

**NBA Players and Teams Relational Database**

Mingkuan Pang, Mingzan Liu, Yilong Wang, Tien Tran, and Johnny Dean

CSCI 126

Dr. David Ruby

April 26, 2022

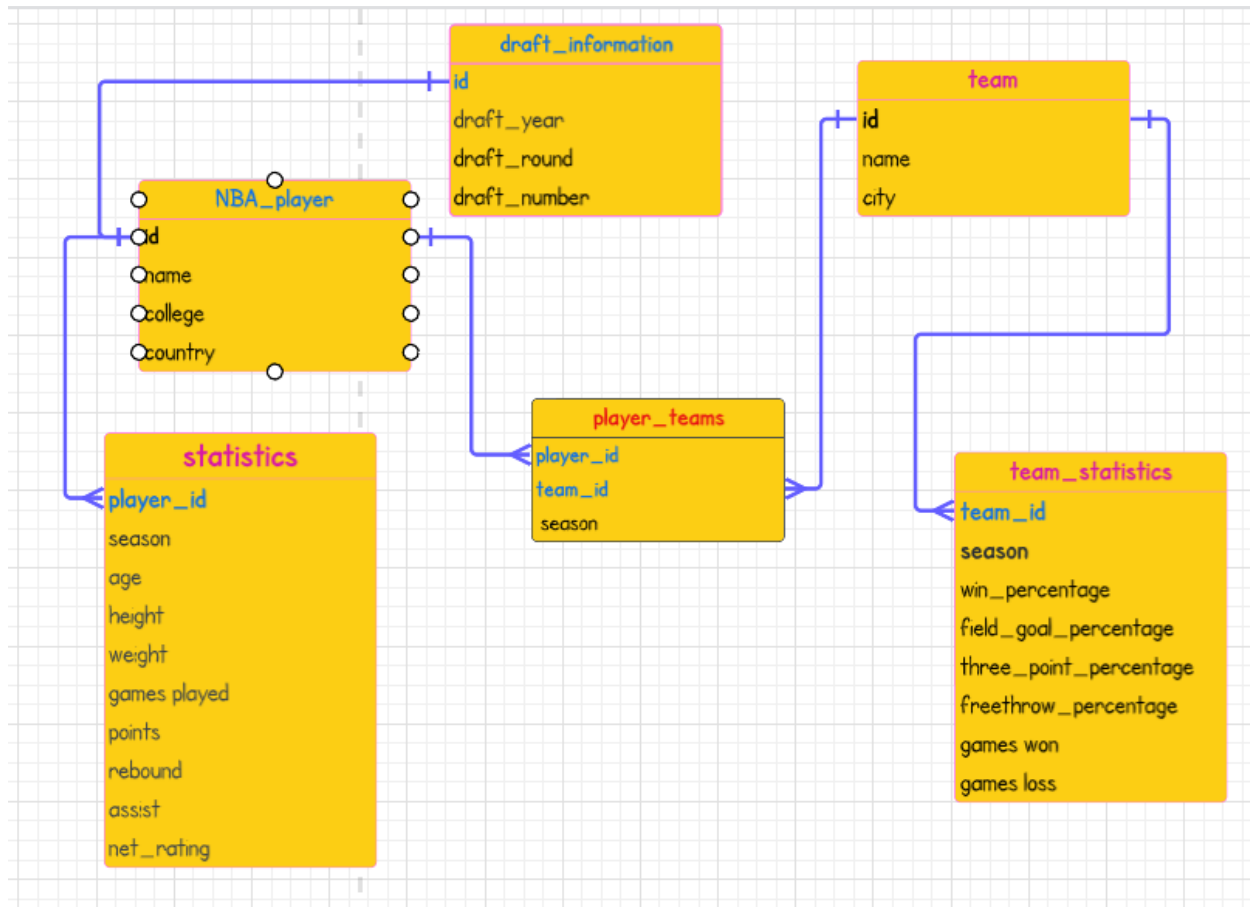
### **Introduction:**

This project involves the search and augmentation of data. Using the data, design a database that includes three to five relations using Entity-Relationship Diagrams (ERD). Which includes one-to-one relations, many-to-many relations, one-to-many relations, and weak entity sets or ISA hierarchy. After designing the database, develop the design into a database schema and include key definitions and referential integrity constraints (foreign key). Then, identify the functional dependencies or multi-valued dependencies for the database and determine whether it is free from violations for the 3rd Normal Form, Boyce-Codd Normal Form, and 4th Normal Form. Describe queries that can be used to interact with the database, such as subqueries, aggregation, insert queries, updates queries, etc. Lastly, include information about the type of applications that will run on the designed database, as well as include the use cases, not just support application functionality.

