**RASPAGEM DE DADOS - BRASILEIRÃO 2023**

Por: Leandro Dias Vieira

- Instalando bibliotecas

In [63]:
```
!pip install Beautifulsoup4
```

Requirement already satisfied: Beautifulsoup4 in /usr/local/lib/python3.10/dist-pa
ckages (4.11.2)
Requirement already satisfied: soupsieve>1.2 in /usr/local/lib/python3.10/dist-pac
kages (from Beautifulsoup4) (2.5)

In [64]:
```
!pip install requests
```

Requirement already satisfied: requests in /usr/local/lib/python3.10/dist-packages
(2.31.0)
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.
10/dist-packages (from requests) (3.3.2)
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-pack
ages (from requests) (3.6)
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.10/dis
t-packages (from requests) (2.0.7)
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.10/dis
t-packages (from requests) (2023.11.17)

- Importando Bibliotecas

In [65]:
```
import requests
from bs4 import BeautifulSoup
import pandas as pd
import numpy as np
```

- Pegando o Link

In [66]:
```
url = "https://www.espn.com.br/futebol/classificacao/_/liga/bra.1"
```

- Passando parametros para acessar site

In [67]:
```
headers = {"User-Agent" : "Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/53
```

- Testando requisicao

```
In [68]: requisicao = requests.get(url, headers = headers)
```

```
In [69]: print(requisicao)
```

```
<Response [200]>
```

```
In [70]: ## print(requisicao.text)
```

- Parseando o site

```
In [71]: site = BeautifulSoup(requisicao.text, "html.parser")
```

- Função Prettify(), para arrumar HTML e Concatenar

```
In [72]: site2 = site.prettify()
```

```
In [73]: ## print(site2)
```

