**BUAN 6320.S01 DATABASE FOUNDATIONS FOR BUSINESS ANALYTICS.**

**SEMESTER PROJECT PHASE-2.**

**DONE BY: VIRAJA NELAGI.**
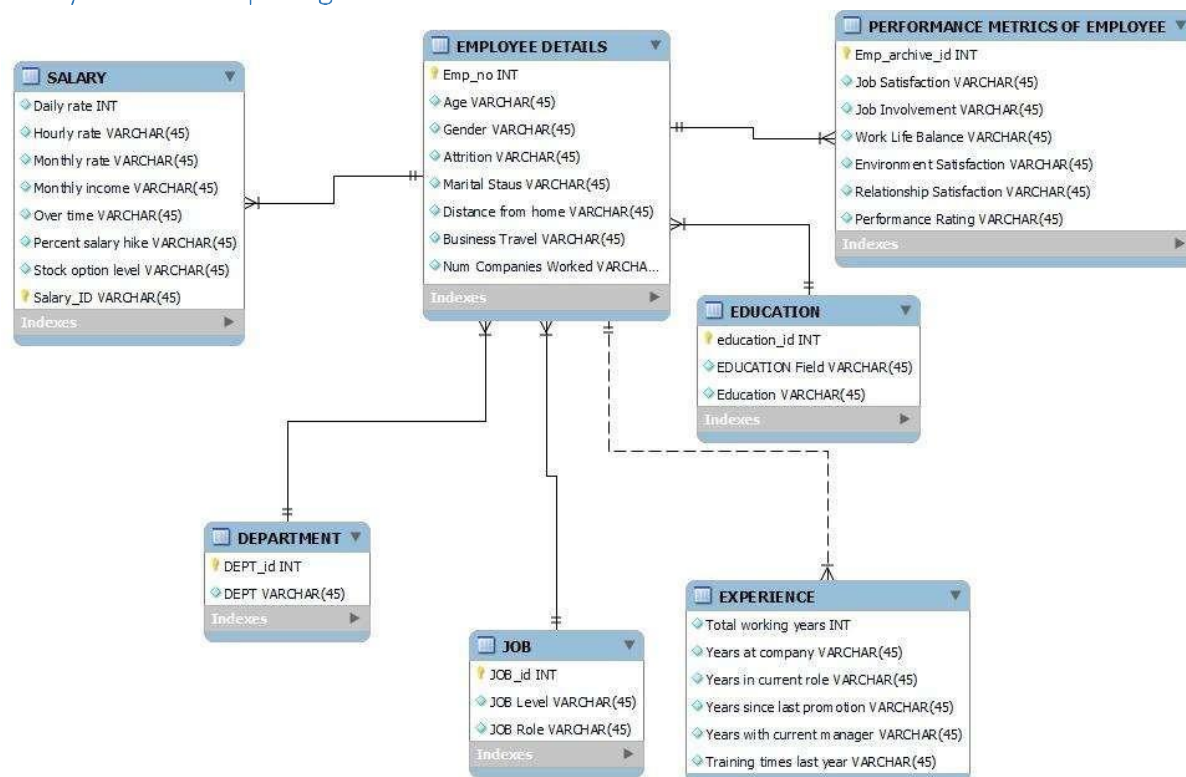
# Contents

# Relational Data Model

## Assumptions/Notes About Data Entities and Relationships

➢ Emp_no serves as the primary key for multiple tables and as a foreign key.

➢ The field of education's varied value is influenced by the Education and Education field columns. The employee's level of education is represented by their education, and their area of specialization within that field is represented by their education field and there are six various education fields in this database.

➢ A many-to-many link between JobLevel and JobRole. At each employment level, there are numerous job roles and there are nine different job roles at different levels in this database.

➢ The pay rates for various employee categories are represented by the hourly rate, overtime, daily rate, monthly rate, and monthly income. It also covers whether the worker has put in extra time.

➢ Dep_Id serves as primary key to department table.In the same way Education_ID ,Salary_ID and Job_ID serves as primary keys in their respective education ,salary and job tables.

➢ All the fields are completely dependent on the primary key of the respective table. There are no partial dependencies, duplicate data, and no transitive dependencies. Hence the model is in 3NF.

