

Welcome to the module 6a coding part: Loading and preparing Data!

This notebook was created at San Francisco State University (SFSU) for the Promoting InClusivity and Computing (PINC) and gSTAR programs by Dr. Pleuni Pennings (SFSU biology professor), Lucy Moctezuma Tan (California State University, East Bay CSUEB master student) and Lorena Benitez-Rivera (SFSU master students). All members of the C0de to understand Drug resistance Evolution (CODE) lab in 2023.

OBJECTIVE OF THIS NOTEBOOK:

In this notebook we are going to load and prepare the antibiotic resistance data to use in the next notebook and reproduce the Moradigaravand 2018 paper machine learning models.

The main goal is to explore and undertand the data we will be using as features to predict **Resistance (R)** or **Susceptibility (S)** in *E.coli* bacteria for several drugs. In addition we will learn how to clean up and preprocess our data before feeding it into different machine learning algorithms.

This tutorial will use the features of:

- **Year of isolation (Y)**
- **Gene absence or presence (G)**
- **Population Structure (S)**

Data origin.

Publication: Moradigaravand D, Palm M, Farewell A, Mustonen V, Warringer J, Parts L (2018) Prediction of antibiotic resistance in *Escherichia coli* from large-scale pan-genome data. PLoS Comput Biol 14(12): e1006258. <https://doi.org/10.1371/journal.pcbi.1006258>

Github page: <https://github.com/DaneshMoradigaravand/PanPred>

WHY ANTIBIOTIC RESISTANCE?

Antibiotic resistance has become a global public health concern. Bacteria are evolving resistance to the current arsenal of prescribed antibiotics resulting in strains that are developing multi-drug resistance. Currently, doctors and clinics often perform traditional culture-based assays (i.e., testing whether a drug would kill the bacteria in a petri dish) to determine antibiotic resistance in bacterial strains.

Alternatively, clinics can also choose to sequence these strains. These sequences can then be analyzed in order to predict if it will show resistance to a certain drug. There are different ways to perform the analysis and machine learning is one of them. This series of tutorials intends to help students understand how to do such analysis.

We will process publically available whole genome sequences of *E. coli* strains to create:

1. **Gradient Boosted Trees**
2. **Random Forest**
3. **Logistic Regression**

These three models are going to be used to predict **Resistance (R)** and **Susceptibility (S)** for each strain. The strains have already been tested in the lab, so we will later be able to compare the predictions made by our machine learning models with the traditional culture-based assays results in order to determine the performance for each of these models.

Step 1) Importing all necessary packages for dataframe manipulation

The code below will allow you to import the packages needed to load pre-process the data used for our models. In addition we will be loading up the "os" package, which will help us use the google drive much like how we create and save files in our own computer.

NOTE: Please allow access to your google drive when prompted, this will let you create and store the files in your drive to be accessed later by subsequent notebooks as we make progress towards getting our final results.

```
# Data Wrangling Imports
import pandas as pd
import numpy as np
from functools import reduce

# Data visualization
import matplotlib.pyplot as plt
import seaborn as sns

# File Manipulation Imports
import os as os
from google.colab import drive
drive.mount('/content/drive')
```

Mounted at /content/drive

Step 2) Loading all feature datasets used for predictions

The features from the paper we will use for this workshop are:

a) Metadata: Year of isolation and results from Antimicrobial Susceptibility Testing.

b) Gene presence and absence: Coding regions converted to proteins and ortholog accessory genes.

c) Population structure: Clusters based on SNPs distance (number of differing sites) for core genome.

```
# assign to url variable for each csv file
metadata_url = 'https://raw.githubusercontent.com/DaneshMoradigaravand/PanPred/master/test_data/Metadata.csv'
gene_presence_url = 'https://raw.githubusercontent.com/DaneshMoradigaravand/PanPred/master/test_data/AccessoryGene.csv'
pop_struc_url = 'https://raw.githubusercontent.com/DaneshMoradigaravand/PanPred/master/test_data/PopulationStructure.csv?labelencoded.csv'

# load csv files into the notebook
metadata = pd.read_csv(metadata_url)
gene_presence_data = pd.read_csv(gene_presence_url)
pop_struc_data = pd.read_csv(pop_struc_url)
```

a) Metadata:

