

OLYMPIC GAMES STATS PROJECT

Data Analysis Project Proposal

28/01/2022

Leonardo Saviane

Why Sports

- I want to retrieve some insights on the best athletes that the world have ever known, to undestand whether some characteristic, such as height, weight, age and country, have some impact on achieving a medal in some sport.
- This can be useful for professional sport institutes who wants to push their athletes more towards some disciplines than another. Eg. If someone is tall maybe is more likely to perform better in sports like high jump or basket.



The Dataset and the cleaning

- The dataset used is SportsStats, which is a sports analysis firm partnering with local news and elite personal trainers to provide "interesting" insights to help their partners. Insights could be patterns/trends highlighting certain groups/events/countries, etc.
- Some variable such as "season" (winter and summer, "sex" (male or female) will be transformed in binary variables.
- "Medals" variable, so if the athletes has won a bronze/silver/ore medal, will be transformed in a binary variable as well. My purpose it's to estimate the likelyhood of winning a medal given the other factors.
- Some data are missing, such as height, weight and age for some athletes. Therefore, those
 value will be dropped to avoid biased estimate. Putting the average in the missing value would
 lead to possible biased estimate. Since the sample is large enough loosing more or less 100000
 observations is ok.
- In the following slide the ERD and an anticipation of the dataset, then the questions I want to answer ->

ENTITY RELATIONAL DIAGRAM

Athletes_events

- ID
- Name
- Sex
- Age
- Height
- Weight
- Team
- NOC
- Games
- Year
- Season
- City
- Sport
- Event
- Medal

NOC_regions

- NOC
- Region
- notes

271116



N Summer Sports

35

66

230

N Winter Sports

17

Short description of the table 'Athletes_events':

134731

SECTION 1: QUESTIONS TO ANSWER

- Exists a correlation between medals won and country?
- Are Height and Age more relevant for some sports than other?
- Does these discovery change between summer sport and winter sports?



SECTION 2: INITIAL HYPOTHESIS

- a. Not all the countries are the same. Some country tends to have a younger population hence more young athletes and may be advantaged on older population country. (eg. African's country)
- b. Some other country such as china have an average height smaller than global avarage therefore they will be in disadavantage on sports like: volleyball,basketball, etc.
- c. Some country will be less likely to win winter sports because the lack of mountain and snow to practice.

SECTION 3.1: DATA ANALYSIS APPROACH

- a. The first step is to separate Male athletes from Female athletes, since they do not have the same abilities.
- b. The second step is to separate winter sports from summer sports, because since some country can practice freely snow sports, without particular equipments, the latter will be on advantage on the other. Including 'no Snow Country' (NOS) on the hight (taller than 175cm) and age (above 30) likelyhood to win a medal could lead to bias results.
- c. Analysis by continent, to understand which continent perform better on some sport.
- d. Analysis by sport, to understand which country perform better on some sport

SECTION 3.2: DATA ANALYSIS APPROACH

- a. Probability: What it the likelyhood that having certain data allows you to win a medal
- A/B Testing and (probably) Regression to understand and test wheter some charateristic are usefull to win a medal
- c. Graphs will be needed to understand if some country are improving their overall performance and are improving the overall performance (number of medal won increased).



Milestone 2: Descriptive Stats

Data Analysis Project Proposal

30/01/2022

Leonardo Saviane

The Sample transformation

As we can see our dataset has mostly men (70%) and the average age is around 25 years old, max age is around 71 years old this may suggest that further cleaning is needed, for example removing atheletes above a certain age threshold (eg.50) to avoid bias estimate and because up to a certain age you cannot compete with younger athletes.

	Age	Height	Weight	Sex	Win
count	206165.000000	206165.000000	206165.000000	206165.000000	206165.000000
mean	25.055509	175.371950	70.688337	0.676419	0.146392
std	5.483096	10.546088	14.340338	0.467843	0.353500
min	11.000000	127.000000	25.000000	0.000000	0.000000
25%	21.000000	168.000000	60.000000	0.000000	0.000000
50%	24.000000	175.000000	70.000000	1.000000	0.000000
75%	28.000000	183.000000	79.000000	1.000000	0.000000
max	71.000000	226.000000	214.000000	1.000000	1.000000

Sex indicate with 1 if the athletes is male and 0 if she is female; Win is a binary variable with 1 if the athlete won a medal and 0 if he didn't

The Sample transformation

- Some athletes in 1920 competed at the age of 72(google research), these are clearly outlier that bias the estimate.
- The way people compete change over time, therefore to have interesting insights we should reduce the sample, taking the athletes who competed at games after 1950, without imposing an arbitrary threshold on age directly
- Note that the dataset is parented with local news to they do not take into consideration solely the Olympic games, therefore we should take this into consideration.
- Hence I want to see the sample after the winter Olympics of Moscow in 1980
- I will take into consideration the following Olympics: Rio de Janeiro, London, Beijing, Athina(Athene), Sydney, Atlanta, Barcelona, Seoul, Los Angeles, Moscow, Sochi, Vancouver, Salt Lake City, Nagano, Lillehammer, Albertville, Calgary, Sarajevo, Lake Placid

```
sport_80=pysqldf('''select * from sport where Year>=1980
and City in ('Rio de Janeiro', 'London', 'Beijing', 'Athina', 'Sydney', 'Atlanta', 'Barcelona
sport_m=pysqldf('''select * from sport_80 where Sex=1 ''')
sport_f=pysqldf('''select * from sport_80 where Sex=0 ''')
sport_80[['Age', 'Height', 'Weight', 'Sex', 'Win']].describe()
```

	Age	Height	Weight	Sex	Win
count	136875.000000	136875.000000	136875.000000	136875.000000	136875.000000
mean	25.285728	175.841571	70.983039	0.622466	0.142729
std	5.387866	10.893256	14.992473	0.484772	0.349797
min	12,000000	127.000000	28.000000	0.000000	0.000000
25%	22.000000	168.000000	60.000000	0.000000	0.000000
50%	25.000000	176.000000	70.000000	1.000000	0.000000
75%	28.000000	183.000000	80.000000	1.000000	0.000000
max	71.000000	226.000000	214.000000	1.000000	1.000000

NB. Mean age I expected to decrease, instead it increased

Checking differences between Sexes

wemen

sport f[['Age', 'Height', 'Weight', 'Sex', 'Win']].describe()

