLE layoffs2
103    MODIFY COLUMN `date` DATE;
```

## 5. Handling Null & Missing Values :

```
111 •    SELECT *
112      FROM layoffs2
113      WHERE industry IS NULL
114      OR industry = '';

116 •    SELECT t1.industry, t2.industry
117      FROM layoffs2 t1
118      JOIN layoffs2 t2
119          ON t1.company = t2.company
120          WHERE t1.industry IS NULL
121          OR t1.industry = ''
122          AND t2.industry IS NOT NULL;
123
124 •    UPDATE layoffs2
125      SET industry = NULL
126      WHERE industry = '';

128 •  UPDATE layoffs2 t1
129    JOIN layoffs2 t2
130        ON t1.company = t2.company
131    SET t1.industry = t2.industry
132        WHERE t1.industry IS NULL
133        AND t2.industry IS NOT NULL;
```

**Exploratory Data Analysis (EDA) :**

## 5   Maximum layoffs :

```
155     -- EXPLATORY DATA ANALYSIS
156 •   SELECT *
157     FROM layoffs2;
158
159 •   SELECT MAX(total_laid_off), MAX(percentage_laid_off)
160     FROM layoffs2;
```

## 6   Companies with 100% layoffs :

```
162 •   SELECT  *
163     FROM layoffs2
164     WHERE percentage_laid_off = 1
165     ORDER BY total_laid_off DESC;
166
```

## 7   Companies with most layoffs :

```
172 •   SELECT company, SUM(total_laid_off)
173     FROM layoffs2
174     GROUP BY company
175     ORDER BY 2 DESC;
```

## 8   Layoffs by Year :

```
177 •   SELECT YEAR(`date`), SUM(total_laid_off)
178     FROM layoffs2
179     GROUP BY YEAR(`date`)
180     ORDER BY YEAR(`date`) DESC
181     ;
```

## 9   Monthly Trends :

```
183 •   SELECT SUBSTRING(`date`, 1,7) AS `MONTH`, SUM(total_laid_off) AS total_off
184     FROM layoffs2
185     WHERE SUBSTRING(`date`, 1,7) IS NOT NULL
186     GROUP BY `MONTH`
187     ORDER BY 1 ASC
188     ;
189
```

## 10   Rolling Total of Layoffs :

```
190 •   WITH ROLLING_TOTAL AS
191   ⊖ (
192     SELECT SUBSTRING(`date`, 1,7) AS `MONTH`, SUM(total_laid_off) AS total_off
193     FROM layoffs2
194     WHERE SUBSTRING(`date`, 1,7) IS NOT NULL
195     GROUP BY `MONTH`
196     ORDER BY 1 ASC
197   )
198     SELECT `MONTH`, SUM(total_off) OVER(ORDER BY `MONTH`) AS rolling_total
199     FROM ROLLING_TOTAL;
200
201 •   SELECT company, YEAR(`date`), SUM(total_laid_off)
202     FROM layoffs2
203     GROUP BY company, YEAR(`date`)
204     ORDER BY company ASC
205     ;
206
```

## 11 Top Companies by Year :

```sql
207    WITH COMPANY_YEAR (company, years, total_laid_off)AS
208    (
209     SELECT company, YEAR(`date`), SUM(total_laid_off)
210     FROM layoffs2
211     GROUP BY company, YEAR(`date`)
212     ORDER BY company ASC
213    ),
214    COMPANY_YEAR_RANK AS
215    (SELECT *, DENSE_RANK() OVER(PARTITION BY years ORDER BY total_laid_off DESC) AS DENSE_RANKING
216     FROM COMPANY_YEAR
217     WHERE YEARS IS NOT NULL
218    )
219    SELECT *
220    FROM COMPANY_YEAR_RANK
221    WHERE DENSE_RANKING <= 5;
222
```

# <u>Conclusion</u>:

This project demonstrated the importance of structured data cleaning and exploratory analysis in deriving meaningful insights from raw datasets. By addressing duplicates, inconsistencies, and missing values, the data was transformed into a reliable foundation for analysis. The EDA highlighted key layoff patterns, industry impacts, and temporal trends, offering valuable perspectives for strategic decision-making. Overall, the study reinforces that clean, well-analyzed data is critical to driving informed business outcomes.

----------------------------------------------------------------------------------------------------