Columns Summary:

- **Isolate number:** Unique number assigned to identify a particular strain of *E. coli* Bacteria. Thus we will refer to each row of our dataset as an "isolate" from now on.
- **Year of isolation:** The year in which a particular bacteria strain was separated from it's natural environment, thus the name isolate.
- **Antibiotic drug:** There are **12 antibiotic drug** columns named after their **3 letter abbreviation** adopted from the 'British Society of Antimicrobial Chemotherapy'.

Abreviation	Class: Subclass	Full name
CTZ	Beta-lactams: Cephalosporins	Ceftazidime
CTX	Beta-lactams: Cephalosporins	Cefotaxime
CXM	Beta-lactams: Cephalosporins	Cefuroxime
CET	Beta-lactams: Cephalosporins	Cephlothin
AMP	Beta-lactams: Penicillin	Ampicillin
AMX	Beta-lactams: Penicillin	Amoxicillin
AMC	Beta-lactams: Penicillin	Amoxicillin + Clavulanate potassium
TZP	Beta-lactams: Piperacillin	Tazobactam
GEN	Aminoglycosides	Gentamicin
TBM	Aminoglycosides	Tobramycin
TMP	Antifolate	Trimethoprim
CIP	Fluoroquinolones	Ciprofloxacin

```
# Visualize the first 5 rows of our dataframe
metadata.head()
```

Use metadata生成程式碼

查看建議的圖表

New interactive sheet

	Isolate	Year	CTZ	CTX	AMP	AMX	AMC	TZP	CXM	CET	GEN	TBM	TMP	CIP
0	11657_5#10	2010.0	S	S	S	NaN	S	S	S	S	S	S	S	
1	11657_5#11	2010.0	S	S	R	NaN	R	S	S	S	S	S	R	R
2	11657_5#12	2010.0	S	S	S	NaN	S	S	S	S	S	S	S	S
3	11657_5#13	2010.0	S	S	R	NaN	R	S	S	S	S	S	S	R
4	11657_5#14	2010.0	S	S	R	NaN	S	S	S	S	S	S	R	S

Determination of Resistant (R) and Susceptible (S) labels:

Antimicrobial Susceptibility Testing: these are laboratory tests, where particular concentrations of a drug are tested for each bacteria isolate to determine whether it is **Resistant (R)**, **Susceptible (S)** or **Intermederate (I)**.

- **Clinical Breakpoints:** Each of the drugs listed in the chart have a different concentration (clinical breakpoint) used to determined the resistance or susceptability of *E. coli* to that drug. These concentrations are agreed upon by pubic health agencies based on different factors and can vary year to year given new discoveries. They serve as a guiding tool for clinicians to prescribe a particular antibiotic drug and appropriate dose for patients.
- **Resistant (R) & Susceptible (S)** results from the laboratory tests were determined based on the series of guidelines from the [European Committee on Antimicrobial Susceptibility Testing \(EUCAST\)](#) on 25/01/2017. For this study, isolates that were classified as **Intermediate (I)** were lumped together with the **Resistant (R)** ones.

Note: NaN was used to mark missing data, for those isolates that were not tested with a particular drug.

```
# Observe the shape and size of the meta data dataframe
metadata.shape
```

(1936, 14)

The code above can show us the shape of our entire dataframe in the following format: **(row count, and column count)**

a) Row count: There is a total of **1,936 isolates** of *E.coli* bacteria.

b) Column count: There are 14 features for the metadata = Isolate number(1) + year of isolation (1) + antibiotic drugs (12).

b) Gene presence and absence:

```
# Observe the shape and size of the gene_presence_data dataframe
gene_presence_data.shape
```

(2033, 17199)

a) Row count: There is a total of **2,033 isolates** of *E. coli* bacteria.