	Age	Height	Weight	Sex	Win
count	51675.000000	51675.000000	51675.000000	51675.0	51675.000000
mean	24.397775	168.552143	60.574446	0.0	0.155530
std	5.511821	8.932364	10.523527	0.0	0.362412
min	12.000000	132.000000	28.000000	0.0	0.000000
25%	21.000000	163.000000	54.000000	0.0	0.000000
50%	24.000000	168.000000	60.000000	0.0	0.000000
75%	28.000000	174.000000	66.000000	0.0	0.000000
max	63.000000	213.000000	167.000000	0.0	1.000000

men

sport_m[['Age', 'Height', 'Weight', 'Sex', 'Win']].describe()

	Age	Height	Weight	Sex	Win
count	85200.000000	85200.000000	85200.000000	85200.0	85200.000000
mean	25.824284	180.262711	77.295998	1.0	0.134965
std	5.238486	9.511485	13.724913	0.0	0.341688
min	12,000000	127.000000	37.000000	1.0	0.000000
25%	22.000000	174.000000	68.000000	1.0	0.000000
50%	25.000000	180.000000	76.000000	1.0	0.000000
75%	28.000000	186.000000	85.000000	1.0	0.000000
max	71.000000	226.000000	214.000000	1.0	1.000000



With respect to previous analysis (not reported in the slides) mean age for men and wemen increased, although median and quantiles are constant.







Checking differences between Sexes with A/B TEST

wemen

sport_f[['Age','Height','Weight','Sex','Win']].describe()

	Age	Height	Weight	Sex	Win
count	51675.000000	51675.000000	51675.000000	51675.0	51675.000000
mean	24.397775	168.552143	60.574446	0.0	0.155530
std	5.511821	8.932364	10.523527	0.0	0.362412
min	12.000000	132.000000	28.000000	0.0	0.000000
25%	21.000000	163.000000	54.000000	0.0	0.000000
50%	24.000000	168.000000	60.000000	0.0	0.000000
75%	28.000000	174.000000	66.000000	0.0	0.000000
max	63.000000	213.000000	167.000000	0.0	1.000000

men

sport_m[['Age','Height'	','Weight'	,'Sex','	<pre>Win']].describe()</pre>
-------------------------	------------	----------	------------------------------

	Age	Height	Weight	Sex	Win
count	85200.000000	85200.000000	85200.000000	85200.0	85200.000000
mean	25.824284	180.262711	77.295998	1.0	0.134965
std	5.238486	9.511485	13.724913	0.0	0.341688
min	12,000000	127.000000	37.000000	1.0	0.000000
25%	22.000000	174.000000	68.000000	1.0	0.000000
50%	25.000000	180.000000	76.000000	1.0	0.000000
75%	28.000000	186.000000	85.000000	1.0	0.000000
max	71.000000	226.000000	214.000000	1.0	1.000000

With respect to previous analysis (not reported in the slides) mean age for men and wemen increased, although median and quantiles are constant. Moreover as you can see below the two

groups are heterogenous in weight and height Therefore..perfmance.

```
t_stat, p_val= ss.ttest_ind(weight_m['Weight'],weight_f['Weight'])
print('Difference in Weight; t_stat: ',t_stat, 'p_value: ', p_val)

t_stat, p_val= ss.ttest_ind(height_m['Height'],height_f['Height'])
print('Difference in Height; t_stat: ',t_stat, 'p_value: ', p_val)

Difference in Weight; t_stat: 237.78583480326142 p_value: 0.0

Difference in Height; t_stat: 225.90664378364585 p_value: 0.0
```

Are there State differences?

- We may ask ourself if there are any differences between two continent or state. Here i take two sample countries that i think are representative of their own continent.
- By running an A/B test the null hyphotetis is rejected, therefore there is a difference between the two countries.
- I assume that such difference exists between each continent and may exists between countries.
- Thefore I will create a set of dummies for each state, it will be useful for the regression

```
height ita=pysqldf('''select Height from sport 80 where NOC="ITA" and Sex=1 ''')
height chn=pysqldf('''select Height from sport 80 where NOC="CHN" and Sex=1 ''')
height ita.describe()
                                    height_chn.describe()
           Height
                                               Height
 count 3369,000000
                                    count 2121,000000
        180,266845
                                            177,689769
          8.870825
                                            11.732913
        152.000000
                                            148.000000
        174.000000
                                            170,000000
        180.000000
                                           178,000000
        186.000000
                                            185,000000
       215.000000
                                           226,000000
t_stat, p_val= ss.ttest_ind(height_ita['Height'],height_chn['Height'])
t stat, p val
```

(9.229752559109919, 3.791512826226679e-20)

Are there State differences? Box-plot ita-chn

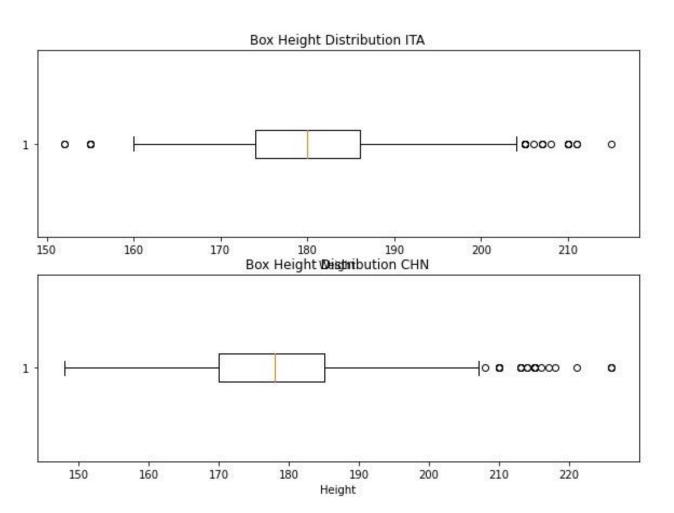
```
fig, axs = plt.subplots(2,figsize=(10,7))

##BOX PLOT
axs[0].boxplot(height_ita['Height'],vert=False)

axs[0].set_xlabel('Weight') # add to x-label to the plot
axs[0].set_title('Box Height Distribution ITA') # add title to the plot
axs[1].boxplot(height_chn['Height'],vert=False)

axs[1].set_xlabel('Height') # add to x-label to the plot
axs[1].set_title('Box Height Distribution CHN') # add title to the plot
```

Codes and results



Weight and Height differences beween Winter and Summer Olympics? A/B Test and distributions graphs

