



EXAMEN PARCIAL PYTHON

GBI6-2021II: BIOINFORMÁTICA

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

03-08-2022

Jose Chedz Martinez

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

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 la función.

In [1]:

```
# Escriba aquí su código para el ejercicio 1
import miningscience_goi
import miningscience_goi as msc
help(msc.download_pubmed)
help(msc.science_plots)
```

Ejercicio 2 [2 puntos]

Utilice dos veces la función `download_pubmed` para:

- 2 Keyword

- 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
paper = msc.download_pubmed("Eumenice")
print("El número de artículos para Eumenice es:", len(paper))

artic = msc.download_pubmed("termit")
print("El número de artículos para termit es:", len(artic))
```

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`

escriba aquí su código para el ejercicio 3

```
Eumeninae = msc.science-plots (paper)
with open ("img/Eumeninae.png", "w") as jpg:
    Eumeninae
```

```
termite = msc.science-plots (article)
with open ("img/termite.png", "w") as jpg:
    termite
```

Ejercicio 4 [1 punto]

Interprete los resultados de las figuras del **ejercicio 3**

En el primer diagrama de pastel se pueden observar los cinco países con más publicaciones referenciadas a la especie *Eumeninae*, donde el porcentaje más alto de artículos, con investigaciones respecto a avispa atrevida, ha sido China con 41,5%, seguido de Alemania con 29,2% y en tercer lugar Brasil con el 15,4%.

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. *- Solo 15*
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
with open('Data/sequence.seq') as file:
    text = file.read()
text = text.split('\n')
handle = Entrez.efetch(db="nucleotide", rettype="gb", retmode="text", id=text)
records = SeqIO.parse("Data/sequence.gb", "genbank")
clustalw_exe = r"C:/Program Files (x86)/Bioinformatics/Clustalw.exe"
clustalw_cli = clustalwcommandline(clustalw_exe, infile="Data/sequence.fasta")
assert os.path.isfile(clustalw_exe), "clustalw executable is missing or not found"
status, stderr = clustalw_cli.run()
clustal_align = AlignIO.read("Data/sequence.aln", clustal)
with open("Data/sequence.aln", "w") as aln:
    alignment = AlignIO.write(aln, "clustal")

calculator = DistanceCalculator('identity')
distance_matrix = calculator.get_distance(alignment)
tree_constructor = DistanceTreeConstructor(calculator)
tree = tree_constructor.build_tree(alignment)
tree.rooted = True
fig = plt.figure(figsize=(10, 12), dpi=200)
matplotlib.rcParams['font.size'] = 12
matplotlib.rcParams['xtick.labelsize'] = 20
matplotlib.rcParams['ytick.labelsize'] = 20
axes = fig.add_subplot(1, 1, 1)
Phylo.draw(tree, axes=axes)
fig.savefig("img/arbol.jpg")
```

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

...

In []:

Nombre [Apellido, Nombre]:

Construya las funciones del módulo miningscience.py

```
def download_pubmed(Keyword):  
    """  
    Este comando me ayudara buscar articulos en pubmed a través de una palabra  
    clave """  
    Entrez.email = "guadalupe.morales@gmail.com"  
    busq = Entrez.read(Entrez.esearch(db="pubmed", term=Keyword, usehistory="y"))  
    webenv = busq["webenv"]  
    query_key = busq["QueryKey"]  
    handle = Entrez.efetch(db="pubmed", rettype="medline", retmode="text", retstart=0,  
                           retmax=543, webenv=webenv, query_key=query_key)  
    data = handle.read()  
    data_exp = re.sub(r'\n\s{6}', '\n', data)
```

Nombre [Apellido, Nombre]:

```
def science_plots( doc
```

Este función me ayuda a contar los países de autores que producen artículos relacionados a la palabra clave que busq'o

```
Correc = re.sub(r'\s[%\~\+]\s+[\w-]+\s+[a-zA-Z]{1,4}', '', doc)
```

```
Coord = re.sub(r'\s\d\d\d\s+', '', Correc)
```

```
numb = re.sub(r'\s\d\d\d\s+', '', Coord)
```

```
x = numb[1:] . split('PHID-')
```

```
Countries_A = []
```

```
For PHID in x:
```

```
q = PHID . split('\n')
```

```
For fila in q:
```

```
w = fila . split(' ')
```

```
if w[0] == "AD":
```

```
e = fila . split(' ', 1)
```

```
Countries_A.append(e[-1])
```

```
a=0
```

```
Countries_B = copy.copy(Countries_A)
```

```
For ls in Countries_A:
```

```
bytes(ls, encoding = "utf8")
```

```
if ls != '':
```

```
w = ls
```

```
if w[0] == '':
```

```
w = re.sub(r'^\s+', '', w)
```

```
if w[-1] == '':
```

```
w = re.sub(r'\s$', '', w)
```

```
w = re.sub(r'\sB', '', w)
```

```
w = re.sub(r'\sB', '', w)
```

```
Countries_B[a] = w
```

```
a = a + 1
```

```
Countries_all = [ todos los países ]
```

```
Countries - C = Countries - B
```

```
k = Countries - all
```

```
f = len(h)
```

```
Countries Count = [0] * f
```

```
k = 0
```

```
for elem in h:
```

```
    d = 0
```

```
    for comp in Countries - C:
```

```
        if elem == str(comp):
```

```
            Go  
            Countries Count [k] = d
```

```
            k = k + 1
```

```
Countries - D = []
```

```
Counter = []
```

```
o = 0
```

```
For line in Countries Count
```

```
    If str(line) != '0':
```

```
        Counter.append(line)
```

```
        m = Countries - all [0]
```

```
        Countries - D.append(m)
```

```
        o = o + 1
```

```
table - A = pd.DataFrame({'Country': Countries - D, 'num_auth': Counter})  
Order = table - A.sort_values(by = 'num_auth', ascending = [False])
```

```
taken = Order.iloc[0:5]
```

```
soma = taken['num_auth'].sum()
```

```
io = taken.iloc[:, 0]
```

```
so = pd.Series(io)
```

```
i = taken.iloc[:, 3]
```

```
sa = pd.Series(i)
```

```
s = sa.tolist()
```

```
prom = []
```

```
for number in s:
```

```
    xa = (number / soma) * 100
```

```
    prom.append(xa)
```

```
table - B = pd.DataFrame({'Country': so, 'percent': prom})
```

```
fig, ax1 = plt.subplots(1)
```

```
ax1.pie(prom, labels = s, autopct = '%1.1f%%',  
        shadow = True, startangle = 90)
```

```
ax1.axis('equal')
```

```
plt.show()
```

```
plt.savefig('img/Grafica de premissas.jpg', dpi = 300) se return (table - B)
```