# Machine Learning Project Summary: Predicting Football Players' Salaries

## Project Objectives

The project aimed to develop a machine learning model capable of predicting football players' salaries based on characteristics such as age, club, and position. The goal was to explore different modeling techniques to find the most effective method in terms of predictive accuracy.

## Methodology

Data Preparation: Data cleaning, converting salaries into a numerical format, encoding categorical variables. Machine Learning Models Used: Linear Regression, Ridge Regression, Random Forest. Model Evaluation: Using the mean squared error (MSE) to evaluate the performance of each model.

## Code for Data Preparation

```python
import pandas as pd
import numpy as np

# Load data
data = pd.read_csv('/path/to/data.csv')

# Clean data
data['Salary'] = data['Salary'].replace('[\€, p/a]', '', regex=True).astype(float)
data['Age'] = data['Age'].astype(int)

# Encode categorical variables
data = pd.get_dummies(data, columns=['Club', 'Position'])
```

## Code for Model Training and Evaluation

```python
from sklearn.model_selection import train_test_split

from sklearn.linear_model import LinearRegression, Ridge
```

```
from sklearn.ensemble import RandomForestRegressor
from sklearn.metrics import mean_squared_error

# Split data
X_train, X_test, y_train, y_test = train_test_split(data.drop('Salary', axis=1),
data['Salary'], test_size=0.2, random_state=0)

# Linear Regression
lr_model = LinearRegression().fit(X_train, y_train)
lr_pred = lr_model.predict(X_test)

# Ridge Regression
ridge_model = Ridge().fit(X_train, y_train)
ridge_pred = ridge_model.predict(X_test)

# Random Forest
rf_model = RandomForestRegressor().fit(X_train, y_train)
rf_pred = rf_model.predict(X_test)

# Evaluate models
lr_mse = mean_squared_error(y_test, lr_pred)
ridge_mse = mean_squared_error(y_test, ridge_pred)
rf_mse = mean_squared_error(y_test, rf_pred)
```

## Results

The various models tested produced high MSEs, suggesting difficulties in accurately predicting salaries using the available variables. The random forest offered a slight improvement over linear and Ridge regression, but the results were not substantially better.

## Visualizations

Salary histogram, Age vs salary scatter plot, Salary box plots by position, Feature importance graph (fictitious example)

## Lessons Learned

The project highlighted the importance of a deep understanding of the data and proper preparation. More complex models did not necessarily lead to better predictions, indicating the need to possibly explore other features or improve data preprocessing.

## Next Steps

Explore advanced techniques such as more extensive feature engineering and the use of deep learning models to improve prediction accuracy.