



EXAMEN PARCIAL PYTHON

GBI6-2021II: BIOINFORMÁTICA

Apellidos, Nombres <--- CAMBIE POR LOS QUE CORRESPONDA A SUS DATOS

03-08-2022

Color de texto

REQUERIMIENTOS PARA EL EXAMEN

Utilice de preferencia Jupyter de Anaconda, dado que tienen que hacer un control de cambios en cada pregunta.

Para este examen se requiere dos documentos:

1. Archivo `miningscience.py` donde tendrá dos funciones:
2. Archivo `2022I_GBI6_ExamenPython` donde se llamará las funciones y se obtendrá resultados.

Ejercicio 0 [0.5 puntos]

Realice cambios al cuaderno de jupyter:

- Agregue el logo de la Universidad
- Coloque sus datos personales
- Escriba una **tabla** con las características de su computador

```

```

```
**Nombre:** Jennyffer Elizalde
**Carrera:** Biotecnología
**Provincia:** Loja
**Etnia:** mestiza
```

```
import pandas as pd
Datos_c = pd.DataFrame({"Procesador": ["Intel(R)", "core (TM)", "2.80 GHz"], "RAM instalada": ["8.00 GB", "None", "None"], "Tipo de sistema": ["64 bits", "None", "None"]})
```

Ejercicio 1 [2 puntos]

Cree el archivo `miningscience.py` con las siguientes dos funciones:

i. `download_pubmed` : para descargar la data de PubMed utilizando el **ENTREZ** de Biopython. El parámetro de entrada para la función es el `keyword`.

ii. `science_plots` : la función debe

- utilizar como argumento de entrada la data descargada por `download_pubmed`
- ordenar los conteos de autores por país en orden ascendente y
- seleccionar los cinco más abundantes. Con esta selección debe graficar un `pie_plot`. Como guía para el conteo por países puede usar el ejemplo de [MapOfScience \(https://github.com/CSB-book/CSB/blob/master/regex/solutions/MapOfScience_solution.ipynb\)](https://github.com/CSB-book/CSB/blob/master/regex/solutions/MapOfScience_solution.ipynb).

iii Cree un `docstring` para cada función.

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

In [1]:

Escriba aquí su código para el ejercicio 1

```
import miningscience_gol 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 os
import re
x = msc.download_pubmed("Dogs")
y = len(x)
print("El número artículos para KEYWORD es: ", y)
with open("../Data/Dogs.txt", "w") as txt:
    txt.write(x)
j = msc.download_pubmed("virus")
k = len(j)
print("El número artículos para KEYWORD es: ", k)
with open("../Data/virus.txt", "w") as txt:
    txt.write(j)
```

Ejercicio 3 [1.5 puntos]

Utilice dos veces la función `science_plots` para:

- Visualizar un pie_plot para cada data descargada en el ejercicio 2.
- Guardar los pie_plot en la carpeta `img`

In [4]:

Escriba aquí su código para el ejercicio 3

```
Dogs = msc.science-plots(x)
with open("../img/Dogs.jpg", "w") as jpg:
    Dogs.
```

```
Virus = msc.science-plots(j)
with open("../img/DogsVirus.jpg", "w") as jpg:
    Virus.
```

Ejercicio 4 [1 punto]

Interprete los resultados de las figuras del ejercicio 3

Escriba la respuesta del ejercicio 5.

Los cinco países en los que se cuenta con más autores que han realizado investigaciones con la keyword "Dogs" son Japón, Italia, China, Brasil y Alemania.

Los cinco países en los que se cuenta con más autores que han realizado investigaciones con la keyword "Virus" son China, Francia, Alemania, Australia y Brasil.

Ejercicio 5 [2 puntos]

Para algún gen de las enzimas que intervienen en la ruta metabólica de la gluconeogenesis (Lista de genes por tipología (<https://www.genome.jp/pathway/map00010+C00068>)), realice lo siguiente:

1. Una búsqueda en la página del NCBI nucleotide (<https://www.ncbi.nlm.nih.gov/nucleotide/>).
2. Descargue el Accession List de su búsqueda y guarde en la carpeta `data`.
3. Cargue el Accession List en este notebook y haga una descarga de las secuencias de los quince primeros IDs de la accesión.
4. Arme un árbol filogenético para los resultados del paso 3.
5. Guarde su árbol filogenético en la carpeta `img`.
6. Interprete el árbol del paso 4.

