# Lab 5: Dissecting the Role of Algorithms in Fragment-based Drug Discovery

Review the lab material and go through the entire notebook. The lab contains 5 exercises for you to solve. The entire lab is worth 2.5% of your final grade and each exercise is worth 0.4% of your final grade. Going through the full notebook is worth 0.5% of your final grade. Any extra credit or bonus exercises are worth an additional 0.4%.

Labs are due by Friday at 11:59 PM PST and can be submitted on BCourses assignment page for the corresponding lab.

## Setup

Run the below steps at the beginning of lab to set up required packages.

```
1 !pip install rdkit
```

```
Collecting rdkit
  Downloading rdkit-2025.9.1-cp312-cp312-manylinux_2_28_x86_64.whl.metadata (4.1 kB)
Requirement already satisfied: numpy in /usr/local/lib/python3.12/dist-packages (from rdkit) (2.0.2)
Requirement already satisfied: Pillow in /usr/local/lib/python3.12/dist-packages (from rdkit) (11.3.0)
Downloading rdkit-2025.9.1-cp312-cp312-manylinux_2_28_x86_64.whl (36.2 MB)
━━━━━━━━━━━━━━━━━━━━━━━━━ 36.2/36.2 MB 52.0 MB/s eta 0:00:00
Installing collected packages: rdkit
Successfully installed rdkit-2025.9.1
```

```
 1 import itertools, math, os
 2 import numpy as np
 3 import matplotlib.pyplot as plt
 4 import pandas as pd
 5 import random
 6 import seaborn as sns
 7
 8 from rdkit import Chem
 9 from rdkit.Chem import AllChem, ChemicalFeatures, DataStructs, MolToSmiles, PandasTools, rdFingerprintGenera
10 from rdkit.Chem.Scaffolds import MurckoScaffold
11 from rdkit.ML.Cluster import Butina
12 from rdkit.RDPaths import RDDataDir
13 from sklearn.cluster import AgglomerativeClustering
14 from scipy.cluster.hierarchy import dendrogram, linkage, fcluster
15 from scipy.optimize import curve_fit
16 from scipy.spatial.distance import cdist, pdist, squareform
17 from tqdm import tqdm
18 from typing import List, Tuple
19
20 RANDOM_SEED = 42
21 np.random.seed(RANDOM_SEED)
22
23 colors = ["#A20025", "#6C8EBF"]
24 sns.set_palette(sns.color_palette(colors))
```

## COVID Moonshot Campaign

In this lab, we'll focus on an application stemming from the COVID Moonshot campaign [1]. The COVID Moonshot campaign is an open-science initiative to accelerate discovery of antiviral drugs targeting inhibition of the SARS-CoV-2 main protease (Mpro). Inhibition of Mpro prevents viral replication of SARS-CoV-2, with precedence for clinical success set by Paxlovid and Ensitrelvir [2,3]. The campaign used high-throughput screening (HTS) techniques to evaluate thousands of compounds with respect to Mpro inhibition, optimizing promising compounds through iterative cycles of design, synthesis, and testing.

The COVID Moonshot incorporated many components, from crowdsourcing over 18,000 designs and judging those designs based on synthetic feasibility and estimated activity, to synthesizing over 2000 of those designs and conducting several types of high-throughput assays and crystallography experiments to determine their structure and viability as leads. We will focus solely on the set of 2,062 synthesized compounds whose IC50 was experimentally measured via HTS.

IC50 is the concentration of a compound required to inhibit 50% of protein activity. Ideally, a compound is effective even at low concentrations to avoid undesired interactions with other proteins.

The following code loads in the data set, which we have already processed by:

1. Computing Morgan fingerprints (2048 bits at radius of 2) for each compound
2. Using a cut-off of 5 μM IC50, based on literature that we'll discuss later, to define active and inactive compounds. Compounds that did not have a measured IC50 are still present in the data set but are labeled as "Unknown."

Make sure that you create a "data" folder and upload both files to the folder:

- LO2_activity_data.csv
- LO2_morgan_fingerprints.npy

```
1 data = pd.read_csv('data/L02_activity_data.csv')
2 fingerprints = np.load('data/L02_morgan_fingerprints.npy')
3 label_names = {1:'Active', 2:'Inactive', 3:'Unknown'}
```

```
1 data.head()
```

| | SMILES | CID | canonical_CID | r_inhibition_at_20_uM | r_inhibition_at_ |
|---|---|---|---|---|---|
| 0 | CCNC(=O)CN1CC2(CCN(c3cncc4ccccc34)C2=O)c2cc(Cl... | LUO-POS-e1dab717-11 | LUO-POS-e1dab717-11 | NaN | |
| 1 | O=C(CN1CC2(CCN(c3cncc4ccccc34)C2=O)c2cc(Cl)ccc... | LUO-POS-e1dab717-12 | LUO-POS-e1dab717-12 | NaN | |
| 2 | CNC(=O)C1(N2C[C@]3(CCN(c4cncc5ccccc45)C3=O)c3c... | MAT-POS-e48723dc-1 | MAT-POS-e48723dc-1 | NaN | |
| 3 | CNC(=O)C1(N2C[C@@]3(CCN(c4cncc5ccccc45)C3=O)c3... | MAT-POS-e48723dc-2 | MAT-POS-e48723dc-2 | NaN | |
| 4 | CNC(=O)CN1C[C@@]2(CCN(c3cncc4ccccc34)C2=O)c2cc... | LUO-POS-9931618f-2 | LUO-POS-9931618f-2 | NaN | |

Next steps: ( Generate code with data ) ( New interactive sheet )

## ⌄ High Throughput Screening (HTS) Data

In this scenario, we want to identify a new antiviral that can stop COVID-19. With advancements in robotic and miniaturization techniques, HTS is capable of handling a high volume of thousands of tiny tests at once and in an automated fashion. In early stages of HTS, single-shot inhibition data provides a quick assessment of each compound's activity against the target. A single concentration, which is relatively high to maximize the chance of detecting any inhibition, is selected. Each compound from a large library of compounds is tested at this concentration, and the inhibition of the target's activity is measured. Compounds that exhibit significant inhibition at the single concentration are flagged as potential hits for follow-up studies.

If resource and time requirements permit, we can measure each compound's activity at multiple concentrations to construct dose response curves. Dose response curves plot the compound's inhibitory activity on the target against different concentrations of the compound, from which we can determine key parameters such as IC50. For example, the COVID Moonshot program determined IC50 from dose response curves that used an 11-point range of 100—0.0017 µM inhibitor concentrations, ensuring that both high and low efficacy could be measured.

To determine the IC50 for each compound, the Moonshot program used the Levenberg-Marquardt algorithm to fit a restrained Hill equation to each compound's dose response curve. Let's break these concepts down. The Hill equation is a mathematical model to describe the relationship between the compound's concentration and the expected response:

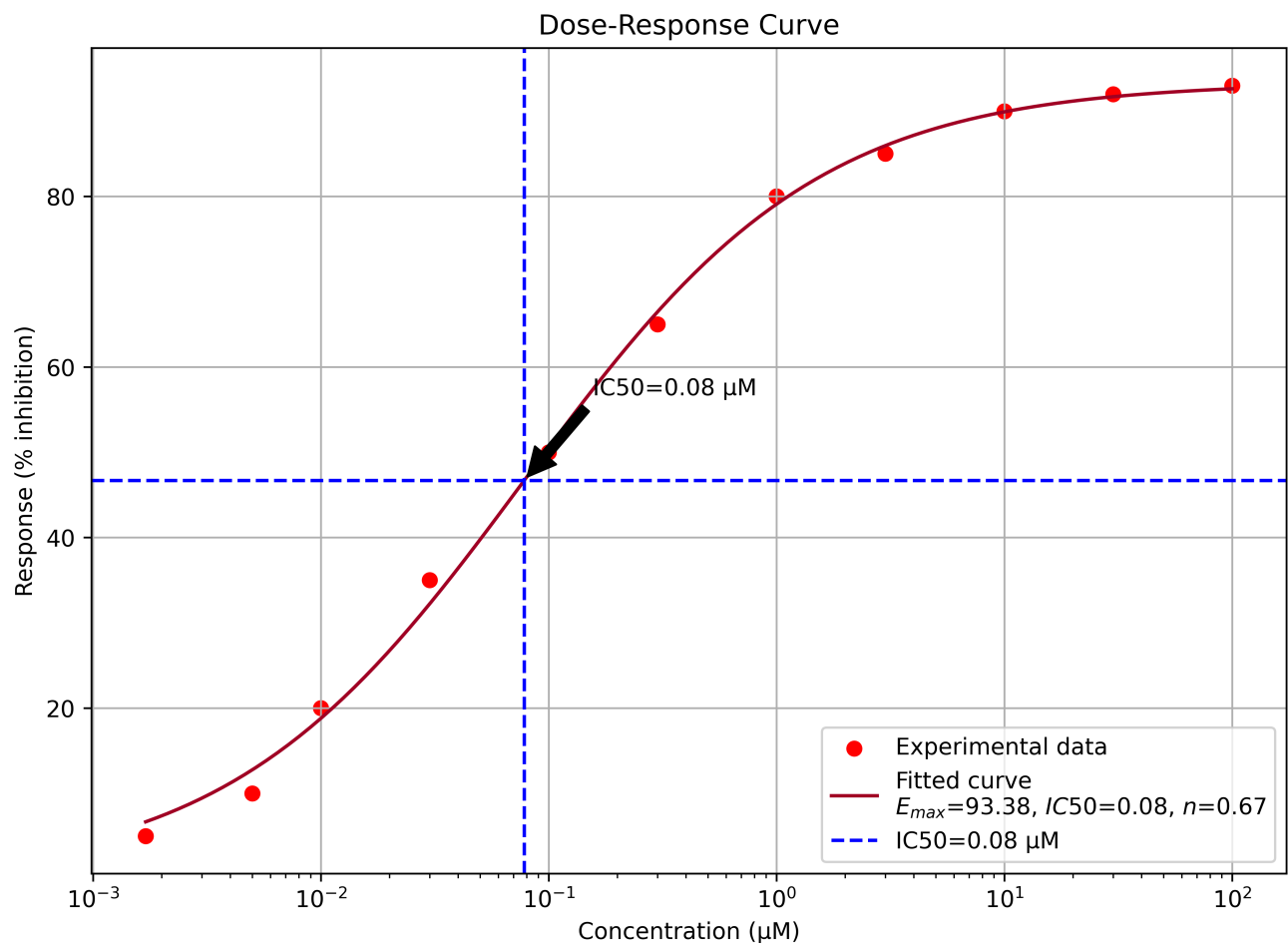$$E = \frac{(E_{max}[I]^n)}{(IC_{50}^n + [I]^n)}$$

- $E$ is the observed response.
- $E_{max}$ is the maximum response.
- $[I]$ is the concentration of the compound.
- $IC_{50}$ is the concentration at which the response is half of $E_{max}$.
- $n$ is the Hill coefficient, which describes the slope of the curve.

We can use the Levenberg-Marquardt algorithm to fit a non-linear model of our HTS data. In the context of the Hill equation, the Levenberg-Marquardt algorithm is useful since it accounts for the non-linear nature of the dose-response relationship. At a high-level, we start with initial guesses for the parameters in the Hill equation ($E_{max}$, $[I]$, $n$) and we iteratively adjust the parameters to minimize the error between the observed data points and the model's predictions via a combination of gradient descent and the Gauss-Newton method. A restrained fit means that certain parameters are kept within specified, realistic bounds to ensure that the resultant IC50 value is reliable and biologically meaningful.

