

Luego de crear las funciones, cargue el módulo `miningscience` como `msc` e imprima documentación de la función.

In [1]:

```
# Escriba aquí su código para el ejercicio 1
```

```
import miningscience as msc
help(msc.download_pubmed)
help(msc.science_plots)
```

Ejercicio 2 [2 puntos]

Utilice dos veces la función `download_pubmed` para:

- Descargar la data, utilizando los keyword de su preferencia.
- Guardar el archivo descargado en la carpeta `data`.

Para cada corrida, imprima lo siguiente:

'El número artículos para KEYWORD es: XX' # Que se cargue con inserción de texto o valor que correspondea KEYWORD y XX

In [2]:

```
# Escriba aquí su código para el ejercicio 2
```

```
import miningscience as msc
import re

my_data2 = msc.download_pubmed("Buerger[Title/Abstract]")
contador = len(my_data2)
print("El número de artículos para el keyword (Buerger) es: ", contador)

with open("data/Buerger.txt", "w") as txt:
    txt.write(my_data2)

my_data3 = msc.download_pubmed("Gastritis[Title/Abstract]")
contador = len(my_data3)
print("El número de artículos para el keyword (Gastritis) es: ", contador)

with open("data/Gastritis.txt", "w") as txt:
    txt.write(my_data3)
```

Ejercicio 3 [1.5 puntos]


```
data_Frame = msc.science_plots(my_data2)
data_Frame1 = data_Frame.sort_values(by=['Numero de autores'], ascending=False)
data_Frame2 = data_Frame1.iloc[0:5]
data_Frame2
```

```
data_Frame_cdi = msc.science_plots(my_data3)
data_Frame_cdi1 = data_Frame_cdi.sort_values(by=['Numero de autores'],
ascending=False)
data_Frame_cdi2 = data_Frame_cdi1.iloc[0:5]
data_Frame_cdi2
```

```
import matplotlib.pyplot as plt.
```

```
labels = 'Japan', 'France', 'USA', 'Turkey', 'Italy'
```

```
sizes = [81, 25, 24, 21, 19]
```

```
fig1, x1 = plt.subplots(1)
```

```
x1.pie(sizes, labels=labels, autopct='%1.1f%%', startangle=10)
```

```
plt.title("Keyword-Burger")
```

```
plt.savefig("img/autoros-Burger.jpg")
```

```
import matplotlib.pyplot as plt.
```

```
labels = 'china', 'Japan', 'Italy', 'USA', 'Germany'
```

```
sizes = [3127, 1290, 905, 850, 252]
```

```
fig1, x1 = plt.subplots(1)
```

```
x1.pie(sizes, labels=labels, autopct='%1.1f%%', startangle=10)
```

```
plt.title("Keyword-Gostritos")
```

```
plt.savefig("img/autoros-Gostritos.jpg")
```


④

Gráfico de la búsqueda Keyword con Boerger se presenta que la mayor cantidad de autores son de Japon. con un porcentaje de 47.6%. seguido de Francia con un porcentaje de 14.7%.