### **Data Selection:**

For the NBA players and teams database, data are drawn from two datasets. One dataset contains the profile information and statistics of the NBA players from the 2019 - 2020 season and the 2020 - 2021 season. The other dataset contains the profile and the statistics of a team from the 2019 - 2020 season and the 2020 - 2021 season. Specifically from the NBA players dataset, the name, season, college, country, team, age, height, weight, games played, points, rebound, assist, net\_rating, draft\_year, draft\_round, and draft\_number are selected. From the NBA teams dataset, the name, city, season, win percentage, field goal percentage, three-point percentage, free-throw percentage, games won, and games loss are selected. From these data points of the two datasets, a database is designed and visualized through the Entity-Relationship Diagrams. The mentioned data is linked in the Links section.

## Database Design:



In the Entity-Relationship Diagram above, the one-to-one relation, many-to-many relations, one-to-many relations, and weak entity set are marked. (1) There is a one-to-one relation between the **NBA\_player** entity and the **draft\_information** weak entity. Specifically, the **id** of the **draft\_information** weak entity is a foreign key, and the **id** of the **NBA\_player** entity is the key that the **draft\_information** weak entity depends on. (2) The **NBA\_player** entity has a one-to-many relation with the **statistics** weak entity. **Statistics** weak entity has a foreign key **player\_id** that depends on the **id** of the **NBA\_player** entity. Furthermore, one **NBA\_player** can have multiple statistics depending on the candidate key **player\_id** and **season** in the **statistics** weak entity; therefore, the **NBA\_player** entity is a one-to-many relation with the statistic weak entity. The **team** entity has a one-to-many relation with the **team\_statistic** weak entity. The

**team\_id** in the team\_statistic weak entity depends on the key in the team entity. Furthermore, the team can have multiple team\_statistic records depending on the candidate key **team\_id** and **season** in the team\_statistics weak entity; thus, making this relation a one-to-many relation. This Entity-Relationship Diagram also has a many-to-many relation. The NBA\_player entity has a one-to-many relation with the player\_teams weak entity because one unique **player\_id** can have multiple **team\_id** depending on the season attribute. The team entity has a one-to-many relation with the player\_teams weak entity because one unique **team\_id** can have many different **player\_id**. As a result, the player\_teams entity contains multiple unique **player\_id** for a **team\_id**, and a **player\_id** can have multiple **team\_id** depending on the season; creating a many-to-many relation. As mentioned earlier, team\_statistics, player\_teams, and statistics are weak entities because they can not determine unique attributes without a foreign key, or a key from another table.

### Database Schema:

From the Entity-Relationship Diagram in the Database Design section, the schema NBA is created in the database. The tables NBA\_player, statistics, draft\_information, player\_teams, team, and team\_statistics are created within the NBA schema. The table NBA\_player contains the primary key **id** and attributes **name**, **college**, and **country**. The statistics table contains the referential integrity constraint **player\_id**, which is the primary key **id** of the NBA\_player table; and attributes **season**, **age**, **height**, **weight**, **games played**, **points**, **rebound**, **assist**, and **net\_rating**. The draft\_information table contains the referential integrity constraint **id** that is the primary key **id** of the NBA\_player table. It also has the attributes **draft\_year**, **draft\_round**, and **draft\_number**. The team table contains the primary key **id** and attributes **name** and **city**. The team\_statistics table contains a referential integrity constraint **team\_id** that is the primary key **id**

of the team table. It also has attributes **season**, **win\_percentage**, **field\_percentage**, **three\_point\_percentage**, **free-throw\_percentage**, **games won**, and **games loss**. The player\_teams table contains two referential integrity constraints. The referential integrity constraint **player\_id** is the primary key **id** of the NBA\_player table. The referential integrity constraint **team\_id** is the primary key **id** of the team table. The relations used in the Entity-Relationship Diagram persist in the created schema through the referential integrity constraints and primary keys of the tables. The schema and table creation SQL queries are linked in the Links section.

