

Introduction

During the last couple of weeks we learned about how to read data, do exploratory data analysis (EDA) and prepare data for training and training a ML model. However, we did not specifically discuss the typical ML pipeline. In this lab, we will go through a typical ML model development process using a classification task as an example. Specifically, we will learn more about the machine learning pipeline, including examining and performing basic data cleaning. We then examine how to perform logistic regression, learn two basic metrics to evaluate this, and perform basic parameter tuning to demonstrate how it can be done. We will apply it to predict whether NBA rookies will play five years or more.

The lab assumes that you have completed the labs for week 2-4. If you havent yet, please do so before attempting this lab.

⚠ **Warning:** Starting this week, we will progressively provide less code, and would like you to use previous labs and what you know to perform the tasks. This will help you to become proficient at this.

The lab can be executed on either your own machine (with anaconda installation) or computer lab.

- Please refer canvas for instructions on installing anaconda python

Objective

- Continue to familiarise with Python and other ML packages
- Perform basic data preparation
- Practice performing logistic regression
- Learn how to perform basic parameter tuning

Dataset

In this lab, we will be using a dataset of NBA rookies, some of their stats and trying to predict whether they will still be playing after 5 years. You can download the data from Canvas.

First, ensure the data file `nbaRookies.csv` is located within the Jupyter workspace.

- If you are on the local machine copy the data file (`nbaRookies.csv`) to your current folder.

In this course we mostly use datasets that are collected by a third party. If you are interested in collecting your own data for your project, some useful information can be found at: [Introduction to Constructing Your Dataset](#)

Problem Formulation

The first step in developing a model is to formulate the problem in a way that we can apply machine learning. To reiterate, the task in the nbaRookies dataset is to predict whether NBA rookies will play five years or more, using some attributes of rookies.

◆ Observe the data and see if there is a pattern that would allow us to predict whether NBA rookies will play five years or more using the attributes given? You can use the observations from the EDA for this.

❖ What category does the task belong to?

✓ **Task category:**

- supervised, univariate/multivariate regression
- We should use the insights gained from observing the data (EDA) in selecting the performance measure. e.g. are there outliers in target?

- ✓ Data Pre-processing

We will first study how to perform some basic data pre-processing. First import pandas, sklearn, numpy and matplotlib.pyplot. You may also want to import seaborn for drawing more beautiful graphs.

Load dataset

We want to load the dataset 'nbaRookies.csv' into a Pandas dataframe (call it nbaDf). Remember to check if your dataframe was loaded correctly by print out the first few records or output some summary information about the dataset.

Start a new Jupyter notebook session. Load the dataset `nbaRookies.csv` in `nbaDf`.

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import numpy as np
4
5 ## TODO
6 nbaDf = pd.read_csv(r"Lab/nbaRookies.csv")##, delimiter="\s+")
7 print(nbaDf)
8
```

	Name	GP	MIN	PTS	FGM	FGA	FG%	3P	Made	3PA	3P%	...	\
0	Brandon Ingram	36	27.4	7.4	2.6	7.6	34.7	0.5	2.1	25.0	...		
1	Andrew Harrison	35	26.9	7.2	2.0	6.7	29.6	0.7	2.8	23.5	...		
2	JaKarr Sampson	74	15.3	5.2	2.0	4.7	42.2	0.4	1.7	24.4	...		
3	Malik Sealy	58	11.6	5.7	2.3	5.5	42.6	0.1	0.5	22.6	...		
4	Matt Geiger	48	11.5	4.5	1.6	3.0	52.4	0.0	0.1	0.0	...		
...		
1335	Chris Smith	80	15.8	4.3	1.6	3.6	43.3	0.0	0.2	14.3	...		
1336	Brent Price	68	12.6	3.9	1.5	4.1	35.8	0.1	0.7	16.7	...		
1337	Marlon Maxey	43	12.1	5.4	2.2	3.9	55.0	0.0	0.0	0.0	...		
1338	Litterial Green	52	12.0	4.5	1.7	3.8	43.9	0.0	0.2	10.0	...		
1339	Jon Barry	47	11.7	4.4	1.6	4.4	36.9	0.4	1.3	33.3	...		
	FTA	FT%	OREB	DREB	REB	AST	STL	BLK	TOV	TARGET_5Yrs			
0	2.3	69.9	0.7	3.4	4.1	1.9	0.4	0.4	1.3	0.0			
1	3.4	76.5	0.5	2.0	2.4	3.7	1.1	0.5	1.6	0.0			