# Importando a primeira tabela

- Encontrando a tag da primeira tabela

```
In [74]: tabela = site.find("table")
```

```
In [75]: print(tabela)
```

```html
<table class="Table Table--align-right Table--fixed Table--fixed-left" style="border-collapse:collapse;border-spacing:0"><colgroup class="Table__Colgroup"><col class="Table__Column"/></colgroup><thead class="Table__header-group Table__THEAD"><tr class="Table__sub-header Table__TR Table__even"><th class="subHeader__item--content Table__TH" title=""><span class="fw-medium w-100 dib subHeader__item--content" title="2023">2023</span></th></tr></thead><tbody class="Table__TBODY"><tr class="Table__TR Table__TR--sm Table__even" data-idx="0"><td class="Table__TD"><div class="team-link flex items-center clr-gray-03"><span class="team-position ml2 pr3">1</span><span class="pr4 TeamLink__Logo"><a class="AnchorLink" data-clubhouse-uid="s:600~t:2029" href="/futebol/time/_/id/2029/palmeiras" tabindex="0"><img alt="PAL" class="Image Logo Logo__sm" data-mptype="image" src="data:image/gif;base64,R0lGODlhAQABAIAAAAAAAP///yH5BAEAAAAALAAAAAABAAEAAAIBRAA7" title="PAL"/></a></span><span class="dn show-mobile"><a class="AnchorLink" data-clubhouse-uid="s:600~t:2029" href="/futebol/time/_/id/2029/palmeiras" tabindex="0"><abbr data-clubhouse-uid="s:600~t:2029" style="text-decoration:none" title="Palmeiras">PAL</abbr></a></span><span class="hide-mobile"><a class="AnchorLink" data-clubhouse-uid="s:600~t:2029" href="/futebol/time/_/id/2029/palmeiras" tabindex="0">Palmeiras</a></span></div></td></tr><tr class="filled Table__TR Table__TR--sm Table__even" data-idx="1"><td class="Table__TD"><div class="team-link flex items-center clr-gray-03"><span class="team-position ml2 pr3">2</span><span class="pr4 TeamLink__Logo"><a class="AnchorLink" data-clubhouse-uid="s:600~t:6273" href="/futebol/time/_/id/6273/gremio" tabindex="0"><img alt="GRE" class="Image Logo Logo__sm" data-mptype="image" src="data:image/gif;base64,R0lGODlhAQABAIAAAAAAAP///yH5BAEAAAAALAAAAAABAAEAAAIBRAA7" title="GRE"/></a></span><span class="dn show-mobile"><a class="AnchorLink" data-clubhouse-uid="s:600~t:6273" href="/futebol/time/_/id/6273/gremio" tabindex="0"><abbr data-clubhouse-uid="s:600~t:6273" style="text-decoration:none" title="Grêmio">GRE</abbr></a></span><span class="hide-mobile"><a class="AnchorLink" data-clubhouse-uid="s:600~t:6273" href="/futebol/time/_/id/6273/gremio" tabindex="0">Grêmio</a></span></div></td></tr><tr class="Table__TR Table__TR--sm Table__even" data-idx="2"><td class="Table__TD"><div class="team-link flex items-center clr-gray-03"><span class="team-position ml2 pr3">3</span><span class="pr4 TeamLink__Logo"><a class="AnchorLink" data-clubhouse-uid="s:600~t:7632" href="/futebol/time/_/id/7632/atletico-mg" tabindex="0"><img alt="CAM" class="Image Logo Logo__sm" data-mptype="image" src="data:image/gif;base64,R0lGODlhAQABAIAAAAAAAP///yH5BAEAAAAALAAAAAABAAEAAAIBRAA7" title="CAM"/></a></span><span class="dn show-mobile"><a class="AnchorLink" data-clubhouse-uid="s:600~t:7632" href="/futebol/time/_/id/7632/atletico-mg" tabindex="0"><abbr data-clubhouse-uid="s:600~t:7632" style="text-decoration:none" title="Atlético-MG">CAM</abbr></a></span><span class="hide-mobile"><a class="AnchorLink" data-clubhouse-uid="s:600~t:7632" href="/futebol/time/_/id/7632/atletico-mg" tabindex="0">Atlético-MG</a></span></div></td></tr><tr class="filled Table__TR Table__TR--sm Table__even" data-idx="3"><td class="Table__TD"><div class="team-link flex items-center clr-gray-03"><span class="team-position ml2 pr3">4</span><span class="pr4 TeamLink__Logo"><a class="AnchorLink" data-clubhouse-uid="s:600~t:819" href="/futebol/time/_/id/819/flamengo" tabindex="0"><img alt="FLA" class="Image Logo Logo__sm" data-mptype="image" src="data:image/gif;base64,R0lGODlhAQABAIAAAAAAAP///yH5BAEAAAAALAAAAAABAAEAAAIBRAA7" title="FLA"/></a></span><span class="dn show-mobile"><a class="AnchorLink" data-clubhouse-uid="s:600~t:819" href="/futebol/time/_/id/819/flamengo" tabindex="0"><abbr data-clubhouse-uid="s:600~t:819" style="text-decoration:none" title="Flamengo">FLA</abbr></a></span><span class="hide-mobile"><a class="AnchorLink" data-clubhouse-uid="s:600~t:819" href="/futebol/time/_/id/819/flamengo" tabindex="0">Flamengo</a></span></div></td></tr><tr class="Table__TR Table__TR--sm Table__even" data-idx="4"><td class="Table__TD"><div class="team-link flex items-center clr-gray-03"><span class="team-position ml2 pr3">5</span><span class="pr4 TeamLink__Logo"><a class="AnchorLink" data-clubhouse-uid="s:600~t:6086" href="/futebol/time/_/id/6086/botafogo" tabindex="0"><img alt="BOT" class="Image Logo Logo__sm" data-mptype="image" src="data:image/gif;base64,R0lGODlhAQABAIAAAAAAAP///yH5BAEAAAAALAAAAAABAAEAAAIBRAA7" title="BOT"/></a></span><span class="dn show-mobile"><a class="AnchorLink" data-clubhouse-uid="s:600~t:6086" href="/futebol/time/_/id/6086/botafogo" tabindex="0"><abbr data-clubhouse-uid="s:600~t:6086" style="text-decoration:none" title="Botafogo">BOT</abbr></a></span><span class="hide-mobile"><a class="AnchorLink" data-clubhouse-uid="s:600~t:6086" href="/futebol/time/_/id/6086/botafogo" tabindex="0">Botafogo</a></span></div></td></tr><tr class="filled Table__TR Table__TR--sm Table__even" data-idx="5"><td class="Table__TD"><div class="team-link flex items-center clr-gray-03"><span class="team-position ml2 pr3">6</span><span
```

```
class="pr4 TeamLink__Logo"><a class="AnchorLink" data-clubhouse-uid="s:600~t:6079"
href="/futebol/time/_/id/6079/red-bull-bragantino" tabindex="0"><img alt="BRA" cla
ss="Image Logo Logo__sm" data-mptype="image" src="data:image/gif;base64,R0lGODlhAQ
ABAIAAAAAAAP///yH5BAEAAAAALAAAAAABAAEAAAIBRAA7" title="BRA"/></a></span><span clas
s="dn show-mobile"><a class="AnchorLink" data-clubhouse-uid="s:600~t:6079" href="/
futebol/time/_/id/6079/red-bull-bragantino" tabindex="0"><abbr data-clubhouse-uid
="s:600~t:6079" style="text-decoration:none" title="Red Bull Bragantino">BRA</abbr
></a></span><span class="hide-mobile"><a class="AnchorLink" data-clubhouse-uid="s:
600~t:6079" href="/futebol/time/_/id/6079/red-bull-bragantino" tabindex="0">Red Bu
ll Bragantino</a></span></div></td></tr><tr class="Table__TR Table__TR--sm Table__
even" data-idx="6"><td class="Table__TD"><div class="team-link flex items-center c
lr-gray-03"><span class="team-position ml2 pr3">7</span><span class="pr4 TeamLink_
_Logo"><a class="AnchorLink" data-clubhouse-uid="s:600~t:3445" href="/futebol/tim
e/_/id/3445/fluminense" tabindex="0"><img alt="FLU" class="Image Logo Logo__sm" da
ta-mptype="image" src="data:image/gif;base64,R0lGODlhAQABAIAAAAAAAP///yH5BAEAAAAL
AAAAAABAAEAAAIBRAA7" title="FLU"/></a></span><span class="dn show-mobile"><a class
="AnchorLink" data-clubhouse-uid="s:600~t:3445" href="/futebol/time/_/id/3445/flum
inense" tabindex="0"><abbr data-clubhouse-uid="s:600~t:3445" style="text-decoratio
n:none" title="Fluminense">FLU</abbr></a></span><span class="hide-mobile"><a class
="AnchorLink" data-clubhouse-uid="s:600~t:3445" href="/futebol/time/_/id/3445/flum
inense" tabindex="0">Fluminense</a></span></div></td></tr><tr class="filled Table_
_TR Table__TR--sm Table__even" data-idx="7"><td class="Table__TD"><div class="team
-link flex items-center clr-gray-03"><span class="team-position ml2 pr3">8</span><
span class="pr4 TeamLink__Logo"><a class="AnchorLink" data-clubhouse-uid="s:600~t:
3458" href="/futebol/time/_/id/3458/athletico-pr" tabindex="0"><img alt="CAP" clas
s="Image Logo Logo__sm" data-mptype="image" src="data:image/gif;base64,R0lGODlhAQA
BAIAAAAAAAP///yH5BAEAAAAALAAAAAABAAEAAAIBRAA7" title="CAP"/></a></span><span class
="dn show-mobile"><a class="AnchorLink" data-clubhouse-uid="s:600~t:3458" href="/f
utebol/time/_/id/3458/athletico-pr" tabindex="0"><abbr data-clubhouse-uid="s:600~
t:3458" style="text-decoration:none" title="Athletico-PR">CAP</abbr></a></span><sp
an class="hide-mobile"><a class="AnchorLink" data-clubhouse-uid="s:600~t:3458" hre
f="/futebol/time/_/id/3458/athletico-pr" tabindex="0">Athletico-PR</a></span></div
></td></tr><tr class="Table__TR Table__TR--sm Table__even" data-idx="8"><td class
="Table__TD"><div class="team-link flex items-center clr-gray-03"><span class="tea
m-position ml2 pr3">9</span><span class="pr4 TeamLink__Logo"><a class="AnchorLink"
data-clubhouse-uid="s:600~t:1936" href="/futebol/time/_/id/1936/internacional" tab
index="0"><img alt="INT" class="Image Logo Logo__sm" data-mptype="image" src="dat
a:image/gif;base64,R0lGODlhAQABAIAAAAAAAP///yH5BAEAAAAALAAAAAABAAEAAAIBRAA7" title
="INT"/></a></span><span class="dn show-mobile"><a class="AnchorLink" data-clubhou
se-uid="s:600~t:1936" href="/futebol/time/_/id/1936/internacional" tabindex="0"><a
bbr data-clubhouse-uid="s:600~t:1936" style="text-decoration:none" title="Internac
ional">INT</abbr></a></span><span class="hide-mobile"><a class="AnchorLink" data-c
lubhouse-uid="s:600~t:1936" href="/futebol/time/_/id/1936/internacional" tabindex
="0">Internacional</a></span></div></td></tr><tr class="filled Table__TR Table__TR
--sm Table__even" data-idx="9"><td class="Table__TD"><div class="team-link flex it
ems-center clr-gray-03"><span class="team-position ml2 pr3">10</span><span class
="pr4 TeamLink__Logo"><a class="AnchorLink" data-clubhouse-uid="s:600~t:6272" href
="/futebol/time/_/id/6272/fortaleza" tabindex="0"><img alt="FOR" class="Image Logo
Logo__sm" data-mptype="image" src="data:image/gif;base64,R0lGODlhAQABAIAAAAAAAP///
yH5BAEAAAAALAAAAAABAAEAAAIBRAA7" title="FOR"/></a></span><span class="dn show-mobi
le"><a class="AnchorLink" data-clubhouse-uid="s:600~t:6272" href="/futebol/time/_/
id/6272/fortaleza" tabindex="0"><abbr data-clubhouse-uid="s:600~t:6272" style="tex
t-decoration:none" title="Fortaleza">FOR</abbr></a></span><span class="hide-mobil
e"><a class="AnchorLink" data-clubhouse-uid="s:600~t:6272" href="/futebol/time/_/i
d/6272/fortaleza" tabindex="0">Fortaleza</a></span></div></td></tr><tr class="Tabl
e__TR Table__TR--sm Table__even" data-idx="10"><td class="Table__TD"><div class="t
eam-link flex items-center clr-gray-03"><span class="team-position ml2 pr3">11</sp
an><span class="pr4 TeamLink__Logo"><a class="AnchorLink" data-clubhouse-uid="s:60
0~t:2026" href="/futebol/time/_/id/2026/sao-paulo" tabindex="0"><img alt="SAO" cla
ss="Image Logo Logo__sm" data-mptype="image" src="data:image/gif;base64,R0lGODlhAQ
ABAIAAAAAAAP///yH5BAEAAAAALAAAAAABAAEAAAIBRAA7" title="SAO"/></a></span><span clas
s="dn show-mobile"><a class="AnchorLink" data-clubhouse-uid="s:600~t:2026" href="/
futebol/time/_/id/2026/sao-paulo" tabindex="0"><abbr data-clubhouse-uid="s:600~t:2
026" style="text-decoration:none" title="São Paulo">SAO</abbr></a></span><span cla
```