b) Column count: There are 17,199 features = isolate number (1) + coding_genes (3,815) + other_genes (13,383)

Genes: Are a chunks of DNA that can have different sizes. Some genes code for specific proteins and their function is fairly well understood, whereas other genes function are not well understood yet. Genes can be generally classified into:

- **Core genome:** These are the genes that are present in almost all individuals for a particular bacterial species. In our case, they are the genes that all our *E.coli* isolates have in common.
- **Accessory genes:** These are the genes that might be found in one individual but may not in another individual of the same species. These would be the genes that are unique to each of our *E. coli* isolate.
- **Pan-genome:** These are all the possible genes that can be found in a particular species. That is they are all the gene presents present in our *E. coli* isolates (Pan-genome = core genome + accesory genes).

The code below will show us the dataframe containing the presence and absence of all the genes detected in for each isolate.

- **0** = Absence of the gene
- **1** = Presence of the gene.

```
# rename all the df with "Unnamed: 0" columns as "Isolate"
gene_presence_data = gene_presence_data.rename(columns={'Unnamed: 0' : 'Isolate'})
# Visualize the first 5 rows of dataframe
gene_presence_data.head()
```


- For example if we increase the cutoff value to say **cutoff_27236**, we can observe that most isolates (1932) would fall in "cluster 0" with only a few isolates dispersed in other clusters.

```
# code will count how many isolates belong to a particular cluster
cutoff_27236= pd.DataFrame(pop_struct_data['cutoff_27236'].value_counts() )
cutoff_27236
```



後續步驟: [使用 cutoff_27236生成程式碼](#) [查看建議的圖表](#) [New interactive sheet](#)

Because there is no information about which cutoff value would yield us the best way to cluster our isolates or which cluster memberships are necessarily related to resistance (R) or susceptibility (S), the author included 1072 different cutoff values, each with their own clusters. Then we will trying to predict if a particular cluster membership will help us predict R or S for each cutoff value.

```
# Observe the shape and size of the dataframe
pop_struct_data.shape
```

(1936, 1072)

Step 3) Preparing each dataset in order to be combined

a) Changing Metadata's Year column to One Hot encoded variables

This is because even though year is a number, it is actually considered a categorical variable. The ratio between two years is not meaningful, so it is incorrect to keep it as a purely quantitative variable.

```
# This prints the years the strains were isolated
year_list = metadata['Year'].unique()
year_list = year_list[np.logical_not(np.isnan(year_list))]
print(sorted(year_list))
```

[1970.0, 1977.0, 1994.0, 1997.0, 1998.0, 1999.0, 2001.0, 2002.0, 2003.0, 2004.0, 2005.0, 2006.0, 2007.0, 2008.0, 2009.0, 2010.0, 2011.0, 2012.0, 2013.0, 2014.0, 2015.0, 2017.0]

By creating one hot encoded variables from the Year column of the metadata, each of their rows are coded the following way:

- "0": Isolate was performed in a particular year
- "1": Isolate was not performed in a particular year

```
# creating one hot encoded variables from "Year" column to create a matrix of years
metadata_d = pd.get_dummies(metadata, columns=["Year"], dummy_na=False)
metadata_d
```

	Isolate	CTZ	CTX	AMP	AMX	AMC	TZP	CMX	CET	GEN	...	Year_2007.0	Year_2008.0	Year_2009.0	Year_2010.0	Year_2011.0	Year_2012.0	Year_2013.0	Year_2014.0	Year_2015.0	Year_2017.0
0	11657_5#10	S	S	S	NaN	S	S	S	S	S	...	False	False	False	True	False	False	False	False	False	False
1	11657_5#11	S	S	R	NaN	R	S	S	S	S	...	False	False	False	True	False	False	False	False	False	False
2	11657_5#12	S	S	S	NaN	S	S	S	S	S	...	False	False	False	True	False	False	False	False	False	False
3	11657_5#13	S	S	R	NaN	R	S	S	S	S	...	False	False	False	True	False	False	False	False	False	False
4	11657_5#14	S	S	R	NaN	S	S	S	S	S	...	False	False	False	True	False	False	False	False	False	False
...
1931	24742_1#96	S	S	S	NaN	NaN	NaN	S	S	S	...	False	False	False	False	False	False	False	False	False	False
1932	24742_1#97	S	S	S	NaN	NaN	NaN	S	S	S	...	False	False	False	False	False	False	False	False	False	False
1933	24742_1#98	S	S	R	NaN	NaN	NaN	S	S	S	...	False	False	False	False	False	False	False	False	False	False
1934	24742_1#99	S	S	R	NaN	NaN	NaN	S	S	S	...	False	False	False	False	False	False	False	False	False	False
1935	24742_1#0	S	S	R	NaN	NaN	NaN	S	S	S	...	False	False	False	False	False	False	False	False	False	True

