

## COMBINING TABLES

Combining tables

```
CREATE TABLE
`harveaspace.data.newhr_table`
AS
SELECT DISTINCT
    e.Employee_ID,
    e.First_Name,
    e.Last_Name,
    e.SSN,
    e.Birth_Date,
    e.Sex,
    e.Address,
    e.Job_ID,
    e.Salary,
    d.Department_ID,
    d.Department_Name,
    d.Location_ID,
    jh.Start_Date,
    jo.Job_Title,
    jo.Minimum_Salary,
    jo.Maximum_Salary
FROM
    `harveaspace.data.Employee` AS e
LEFT JOIN
    `harveaspace.data.departments` AS d
ON
    e.Department_ID = d.Department_ID
LEFT JOIN
    `harveaspace.data.job history` AS jh
ON
    e.Employee_ID = jh.Employee_ID
LEFT JOIN
    `harveaspace.data.jobs` AS jo
ON
    jh.Job_ID = jo.Job_ID;
```

READING DATA AND CLEANING	
Overview of table	<pre>SELECT * FROM `harveaspace.data.newhr_table`</pre>
Adding Rows	<pre>INSERT INTO `harveaspace.data.newhr_table` (Employee_ID, First_Name, Last_Name, SSN, Birth_Date, Sex, Address, Job_ID, Salary, Department_ID, Department_Name, Location_ID, Start_Date, Job_Title, Minimum_Salary, Maximum_Salary) VALUES (1, 'Aarav', 'Sharma', '234-56-7890', '1992-11-25', 'M', '321 Pine Ave', 104, 55000, 1, 'Marketing', 1, '2018-08-01', 'Marketing Coordinator', 50000, 70000), (2, 'Neha', 'Patel', '876-54-3210', '1987-07-10', 'F', '654 Maple Dr', 105, 60000, 1, 'Marketing', 1, '2017-12-20', 'Marketing Assistant', 55000, 75000);</pre>
Checking the Years of Experience by Employees	<pre>SELECT Employee_ID, First_Name, Last_Name, Start_Date, DATE_DIFF(CURRENT_DATE(), Start_Date, YEAR) AS Total_Years FROM `harveaspace.data.newhr_table`</pre>
Creating a new column named "Total year"	<pre>CREATE OR REPLACE TABLE harveaspace.data.newhr_table_new AS SELECT Employee_ID, First_Name, Last_Name, Start_Date,</pre>

	<pre> (CASE     WHEN Employee_ID IN (4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26) THEN Total_Years     ELSE DATE_DIFF(CURRENT_DATE(), Start_Date, YEAR) END) AS Total_Years FROM harveaspace.data.newhr_table; </pre>
CORRECTING INCONSISTENT OR ERRONEOUS DATA	
Updated Employee Neha Patel's Address	<pre> UPDATE `harveaspace.data.newhr_table` SET Address = '123 Elm St' WHERE First_Name = 'Neha' AND Last_Name = 'Patel'; </pre>
HANDLING MISSING VALUES	
Missing data checks	<pre> SELECT COUNTIF(Employee_ID IS NULL) AS Missing_Employee_ID, COUNTIF(First_Name IS NULL) AS Missing_First_Name, COUNTIF&gt;Last_Name IS NULL) AS Missing_Last_Name, COUNTIF(SSN IS NULL) AS Missing_SSN, COUNTIF(Birth_Date IS NULL) AS Missing_Birth_Date, COUNTIF(Sex IS NULL) AS Missing_Sex, COUNTIF(Address IS NULL) AS Missing_Address, COUNTIF(Job_ID IS NULL) AS Missing_Job_ID, COUNTIF(Salary IS NULL) AS Missing_Salary, COUNTIF(Department_ID IS NULL) AS Missing_Department_ID, </pre>

	<pre> COUNTIF(Department_Name IS NULL) AS Missing_Department_Name, COUNTIF(Location_ID IS NULL) AS Missing_Location_ID, COUNTIF(Start_Date IS NULL) AS Missing_Start_Date, COUNTIF(Job_Title IS NULL) AS Missing_Job_Title, COUNTIF(Minimum_Salary IS NULL) AS Missing_Minimum_Salary, COUNTIF(Maximum_Salary IS NULL) AS Missing_Maximum_Salary FROM `harveaspace.data.newhr_table` </pre>
	<pre> -- Check if all Location_ID values exist in the Employee_ID column SELECT DISTINCT t.Location_ID FROM `harveaspace.data.newhr_table` t LEFT JOIN `harveaspace.data.newhr_table` e ON t.Location_ID = e.Employee_ID WHERE e.Employee_ID IS NULL </pre>
Manage Missing values	<pre> -- Replace missing values in Address column with 'Unknown' UPDATE `harveaspace.data.newhr_table` SET Address = 'Unknown' WHERE Address IS NULL </pre>
UNIQUE CONSTRAINT	
Duplicate checks	<pre> SELECT * FROM ( SELECT *,     ROW_NUMBER() OVER (PARTITION BY Employee_ID, First_Name, Last_Name, SSN, Birth_Date, Sex, Address, Job_ID, Salary, Department_ID, Department_Name, </pre>