```
2      1.3 67.0 0.5 1.7 2.2 1.0 0.5 0.3 1.0      0.0
3      1.3 68.9 1.0 0.9 1.9 0.8 0.6 0.1 1.0      1.0
4      1.9 67.4 1.0 1.5 2.5 0.3 0.3 0.4 0.8      1.0
...    ...  ...  ...  ...  ...  ...  ...  ...  ...
1335   1.5 79.2 0.4 0.8 1.2 2.5 0.6 0.2 0.8      0.0
1336   1.0 79.4 0.4 1.1 1.5 2.3 0.8 0.0 1.3      1.0
1337   1.6 64.3 1.5 2.3 3.8 0.3 0.3 0.4 0.9      0.0
1338   1.8 62.5 0.2 0.4 0.7 2.2 0.4 0.1 0.8      1.0
1339   1.0 67.3 0.2 0.7 0.9 1.4 0.7 0.1 0.9      1.0
```

[1340 rows x 21 columns]

▼ Data pre-processing

Let's plot a series of histogram to understand the distribution of the data more. Is there anything that captures your interest?

```
1 nbaDf.head()
```

	Name	GP	MIN	PTS	FGM	FGA	FG%	3P Made	3PA	3P%	...	FTA	FT%	OREB	DREB	REB	AST	STL	BLK	TOV	TARGET_5Yrs
0	Brandon Ingram	36	27.4	7.4	2.6	7.6	34.7	0.5	2.1	25.0	...	2.3	69.9	0.7	3.4	4.1	1.9	0.4	0.4	1.3	0.0
1	Andrew Harrison	35	26.9	7.2	2.0	6.7	29.6	0.7	2.8	23.5	...	3.4	76.5	0.5	2.0	2.4	3.7	1.1	0.5	1.6	0.0
2	JaKarr Sampson	74	15.3	5.2	2.0	4.7	42.2	0.4	1.7	24.4	...	1.3	67.0	0.5	1.7	2.2	1.0	0.5	0.3	1.0	0.0
3	Malik Sealy	58	11.6	5.7	2.3	5.5	42.6	0.1	0.5	22.6	...	1.3	68.9	1.0	0.9	1.9	0.8	0.6	0.1	1.0	1.0
4	Matt Geiger	48	11.5	4.5	1.6	3.0	52.4	0.0	0.1	0.0	...	1.9	67.4	1.0	1.5	2.5	0.3	0.3	0.4	0.8	1.0

5 rows x 21 columns

The target column is **TARGET_5Yrs** and all the other columns are attributes.

