# Capstone Project : FIFA Market Value & Wage

## Problem definition

Nowadays, Football is the most popular sport in the world. The financial in this industry is very massive and rapid growth. So, it's very interesting to analyze and consider for the subject. Among the cost in the industry, the most spending is related the market transfer value and wage for the football player. As a result, most of the clubs in the world would need to reduce this cost and spend only on the proper contract. Moreover, it's become serious when the FIFA announce the 'Financial Fair Play' rules to limit the spending on buying or hiring.

We craft this project to find the insight of average Player's Market value and Wage base on each dimension such as the Football League or Player Position to help the Management level on making decision for the contract agreement.

## Data source

FIFA 21 complete player dataset :

https://www.kaggle.com/datasets/stefanoleone992/fifa-21-complete-player-dataset

## Raw data

players_21.csv :



** The raw data file is stored in the AWS S3

## Data model

Datalake :



**Leagues** : Collect the football league information.

**Clubs** : Collect the football club information. Appear the relationship with Leagues because each club must be in a league.

**Nationalities** : Collect the nationality which belong to existing football player from source data.

**Positions** : Collect the player position, played for the club.

**Players** : Collect the player personal information including transfer value and wage. Appear the relationship with Positions, Nationalities and Clubs to identify the dimensions for this player.

\*\* Datalake tables are stored in csv format in AWS S3 with partitioning on **"date_oprt"**

Amazon S3 > Buckets > jaochin-dataset-fifa > cleaned/ > clubs/ > date_oprt=2022-12-17/

## date_oprt=2022-12-17/

**Objects** | Properties

**Objects** (1)

Objects are the fundamental entities stored in Amazon S3. You can use Amazon S3 inventory ↗ to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions. Learn more ↗

Copy S3 URI | Copy URL | Download | Open ↗ | Delete | Actions ▼ | Create folder | Upload

Find objects by prefix

| Name | Type | Last modified | Size |
|------|------|---------------|------|
| part-00000-206c497d-0ec2-402e-bee6-958f0426fbdb.c000.csv | csv | December 18, 2022, 00:33:04 (UTC+07:00) | 14.2 KB |

part-00000-206c497d-0ec2-402e-bee6-958f0426fbdb.c000.csv ✕

Users > JC > Downloads > part-00000-206c497d-0ec2-402e-bee6-958f0426fbdb.c000.csv

```
1    club_id,club_name,league_id
2    1,1. FC Heidenheim 1846,22
3    2,1. FC Kaiserslautern,23
4    3,1. FC Köln,21
5    4,1. FC Magdeburg,23
```

## Datawarehouse :

| player_value_wage | |
|---|---|
| player_id | bigint |
| player_name | text |
| player_age | int |
| player_overall | int |
| player_value | decimal |
| player_wage | decimal |
| position_name | text |
| club_name | text |
| nationality_name | text |
| league_name | text |
| date_oprt | date |

dbdiagram.io

**Player_value_wage** : Collect the information of player including necessary information such as position, club. The purpose of this datawarehouse table is to serve the OLAP processing, Tableau for example.