ss="hide-mobile"><a class="AnchorLink" data-clubhouse-uid="s:600~t:2026" href="/fu
tebol/time/_/id/2026/sao-paulo" tabindex="0">São Paulo</a></span></div></td></tr><
tr class="filled Table__TR Table__TR--sm Table__even" data-idx="11"><td class="Tab
le__TD"><div class="team-link flex items-center clr-gray-03"><span class="team-pos
ition ml2 pr3">12</span><span class="pr4 TeamLink__Logo"><a class="AnchorLink" dat
a-clubhouse-uid="s:600~t:17313" href="/futebol/time/_/id/17313/cuiaba" tabindex
="0"><img alt="CUI" class="Image Logo Logo__sm" data-mptype="image" src="data:imag
e/gif;base64,R0lGODlhAQABAIAAAAAAAP///yH5BAEAAAAALAAAAAABAAEAAAIBRAA7" title="CU
I"/></a></span><span class="dn show-mobile"><a class="AnchorLink" data-clubhouse-u
id="s:600~t:17313" href="/futebol/time/_/id/17313/cuiaba" tabindex="0"><abbr data-
clubhouse-uid="s:600~t:17313" style="text-decoration:none" title="Cuiabá">CUI</abb
r></a></span><span class="hide-mobile"><a class="AnchorLink" data-clubhouse-uid
="s:600~t:17313" href="/futebol/time/_/id/17313/cuiaba" tabindex="0">Cuiabá</a></s
pan></div></td></tr><tr class="Table__TR Table__TR--sm Table__even" data-idx="12">
<td class="Table__TD"><div class="team-link flex items-center clr-gray-03"><span c
lass="team-position ml2 pr3">13</span><span class="pr4 TeamLink__Logo"><a class="A
nchorLink" data-clubhouse-uid="s:600~t:874" href="/futebol/time/_/id/874/corinthia
ns" tabindex="0"><img alt="COR" class="Image Logo Logo__sm" data-mptype="image" sr
c="data:image/gif;base64,R0lGODlhAQABAIAAAAAAAP///yH5BAEAAAAALAAAAAABAAEAAAIBRAA7"
title="COR"/></a></span><span class="dn show-mobile"><a class="AnchorLink" data-cl
ubhouse-uid="s:600~t:874" href="/futebol/time/_/id/874/corinthians" tabindex="0"><
abbr data-clubhouse-uid="s:600~t:874" style="text-decoration:none" title="Corinthi
ans">COR</abbr></a></span><span class="hide-mobile"><a class="AnchorLink" data-clu
bhouse-uid="s:600~t:874" href="/futebol/time/_/id/874/corinthians" tabindex="0">Co
rinthians</a></span></div></td></tr><tr class="filled Table__TR Table__TR--sm Tabl
e__even" data-idx="13"><td class="Table__TD"><div class="team-link flex items-cent
er clr-gray-03"><span class="team-position ml2 pr3">14</span><span class="pr4 Team
Link__Logo"><a class="AnchorLink" data-clubhouse-uid="s:600~t:2022" href="/futebo
l/time/_/id/2022/cruzeiro" tabindex="0"><img alt="CRU" class="Image Logo Logo__sm"
data-mptype="image" src="data:image/gif;base64,R0lGODlhAQABAIAAAAAAAP///yH5BAEAAAA
ALAAAAAABAAEAAAIBRAA7" title="CRU"/></a></span><span class="dn show-mobile"><a cla
ss="AnchorLink" data-clubhouse-uid="s:600~t:2022" href="/futebol/time/_/id/2022/cr
uzeiro" tabindex="0"><abbr data-clubhouse-uid="s:600~t:2022" style="text-decoratio
n:none" title="Cruzeiro">CRU</abbr></a></span><span class="hide-mobile"><a class
="AnchorLink" data-clubhouse-uid="s:600~t:2022" href="/futebol/time/_/id/2022/cruz
eiro" tabindex="0">Cruzeiro</a></span></div></td></tr><tr class="Table__TR Table__
TR--sm Table__even" data-idx="14"><td class="Table__TD"><div class="team-link flex
items-center clr-gray-03"><span class="team-position ml2 pr3">15</span><span class
="pr4 TeamLink__Logo"><a class="AnchorLink" data-clubhouse-uid="s:600~t:3454" href
="/futebol/time/_/id/3454/vasco-da-gama" tabindex="0"><img alt="VAS" class="Image
Logo Logo__sm" data-mptype="image" src="data:image/gif;base64,R0lGODlhAQABAIAAAAAA
AP///yH5BAEAAAAALAAAAAABAAEAAAIBRAA7" title="VAS"/></a></span><span class="dn show
-mobile"><a class="AnchorLink" data-clubhouse-uid="s:600~t:3454" href="/futebol/ti
me/_/id/3454/vasco-da-gama" tabindex="0"><abbr data-clubhouse-uid="s:600~t:3454" s
tyle="text-decoration:none" title="Vasco da Gama">VAS</abbr></a></span><span class
="hide-mobile"><a class="AnchorLink" data-clubhouse-uid="s:600~t:3454" href="/fute
bol/time/_/id/3454/vasco-da-gama" tabindex="0">Vasco da Gama</a></span></div></td>
</tr><tr class="filled Table__TR Table__TR--sm Table__even" data-idx="15"><td clas
s="Table__TD"><div class="team-link flex items-center clr-gray-03"><span class="te
am-position ml2 pr3">16</span><span class="pr4 TeamLink__Logo"><a class="AnchorLin
k" data-clubhouse-uid="s:600~t:9967" href="/futebol/time/_/id/9967/bahia" tabindex
="0"><img alt="BAH" class="Image Logo Logo__sm" data-mptype="image" src="data:imag
e/gif;base64,R0lGODlhAQABAIAAAAAAAP///yH5BAEAAAAALAAAAAABAAEAAAIBRAA7" title="BA
H"/></a></span><span class="dn show-mobile"><a class="AnchorLink" data-clubhouse-u
id="s:600~t:9967" href="/futebol/time/_/id/9967/bahia" tabindex="0"><abbr data-clu
bhouse-uid="s:600~t:9967" style="text-decoration:none" title="Bahia">BAH</abbr></a
></span><span class="hide-mobile"><a class="AnchorLink" data-clubhouse-uid="s:600~
t:9967" href="/futebol/time/_/id/9967/bahia" tabindex="0">Bahia</a></span></div></
td></tr><tr class="Table__TR Table__TR--sm Table__even" data-idx="16"><td class="T
able__TD"><div class="team-link flex items-center clr-gray-03"><span class="team-p
osition ml2 pr3">17</span><span class="pr4 TeamLink__Logo"><a class="AnchorLink" d
ata-clubhouse-uid="s:600~t:2674" href="/futebol/time/_/id/2674/santos" tabindex
="0"><img alt="SAN" class="Image Logo Logo__sm" data-mptype="image" src="data:imag
e/gif;base64,R0lGODlhAQABAIAAAAAAAP///yH5BAEAAAAALAAAAAABAAEAAAIBRAA7" title="SA