- Is there a difference weight and height between winter olympcs and summer olympics?
- Since we have seen that there is a difference between man and wemen, I further explore this analysis by diving the two sexes
- A/B test shows that between
 Season=1(summer) and winter there are weight and height differences! exept for male weight which is the same in the two groups!
- You will see many outlier in the box plot for male in summer Olympics

```
height_count_m_s=pysqldf('''select Height,Count(Height) as 'Height_Count' from sport_80 where Sex=1 and Season = 1 group by Height ''')
height_count_f_s=pysqldf('''select Height,Count(Height) as 'Height_Count' from sport_80 where Sex=0 and Season = 1 group by Height ''')
height_m_s=pysqldf('''select Height from sport_80 where Sex=1 and Season = 1 ''')
height_f_s=pysqldf('''select Height from sport_80 where Sex=0 and Season = 1 ''')
height_count_m_w=pysqldf('''select Height,Count(Height) as 'Height_Count' from sport_80 where Sex=1 and Season = 0 group by Height ''')
height_m_w=pysqldf('''select Height,Count(Height) as 'Height_Count' from sport_80 where Sex=0 and Season = 0 group by Height ''')
height_m_w=pysqldf('''select Height from sport_80 where Sex=1 and Season = 0 ''')
height_f_w=pysqldf('''select Height from sport_80 where Sex=0 and Season = 0 ''')
```

A/B TESTING

```
t_stat, p_val= ss.ttest_ind(height_m_s['Height'],height_m_w['Height'])
t_stat, p_val

(10.537528545769307, 6.017725959006674e-26)

t_stat, p_val= ss.ttest_ind(height_f_s['Height'],height_f_w['Height'])
t_stat, p_val

(20.84581328231301, 4.1439688067478515e-96)

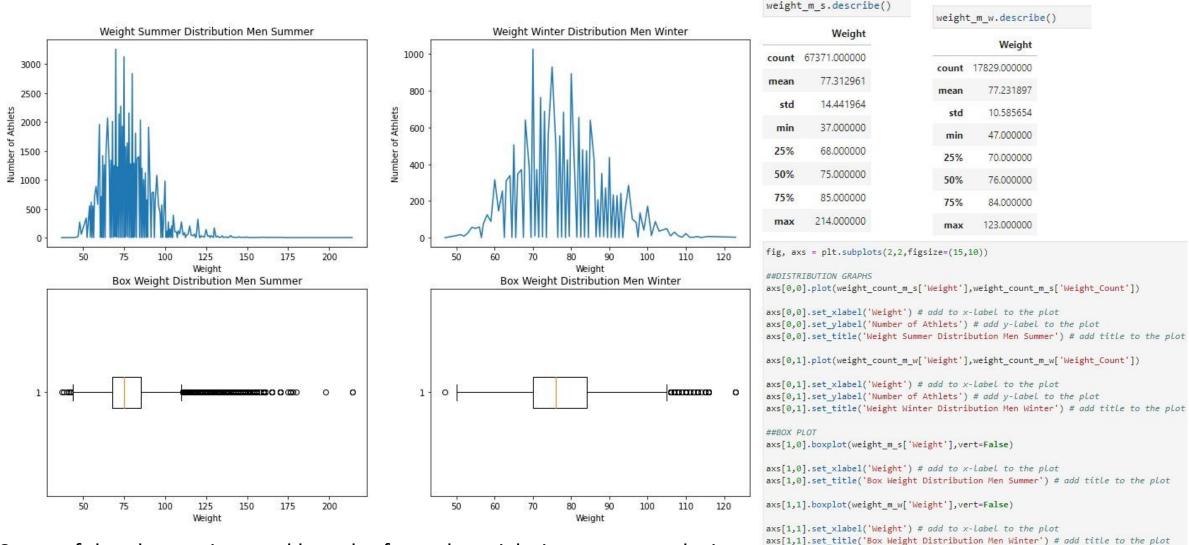
t_stat, p_val= ss.ttest_ind(weight_m_s['Weight'],weight_m_w['Weight'])
t_stat, p_val

(0.7012859803826186, 0.4831264765792309)

t_stat, p_val= ss.ttest_ind(weight_f_s['Weight'],weight_f_w['Weight'])
t_stat, p_val= ss.ttest_ind(weight_f_s['Weight'],weight_f_w['Weight'])
t_stat, p_val= ss.ttest_ind(weight_f_s['Weight'],weight_f_w['Weight'])
t_stat, p_val= ss.ttest_ind(weight_f_s['Weight'],weight_f_w['Weight'])
```

SQL CODE ABOVE

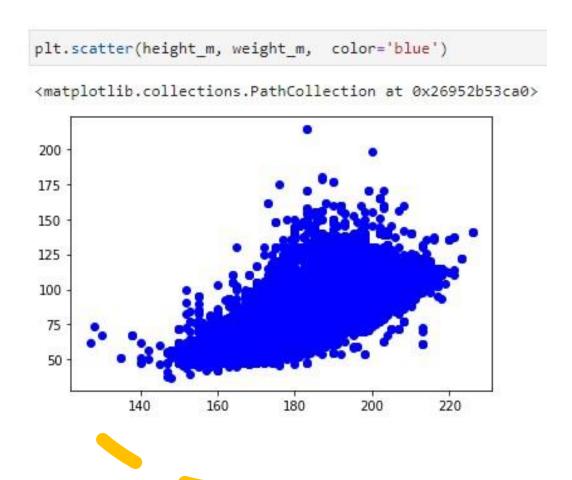
Weight and Height differences beween Winter and Summer Olympics? A/B Test and distributions graphs

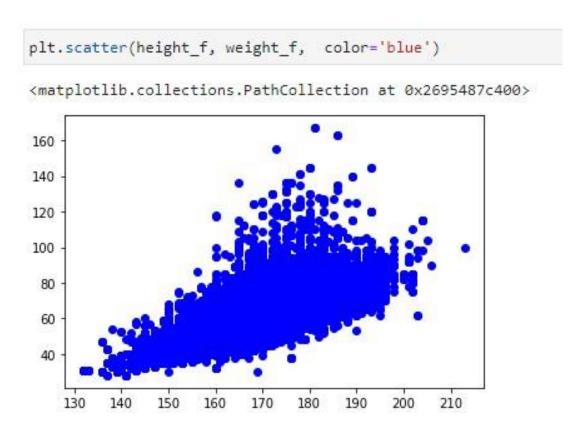


Count of the observations and box plot for male weight in summer and winter

EXTRA: Scatter plot between height and weight

The following graphs show a positive correlation between Height and Weight, respectively for Man and Women