1936 rows × 35 columns

Step 4) Making a single dataframe using all 3 sources of data

a) Joining all data sources into a single dataframe

Note: The function **merge()**, allows to pass the parameter (**on = "Isolate"**), which will ensure that each isolate number is correctly matched for all 3 data sources, the parameter (**how="inner"**) will make sure that isolates without a match are not included in the final dataframe.

```
# List of all 3 data sources
df_list = [metadata_d, gene_presence_data, pop_struct_data]
```

```
# creating a single dataframe with all drugs and features available
Drug_df = reduce(lambda left, right: pd.merge(left, right, on=["Isolate"], how="inner"), df_list)
Drug_df.head()
```

	Isolate	CTZ	CTX	AMP	AMX	AMC	TZP	CMX	CET	GEN	...	cutoff_25459	cutoff_25654	cutoff_25772	cutoff_25979	cutoff_26792	cutoff_27119	cutoff_27236	cutoff_27248	cutoff_27690	cutoff_45092
0	11657_5#10	S	S	S	NaN	S	S	S	S	S	...	0	0	0	0	0	0	0	0	0	0
1	11657_5#11	S	S	R	NaN	R	S	S	S	S	...	0	0	0	0	0	0	0	0	0	0
2	11657_5#12	S	S	S	NaN	S	S	S	S	S	...	0	0	0	0	0	0	0	0	0	0
3	11657_5#13	S	S	R	NaN	R	S	S	S	S	...	0	0	0	0	0	0	0	0	0	0
4	11657_5#14	S	S	R	NaN	S	S	S	S	S	...	0	0	0	0	0	0	0	0	0	0

5 rows × 18304 columns

- Notice that the number of rows is now correctly matched with Isolate number, yielding a total of 1936 rows as in the metadata.
- Also we now have a bunch of columns that currently include:
 - 1 isolate number column these are the unique tags for each of our isolates.
 - 12 labels, one for each drug we will try to make predictions for.
 - 18291 features that we will be using to make prediction for the labels (isolation year, gene presence or absence and population structure).

```
# Check out all the columns included in the final dataframe and the final shape it takes
print(Drug_df.columns) # contains all labels (drug abbreviation column names)
# and all features (year, gene presence absence and population structure)
```

```
print("Final shape of combined dataframe", Drug_df.shape)
```

```
Index(['Isolate', 'CTZ', 'CTX', 'AMP', 'AMX', 'AMC', 'TZP', 'CMX', 'CET', 'GEN', '...', 'cutoff_25459', 'cutoff_25654', 'cutoff_25772', 'cutoff_25979', 'cutoff_26792', 'cutoff_27119', 'cutoff_27236', 'cutoff_27248', 'cutoff_27690', 'cutoff_45092'],
      dtype='object', length=18304)
Final shape of combined dataframe (1936, 18304)
```

b) Convert the Dataframe into a CSV and save it in a folder

After running the code below, feel free to check your Drive to make sure that you have a folder named **"EColi_ML_CSV_files"** and that inside you have a csv called **"EColi_Merged_dfs.csv"**

NOTE: The code below creates a directory in your Drive that allows you to save your results. If you have previously created one already, it will produce an error message, that you can just ignore, as its just telling you that it already exists.

```
# makes a directory to save all your csv's
os.mkdir('/content/drive/My Drive/EColi_ML_CSV_files')
```

```
# path where we will store csv data #change to any path you want
path = '/content/drive/My Drive/EColi_ML_CSV_files/'
```

```
# this code exports the dataframe into a CSV file
Drug_df.to_csv(path+"EColi_Merged_dfs.csv", index= False)
```

DATA VISUALIZATION EXERCISES:

Below we have a couple of data visualization tasks for you to try!

