

iCFN Manual

Mostafa Karimi, mostafa_karimi@tamu.edu

November 8, 2018

Contents

| | | |
|----------|--|-----------|
| 1 | iCFN Download link and its Prerequisites | 3 |
| 2 | Data format | 4 |
| 2.1 | Input data for each protein | 4 |
| 2.2 | Multistate protein design with ensemble of substates | 6 |
| 2.3 | Toy example | 6 |
| 3 | Parameters in the iCFN | 10 |
| 4 | Replication of TCR design | 11 |
| 5 | iCFN for new design | 12 |
| 5.1 | Running protein design based on OSPREY energy calculation | 13 |
| 6 | iCFN output | 15 |

1 iCFN Download link and its Prerequisites

Binary executable version of iCFN is available in <https://shen-lab.github.io/software/iCFN/>. iCFN has been developed and tested in:

- C programming
- GCC 4.4.7 compiler
- Centos 6.7 Linux

2 Data format

2.1 Input data for each protein

For any protein, iCFN needs four files:

- Rot_positions: There should be a text file with 4 columns which explains about the design of this protein with each correspond to a mutable or flexible residue in the design:
 - First column: chain of the residue
 - Second column: residue number
 - Third column: Wild type amino acid of residue with three letter codes given in the below table.
 - Fourth column: number of rotamers for the residue
- Constant.dat: There should be a binary file with the constant energy (energy between fixed residue and the mutable/flexible with the fixed ones). (It should be in the format of double)
- Amino.txt: There should be a text file with each row correspond to one rotamer with their amino acid code based on the following table. They should follow the same order with the Rot_positions file.

| Amino acid | Code |
|------------|------|
| ALA | 0 |
| ARG | 1 |
| ASN | 2 |
| AS1 | 3 |
| AS2 | 4 |
| ASP | 5 |
| CYS | 6 |
| GLN | 7 |
| GL1 | 8 |
| GL2 | 9 |
| GLU | 10 |
| GLY | 11 |
| HIS | 12 |
| HSE | 13 |
| HSD | 14 |
| HSP | 15 |
| ILE | 16 |
| LEU | 17 |
| MET | 18 |
| PHE | 19 |
| PRO | 20 |
| SER | 21 |
| THR | 22 |
| TRP | 23 |
| TYR | 24 |
| VAL | 25 |
| LYS | 26 |

- Energies.dat: A binary file contains both singleton and pairwise energy terms for the protein. (All the values should be in double format). All the singleton energy terms should be provided at first and then the all the pairwise energy terms should be provided. Singleton energy terms should be written in the energy.dat binary file in the same order with Rot_positions file. Then, just the half rectangle of pairwise energy terms should be provided in the same order.

2.2 Multistate protein design with ensemble of substates

Suppose for the positive state there are P proteins and for negative state there are Q proteins. At first, for each of the protein, the 4 mentioned files should be provided. Moreover, following files should be provided:

- Path of all amino.txt: Two files should be created separately for the positive and negative state. Each of which should contain the absolute path to the P (Q) Amino.txt files for positive (negative) state.
- Path of all rot_positions: Two files should be created separately for the positive and negative state. Each of which should contain the absolute path to the P (Q) rot_positions files for positive (negative) state.
- Path of all energy.dat: Two files should be created separately for the positive and negative state. Each of which should contain the absolute path to the P (Q) energy.dat files for positive (negative) state.
- Path of all constant.dat: Two files should be created separately for the positive and negative state. Each of which should contain the absolute path to the P (Q) constant.dat files for positive (negative) state.
- Rot_input: There should be a file for the general design problem which contains all of the mutable or flexible residues with wild-type and mutant amino acids. (All proteins should have same mutable/flexible residues with the same wild-type amino acid). It should contains at least 3 columns with each corresponds to a mutable or flexible residues:
 - First column: chain of the each residue
 - Second column: residue number
 - Third column: Wild type amino acid of residue with three letter codes given in the above table.
 - Fourth and so on: For any residue which is mutable, you should list all the amino acid it should be mutate to with a space separated.

2.3 Toy example

I will explain about the detail of providing data for iCFN with a very simple toy example. Suppose there are two state in which positive state has two sub-states and the negative state has just one sub-state. (Therefore there are 3 proteins) Suppose there are 3 mutable or flexible residues in this toy

example.