## Exercise 1: Levenberg-Marquardt Algorithm

The following code block will plot a figure that demonstrates the sigmoidal shape common to dose-response curves, a fitted Hill equation, and which point on the fitted curve corresponds to the IC50 value. This figure is based on an arbitrary, toy data set.

Your task is to complete the missing sections of the code block (where a `TODO` comment is present). Once you have added the correct code, it will output the same figure as shown here:



### Dose-Response Curve

## Student Solution to Exercise 1

Provide your solution to the above exercise in this cell and/or immediately following cells.

```python
1 # Synthetic dose-response data
2 concentrations = np.array([0.0017, 0.005, 0.01, 0.03, 0.1, 0.3, 1, 3, 10, 30, 100])  # in µM
3 responses = np.array([5, 10, 20, 35, 50, 65, 80, 85, 90, 92, 93])  # arbitrary units, % inhibition for examp
4
5 # Initial guesses for E_max, IC50, and Hill coefficient (n)
6 initial_guesses = [100, 1, 1]
7
8 # Constraints: E_max between 0 and 100, IC50 between 0 and 100, n between 0.1 and 3
9 param_bounds = ([0, 0, 0.1], [100, 100, 3])
10
11 # Hill equation function
12 def hill_equation(conc, E_max, IC50, n):
13     """
14     Calculate response using the Hill equation
15     conc: drug concentration
16     E_max: maximum effect
17     IC50: concentration producing 50% of max effect
18     n: Hill coefficient (steepness)
19     """
20     # TODO: Implement the Hill hill_equation
21     E = E_max * conc ** (n) / ((IC50 ** (n)) + (conc ** (n)))
22     return E
23
24
25 # Fit the Hill equation to the data using the Levenberg-Marquardt algorithm
26 popt, pcov = curve_fit(hill_equation, concentrations, responses,
```
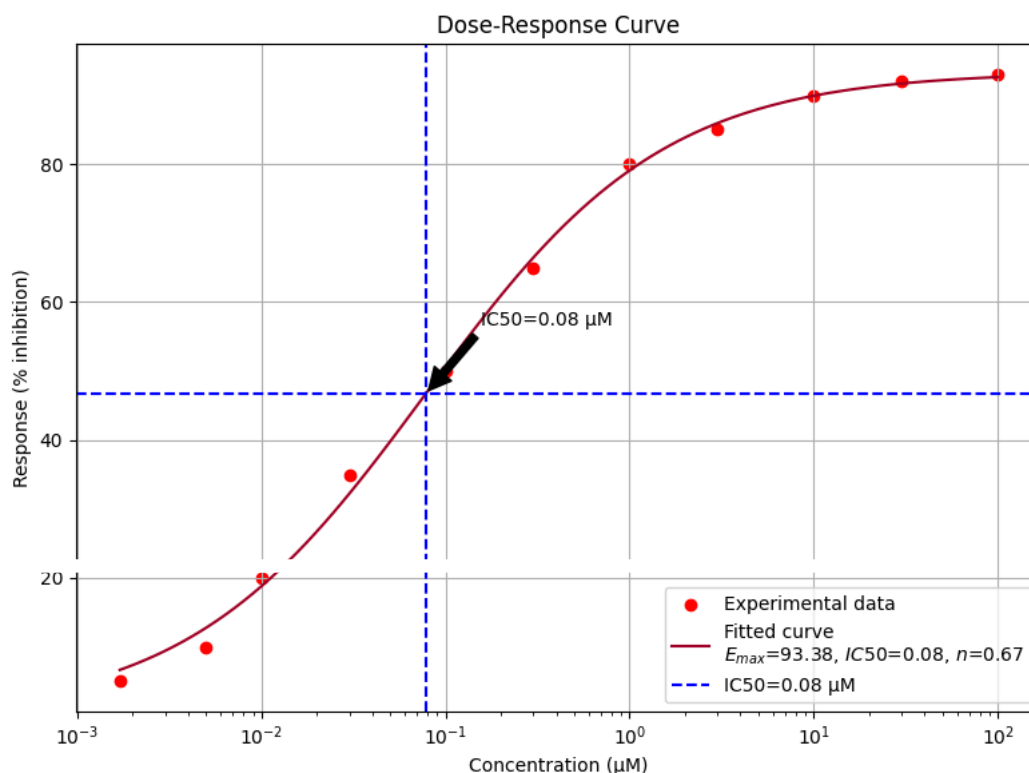
```
27                              p0=initial_guesses,
28                              bounds=param_bounds,
29                              method='trf')
30
31 # Extract the optimal parameters
32 E_max_opt, IC50_opt, n_opt = popt
33
34 # Generate data for the fitted curve, then plot it
35 concentration_range = np.logspace(np.log10(concentrations.min()), np.log10(concentrations.max()), 100)
36 fitted_responses = hill_equation(concentration_range, *popt)
37
38 plt.figure(figsize=(8, 6))
39 plt.scatter(concentrations, responses, color='red', label='Experimental data')
40 plt.plot(concentration_range, fitted_responses, label=f'Fitted curve\n$E_{{max}}$={E_max_opt:.2f}, $IC50$={I
41 plt.xscale('log')
42 plt.xlabel('Concentration (µM)')
43 plt.ylabel('Response (% inhibition)')
44 plt.title('Dose-Response Curve')
45 plt.legend()
46 plt.grid(True)
47
48 # To get the IC50 response, return the output if hill_equation using a combination of the parameters IC50_op
49 ic50_response = hill_equation(IC50_opt, E_max_opt, IC50_opt, n_opt)
50 plt.axhline(y=ic50_response, color='blue', linestyle='--', label=f'IC50={IC50_opt:.2f} µM')
51 plt.axvline(x=IC50_opt, color='blue', linestyle='--')
52 plt.annotate(f'IC50={IC50_opt:.2f} µM',
53                 xy=(IC50_opt, ic50_response),
54                 xytext=(IC50_opt*2, ic50_response+10),
55                 arrowprops=dict(facecolor='black', shrink=0.05))
56
57 plt.legend()
58 plt.tight_layout()
59
60 print(f"Fitted parameters:\nE_max: {E_max_opt}\nIC50: {IC50_opt}\nHill coefficient (n): {n_opt}")
```

```
Fitted parameters:
E_max: 93.38453452061374
IC50: 0.07799836455235021
Hill coefficient (n): 0.6707229452417185
```



HTS presents a high volume of data at our doorstep, but with additional confounders in the form of greater levels of experimental noise and greater diversity in the chemical structures within the compound library. We'd like to extrapolate the HTS data on these compounds to novel compounds. In the next section, we'll consider a few approaches to generating a compound library with novel compounds.

## ⌄ Curating Diverse Compound Libraries

In lab 1, we used the similar property principle to search for compounds within a screening library that had high structural similar to known antimalarial compounds, with the hope that such compounds would exhibit similar antimalarial properties. In this section, we'll consider how similarity, and diversity, play fundamental roles in the design of compound libraries.

## Diversity and Focus

The concept of creating representative sets of compounds is a critical strategy in the pharmaceutical industry. Consider a scenario where we have conducted a virtual screening campaign but have limited resources to experimentally test only a few compounds in a confirmatory assay. To maximize the information obtained from this screen, we select a diverse set of compounds. This might mean picking one representative from each chemical series in our list of potentially active compounds, ensuring that the selected compounds span the broadest range of chemical diversity.

Alternatively, we might focus on a single chemical series to delve into the structure-activity relationship within that series. By systematically testing compounds within the same series, we can gain detailed insights into how small structural changes drive biological activity, thus informing future optimization and development efforts. When designing compound libraries, we strike a balance between diversity and focus:

- Diversity, i.e., maximizing exploration: More structurally diverse libraries provide greater coverage over the expanse of chemical space that we want to explore, maximizing the number of compounds with different activities and minimizing the number of redundant compounds with similar activities. We place greater emphasis on diversity when we expect our compound library to be used to screen against a broad set of biological targets or when we don't have much prior information about the therapeutic target of interest.
- Focus, i.e., maximizing exploitation: Focused libraries restrict our navigation of chemical space to a narrower region. We place greater emphasis on focused libraries when we know more about the downstream use case, incorporating as much information about the therapeutic target (such as its 3D structure, if available) to maximize our hit rate. Focused libraries provide higher hit rates, but these starting points may have greater redundancy in activity or function.

There is not a single, authoritative definition or quantification of chemical diversity. There are many methods for assessing the internal diversity of the molecules in our library, as well as the diversity of those molecules compared to external data sets from prior screens. We can quantify diversity as a function of any molecular descriptors or scaffolds, and how we measure diversity will also impact how we sample which molecules are to be included in the library.
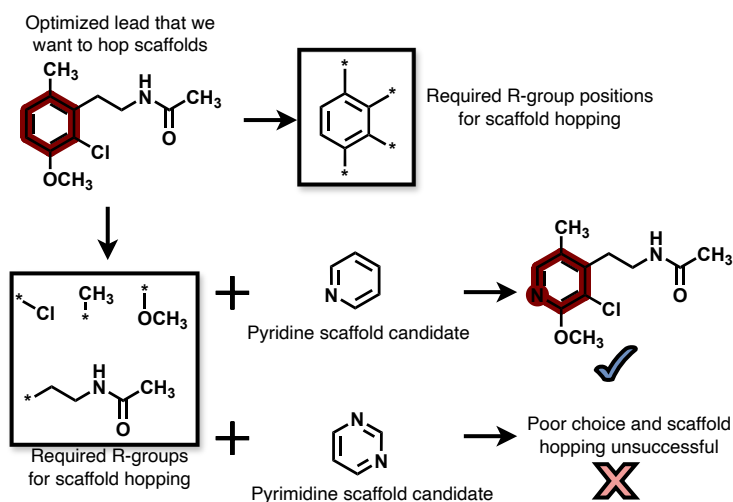
Once we've settled on which descriptors we want to use to characterize our molecules and how we want to measure the aggregate diversity of the molecules we are considering, we can construct each subset of size K from the total number of molecules, N, that we are considering, measure each subset's diversity, and keep the subset with greatest diversity. If we want to construct a library of 100 molecules from an initial set of 500 molecules, we'd need to evaluate $\binom{N}{K} = \binom{500}{100} = 2 * 10^{107}$ subsets. This is not realistic.

Later in this lab, we'll lean on machine learning to help resolve this sampling problem. Note that you aren't expected to know the mechanics of how any machine learning algorithms work, we are just using them to move the lab forward!

## ⌄ Molecular Fragments

Before we can select a diverse subset of compounds, we have to establish the data source that these compounds are being selected from. Such data sources, typically large compound libraries, do not have to be pre-existing collections like the Specs library from lab 1. In the present, we'll cover two classic approaches used to construct synthetic compound libraries, from which we can use as a basis to sample a diverse subset of compounds. If you find this subject interesting, you may want to explore the use of generative ML models to create synthetic compound libraries outside of this class.

The two classic approaches we will cover rely on understanding what a molecular fragment is. When we use the term "molecular fragment," we are typically referring to a scaffold or an R-group. During drug design, medicinal chemists might identify a scaffold, or core, of their compounds that is essential to achieve the desired pharmacological effect. Compounds with the same scaffold often constitute a chemical series.

Optimized lead that we want to hop scaffolds

Required R-group positions for scaffold hopping

Required R-groups for scaffold hopping

Pyridine scaffold candidate

Pyrimidine scaffold candidate

Poor choice and scaffold hopping unsuccessful

During optimization, the scaffold is typically kept constant while the substituents attached to the scaffold (the R-groups) are varied to improve the compound's properties. Different R-groups can be attached to the same scaffold to create a library of analogs, which we can also use to explore how structural modifications to the compound affect its activity. This process is referred to as scaffold decoration.