```
1 nbaDf.shape
```

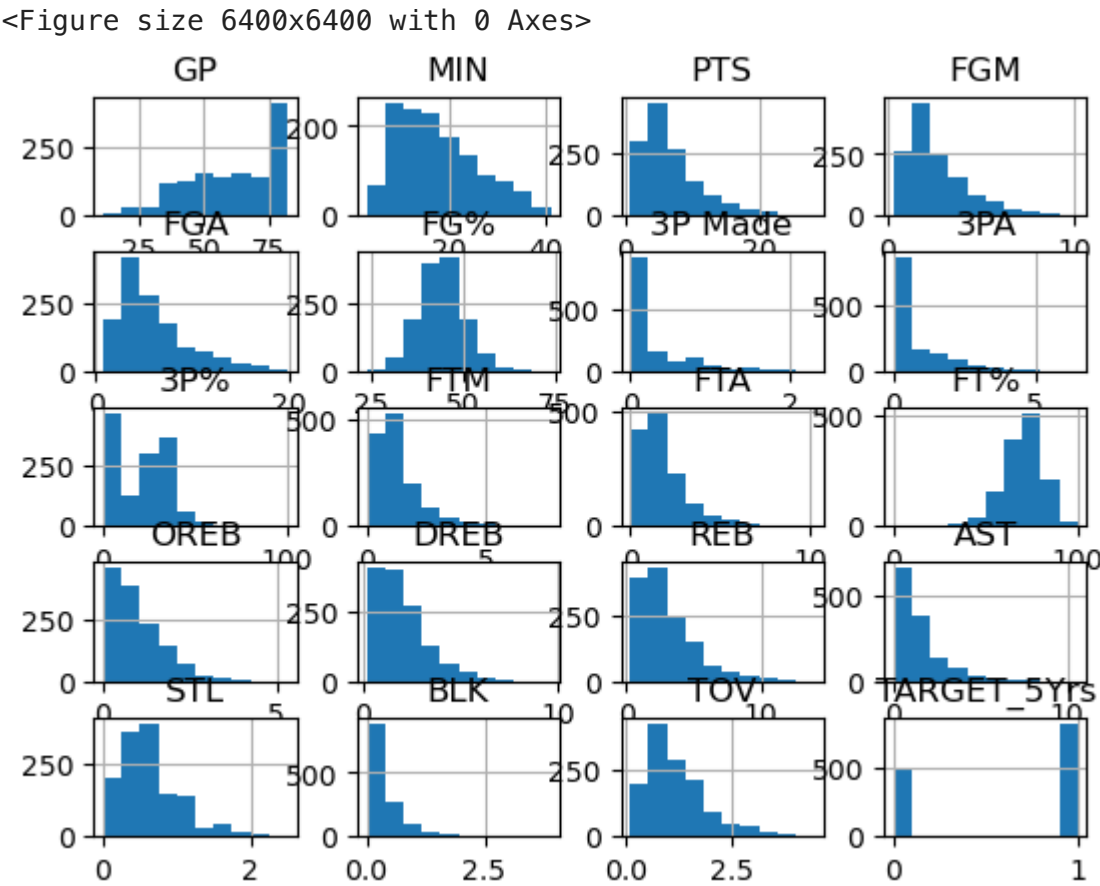
(1340, 21)