In [3]:

Escriba aquí su código para el ejercicio 6

```

from Bio import Entrez
from Bio import Phylo
import matplotlib.pyplot as plt
from Bio.Phylo.TreeConstruction import DistanceTreeConstructor
from Bio import SeqIO
from Bio.Phylo.TreeConstruction import DistanceCalculator
from Bio.Align.Applications import ClustalwCommandline
import os

with open("../data/sequence.fasta") as file:
    text = file.read()
text = text.split('\n')
text = ''.join(text[1:5])
handle = Entrez.efetch(db="nucleotide", rettype="gb", retmode="text",
id = text)
print(handle.url)
records = SeqIO.parse("../data/sequence.fasta", "gb")
count = SeqIO.write(records, "../data/sequence.fasta", "fasta")

```

Escriba aquí la interpretación del árbol

Ejercicio 6 [1 punto]

1. Cree en GitHub un repositorio de nombre GBI6_ExamenPython.
2. Cree un archivo Readme.md que debe tener lo siguiente:
 - Datos personales
 - Características del computador
 - Versión de Python/Anaconda y de cada uno de los módulos/paquetes y utilizados
 - Explicación de la data utilizada
 - Un diagrama de procesos del módulo miningscience
3. Asegurarse que su repositorio tiene las carpetas data e img con los archivos que ha ido guardando en las preguntas anteriores.
4. Realice al menos 1 control de la versión (commits) por cada ejercicio (del 1 al 5), con un mensaje que inicie como:

Carlitos Alimaña ha realizado el ejercicio 1

Carlitos Alimaña ha realizado el ejercicio 2

...

```

Clustalw_exe = r"C:/program files (x86)/Clustalw2/clustalw2.exe"
Clustalw_cline = ClustalwCommandline(Clustalw_exe, infile = "../data/sequence.fasta")
assert os.path.isfile(Clustalw_exe), infile = "../data/sequence.fasta"

```

```

stdout, stderr = clustalw_cline()
print(stdout)
ClustalAlign = AlignIO.read(aln, "clustal")
Calculator = DistanceCalculator("identity")
distance_matrix = calculator.get_distance(alignment)

```

In []: Constructor = DistanceTreeConstructor(Calculator)

```

tree = constructor.build_tree(alignment)
tree = rooted = True
fig = plt.figure(figsize=(10, 12), dpi=200)
plt.rc("font", size=12)
plt.rc("xtick", labels=20)
plt.rc("ytick", labels=20)
axes = fig.add_subplot(1, 1, 1)
Phylo.draw(tree, axes=axes)
fig.savefig("../img/arbolfilogenético.jpg")

```

Descripción:

Inner 1 es antecesor de Xm-029014683.1 y Xm-039427779.1, Inner 2 es antecesor de Inner 1 y Xm-024718519.1, Inner 3 es antecesor de Inner 2 y Xm-040253125.1


```

1 # NOMBRE (Elizalde, Joseyffer):
2
3 from Bio import Entrez
4 from Bio import SeqIO
5 from Bio import GenBank
6 import re
7 import pandas as pd
8 import matplotlib.pyplot as plt
9 import csv as csv
10
11 def download_pubmed(keyword):
12     """
13     Esta función permite la búsqueda científica en la base de datos PubMed
14     """
15
16     Entrez.email = "A.N.Other@example.com"
17     busq = Entrez.read(Entrez.esearch(db="pubmed",
18                                     term=keyword,
19                                     usehistory="y"))
20     webenv = busq["webEnv"]
21     query_key = busq["QueryKey"]
22     handle = Entrez.efetch(db="pubmed",
23                           rettype="medline",
24                           retmode="text",
25                           retstart=0,
26                           retmax=543, webenv=webenv, query_keys=query_key)
27     data = handle.read()
28     dataexp = re.sub(r'\n{5}', '', data)
29     return dataexp
30
31
32 def science_plots(archivo):
33     """
34     Esta función permite ordenar ascendentemente el número de autores que realizaron un estudio científico de acuerdo a la búsqueda en la
35     función 1 y seleccionar los 5 más abundantes.
36     """
37
38     a = re.sub(r'^(?!(\w_)+)([^\w.-]+)\.\.[a-zA-Z]{1,4}$', '', archivo)
39     b = re.sub(r'^..[a-z_]*$', '', a)
40     c = re.sub(r'^..[a-z_]*$', '', b)
41     wc[c:].split('PMID-')
42
43     Palres1=[]
44     for PMID in c:
45         q=PMID.split('\n')
46         for fila in q:
47             w=fila.split(' ')
48             if w[0] == "AD":
49                 w=fila.split(',')
50                 Palres1.append(w[-1])
51
52     go0
53     Palres2 = [0]*len(Palres1)
54     for lis in Palres1:
55         bytes(lis,encoding="utf8")
56         if lis != '':
57             w=lis
58             if w[0] == ' ':
59                 w = re.sub(r'^\s','',w)
60             if w[-1] == '.':
61                 w = re.sub(r'\.$','',w)
62             u = re.sub(r'^\.$','',w)
63             u = re.sub(r'^\s$','',w)
64             Palres2[u]=w
65     print