In some cases, fundamental flaws inherent to the scaffold core arise, such as toxicity or patentability concerns. Scaffold hopping involves identifying entirely new scaffolds that exhibit similar biological activity as our current scaffold, while obviating their limitations. Note that we exchange the entire scaffold of the molecule while preserving the R-groups that were originally attached to it. Searching for alternative, novel scaffolds is harder than just adjusting the R-groups. From a modeling perspective, we must also consider that models exposed only to compounds with the old scaffold may not translate to reliable predictions on compounds with the new scaffold.

One application of scaffold hopping is to resolve intellectual property issues in drug repurposing. The original patent holders may have exclusive rights to the drug being repurposed, which can complicate IP ownership unless new patents for new indications can be obtained. If this interests you, check out Markush structures, which is a representation in chemical patents and IP claims that defines a broad class of chemical compounds by specifying variable parts of a molecule, allowing for the inclusion of multiple possible substituents or structural variations.

## Retrosynthetic Library Design

The COVID Moonshot project's methodology was primarily concerned with developing a therapeutic as quickly as possible. In an (at the time) on-going pandemic, speed matters! Their aim was to make as many different ideas as possible for experimental follow ups, and so they prioritized compounds that could be rapidly synthesized by using ML to design synthetic schemes with uncertainty estimates.

Retrosynthetic library design safeguards against generation of compounds that are not synthesizable by baking in consideration of synthetic feasibility. We can use RDKit's BRICS (Breaking of Retrosynthetically Interesting Chemical Substructures) module to both fragment compounds and generate new compounds by recombining these fragments based on common chemical reaction rules. Now we will construct a compound library based on retrosynethic library design. We will then use a model to screen against these synthetically generated compounds.

```
 1 from rdkit.Chem import BRICS
 2
 3 data['mols'] = data['SMILES'].apply(Chem.MolFromSmiles)
 4 frags = set([])
 5 for mol in data['mols']:
 6     if mol is None: continue
 7     fragments = BRICS.BRICSDecompose(mol)    #A
 8     frags.update(fragments)
 9
10 frag_mols = []
11 for frag in frags:
12     frag_mol = Chem.MolFromSmiles(frag)
13     if frag_mol:
14         AllChem.Compute2DCoords(frag_mol)
15         frag_mols.append(frag_mol)
16
17 rld_new_mols = []
18 for mol in BRICS.BRICSBuild(frag_mols, maxDepth=1):    #B
19     if len(rld_new_mols) > 10000:    #C
20         break
21     try:
22         Chem.SanitizeMol(mol)
23         AllChem.Compute2DCoords(mol)
24         rld_new_mols.append(mol)
```

```
25     except:
26         continue    #D
27 #A Decompose molecules into BRICS fragments and store unique fragments.
28 #B Call generator to build new molecules from the BRICS fragments with maximum depth of 1. The maximum depth
29 #C Maximum depth greater than 3 will explode the number of potential molecules. Even with a depth of 1, we b
30 #D Skip molecules that cannot be sanitized.
31
```

## Exercise 2: Combinatorial & Retrosynthetic Library Design

Finish the implementation of a basic retrosynthetic analysis function that suggests possible precursors for a given string.

Test cases are provided to illustrate example inputs and outputs.

### ⌄ Student Solution to Exercise 2

Provide your solution to the above exercise in this cell and/or immediately following cells.

```
 1 def retrosynthetic_analysis(molecule):
 2     # TODO: your implementation
 3     """
 4     Performs a simplified retrosynthetic analysis on a given molecule string.
 5     Returns all possible two-fragment combinations that could form the target molecule.
 6
 7     Args:
 8         molecule (str): A string representing the target molecule
 9
10     Returns:
11         list: A list of tuples, where each tuple contains two strings representing
12             possible precursor fragments (A, B) that could form the target molecule
13
14     Example:
15         >>> retrosynthetic_analysis("ABC")
16         [('A', 'BC'), ('AB', 'C')]
17     """
18     if not molecule or len(molecule) < 2:
19         return []
20
21     # Generate all possible two-fragment combinations
22     precursors = []
23
24     # Iterate through all possible splitting points
25     for i in range(1, len(molecule)):
26         # Split molecule into two fragments
27         fragment1 = molecule[:i]
28         fragment2 = molecule[i:]
29
30         # Add the precursor combination to our list
31         precursors.append((fragment1, fragment2))
32
33     return precursors
```

```
 1     # Test case 1: Basic splitting
 2     assert retrosynthetic_analysis("ABC") == [('A', 'BC'), ('AB', 'C')], "Test case 1 failed"
 3
 4     # Test case 2: Empty string
 5     assert retrosynthetic_analysis("") == [], "Test case 2 failed"
 6
 7     # Test case 3: Single character
 8     assert retrosynthetic_analysis("A") == [], "Test case 3 failed"
 9
10     # Test case 4: Two characters
11     assert retrosynthetic_analysis("AB") == [('A', 'B')], "Test case 4 failed"
12
13     # Test case 5: Longer string
14     assert retrosynthetic_analysis("ABCD") == [
15         ('A', 'BCD'),
16         ('AB', 'CD'),
17         ('ABC', 'D')
18     ], "Test case 5 failed"
19
20     print("All test cases passed!")

All test cases passed!
```

### ⌄ An Introduction to Clustering

In this section, we will explore an important category of unsupervised machine learning models: clustering algorithms. Clustering algorithms can be used to produce an optimal division of our compound library into different clusters, where compounds within a cluster are more similar to each other than to compounds in other clusters.

In clustering, our model segments the data into similar subsets (where each subset is a cluster). The clustering model learns its own definition of classes, such that each cluster corresponds to a different class, and it learns which cluster to assign each data instance to.

One or more representative compounds can then be selected from each cluster to form our final, diverse compound library. Key aspects of cluster-based compound selection include:

1. How we choose to represent the compounds, e.g., descriptors or fingerprints.
2. What metric we use to assess the distance between each compound, given its representation.
3. Our choice of clustering algorithm to segment the compounds.
4. How we select a subset of compounds from each cluster depending on our problem requirements.

For simplicity, we'll default to Morgan fingerprints of radius 2 and we'll use the Tanimoto similarity coefficient, which we can use as a distance by subtracting the similarity coefficient from 1. A popular clustering algorithm for cheminformatics applications is Taylor-Butina clustering. We'll use both Taylor-Butina clustering and agglomerative clustering throughout the next section of the lab. The details of Taylor-Butina clustering and agglomerative clustering are out-of-scope for this course. Below, we use agglomerative clustering and Taylor-Butina clustering to sample a subset of molecules from the retrosynthetically constructed library.

The following code block defines useful functions for:

1. Converting molecules to fingerprints

   - The description by Butina for the Taylor-Butina clustering algorithm uses Daylight fingerprints, whose closest analog in the RDKit library is RDKit fingerprints of max path 5. This is why we use a different type of fingerprints.

2. Sampling compounds from each cluster

   - We will sample an equivalent number of compounds from each cluster to ensure that the sampled subset is representative of the original data set distribution. Within each cluster, we will sample the `n_samples` compounds that are closest to the cluster center, producing groupings of similar molecules within the sampled subset that will be easier for a downstream model to provide insights into structure-activity relationships.

3. Taylor-Butina clustering

```
1 # Convert RDKit molecules to Morgan fingerprints
2 def molecules_to_fingerprints(mols, radius=2, nBits=2048):
3     rdkit_gen = rdFingerprintGenerator.GetRDKitFPGenerator(maxPath=5) # Daylight FP approximation
4     fingerprints = [rdkit_gen.GetFingerprint(mol) for mol in mols]
5
6     np_fps = []
7     for fp in fingerprints:
8         arr = np.zeros((1,))
9         DataStructs.ConvertToNumpyArray(fp, arr)
10         np_fps.append(arr)
11     return np.array(np_fps), fingerprints
12
13 def sample_from_clusters(fps, clusters, n_samples=3):
14     sampled_mols = []
15     cluster_centers = []
16
17     for cluster in clusters:
18         if len(cluster) < n_samples:
19             sampled_mols.extend(cluster)
20         else:
21             cluster_fps = [fps[i] for i in cluster]
22             cluster_center = np.mean(cluster_fps, axis=0).reshape(1, -1)
23             cluster_centers.append(cluster_center)
24
25             distances = cdist(cluster_fps, cluster_center, metric='euclidean')
26             closest_indices = np.argsort(distances.flatten())[:n_samples]
27             sampled_mols.extend([cluster[i] for i in closest_indices])
28
29     return sampled_mols
30
31 def butina_clustering(fingerprints, cutoff=0.2):
32     dists = []
33     nfps = len(fingerprints)
34     for i in range(1, nfps):
35         sims = DataStructs.BulkTanimotoSimilarity(fingerprints[i], fingerprints[:i])
36         dists.extend([1 - x for x in sims])
37     cluster_data = Butina.ClusterData(dists, nfps, cutoff, isDistData=True)
38     return cluster_data
```

The next two code blocks fit agglomerative clustering and Taylor-Butina clustering models, respectively, and sample separate sets of compounds for each methodology. Each code block may take a couple of minutes.

```
1 fps, rdkit_fps = molecules_to_fingerprints(rld_new_mols)
2
3 agg_model = AgglomerativeClustering(n_clusters=16, linkage='ward')
4 agg_model.fit(fps)
5 agg_clusters = [np.where(agg_model.labels_ == i)[0].tolist() for i in np.unique(agg_model.labels_)]
6
7 agg_sampled_indices = sample_from_clusters(fps, agg_clusters)
8 agg_sampled_molecules = [rld_new_mols[i] for i in agg_sampled_indices]
```

```
1 clusters = butina_clustering(rdkit_fps, cutoff=0.4)
2 cluster_labels = np.zeros(len(rld_new_mols))
3 for cluster_id, cluster in enumerate(clusters):
4     for mol_id in cluster:
5         cluster_labels[mol_id] = cluster_id
6
7 butina_sampled_indices = sample_from_clusters(fps, clusters)
8 butina_sampled_molecules = [rld_new_mols[i] for i in butina_sampled_indices]
9
```

## Exercise 3: Sorting in Taylor-Butina Clustering

In Taylor-Butina clustering, the following steps play critical roles in the algorithms performance:

1. Computing the pairwise similarity between all compounds.
2. Calculting the number of neighbors for each compound, where a neighbor is defined as another compound with a similarity that is above a specified threshold.
3. Sorting the compounds by their number of neighbors.

In this exercise, we implement two key functions:

1. tanimoto_similarity: Calculate the Tanimoto similarity between two binary fingerprints.
2. merge_sort: Implement the merge sort algorithm to sort fingerprints by their number of neighbors.

We then use these functions to analyze a set of molecular fingerprints.

**Part 1: Tanimoto Similarity**

Review the already implemented tanimoto_similarity function to calculate the similarity between two binary fingerprints. Recall that Tanimoto similarity of two binary fingerprints is equal to: (intersection of bits) / (union of bits)

**Part 2: Merge Sort**

Implement the merge_sort function to sort fingerprints by their number of neighbors in descending order. Hint: The merge function has already been implemented, which you can then use in the recursive merge_sort function.

**Part 3: Analysis**

Use your implemented functions to analyze the given set of molecular fingerprints by combining the Tanimoto similarity calculation with sorting. Specifically, calculate the number of neighbors for each fingerprint and sort them using your merge sort function. This functionality is already implemented for you.