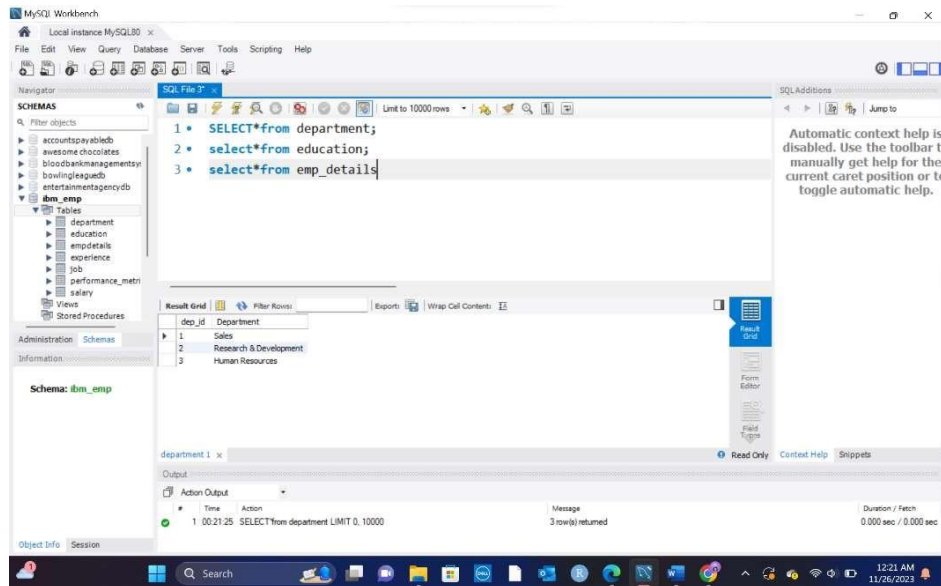## Entity-Relationship Diagram

# Physical MySQL Database

• Since the EmployeeCount column was not helping the model, it was eliminated. Its value remains constant at 1.

• Data in the Over-18 column is unnecessary. The age column can be used to deduce information from the over-18 column. Thus, the column was eliminated.

• The StandardHours column was eliminated as every employee's value was the same.
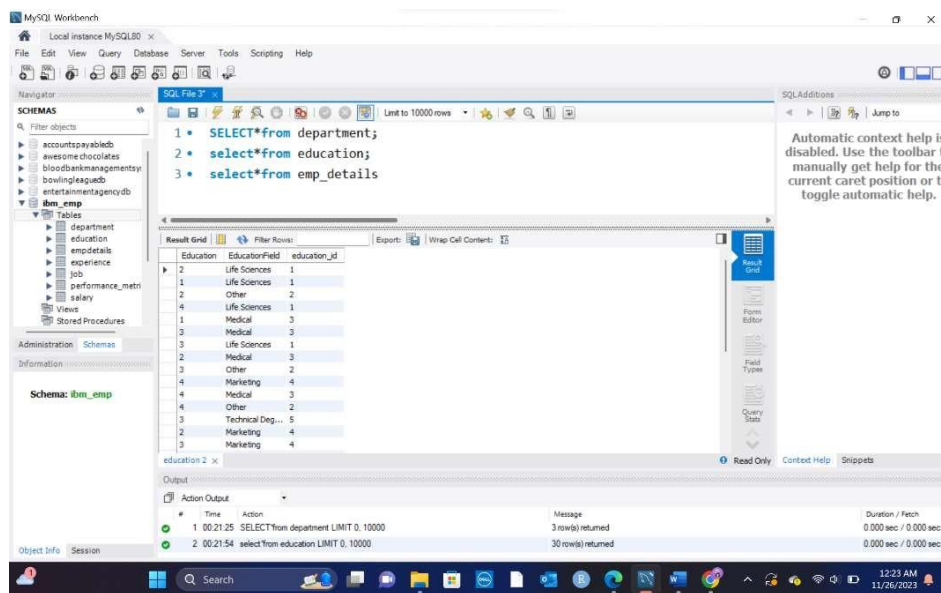
## Data in the Database

| Table Name | Primary Key | Foreign Key | # of Rows in Table |
|---|---|---|---|
| Department | Department_ID | | 3 |
| Education | Education_ID | | 30 |
| Employeedetails | Emp_no | Education_ID Job_ID DepartmentID | 1470 |
| Experience | | Emp_no | 1259 |
| _metrics_employee | Emp_archive ID | ➢ Emp_no | 1470 |
| Salary | Salary ID | ➢ Emp_no | 1470 |
| Job | Job ID | | 26 |

## Screen shot of Physical Database objects.

**Department Table:**

**Education table:**



**Employeedetails table:**

**Experience Table:**



**Performance_metrics_emp table:**

**Salary table:**



**Job Table:**

In Phase 2, we imported our data into MongoDB and wrote 3 queries, 5 queries from SQL.