N"/></a></span><span class="dn show-mobile"><a class="AnchorLink" data-clubhouse-u
id="s:600~t:2674" href="/futebol/time/_/id/2674/santos" tabindex="0"><abbr data-cl
ubhouse-uid="s:600~t:2674" style="text-decoration:none" title="Santos">SAN</abbr>
</a></span><span class="hide-mobile"><a class="AnchorLink" data-clubhouse-uid="s:6
00~t:2674" href="/futebol/time/_/id/2674/santos" tabindex="0">Santos</a></span></d
iv></td></tr><tr class="filled Table__TR Table__TR--sm Table__even" data-idx="17">
<td class="Table__TD"><div class="team-link flex items-center clr-gray-03"><span c
lass="team-position ml2 pr3">18</span><span class="pr4 TeamLink__Logo"><a class="A
nchorLink" data-clubhouse-uid="s:600~t:3395" href="/futebol/time/_/id/3395/goias"
tabindex="0"><img alt="GOI" class="Image Logo Logo__sm" data-mptype="image" src="d
ata:image/gif;base64,R0lGODlhAQABAIAAAAAAAP///yH5BAEAAAAALAAAAAABAAEAAAIBRAA7" tit
le="GOI"/></a></span><span class="dn show-mobile"><a class="AnchorLink" data-clubh
ouse-uid="s:600~t:3395" href="/futebol/time/_/id/3395/goias" tabindex="0"><abbr da
ta-clubhouse-uid="s:600~t:3395" style="text-decoration:none" title="Goiás">GOI</ab
br></a></span><span class="hide-mobile"><a class="AnchorLink" data-clubhouse-uid
="s:600~t:3395" href="/futebol/time/_/id/3395/goias" tabindex="0">Goiás</a></span>
</div></td></tr><tr class="Table__TR Table__TR--sm Table__even" data-idx="18"><td
class="Table__TD"><div class="team-link flex items-center clr-gray-03"><span class
="team-position ml2 pr3">19</span><span class="pr4 TeamLink__Logo"><a class="Ancho
rLink" data-clubhouse-uid="s:600~t:3456" href="/futebol/time/_/id/3456/coritiba" t
abindex="0"><img alt="CFC" class="Image Logo Logo__sm" data-mptype="image" src="da
ta:image/gif;base64,R0lGODlhAQABAIAAAAAAAP///yH5BAEAAAAALAAAAAABAAEAAAIBRAA7" titl
e="CFC"/></a></span><span class="dn show-mobile"><a class="AnchorLink" data-clubho
use-uid="s:600~t:3456" href="/futebol/time/_/id/3456/coritiba" tabindex="0"><abbr
data-clubhouse-uid="s:600~t:3456" style="text-decoration:none" title="Coritiba">CF
C</abbr></a></span><span class="hide-mobile"><a class="AnchorLink" data-clubhouse-
uid="s:600~t:3456" href="/futebol/time/_/id/3456/coritiba" tabindex="0">Coritiba</
a></span></div></td></tr><tr class="filled Table__TR Table__TR--sm Table__even" da
ta-idx="19"><td class="Table__TD"><div class="team-link flex items-center clr-gray
-03"><span class="team-position ml2 pr3">20</span><span class="pr4 TeamLink__Log
o"><a class="AnchorLink" data-clubhouse-uid="s:600~t:6154" href="/futebol/time/_/i
d/6154/america-mg" tabindex="0"><img alt="AMG" class="Image Logo Logo__sm" data-mp
type="image" src="data:image/gif;base64,R0lGODlhAQABAIAAAAAAAP///yH5BAEAAAAALAAAA
ABAAEAAAIBRAA7" title="AMG"/></a></span><span class="dn show-mobile"><a class="Anc
horLink" data-clubhouse-uid="s:600~t:6154" href="/futebol/time/_/id/6154/america-m
g" tabindex="0"><abbr data-clubhouse-uid="s:600~t:6154" style="text-decoration:non
e" title="América-MG">AMG</abbr></a></span><span class="hide-mobile"><a class="Anc
horLink" data-clubhouse-uid="s:600~t:6154" href="/futebol/time/_/id/6154/america-m
g" tabindex="0">América-MG</a></span></div></td></tr></tbody></table>

