# Fall 2025 Data Science Bootcamp

## Week 4: SQL Fundamentals

## Take-home assignment

**Section 1 (Mondays): 2:00PM-4:00PM EST**

**Instructor: Shaahid Ahmed Nadeem**

**Made by: Mahdi (Matt) Ghadimi**

**N14369173**

mg8786@nyu.edu

NYU | TANDON SCHOOL OF ENGINEERING

## 1. Actors and Directors Who Cooperated At Least Three Times:

SQL:

```
1. # Write your MySQL query statement below
2. SELECT actor_id, director_id
3. FROM ActorDirector
4. GROUP BY actor_id, director_id
5. HAVING COUNT(*) >= 3;
```

Pandas:

```
1. import pandas as pd
2.
3. def actors_and_directors(actor_director: pd.DataFrame) -> pd.DataFrame:
4.     result = (
5.         actor_director
6.         .groupby(['actor_id', 'director_id'])
7.         .size()
8.         .reset_index(name='cnt')
9.         .query('cnt >= 3')[['actor_id', 'director_id']]
10.    )
11.    return result
```

## 2. Fix Names in a Table

SQL:

```
1. SELECT
2.     user_id,
3.     CONCAT(UPPER(LEFT(name, 1)), LOWER(SUBSTRING(name, 2))) AS name
4. FROM Users
5. ORDER BY user_id;
```

Pandas:

```
1. import pandas as pd
2.
3. def fix_names(users: pd.DataFrame) -> pd.DataFrame:
4.     users['name'] = users['name'].str.capitalize()
5.     return users.sort_values('user_id')
```

## 3. Combine Two Tables

SQL:

```
1. SELECT
2.     Person.firstName,
3.     Person.lastName,
4.     Address.city,
5.     Address.state
6. FROM Person
7. LEFT JOIN Address
8.     ON Person.personId = Address.personId;
```

Pandas:

```
1. import pandas as pd
2.
3. def combine_two_tables(person: pd.DataFrame, address: pd.DataFrame) -> pd.DataFrame:
4.     result = person.merge(
5.         address[['personId', 'city', 'state']],
```

NYU | TANDON SCHOOL OF ENGINEERING

```
6.          on='personId',
7.          how='left'
8.      )
9.      return result[['firstName', 'lastName', 'city', 'state']]
```

## 4. Second Highest Salary

SQL:

```
1. SELECT
2.      (SELECT DISTINCT salary
3.       FROM Employee
4.       ORDER BY salary DESC
5.       LIMIT 1 OFFSET 1) AS SecondHighestSalary;
```

Pandas:

```
1. import pandas as pd
2.
3. def second_highest_salary(employee: pd.DataFrame) -> pd.DataFrame:
4.      distinct_salaries = employee['salary'].drop_duplicates().sort_values(ascending=False)
5.      second = distinct_salaries.iloc[1] if len(distinct_salaries) > 1 else None
6.      return pd.DataFrame({'SecondHighestSalary': [second]})
```

## 5. List the Products Ordered in a Period

SQL:

```
1. SELECT
2.      p.product_name,
3.      SUM(o.unit) AS unit
4. FROM Orders o
5. JOIN Products p
6.      ON o.product_id = p.product_id
7. WHERE o.order_date >= '2020-02-01'
8.   AND o.order_date < '2020-03-01'
9. GROUP BY p.product_name
10. HAVING SUM(o.unit) >= 100;
```

Pandas:

```
1. import pandas as pd
2.
3. def list_products(products: pd.DataFrame, orders: pd.DataFrame) -> pd.DataFrame:
4.      # Filter orders in Feb 2020
5.      feb_orders = orders[
6.          (orders['order_date'] >= '2020-02-01') &
7.          (orders['order_date'] < '2020-03-01')
8.      ]
9.
10.      # Sum units by product_id
11.      summed = feb_orders.groupby('product_id', as_index=False)['unit'].sum()
12.
13.      # Filter products with at least 100 units
14.      filtered = summed[summed['unit'] >= 100]
15.
16.      # Join with products to get product_name
17.      result = filtered.merge(products[['product_id', 'product_name']], on='product_id')
18.
19.      return result[['product_name', 'unit']]
```