Gráfico de la búsqueda Keyword con Gostitis se obtuvo que la mayor parte de autores son de China con un porcentaje de 48.7%. seguido de Japon con 20.4%. Por lo que se puede concluir que Japon es uno de los países que posee una gran cantidad de autores que han estudiado enfermedades como la Gostitis y la enfermedad de Boerger (Keyword).


```

from Bio.Phylo.TreeConstruction import DistanceTreeConstructor
from Bio.Phylo.TreeConstruction import DistanceCalculator
from Bio.Align.Applications import ClustalwCommandline
from Bio import AlignIO
from Bio import Phylo
from Bio import Entrez
from Bio import SeqIO
import Bio
import warnings
warnings.filterwarnings('ignore')
import os
import matplotlib
import matplotlib.pyplot as plt

with open("data/sequence_eulasa.seq", "r") as f:
    seq_eulasa = f.read()
    seq_eulasaList = seq_eulasa.split('\n')
    seq_eulasaListT = []
    count = 0
    for i in seq_eulasaList:
        if count < 15:
            seq_eulasaListT.append(i[:1])
            count += 1

Entrez.email = "guelapero.mois@ gmail.com"
ofile = open("data/sequence_eulasa.gb", 'w')
with Entrez.efetch(db="nucleotide", rettype="gb", retmode="text", id=seq_eulasaListT)
    as handle:
        for seq_record in SeqIO.parse(handle, "gb"):
            ofile.write(">" + str(seq_record.description[:50]) + '\n')
            ofile.write(str(seq_record.id) + '\n')
            ofile.write('\n')

ffile = open("data/sequence_eulasa.fasta", "w")
with open("data/sequence_eulasa.gb", 'r') as genbank:
    c = genbank.read()
    for line in c:
        ffile.write(str(line))

Clustalw_exe = r"C:\Program Files (x86)\Clustal W2\clustalw2.exe"
clustal_cline = ClustalwCommandline(clustalw_exe, infile="data/sequence_eulasa.fasta")
assert os.path.isfile(clustalw_exe), "Clustal-W executable is missing or not found"
stdout, stderr = Clustalw_cline()
ClustalAlign = AlignIO.read("data/sequence_eulasa.aln", "clustal")

```



```
calculator = DistanceCalculator('identity')
distance_matrix = calculator.get_distance(ClustalAlign)
constructor = DistanceTreeConstructor(calculator)
Data_tree = constructor.buid_tree(ClustalAlign)
Data_tree.rooted = True
Phylo.write(Data_tree, "data/sequence-enolasa.xml", "phylxml")
```

```
fig = plt.figure(figsize=(10, 25), dpi=200)
matplotlib.rc('font', size=12)
matplotlib.rc('xtick', labelsizex=10)
matplotlib.rc('ytick', labelsizex=10)
axes = fig.add_subplot(1, 1, 1)
Phylo.draw(Data_tree, axes=axes)
print("\t\t\t\t\t\t\t Fig 4. Árbol Filogenético enzima enolasa\n")
fig.savefig("img/sequence-enolasa.jpg")
```


Nombre [Apellido, Nombre]:

Construya las funciones del módulo miningscience.py

```
def download_pubmed(Keyword
```

): from Bio import Entrez.
import re.

Con esta función se buscara en pubmed los
articulos asociados con la palabra que se colocara como keyword, siendo esta la
referencia al tema del que se hara la busqueda de los articulos.

```
"""
```

```
Entrez.email = "gualapuro.moises@gmail.com"
```

```
handle = Entrez.esearch(db = "pubmed",  
                        term = keyword,  
                        usehistory = "y")
```

```
record = Entrez.read(handle)
```

```
id_list = record["IdList"]
```

```
webenv = record["WebEnv"]
```

```
query_key = record["QueryKey"]
```

```
handle = Entrez.efetch(db = "pubmed"
```

```
rettype = "medline",
```

```
retmode = "text",
```

```
retstart = 0,
```

```
retmax = 1500,
```

```
webenv = webenv,
```

```
query_key = query_key)
```

```
my_data = handle.read()
```

```
my_data2 = re.sub(r'\n\s{6}', ' ', my_data)
```

```
return(my_data2)
```


Nombre [Apellido, Nombre]:

```
def science_plots(my_data2):
```

```
):
```

```
import pandas as pd  
import csv  
import itertools
```

