



# EXAMEN PARCIAL PYTHON

## GBI6-2021II: BIOINFORMÁTICA

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

Calva Moreno Ibeth Calva  
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

### 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 as msc
help(download_pubmed)
help(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
journal = download_pubmed("Ecuador SARS-CoV-2")
journal1 = (download_pubmed("Ecuador shrimp"))
print("El número artículos para Ecuador SARS-CoV-2 es: ", len(journal))
print("El número artículos para Ecuador shrimp: ", len(journal1))
```

## 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`



n [4]:

# Escriba aquí su código para el ejercicio 3

```
science-plots("Ecuador SARS-CoV-2")
science-plots("Ecuador shrimp")
```

## Ejercicio 4 [1 punto]

Interprete los resultados de las figuras del ejercicio 3

Escriba la respuesta del ejercicio 5.

- En la primera búsqueda se obtuvieron 126 artículos de SARS-CoV2, es un número alto dado el descubrimiento del virus a finales del 2019. Sin embargo, es el esperado debido a los años de pandemia que soportó la población mundial. En el gráfico se puede observar que por "Keyword" elegido el país que presente un mayor número de frecuencias fue Ecuador, los siguientes países fueron UK, Spain y Usa, los motivos de la frecuencia puede deberse a múltiples factores, tales como la producción de fármacos en UK y USA que hicieron que realicen investigaciones en colaboración en Ecuador y puede el intercambio de datos para identificar posibles brotes.
- En la segunda búsqueda se obtuvieron 37 artículos de "Ecuador shrimp", siendo un número bajo dado la importancia del camarón en la industria productora del Ecuador. Dado el "keyword" se espera que el país que presente mayor frecuencia sea Ecuador. El segundo país fue USA que podría deberse a que es uno de los más importantes importadores del camarón ecuatoriano. Finalmente los países latinos restantes, tales como México, Colombia y Chile puede deberse a la cooperación científica que existe en áreas investigativas de interés.

## 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.

```
from Bio import Phylo
from Bio import SeqIO
from Bio import AlignIO
from Bio.Phylo.TreeConstruction import Distance Calculator
from Bio.Phylo.TreeConstruction import Distance TreeConstructor
from Bio import Entrez
import re
import os
from Bio.Align.Applications import ClustalwCommandline

with open("sequence.seq") as f:
    data = f.readlines()[0:15] ## Seleccionamos las primeras 15 secuencias
out_sequence = open("secuencias.fasta", "w")
for linea in data:
    Entrez_email = "ibeth.calva@est.ikiam.edu.ec"
    handle = Entrez.efetch(db="nucleotide", id=linea, rettype="fasta", retmode="text")
    data = (handle.read())
    out_sequence.write(data)
out_sequence.close()
localhost:8888/notebooks/GDrive/IKIAM/CLASES/2022I/20221_GBI6/20221_GBI6_Examen_Python/20221_GBI6G01_ExamenPython.ipynb
```



In [3]:

# Escriba aquí su código para el ejercicio 6

```

Crustalw_exe = r"C:\Program Files (x86)\CrustalW2\crustalw2.exe"
Crustalw_cline = CrustalwCommandline(crustalw_exe, infile = "secuencias.fasta")
assert os.path.isfile(crustalw_exe), "crustalw executable is missing or not found"
stdout, stderr = crustalw_cline()
Print(crustalw_cline)
CrustalAlign = AlignIO.read("secuencias.aln", "crustal")
print(CrustalAlign)

from Bio import Phylo
tree = Phylo.read("secuencias.dnd", "newick")

with open("secuencias.dn", "r") as aln:
    alignment = AlignIO.read(aln, "crustal")

from Bio.Phylo.TreeConstruction import DistanceCalculator
calculator = DistanceCalculator('identity')
distance_matrix = calculator.get_distance(alignment)
from Bio.Phylo.TreeConstruction import DistanceTreeConstructor
constructor = DistanceTreeConstructor(calculator)
# construir el árbol
glut_1 = constructor.build_tree(alignment)
glut_1.rooted = True
Phylo.write(glut_1, "tree-glut-1.xml", "phyloxml")
glut_1 = Phylo.read(file="tree-glut-1.xml", format="phyloxml")