Which is the sport that counts the highest number of athletes?

```
obv_sport=pysqldf('''select Sport,Count(Sport) as 'Sport_Count' from sport_80 group by Sport ''')

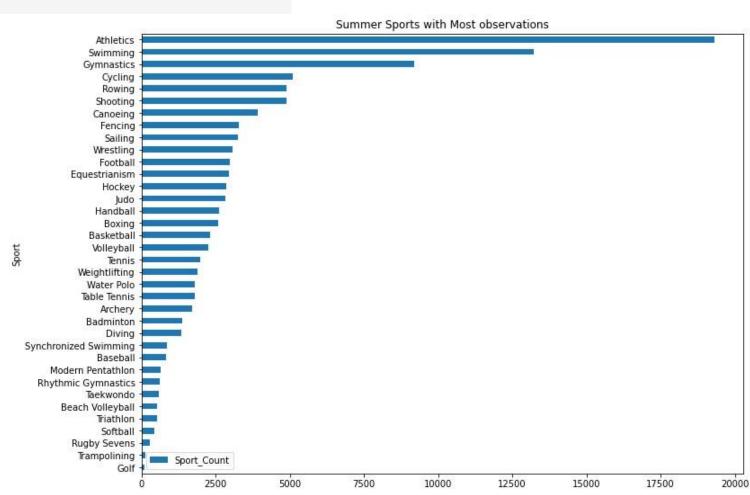
obv_sport_s=pysqldf('''select Sport,Count(Sport) as 'Sport_Count' from sport_80 Where Season = 1 group by Sport ''')

obv_sport_s.sort_values(by='Sport_Count', ascending=True, inplace=True)

obv_sport_s.set_index('Sport', inplace=True)

obv_sport_s.plot.barh(figsize=(12,9)).set_title('Summer Sports with Most
```

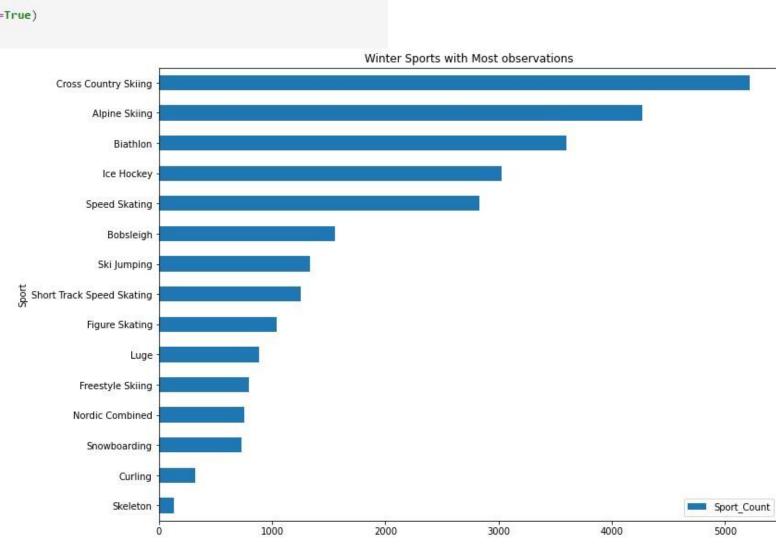
Se we can see the most popular sport among the Olympic summer games is Athelitics with counted more and 20000 atheltes throughout the years



Which is the sport that counts the highest number of athletes?

```
obv_sport_w=pysqldf('''select Sport,Count(Sport) as 'Sport_Count' from sport_80 Where Season = 0 group by Sport ''')
obv_sport_w.sort_values(by='Sport_Count', ascending=True, inplace=True)
obv_sport_w.set_index('Sport', inplace=True)
obv_sport_w.plot.barh(figsize=(12,9)).set_title('Winter Sports wi
```

Se we can see the most popular sport among the Olympic Winter games is Cross country Skiing with counted more and 5000 atheltes throughout the years



LAST: Is there difference in Height and Weight between Sport? A/B Testing

- If there exists a difference in height and weight between each sport, the group of athletes must be separated in the regression becouse otherwise they will bias the estimate.
- AB test works perfectly here, in the subcategory of the men I will take two sports randomly(Gymnastics and Athletics) and perfom the test on height and weight.
- I will assume that the results it's rapresentative of each category.

Weight T-stat: -53.347663315055506 Weight p-value: 0.0

```
Sport_tt1=pysqldf('''select Height, Weight from sport_80 Where Sport = 'Gymnastics' and Sex=1 ''')
Sport_tt2=pysqldf('''select Height, Weight from sport_80 Where Sport = 'Athletics' and Sex=1 ''')

t_stat, p_val= ss.ttest_ind(Sport_tt1['Height'], Sport_tt2['Height'])
print('Height T-stat: ', t_stat, 'Height p-value: ', p_val)

t_stat, p_val= ss.ttest_ind(Sport_tt1['Weight'], Sport_tt2['Weight'])
print('Weight T-stat: ', t_stat, 'Weight p-value: ', p_val)

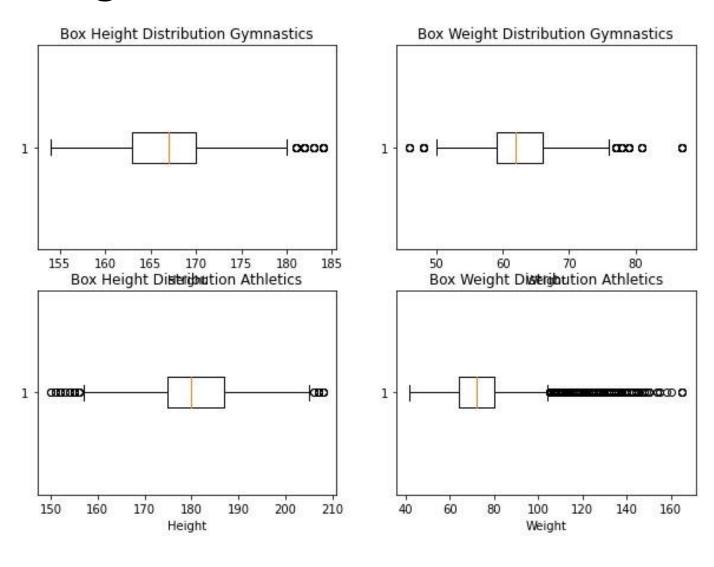
Height T-stat: -111.3967094784561 Height p-value: 0.0
```

The two groups are strongly different! Null hyphotesis is rejected!