- Criando DataFrame

```
In [76]: df = pd.DataFrame(columns=["Time"])
```

```
In [77]: df
```

Out[77]:

**Time**

- Fazendo a Inserção das Linhas

```
In [78]: for linha in tabela.tbody.find_all("tr"):
           for coluna in linha.find_all("td"):
             for div in coluna.find_all("div"):
               span = div.find("span", class_ ="hide-mobile")
               nometime = span.find_all("a")
             if (nometime != []):
               time = nometime[0].text.strip(" ")
               df = pd.concat([df, pd.DataFrame.from_records([{"Time": time}])])
```

In [79]: `df`

Out[79]:

| | Time |
|---|---|
| **0** | Palmeiras |
| **0** | Grêmio |
| **0** | Atlético-MG |
| **0** | Flamengo |
| **0** | Botafogo |
| **0** | Red Bull Bragantino |
| **0** | Fluminense |
| **0** | Athletico-PR |
| **0** | Internacional |
| **0** | Fortaleza |
| **0** | São Paulo |
| **0** | Cuiabá |
| **0** | Corinthians |
| **0** | Cruzeiro |
| **0** | Vasco da Gama |
| **0** | Bahia |
| **0** | Santos |
| **0** | Goiás |
| **0** | Coritiba |
| **0** | América-MG |

# Importando a segunda tabela

- Encontrando a tag da segunda tabela

In [80]: 
```
tabela2 = site.find("table", class_="Table Table--align-right")
```