# DS525 – Capstone Project - Chin Lertvipada - 64199130039

** Datawarehouse table is stored in AWS Redshift with **"date_oprt"** as execution date

## Project workflow

1. The AWS Cloud environments setup
   a. S3 to store the *"raw data"* and *"cleaned data"* (Datalake)
   b. Redshift to store the *"OLAP data"* (Datawarehouse)
   c. AWS Credentials for application to access AWS
2. The source data (raw csv file) will be stored in the AWS S3
   a. The source data will be loaded into S3 with other process or team
   b. S3 repository : "jaochin-dataset-fifa/landing/"
3. The Datalake process will load the *"raw data"* and produce the *"cleaned data"* in AWS S3
   a. Data transformation & cleansing by PySpark
   b. Produce 5 tables : Clubs, Leagues, Positions, Nationalities, Players
   c. Output tables are in csv format
   d. Each table (csv file) is partitioned by *"date_oprt"* (execution date)
   e. Example S3 repository : "jaochin-dataset-fifa/cleaned/clubs/date_oprt=2022-12-17/"
4. The Datawarehouse process will load the *"cleaned data"* and produce the *"OLAP data"* in AWS Redshift
   a. Data transformation & load by Python
   b. Monthly schedule at $1^{st}$ of each month
   c. Load *"cleaned data"* from S3 into Redshift
   d. Load *"cleaned data"* with filtering *"date_oprt"* = execution date (current date)
   e. Merge and transform *"cleaned data"* to produce *"OLAP data"*
   f. Produce datawarehouse table : Player_value_wage
   g. The datawaehouse table is partitioned by *"date_oprt"* (execution date)
5. The dashboard for data visualization using Tableau Desktop
   a. Connect Tableau Desktop to Redshift with information in step **1.b**
   b. Select the database and data which collect the necessary information
   c. Create the dashboard to answer the problem or question definition
   d. Publish the dashboard to the Tableau Public for online access
6. The workflow orchestration using Airflow
   a. The step **1 - 4** should be fully automated and controlled by Airflow
   b. However, due to configuration in many steps are very complicate and consume a lot of time. So, for the moment we decide to use the Airflow to only automate and control for step **4** Datawarehouse process.
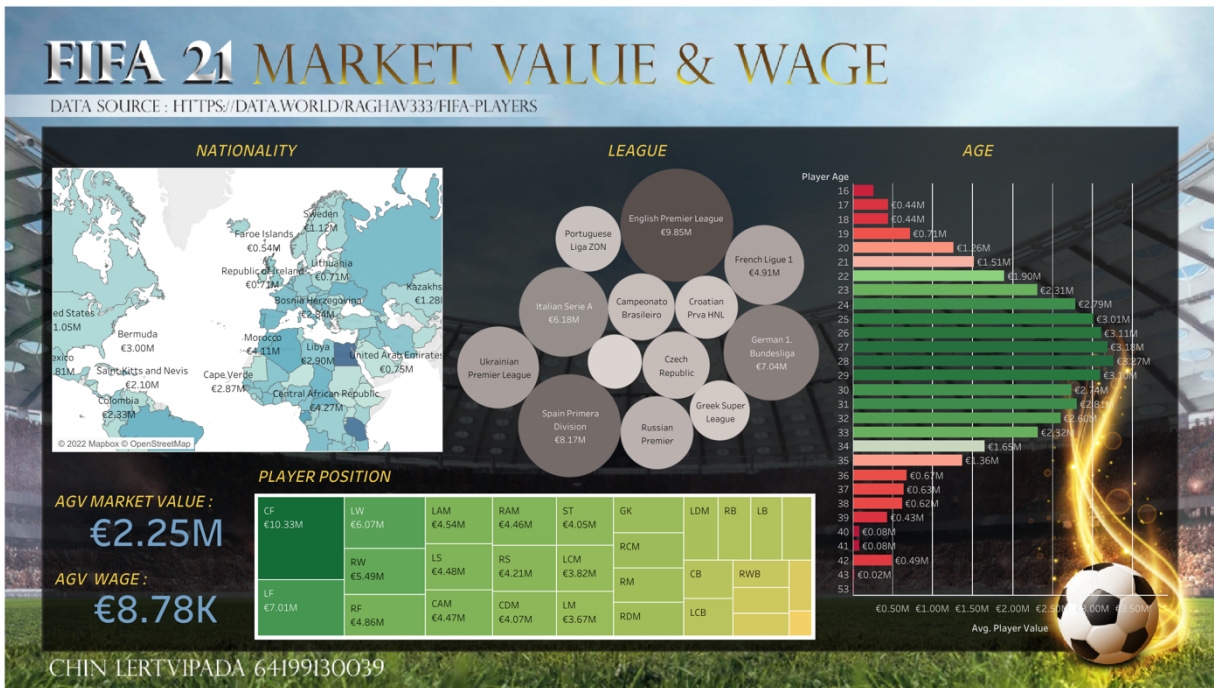
# Project implementation instruction

For the full implementation instruction (step-by-step), please find the information here :

https://github.com/chin-lertvipada/swu-ds525/tree/main/chin-capstone#readme

# Dashboard

Please find the dashboard online here :

https://public.tableau.com/app/profile/chin.lertvipada/viz/Capstone_csv/FootballMarketValue



The dashboard visualizes "Average Market value & Wage" for football player in many dimensions to answer the questions and solve the problems

- **Average Market value & Wage** : show the average value & wage among all football player
- **Nationality** : show the average value & wage for player from their nationality
- **League** : show the average value & wage by the football league
- **Age** : show the average value & wage by span of football player age
- **Player Position** : show the average value & wage by the position that they play for a club

## Summary

From the raw data, we create the process and workflow to extract, transform and load (ETL) the data to produce the data standardization for utilization of the data.

We make the flow thru the Datalake process and the Datawarehouse process to build to "OLAP data" which we believe that it will solve the problems and questions.

After that we create the Dashboard to visualize the data to help Management level or any stakeholder for the decision-making process.