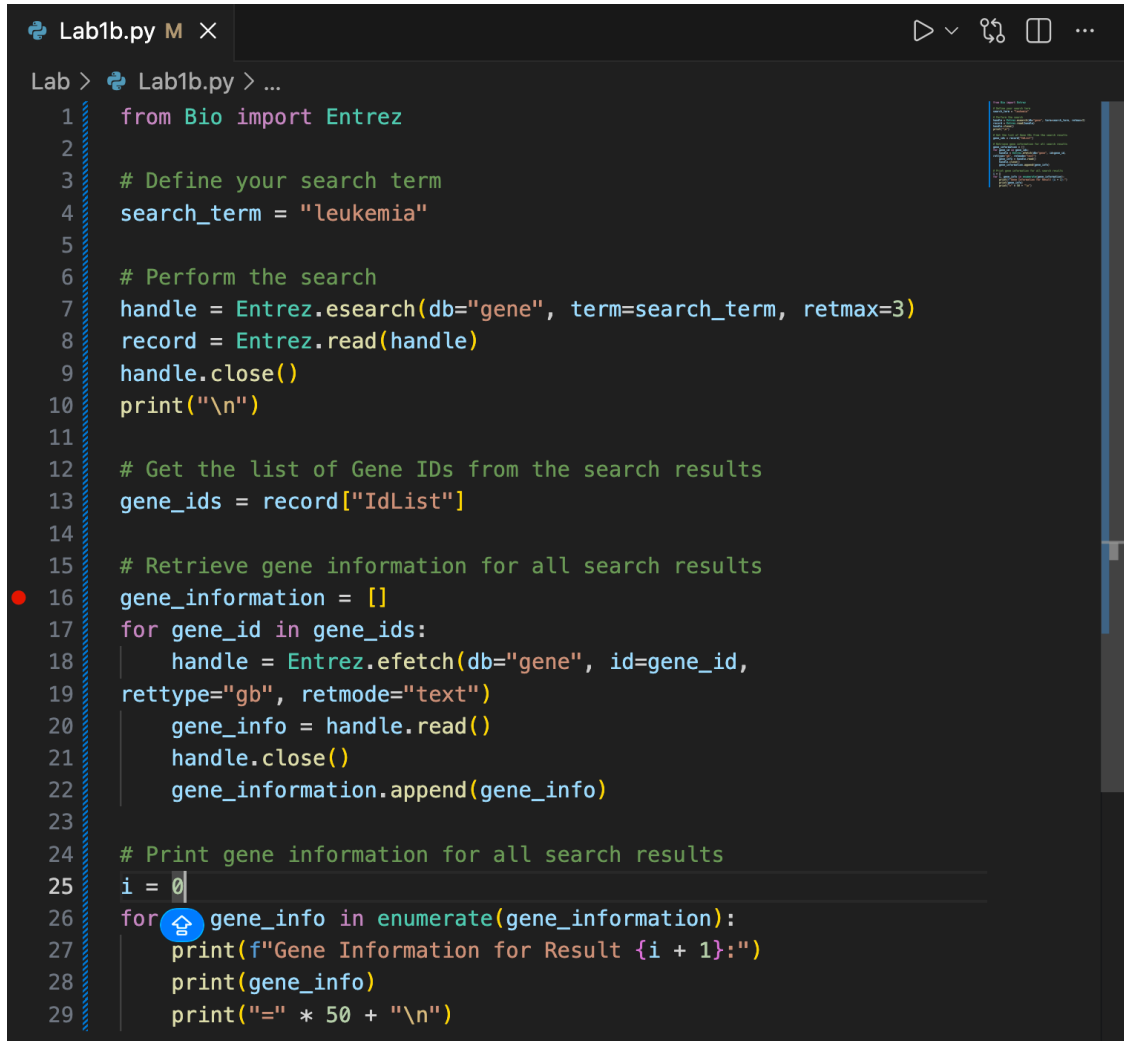


PART B

1.