```
1 nbaDf.info()
```

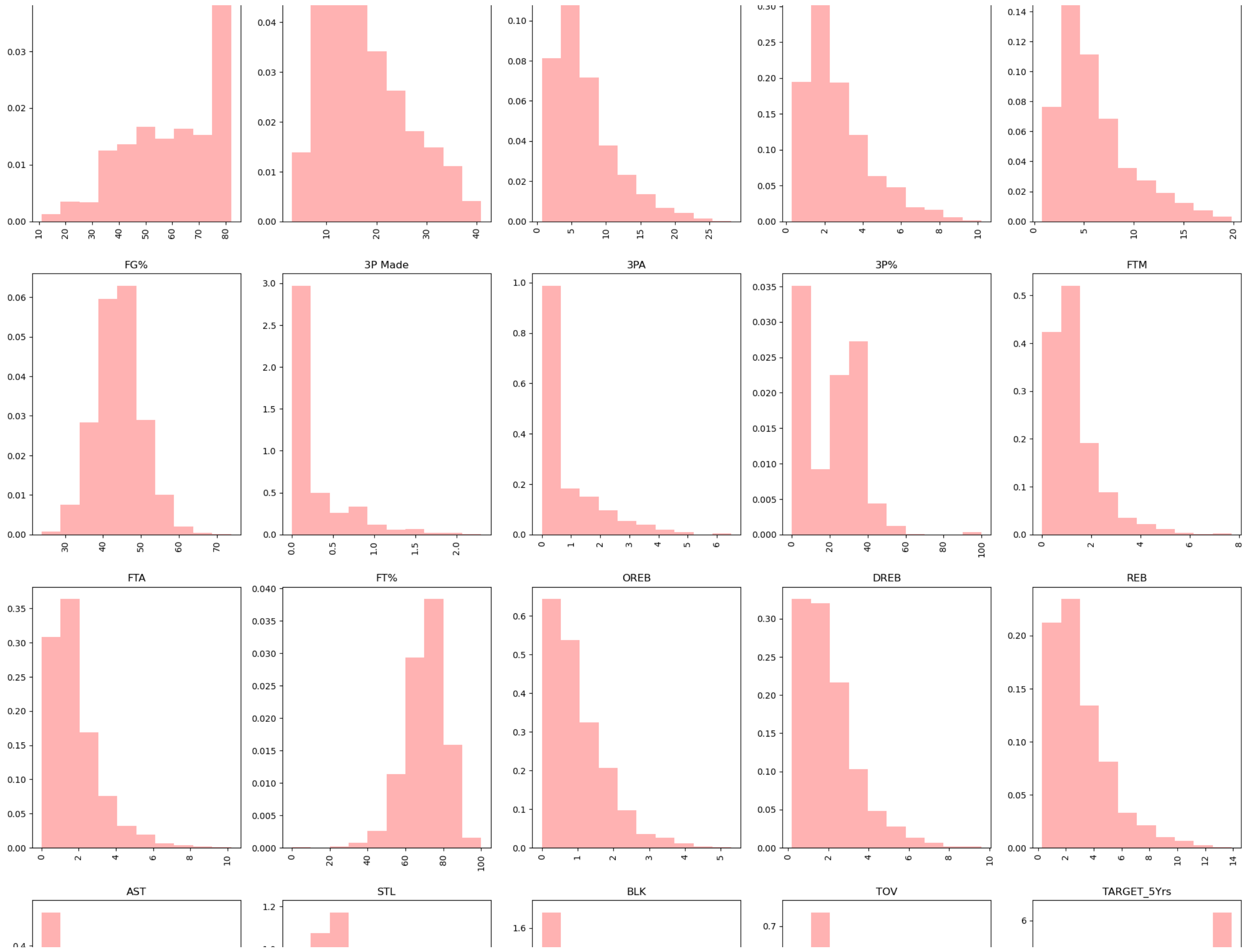
```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1340 entries, 0 to 1339
Data columns (total 21 columns):
#   Column          Non-Null Count  Dtype
---  ---
0   Name            1340 non-null   object
1   GP              1340 non-null   int64
2   MIN             1340 non-null   float64
3   PTS             1340 non-null   float64
4   FGM             1340 non-null   float64
5   FGA             1340 non-null   float64
6   FG%             1340 non-null   float64
7   3P Made         1340 non-null   float64
8   3PA             1340 non-null   float64
9   3P%             1329 non-null   float64
10  FTM             1340 non-null   float64
11  FTA             1340 non-null   float64
