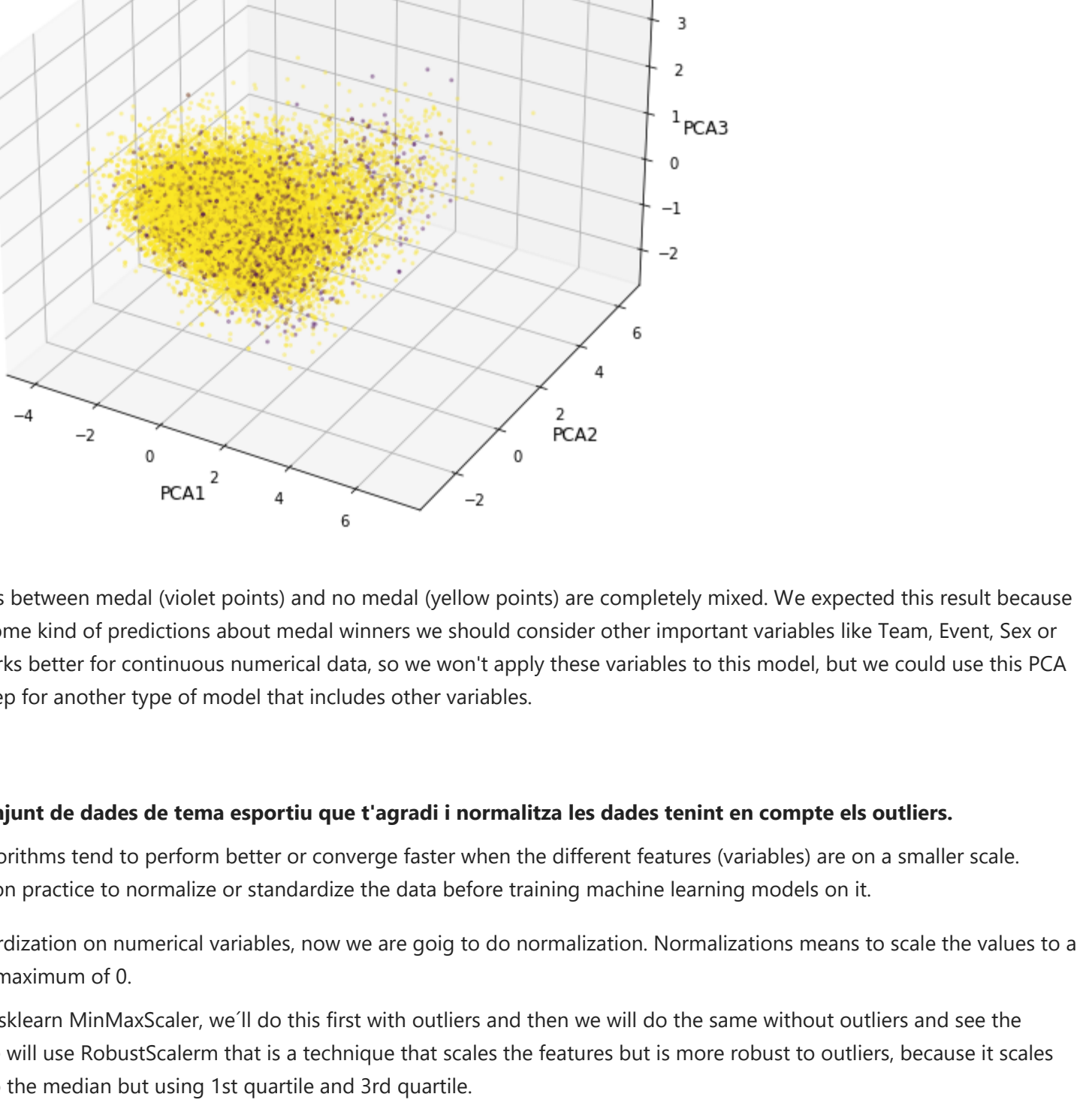






Principal Component Analysis of Athletics athletes in the Olympic Games (3 components)



We can see the points between medal (violet points) and no medal (yellow points) are completely mixed. We expected this result because to be able to make some kind of predictions about medal winners we should consider other important variables like Team, Event, Sex or even Games. PCA works better for continuous numerical data, so we won't apply these variables to this model, but we could use this PCA approach as a first step for another type of model that includes other variables.