In [81]: 
```
print(tabela2)
```

| J | V | E | D | GP | GC | SG | PTS |
|---|---|---|---|----|----|----|-----|
| 38 | 20 | 10 | 8 | 64 | 33 | +31 | 70 |
| 38 | 21 | 5 | 12 | 63 | 56 | +7 | 68 |
| 38 | 19 | 9 | 10 | 52 | 32 | +20 | 66 |
| 38 | 19 | 9 | 10 | 56 | 42 | +14 | 66 |
| 38 | 18 | 10 | 10 | 58 | | | |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| | | | | | 37 | +21 | 64 |
| 38 | 17 | 11 | 10 | 49 | 35 | +14 | 62 |
| 38 | 16 | 8 | 14 | 51 | 47 | +4 | 56 |
| 38 | 14 | 14 | 10 | 51 | 43 | +8 | 56 |
| 38 | 15 | 10 | 13 | 46 | 45 | +1 | 55 |
| 38 | 15 | 9 | 14 | 45 | 44 | +1 | 54 |
| 38 | 14 | 11 | 13 | 40 | 38 | +2 | 53 |
| 38 | 14 | 9 | 15 | 40 | 39 | +1 | 51 |
| 38 | 12 | 14 | 12 | 47 | 48 | -1 | 50 |
| 38 | 11 | 14 | 13 | 35 | 32 | +3 | 47 |

```
s="stat-cell">38</span></td><td class="Table__TD"><span class="stat-cell">12</span
></td><td class="Table__TD"><span class="stat-cell">9</span></td><td class="Table_
_TD"><span class="stat-cell">17</span></td><td class="Table__TD"><span class="stat
-cell">41</span></td><td class="Table__TD"><span class="stat-cell">51</span></td><
td class="Table__TD"><span class="stat-cell">-10</span></td><td class="Table__TD">
<span class="stat-cell">45</span></td></tr><tr class="filled Table__TR Table__TR--
sm Table__even" data-idx="15"><td class="Table__TD"><span class="stat-cell">38</sp
an></td><td class="Table__TD"><span class="stat-cell">12</span></td><td class="Tab
le__TD"><span class="stat-cell">8</span></td><td class="Table__TD"><span class="st
at-cell">18</span></td><td class="Table__TD"><span class="stat-cell">50</span></td
><td class="Table__TD"><span class="stat-cell">53</span></td><td class="Table__T
D"><span class="stat-cell">-3</span></td><td class="Table__TD"><span class="stat-c
ell">44</span></td></tr><tr class="Table__TR Table__TR--sm Table__even" data-idx
="16"><td class="Table__TD"><span class="stat-cell">38</span></td><td class="Table
__TD"><span class="stat-cell">11</span></td><td class="Table__TD"><span class="sta
t-cell">10</span></td><td class="Table__TD"><span class="stat-cell">17</span></td>
<td class="Table__TD"><span class="stat-cell">39</span></td><td class="Table__TD">
<span class="stat-cell">64</span></td><td class="Table__TD"><span class="stat-cel
l">-25</span></td><td class="Table__TD"><span class="stat-cell">43</span></td></tr
><tr class="filled Table__TR Table__TR--sm Table__even" data-idx="17"><td class="T
able__TD"><span class="stat-cell">38</span></td><td class="Table__TD"><span class
="stat-cell">9</span></td><td class="Table__TD"><span class="stat-cell">11</span>
</td><td class="Table__TD"><span class="stat-cell">18</span></td><td class="Table_
_TD"><span class="stat-cell">36</span></td><td class="Table__TD"><span class="stat
-cell">53</span></td><td class="Table__TD"><span class="stat-cell">-17</span></td>
<td class="Table__TD"><span class="stat-cell">38</span></td></tr><tr class="Table_
_TR Table__TR--sm Table__even" data-idx="18"><td class="Table__TD"><span class="st
at-cell">38</span></td><td class="Table__TD"><span class="stat-cell">8</span></td>
<td class="Table__TD"><span class="stat-cell">6</span></td><td class="Table__TD"><
span class="stat-cell">24</span></td><td class="Table__TD"><span class="stat-cel
l">41</span></td><td class="Table__TD"><span class="stat-cell">73</span></td><td c
lass="Table__TD"><span class="stat-cell">-32</span></td><td class="Table__TD"><spa
n class="stat-cell">30</span></td></tr><tr class="filled Table__TR Table__TR--sm T
able__even" data-idx="19"><td class="Table__TD"><span class="stat-cell">38</span>
</td><td class="Table__TD"><span class="stat-cell">5</span></td><td class="Table__
TD"><span class="stat-cell">9</span></td><td class="Table__TD"><span class="stat-c
ell">24</span></td><td class="Table__TD"><span class="stat-cell">42</span></td><td
class="Table__TD"><span class="stat-cell">81</span></td><td class="Table__TD"><spa
n class="stat-cell">-39</span></td><td class="Table__TD"><span class="stat-cell">2
4</span></td></tr></tbody></table>
```

- Criando outro DF para armazenar resultado

```python
In [82]: df2 = pd.DataFrame(columns=["Jogos", "Vitórias","Empates", "Derrotas",  "Gols Pró",
```

- Fazendo a inserção das linhas

```python
In [83]: for linha in tabela2.tbody.find_all("tr"):
    coluna = linha.find_all("td")
    if (coluna != []):
        jogos = coluna[0].text.strip(' ')
        vitorias = coluna[1].text.strip(' ')
        empates = coluna[2].text.strip(" ")
        derrotas = coluna[3].text.strip(" ")
        gols_pro = coluna[4].text.strip(" ")
        gols_contra = coluna[5].text.strip(" ")
        saldo_de_gols = coluna[6].text.strip(" ")
        pontos = coluna[7].text.strip(" ")
        df2 = pd.concat([df2, pd.DataFrame.from_records([{"Jogos": jogos ,"Vitórias": v
```

```
In [84]: df2
```