12  FT%             1340 non-null   float64
13  OREB            1340 non-null   float64
14  DREB            1340 non-null   float64
15  REB             1340 non-null   float64
16  AST             1340 non-null   float64
17  STL             1340 non-null   float64
18  BLK             1340 non-null   float64
19  TOV             1340 non-null   float64
20  TARGET_5Yrs     1340 non-null   float64
dtypes: float64(19), int64(1), object(1)
memory usage: 220.0+ KB
```

▼ Let's plot a series of histogram to understand the distribution of the data more. Is there anything that captures your interest?

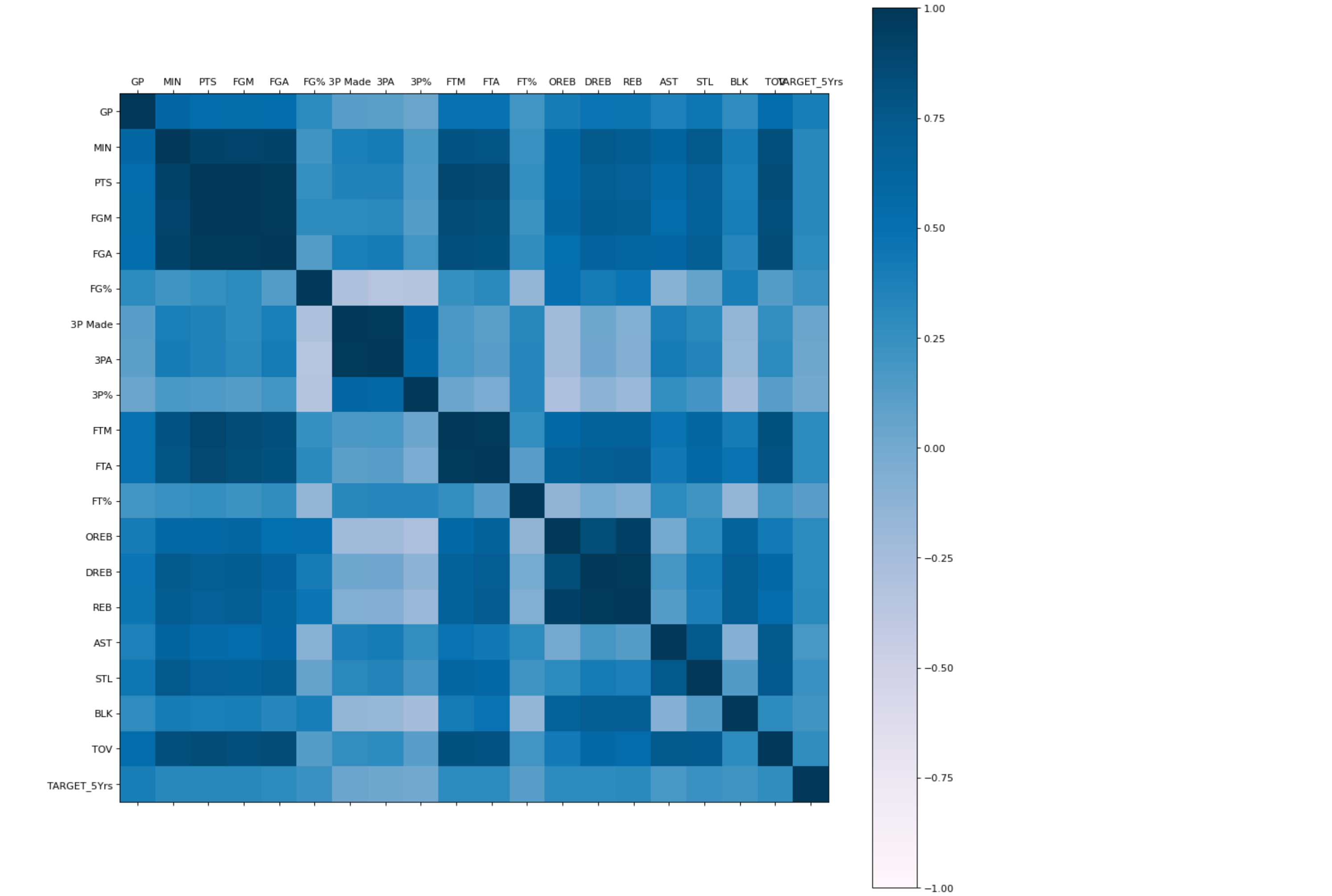
```
1 plt.figure(figsize=(40, 40), dpi=160)
2 nbaDf.hist()
3 plt.show()
```