Rot_input format: The toy example, contains 3 residues with 2 mutable residues (26, 54) and one flexible (32). The first residue is ASP in wild type which can be mutated to ARG/ASN. Second residue is GLY which is just flexible. Third residue is VAL in wild type which can be mutate to GLU.

| | | | | |
|---|----|-----|-----|-----|
| A | 26 | ASP | ARG | ASN |
| A | 32 | GLY | | |
| A | 54 | VAL | GLU | |

Rot_positions format: It should be in format similar to the following. Suppose first sub-state in positive state: It shows that the first residue in total has 5 rotamers (considering rotamers of wild type and mutates). Second residue has just one rotamer. Third residue has in total 5 rotamers.

| | | | |
|---|----|-----|---|
| A | 26 | ASP | 5 |
| A | 32 | GLY | 1 |
| A | 54 | VAL | 5 |

Second sub-state in positive state:

| | | | |
|---|----|-----|---|
| A | 26 | ASP | 4 |
| A | 32 | GLY | 1 |
| A | 54 | VAL | 3 |

First sub-state in negative state:

| | | | |
|---|----|-----|---|
| A | 26 | ASP | 3 |
| A | 32 | GLY | 1 |
| A | 54 | VAL | 4 |

The number of rotamers for a given amino acid in different sub-state can vary as the toy example shows.

Amino files: I will give an example for the first sub-state in the positive state and the other two protein are very similar. Suppose for the first residue, there 2 rotamers for ASP (with code 5), 2 for ARG (with code 1) and 1 for ASN (with code 2). For the second residue, there are 1 rotamer for the GLY (with code 11). There are 2 rotamers for VAL (with code 25) and 3 for GLU (with code 10). Therefore, the amino.txt for this protein will be:

5
5
1
1
2
11
25
25
10
10
10

Energy file: Suppose for the above example (just one of the proteins), singleton energy terms are as follow in which r_j^i corresponds to the i^{th} rotamer in j^{th} residue:

$$\begin{aligned} E(r_1^1) &= -1 & E(r_1^2) &= -2 & E(r_1^3) &= 0 & E(r_1^4) &= -10 & E(r_1^5) &= -1.5 \\ E(r_2^1) &= -2.3 \\ E(r_3^1) &= -1 & E(r_3^2) &= -6 & E(r_3^3) &= -1.7 & E(r_3^4) &= -9 & E(r_3^5) &= -3 \end{aligned}$$

Pairwise energy terms for the first and second rotamer of the first residue with rotamer of all the other residue are given as following:

$$\begin{aligned} E(r_1^1, r_2^1) &= -0.6 & E(r_1^1, r_3^1) &= -1 & E(r_1^1, r_3^2) &= 0 & E(r_1^1, r_3^3) &= 10000 \\ E(r_1^1, r_3^4) &= -1.5 & E(r_1^1, r_3^5) &= 3 & E(r_1^2, r_2^1) &= 10 & E(r_1^2, r_3^1) &= -9 \\ E(r_1^2, r_3^2) &= 5 & E(r_1^2, r_3^3) &= 10000 & E(r_1^2, r_3^4) &= 10000 & E(r_1^2, r_3^5) &= 2.4 \end{aligned}$$

Therefore, the energy file should be in the following order but it should be converted to binary file (double format):

-1
 -2
 0
 -10
 -1.5
 -2.3
 -1
 -6
 -1.7
 -9
 -3
 -0.6
 -1
 0
 10000
 -1.5
 3
 10
 -9
 5
 10000
 10000
 2.4

At first all the singleton energy term should be provided with same order as amino.txt. First 5 values correspond to the singleton energy terms of the first residue, the next one value to the second residue and the next 5 values correspond to third residue. After all the singleton energy terms are provided, the pairwise energy terms between first rotamer in first residue with rotamer of the succeeding residues should be provided. (For example the next 6 values correspond to pairwise energy terms of the first rotamer of the first residue with one rotamer of the second residue and the 5 rotamer of the third residue). Then, all the pairwise energy terms between the second rotamer of the first residue with rotamer of the succeeding residues should be provided. We goes on until we have the pairwise energy terms between all rotamers of the first residue with all the rotamers of succeeding residues. Then, pairwise energy terms of first rotamer of second residue with all the rotamer of succeeding residues (third, fourth, \dots) should be provided.

