

Conservation analysis on PDB fragments binding mode

The following scripts underline the main analysis on the evaluation of binding mode conservation of PDB fragments. In a first step, the subset files are generated considering the initial filtering thresholds and the fragments subsets at the three possible conservation levels. Then a parallel pairwise comparison for the interactions fingerprints is carried out, to finalize with a global conservation score calculation.

Example code used in test case

Main script to generate the conservation subsets according to the initial filters and the different levels of conservation: fragments in all ligands binding all targets, fragments binding the same target, and fragments binding the same ligand. Independent files are generated after each step for further analysis and easy track back of data generation.

```
1 ### Parameters to filter the dataset ###
2 MWmin = 40 #Defines the Threshold for fragments size. (less than 40 too small)
3 MWmax = 300 #According to definition of fragments
4 ligs = 4 #Number of min different ligands the fragment is found. (according to scatter plot and mean)
5 uniireps = 9 #Number os min different targets the fragments is found. (according to scatter plot and mean)
6
7 def write_subset(dataset, subset_path, MWmin, MWmax, ligs, uniireps):
8     """Read fragments data and generates a subset according to the filter"""
9     #parameters: fragmentation report, output folder, min molsize, max molsize, Nr ligs, Nr targets
10    #Reads the fragmentation report
11    report = pd.read_csv(dataset, sep = '\t', index_col=None, dtype=None)
12    #Filters the data according to type of data and cutoff of MW
13    reportf1 = report[(report['TYPE']=='FR') & (report['MW']>=MWmin) & (report['MW']<=MWmax)]
14    #Group the data by fragments and gets count of unique ligands and unique proteins
15    groups = reportf1.groupby('FRAGINCHI')['LIGINCHI', 'UNIREP'].nunique(dropna=False).reset_index()
16    #Gets the subset of fragments according to cutoff of ligands and proteins
17    fragments = groups[(groups['LIGINCHI']>ligs) & (groups['UNIREP']>uniireps)]
18    #Subsets the original report by fraginchis in fragments subset
19    fil_subset = reportf1[reportf1['FRAGINCHI'].isin(fragments['FRAGINCHI'])]
20    #writes the subset into a file
21    fil_subset.to_csv(subset_path, sep='\t', index=False)
22
23
24 def generate_frag_files(subset_path, frags_subsets_path):
25     """Reads the filtered subset and for each fragment gets a new subset"""
26     #parameters: the path to the filtered subset and the path to folder that will contain the independent fragments
27     subsets
28     fil_subset = pd.read_csv(subset_path, sep = '\t', index_col=None, dtype = str)
29     #Gets a list of unique fragments among the filtered subset
30     fragments = fil_subset["FRAGINCHI"].unique().tolist()
31     for frag in fragments:
32         #for a given fragment inchikey gets the full subset
33         frag_subset = fil_subset[fil_subset["FRAGINCHI"] == frag]
34         #writes the fragments subset into an independent file
35         frag_subset.to_csv(frags_subsets_path+frag+".csv", sep='\t', index=False)
36
37
38 def generate_same_tar_files(subset_path, frag-tar_subsets_path):
39     """Reads the filtered subset and for each fragment and target combinations generates a new subset"""
40     #parameters: the path to the filtered subset and the path to folder that will contain the independent fragments-
41     targets subsets
42     fil_subset = pd.read_csv(subset_path, sep = '\t', index_col=None, dtype = str)
43     setcode = []
44     for index, row in fil_subset.iterrows():
45         #for each fragment-target combination generates a code
46         frag = str(row["FRAGINCHI"])
47         tar = str(row["UNIREP"])
48         #appends the code to a list
49         setcode.append('~'.join((frag,tar)))
50     #Adds the combination code list as new column into the subset
51     fil_subset["SETCODE"] = setcode
52     #obtains the unique codes as a list
53     setcodes = fil_subset["SETCODE"].unique().tolist()
```

```

51     for setcode in setcodes:
52         #for each combination code, extract the subset of complexes
53         set_subset = fil_subset[fil_subset["SETCODE"] == setcode]
54         #write the combination code subset into an independent file
55         set_subset.to_csv(frag-tar_subsets_path+setcode+".csv", sep='\t', index=False)
56
57 def generate_same_lig_files(subset_path, frag-lig_subsets_path):
58     """Reads the filtered subset and for each fragment and ligand combinations generates a new subset"""
59     #parameters: the path to the filtered subset and the path to folder that will contain the independent fragments-
    ligands subsets
60     fil_subset = pd.read_csv(subset_path, sep = '\t', index_col=None, dtype = str)
61     setcode = []
62     for index, row in fil_subset.iterrows():
63         #for each fragment-ligand combination generates a code
64         frag = str(row["FRAGINCHI"])
65         lig = str(row["LIGINCHI"])
66         #appends the code to a list
67         setcode.append('~'.join((frag,lig)))
68     #Adds the combination code list as new column into the subset
69     fil_subset["SETCODE"] = setcode
70     #obtains the unique codes as a list
71     setcodes = fil_subset["SETCODE"].unique().tolist()
72     for setcode in setcodes:
73         #for each combination code, extract the subset of complexes
74         set_subset = fil_subset[fil_subset["SETCODE"] == setcode]
75         #write the combination code subset into an independent file
76         set_subset.to_csv(frag-lig_subsets_path+setcode+".csv", sep='\t', index=False)