Task 1: visualize how many of the samples are resistant per year and per drug.

Plot stacked barcharts for the percentage of Resitant versus Susceptible throughout the years, using the metadata dataframe for each of the 12 drugs. Tip: use the **groupby()** function to summarize the data first.

```
# prompt: Plot stacked barcharts for the percentage of Resitant versus Susceptible throughout the years, using the #
```

```
# Group the data by year and drug, then count the resistant and susceptible isolates
grouped_data = metadata.groupby(['Year'])
```

```
# Create a figure and axes for the plot
fig, axes = plt.subplots(nrows=6, ncols=2, figsize=(15, 30))
```

```
# Define the drugs for which to create plots
drugs = ['CTZ', 'CTX', 'CMX', 'CET', 'AMP', 'AMX', 'AMC', 'TZP', 'GEN', 'TBM', 'TMP', 'CIP']
```

```
# Iterate over the drugs and create stacked bar plots
```

```
for i, drug in enumerate(drugs):
    # Get the subplot for the current drug
    row = i // 2
    col = i % 2
    ax = axes[row, col]
```

```
# Filter data for the current drug
drug_data = metadata[metadata[drug].isna()==0]
drug_data_grouped = drug_data.groupby(['Year', drug]).size().unstack(fill_value=0)
```

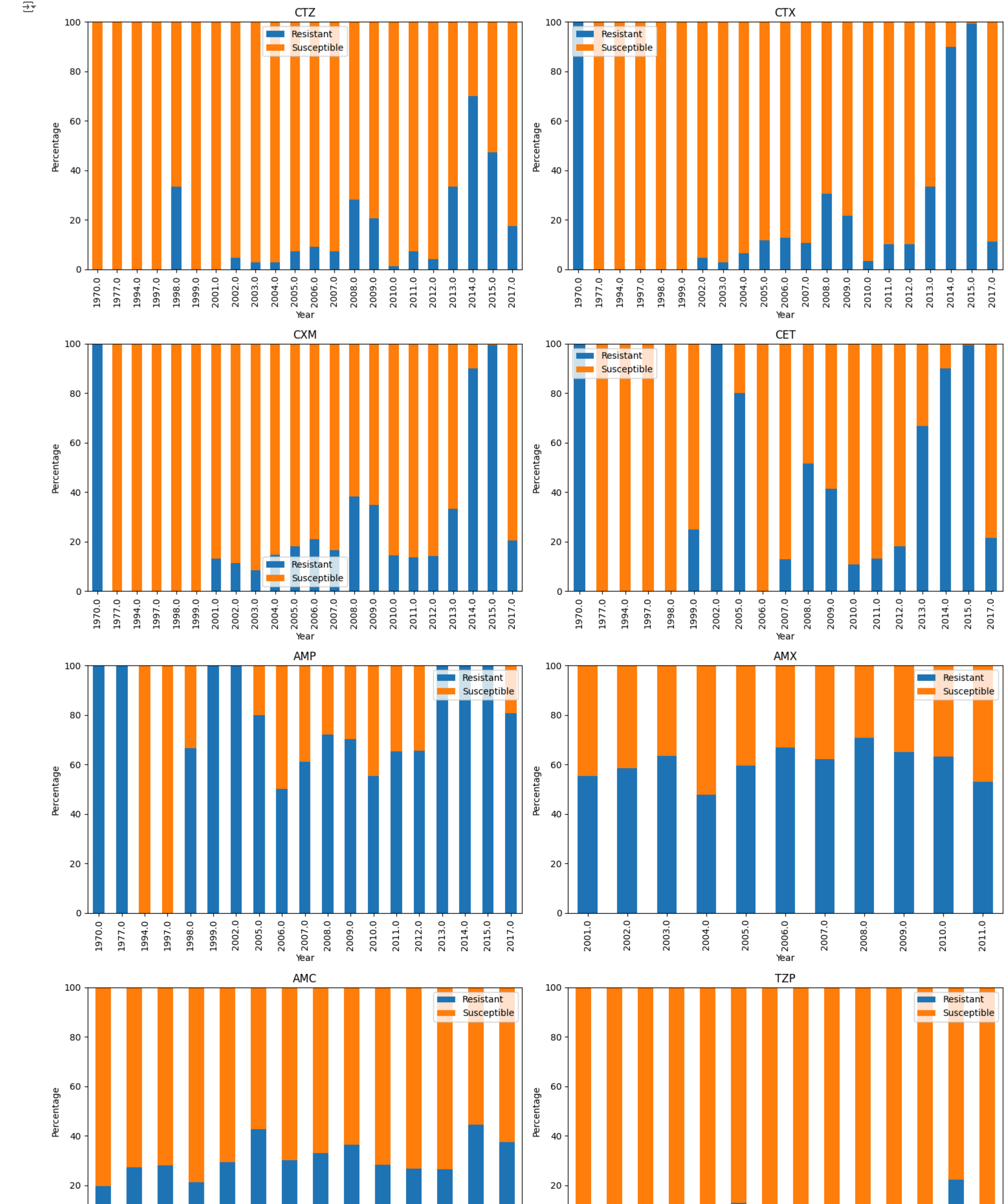
```
# Calculate the total number of isolates per year
drug_data_grouped['Total'] = drug_data_grouped.sum(axis=1)
```

```
# Calculate the percentage of resistant and susceptible isolates per year
drug_data_grouped['R_percentage'] = (drug_data_grouped['R'] / drug_data_grouped['Total']) * 100
drug_data_grouped['S_percentage'] = (drug_data_grouped['S'] / drug_data_grouped['Total']) * 100
```

```
# Plot stacked bar chart
drug_data_grouped[['R_percentage', 'S_percentage']].plot(kind='bar', stacked=True, ax=ax)
```

```
# Customize the plot
ax.set_title(f'{drug}')
ax.set_xlabel('Year')
ax.set_ylabel('Percentage')
ax.legend(['Resistant', 'Susceptible'])
ax.set_ylim(0, 100)
```

```
# Adjust layout for better visualization
plt.tight_layout()
plt.show()
```

- Task 2: visualize the clustering of the strains.