Con esta función se va a leer la fuente que se ha empleado en el keyword por la búsqueda en Pubmed, por medio del uso de expresiones regulares para delimitar cada posible variante de lectura, con el fin de que tome el país y el autor de cada registro. Además genero un dataframe relacionando los coordenados y el país con el uso de la data coordenados.txt / relacionando que luego se tendrá que hacer un arreglo para que se grafique solo los mismos o bandantes

```
AD = []
```

```
países1 = []
```

```
países2 = []
```

```
países3 = []
```

```
países4 = []
```

```
países5 = []
```

```
países6 = []
```

```
países7 = []
```

```
países8 = []
```

```
países9 = []
```

```
paísesT = []
```

```
for line in my_data2.splitlines():
```

```
if line.startswith("AD -"):
```

```
AD.append(line[:])
```

```
for line in my_data2.splitlines():
```

```
if line.startswith("AD -"):
```

```
AD = line[:]
```

```
país1 = re.findall(r'\s(\w{2,16})\s', AD)
```

```
países1.append(país1)
```

```
país2 = re.findall(r'\s(\w{2,16})\s', AD)
```

```
países2.append(país2)
```



```

pairs3 = re.findall(r'\s(\w{3,16}[\^o-q\,]\s\w{2,3}[\^o-q\,]\s\w{3,16}[\^o-q\,])\.', AD)

```

```

pairs3.append(pairs3)

```

```

pairs4 = re.findall(r'\s(\w{2,16})\.\s[a-z0-9\.-]+\@[a-z\.-]+\.[a-z\.-]{2,6}', AD)

```

```

pairs4.append(pairs4)

```

```

pairs5 = re.findall(r'\s(\w{2,16}[\^o-q\,]\s\w{2,16}[\^o-q\,])\.\s[a-z0-9\.-]+\@[a-z\.-]+\.[a-z\.-]{2,6}', AD)

```

```

pairs5.append(pairs5)

```

```

pairs6 = re.findall(r'\s(\w{3,16}[\^o-q\,]\s\w{2,3}[\^o-q\,]\s\w{3,16}[\^o-q\,])\.\s[a-z0-9\.-]+\@[a-z\.-]+\.[a-z\.-]{2,6}', AD)

```

```

pairs6.append(pairs6)

```

```

pairs7 = re.findall(r'\s(\w{2,16})\.\sElectronic address:\s[a-z0-9\.-]+\@[a-z\.-]+\.[a-z\.-]{2,6}', AD)

```

```

pairs7.append(pairs7)

```

```

pairs8 = re.findall(r'\s(\w{2,16}[\^o-q\,]\s\w{2,16}[\^o-q\,])\.\sElectronic address:\s[a-z0-9\.-]+\@[a-z\.-]+\.[a-z\.-]{2,6}', AD)

```

```

pairs8.append(pairs8)

```

```

pairs9 = re.findall(r'\s(\w{3,16}[\^o-q\,]\s\w{2,3}[\^o-q\,]\s\w{3,16}[\^o-q\,])\.\sElectronic address:\s[a-z0-9\.-]+\@[a-z\.-]+\.[a-z\.-]{2,6}', AD)

```

```

pairs9.append(pairs9)

```

```

pairsT = pairs1 + pairs2 + pairs3 + pairs4 + pairs5 + pairs6 + pairs7 + pairs8 + pairs9

```

```

pairsT = list(itertools.chain.from_iterable(pairsT))

```

```

len(pairsT)

```

```

unique_pairsT = list(set(pairsT))

```

```

unique_pairsT.sort()
len(unique_pairsT)

```

```

import csv

```

```

coordinates = {}

```

```

with open('coordinates.txt') as f:

```

```

    csvr = csv.DictReader(f)

```

```

    for row in csvr:

```

```

        coordinates[row['Name']] = [row['latitud'], row['longitud']]
        pairs.append(
            lon = [],
            lat = [],
            c = []
        )
        for z in unique_pairsT:
            if z in coordinates.keys():
                pairs.append(z)
                lat.append(float(coordinates[z][0]))
                lon.append(float(coordinates[z][1]))

```



```
C.append(paises T.count(Z))  
tablpaises T.count(Z)  
tablpaises = pd.DataFrame()
```

```
tablpaises ["pais"] = pais
```

```
tablpaises ["Numero de autores"] = C
```

```
return(tablpaises)
```