Out[84]:

| | Jogos | Vitórias | Empates | Derrotas | Gols Pró | Gols Contra | Saldo de Gols | Pontos |
|---|---|---|---|---|---|---|---|---|
| **0** | 38 | 20 | 10 | 8 | 64 | 33 | +31 | 70 |
| **0** | 38 | 21 | 5 | 12 | 63 | 56 | +7 | 68 |
| **0** | 38 | 19 | 9 | 10 | 52 | 32 | +20 | 66 |
| **0** | 38 | 19 | 9 | 10 | 56 | 42 | +14 | 66 |
| **0** | 38 | 18 | 10 | 10 | 58 | 37 | +21 | 64 |
| **0** | 38 | 17 | 11 | 10 | 49 | 35 | +14 | 62 |
| **0** | 38 | 16 | 8 | 14 | 51 | 47 | +4 | 56 |
| **0** | 38 | 14 | 14 | 10 | 51 | 43 | +8 | 56 |
| **0** | 38 | 15 | 10 | 13 | 46 | 45 | +1 | 55 |
| **0** | 38 | 15 | 9 | 14 | 45 | 44 | +1 | 54 |
| **0** | 38 | 14 | 11 | 13 | 40 | 38 | +2 | 53 |
| **0** | 38 | 14 | 9 | 15 | 40 | 39 | +1 | 51 |
| **0** | 38 | 12 | 14 | 12 | 47 | 48 | -1 | 50 |
| **0** | 38 | 11 | 14 | 13 | 35 | 32 | +3 | 47 |
| **0** | 38 | 12 | 9 | 17 | 41 | 51 | -10 | 45 |
| **0** | 38 | 12 | 8 | 18 | 50 | 53 | -3 | 44 |
| **0** | 38 | 11 | 10 | 17 | 39 | 64 | -25 | 43 |
| **0** | 38 | 9 | 11 | 18 | 36 | 53 | -17 | 38 |
| **0** | 38 | 8 | 6 | 24 | 41 | 73 | -32 | 30 |
| **0** | 38 | 5 | 9 | 24 | 42 | 81 | -39 | 24 |

# Concatenando as tabelas

```
In [85]: tabela_final = pd.concat([df,df2], axis=1)
         tabela_final
```

| | Time | Jogos | Vitórias | Empates | Derrotas | Gols Pró | Gols Contra | Saldo de Gols | Pontos |
|---|---|---|---|---|---|---|---|---|---|
| **0** | Palmeiras | 38 | 20 | 10 | 8 | 64 | 33 | +31 | 70 |
| **0** | Grêmio | 38 | 21 | 5 | 12 | 63 | 56 | +7 | 68 |
| **0** | Atlético-MG | 38 | 19 | 9 | 10 | 52 | 32 | +20 | 66 |
| **0** | Flamengo | 38 | 19 | 9 | 10 | 56 | 42 | +14 | 66 |
| **0** | Botafogo | 38 | 18 | 10 | 10 | 58 | 37 | +21 | 64 |
| **0** | Red Bull Bragantino | 38 | 17 | 11 | 10 | 49 | 35 | +14 | 62 |
| **0** | Fluminense | 38 | 16 | 8 | 14 | 51 | 47 | +4 | 56 |
| **0** | Athletico-PR | 38 | 14 | 14 | 10 | 51 | 43 | +8 | 56 |
| **0** | Internacional | 38 | 15 | 10 | 13 | 46 | 45 | +1 | 55 |
| **0** | Fortaleza | 38 | 15 | 9 | 14 | 45 | 44 | +1 | 54 |
| **0** | São Paulo | 38 | 14 | 11 | 13 | 40 | 38 | +2 | 53 |
| **0** | Cuiabá | 38 | 14 | 9 | 15 | 40 | 39 | +1 | 51 |
| **0** | Corinthians | 38 | 12 | 14 | 12 | 47 | 48 | -1 | 50 |
| **0** | Cruzeiro | 38 | 11 | 14 | 13 | 35 | 32 | +3 | 47 |
| **0** | Vasco da Gama | 38 | 12 | 9 | 17 | 41 | 51 | -10 | 45 |
| **0** | Bahia | 38 | 12 | 8 | 18 | 50 | 53 | -3 | 44 |
| **0** | Santos | 38 | 11 | 10 | 17 | 39 | 64 | -25 | 43 |
| **0** | Goiás | 38 | 9 | 11 | 18 | 36 | 53 | -17 | 38 |
| **0** | Coritiba | 38 | 8 | 6 | 24 | 41 | 73 | -32 | 30 |
| **0** | América-MG | 38 | 5 | 9 | 24 | 42 | 81 | -39 | 24 |

# Arrumando indices

- Começando por 1, e nao por 0

In [86]:
```python
tabela_final.index = np.arange(1, len(tabela_final)+1)
```

In [87]:
```python
tabela_final
```

| | Time | Jogos | Vitórias | Empates | Derrotas | Gols Pró | Gols Contra | Saldo de Gols | Pontos |
|---|---|---|---|---|---|---|---|---|---|
| 1 | Palmeiras | 38 | 20 | 10 | 8 | 64 | 33 | +31 | 70 |
| 2 | Grêmio | 38 | 21 | 5 | 12 | 63 | 56 | +7 | 68 |
| 3 | Atlético-MG | 38 | 19 | 9 | 10 | 52 | 32 | +20 | 66 |
| 4 | Flamengo | 38 | 19 | 9 | 10 | 56 | 42 | +14 | 66 |
| 5 | Botafogo | 38 | 18 | 10 | 10 | 58 | 37 | +21 | 64 |
| 6 | Red Bull Bragantino | 38 | 17 | 11 | 10 | 49 | 35 | +14 | 62 |
| 7 | Fluminense | 38 | 16 | 8 | 14 | 51 | 47 | +4 | 56 |
| 8 | Athletico-PR | 38 | 14 | 14 | 10 | 51 | 43 | +8 | 56 |
| 9 | Internacional | 38 | 15 | 10 | 13 | 46 | 45 | +1 | 55 |
| 10 | Fortaleza | 38 | 15 | 9 | 14 | 45 | 44 | +1 | 54 |
| 11 | São Paulo | 38 | 14 | 11 | 13 | 40 | 38 | +2 | 53 |
| 12 | Cuiabá | 38 | 14 | 9 | 15 | 40 | 39 | +1 | 51 |
| 13 | Corinthians | 38 | 12 | 14 | 12 | 47 | 48 | -1 | 50 |
| 14 | Cruzeiro | 38 | 11 | 14 | 13 | 35 | 32 | +3 | 47 |
| 15 | Vasco da Gama | 38 | 12 | 9 | 17 | 41 | 51 | -10 | 45 |
| 16 | Bahia | 38 | 12 | 8 | 18 | 50 | 53 | -3 | 44 |
| 17 | Santos | 38 | 11 | 10 | 17 | 39 | 64 | -25 | 43 |
| 18 | Goiás | 38 | 9 | 11 | 18 | 36 | 53 | -17 | 38 |
| 19 | Coritiba | 38 | 8 | 6 | 24 | 41 | 73 | -32 | 30 |
| 20 | América-MG | 38 | 5 | 9 | 24 | 42 | 81 | -39 | 24 |