- a. I have changed the `search_term` to “leukemia” and then to “cystic fibrosis” as the topic of interest.



```
Lab1b.py M X
Lab > Lab1b.py > ...
1  from Bio import Entrez
2
3  # Define your search term
4  search_term = "leukemia"
5
6  # Perform the search
7  handle = Entrez.esearch(db="gene", term=search_term, retmax=3)
8  record = Entrez.read(handle)
9  handle.close()
10 print("\n")
11
12 # Get the list of Gene IDs from the search results
13 gene_ids = record["IdList"]
14
15 # Retrieve gene information for all search results
16 gene_information = []
17 for gene_id in gene_ids:
18     handle = Entrez.efetch(db="gene", id=gene_id,
19 rettype="gb", retmode="text")
20     gene_info = handle.read()
21     handle.close()
22     gene_information.append(gene_info)
23
24 # Print gene information for all search results
25 i = 0
26 for gene_info in enumerate(gene_information):
27     print(f"Gene Information for Result {i + 1}:")
28     print(gene_info)
29     print("=" * 50 + "\n")
```

- b. After repeating step 1 twice, I found that different diseases (**search_term**) return different genes in different species. Each search focuses on genes related to that specific disease condition. Another interesting observation is the genes returned for **search_term = "cystic fibrosis"**: EGFR from humans do has some connection to leukemia (the first **search_term**).

Gene Information for Result 1:

1. LOC144531885
leukemia inhibitory factor receptor-like [Sander vitreus (walleye)]
Other Designations: leukemia inhibitory factor receptor-like
Chromosome: 16
Annotation: Chromosome 16 NC_135870.1 (28581425..28602849)
ID: 144531885

Gene Information for Result 2:

1. pbx4
pre-B-cell leukemia transcription factor 4 [Sander vitreus (walleye)]
Other Designations: pre-B-cell leukemia transcription factor 4
Chromosome: 15
Annotation: Chromosome 15 NC_135869.1 (16032391..16060199, complement)
ID: 144530797

Gene Information for Result 3:

1. pbxip1a
pre-B-cell leukemia homeobox interacting protein 1a [Sander vitreus (walleye)]
Other Designations: pre-B-cell leukemia transcription factor-interacting protein 1
Chromosome: 14
Annotation: Chromosome 14 NC_135868.1 (11971962..11982191)
ID: 144528519

Gene Information for Result 1:

1. cftr
CF transmembrane conductance regulator [Sander vitreus (walleye)]
Other Designations: cystic fibrosis transmembrane conductance regulator
Chromosome: 8
Annotation: Chromosome 8 NC_135862.1 (5473012..5522248, complement)
ID: 144521857

Gene Information for Result 2:

1. cftr
CF transmembrane conductance regulator [Mustelus asterias (starry smooth-hound)]
Other Designations: cystic fibrosis transmembrane conductance regulator
Chromosome: 9
Annotation: Chromosome 9 NC_135809.1 (27249470..27379241, complement)
ID: 144498565

Gene Information for Result 3:

1. EGFR
Official Symbol: EGFR and Name: epidermal growth factor receptor [Homo sapiens (human)]
Other Aliases: ERBB, ERBB1, ERBP, HER1, NISBD2, NNCIS, PIG61, mENA
Other Designations: epidermal growth factor receptor; EGFR vIII; avian erythroblastic leukemia viral (v-erb-b) oncogene homolog; cell growth inhibiting protein 40; cell proliferation-inducing protein 61; epidermal growth factor receptor tyrosine kinase domain; erb-b2 receptor tyrosine kinase 1; proto-oncogene c-ErbB-1; receptor tyrosine-protein kinase erbB-1
Chromosome: 7; Location: 7p11.2
Annotation: Chromosome 7 NC_000007.14 (55019017..55211628)
MIM: 131550
ID: 1956