	<pre> Location_ID, Start_Date, Job_Title, Minimum_Salary, Maximum_Salary ORDER BY Employee_ID) AS row_num FROM `harveaspace.data.newhr_table` ) WHERE row_num &gt; 1; </pre>
Duplicate checking in specific columns	<pre> -- Check for duplicate Employee_IDs SELECT Employee_ID, COUNT(*) AS Duplicate_Count FROM `harveaspace.data.newhr_table` GROUP BY Employee_ID HAVING COUNT(*) &gt; 1 </pre>
VALIDATING DATA INTEGRITY	
actual minimum and maximum values for the salary range you want to validate	<pre> SELECT * FROM `harveaspace.data.newhr_table` WHERE Salary BETWEEN 50000 AND 80000; </pre>
This statement created a new table named newhr_table_distinct.	<pre> CREATE OR REPLACE TABLE `harveaspace.data.newhr_table_distinct` AS SELECT DISTINCT * FROM `harveaspace.data.newhr_table`; </pre>
STANDARDIZING DATA FORMATS	
Department_Name length	<pre> SELECT LENGTH(Department_Name) AS Department_Name_Length FROM `harveaspace.data.newhr_table_distinct` </pre>
Find the department more than 8 letters	<pre> SELECT Department_Name FROM `harveaspace.data.newhr_table_distinct` WHERE LENGTH(Department_Name) &gt; 8 </pre>
Change "marketing" to "marketing by filed"	<pre> UPDATE `harveaspace.data.newhr_table_distinct` </pre>

	<pre> SET Department_Name = 'Marketing by Field' WHERE Department_Name = 'Marketing' </pre>
Removing spaces	<pre> -- Remove leading/trailing spaces from First_Name column  UPDATE `harveaspace.data.newhr_table_distinct` SET First_Name = TRIM(First_Name) WHERE First_Name IS NOT NULL </pre>
Date Format changing	<pre> ---To change the display format of the Birth_Date column from 'YYYY/MM/DD' to 'DD/MM/YYYY' without modifying the underlying data type,  SELECT Employee_ID, First_Name, Last_Name, FORMAT_DATE('%d/%m/%Y', Birth_Date) AS Formatted_Birth_Date, Sex, Address, Job_ID, Salary, Department_ID, Department_Name, Location_ID, Start_Date, Job_Title, Minimum_Salary, Maximum_Salary FROM `harveaspace.data.newhr_table_distinct` </pre>
DATA TYPE VALIDATION	

Data Type Validation:	<pre>-- Check if Birth_Date is a valid date SELECT Birth_Date FROM `harveaspace.data.newhr_table_distinct` WHERE SAFE_CAST(Birth_Date AS DATE) IS NULL</pre>
HANDLING OUTLIERS	
Outlier Detection:	<pre>-- Identify outliers in Maximum Salary column using z-score SELECT * FROM (   SELECT *,     ABS((Maximum_Salary - AVG(Maximum_Salary) OVER ()) / STDDEV(Maximum_Salary) OVER ()) AS z_score FROM `harveaspace.data.newhr_table_distinct` ) AS subquery WHERE Maximum_Salary IS NOT NULL AND z_score &gt; 3</pre>
Removing Irrelevant Data:	<pre>-- Remove the Location ID column from the table ALTER TABLE `harveaspace.data.newhr_table_distinct` DROP COLUMN Location_ID</pre>

## DATA ANALYSIS ABOUT SALARY

Questions	Execution Codes
How many employees are there in the dataset?	<pre>SELECT COUNT(*) AS Total_Employees FROM `harveaspace.data.newhr_table_distinct`</pre>

How many employees are there in each department?	<pre>SELECT Department_Name, COUNT(*) AS Total_Employees FROM `harveaspace.data.newhr_table_distinct` GROUP BY Department_Name</pre>
What is the employee count by job title?	<pre>SELECT Job_Title, COUNT(*) AS Total_Employees FROM `harveaspace.data.newhr_table_distinct` GROUP BY Job_Title</pre>
Most High Paying Departments:	<pre>SELECT     Department_Name,     MAX(Salary) AS Highest_Salary FROM     `harveaspace.data.newhr_table_distinct` GROUP BY     Department_Name ORDER BY     Highest_Salary DESC;</pre>
Most Low Paying Departments:	<pre>SELECT     Department_Name,     MIN(Salary) AS Lowest_Salary FROM     `harveaspace.data.newhr_table_distinct` GROUP BY     Department_Name ORDER BY     Lowest_Salary ASC;</pre>
Which are the jobs earn between 35000 to 50000?	<pre>SELECT DISTINCT Job_Title FROM `harveaspace.data.newhr_table_distinct` WHERE Salary &gt;= 35000 AND Salary &lt;= 50000</pre>
Most High Paying Job:	<pre>SELECT     Job_Title,     MAX(Salary) AS Highest_Salary FROM     `harveaspace.data.newhr_table_distinct` GROUP BY     Job_Title ORDER BY</pre>