# Ajustando Tipagem dos Dados

```python
tabela_final["Jogos"] = pd.to_numeric(tabela_final["Jogos"])
tabela_final["Vitórias"] = pd.to_numeric(tabela_final["Vitórias"])
tabela_final["Empates"] = pd.to_numeric(tabela_final["Empates"])
tabela_final["Derrotas"] = pd.to_numeric(tabela_final["Derrotas"])
tabela_final["Gols Pró"] = pd.to_numeric(tabela_final["Gols Pró"])
tabela_final["Gols Contra"] = pd.to_numeric(tabela_final["Gols Contra"])
tabela_final["Saldo de Gols"] = pd.to_numeric(tabela_final["Saldo de Gols"])
tabela_final["Pontos"] = pd.to_numeric(tabela_final["Pontos"])
tabela_final.dtypes
```

```
Out[88]:  Time              object
          Jogos              int64
          Vitórias           int64
          Empates            int64
          Derrotas           int64
          Gols Pró           int64
          Gols Contra        int64
          Saldo de Gols      int64
          Pontos             int64
          dtype: object
```

# Top 5

- Top 5 times com menos vitórias

```
In [89]:  tabela_final.sort_values(by="Vitórias", ascending= True).head(5)
```

Out[89]:

| | Time | Jogos | Vitórias | Empates | Derrotas | Gols Pró | Gols Contra | Saldo de Gols | Pontos |
|---|---|---|---|---|---|---|---|---|---|
| **20** | América-MG | 38 | 5 | 9 | 24 | 42 | 81 | -39 | 24 |
| **19** | Coritiba | 38 | 8 | 6 | 24 | 41 | 73 | -32 | 30 |
| **18** | Goiás | 38 | 9 | 11 | 18 | 36 | 53 | -17 | 38 |
| **17** | Santos | 38 | 11 | 10 | 17 | 39 | 64 | -25 | 43 |
| **14** | Cruzeiro | 38 | 11 | 14 | 13 | 35 | 32 | 3 | 47 |

- Top 5 times com mais empates

```
In [90]:  tabela_final.sort_values(by="Empates", ascending= False).head(5)
```

Out[90]:

| | Time | Jogos | Vitórias | Empates | Derrotas | Gols Pró | Gols Contra | Saldo de Gols | Pontos |
|---|---|---|---|---|---|---|---|---|---|
| **14** | Cruzeiro | 38 | 11 | 14 | 13 | 35 | 32 | 3 | 47 |
| **13** | Corinthians | 38 | 12 | 14 | 12 | 47 | 48 | -1 | 50 |
| **8** | Athletico-PR | 38 | 14 | 14 | 10 | 51 | 43 | 8 | 56 |
| **11** | São Paulo | 38 | 14 | 11 | 13 | 40 | 38 | 2 | 53 |
| **18** | Goiás | 38 | 9 | 11 | 18 | 36 | 53 | -17 | 38 |

- Top 5 times que mais fizeram gols

```
In [91]:  tabela_final.sort_values(by="Gols Pró", ascending= False).head(5)
```

| | Time | Jogos | Vitórias | Empates | Derrotas | Gols Pró | Gols Contra | Saldo de Gols | Pontos |
|---|---|---|---|---|---|---|---|---|---|
| **1** | Palmeiras | 38 | 20 | 10 | 8 | 64 | 33 | 31 | 70 |
| **2** | Grêmio | 38 | 21 | 5 | 12 | 63 | 56 | 7 | 68 |
| **5** | Botafogo | 38 | 18 | 10 | 10 | 58 | 37 | 21 | 64 |
| **4** | Flamengo | 38 | 19 | 9 | 10 | 56 | 42 | 14 | 66 |
| **3** | Atlético-MG | 38 | 19 | 9 | 10 | 52 | 32 | 20 | 66 |

- Top 5 times com mais saldo de gols

```python
tabela_final.sort_values(by="Saldo de Gols", ascending= False).head(5)
```

| | Time | Jogos | Vitórias | Empates | Derrotas | Gols Pró | Gols Contra | Saldo de Gols | Pontos |
|---|---|---|---|---|---|---|---|---|---|
| **1** | Palmeiras | 38 | 20 | 10 | 8 | 64 | 33 | 31 | 70 |
| **5** | Botafogo | 38 | 18 | 10 | 10 | 58 | 37 | 21 | 64 |
| **3** | Atlético-MG | 38 | 19 | 9 | 10 | 52 | 32 | 20 | 66 |
| **4** | Flamengo | 38 | 19 | 9 | 10 | 56 | 42 | 14 | 66 |
| **6** | Red Bull Bragantino | 38 | 17 | 11 | 10 | 49 | 35 | 14 | 62 |