LAST: Is there difference in Height and Weight between Sport? A/B Testing

- As we can see from the box plot the median value differes a lot from the other.
- Athletics weights have a lot of outlier which may bias estimate

In the next slide the findings to implement in the last part!



CONCLUSION Milestone 2: Descriptive Stats

Codes:

https://github.com/LeonSavi/Projects_Leonardo_Saviane.git

- Most of the hypothesis i had made at beginning can be tested only in the next part.
- AB test comes in handy for testing differences between groups and it have been used to test if the groups are similar
- Since there is a difference between sexes they should be tested separately in the regression to avoid bias estimate
- Similarly for Winter and Summer Sport, athletes have different weight and height therefore we shouldn't confound them.
- Again, I tested difference between athletes of different sport and there is a statistically significative difference.
- Since are many sports i will explore only the 4 most popular sport among the athletes: Cross Country Skiing (most popular among winter sports), Gymnastics(third to last most popular sport among summer sports), Swimming(second to last most popular sport among summer sports), Atheletics(most popular sport among summer sports).
- To see if someone is a potential athletes we have to take into consideration all these things in the next part.
- Since are many sports i will explore only the 4 most popular sport among the athletes: Cross Country Skiing (most popular among winter sports), Gymnastics (third to last most popular sport among summer sports), Swimming (second to last most popular sport among summer sports), Atheletics (most popular sport among summer sports).



Milestone 3: Beyond Descriptive Stats

Data Analysis Project Proposal

30/01/2022

Leonardo Saviane

CODES:

https://github.com/LeonSavi/Projects_Leonardo_Saviane.git

PREPERARING FOR REGRESSION

- One of the first thing that we have to do for our regression is to create a dummy for each country, which means if a an athletes belong to country(=1) or not(=0)
- Then through SQL we get the data we need to put in our regression and watch the coefficient
- NOC_ITA will be removed to avoid collinearity, the intercept/costant will capture italy

Now I want to create binary variable for each NOC, in order to see the influence of each country

the following variables will be removed from sport_80: ID, Name, City, Year, NOC_ITA(to avoid collinearity), costant will capture ITALIAN observations



	Age	Height	Weight	Sport	Sex	Season	Win	NOC_AFG	NOC_AHO	NOC_ALB		NOC_VAN	NOC_VEN	NOC_VIE	NOC_VIN	NOC_YAR	NOC_YEM	NOC_YMD	NOC_YUG	NOC_ZAM	NOC_ZIM
0	24.0	180.0	80.0	Basketball	1	1	0	0	0	0	in	0	0	0	0	0	0	0	0	0	0
1	23,0	170.0	60.0	Judo	1	1	0	0	0	0	***	0	0	0	0	0	0	0	0	0	0
2	21.0	185.0	82.0	Speed Skating	0	0	0	0	0	0		0	0	0	0	0	0	0	0	0	0
3	21,0	185.0	82.0	Speed Skating	0	0	0	0	0	0	***	0	0	0	0	0	0	0	0	0	0
4	25.0	185.0	82.0	Speed Skating	0	0	0	0	0	0	in	0	0	0	0	0	0	0	0	0	0

5 rows × 224 columns



```
1 reg1_m_x=pysqldf('''select * from dataset Where Sport = 'Cross Country Skiing' and Sex = 1 ''')
        2 reg1 m y=pysqldf('''select Win from dataset Where Sport = 'Cross Country Skiing' and Sex = 1 ''')
                  reg1_f_x=pysqldf('''select * from dataset Where Sport = 'Cross Country Skiing' and Sex = 0 ''')
                 reg1 f y=pysqldf('''select Win from dataset Where Sport = 'Cross Country Skiing' and Sex = 0 ''')
                 reg1_m_x.drop(['Win', 'Sport', 'Season'], axis=1, inplace=True)
        8 reg1_f_x.drop(['Win', 'Sport', 'Season'], axis=1, inplace=True)
         9 reg1 m x.head()
          Age Height Weight Sex NOC_AFG NOC_AHO NOC_ALB NOC_ALG NOC_AND NOC_ANG ... NOC_VAN NOC_VEN NOC_VIE NOC_VIN NOC_YAR NOC_YEM NOC_
 0 31.0
                                                           75.0
1 31.0 188.0
                                                           75.0
                                                                                                                                                                                                                                                   0
                                                                                                                                                                                                                                                                                      0 ...
                                                                                                                                                                                                                                                                                                                                0
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              0
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                  0
2 31.0
                          188.0
                                                           75.0
3 31.0 188.0
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              0
4 33.0 188.0
                                                           75.0
```

5 rows × 221 columns

REGRESSION SET-UP(1)

- While running the regression all country variables seemed to not be statistically significant, R2 was low and F statistic as well. This means that there were too many variables.
- To explore country effect on the probability of winning a medal more complex model are need. Some model that can exploit the difference between groups.
- In this project I stick with multiple linear regression (OLS)
- To avoid the problem I said i take into consideration only some of the variables such as Age, Weight, Height. The intercept is added in the model.

```
reg1_m_x=pysqldf('''select Age,Weight,Height from dataset Where Sport = 'Cross Country Skiing' and Sex = 1 ''')

## X is the input variables (or independent variables)

X = reg1_m_x

## y is the target/dependent variable

y = reg1_m_y

## add an intercept (beta_0) to our model

X = sm.add_constant(X)

model = sm.OLS(y, X).fit()
predictions = model.predict(X)

## Print out the statistics
model.summary()
```

REGRESSION SET-UP(2)

- This is the correction I made for each sex and sport.
- And these are the code used for the OLS regression

MEN WOMEN

		OLS F	Regressi	on Resu	lts		
Dep	. Variable	:	W	in	R-sq	uared:	0.01
	Model	:	O	LS A	dj. R-sq	uared:	0.01
	Method	: Leas	t Squar	es	F-sta	atistic:	8.54
	Date	: Sun, 30	Jan 202	22 Pro	b (F-sta	tistic):	1.21e-0
	Time	:	17;28:4	40 L o	og-Likeli	ihood:	-389.4
No. Obs	ervations	:	229	95		AIC:	786.
Df	Residuals		229	91		BIC:	809.
	Df Model	:		3			
Covaria	nce Type	: r	onrobu	st			
	coef	std err	t	P> t	[0.025	0.975]	Į.
const	-0.2911	0.206	-1.412	0.158	-0.695	0.113	Ú.
Age	0.0065	0.001	4.743	0,000	0.004	0.009	
Weight	0.0014	0.002	0.760	0,448	-0.002	0.005	i.
Height	0.0008	0.002	0.522	0.602	-0.002	0.004	
0	mnibus:	1264,762	Durl	bin-Wa	tson:	1,508	
Prob(On	nn <mark>i</mark> bus):	0.000	Jarqu	e-Bera	(JB): 6	288,098	
	Skew:	2.785		Prob	(JB):	0.00	
K	Curtosis:	8.893		Cond	. No. 6	.15e+03	