## 6. Replace Employee ID With The Unique Identifier

SQL:

```
1. SELECT
2.     eui.unique_id,
3.     e.name
4. FROM Employees e
5. LEFT JOIN EmployeeUNI eui
6.     ON e.id = eui.id;
```

Pandas:

```
1. import pandas as pd
2.
3. def replace_employee_id(employees: pd.DataFrame, employee_uni: pd.DataFrame) -> pd.DataFrame:
4.     result = employees.merge(
5.         employee_uni[['id', 'unique_id']],
6.         on='id',
7.         how='left'
8.     )
9.     return result[['unique_id', 'name']]
```

## 7. Game Play Analysis IV

SQL:

```
1. WITH first_login AS (
2.     SELECT player_id, MIN(event_date) AS first_date
3.     FROM Activity
4.     GROUP BY player_id
5. )
6. SELECT ROUND(
7.     SUM(
8.         CASE
9.             WHEN EXISTS (
10.                SELECT 1
11.                FROM Activity a2
12.                WHERE a2.player_id = f.player_id
13.                  AND a2.event_date = DATE_ADD(f.first_date, INTERVAL 1 DAY)
14.            ) THEN 1 ELSE 0
15.        END
16.    ) * 1.0 / COUNT(*), 2
17. ) AS fraction
18. FROM first_login f;
```

Pandas:

```
1. import pandas as pd
2.
3. def gameplay_analysis(activity: pd.DataFrame) -> pd.DataFrame:
4.     activity['event_date'] = pd.to_datetime(activity['event_date'])
5.
6.     first_login = activity.groupby('player_id')['event_date'].min().reset_index()
7.
8.     def logged_next_day(row):
9.         return ((activity['player_id'] == row['player_id']) &
10.                (activity['event_date'] == row['event_date'] + pd.Timedelta(days=1))).any()
11.
12.     count_next_day = first_login.apply(logged_next_day, axis=1).sum()
13.     total_players = first_login.shape[0]
14.
15.     fraction = round(count_next_day / total_players, 2)
```

```
16.
17.     return pd.DataFrame({'fraction': [fraction]})
```

## 8. Project Employees I

SQL:

```
1. SELECT
2.     p.project_id,
3.     ROUND(AVG(e.experience_years), 2) AS average_years
4. FROM Project p
5. JOIN Employee e
6.     ON p.employee_id = e.employee_id
7. GROUP BY p.project_id;
```

Pandas:

```
1. import pandas as pd
2.
3. def project_employees_i(project: pd.DataFrame, employee: pd.DataFrame) -> pd.DataFrame:
4.     merged = project.merge(employee[['employee_id', 'experience_years']], on='employee_id')
5.
6.     result = merged.groupby('project_id', as_index=False)['experience_years'].mean()
7.
8.     result['average_years'] = result['experience_years'].round(2)
9.
10.     return result[['project_id', 'average_years']]
```

## 9. Department Top Three Salaries

SQL:

```
1. WITH RankedSalaries AS (
2.     SELECT
3.         e.id,
4.         e.name AS Employee,
5.         e.salary AS Salary,
6.         d.name AS Department,
7.         DENSE_RANK() OVER(PARTITION BY e.departmentId ORDER BY e.salary DESC) AS salary_rank
8.     FROM Employee e
9.     JOIN Department d
10.        ON e.departmentId = d.id
11. )
12. SELECT Department, Employee, Salary
13. FROM RankedSalaries
14. WHERE salary_rank <= 3;
```

Pandas:

```
1. import pandas as pd
2.
3. def top_three_salaries(employee: pd.DataFrame, department: pd.DataFrame) -> pd.DataFrame:
4.     merged = employee.merge(department, left_on='departmentId', right_on='id', suffixes=('',
'_dept'))
5.
6.     merged['salary_rank'] = merged.groupby('departmentId')['salary'] \
7.                             .rank(method='dense', ascending=False)
8.
9.     top_earners = merged[merged['salary_rank'] <= 3]
10.
11.     return top_earners[['name_dept', 'name', 'salary']].rename(
12.         columns={'name_dept': 'Department', 'name': 'Employee', 'salary': 'Salary'}
13.     )
```