

**Assignment – Relational and Dimensional Model, ETL****(20% of final exam marks)***Due Date: Monday 4th January (midnight)*

In this assignment you are required to create a relational database, populate it, create a dimensional model and perform an ETL process at two different points in time.

Your dataset is about players' statistics during the 2011/2012 Premier League season. For each player, the dataset contains match-level stats such as number of passes, goals, shots...

The relational database you are asked to create must have the following schema:

Table Player

<u>Player_ID</u>	Pl_name	Pl_surname	TeamID (FK)
------------------	---------	------------	-------------

Table Stadium

<u>Stadium_ID</u>	St_name	City	Capacity	TeamID (FK)
-------------------	---------	------	----------	-------------

Table Teams

<u>Team_ID</u>	Team_name	YearOfFound
----------------	-----------	-------------

Table Matches

<u>Team_A</u>	<u>Team_B</u>	<u>M_date</u>
---------------	---------------	---------------

Table Players\_stats

<u>Team_A_ID</u>	<u>Team_B_ID</u>	<u>M_date</u>	<u>Player_ID</u> (FK)	MinPlayed	Goals	Shot_on	Shot_off	Penalty	Pass OK	Pass KO
Composite FK to table Matches										

The data to be loaded into this relational DB are in the script insert.sql and in a csv file (premier.csv).

The script Insert.sql populates the table Teams and the table Stadiums.

The csv file has the following format:

m\_date,player\_id,pl\_surname,pl\_name,team\_name,team\_id,opposition\_name,opposition\_id, min\_played,shot\_on,shot\_off,penalty,pass\_ok,pass\_ko

The meaning of each field is as follows:

Goal = number of goals scored during the match

m\_date = date of the match

player\_id = unique id of the player

pl\_surname = surname of the player

pl\_name = name of the player

team\_name = name of the player's team

team\_id = id of the player's team

opposition\_name = name of the opposing team (the one where the player is not playing in)

opposition\_id = id of the opposing team (the one where the player is not playing in)

min\_played = number of minutes played by the player in the match

shot\_on = number of shots on target by the player

shot\_off = number of shots off target by the player

penalty = number of goals scored by the player with a penalty

pass\_ok = number of passes completed by the player

pass\_ko = number of unsuccessful passes by the player

Note that the team names and team ids used in the sql script and in the csv file do match!

In the rest of the assignment, assume the following:

1. Player name + team name are unique in all the database, players do not change team during the season!
2. A team has only one home stadium

### **Deliverables**

- a. You are required to create the above relational model using Oracle or MySQL. Remember to implement the primary and foreign key constraints.
- b. You are then required to populate the data by loading the data from the sql script and from the csv into the relational model. Suggestion: you could load all the csv into a big denormalized table, and then move the data into the appropriate DB tables.
- c. Once the database has been loaded, you are required to perform an initial ETL process to move data into the data warehouse. You are required to define a sql script (that you might bundle into a procedure(s) or/and triggers) that automatically executes the process, creating the required stage tables and all the data matching procedures.

- d. Finally, you are required to perform a second ETL using the additional data contained in the etl2.csv files. This file contains a number of new records describing new statistics for new and old players (but no new matches). First you need to load the data into some staging table and then perform the second ETL reusing some of the stage tables and procedures used for the first ETL.
- e. Provide 2 sample queries that could be executed over the dimensional model to create reports about teams and players.
- f. Document all your work in a document containing the commented script and a detailed explanation of how your code works, your design choices, the instructions about how to run your code and some evidence/test that the script is working correctly. The scripts have to run in ORACLE or MYSQL! You will demo your work on week 13.

#### IMPORTANT:

The data warehouse has the following dimensional model:

##### Fact Table

Fact\_Stats(**date\_sk**,**player\_sk**,**team\_sk**,**opponent\_sk**,**stadium\_sk**,min\_played,goals,shot\_on,shot\_off,penalty,pass\_ok,pass\_ko)

##### Dimension Player

DimPlayer(**player\_sk**,player\_name,player\_surname,player\_age)

##### Dimension Team (used for the team and the opponent)

DimTeam(**team\_sk**, team\_name, year\_of\_foundation)

##### Dimension Time

DimTime(**date\_sk**,year,month,day)

##### Dimension Stadium

DimStadium(**stadium\_sk**,stadium\_name,stadium\_city,capacity)