**5) An employee in Sales department has complained to HR saying that females are paid less than males in the company, in all departments. What insight can you provide to prove or disprove that statement?**

**Result: Except in sales department, females are paid less in other departments of Human Resources and Research and Development. Hence, we disapprove the statement.**

**Translation**: Select gender and average of HourlyRate from salary joined on tables with empdetails table on s.Emp_no = ed.Emp_no joined with department ~~table~~ on ed.dep_id = d.dep_id. Group~~ed~~ by department and gender. Order by gender in ascending.

**Cleanup**: Select gender ~~and~~ average of hourly rate from salary join~~ed on tables~~ empdetails, salary matched ~~with~~ empno ~~in~~

salary ~~to~~ empno ~~in~~ emp_details ~~and~~ dep_id ~~in~~ empdetails ~~to~~ dep_id in department. Group~~ed~~ by department ~~and~~ gender, ~~sorted in ascending~~ order ~~on~~ gender.



**6)**

**A press article in a business magazine has said that at this company, married men have higher performance ratings than divorced or single men. What initial finding can you obtain from the data to help articulate the company's response in this regard?**

**Result: Married men have higher performance rating than Divorced and Single set of males. Hence, The press article's claim is correct.**

**Translation**: Select Gender, MaritalStatus, count of PerformanceReporting joined table performance_metrics_emp on ed.Emp_no=pm.Emp_no where ed.gender=male and pm.PerformanceRating > 3. Group by Gender and MaritalStatus.

**Cleanup**: Select Gender, MaritalStatus, count ~~of~~ PerformanceReporting joined ~~table~~ performance_metrics_emp on ed.Emp_no=pm.Emp_no where ed.gender=male and pm.PerformanceRating > 3. Group by Gender and MaritalStatus.

7)



**If the company wants to cut travel costs, which department should the company focus on?**

**Result: The company should focus on Research and Development department to cut down their costs as the amount of people who travel frequently are higher than other departments.**

**Translation**: Select Department, count of dep_id from department joined table empdetails on dep_id = ed.dep_id and ed.BusinessTravel = 'travel_frequently'. Grouped by department.
.
**Cleanup**: Select Department, count of dep_id from department joined ~~table~~ empdetails on dep_id = ed.dep_id and ed.BusinessTravel = 'travel_frequently'. Group~~ed~~ by department.



**4 – From the output of the query it is evident that the Research and Development department have more average job**

**satisfaction.**

**10- From the output of the query it is evident that the environment satisfaction score of HR is higher than sales but HR job satisfaction score is lower than Research & Development. The HR department is right.**

**Translation**: Select Department, average of JobSatisfaction, average of EnvironmentSatisfaction from performance_metrics_emp joined table empdetails on pm.Emp_no=ed.Emp_no and joined table department on ed.dep_id = d.dep_id. Grouped by department.

**Cleanup**: Select Department, average of JobSatisfaction, average of EnvironmentSatisfaction from performance_metrics_emp joined table empdetails on pm.Emp_no=ed.Emp_no and joined table department on ed.dep_id = d.dep_id. Grouped by department.

# Data Review for MongoDB
## Assumptions/Notes About Data Collections, Attributes and Relationships between Collections

Collection Structure: MongoDB stores data in a hierarchical format, complex and nested data structures are possible. The only collection in my database with the name "ibm_emp" is called "employee."

Data organization: The original MySQL database's tables are all combined into one collection by the Employee collection. This indicates that all the attributes from the original tables are present in every document in the collection.

Attributes: The attributes Over18 and EmployeeCount are discarded while loading the dataset into the MongodB as

o Over18 can be easily found by performing logic {$gte:18}

o Employeecount is just counting the same employee alone and has no meaning.

Relationships: There is only one collection.

# Physical Mongo Database
## Assumptions/Notes About Data Set

Data completeness: It is assumed that the data set is comprehensive and includes all pertinent details regarding the employees' job descriptions, pay scales, levels of job satisfaction, and other relevant characteristics.

Data consistency: It is assumed that the data set is consistent, which means that each attribute has a set format, and that the data is constant across all records.

Data reliability: It is assumed that the data set is trustworthy, and that the data can be relied upon for making decisions.

Timeframe: Any analysis or inferences made from the data should consider the time period for which the data set is assumed to be representative.

Data integrity: The extraction, transformation, and loading (ETL) process are assumed to have preserved the data set's integrity, preventing any data loss or corruption.

Data confidentiality: It is assumed that the data set will be handled with the proper degree of confidentiality and that any sensitive employee information will be adequately safeguarded.