# Arbol elemental en Matplotlib
import matplotlib
import matplotlib.pyplot as plt
# fig = Phylo.draw(aln_tree)
fig = plt.figure(figsize=(30,40), dpi=100) # create figure y set the size
matplotlib.rcParams['font', size=30] # font size of the leaf and node labels
matplotlib.rcParams['xtick', labels=20] # font size of the tick labels
matplotlib.rcParams['ytick', labels=20] # font size of the tick labels
axes = fig.add_subplot(1, 1, 1)
Phylo.draw(glut_1, axes=axes)
fig.savefig("img/tree-glut-1.jpg")

```

Escriba aquí la interpretación del árbol

El árbol filogenético se puede distinguir a 13 clados definidos del gen GLUT-1 que codifica a la glucosa transporter protein. Las secuencias con el código NM\_174602.2, NM\_001103687.1 fueron las que presentaron mayor similitud y compartieron 9 clados, pudiendo indicar que ambas secuencias provienen de un mismo organismo.

## 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]:

```
import Bio
from Bio.Seq import Seq
from Bio import Entrez
import re
```

Construya las funciones del módulo miningscience.py

```
def download_pubmed(Keyword):
```

```
):
```

Función que entrada pide la frase clave de búsqueda y en output guarda un archivo que contiene los resultados de la búsqueda, considerando los títulos.

```
    """
    Entrez.email = "ibeth.calva@est.ikiam.edu.ec"
    handle = Entrez.esearch(db="pubmed",
                           term=Keyword+"[Title]",
                           retmax=1000,
                           usehistory="y")

    record = Entrez.read(handle)
    id_list = record["Idlist"]
    webenv = record["Webenv"]
    query_key = record["Query Key"]
    handle = Entrez.efetch(db="pubmed",
                           rettype="medline",
                           retmode="text",
                           webenv=webenv,
                           query_key=query_key)
    out_handle = open("data/"+Keyword, "w")
    data = handle.read()
    (id_list)
    handle.close()
    out_handle.write(data)
    out_handle.close()
    return id_list
```

```
import re
import csv
import matplotlib.pyplot as plt
from collections import Counter
```

```
def science_plots(data):
```

Función que pide como entrada la búsqueda de la función `download_pubmeds` y como resultado muestra un gráfico de pastel indicando a los cinco países de origen de autores que aparecieron más veces.

```
    """
    with open("data/"+data, errors="ignore") as f:
        texto = f.read()
        texto = re.sub(r"\n\s{6}", " ", texto)
        countries_1 = re.findall(r"AD\s{2}\s[A-Za-z].*,\s([A-Za-z]*)\s", texto)
        unique_countries = list(set(countries_1))
        conteo = Counter(countries_1)
        resultado = {} ## creamos un diccionario vacío
        ## Bucle agregamos los valores del diccionario
        for clave in conteo:
            valor = conteo[clave]
            if valor > 1:
                resultado[clave] = valor
        ## Ordenamos de forma ascendente
        ## ordenamos de forma ascendente
        sort = (sorted(resultado.values()))
        sort.sort(reverse=True)
        import operator
        ## creamos dos listas vacías
        país = []
        frecuencia = []
        ## bucle que añade los valores país y frecuencia a la listas vacías país y frecuencia
        reverse = sorted(resultado.items(), key=operator.itemgetter(1), reverse=True)
        for name in enumerate(reverse):
            país.append(name[1][0])
            frecuencia.append(resultado[name[1][0]])
        max_país = país[0:5] ## Seleccionamos los cinco primeros países
        max_frec = frecuencia[0:5] ## Seleccionamos los cinco primeros frecuencia respecto a los países
        fig = plt.figure(figsize=(10, 8))
        plt.pie(max_frec, labels=max_país)
        (plt.savefig("img/"+data, dpi=100, bbox_inches='tight'))
        plt.show()
```