		OLS F	Regressi	on Resu	lts		
Dep	. Variable	:	W	in	R-sq	juared:	0.022
	Model	:	Ol	S A	dj. R-sq	juared:	0.021
	Method	: Leas	t Square	es	F-st	atistic:	22.14
	Date	: Sun, 30	Jan 202	22 Pro	b (F-sta	itistic):	3.57e-14
	Time	:	17:28:1	14 Lo	g-Likel	lihood:	- <mark>1</mark> 46.30
No. Obs	ervations	:	291	18		AIC:	300.6
Df	Residuals	:	291	14		BIC:	324.5
	Df Model	:		3			
Covaria	nce Type	: r	onrobu	st			
	coef	std err	t	P> t	[0.025	0.975]
const	-0.8380	0.169	-4.962	0.000	-1.169	-0.507	7
Age	0.0058	0.001	5.279	0.000	0.004	0.008	3
Weight	0.0009	0.001	0.703	0.482	-0.002	0.003	3
Height	0.0039	0.001	3.006	0.003	0.001	0.006	5
O	mnibus:	1894.144	Durl	bin-Wa	tson:	1.74	4
Prob(On	nnibus):	0.000	Jarqu	e-Bera	(JB):	14206.97	7
	Skew:	3.218		Prob	(JB):	0.0	0
K	urtosis:	11.684		Cond	No.	6.99e+0	3

Results Cross Country Skiing

- As we can see for man, weight and height seems to not be statistically significant this can be due to small sample (2500 observation each).
- R-squared and F-test are not very encuraging on the results.
- Probably a non linear regression would give better results

Results Gymnastics

- Except for the age the other variable are statistically significant
- As we expected smaller people are more likely to win a medal (coefficient negative on weight and height)
- The smaller you are the better

MEN WOMEN

OLS	Regressio	on Resu	lts					010			To .		
	Wi	n	R-squ	uared:	0.009	Den	Variable		- F			nuared:	0.008
	OL	S A	dj. R-sq	uared:	0.008	J.C.P.							0.008
Leas	st Square	:5	F-sta	atistic:	15.40		14100000000	-					10.84
Sun, 30) Jan 202	2 Pro	b (F-sta	tistic):	5.67e-10								1,7076
	17:30:5	6 Lo	g-Likeli	hood:	-116.19						N. S. Office Co.		-171.85
	531	3		AIC:	240.4	No. Obs						AIC:	351.7
	530	9		BIC:	266.7	Df	Residuals	:	387	77		BIC:	376.7
		3				1	Df Model	:		3			
1	nonrobus	st				Covaria	nce Type	: 1	nonrobu	st			
std err	t	P> t	[0.025	0.975	L		coef	std err	t	P> t	[0.025	5 0.975	1
0.129	5.372	0.000	0.441	0.949)	const	0.3035	0.112	2.703	0.007	0.083	3 0.524	4
0.001	-1.023	0.306	-0.003	0.001		Age	0.0015	0.002	1.001	0.317	-0.00	1 0.004	4
0.001	-1.688	0.092	-0.003	0.000)	Weight	-0.0031	0.001	-3.055	0.002	-0.00	5 -0.00	1
0.001	-3.131	0.002	-0.005	-0.001		Height	-0.0008	0.001	-0,883	0.377	-0,003	3 0.00	1
3673.557	Durb	in-Wa	tson:	1.74	2	O	mnibus:	2608,201	Durl	oin-Wa	tson:	1.56	1
0.000	Jarque	e-Bera	(JB): 3	2838.16	1	Prob(On	nnibus):	0.000	Jarqu	e-Bera	(JB):	21328.32	9
	79						Skew:	3,347		Prob	(JB):	0.0	0
			- T. C.			К	urtosis:	12.331		Cond	. No.	4.46e+0	3
	Lea: Sun, 30 std err 0.129 0.001 0.001 0.001 3673.557 0.000 3.450	Wi OL Least Square Sun, 30 Jan 202 17:30:5 531 530 nonrobus std err t 0.129 5.372 0.001 -1.023 0.001 -1.688 0.001 -3.131 3673.557 Durb	Win OLS A Least Squares Sun, 30 Jan 2022 Pro 17:30:56 Lo 5313 5309 3 nonrobust std err t P> t 0.129 5.372 0.000 0.001 -1.023 0.306 0.001 -1.688 0.092 0.001 -3.131 0.002 3673.557 Durbin-War 0.000 Jarque-Bera 3.450 Prob	OLS Adj. R-squ Least Squares F-sta Sun, 30 Jan 2022 Prob (F-star 17:30:56 Log-Likeli 5313 5309 3 nonrobust std err t P> t [0.025 0.129 5.372 0.000 0.441 0.001 -1.023 0.306 -0.003 0.001 -1.688 0.092 -0.003 0.001 -3.131 0.002 -0.005 3673.557 Durbin-Watson: 0.000 Jarque-Bera (JB): 3 3.450 Prob(JB):	Win R-squared: OLS Adj. R-squared: Least Squares F-statistic: Sun, 30 Jan 2022 Prob (F-statistic): 17:30:56 Log-Likelihood: 5313 AIC: 5309 BIC: 3 nonrobust std err t P> t [0.025 0.975] 0.129 5.372 0.000 0.441 0.949 0.001 -1.023 0.306 -0.003 0.001 0.001 -1.688 0.092 -0.003 0.000 0.001 -3.131 0.002 -0.005 -0.001 3673.557 Durbin-Watson: 1.744 0.000 Jarque-Bera (JB): 32838.16 3.450 Prob(JB): 0.00	Win R-squared: 0.009 OLS Adj. R-squared: 0.008 Least Squares F-statistic: 15.40 Sun, 30 Jan 2022 Prob (F-statistic): 5.67e-10 17:30:56 Log-Likelihood: -116.19 5313 AIC: 240.4 5309 BIC: 266.7 3 nonrobust std err t P> t [0.025 0.975] 0.129 5.372 0.000 0.441 0.949 0.001 -1.023 0.306 -0.003 0.001 0.001 -1.688 0.092 -0.003 0.000 0.001 -3.131 0.002 -0.005 -0.001 3673.557 Durbin-Watson: 1.742 0.000 Jarque-Bera (JB): 32838.161 3.450 Prob(JB): 0.00	Win R-squared: 0.009 OLS Adj. R-squared: 0.008 Least Squares F-statistic: 15.40 Sun, 30 Jan 2022 Prob (F-statistic): 5.67e-10 17:30:56 Log-Likelihood: -116.19 5313 AIC: 240.4 No. Obsection 5309 BIC: 266.7 Df 3 nonrobust Covaria std err t P> t [0.025 0.975] 0.129 5.372 0.000 0.441 0.949 const 0.001 -1.023 0.306 -0.003 0.001 Age Weight 0.001 -3.131 0.002 -0.005 -0.001 Height 3673.557 Durbin-Watson: 1.742 Prob(On 0.000 Jarque-Bera (JB): 32838.161 3.450 Prob(JB): 0.000	Win R-squared: 0.009 Dep. Variable OLS Adj. R-squared: 0.008 Model Least Squares F-statistic: 15.40 Method Sun, 30 Jan 2022 Prob (F-statistic): 5.67e-10 Date 17:30:56 Log-Likelihood: -116.19 Time 5313 AIC: 240.4 No. Observations 5309 BIC: 266.7 Df Residuals 3 Df Model Covariance Type std err t P> t [0.025 0.975] coef 0.129 5.372 0.000 0.441 0.949 const 0.3035 0.001 -1.688 0.092 -0.003 0.001 Age 0.0015 0.001 -3.131 0.002 -0.005 -0.001 Height -0.0008 3673.557 Durbin-Watson: 1.742 Omnibus: Prob(Omnibus): Skew: 3.450 Prob(JB): 0.00 Kurtosis: Kurtosis:	Win R-squared: 0.009 Dep. Variable:	Win R-squared: 0.009 Dep. Variable: Wind No. Ol.	No. Observations: Std err t P > t	No. Observations: Sun, 30 Jan 2022 Prob (F-statistic) 5.67e-10 Date: Sun, 30 Jan 2022 Prob (F-statistic) 5.67e-10 Date: Sun, 30 Jan 2022 Prob (F-statistic) 5.67e-10 Date: Sun, 30 Jan 2022 Prob (F-statistic) F-s	Win R-squared: 0.009 Dep. Variable: Win R-squared: 0.008 Least Squares F-statistic: 15.40 Method: Least Squares F-statistic: 5.67e-10 Date: Sun, 30 Jan 2022 Prob (F-statistic): 5.67e-10 Date: Sun, 30 Jan 2022 Prob (F-statistic): F-statistic: Sun, 30 Jan 2022 Prob (F-statistic): Least Squares F-statistic: Sun, 30 Jan 2022 Prob (F-statistic): Date: Sun, 30 Jan 2022 Prob (F-statistic): Log-Likelihood: Inc. Sun, 30 Jan 2022 Prob (F-statistic): Log-Likelihood: Inc. Inc.