⌄ Student Solution to Exercise 3

Provide your solution to the above exercise in this cell and/or immediately following cells.

```
1 import random
2
3 # Helper function to generate random fingerprints of a more manageable length
4 def generate_fingerprint(length=64):
5     return ''.join(random.choice('01') for _ in range(length))
6
7 # Part 1: Implement Tanimoto Similarity
8 def tanimoto_similarity(fp1, fp2):
9     """
10    Calculate the Tanimoto similarity between two binary fingerprints.
11
12    :param fp1: First fingerprint (string of '0's and '1's)
13    :param fp2: Second fingerprint (string of '0's and '1's)
14    :return: Tanimoto similarity (float between 0 and 1)
15    """
16    if len(fp1) != len(fp2):
17      raise ValueError("Fingerprints must be of equal length")
18
```

```python
19      # Convert binary strings to sets of indices where bits are 1
20      set1 = set(i for i, bit in enumerate(fp1) if bit == '1')
21      set2 = set(i for i, bit in enumerate(fp2) if bit == '1')
22
23      # Calculate intersection and union
24      intersection = len(set1 & set2)
25      union = len(set1 | set2)
26
27      # Handle edge case where both fingerprints are all zeros
28      if union == 0:
29          return 1.0 if fp1 == fp2 else 0.0
30
31      return intersection / union
32
33  # Part 2: Implement Merge Sort
34  def merge(left, right, key_func):
35      """
36      Merge two sorted lists into a single sorted list.
37
38      :param left: Left sorted list
39      :param right: Right sorted list
40      :param key_func: Function to compute the sort key
41      :return: Merged sorted list
42      """
43      result = []
44      i = j = 0
45
46      while i < len(left) and j < len(right):
47          if key_func(left[i]) >= key_func(right[j]):
48              result.append(left[i])
49              i += 1
50          else:
51              result.append(right[j])
52              j += 1
53
54      # Add remaining elements
55      result.extend(left[i:])
56      result.extend(right[j:])
57      return result
58
59  def merge_sort(arr, key_func):
60      """
61      Sort the given array using the merge sort algorithm.
62
63      :param arr: Array to be sorted
64      :param key_func: Function to compute the sort key
65      :return: Sorted array
66      """
67      # TODO: Implement merge sort algorithm
68      # Base case: arrays of length 0 or 1 are already sorted
69      if len(arr) <= 1:
70          return arr
71
72      # Split array into two halves
73      middle = len(arr) // 2
74
75      # Recursively sort both halves
76      left = merge_sort(arr[:middle], key_func)
77      right = merge_sort(arr[middle:], key_func)
78
79      # Merge the sorted halves
80
81      return merge(left, right, key_func)
82
83
84  # Part 3: Analysis
85  def analyze_fingerprints(fingerprints, similarity_threshold):
86      """
87      Analyze the given fingerprints using Tanimoto similarity and sort them by number of neighbors.
88
89      :param fingerprints: List of fingerprints to analyze
90      :param similarity_threshold: Threshold for considering fingerprints as neighbors
91      :return: Sorted list of (fingerprint, neighbor_count) tuples
92      """
93      neighbor_counts = []
94      for i, fp1 in enumerate(fingerprints):
95          count = sum(1 for j, fp2 in enumerate(fingerprints)
96                      if i != j and tanimoto_similarity(fp1, fp2) >= similarity_threshold)
97          neighbor_counts.append((fp1, count))
98
99      return merge_sort(neighbor_counts, key_func=lambda x: x[1])
```

```python
1  # Test the implementation
2  def run_tests():
3      print("Running tests...")
4
5      # Test Tanimoto similarity
6      assert abs(tanimoto_similarity('1010', '1100') - 0.3333) < 0.0001, "Tanimoto similarity test failed"
7      assert tanimoto_similarity('11110000', '00001111') == 0, "Tanimoto similarity test failed"
8      assert tanimoto_similarity('11111111', '11111111') == 1, "Tanimoto similarity test failed"
9      assert tanimoto_similarity('10101010', '10101010') == 1, "Tanimoto similarity test failed"
10     assert abs(tanimoto_similarity('10101010', '11110000') - 0.3333) < 0.0001, "Tanimoto similarity test fai
11
12     # Test merge sort
13     test_arr = [(1, 5), (2, 3), (3, 8), (4, 1)]
14     sorted_arr = merge_sort(test_arr, key_func=lambda x: x[1])
15     assert sorted_arr == [(3, 8), (1, 5), (2, 3), (4, 1)], "Merge sort test failed"
16
17     # Test with strings
18     test_arr = [('a', 3), ('b', 1), ('c', 4), ('d', 2)]
19     sorted_arr = merge_sort(test_arr, key_func=lambda x: x[1])
20     assert sorted_arr == [('c', 4), ('a', 3), ('d', 2), ('b', 1)], "Merge sort test with strings failed"
21
22     # Test analyze_fingerprints
23     test_fingerprints = ['1010', '1100', '0011', '1111']
24     result = analyze_fingerprints(test_fingerprints, 0.5)
25     assert result == [('1111', 3), ('1010', 1), ('1100', 1), ('0011', 1)], "Analyze fingerprints test failed
26
27     print("All tests passed!")
28
29 # Generate a set of random fingerprints
30 num_fingerprints = 20
31 fingerprint_length = 64
32 fingerprints = [generate_fingerprint(fingerprint_length) for _ in range(num_fingerprints)]
33
34 # Set similarity threshold
35 similarity_threshold = 0.7
36
37 # Analyze fingerprints
38 result = analyze_fingerprints(fingerprints, similarity_threshold)
39
40 # Print results
41 print(f"Fingerprints sorted by number of neighbors (similarity threshold: {similarity_threshold}):")
42 for fp, count in result:
43     print(f"Fingerprint: {fp}, Neighbors: {count}")
44
45 run_tests()
```

```
Fingerprints sorted by number of neighbors (similarity threshold: 0.7):
Fingerprint: 1111111010100101100110110111100001011011000011001101100110011001, Neighbors: 0
Fingerprint: 1010101110000010110110011011001100001010000000011100110001001111000, Neighbors: 0
Fingerprint: 1010100111001000011100010100001010010001111111000100000101100010, Neighbors: 0
Fingerprint: 1000100011011100011011011110000110000001001000111000010000111011, Neighbors: 0
Fingerprint: 1111100100011111001001000011110011010101000010010000010100101100, Neighbors: 0
Fingerprint: 1011101110010010001000011001111011111101010100011010101000100100, Neighbors: 0
Fingerprint: 1110110001110011011100011110101001000000100101101000100100110111, Neighbors: 0
Fingerprint: 1001011111101100011011111000010000100001001010011000111011110, Neighbors: 0
Fingerprint: 0001001111101011011011010110010000100000111101100100000001001100, Neighbors: 0
Fingerprint: 0001011000010100101000011001111101110100010000011010000110000010, Neighbors: 0
Fingerprint: 1110111011010011010111011110011111001010000111111000011001101101011, Neighbors: 0
Fingerprint: 0000001101010110010011110110000110001000000101101110100000101100, Neighbors: 0
Fingerprint: 1101111111001111000001111100110101110001001100010101001011101001, Neighbors: 0
Fingerprint: 1011110101110101000000000001101000010011100100011010011010101001, Neighbors: 0
Fingerprint: 1111110100010001100100101010010011110010000000001001111101001100010, Neighbors: 0
Fingerprint: 1000001101110001100101111000011101100000110100011110110010110, Neighbors: 0
Fingerprint: 1101110000000001111101011000101011111100110000111010101100010100, Neighbors: 0
Fingerprint: 0011011110011000001100101011001111010011111110100100100000011110, Neighbors: 0
Fingerprint: 0011100110011001011010100111110011000011111000001111110101111000, Neighbors: 0
Fingerprint: 1001001101111110110101010010110100000001110010010001101110001111, Neighbors: 0
Running tests...
All tests passed!
```

## ⌄ Comparing Clustering Methods

Starting with our retrosynthetically constructed library, we've sampled two subsets of compounds with agglomerative clustering and Taylor-Butina clustering. The following code block compares the two subsets based on two dimensions.

As stated several times, we really care about the sampled sets' diversity. To measure diversity, we calculate the average pairwise Tanimoto similarity for the sampled sets. The lower the average similarity, the more diverse the set. The diversity of the Taylor-Butina sampled set was greater than sampled via agglomerative clustering.

We also care about redundancy, as compounds that are likely to share the same activity are redundant and not as valuable for our downstream purposes. Following from literature, compound pairs with greater than 0.85 Tanimoto similarity have an 80% chance of

sharing same activity. We calculate the redundancy ratio for each set, which is the ratio of the number of redundant pairs (pairs with a Tanimoto similarity greater than 0.85) to the total number of unique pairs in the set. The lower the redundancy ratio, the less redundant the set is. The redundancy of the Taylor-Butina sampled set was much lower than that achieved via agglomerative clustering.

Note: Our source for the 0.85 cut off is not perfect since it was calculated only with respect to Daylight fingerprints. We use it more for convenience, with the caveat that there is no shortage of literature that examines different trends of redundancy (pairs with X similarity have Y% chance of being active given some fingerprint representation, similarity metric, and benchmark data set)

```
1 # Calculate Tanimoto distance matrix
2 def calculate_tanimoto_distances(fingerprints):
3     n = len(fingerprints)
4     tanimoto_distances = np.zeros((n, n))
5     for i in range(n):
6         for j in range(i + 1, n):
7             tanimoto_distances[i, j] = 1 - DataStructs.FingerprintSimilarity(fingerprints[i], fingerprints[j
8             tanimoto_distances[j, i] = tanimoto_distances[i, j]
9     return tanimoto_distances
10
11 # Calculate the diversity of the sampled sets based on pairwise Tanimoto similarity
12 def calculate_diversity(sampled_indices, fingerprints):
13     sampled_fps = [fingerprints[i] for i in sampled_indices]
14     tanimoto_distances = calculate_tanimoto_distances(sampled_fps)
15     avg_tanimoto_similarity = 1 - np.mean(tanimoto_distances)
16     return avg_tanimoto_similarity
17
18 agg_diversity = calculate_diversity(agg_sampled_indices, rdkit_fps)
19 butina_diversity = calculate_diversity(butina_sampled_indices, rdkit_fps)
20
21 print(f"Agglomerative clustering sampled set diversity (lower is better): {agg_diversity}")
22 print(f"Taylor-Butina clustering sampled set diversity (lower is better): {butina_diversity}")
23
24 if agg_diversity < butina_diversity:
25     print("The Agglomerative clustering sampled set is more diverse.")
26 else:
27     print("The Taylor-Butina clustering sampled set is more diverse.")
```

```
Agglomerative clustering sampled set diversity (lower is better): 0.24658869026751318
Taylor-Butina clustering sampled set diversity (lower is better): 0.15754030228878224
The Taylor-Butina clustering sampled set is more diverse.
```

```
1 # Calculate the redundancy of the sampled sets based on pairwise Tanimoto similarity
2 def calculate_redundancy(sampled_indices, fingerprints, threshold=0.85):
3     sampled_fps = [fingerprints[i] for i in sampled_indices]
4     tanimoto_matrix = np.zeros((len(sampled_fps), len(sampled_fps)))
5     for i, fp1 in enumerate(sampled_fps):
6         for j, fp2 in enumerate(sampled_fps):
7             if i != j:
8                 similarity = DataStructs.TanimotoSimilarity(fp1, fp2)
9                 tanimoto_matrix[i, j] = similarity
10    redundant_pairs = np.where(tanimoto_matrix > threshold)
11    num_redundant_pairs = len(redundant_pairs[0])
12    num_unique_pairs = len(sampled_fps) * (len(sampled_fps) - 1) / 2
13    redundancy_ratio = num_redundant_pairs / num_unique_pairs
14    return redundancy_ratio
15
16 agg_redundancy = calculate_redundancy(agg_sampled_indices, rdkit_fps)
17 butina_redundancy = calculate_redundancy(butina_sampled_indices, rdkit_fps)
18
19 print(f"Agglomerative clustering sampled set redundancy ratio (lower is better): {agg_redundancy}")
20 print(f"Taylor-Butina clustering sampled set redundancy ratio (lower is better): {butina_redundancy}")
21
22 if agg_redundancy < butina_redundancy:
23     print("The Agglomerative clustering sampled set is less redundant.")
24 else:
25     print("The Taylor-Butina clustering sampled set is less redundant.")
```

```
Agglomerative clustering sampled set redundancy ratio (lower is better): 0.06382978723404255
Taylor-Butina clustering sampled set redundancy ratio (lower is better): 0.0007370316340075438
The Taylor-Butina clustering sampled set is less redundant.
```

## Big O in Machine Learning

Another comparison point for ML algorithms is their time and memory complexity. To choose methods that can handle large data sets, especially when processing time and memory usage are constrained, we can use Big O notation to profile how ML algorithms will scale with increasing data size or dimensionality. Let's breakdown the complexity for agglomerative clustering, K-means clustering, and Taylor-Butina clustering. In the below, $n$ is the number of data points, $d$ is the number of dimensions (features), $i$ is the number of iterations the algorithm runs for, and $k$ is the number of centroids.