### Normal Forms:

The functional dependencies and multi-valued dependencies of each table are tested for the 3rd Normal Form, Boyce-Codd Normal Form, and 4th Normal Form violations.

**Table 1: NBA\_player** NBA\_player = {id, name, college, country}

**Functional Dependencies:** id -> {name, college, country}

**3rd form:** No violation.

**Explanation:** (1) The candidate key, **id**, can uniquely determine the attributes: **name**, **college**, and **country**. (2) This relation does not contain any partial dependency and it conforms to 1st normal form.

**BCNF:** No violation.

**Explanation:** (1) The candidate key, **id**, can uniquely determine the attributes: **name**, **college**, and **country**. Since the determinant of functional dependency **id** is a key, it is

conforming to BCNF. (2) Since it is already conforming to the 3rd form, and the determinant **id** of the functional dependency is a key, it conforms to BCNF.

**4th form:** Violation.

**Explanation:** (1) The functional dependency  $id \rightarrow \{name, college, country\}$  is a multivalued dependency.  $id \twoheadrightarrow name$ ,  $id \twoheadrightarrow college$ , and  $id \twoheadrightarrow country$ . (2) This table conforms with case 1: conforming to BCNF. But it contains multivalued dependencies, therefore, this table is not in 4th form.

**Table 2: Statistics** Statistics = {player\_id, season, age, height, weight, games played, points, rebound, assist, net\_rating}

**Functional dependencies:** {player\_id, season}  $\rightarrow$  {age, height, weight, games played, points, rebound, assist, net\_rating}

**3rd form:** No violation.

**Explanation:** (1) The candidate key {player\_id, season} can uniquely determine the attributes: **age, height, weight, games played, points, rebound, assist, and net\_rating**. (2) This relation does not contain any partial dependency and it conforms to 1st normal form.

**BCNF:** No violation.

**Explanation:** (1) The candidate key {player\_id, season} can uniquely determine the attributes: **age, height, weight, games played, points, rebound, assist, and net\_rating**. Since the determinant of functional dependency {player\_id, season} is a key, it is

conforming to BCNF. (2) Since it is already conforming to the 3rd form, and the determinant **{player\_id, season}** of the functional dependency is a key, it conforms to BCNF.

**4th form:** Violation.

**Explanation:** (1) The functional dependency {player\_id, season} -> {age, height, weight, games played, points, rebound, assist, net\_rating} has multivalued dependencies.

**{player\_id, season} ->> age, {player\_id, season} ->> height, {player\_id, season} ->> weight, {player\_id, season} ->> games played, {player\_id, season} ->> points, {player\_id, season} ->> rebound, {player\_id, season} ->> assist, and {player\_id, season} ->> net\_rating.** (2) This table conforms with case 1: conforming to BCNF. But it contains multivalued dependencies, therefore, this table is not in 4th form.

**Table 3: Team** team = {id, name, city}

**Functional dependencies:** id -> {name, city}

**3rd form:** No violation

**Explanation:** (1) The candidate key, **id**, can uniquely determine the attributes: **name**, and **city**. (2) This relation does not contain any partial dependencies, and it conforms to 1st normal form.

**BCNF:** No violation

**Explanation:** (1) The candidate key, **id**, can uniquely determine the attributes: **name**, **city**. (2) Since the determinant **id** of the functional dependency is a key, and it conforms to the 3rd form; it conforms to BCNF.

**4th form:** Violation

**Explanation:** (1) For the functional dependency  $id \rightarrow \{name, city\}$ , it has multivalued dependencies.  $id \twoheadrightarrow name$ , and  $id \twoheadrightarrow city$ . Multivalued dependencies are not allowed in the 4th normal form, therefore it violates the 4th form.

**Table 4: Player\_teams** Player\_teams = {player\_id, team\_id, season}

**Function dependencies:** {season, player\_id}  $\rightarrow$  team\_id

**3rd form:** No violation

**Explanation:** (1) The candidate key {season, player\_id} can uniquely determine the attributes: **team\_id**. (2) This relation does not contain any partial dependencies, and it conforms to 2nd normal form.

**BCNF:** No violation

**Explanation:** (1) The candidate key {season, player\_id} can uniquely determine the **team\_id**. (1) This functional dependency conforms to BCNF because the determinant {season, player\_id} is a key; and it conforms to the 3rd normal form.