3 Parameters in the iCFN

In this section, I will explain about the parameters in the model:

- ϵ : This parameter is related to the specificity cutoff. This parameter is for the desired ϵ -optima (from the global optimum) of the overall objective function.
- τ : This parameter is related to the positive-substate stability cutoff. This parameter is for a constraint on each positive substate function (stability here) so that they are no worse than the wild-type value plus τ . This constraint is only needed as the application demands.
- δ : In the paper, this parameter is used for Both within-substate and across-substate energy cutoff since we set them equally in the paper. This parameter is for δ -optima of each substate function.

To guarantee the ϵ -optima of the overall objective function, one needs to set δ and ϵ so that δ is at least half of ϵ . This could be relaxed for practical reasons when small δ makes a design problem tractable, especially considering that δ is the most "inner"-level cutoff.

4 Replication of TCR design

To easily replicate our results in the paper for the T Cell Receptor (TCR) design, we provided a shell script with all needed files in correct directories with the path given in the shell script. With this shell scripts, one can run our single point mutation designs for TCR.

To run it:

```
$ ./run_protein_design_TCR_Multistate.sh 26 2 2 5 10 > log_output &
```

in which the arguments correspond to:

- First argument: the position we want to mutate. In our design there are four positions 26, 28, 98 and 100.
- Second argument: Energy cutoff for side chain packing and pruning by type dependent DEE. In the above example it is 2 Kcal/Mol.
- Third argument: Energy cutoff for across substate type dependent DEE. In the above example it is 2 Kcal/Mol.
- Fourth argument: Stability energy cutoff for positive state from wild type. In the above example it is 5 Kcal/Mol.
- Fifth argument: Specificity energy cutoff for ensemble of sequences. In the above example it is 10 Kcal/Mol.

To check if one have run the code correctly, we provided the correct results in the folder `outputs/correct_output_MSD`.

5 iCFN for new design

To run iCFN for a new design, we have provided a shell script. At first, one need to provide the files for each of the proteins in two states. When the files are ready, We create folder "new_data", in which has two empty folders of "positive" (for the protein correspond to positive state) and "negative" (for the protein correspond to negative state). One should do the following:

1. `$ cd new_data`
2. Put the rot_input, with exactly "rot_input" name inside of the "new_data" folder.
3. `$ cd positive`
4. Put all the proteins correspond to positive state and their files in "positive" directory. Provide the following as well:
 - Path of all rot_positions files with the name "rot_pos_positive_files.txt"
 - Path of all binary energy files with the name "energies_positive_files.txt"
 - Path of all amino files with the name "amino_positive_files.txt"
 - Path of all constant energy files with the name "const_positive_files.txt"
5. `$ cd negative`
6. Similarly, put all the proteins correspond to negative state and their files in "negative" directory. Provide the following as well:
 - Path of all rot_positions files with the name "rot_pos_negative_files.txt"
 - Path of all binary energy files with the name "energies_negative_files.txt"
 - Path of all amino files with the name "amino_negative_files.txt"
 - Path of all constant energy files with the name "const_negative_files.txt"
7. `$ cd ../../src`
8. Run
`./run_protein_design_Multistate.sh`
with the following arguments:
 - First argument: Energy cutoff for side chain packing and pruning by type dependent DEE

- Second argument: Energy cutoff for across substate type dependent DEE.
- Third argument: Stability energy cutoff for positive state from wild type.
- Fourth argument: Specificity energy cutoff for ensemble of sequences.
- Fifth argument: Maximum number of conformation for each substate per amino acid assignment (Integer value)
- Sixth argument: Maximum number of discrepancy, sequences can have from wild type. (Integer value). If you don't want to constraint it, you can put arbitrary large number. For example at least equal the number of mutable residues.
- seventh argument: The output name

9. Then the outputs will be generated in the "outputs" directory.

5.1 Running protein design based on OSPREY energy calculation

If you have already calculated the energy based on OSPREY, we have provided a source code to create all the necessary files for a given pdb. This code, will create the float version of binary energies, therefore you need to use our **float** version binary code. To do that:

1. `$ cd OSPREY-conversion`

2. Run

`$ make_io.o`

with following arguments:

- First argument: OSPREY energy file for a given pdb
- Second argument: output name of Rot_position file
- Third argument: output name of energy file
- Fourth argument: output name of amino file
- Fifth argument: output name of fleximer file (An internal file might not need it later)
- Sixth argument: output name of constant energy term file

Then the files will be ready based on the names you gave in the arguments. There are two example of energy calculation from OSPREY in the "OSPREY-conversion/example" directory with the names 3K75_unbound (as negative state) and 3K75_bound (as positive state). Try to run our code for this example!!