Agglomerative Clustering: We compute pairwise distances between each pair of points ($O(d)$), which is $n^2$ computations for the initial distance matrix and, in worst-cast, an additional computation in each of the $n-1$ merges. The pairwise distance matrix takes up $n^2$ space regardless of dimensionality. - Time complexity: $O(n^3 * d)$ - Memory complexity: $O(n^2)$

K-means Clustering: Distances are computed between each centroid and each data point, and each distance calculation takes $O(d)$ time. We compute these distances for each of $i$ iterations. We store each data point and each centroid in memory. - Time complexity: $O(k * n * i * d)$ - Memory complexity: $O((n + k) * d)$

Taylor-Butina Clustering: We computer pairwise distances between each pair of points. Each pairwise distance calculation takes $O(d)$ time. It only needs to store each of the $n$ data points and their $d$ values. - Time complexity: $O(n^2 * d)$ - Memory complexity: $O(n * d)$

Note that these are general complexities and can vary with specific implementations or optimizations. For instance, agglomerative clustering under certain linkage criteria has a time complexity of $O(n^2 log n * d)$.

## Fragment-based Drug Discovery

The application of fragments extends beyond compound library construction, encompassing fragment-based drug discovery (FBDD). Compared to the traditional approach of screening against a compound library to identify hits and optimize those hits into leads, FBDD screens against a fragment library to identify multiple fragments with binding affinity to a target. Once fragments are identified, they are expanded, merged, and optimized to yield lead compounds with stronger binding profiles.
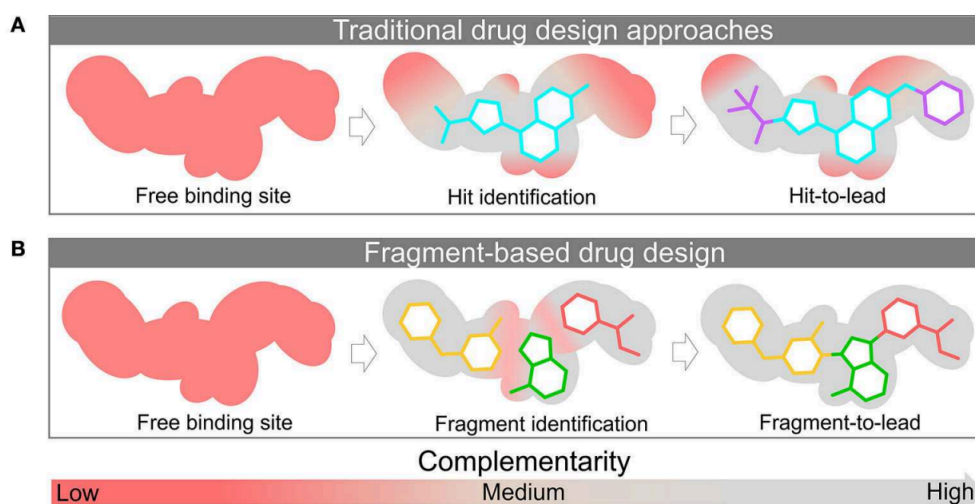
Screening against a fragment library implies that there are some rules for what properties qualify a molecule as a fragment. Just as we have rule-based guidelines (e.g., Lipinski's Ro5) for whether a molecule is "drug-like," fragment libraries generally conform to a "Rule of 3":

1. Less than 300 Da molecule weight,
2. 3 or less H bond donors,
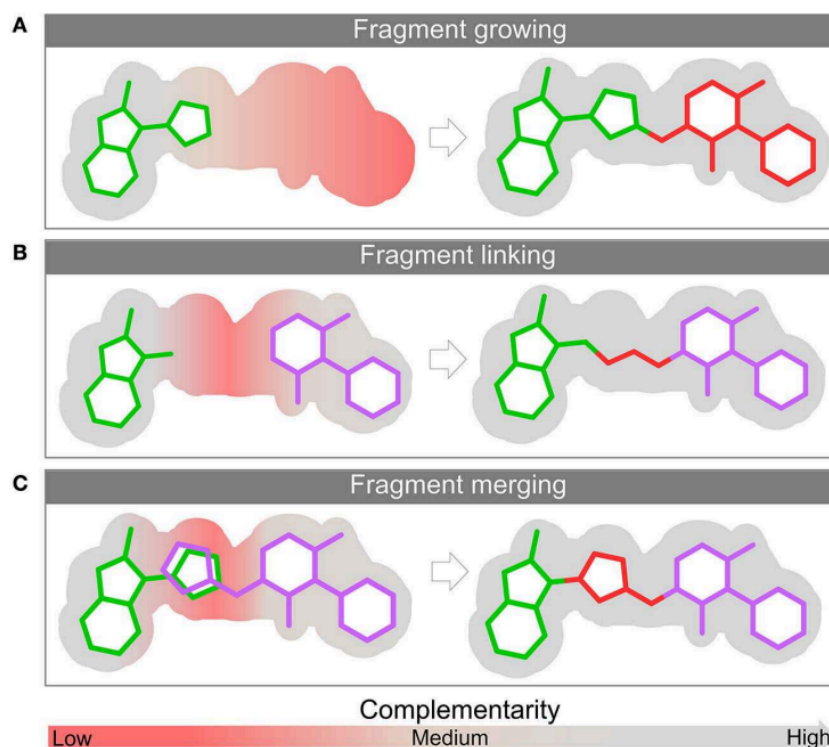3. 3 or less H bond acceptors,
4. 3 or less CLogP

Compared to compound screening libraries, which constitute millions of compounds, fragment libraries tend towards less than 1,000 fragments. These fragments are less complex than the larger, drug-like compounds they are derived from, which is advantageous because the fragments can access and explore more of the potential binding sites of the target protein. Thus, even though the library contains only 1,000 fragments, each fragment is able to probe the target protein's structure more extensively than a full-fledged compound. The resulting higher hit rate and increased coverage of chemical space contribute to greater monetary efficiency.

Each fragment is weakly binding, so once our screen has run its course and identified fragments, we need conflate these fragments into actionable leads with high binding affinity. Several methods for doing so are visualized below (figure source: de Souza Neto, L. R., Moreira-Filho, J. T., Neves, B. J., Maidana, R. L. B. R., Guimarães, A. C. R., Furnham, N., Andrade, C. H., & Silva, F. P., Jr (2020). In silico Strategies to Support Fragment-to-Lead Optimization in Drug Discovery. Frontiers in chemistry, 8, 93. https://doi.org/10.3389/fchem.2020.00093).



**FIGURE 2 |** Discovery and structural-optimization of drug-like molecules **(A)** and fragments **(B)** using protein target information. The surface represents the binding site. The red and gray colors represent the level of complementarity of ligand with the active site. Pockets with low complementarity with ligand are colored in red; pockets with high complementarity with ligand are highlighted in gray.

**FIGURE 3 |** Fragment optimization approaches: fragment growing **(A)**, fragment linking **(B)**, and fragment merging **(C)**. The surface of the binding site is depicted in gray. The red and gray colors represent the level of complementarity of ligand with the active site. Pockets with low complementarity with ligand are colored in red; pockets with high complementarity with ligand are highlighted in gray.

Once we've expanded the fragment hits into a compound, are we done? It's not that simple. We want the resulting compound from growing, linking, or merging the fragments to have a better binding profile than the individual fragments. For the majority of the compounds we'll derive from fragments, this won't be the case and they won't be suitable leads. The fragments are, individually, weak binders and their experimental binding is often with respect to a crystallized, static structure that doesn't reflect how that structure, and resulting binding poses, change at non-cryogenic conditions (e.g., room temperature). Merging the fragments doesn't make up for the unreliable, weak binding of each fragment in isolation.
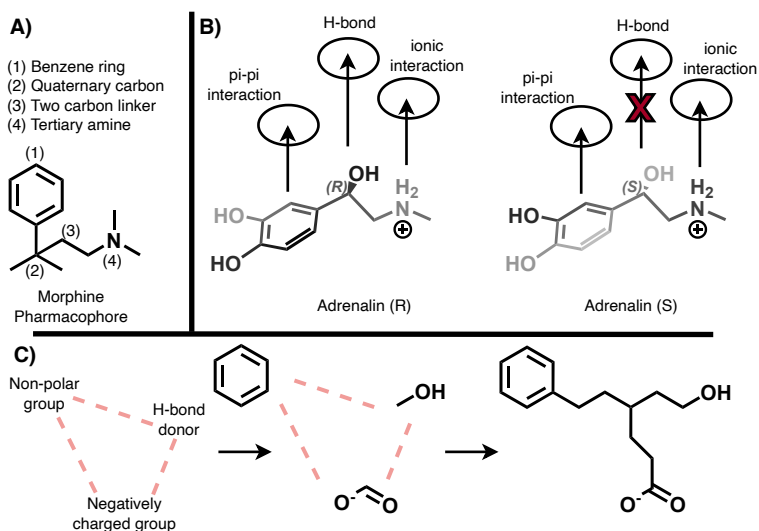
So now what?

## Pharmacophore Modeling

Let's take a step back and think about drug-target complementarity. Intermolecular forces, such as hydrogen bonding, hold a drug to its target. Intermolecular forces will not form unless the drug is positioned such that its functional groups interact with the target's functional groups in a way that facilitates binding. A good drug candidate that binds strongly and is biologically active to our target needs, at minimum, (1) the right functional groups with (2) the correct geometry (distances, angles, orientations).