- Exercici 3

Continua amb el conjunt de dades de tema esportiu que t'arregadi i normalitza les dades tenint en compte els outliers.

Machine learning algorithms tend to perform better or converge faster when the different features (variables) are on a smaller scale. Therefore it is common practice to normalize or standardize the data before training machine learning models on it.

Before we did standardization on numerical variables, now we are going to do normalization. Normalizations means to scale the values to a minimum of 1 and a maximum of 10.

We will do this using sklearn MinMaxScaler, we'll do this first with outliers and then we will do the same without outliers and see the differences. Lastly, we will use RobustScaler that is a technique that scales the features but is more robust to outliers, because it scales data not according to the median but using 1st quartile and 3rd quartile.

- #### Normalization without outliers

```
In [97]: from sklearn.preprocessing import MinMaxScaler  
# use the athletics_df that has the original values  
athletics_df.sample(3)
```

Out[97]:

	ID	Name	Sex	Age	Height	Weight	Year	Medal.0	Medal	Bronze	Medal	Gold	Medal	Silver	Event	Athletics	Men's 1,500	metres	Wal	Event	Athletics	Men's 1,500	metres	Wal
26509	132731	Fouad Wadid Yazgi	M	0	19.00	179.77	73.84	1952	1		0		0		0									
11068	48876	Moa Elin Mårtensson	F	1	22.00	172.00	60.00	2012	1		0		0		0									
19347	95497	Andrey Plotnikov	M	0	28.00	186.00	77.00	1996	1		0		0		0									

```
In [98]: # Select the numerical columns we want to normalize:  
X = athletics_df[['Age', 'Height', 'Weight']]  
mms = MinMaxScaler()  
# apply transformation on selected columns:  
Xnorm = mms.fit_transform(X)
```

```
In [99]: # look for the min and max ranges for the normalized matrix:  
print(Xnorm.min(axis=0))  
print(Xnorm.max(axis=0))  
  
[0. 0. 0.]  
[1. 1. 1.]
```

```
In [100]: print(Xnorm.std(axis=0))  
  
[0.10815344 0.13233521 0.11008015]
```

Since all three columns are in the same scale now, we can compare the values between them. Since standard deviation is similar in the three columns, we can say that the data is similarly dispersed for the three variables.

We will convert this again to a dataframe:

```
In [128]: norm_df = pd.DataFrame(Xnorm, columns = ['Age', 'Height', 'Weight'])
```

- #### Normalization without outliers

Now we are going to do the same but first we will remove outliers, we'll use the IQR approach (IQR for Inter Quartile Range).

IQR = Quartile3 - Quartile1

To define an outlier we must set an Upper and Lower bound, using IQR we'll define this bounds and the remove all the values that don't fit this:

- upper bound = Quartile3 + 1.5 \* IQR
- lower bound = Quartile1 - 1.15 \* IQR

```
In [101]: # we'll define the columns  
cols = ['Age', 'Height', 'Weight']  
  
# define quartile 1  
Q1 = athletics_df[cols].quantile(0.25) # 25% of the observations in this quartile  
Q3 = athletics_df[cols].quantile(0.75) # 75% of the observations in this quartile  
# define inter quartile range  
IQR = Q3 - Q1  
# new dataframe with observations from outliers removed (removes whole row if finds any outlier)  
athletics_df_no_outliers = athletics_df[(athletics_df[cols] < (Q1 + 1.5 * IQR)) && (athletics_df[cols] > (Q3 - 1.5 * IQR)).any(axis=1)]
```

```
In [102]: print('We have removed ', len(athletics_df) - len(athletics_df_no_outliers), ' observations from our dataset.')  
print('After removing outliers we have ', len(athletics_df_no_outliers), ' rows in the dataframe.')
```

We have removed 2385 observations from our dataset.  
After removing outliers we have 36239 rows in the dataframe.