MEN WOMEN

		OLS I	edicasi	on nesu	11.2		
Dep	. Variable	:	W	in	R-s	quared:	0.033
	Model	:	OI	LS A	dj. R-s	quared:	0.032
	Method	: Leas	t Squar	es	F-s	tatistic:	67.15
	Date	: Sun, 30	Jan 202	22 Pro	b (F-st	atistic):	1.06e-42
	Time	:	17:40:5	52 L o	g-Like	elihood:	-1886.
No. Obs	ervations	:	599	93		AIC:	3781
Df	Residuals	:	598	39		BIC:	3808
1	Df Model	:		3			
Covaria	nce Type	: r	nonrobu	st			
	coef	std err	t	P> t	[0.02	5 0.975	1
const	-1.0801	0.123	-8.782	0.000	-1.32	1 -0.839)
Age	0.0039	0.001	3.475	0.001	0.00	2 0.006	ō
Weight	0.0038	0.001	3.907	0.000	0.00	2 0.006	5
Height	0.0052	0.001	5.586	0.000	0.00	3 0.007	13
0	mnibus:	2223.010	Durl	bin-Wa	tson:	1,396	
Prob(Or	nnibus):	0.000	Jarqu	e-Bera	(JB):	6003.746	
	Skew:	2,080		Prob	(JB):	0.00	
H	Curtosis:	5.597		Cond	No.	5.32e+03	

OLS Regression Results

		OLS F	Regressi	on Resu	ılts		
Dep	. Variable	:	in	R-s	quared:	0.041	
	Model	:	Ol	S A	dj. R-s	quared:	0.040
	Method	: Leas	t Square	es F-s		tatistic:	102.5
	Date	: Sun, 30	Jan 202	2022 Prob		atistic):	5.30e-65
	Time	:	17:40:3	31 L	og-Like	elihood:	-1820.6
No. Obs	ervations	:	722	24		AIC:	3649.
Df	Residuals	:	7220			BIC:	3677.
	Df Model	:		3			
Covaria	nce Type	: r	nonrobu	st			
	coef	std err	t	P> t	[0.02	5 0.975]	
const	-0.8752	0.107	-8.202	0.000	-1.08	4 -0.666	
Age	0.0024	0.001	2.197	0.028	0.00	0.005	
Weight	0.0053	0.001	7.580	0.000	0.00	4 0.007	ŧ
Height	0.0028	0.001	3.625	0.000	0.00	1 0.004	
0	mnibus:	3070.015	Durl	oin-Wa	tson:	1.415	5
Prob(Omnibus):		0.000	Jarqu	e-Bera	(JB):	10083.866	5
	Skew:	2.278		Prob	(JB):	0.00)
K	urtosis:	6,569	Cond	No.	5.94e+03	3	

Results Swimming

- All coefficient are statistically significant
- F-statistics are high and encouraging
- The hight and weight have a posive impact on the probability of winning a medal this suggest that taller people are in advantage wrt to the other.
- Although i think that this relationship is non-linear