A drug's pharmacophore represents the proper functionality and geometry for the drug to be biologically active. Pharmacophores are a useful abstraction, representing the ensemble of functional groups and their ideal spacing that is necessary to ensure optimal interaction with a specific biological target structure to trigger (or block) its biological response.
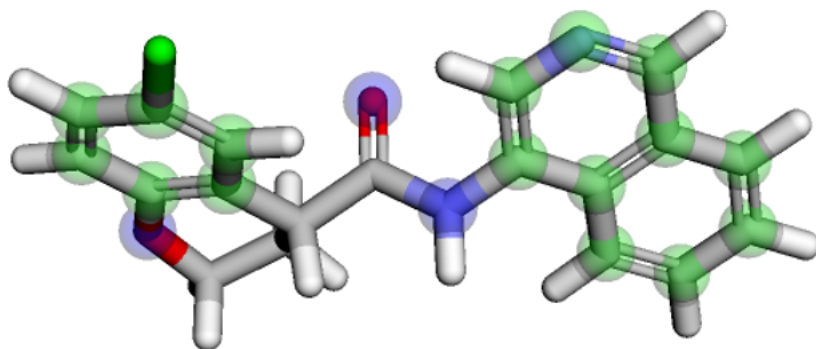
Consider the figure below. (A) shows the pharmacophore for morphine and (B) illustrates the explanation for why one enantiomer of epinephrine is much more potent in humans (adrenalin) while the other has weak potency since its orientation does not enable H-bond contact.

**A)**
(1) Benzene ring
(2) Quaternary carbon
(3) Two carbon linker
(4) Tertiary amine

Morphine Pharmacophore

**B)**
Adrenalin (R)

Adrenalin (S)

**C)**
Non-polar group
H-bond donor
Negatively charged group

We can intuit pharmacophoric features as describing a compound's preferences for specific molecular interactions, such as H-bond acceptors, H-bond donors, hydrophobic areas, aromatic interactions, coordination to metal ions, and so on. As we've done in the past with other descriptors or features, we can annotate molecules with all interaction features possible, either in a vacuum (ligand-based) or with respect to binding to a specific target (structure-based). With enough data, pharmacophore modeling can identify essential interaction features that are responsible for activity and use important pharmacophores to score compounds from a screening database.

How do pharmacophores relate to FBDD? Consider the design of a pharmacophore as shown in (C) of the previous figure. To satisfy the pharmacophore's requirements, we might select the shown functional groups, and connect them to produce a compound that contains the pharmacophore. This is analogous to how we might screen fragment hits and consider ways to expand or merge them into a compound containing those fragments.

Pharmacophores include many key functional groups, often more than 3 and not coplanar, so they exist in three-dimensional space (though we can compute either 2D or 3D pharmacophoric features). This relates back to the importance of conformational flexibility. For instance, a compound with a high number of rotatable bonds may have many probable conformations with improper spacing of the functional groups, which minimizes the probability that the functional groups will match the desired pharmacophore. The next figure highlights 3D pharmacophoric features (H-bond donors in red, H-bond acceptors in blue, aromatics in green) for one possible conformer sampled from an arbitrary molecule.



For a given fragment hit, we can decompose it into its pharmacophores. If we repeat this for every fragment that was a hit during experimental screening, we will now have a probability distribution of pharmacophores in relation to the target structure of interest, which is advantageous compared to the traditional merging approaches that think of each fragment one-by-one. We can then use unsupervised learning (density estimation) to identify significant pharmacophores from the ensemble of fragment hits and filter out noise from weakly binding fragments.

## Density Estimation

Our plan is to assess the pharmacophoric features of all of the fragments that were hits in the experimental screening step. We'll then develop a model to estimate the probability distribution underlying the influence of different pharmacophores at different locations in the protein's binding site. If we have enough data points for the model to adequately learn this distribution, we can take a vast library of compounds, annotate their pharmacophores, and score how well those compounds fit the estimated pharmacophore distribution with respect to the protein binding site (i.e., structure-based virtual screening).

Density estimation is an unsupervised learning technique that will help us model the unknown probability distribution from which a given data set has been sampled. Imagine we were trying to figure out how crowded different parts of a city are. We might look at a
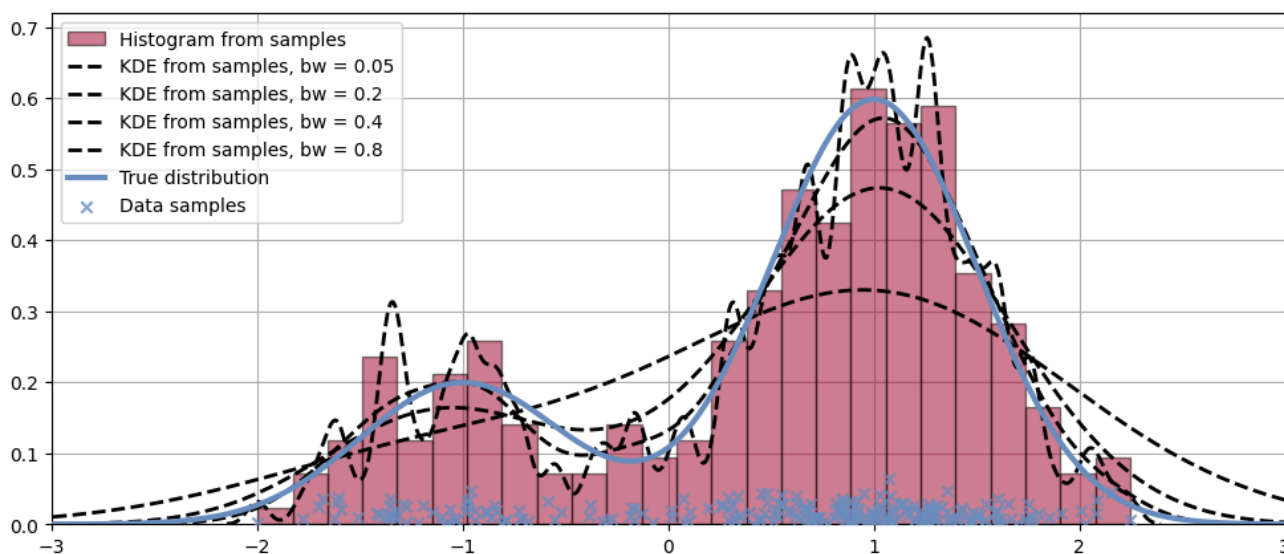
map and estimate that the city is more densely populated in areas where high-rise offices or apartment buildings are packed closely together, compared to single-family homes spread out in suburbs on the city's outer edge. We might then use these estimates to decide the appropriate location to build an upscale restaurant that caters to wine-and-dine business expenses and wants to maximize proximity to corporate clientele.

Instead of buildings, we have data points comprised of individual compounds. Likewise, instead of the building's height, we have features such as a compound's regions of intermolecular interactions, and instead of the building's location in the city, we have the spatial arrangement and geometry of those intermolecular interactions (which, combined, define the compound's pharmacophores). In summary, density estimation models how "crowded" different segments of our data set are.

A detailed description of density estimation is outside of the scope of this course. However, if you've ever worked with a histogram, that is an example of a density estimator. A histogram groups data into defined intervals and represents the number of data points that fall into each interval as a vertical bar whose height corresponds to the number of data points in that interval. Problematically, histograms can change significantly depending on the number of intervals and the length of the intervals, resulting in different interpretations of the same data set. Histograms also discretize the data into intervals, which obscures the underlying continuous nature of the data distribution into a jagged, step-like representation and may introduce artifacts at the edges of the interval boundaries.

While we won't explain the mechanics, we will use something similar to a histogram density estimator -- called a kernel density estimator (KDE). However, instead of the jagged, step-like estimation of a histogram, KDE provides a more continuous estimation as visualized in the figure below.



## Exercise 4: Sliding Windows and KDE

The sliding window pattern is a technique in data processing where a fixed-size "window" moves over a dataset, analyzing or transforming only the data within that window at each step. It's particularly useful for time series analysis, signal processing, and sequence data. In relation to KDE, the sliding window concept is similar to how KDE uses a kernel function to estimate the probability density at each point, with the kernel effectively acting as a window that weighs nearby points more heavily. Both techniques involve local calculations that move across the entire dataset, allowing for efficient analysis of large datasets by focusing on subsets at a time.

In this exercise, you will implement a sliding window algorithm to estimate the local density of data points in a time series. This technique is useful in various scientific applications, including molecular dynamics simulations and signal processing.

The task can decomposed into the following steps, with several steps (partially) implemented for you already:

1. Implement the `sliding_window_density function` that uses a sliding window to estimate the local density of data points.
2. The function should take a list of (timestamp, value) tuples, a window size, and a step size as input.
3. For each window, calculate the density as the number of points in the window divided by the time span of the window.
4. Return a list of (window_center_time, density) tuples.

⌄   Student Solution to Exercise 4

Provide your solution to the above exercise in this cell and/or immediately following cells.

```
1 def generate_sample_data(n: int, time_range: Tuple[float, float], value_range: Tuple[float, float]) -> List
2     """Generate sample time series data."""
3     return sorted([
4         (random.uniform(*time_range), random.uniform(*value_range))
5         for _ in range(n)
6     ])
```

```
 7
 8 def sliding_window_density(data: List[Tuple[float, float]], window_size: float, step_size: float) -> List[T
 9     """
10     Estimate local density using a sliding window.
11
12     :param data: List of (timestamp, value) tuples
13     :param window_size: Size of the sliding window
14     :param step_size: Step size for sliding the window
15     :return: List of (window_center_time, density) tuples
16     """
17
18     # Validate inputs
19     if window_size <= 0:
20         raise ValueError("Window size must be positive")
21     if step_size <= 0:
22         raise ValueError("Step size must be positive")
23
24     # Handle empty data case
25     if not data:
26         return []
27
28     # Handle single data point case
29     if len(data) == 1:
30         return [(data[0][0], 1/window_size)] if window_size > 0 else [(data[0][0], float('inf'))]
31
32     result = []
33     start_time, end_time = data[0][0], data[-1][0]
34
35     window_start = start_time
36
37     # TODO: Main loop executes while window_start + window_size is <= end_time
38     while window_start + window_size <= end_time:
39     # TODO: Within the loop, conduct the following steps
40     #           Define window_end for the current window, which is the start of the window plus the window
41     #           Compute the window_center using window_start and window_end
42     #           Count the points in the current window, points_in_window. To do so, iterate over each (time
43     #           Calculate the density, which is equal to the points in the window divided by the window siz
44     #           Append the (window_center, density) tuple to result
45     #           Increment window_start by the step_size
46         window_end = window_start + window_size
47         window_center = (window_start + window_end) / 2
48         count = 0
49
50         for timestamp, value in data:
51             if timestamp >= window_start and timestamp < window_end: # window_start <= timestamp < window_end
52                 count += 1
53
54         density = count / window_size
55         result.append((window_center, density))
56
57         window_start += step_size
58
59     return result
60
61 def run_tests():
62     # Test case 1: Basic functionality
63     data1 = [(1, 2), (2, 3), (3, 4), (4, 5), (5, 6)]
64     result1 = sliding_window_density(data1, window_size=2, step_size=1)
65     expected1 = [(2.0, 1.0), (3.0, 1.0), (4.0, 1.0)]
66     assert result1 == expected1, f"Test case 1 failed. Expected {expected1}, but got {result1}"
67
68     # Test case 2: Varying density
69     data2 = [(1, 1), (1.1, 2), (1.2, 3), (2, 4), (3, 5), (4, 6)]
70     result2 = sliding_window_density(data2, window_size=1, step_size=0.5)
71     expected2 = [(1.5, 3.0), (2.0, 1.0), (2.5, 1.0), (3.0, 1.0), (3.5, 1.0)]
72     assert result2 == expected2, f"Test case 2 failed. Expected {expected2}, but got {result2}"
73
74     # Test case 3: Empty data
75     data3 = []
76     result3 = sliding_window_density(data3, window_size=1, step_size=0.5)
77     expected3 = []
78     assert result3 == expected3, f"Test case 3 failed. Expected {expected3}, but got {result3}"
79
80     # Test case 4: Single data point
81     data4 = [(1, 1)]
82     result4 = sliding_window_density(data4, window_size=1, step_size=0.5)
83     expected4 = [(1.0, 1.0)]
84     assert result4 == expected4, f"Test case 4 failed. Expected {expected4}, but got {result4}"
85
86     # Test case 5: Window larger than data range
87     data5 = [(1, 1), (2, 2), (3, 3)]
88     result5 = sliding_window_density(data5, window_size=5, step_size=1)
```

```
89      expected5 = []
90      assert result5 == expected5, f"Test case 5 failed. Expected {expected5}, but got {result5}"
91
92      print("All test cases passed!")
93
94 sample_data = generate_sample_data(100, (0, 10), (0, 1))
95 density_estimation = sliding_window_density(sample_data, window_size=2, step_size=0.5)
96 print("Sample density estimation results:")
97 for time, density in density_estimation[:10]:  # Print first 10 results
98      print(f"Time: {time:.2f}, Density: {density:.2f}")
99
100 # Test the implementation
101 run_tests()
```

```
Sample density estimation results:
Time: 1.16, Density: 10.00
Time: 1.66, Density: 9.50
Time: 2.16, Density: 9.50
Time: 2.66, Density: 8.50
Time: 3.16, Density: 9.50
Time: 3.66, Density: 10.00
Time: 4.16, Density: 12.50
Time: 4.66, Density: 14.00
Time: 5.16, Density: 10.50
Time: 5.66, Density: 10.50
All test cases passed!
```