	Highest_Salary DESC;
Most Low Paying Job:	<pre> SELECT     Job_Title,     MIN(Salary) AS Lowest_Salary FROM     `harveaspace.data.newhr_table_distinct` GROUP BY     Job_Title ORDER BY     Lowest_Salary ASC; </pre>
Most High Paying 5 Employees:	<pre> SELECT Employee_ID, First_Name, Last_Name, Salary FROM `harveaspace.data.newhr_table_distinct` ORDER BY Salary DESC LIMIT 5 </pre>
Most Low Paying 5 Employee:	<pre> SELECT     First_Name,     Last_Name,     Salary FROM     `harveaspace.data.newhr_table_distinct` ORDER BY     Salary ASC LIMIT 5; </pre>
High Paying Person and Experience:	<pre> SELECT     First_Name,     Last_Name,     Salary,     Start_Date,     DATEDIFF(NOW(), Start_Date) AS ExperienceDays FROM     `harveaspace.data.newhr_table_distinct` ORDER BY     Salary DESC LIMIT 1; </pre>

Low Paying Person and Experience:	<pre> SELECT     First_Name,     Last_Name,     Salary,     Start_Date,     DATEDIFF(NOW(), Start_Date) AS ExperienceDays FROM     `harveaspace.data.newhr_table_distinct` ORDER BY     Salary ASC LIMIT 1; </pre>
return the name,jobtitle,department and salary of employee that have a name of "ann".	<pre> SELECT     Department_Name,     Salary FROM     `harveaspace.data.newhr_table_distinct` WHERE     Department_Name IN (         SELECT             Department_Name         FROM             `harveaspace.data.newhr_table_distinct`         WHERE             Firstname = "Ann"     ) </pre>
How many employees have a salary above a certain threshold (e.g., \$80,000)?	<pre> SELECT COUNT(*) AS Employees_Above_Threshold FROM `harveaspace.data.newhr_table_distinct` WHERE Salary &gt; 80000 </pre>
What is the average salary by department?	<pre> SELECT Department_Name, AVG(Salary) AS Average_Salary FROM `harveaspace.data.newhr_table_distinct` GROUP BY Department_Name </pre>
Compare average,minmum and maximum salary range in department and salary	<pre> --Grouping and aggregating--- SELECT     Department_Name,     SUM(Salary) AS Total_Salary, </pre>

	<pre> AVG(Salary) AS Average_Salary, MIN(Salary) AS Minimum_Salary, MAX(Salary) AS Maximum_Salary FROM `harveaspace.data.newhr_table_distinct` GROUP BY Department_Name </pre>
calculates the average salary for each department by partitioning the data based on the Department_Name column.	<pre> --aggregation-- SELECT Department_Name, Salary, AVG(Salary) OVER (PARTITION BY Department_Name) AS Average_Salary FROM `harveaspace.data.newhr_table_distinct` </pre>
Sorting salary by department name	<pre> SELECT Department_Name, Salary FROM `harveaspace.data.newhr_table_distinct` ORDER BY Department_Name, Salary DESC </pre>
What is the salary range for each job title?	<pre> SELECT Job_Title, MIN(Salary) AS Minimum_Salary, MAX(Salary) AS Maximum_Salary FROM `harveaspace.data.newhr_table_distinct` GROUP BY Job_Title </pre>
What is the distribution of employees by gender?	<pre> SELECT Sex, COUNT(*) AS Total_Count FROM `harveaspace.data.newhr_table_distinct` GROUP BY Sex </pre>

## ANALYSIS ABOUT RECRUITMENT

How many employees were hired in each year?	<pre>SELECT EXTRACT(YEAR FROM Start_Date) AS Hire_Year, COUNT(*) AS Total_Employees FROM `harveaspace.data.newhr_table_distinct` GROUP BY Hire_Year ORDER BY Hire_Year</pre>
---	--

## ANALYSIS ABOUT COMPANY

Total payment company spends in one month?	<pre>SELECT SUM(Minimum_Salary) AS total_company_spends FROM `harveaspace.data.newhr_table_distinct`</pre>
Total payment company spends in one month for each departments?	<pre>SELECT Department_Name, SUM(Salary) AS Total_Salary FROM  `harveaspace.data.newhr_table_distinct` GROUP BY Department_Name;</pre>