MEN WOMEN

OLS Regression Results							OLS Regression Results								
Dep. Variable: Win		in	R-squared:		0.006	Dep	Dep. Variable:		Win		R-squared:		0.004		
Model: OLS		S A	Adj. R-squared:		0.006		Model	l:	0	OLS Adj. R-sq		quared:	0.004		
Method: Least Square		es	F-statisti		22.79		Method	l: Leas	Least Squares		F-statistic:		11.24		
Date: Sun, 30 Jan 2022		22 Pro	Prob (F-statistic):		1.05e-14		Date	: Sun, 30	Sun, 30 Jan 202		2 Prob (F-statistic)		2.35e-07		
Time: 17:4		17:42:3	39 L o	Log-Likelihood:		-1040.9		Time:		17;43:01		Log-Likelihood:		-1556.0	
No. Observations: 111		39		AIC:	2090.	No. Obs	ervations	:	813	35		AIC:	3120.		
Df Residuals: 11		1118	35		BIC:	2119.	Df	Residuals	:	813	31		BIC:	3148.	
Df Model:		:		3				1	Df Model	l:		3			
Covariance Type:		: 1	nonrobu	st				Covaria	nce Type	: 1	nonrobu	ıst			
	coef	std err	t	P> t	[0.025	0.975]	ĺ		coef	std err	t	P> t	[0.02	5 0.975]
const	-0.2283	0.067	-3.399	0.001	-0.360	-0.097		const	-0.2771	0.085	-3.243	0.001	-0.44	5 -0.110	0
Age	-0.0013	0.001	-2.292	0.022	-0.002	-0.000	i .	Age	0.0016	0.001	2.316	0.021	0.00	0.003	3
leight	0.0004	0.000	2.021	0.043	1.3e-05	0.001		Weight	0.0004	0.000	1.055	0.291	-0.00	0.00	1
leight	0.0017	0.000	4.107	0.000	0.001	0.003		Height	0.0018	0.001	3.205	0.001	0.00	1 0.003	3
0	mnibus:	6992.686	Durl	oin-Wa	tson:	1.641		O	mnibus:	4334.156	Dur	bin-Wa	tson:	1.64	1
rob(Or	nnibus):	0.000	Jarqu	e-Bera	(JB): 48	270.535		Prob(On	nnibus):	0.000	Jarqu	ie-Bera	(JB):	20643.25	6
	Skew:	3.148		Prob	(JB):	0.00			Skew:	2.740	ì	Prob	(JB):	0.0	0
H	Curtosis:	10.994		Cond	. No. 5	.30e+03		K	urtosis:	8.556	Ü	Cond	. No.	4.78e+0	3

Results Athletics

- Weight is not statistically significant for the regression.
- This suggest that it should be removed and not taken into consideration
- Fat people tend to be in disadvantage anyway in this sport
- Weirdly age has a negative impact on probability of victory for man but it is positve for women

```
X_train, X_test, y_train, y_test = train_test_split( reg1_m_x, reg1_m_y, test_size=0.2, random_state=4)
print ('Train set:', X_train.shape, y_train.shape)
print ('Test set:', X_test.shape, y_test.shape)

LR = LogisticRegression(C=0.01, solver='liblinear').fit(X_train,y_train.values.ravel())

yhat_lr = LR.predict(X_test)
yhat_lr_prob = LR.predict_proba(X_test)

print('JACCARD: ', jaccard_score(y_test, yhat_lr,pos_label=0)) #not winning a medal
print('F1: ', f1_score(y_test, yhat_lr, average='weighted'))
print('LogLoss: ', log_loss(y_test, yhat_lr_prob))
```

MACHINE LEARNING: PREDICTION MODEL

- I suggest a logistic model which works very good with binary dependent variable
- I will use Sklearn library of python to run this kind of regression
- If the logloss is low the prediction model is good and can be used to forecast!
- I will use a train sample (80% of the sample) and a train set (20%), obtained randomly through «train_test_split»

LOGISTIC REGRESSION MODEL

- For the sake of brevity I will report some of the scores, the results are similar also for the other models!
- As you can see the model seem to predict very well if an Athlete won't win a medal:

ATHLETICS WOMEN

Train set: (8951, 221) (8951, 1)
Test set: (2238, 221) (2238, 1)
JACCARD: 0.9186773905272565
F1: 0.8797395039236328
LogLoss: 0.27210371220593615

GYMNASTIC WOMEN

Train set: (3104, 221) (3104, 1)
Test set: (777, 221) (777, 1)
JACCARD: 0.9395109395109396
F1: 0.9102096693337569
LogLoss: 0.2105092609797363

ATHLETICS MEN

Train set: (6508, 221) (6508, 1)
Test set: (1627, 221) (1627, 1)
JACCARD: 0.9084204056545789
F1: 0.864827928861493
LogLoss: 0.296838593714039

GYMNASTIC MEN

Train set: (4250, 221) (4250, 1)
Test set: (1063, 221) (1063, 1)
JACCARD: 0.9275634995296331
F1: 0.8927063060382804
LogLoss: 0.253696171549855

THE LOG LOSS FUNCTION SI LOW

JACCARD AND F1 SCORE ARE HIGH!

CONCLUSIONS:

- My purpose as stated earlier is to find, for potential gym institute, a model to understand if somone has the potential to win a medal (binary variable) given age, height, weight, country(binary variable for each country whether an athlete belong to that country).
- I found out from a previous descriptive analysis that there is a difference between sexes they should be tested separately in the regression to avoid bias estimate.
- moreover i found out previously (always by mean of a AB test) that each country has athletes with height and weight different from another country.
- The first challenge that i faced was in the linear regression, it seemed that t-statistics for country variables is not statistically significant therefore i remove them and considered solely the following variable: age, weight, height and the intercepet.
- i run 4 regression (linear OLS) for each sex, in every regression i used 4 subsample of athletes, that is: Cross country skii, gymnastic, swimming, athletics athletes.
- the most interesting findings are:
- for gymnastic weight and height coefficient are negative: this is in line with the idea that the smaller you are the better you perform, both for male and female. althought the p-value for the Weight variable is close to 0.95 percentile. age coef is positive (male coef: 0.0015); this is a difficult sport and experince helps improving the performance up to a certain age.
- similar result for cross country skiing; age, height, weight have negative coefficient. p-value of height is statistically significant.
- In swimming all variable are positive and statistically significant!; eg for male the coefficient of weight and height are: 0.0053 and 0.0028, the taller you are the easier is to win.
- In atheltics all variable are statistically significant for both group, variable age for man seems to have a negative impact on the probability for the athelete to win a medal (-0.0013) for women is positive (0.0016)
- I estimated also a *logistic regression* which is very usefull when the dependent variable is binary. I did on the full sample for each of the 8 cases listed above, using the sklearn library and training test and test set (the latter is 20% of the sample). It turned out that the model is very good in predicting if an athlete is not going to win a medal! LOG LOSS value in all the cases is close to 0.3 and jaccard and f1 score is close to 0.9!
- CODES: https://github.com/LeonSavi/Projects_Leonardo_Saviane.git