```

64

65

```

Paises_completos=['Andorra','United Arab Emirates ','Afghanistan','Antigua and Barbuda','Anguilla','Albania','Armenia','Netherlands
Antilles','Angola','Antarctica','Argentina','American Samoa','Austria','Australia','Aruba','Azerbaijan','Bosnia and
Herzegovina','Barbados','Bangladesh','Belgium','Bhutan','Burkina Faso','Bulgaria','Bahrain','Burundi','Benin','Bermuda','Bhutan','Bolivia',
'Brazil','Bahamas','Bhutan','Bouvet Island','Botswana','Belarus','Belize','Canada','Cocos [Keeling] Islands','Congo [DRC]','Central
African Republic','Congo Republic','Switzerland','Côte d'Ivoire','Cook Islands','Chile','Cameroon','China','Colombia','Costa
Rica','Cuba','Cape Verde','Christmas Island','Cyprus','Czech Republic','Germany','Djibouti','Denmark','Dominica','Dominican
Republic','Algeria','Ecuador','Estonia','Egypt','Western Sahara','Eritrea','Spain','Ethiopia','Finland','Fiji','Falkland Islands [Islas
Malvinas]','Micronesia','Faroe Islands','France','Gabon','United Kingdom','Grenada','Georgia','French
Guiana','Guernsey','Ghana','Gibraltar','Greenland','Gambia','Guinea','Guadeloupe','Equatorial Guinea','Greece','South Georgia and the
South Sandwich Islands','Guatemala','Guam','Guinea-Bissau','Guyana','Caza Strip','Hong Kong','Heard Island and McDonald
Islands','Honduras','Croatia','Haiti','Hungary','Indonesia','Ireland','Israel','Isle of Man','India','British Indian Ocean
Territory','Iraq','Iran','Iceland','Italy','Jersey','Jamaica','Jordan',
'Japan','Kenya','Kyrgyzstan','Cambodia','Kiribati','Comoros','Saint Kitts and Nevis','North Korea','South Korea','Kuwait','Cayman
Islands','Kazakhstan','Laos','Lebanon','Saint Lucia','Liechtenstein','Sri Lanka','Liberia','Lesotho','Lithuania','Luxembourg','Latvia',
'Libya','Morocco','Macao','Moldova','Montenegro','Madagascar','Marshall Islands','Macedonia [FYROM]','Mali','Myanmar
[Burma]','Mongolia','Macau','Northern Mariana
Islands','Martinique','Mauritania','Mauritius','Malta','Mauritius','Maldives','Malawi','Mexico','Malaysia','Mozambique','Namibia','New
Caledonia','Niger','Newfolk Island','Nigeria','Nicaragua','The Netherlands','Norway','Nepal','Niue','New
Zealand','Oman','Panama','Peru','French Polynesia','Papua New Guinea','Philippines','Pakistan','Poland','Saint Pierre and Miquelon',
'Pitcairn Islands','Puerto Rico','Palestinian Territories','Portugal','Palau','Paraguay','Qatar','Reunion','Romania','Serbia','Seychelles',
'Swaziland','Saudi Arabia','Solomon Islands','Seychelles','Sudan','Sweden','Singapore','Saint Helena','Slovenia','Svalbard and Jan
Mayen','Slovakia','Sierra Leone','San Marino','Senegal','Somalia','Suriname','São Tomé and Príncipe','El Salvador','Syria','Swaziland',
'Turks and Caicos Islands','Chad','French Southern Territories','Togo','Thailand','Tajikistan','Tahiti','Timor-Leste','Turkmenistan',
'Tunisia','Tuvalu','Turkey','Trinidad and Tobago','Togo','Togo','Togo','Togo','Togo','Togo','Togo','Togo','Togo','Togo','Togo','Togo','Togo',
'United States of America','Uruguay','Uzbekistan','Vatican City','Saint Vincent and the Grenadines','Venezuela','British Virgin Islands','U.S.
Virgin Islands','Vietnam','Vanuatu','Wallis and Futuna','Samoa','Kosovo','Yemen','Mayotte','South Africa','Zambia','Zimbabwe']

```

66

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92

93

94

95

96

97

98

99

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

```

Paises3=Paises2
h=Paises_completos
f=len(h)
CountriesCount=[0]*f
k=0
for i in range(f):
    d=0
    for comp in Paises3:
        if comp == h[i]:
            d+=1
    CountriesCount[k]=d
    k+=1

Paises4=[]
CountriesCount=[0]
for i in range(CountriesCount):
    if CountriesCount[i] != 0:
        Counter.append(i)
        m=Paises_completos[i]
        Paises4.append(m)
    i+=1

Table1 = pd.DataFrame({'Country': Paises4,
                       'sum_auth': Counter})

Order=Table1.sort_values(by=['sum_auth'], ascending=False)
Taken = Order.iloc[0:5]
sum_auth=Taken['sum_auth'].sum()
len_Taken=len(Taken)
smpd.Series(1)
l = Taken.iloc[0,1]
sa = pd.Series(1)
s = sa.tolist()

prom = []
for number in s:
    x=(number/sum_auth)*100
    prom.append(x)
Table2 = pd.DataFrame({'Country': sa,
                       'Percent': prom})

fig1, ax1 = plt.subplots()
ax1.pie(prom, labels=s, autopct='%1.1f%%',
        shadow=True, startangle=90)
ax1.axis('equal')
plt.show()

return (Table2)

```