## ⌄  Combining it all together: FRESCO

We have all the tools we need to estimate the pharmacophore distributions and use them to score molecules for continued investigation. Let's put it all together following the work on the package FRESCO (reference: McCorkindale, W., et al. (2022). Fragment-Based Hit Discovery via Unsupervised Learning of Fragment-Protein Complexes. bioRxiv 2022.11.21.517375; doi: https://doi.org/10.1101/2022.11.21.517375). Using the Fresco demo as a guide, we will discuss each step with attention to the internals that focus on pharmacophore modeling and density estimation. Much of the following code is adapted from FRESCO's GitHub repository: https://github.com/wjm41/fresco/tree/master

We start with pre-processed mpro_frags.sdf provided for demo purposes by FRESCO and note that the following code listings are inspired and modified from FRESCO implementations. The mpro_frags.sdf file contains aligned conformations of 23 fragments bound to SARS-CoV-2 Mpro. We annotate this fragment ensemble with 3D pharmacophore features. The produced pcore_df contains the pharmacophore features, along with 3D coordinates for each molecule. Note that we limit our interest to H-bond donor, H-bond acceptor, and aromatic pharmacophore features.

```
1 sdfFile = 'data/L02_mpro_frags.sdf'
2 df_fragments = PandasTools.LoadSDF(sdfFile, idName='name', smilesName='SMILES', molColName='mol')
```

```
1 def compute_pcore_data(mol, mol_id=0):
2     """
3     Compute pharmacophoric core (pcore) data for a given molecule.
4
5     Parameters:
6     - mol: RDKit molecule object
7     - mol_id: Identifier for the molecule
8
9     Returns:
10     - DataFrame containing pcore information
11     """
12     # Define pharmacophoric cores of interest
13     pcores_of_interest = ['Donor', 'Acceptor', 'Aromatic']
14
15     # Load feature definition file
16     fdefFile = os.path.join(RDDataDir, 'BaseFeatures.fdef')
17     featFactory = ChemicalFeatures.BuildFeatureFactory(fdefFile)
18
19     # Get atom coordinates and pharmacophore features
20     atom_coordinates = mol.GetConformer().GetPositions()
21     pharmacophore_features = featFactory.GetFeaturesForMol(mol)
22
23     # Collect DataFrame rows
24     rows = []
25     for pcore in pharmacophore_features:
26         pharmacophore_name = pcore.GetFamily()
27
28         if pharmacophore_name in pcores_of_interest:
29             atom_ids = pcore.GetAtomIds()
30             xyz = np.mean([atom_coordinates[id] for id in atom_ids], axis=0)
31             rows.append({
```

```
32                'pcore': pharmacophore_name,
33                'smiles': MolToSmiles(mol),
34                'mol_id': mol_id,
35                'coord_x': xyz[0],
36                'coord_y': xyz[1],
37                'coord_z': xyz[2],
38            })
39    return pd.DataFrame(rows)
40
41 mols = df_fragments['mol'].values
42 dfs = []
43 for mol_id, mol in tqdm(enumerate(mols), total=len(mols)):
44     dfs.append(compute_pcore_data(mol, mol_id))
45 pcore_df = pd.concat(dfs)
```

```
100%|██████████| 23/23 [00:00<00:00, 68.66it/s]
```

```
1 pcore_df.head()
```

| | pcore | smiles | mol_id | coord_x | coord_y | coord_z |
|---|---|---|---|---|---|---|
| 0 | Donor | COC(=O)c1ccc(S(N)(=O)=O)cc1 | 0 | 8.8612 | 5.1249 | 22.8984 |
| 1 | Acceptor | COC(=O)c1ccc(S(N)(=O)=O)cc1 | 0 | 13.0863 | -1.0956 | 24.5273 |
| 2 | Acceptor | COC(=O)c1ccc(S(N)(=O)=O)cc1 | 0 | 13.9905 | -0.3540 | 22.5499 |
| 3 | Acceptor | COC(=O)c1ccc(S(N)(=O)=O)cc1 | 0 | 11.1835 | 6.2738 | 23.3719 |
| 4 | Acceptor | COC(=O)c1ccc(S(N)(=O)=O)cc1 | 0 | 9.9526 | 5.4586 | 25.2863 |

Next steps: ( Generate code with pcore_df )   ( New interactive sheet )

We now assess each pair of fragments in the ensemble, and we measure the distance between each pair of pharmacophore features across the fragment pair. Since we are interested in three pharmacophore feature types, there are six possible pairwise pharmacophore feature types, and we calculate the pairwise distances between each instance of a pharmacophore feature pair. The smaller the distance between the pharmacophore feature in one fragment and that in another fragment, the closer those pharmacophores are to each other within the binding site.

```
1 def calculate_pairwise_distances_between_pharmacophores_for_fragment_ensemble(df_of_frag_ensemble_pcores, pc
2     """
3     Calculate pairwise distances between two specified pharmacophores within a fragment ensemble.
4
5     Parameters:
6     – df_of_frag_ensemble_pcores: DataFrame containing pharmacophoric data for fragment ensemble
7     – pcore_a: Name of the first pharmacophore
8     – pcore_b: Name of the second pharmacophore
9
10    Returns:
11    – A NumPy array of pairwise distances between the specified pharmacophores
12    """
13    # Filter DataFrame for the specified pharmacophores
14    df_pcore_a = df_of_frag_ensemble_pcores.query('pcore == @pcore_a')
15    df_pcore_b = df_of_frag_ensemble_pcores.query('pcore == @pcore_b')
16
17    unique_smiles_a = df_pcore_a['smiles'].unique()
18    # Initialize list to store distances
19    distances_for_all_pairs = []
20
21    for smile in unique_smiles_a:
22        # Coordinates of pcore_a in the current fragment
23        coords_a = df_pcore_a[df_pcore_a['smiles'] == smile][['coord_x', 'coord_y', 'coord_z']].to_numpy()
24
25        # Coordinates of pcore_b in different fragments (don't count distances within the same fragment)
26        coords_b = df_pcore_b[df_pcore_b['smiles'] != smile][['coord_x', 'coord_y', 'coord_z']].to_numpy()
27
28        if coords_b.size == 0:
29            continue  # Skip if there are no pcore_b in different fragments
30
31        # Calculate  and store pairwise distances
32        delta_coords = coords_a[:, np.newaxis, :] – coords_b[np.newaxis, :, :]
33        distances_for_this_pair = np.linalg.norm(delta_coords, axis=2).flatten()
34        distances_for_all_pairs.append(distances_for_this_pair)
35
36    return np.hstack(distances_for_all_pairs)
37
38 interesting_pcores = [
39     'Donor–Aromatic', 'Aromatic–Acceptor', 'Aromatic–Aromatic',
40     'Donor–Donor', 'Donor–Acceptor', 'Acceptor–Acceptor'
41 ]
```

```
42 frag_pcore_histogram = {}
43 frag_pcore_weight = {}
44 for pcore_pair in interesting_pcores:
45     core_a, core_b = pcore_pair.split('-')
46     frag_pcore_histogram[pcore_pair] = calculate_pairwise_distances_between_pharmacophores_for_fragment_ense
```

Once we have the distances for each possible pair of fragments and for each possible pairing of pharmacophores, we have enough information to construct a pharmacophore-pharmacophore pairwise distance distribution for each pharmacophore feature. We fit multiple KDE models to these distributions to learn the probability density of a pair of pharmacophores associated with fragment binding having a particular distance from each other in the binding site.