6 iCFN output

iCFN algorithm will provide 5(4) outputs for multi-state protein design (single state protein design). These outputs will include the optimal sequence, ensemble of sequences and ensemble of rotameric conformations for per sequence, substate and state. Output's names are based on the output name has been chosen in the shell script for running the code. For example for the multi-state TCR example, the output's names are based on the mutable position, stability cutoff and specificity cutoff such as "26_output_stab.5_seq.10" (position 26 was the mutable position with the stability cutoff of 5Kcal/-Mol and specificity cutoff of 10Kcal/Mol). In the following, output files are explained in detail based on the multi-state TCR example:

- 26_output_stab.5_seq.10: This file will provide the evolution of best solution has been found by the algorithm until it reaches the optimal sequence. It will contain one or several blocks of the following in the file:

```
*****best score is -12.950394 pos -175.762329 neg -162.811935
pos 0 amino 26 rot pos 2662 rot neg 1887
pos 1 amino 11 rot pos 0 rot neg 0
pos 2 amino 19 rot pos 3 rot neg 11
pos 3 amino 25 rot pos 6 rot neg 5
pos 4 amino 7 rot pos 0 rot neg 132
pos 5 amino 26 rot pos 71 rot neg 0
pos 6 amino 25 rot pos 3 rot neg 2
pos 7 amino 19 rot pos 9 rot neg 7
pos 8 amino 10 rot pos 106 rot neg 131
pos 9 amino 21 rot pos 103 rot neg 3
pos 10 amino 11 rot pos 0 rot neg 0
pos 11 amino 1 rot pos 0 rot neg 473
pos 12 amino 21 rot pos 0 rot neg 48
pos 13 amino 24 rot pos 69 rot neg 201
pos 14 amino 26 rot pos 17 rot neg 1
pos 15 amino 11 rot pos 0 rot neg 0
pos 16 amino 24 rot pos 139 rot neg 223
pos 17 amino 10 rot pos 104 rot neg 99
pos 18 amino 23 rot pos 0 rot neg 0
```

Figure 1: Block of the best sequence found so far in the file

Each of this block, is the best sequence has been found so far and the last block is the optimal sequence. The first row explain the specificity

energy of the sequence which is the difference of the best conformation energy of positive and negative state. In the above example, the energy of best conformation for the positive state is -175.76 and the negative state is -162.81, therefore the specificity is -12.95. The rest of the block, explain the rotameric conformation of the best conformation for each state. In the above above block, the second row, correspond to the first residue ("pos 0"), amino acid with the code 26 ("amino 26"), rotamer 2662 for positive state ("rot pos 2662") and rotamer 1887 for negative state ("rot neg 1887"). The rest of the block can be interpreted similarly. At the end of this file, the running time of each component is explained.

- 26_output_stab.5.seq.10.seq_name: This file will include the ensemble of sequences within the desired range of the optimal sequence (we will always include the wild-type as well). Following explain the columns in the file:
 - First column: Sequence number
 - Second column: the specificity energy for the sequence
 - From third column to the end: Amino acid for the first till the end residue
- 26_output_stab.5.seq.10.seq_code: This file is exactly similar to 26_output_stab.5.seq.10.seq_name in which the amino acid names are replaced with the amino acid codes.
- 26_output_stab.5.seq.10.pos_ensemble: This file will provide the top K best rotameric conformations (If they satisfy the DEE, stability, specificity constraints) for each substate and sequence in the positive state. In following, I will explain the columns in the file:
 - First column: Sequence number which is exactly correspond to the sequence in 26_output_stab.5.seq.10.seq_name file.
 - Second column: The substate or backbone number for positive state.
 - From third column to two columns before the last column: The rotamer for that residue for that specific substate and sequence
 - One column before the last one: Conformation energy
 - Last column: The specificity energy for that conformation

- 26_output_stab.5.seq.10.neg_ensemble: This file is exactly similar to 26_output_stab.5.seq.10.pos_ensemble but for the negative state. If the problem is single state protein design, this file won't exist.