Before applying normalization let's check describe method:

```
In [103]: athletics_df[['Age', 'Height', 'Weight']].describe().round(2)
```

Out[103]:

	Age	Height	Weight
count	36239.00	36239.00	36239.00
mean	24.91	175.99	67.54
std	4.04	8.16	10.29
min	13.00	152.00	39.00
25%	22.00	170.00	60.00
50%	25.00	178.00	68.00
75%	27.00	180.00	73.84
max	37.00	200.00	95.00

If we compare this summary statistics with the ones before removing outliers we can say that most outliers were 'Greater' values, since removing them the means and specially the max values have decreased (minimum values too but not so much).

If we look at standard deviation before removing outliers:

```
In [105]: athletics_df[['Age', 'Height', 'Weight']].std().round(2)
```

Out[105]:

Age	4.33
Height	8.73
Weight	14.31
dtype:	float64

Looking at Standard deviations we notice that they also have decreased a little, specially in the case of 'Weight' attribute, which was the one with most outliers (we saw that when we draw the boxplots). This effect is caused because standard deviation is a measure of dispersion of the data, and removing outliers makes data less dispersed.

Now let's normalize the data:

```
In [106]: X_without_outliers = athletics_df_no_outliers[['Age', 'Height', 'Weight']]  
mms = MinMaxScaler()  
  
Xnorm_without_outliers = mms.fit_transform(X_without_outliers)
```

```
In [107]: Xnorm_without_outliers[:3]
```

```
Out[107]: array([[0.45833333, 0.45833333, 0.55357143],  
[0.41666667, 0.57850631, 0.62212123],  
[0.70833333, 0.72956667, 0.66071429]])
```

```
In [108]: # look for the min and max ranges for the normalized matrix:  
print(Xnorm_without_outliers.min(axis=0))  
print(Xnorm_without_outliers.max(axis=0))  
  
[0. 0. 0.]  
[1. 1. 1.]
```

As before, we see that now all of our values are in the scale between 0 and 1. Now we will convert the array to a dataframe:

```
In [130]: norm_df_no_outliers = pd.DataFrame(Xnorm_without_outliers, columns = ['Age', 'Height', 'Weight'])
```

- Scale using RobustScaler

Another approach to scale variables in the presence of outliers is 'robust data scaling', that ignores the outliers from the calculation of the mean and standard deviation, then uses the calculated values to scale the variable. The result then is not skewed by the outliers.

```
In [117]: from sklearn import preprocessing  
  
robust_scaler = preprocessing.RobustScaler()  
robust_df = robust_scaler.fit_transform(athletics_df[['Age', 'Height', 'Weight']])  
print('Number of observations in new dataframe is', len(robust_df), ', same as the original df:', len(athletics_df))
```

We see that RobustScaler has not removed any values, the way it works is that it preserves outliers but tries to not let them influence the scaling of the non outliers values.

```
In [122]: robust_df = pd.DataFrame(robust_df, columns = ['Age', 'Height', 'Weight'])  
robust_df.describe().round(2)
```

Out[122]:

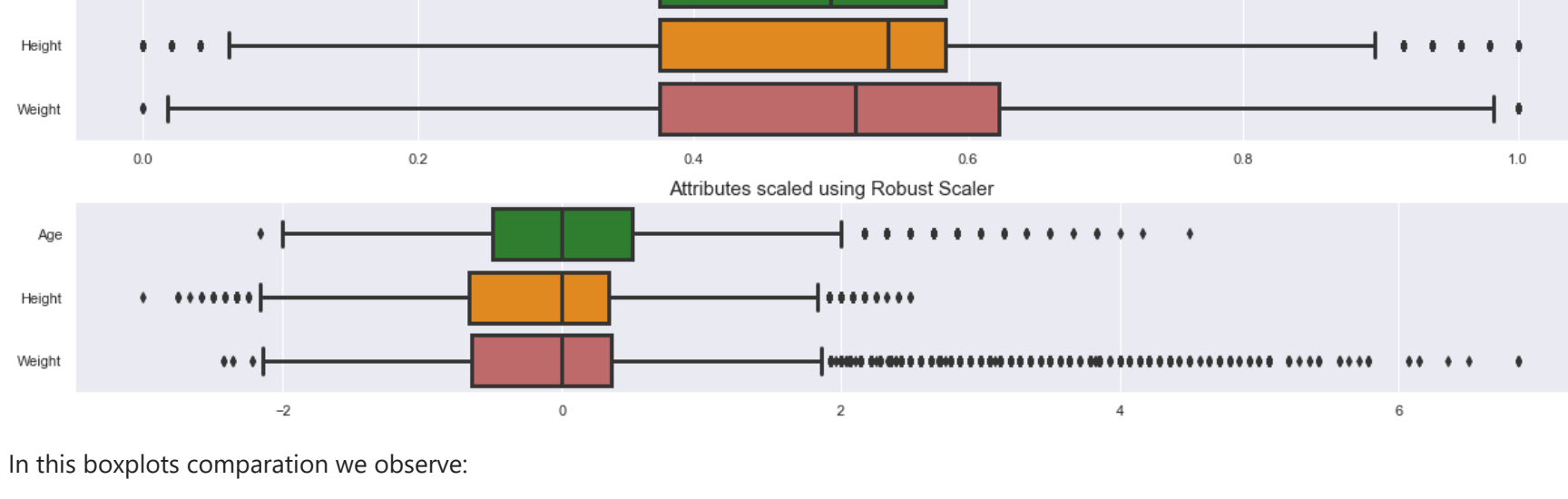
	Age	Height	Weight
count	38624.00	38624.00	38624.00
mean	0.03	-0.12	0.05
std	0.72	0.73	1.02
min	-2.17	-3.00	-2.43
25%	-0.50	-0.67	-0.64
50%	0.00	0.00	0.00
75%	0.50	0.33	0.36
max	4.50	2.50	6.86

We observe that RobustScaler works similar to standardization, the mean is centered around 0 (in the parameters, with centering: boolean: It is True by default), but in the case of Standard Deviation, it's not equal to 1 in all cases (only in the case of weight). It may be because making the scale more robust to outliers may have lowered std.

Let's plot this together and see if we could reach any conclusions. We have scaled the data in four different ways:

- standard\_df: we have used StandardScaler
- norm\_df: we have used MinMaxScaler
- norm\_df\_no\_outliers: first we have removed outliers and then used MinMaxScaler
- robust\_df: we have used RobustScaler