# Screen shot of Physical Database objects (Database, Collections and Attributes)



# Data in the Database

| Collection Name | Relationships With Other Collections (if any) | # of Documents in Collection |
|---|---|---|
| ibm_emp | | 1470 |

# MongoDB Queries/Code

Pick 3 SQL queries and write them in MongoDB

## Question 3
A new employee from a Medical-related education field wants to work in Sales. Do you believe the company might be able to give her a chance to work in Sales? Why or why not?

## Notes/Comments About MongoDB Query/Code and Results (Include # of Documents in Result)
Based on their educational backgrounds, sales department employees' average performance rating has been determined. The medical sector ranks second among all education fields, while other fields make up the first place. As a result, the corporation may offer her a position in sales.

## Translation
Aggregate the emp_attrition collection the Database. In the first stage filter the documents using $match operator to and include only those whose Gender is set to "Female" and Department is set to "Sales". In the second stage group the filtered documents using $group operator, the grouping key is composed by "EducationField". Calculate the average value of PerformanceRating for each grouped documents assign it to the field avgPerformanceRating in the result.

## Screen Shot of MongoDB Query/Code and Results

# Mongo Query 2

## Question 8
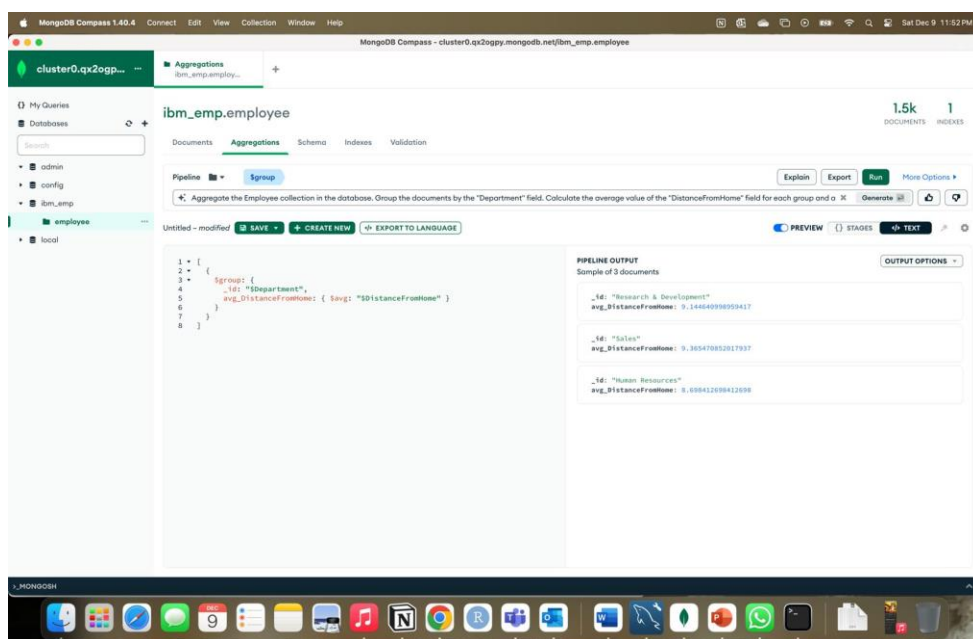Screen Shot of MongoDB Query/Code and Results

## Question 8, from the list
### Notes/Comments About MongoDB Query/Code and Results (Include # of Documents in Result)

After calculating the average distance from home for each department using the Mongo query below, the employees in the "Sales" department will receive the highest wage. The reason for this is that the sales section is the furthest away from home. The average distance from home for each department is listed in three papers in the result set.

## Translation

Aggregate the emp_attrition collection in the database. Group the documents by the "Department" field. Calculate the average value of the "DistanceFromHome" field for each group and assign it to the field "avg_DistanceFromHome" in the result.

## Question 11

An employee from Medical education field working in Sales department has spread a rumor saying that employees with his educational background are paid more in Research & Development than in Sales. What insight can you provide to prove or disprove that statement?

## Notes/Comments About MongoDB Query/Code and Results (Include # of Documents in Result)

Determine the mean income of medical education field personnel working in the "Sales" and "Research & Development" divisions.

•It is evident that the Research & Development Department pays more. This can therefore be used as proof to back up the rumor.

• The outcome contains two documents: one for the medical field sales department and one for Department of Research and Development in the Medical FieldTranslation

## Translation

Aggregate the emp_attrition collection in the database. In the first stage, filter the documents using the $match operator to include only those with the EducationField set to "Medical" and the Department set. to either "Sales" or "Research & Development". In the second stage, group the filtered documents using. the $group operator. The grouping key is composed of the Department and EducationField fields. Calculate the average value of the MonthlyIncome field for each group and assign it to the field. avg_Salary in the result.