c. Codes:

```
Lab1b.py M X
Lab > Lab1b.py > ...
1  from Bio import Entrez
2  import pandas as pd
3  import re
4
5  # Define your search term
6  search_term = "cystic fibrosis"
7
8  # Perform the search
9  handle = Entrez.esearch(db="gene", term=search_term, retmax=3)
10 record = Entrez.read(handle)
11 handle.close()
12 print("\n")
13
14 # Get the list of Gene IDs from the search results
15 gene_ids = record["IdList"]
16
17 # Retrieve gene information for all search results
18 gene_information = []
19 for gene_id in gene_ids:
20     handle = Entrez.efetch(db="gene", id=gene_id,
21 rettype="gb", retmode="text")
22     gene_info = handle.read()
23     handle.close()
24     gene_information.append(gene_info)
25
26 # Extract gene symbol, other aliases, and ID
27 table_data = []
28 for gene_info in gene_information:
29     # Extract gene symbol
30     symbol_match = re.search(r'^\d+\.\s+(\S+)', gene_info, re.MULTILINE)
31     symbol = symbol_match.group(1) if symbol_match else "N/A"
32
33     # Extract other aliases
34     aliases_match = re.search(r'Other Aliases:\s*(.+?)(?:\n|$)', gene_info)
35     aliases = aliases_match.group(1).strip() if aliases_match else "N/A"
36
37     # Extract ID
38     id_match = re.search(r'ID:\s*(\d+)', gene_info)
39     gene_id = id_match.group(1) if id_match else "N/A"
40
41     table_data.append([symbol, aliases, gene_id])
42
43 # Create and display pandas DataFrame
44 df = pd.DataFrame(table_data, columns=["Gene Symbol", "Other Aliases", "ID"])
45 print(f"Gene search results for: {search_term}\n")
46 print(df.to_string(index=False))
```

Output:

Gene search results for: cystic fibrosis

Gene Symbol	Other Aliases	ID
cftr	N/A	144521857
cftr	N/A	144498565
EGFR ERBB, ERBB1, ERRP, HER1, NISBD2, NNCIS, PIG61, mENA		1956

2. Codes:

```
Lab1(2).py U X
Lab > Lab1(2).py > ...
1  from sklearn.datasets import load_iris
2  from sklearn.model_selection import StratifiedKFold
3  import pandas as pd
4
5  # Load the Iris dataset
6  iris = load_iris()
7  X = iris.data
8  y = iris.target
9
10 # Display dataset size before cross-validation
11 print(f"Total samples: {X.shape[0]}")
12 print(f"Number of features: {X.shape[1]}")
13 print(f"Number of classes: {len(set(y))}")
14 print(f"Dataset dimensions: {X.shape}\n")
15 print("=" * 70 + "\n")
16
17 # Create Stratified K-Fold with 10 splits
18 skf = StratifiedKFold(n_splits=10, shuffle=True, random_state=42)
19
20 # Perform stratified k-fold cross-validation
21 fold_results = []
22 for fold, (train_index, test_index) in enumerate(skf.split(X, y), 1):
23     X_train, X_test = X[train_index], X[test_index]
24     y_train, y_test = y[train_index], y[test_index]
25
26     train_size = X_train.shape[0]
27     test_size = X_test.shape[0]
28
29     fold_results.append({
30         'Fold': fold,
31         'Training Set Size': train_size,
32         'Test Set Size': test_size,
33         'Total': train_size + test_size
34     })
35
36 # Create a pandas DataFrame
37 df = pd.DataFrame(fold_results)
38 print(df.to_string(index=False))
```

Output:

```
Total samples: 150
Number of features: 4
Number of classes: 3
Dataset dimensions: (150, 4)

=====

Fold  Training Set Size  Test Set Size  Total
1      135              15    150
2      135              15    150
3      135              15    150
4      135              15    150
5      135              15    150
6      135              15    150
7      135              15    150
8      135              15    150
9      135              15    150
10     135              15    150
```