```
1 plt.figure(figsize=(25,25))
2 for i, col in enumerate(nbaDf.columns[1:]):
3     plt.subplot(4,5,i+1)
4     plt.hist(nbaDf[col], alpha=0.3, color='r', density=True)
5     plt.title(col)
6     plt.xticks(rotation='vertical')
```



```
1 # Select only numeric columns
2 numeric_nbaDf = nbaDf.select_dtypes(include=[np.number])
3
4 # Compute the correlation matrix
5 correlation = numeric_nbaDf.corr()
6
7 # Create a new figure
8 fig = plt.figure(figsize=(16, 16), dpi=80)
9
10 # Add a subplot
11 ax = fig.add_subplot(111)
12
13 # Display the correlation matrix
14 cax = ax.matshow(correlation, vmin=-1, vmax=1, cmap=plt.cm.PuBu)
15
16 # Add a colorbar
17 fig.colorbar(cax)
18
19 # Set the x and y ticks
20 ticks = np.arange(0, len(numeric_nbaDf.columns), 1)
21 ax.set_xticks(ticks)
22 ax.set_yticks(ticks)
23
24 # Set the x and y tick labels
25 ax.set_xticklabels(numeric_nbaDf.columns)
26 ax.set_yticklabels(numeric_nbaDf.columns)
27
28 # Display the plot
29 plt.show()
```



1 nbaDf.describe()

	GP	MIN	PTS	FGM	FGA	FG%	3P Made	3PA	3P%	FTM	FTA	FT%	OREB	DREB	REB	AST	STL	BLK	TOT	TARGET_5Yrs
count	1340.000000	1340.000000	1340.000000	1340.000000	1340.000000	1340.000000	1340.000000	1340.000000	1329.000000	1340.000000	1340.000000	1340.000000	1340.000000	1340.000000	1340.000000	1340.000000	1340.000000	1340.000000	1340.000000	1340.000000
mean	60.414179	17.624627	6.801493	2.629104	5.885299	44.169403	0.247612	0.779179	19.308126	1.297687	1.821940	70.300299	1.009403	2.025746	3.035149	2.237211	0.760164	1.610164	1.610164	1.610164
std	17.433992	8.307964	4.357545	1.683555	3.593488	6.137679	0.383688	1.061847	16.022916	0.987246	1.322984	10.578479	0.777119	1.360008	2.035149	1.360008	0.760164	1.610164	1.610164	1.610164
min	11.000000	3.100000	0.700000	0.300000	0.800000	23.800000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.200000	0.300000	0.700000	0.700000	0.700000	0.700000	0.700000
25%	47.000000	10.875000	3.700000	1.400000	3.300000	40.200000	0.000000	0.000000	0.000000	0.600000	0.900000	64.700000	0.400000	1.000000	1.500000	1.000000	0.700000	1.000000	1.000000	1.000000
50%	63.000000	16.100000	5.550000	2.100000	4.800000	44.100000	0.100000	0.300000	22.400000	1.000000	1.500000	71.250000	0.800000	1.700000	2.500000	1.700000	1.000000	1.700000	1.700000	1.700000
75%	77.000000	22.900000	8.800000	3.400000	7.500000	47.900000	0.400000	1.200000	32.500000	1.600000	2.300000	77.600000	1.400000	2.600000	4.000000	2.600000	1.400000	2.600000	2.600000	2.600000
max	82.000000	40.900000	28.200000	10.200000	19.800000	73.700000	2.300000	6.500000	100.000000	7.700000	10.200000	100.000000	5.300000	9.600000	13.900000	9.600000	5.300000	9.600000	9.600000	9.600000

What observations did you make?

✓ **Observations:**

- We can see that the 3P% column has only 1329 items while other columns all have 1340 items.

If there are missing values in the dataset, they are generally represented as NaN Values.

If we tried to run this with a classifier, we will find it will complaint about NaN values. Let's examine them:

```
1 import pandas as pd
2 pd.isna(nbaDf)
```

	Name	GP	MIN	PTS	FGM	FGA	FG%	3P Made	3PA	3P%	...	FTA	FT%	OREB	DREB	REB	AST	STL	BLK	TOV	TARGET_5Yrs
0	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False	False	False	False	False
1	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False	False	False	False	False
2	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False	False	False	False	False
3	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False	False	False	False	False
4	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False	False	False	False	False
...
1335	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False	False	False	False	False
1336	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False	False	False	False	False
1337	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False	False	False	False	False
1338	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False	False	False	False	False
1339	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False	False	False	False	False

1340 rows x 21 columns

That outputs the whole dataframe and entries with True means the value is NaN or None. Given the size of the dataframe, it is hard to visualise it. Please check up the reference for isna() at <https://pandas.pydata.org/pandas-docs/stable/generated/pandas.DataFrame.isna.html> (Links to an external site.).

Knowing that the function isna() produces a dataframe, can you find a way to summarise how many rows that contain missing data? What are the column(s) that contain missing data, and how many rows? Next, slice the nbaDf dataframe to examine the rows that have missing data.

There are several ways to deal with this, but in this case, we can set the missing data to zeros. Please use the built-in function fillna() of pandas to do this.

This essentially fills all NaN entries with 0 (remember to check the documentation for details of the method). There is another useful function to deal with NaN and missing values called interpolate, that tries to infer values – again check the documentation for details. Another option is to drop the row/instance if it appears the instance might be erroneous or there is no good way to fill or infer.

```
1 pd.isna(nbaDf).sum()
```

```
Name      0
GP         0
MIN        0
PTS        0
FGM        0
FGA        0
FG%        0
3P Made    0
3PA        0
3P%       11
FTM        0
FTA        0
FT%        0
OREB       0
DREB       0
REB        0
AST        0
STL        0
BLK        0
TOV        0
TARGET_5Yrs 0
dtype: int64
```

The 3P% column has 11 NaN values. We can find which instances/rows this corresponds to:

```
1 nbaDf[pd.isna(nbaDf).any(axis=1)]
```

	Name	GP	MIN	PTS	FGM	FGA	FG%	3P Made	3PA	3P%	...	FTA	FT%	OREB	DREB	REB	AST	STL	BLK	TOV	TARGET_5Yrs
338	Ken Johnson	64	12.7	4.1	1.8	3.3	52.8	0.0	0.0	NaN	...	1.3	43.5	1.4	2.4	3.8	0.3	0.2	0.3	0.9	0.0
339	Ken Johnson	64	12.7	4.1	1.8	3.3	52.8	0.0	0.0	NaN	...	1.3	43.5	1.4	2.4	3.8	0.3	0.2	0.3	0.9	0.0
340	Pete Williams	53	10.8	2.8	1.3	2.1	60.4	0.0	0.0	NaN	...	0.8	42.5	0.9	1.9	2.8	0.3	0.4	0.4	0.4	0.0
358	Melvin Turpin	79	24.7	10.6	4.6	9.0	51.1	0.0	0.0	NaN	...	1.8	78.4	2.0	3.8	5.7	0.5	0.5	1.1	1.5	1.0
386	Jim Petersen	60	11.9	3.2	1.2	2.4	48.6	0.0	0.0	NaN	...	1.1	75.8	0.7	1.7	2.5	0.5	0.2	0.5	1.2	1.0
397	Tom Scheffler	39	6.9	1.3	0.5	1.3	41.2	0.0	0.0	NaN	...	0.5	50.0	0.5	1.5	1.9	0.3	0.2	0.3	0.4	0.0
507	Sam Williams	59	18.2	6.1	2.6	4.7	55.6	0.0	0.0	NaN	...	1.5	55.1	1.5	3.7	5.2	0.6	0.8	1.3	1.1	0.0
509	Kurt Nimphius	63	17.2	5.3	2.2	4.7	46.1	0.0	0.0	NaN	...	1.7	58.3	1.5	3.2	4.7	1.0	0.3	1.3	0.9	1.0
510	Pete Verhoeven	71	17.0	4.9	2.1	4.2	50.3	0.0	0.0	NaN	...	1.0	70.8	1.5	2.1	3.6	0.7	0.6	0.3	0.8	1.0
521	Jim Smith	72	11.9	2.9	1.2	2.3	50.9	0.0	0.0	NaN	...	1.2	45.9	1.0	1.5	2.5	0.6	0.3	0.7	0.7	0.0
559	Jeff Wilkins	56	18.9	4.7	2.1	4.6	45.0	0.0	0.0	NaN	...	0.7	67.5	1.1	3.8	4.9	0.7	0.6	0.8	1.1	1.0

11 rows x 21 columns

💡 What are the possible actions we can take?

✔ Actions:

- We can remove the above rows from the dataset. This will lead to loss of some information as we will lose the other attribute information in those rows.
- We can replace the missing values with zero (or the mean of that column with missing values). Need to see if this is reasonable for a given attribute, using nbaDf.fillna(0) or nbaDf.fillna()
- We can use another feature(s) to predict the missing values and use that.

For this problem we can observe that the 3P% and the FTM (or MIN) has a very strong correlation (See EDA results that appear before). Therefore we can use the value of the FTM to replace the missing values of 3P%. Generally we might have to train a ML model to predict the missing attributes (x: FTM , y: 3P%). However for this problem we can even directly replace the missing mode values without building a model.