**4th form:** No violation



**Explanation:** (1) For the functional dependency {season, player\_id} → team\_id, it conforms to the 4th form. (2) The candidate key {season, player\_id} only determines team\_id and nothing else. As well as the functional dependency already conforms to BCNF.

**Table 5: team\_statistics** team\_statistics = {team\_id, season, win\_percentage, field\_goal\_percentage, three\_point\_percentage, freethrow\_percentage, games won, games loss}

**Functional dependencies:** {team\_id, season} → { win\_percentage, field\_goal\_percentage, three\_point\_percentage, freethrow\_percentage, games won, games loss}

**3rd form:** No violation

**Explanation:** (1) The candidate key, {team\_id, season}, can uniquely determine the attributes: win\_percentage, field\_goal\_percentage, three\_point\_percentage, freethrow\_percentage, games won, and games loss. (2) This relation does not contain any partial dependencies, and it conforms to 1st normal form.

**BCNF:** No violation

**Explanation:** (1) The candidate key, {team\_id, season}, can uniquely determine the attributes: win\_percentage, field\_goal\_percentage, three\_point\_percentage, freethrow\_percentage, games won, and games loss. (2) Since the determinant {team\_id, season} of the functional dependency is a key, and it conforms to the 3rd form; it conforms to BCNF.

**4th form:** Violation

**Explanation:** (1) For the functional dependency  $\{\text{team\_id}, \text{season}\} \rightarrow \{\text{win\_percentage}, \text{field\_goal\_percentage}, \text{three\_point\_percentage}, \text{freethrow\_percentage}, \text{games won}, \text{games loss}\}$ , it has multivalued dependencies.  $\{\text{team\_id}, \text{season}\} \twoheadrightarrow \text{win\_percentage}$ ,  $\{\text{team\_id}, \text{season}\} \twoheadrightarrow \text{field\_goal\_percentage}$ ,  $\{\text{team\_id}, \text{season}\} \twoheadrightarrow \text{three\_point\_percentage}$ ,  $\{\text{team\_id}, \text{season}\} \twoheadrightarrow \text{freethrow\_percentage}$ ,  $\{\text{team\_id}, \text{season}\} \twoheadrightarrow \text{games won}$ , and  $\{\text{team\_id}, \text{season}\} \twoheadrightarrow \text{games loss}$ . (2) Despite it conforms to BCNF, Multivalued dependencies are not allowed in the 4th normal form, therefore it violates the 4th form.

**Table 6: draft\_information** draft\_information = {id, college, draft\_round, draft\_number, country}

**Functional dependencies:** id  $\rightarrow$  {draft\_year, draft\_round, draft\_number}

**3rd form:** No violation

**Explanation:** (1) The candidate key, **id**, can uniquely determine the attributes: **draft\_year, draft\_round, draft\_number**. (2) This relation does not contain any partial dependencies, and it conforms to 1st normal form.

**BCNF:** No violation

**Explanation:** (1) The candidate key, **id**, can uniquely determine the attributes: **draft\_year, draft\_round, draft\_number**. (2) Since the determinant **id** of the functional dependency is a key, and it conforms to the 3rd form; it conforms to BCNF.

**4th form:** Violation

**Explanation:** (1) For the functional dependency  $id \rightarrow \{draft\_year, draft\_round, draft\_number\}$ , it has multivalued dependencies.  **$id \twoheadrightarrow draft\_year$ ,  $id \twoheadrightarrow draft\_round$ ,  $id \twoheadrightarrow draft\_number$** . (2) Despite it conforms to BCNF, Multivalued dependencies are not allowed in the 4th normal form, therefore it violates the 4th form.