```

The class `Combis` takes care of the pairwise comparisons. Collects the fingerprints by UIDs and writes down all possible pairs of them. In order to achieve a fair comparison of binding modes, the number of UIDs by fragment subset was limited to 500. The multiple UIDs having the same ligand and target were reduced to one representative and only 500 samples were chosen.

```

1 class Combis():
2     """Defines a Combis object containing the pairwise comparisons of a fragment subset"""
3     def __init__(self, path_combis, path_FPs, frag_subset, frag):
4         #parameters: output path for combis, output path for the fingerprints, the input fragment subset, the fragment
        inchikey
5         self.path_combis = path_combis
6         self.path_FPs = path_FPs
7         self.subset = pd.read_csv(frag_subset, sep = '\t', index_col=None, dtype = str)
8         self.frag = frag
9         self.frag_uids = defaultdict(list)
10        self.uid_combis = defaultdict(list)
11
12    def write_uid_FP_csv(self):
13        """Read subset dataframe and writes csv file with the fingerprint data"""
14        subset_fps = self.subset[['UID', 'FPSIMPLE']]
15        subset_fps.to_csv(self.fol_FPs+self.frag+"_fps.csv", sep='\t', index=False) #Only FPSimple
        data
16
17    def get_rep_complex(self):
18        """For a given fragment, subset the dataset and group by uniprot-liginchi combinations.
19        Then selects only one uid as representative complex of the group and return the full list of uids"""
20        groups = self.subset.groupby(['UNIREP', 'LIGINCHI'])['UID'].apply(list).to_frame() #group by
        uniprot-liginchi and get list of uids
21        uids = [uid_list[0] for uid_list in groups['UID']] #Select the first complex in group and add to
        list
22        if len(uids) > 500:
23            uids = uids[:500]
24        return uids
25
26    def get_uids_for_frag(self):
27        """Creates a dictionary with FRAGINCHI as key and as value a list of all UIDs with such fragment.
28        If the fragment has > 500 uids, the set of UIDs will be reduced"""
29        self.frag_uids = self.subset["UID"].tolist()
30        if len(self.frag_uids) > 500:

```

```

31         uids = self.get_rep_complex()
32         self.frag_uids = uids
33
34     def write_uid_combis(self):
35         """Writes a file containing the list of all uid combinations for each fragment and another file containing
36         each uid combination"""
37         uids = self.frag_uids
38         if len(uids) > 1:
39             rep = open(self.fol_combis+self.frag+".csv", "w")
40             rep.write('UID1\tUID2\tFRAGINCHI\n')
41             for combi in itertools.combinations(uids,2):
42                 rep.write('{0}\t{1}\t{2}\n'.format(combi[0], combi[1], self.frag))
43             rep.close()