Using The Population Structure Dataframe, create 4 different pie charts for the following cutoff values: **cutoff_823**, **cutoff_6934**, **cutoff_9768**, [link text](#) **cutoff_11227**. Write what you notice as cutoff_values increase.

```
# prompt: Using The Population Structure Dataframe, create 4 different pie charts for the following cutoff values: 0.1, 0.2, 0.3, 0.4

import matplotlib.pyplot as plt

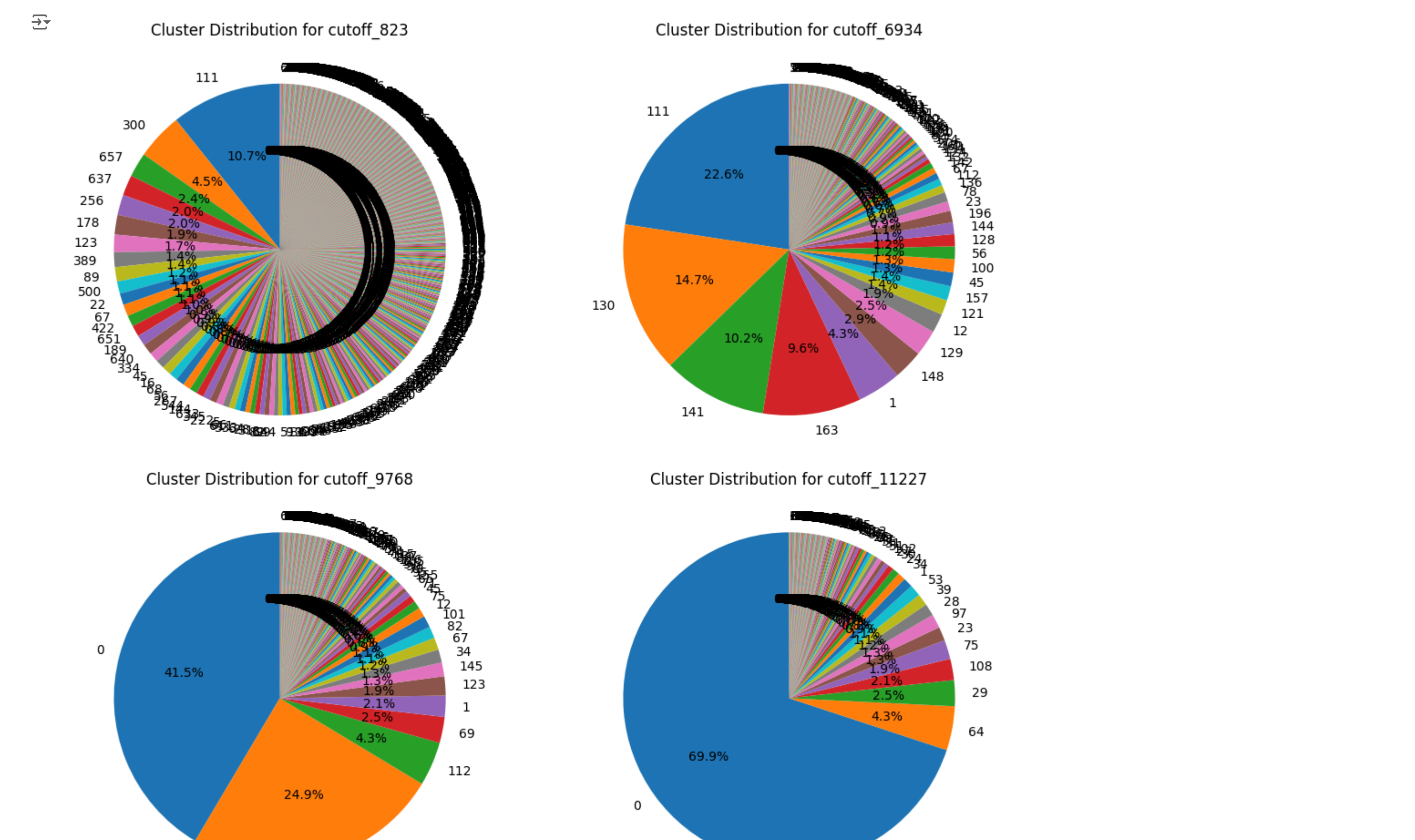
## Assuming pop_struc_data is your dataframe
cutoff_values = ['cutoff_823', 'cutoff_6934', 'cutoff_9768', 'cutoff_11227']

fig, axes = plt.subplots(2, 2, figsize=(12, 10))
axes = axes.flatten()

for i, cutoff in enumerate(cutoff_values):
    cluster_counts = pop_struc_data[cutoff].value_counts()
    axes[i].pie(cluster_counts, labels=cluster_counts.index, autopct='%1.1f%%', startangle=90)
    axes[i].set_title(f'Cluster Distribution for {cutoff}')

plt.tight_layout()
plt.show()

## Observations as cutoff values increase (add your observations here)
## For example:
## As the cutoff values increase, the number of clusters tends to decrease.
## Initially, with lower cutoffs, there's a greater diversity in cluster sizes.
## As the cutoff increases, a dominant cluster emerges, encompassing a larger portion of the isolates.
```



THANKS for making it this far! Now that our data is ready the next notebooks will deal with the creation and training of different algorithms using the dataframe we just created in this notebook!!!