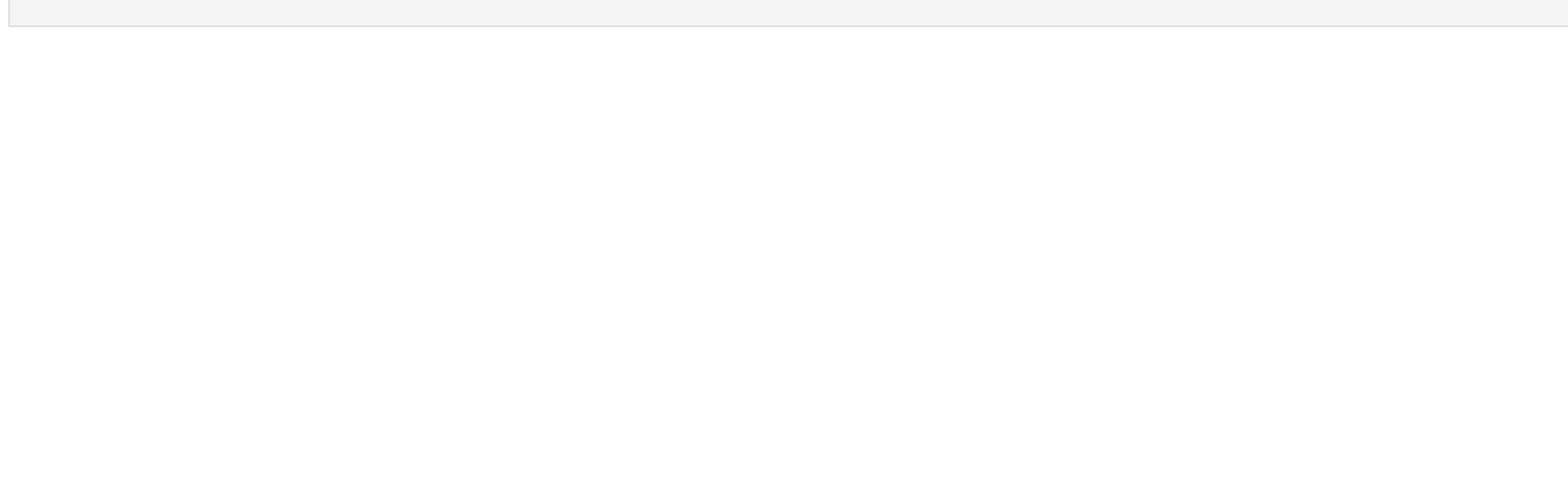
```
In [285]: # plotting the distributions in each case  
fig, axes = plt.subplots(5, 1, figsize=(17,14))  
fig.suptitle('Distribution of Attributes among Olympic Athletes after different Scaling Techniques', fontsize = 20)  
colors = ['#228b22', '#ff8c00', '#cd5c5c']
```



- In the original data, the distributions are far from each other because they are measured in different scales.
- Using StandardScaler Height and Weight look more similar to each other and Age resembles more like a normal distribution.
- MinMaxScaler including outliers makes the three distributions look different from each other, Height has the data more dispersed and Weight is the attribute that has more outliers and is clearly right-skewed.
- MinMaxScaler with outliers removed: now Height and Weight look more similar, the removing of the outliers has made the Weight distribution shift to the right. Age distribution is the one that looks closer to normal and is the most centered one.
- RobustScaler makes the data much more concentrated around the center than the other techniques, the three attributes look closer to each other but the outliers are still there.

Now let's do the boxplots to see if we see anything else:

```
In [284]: fig, axes = plt.subplots(5, 1, figsize=(17,14))  
fig.suptitle('Boxplots of Attributes among Olympic Athletes after different Scaling Techniques', fontsize = 20)  
colors = ['#228b22', '#ff8c00', '#cd5c5c']
```



In this boxplots comparison we observe:

- Using StandardScaler: all boxplots are centered around the same area (median is close to zero). Weight attribute has so many outliers that makes the boxplot more tight compared to the other two. In the three cases we see that outliers mostly come from higher values.
- Using MinMaxScaler: Even though all three attributes are in the same scale, they are centered around different values, we already saw that looking at the distributions, maybe this is the effect the outliers have on each type of data.
- MinMaxScaler after removing outliers: In this case we get a more symmetric distribution of the three attributes, and we observe that Weight boxplot is much more wider in this case than in the previous ones, being the attribute that has more outliers it is the one that changes the most after removing them.
- Use of RobustScaler: here we see that the three attributes are quite similar to each other, their boxplot's size is almost equal between them and their medians are centered around zero. We continue to see the outliers so no information is lost in this transformation.

Conclusions

After using all this techniques of scaling and treatment of outliers, we can reach some conclusions. Removing outliers may not be such a good idea in all cases because outliers may provide with some information about the data or tendencies that we would be missing if we delete them. In our dataset there are a lot of outliers but they do not seem to be measurement errors. Even though we are studying athletes of the same sport (Athletics), there are subcategories in this sport in which some athletes may have different attributes (for instance, athlete that compete in Hammer Throw clearly will not have the same values as the ones running the Marathon). The different types of events we have leads to different attributes in the athletes, that can be a possible explanation of the reason we have so many outliers.

RobustScaler seems like a good technique, it's a quick way to scale the different attributes to the same values but at the same time we keep the information that outliers provide us.

```
In [ ]:
```