```

The class FPSimCalculator defines the calculation of the (cosine or tanimoto) similarities between pairwise fingerprints. For each UID pair, it calculates the fingerprints similarity and writes it down to a file.

```

1 class FPSimCalculator():
2     def __init__(self, uid_pairs, fpsimple):
3         self.uid_pairs = pd.read_csv(uid_pairs, sep = '\t', index_col=None, dtype = str)
4         self.FPsimple_dic = self.get_data_dict(fpsimple, 'UID', 'FPSIMPLE')
5
6     def get_data_dict(self, datafile, key1, key2):
7         """Reads a file as dataframe and constructs a dictionary with column (key1) as dict key and column (key2) as
8         value"""
9         data_dic = {}
10        data = pd.read_csv(datafile, sep = '\t', index_col=None, dtype = str)
11        for index, row in data.iterrows():
12            data_dic[row[key1]] = row[key2]
13        return data_dic
14
15    def get_cosine_sim(self, va, vb):
16        """Takes two fingerprint vectors va and vb and calculates the cosine similarity."""
17        va = list(map(int, va))
18        vb = list(map(int, vb))
19        csim = 1.0-distance.cosine(va, vb)
20        if np.isnan(csim):
21            csim = -1.0
22        return csim
23
24    def get_tanimoto_sim(self, va, vb):
25        """Takes two fingerprint vectors va and vb and calculates the cosine tanimoto coefficient metric"""
26        #common = [a for a, b in zip(va, vb) if a == b]
27        #return float((len(common))/(len(va) + len(vb) - len(common)))
28        common = [a for a, b in zip(va, vb) if a == b and a=='1']
29        binsva = [a for a in va if a=='1']
30        binsvb = [b for b in vb if b=='1']
31        nom = len(common)
32        denom = len(binsva)+len(binsvb)-len(common)
33        if denom != 0:
34            return float(nom/denom)
35        else:
36            return float(1.0)
37
38    def get_fp_sim(self, fp_dic, id1, id2):
39        """Checks if the tuple of ids is different, then calculates the FP similarity. If not, then sim is
40        automatically 1"""
41        if id1 != id2:
42            id1_fp = fp_dic[id1]
43            id2_fp = fp_dic[id2]
44            fp_sim = self.get_tanimoto_sim(list(id1_fp), list(id2_fp))
45        else:
46            fp_sim = 1.0
47        return fp_sim
48
49    def write_fingerprints_sim(self, output):
50        """Iterates over the ID pairs, calls the function to calculates similarity and write a new report"""

```

```

49     report = open(output, 'w')
50     report.write('UID1\tUID2\tFPSIMPLE_SIM\tGROUP\n')
51     for index, row in self.uid_pairs.iterrows():
52         id1 = row['UID1']
53         id2 = row['UID2']
54         group = row['FRAGINCHI'] #In Desaphy analysis is GROUP
55         fpsimple_sim = self.get_fp_sim(self.FPsimple_dic, id1,id2)
56         report.write('{0}\t{1}\t{2:.2f}\t{3}\n'.format(id1,id2, fpsimple_sim, group))
57     report.close()

```

The class DataAnalyzer provides the final report with the global binding mode conservation (similarities) for each fragments subset or group.

```

1 class DataAnalyzer():
2     def __init__(self, sim_data_path, frag):
3         self.sim_data = pd.read_csv(sim_data_path, sep = '\t', index_col=None, dtype = str)
4         self.frag = frag
5         self.sim_fps = self.sim_data['FPSIMPLE_SIM'].tolist()
6
7     def write_fp_sim(self, output):
8         """Writes the MIN/MAX values of fingerprints into a report file"""
9         report = open(output, 'w')
10        report.write("FRAGMENT\tFPSMIN\tFPSMAX\tFPSAVG\n")
11        minfps = min(self.sim_fps)
12        maxfps = max(self.sim_fps)
13        fps = list(map(float, self.sim_fps))
14        avgfps = sum(fps)/len(self.sim_fps)
15        report.write('{0}\t{1}\t{2}\t{3:.2f}\n'.format(self.frag, minfps, maxfps, avgfps))
16        report.close()

```