The following code block fits multiple KDE models and then selects the optimal KDE model. We visualize the final KDE models in a figure following the code block. Note that there is room for improvement and our best model is not perfect. The estimates for Aromatic-Aromatic and Donor-Acceptor are best, and those for Donor-Aromatic and Aromatic-Acceptor are permissible, albeit arguably too smooth. The estimates for Donor-Donor and Acceptor-Acceptor are noisy, with too small of a bandwidth.

```
 1 from sklearn.neighbors import KernelDensity
 2 from sklearn.model_selection import GridSearchCV
 3
 4 def fit_sklearn_pair_kde(data):
 5     params = {'bandwidth': np.logspace(-3, 3, 50)}
 6     grid = GridSearchCV(KernelDensity(kernel='gaussian', rtol=1e-4), params)
 7     grid.fit(data.reshape(-1, 1))
 8     kde = grid.best_estimator_
 9     return kde
10
11 fresco_kdes = {}
12 for pcore_pair in interesting_pcores:
13     fresco_kdes[pcore_pair] = fit_sklearn_pair_kde(frag_pcore_histogram[pcore_pair])
```
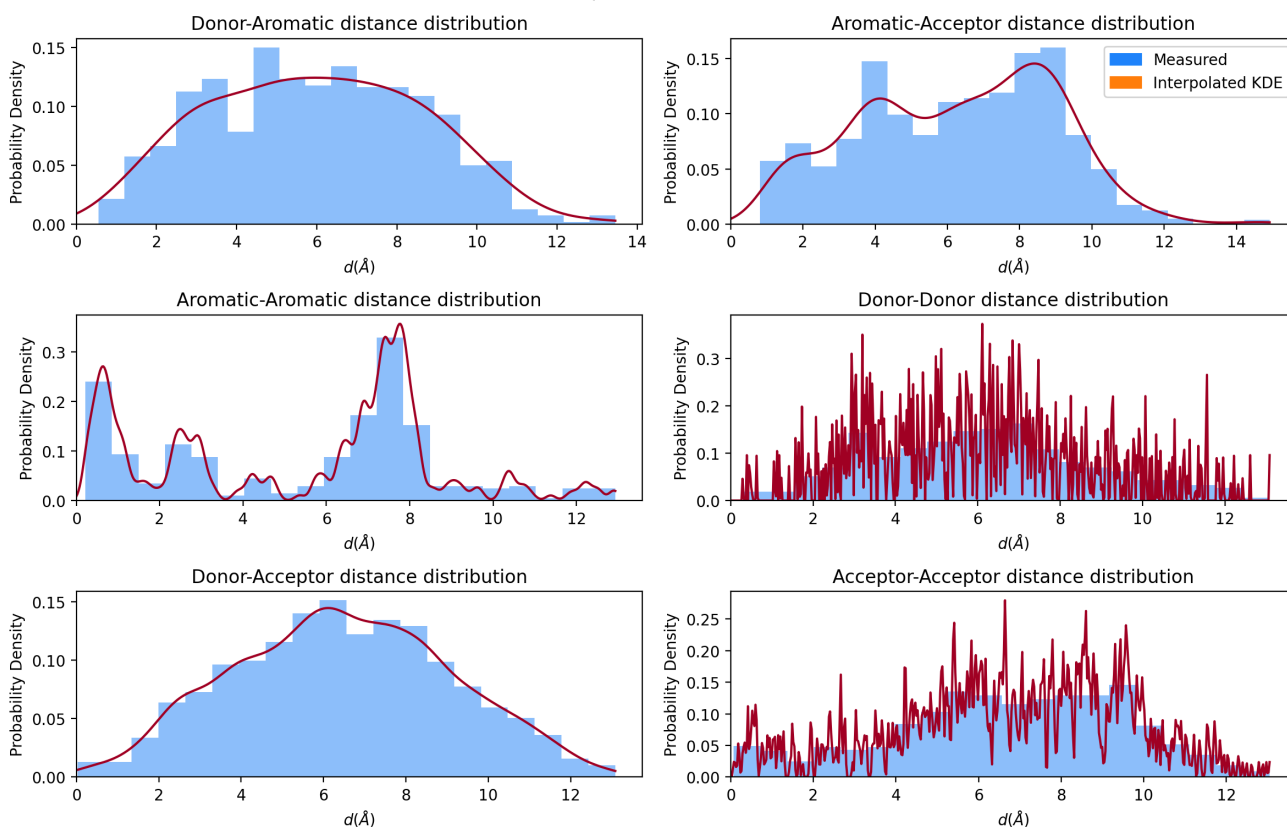
```
 1 fresco_kdes
```

```
{'Donor-Aromatic': KernelDensity(bandwidth=np.float64(0.868511373751352), rtol=0.0001),
 'Aromatic-Acceptor': KernelDensity(bandwidth=np.float64(0.655128556859551), rtol=0.0001),
 'Aromatic-Aromatic': KernelDensity(bandwidth=np.float64(0.1206792640639329), rtol=0.0001),
 'Donor-Donor': KernelDensity(bandwidth=np.float64(0.0071968567300115215), rtol=0.0001),
 'Donor-Acceptor': KernelDensity(bandwidth=np.float64(0.49417133613238334), rtol=0.0001),
 'Acceptor-Acceptor': KernelDensity(bandwidth=np.float64(0.016768329368110083), rtol=0.0001)}
```



Pharmacophore Pair Distributions

Once we have modeled the probability density of a pair of pharmacophores having a particular distance, we can score molecules that we have never seen. For an unseen molecule, we can compute the pairwise distances of its pharmacophores and observe how well its

pairwise pharmacophore distances align with the fitted density estimates. If the unseen molecule's distance measurements closely correspond to the distributions that we observed from the product of our experiment fragment screens, then it is assigned a higher score. If it doesn't fit well, then its score is lower.

As an example, consider the donor-acceptor distance distribution in the previous figure. A donor-acceptor pair that is 6 angstroms apart is much more likely than a donor-acceptor pair that is 10 angstroms apart.

To score a new molecule, we loop over its pharmacophore pairs and calculate the probability of the distance between each pharmacophore pair (with respect to the fitted density estimate). We compute the mean of the log probabilities, which outputs the final score. The next code block shows how to do this to score the new molecules we sampled via Taylor-Butina clustering.

```
 1 def calculate_pairwise_distances_between_pharmacophores_for_a_single_ligand(df_of_pcores_for_single_ligand,
 2     """
 3     Calculate the pairwise distance between pharmacophores pcore_a and pcore_b in the same ligand.
 4
 5     Parameters:
 6     – df: DataFrame containing pharmacophoric data for a single ligand
 7     – pcore_a: Name of the first pharmacophore
 8     – pcore_b: Name of the second pharmacophore
 9
10     Returns:
11     – A NumPy array of pairwise distances between the specified pharmacophores
12     """
13     df_pcore_a = df_of_pcores_for_single_ligand.query('pcore == @pcore_a')
14     coords_a = df_pcore_a[['coord_x', 'coord_y', 'coord_z']].to_numpy()
15     df_pcore_b = df_of_pcores_for_single_ligand.query('pcore == @pcore_b')
16     coords_b = df_pcore_b[['coord_x', 'coord_y', 'coord_z']].to_numpy()
17
18     # If coordinates are missing then we can't calculate distances
19     if len(coords_b) > 0:
20         delta_coords = coords_a[:, np.newaxis] – coords_b
21         distances = np.linalg.norm(delta_coords, axis=2)
22         return distances.flatten()
23
24 def score_mol(kde_dict, pair_distribution, pcore_pairs):
25     """
26     Score a molecule based on its pharmacophoric pair distribution.
27
28     Parameters:
29     – kde_dict: Dictionary containing Kernel Density Estimations (KDE) for each pharmacophore pair
30     – pair_distribution: Dictionary containing pairwise distances for the molecule
31     – pcore_pairs: List of pharmacophore pairs
32
33     Returns:
34     – The score for the molecule
35     """
36     score_df = pd.DataFrame(columns=pcore_pairs)
37     for pcore_combination in pcore_pairs:
38         kde = kde_dict[pcore_combination]
39         pcore_dist = pair_distribution[pcore_combination].reshape(-1, 1)
40         pcore_score = np.max(kde.score_samples(pcore_dist.reshape(-1, 1)))
41         score_df.at[0, pcore_combination] = pcore_score
42
43     return np.nanmean(score_df[pcore_pairs].to_numpy().astype(float))
```

```
1 smiles = 'Cc1ccccc1CNc1ccccc1NC(=O)[C@@H](O)c1cccnc1'
2 mol = Chem.MolFromSmiles(smiles)
3 mol = Chem.AddHs(mol)
4 AllChem.EmbedMolecule(mol)
```

```
0
```

```
1 mol_pcore_df = compute_pcore_data(mol)
2 pair_distribution_for_this_ligand = {}
3 for pcore_pair in interesting_pcores:
4     core_a, core_b = pcore_pair.split('-')
5     pair_distribution_for_this_ligand[pcore_pair] = calculate_pairwise_distances_between_pharmacophores_for_
6         mol_pcore_df, core_a, core_b)
```

```
1 score_for_this_mol = score_mol(fresco_kdes, pair_distribution_for_this_ligand, interesting_pcores)
```

```
1 score_for_this_mol
```

```
np.float64(-1.9437227942680053)
```

And with that, we can score any compounds (e.g., our Taylor-Butina synthetic sample) with regard to their activity against Mpro in the context of our estimated probability distribution. We can then use these scores to prioritize compounds for downstream computational

and experimental studies.

## Exercise 5: Pairwise Distance

In our computation of the 2-body pharmacophore distribution, we calculated pairwise distances of pharmacophores. In a similar, simplified context, implement a pairwise distance calculator. Given a list of points in 2D space, implement a function that calculates the pairwise Euclidean distances between all points.

The function should:

1. Take a list of tuples, where each tuple represents a point (x, y).
2. Calculate the Euclidean distance between every pair of points.
3. Return a list of tuples, where each tuple contains: (point1, point2, distance)

Note: The distance between a point and itself should not be included in the result.

Example: Input: [(0, 0), (3, 4), (1, 1)] Output: [((0, 0), (3, 4), 5.0), ((0, 0), (1, 1), 1.4142135623730951), ((3, 4), (1, 1), 3.605551275463989)]

Hint: The Euclidean distance between two points (x1, y1) and (x2, y2) is: $sqrt((x2 - x1)^2 + (y2 - y1)^2)$

## ⌄ Student Solution to Exercise 5

Provide your solution to the above exercise in this cell and/or immediately following cells.

```python
1  import math
2  from typing import List, Tuple
3
4  def calculate_distance(p1: Tuple[float, float], p2: Tuple[float, float]) -> float:
5      """
6      Calculate Euclidean distance between two points.
7
8      Args:
9          p1: First point (x, y)
10         p2: Second point (x, y)
11     Returns:
12         float: Euclidean distance between the points
13     """
14     # TODO: Implement the function
15     delta_x = math.pow(p2[0] - p1[0], 2)
16     delta_y = math.pow(p2[1] - p1[1], 2)
17     distance = math.sqrt(delta_x + delta_y)
18
19     return (p1, p2, distance)
20
21 def calculate_pairwise_distances(points: List[Tuple[float, float]]) -> List[Tuple[Tuple[float, float], Tuple[
22     """
23     Calculate distances between all pairs of points.
24
25     Args:
26         points: List of (x, y) coordinate tuples
27     Returns:
28         List of tuples containing (point1, point2, distance)
29         Each pair is included only once and point1 comes before point2 in the input list
30     """
31     # Handle empty list or single point cases
32     if len(points) <= 1:
33         return []
34
35     result = []
36
37     # Calculate distances for all pairs
38     for i in range(len(points) - 1):
39         for j in range(i + 1, len(points)):
40             p1 = points[i]
41             p2 = points[j]
42             distance = calculate_distance(p1, p2)
43             result.append((p1, p2, distance))
44
45     return result
46
47 # Test cases
48 def run_tests():
49     test_cases = [
50         ([(0, 0), (3, 4), (1, 1)], 3),
51         ([(0, 0), (1, 1), (2, 2), (3, 3)], 6),
52         ([(0, 0)], 0),
53         ([], 0),
54         ([(1.5, 2.5), (3.5, 4.5), (5.5, 6.5), (7.5, 8.5)], 6)
55     ]
```

```
56
57     for i, (points, expected_count) in enumerate(test_cases):
58         print(f"\nTest case {i + 1}:")
59         print(f"Input: {points}")
60         result = calculate_pairwise_distances(points)
61         print(f"Output:")
62         for pair in result:
63             print(f"  {pair}")
64         print(f"Number of pairs: {len(result)}")
65         assert len(result) == expected_count, f"Expected {expected_count} pairs, but got {len(result)}"
66         print("Test case passed!")
67
68 run_tests()
```

```
Test case 1:
Input: [(0, 0), (3, 4), (1, 1)]
Output:
  ((0, 0), (3, 4), ((0, 0), (3, 4), 5.0))
  ((0, 0), (1, 1), ((0, 0), (1, 1), 1.4142135623730951))
  ((3, 4), (1, 1), ((3, 4), (1, 1), 3.605551275463989))
Number of pairs: 3
Test case passed!

Test case 2:
Input: [(0, 0), (1, 1), (2, 2), (3, 3)]
Output:
  ((0, 0), (1, 1), ((0, 0), (1, 1), 1.4142135623730951))
  ((0, 0), (2, 2), ((0, 0), (2, 2), 2.8284271247461903))
  ((0, 0), (3, 3), ((0, 0), (3, 3), 4.242640687119285))
  ((1, 1), (2, 2), ((1, 1), (2, 2), 1.4142135623730951))
  ((1, 1), (3, 3), ((1, 1), (3, 3), 2.8284271247461903))
  ((2, 2), (3, 3), ((2, 2), (3, 3), 1.4142135623730951))
Number of pairs: 6
Test case passed!

Test case 3:
Input: [(0, 0)]
Output:
Number of pairs: 0
Test case passed!

Test case 4:
Input: []
Output:
```