```
1 nbaDf.loc[pd.isna(nbaDf['3P%']), '3P%'] = nbaDf.loc[pd.isna(nbaDf['3P%']), 'FTM']
```

The `loc` function is used to access a group of rows and columns by label(s) or a boolean array. In this case, it's being used twice: once to identify the rows where '3P%' is NaN (Not a Number), and once to replace those NaN values.

The expression `pd.isna(nbaDf['3P%'])` returns a boolean Series where each element is True if the corresponding value in the '3P%' column is NaN, and False otherwise.

The code `nbaDf.loc[pd.isna(nbaDf['3P%']), '3P%']` then uses this boolean Series to select only the rows in '3P%' column of `nbaDf` where '3P%' is NaN.

The entire line `nbaDf.loc[pd.isna(nbaDf['3P%']), '3P%'] = nbaDf.loc[pd.isna(nbaDf['3P%']), 'FTM']` replaces the NaN values in the '3P%' column with the corresponding values from the 'FTM' column.

This might be done, for example, if you're preparing your data for a machine learning algorithm that cannot handle NaN values, and you've decided that the 'FTM' values are a good substitute for missing '3P%' values.

Check the data again after fill-in NaN values

1 pd.isna(nbaDf).sum()		
Name	0	
GP	0	
MIN	0	
PTS	0	
FGM	0	
FGA	0	
FG%	0	
3P Made	0	
3PA	0	
3P%	0	
FTM	0	
FTA	0	
FT%	0	
OREB	0	
DREB	0	
REB	0	
AST	0	
STL	0	
BLK	0	
TOV	0	
TARGET_5Yrs	0	
dtype:	int64	

> Setting up training and testing data

The final task in this section is to set up the feature/attribute data and the column we are predicting 'TARGET_5Yrs'. We have done this in the previous lab, please do that now.

Similar to last week and we discuss in lectures about evaluation, we will divide our data into a number of testing datasets.

What we want to do is to use the training (data)set to construct the model, then use the validation set to tune the parameters of the model. Then once the parameters + model are tuned, we evaluate it on the testing set. This reduces the risk that we overfit if we use the testing set to tune the parameters (something we will talk about in lectures).

Scikit-learn doesn't have a function to split the data into the three sets. Instead, we can call it twice! First, lets split into training and testing dataset, as per last week (remember to import the relevant packages):

[] ↪ 10 cells hidden

> Baseline model

We need to select a baseline mode to do this task. I am going to select `regularised polynomial logistic regression` for this example.

There are better models than this, however we only know logistic regression technique that can be used for this problem at the moment, therefore out choices are limited and the decision is simple. If we had other options, we need to use our knowledge on those techniques and the EDA to select the best base model.

The polynomial model is justified because in the EDA we can see that a non-linear decision boundary can separate the classes. regularisation is justified because we have correlated attributes and in EDA we also had some features where a linear decision boundary looked appropriate.

[] ↪ 11 cells hidden

> Apply regularisation

When applying regularisation we need to select the lambda value. For this we can use

- 1. Grid search
- 2. Random search

We will do grid search in this example. In grid search we establish a set of lambda values in a frid. Selecting the range of lambda values is a process mostly done with trial and error.