### Queries:

Many queries of use are created for this **nba** database. It includes aggregations, insert queries, subqueries, and update queries. For insert queries, inserting a player, inserting a player's statistics, inserting a player who is at least 18 years old (constraint), and inserting a non-unique player (constraint) are available for use in this database. Inserting a player requires a tuple (**player\_id**, **name**, **college**, and **country**). **player\_id** is required and has to be unique from every entry that is currently in the database. The other attributes are optional and will be set to NULL if left empty. Similar to inserting a player, inserting a player's statistics requires the unique key of **player\_id**, and **season** to uniquely identify a player in a particular season; other attributes (**age**, **height**, **weight**, **games played**, **points**, **rebound**, **assist**, **net\_rating**) are optional and set to NULL if left empty. When inserting a player that is less than the age of 18, which is a constraint created at table creation, it will be invalidated. Similarly, for inserting a non-unique player, it will be invalidated automatically without using the *CHECK constraint* command due to **player\_id** being a unique key. The constraint is checked using *CHECK constraint* command before the insertion. Test queries are provided to check if the player and statistics are inserted successfully and correctly. For update queries, updating a player's team and updating a particular statistic of a player are provided for use in this database. In order to update a player's team, the **player\_id** in the **NBA\_player** table and the **player\_id** in the table **player\_teams** are required. For updating a player's statistic, specifically, increment the rebound by one. The **player\_id** is required to

identify the correct player for a particular season. The rebound incrementation is performed in the *SET* command. Again, test queries are provided to check whether the update query was successful and performed correctly.

For subqueries and aggregations, get players with greater than average height, get teams with greater than average win rate, get teams with an average height greater than the average height across all teams, and get the teams with greater than average **heights** who have a greater than average win rate. For getting players with greater than average **heights**, an aggregation *AVG()* is used inside a subquery in order to get the average **height** in the *player\_statistics* table. Similarly, for getting teams that have a greater than average win rate, an aggregation *AVG()* is used inside a subquery to get the average of the **win\_percentage** in the *team\_statistics* table. For finding the teams with a team **height** average that is above average **height** across all the **NBA\_players**. We used two *AVG()* aggregations, one inside a subquery to find the average player **height**. Another is to see if the team satisfies the condition that the average **height** of the team is greater than the average **height** across all the players. The result of this query is grouped by the **team\_name** and ordered by **win\_percentage** in descending order. For finding teams with greater than average **heights** that also have greater than average win rates, two aggregations, and two subqueries are used. One *AVG()* is used inside a subquery that finds the average **height** of the players in the *player\_statistics* table. Another *AVG()* is used inside a subquery that finds the average **win\_percentage** from the *team\_statistics* table. The result of this query is grouped by the **team\_name** and ordered by **win\_percentage** in descending order.

Further documentation of the queries is made inside their respective files. The mentioned queries are linked in the Links section.

### Application and Use Cases:

This NBA database allows applications such as basketball-related games to determine if a player exists in the NBA during a particular season and is performing well before they decide to create this player inside the game. If it is a less-known player, then it is not worthwhile for the application to create this player inside the game and generate revenue. Also, this database allows basketball-related games to determine the in-game worth of a player based on the information stored inside the database. For example, the database stores a player's **shooting\_percentage** and **net\_rating** of a particular season. The application can use this information to decide the in-game value of this player based on the records stored inside the database. This database can provide functionalities beyond applications. For example, the query for getting teams with a **win\_percentage** above the average **win\_percentage** of a particular season (mentioned in the Queries section). This can provide analysis about which team will likely perform well in the post-season and win the championship as these data are from the regular season, before the championship games. Also, the query for getting the player with above-average **height** (mentioned in the Queries section) can provide useful analysis. In particular, this query can determine if a particular player's **height** affects his team's **win\_percentage**. Similar analysis can be made using this database through small alterations to the query that finds players with above-average **height**. The query can be changed to players with above-average **shooting\_percentage** instead by changing the **height** attribute to the **shoot\_percentage** attribute. By doing so, analysis can be made about if a player with an above-average **shooting\_percentage** affects his team's **win\_percentage**.

**Loading Collected Data Instructions:**

These instructions assume that MySQL is used as the database. Some of these instructions might apply or are similar to other databases.

1. Download the csv and sql files from the Github repository.
2. Create a new connection in MySQL and run the create.sql file inside the new connection to create the schema and tables.
3. Inside the schema, right-click on the tables directory to access the table data import wizard and import the NBA\_player.csv, player\_statistics.csv, draft\_information.csv, team.csv, team\_statistics.csv, and player\_teams.csv in this specific order. During the import, make sure the source columns match the destination columns.
4. Queries in the insert\_and\_update.sql and sub\_queries.sql files can be run after importing the data.

**Links:**

The datasets, Entity-Relationship Diagrams, database and table creation queries, SQL files with testable queries, and documentation are inside the Github repository below.

Link: <https://github.com/Venceyv/DBMS-NBAplayer>