# TopologyNet: Topology based deep convolutional and multi-task neural networks for biomolecular property predictions — Handcrafted auxiliary features for the mutation example

Zixuan Cang and Guowei Wei

For the prediction of protein folding free energy change upon mutation, there are factors that are crucial to protein stability but can not be easily captured by persistent homology analysis, such as electrostatics. Also, there exists valuable sequential data base that can further enhance the performance of the predictor. Therefore, handcrafted features regarding the aforementioned aspects are generated and described in details in this document.

The atomic features are computed for each atom or each atom pair and are then summed over according to their element types and their locations. We consider the following atom types, $e_1 : \mathrm{C}, e_2 : \mathrm{N}, e_3 : \mathrm{O}, e_4 : \mathrm{S}, e_5 : \mathrm{H}, e_6 : \{\mathrm{C}, \mathrm{N}, \mathrm{O}, \mathrm{S}\}$, and $e_7 : \{\mathrm{C}, \mathrm{N}, \mathrm{O}, \mathrm{S}, \mathrm{H}\}$, and the following locations, $l_1$ : at the mutation site, $l_2$: within 10 Å of the mutation site but not at the mutation site, and $l_3$: with a distance greater than 10 Å to the mutation site while smaller than a cutoff of 40 Å. All features except the PSSM and the neighborhood amino acid composition are computed for both the mutant and the wild type. The difference of the features between the two structures are also taken into consideration.

**Surface area** Solvent excluded surface area for each atom is computed. A total of 21 area features are generated for either the mutant or the wild type by summing over combinations of 7 element types $e_1 - e_7$ and 3 locations $l_1 - l_3$. The difference between the mutant and the wild type gives rise to other 21 features. Therefore, there are 63 surface area features. The atomic surface area is computed with our in-house software package ESES.[5] ESES is available at `http://weilab.math.msu.edu/ESES/`.

**van der Waals interaction** The van der Waals interaction is modeled by the Lennard-Johns potential and is computed on each heavy atom by contrasting the atom against all other heavy atoms. A total of 63 features are generated by summing over the combinations of element types $e_1, e_2, e_3, e_4, e_6$ and locations $l_1 - l_3$.

**Partial charge** The atomic partial charges approximate the electron density. A total of 63 charge features are generated by summing over the combinations of element types $e_1 - e_7$ and locations $l_1 - l_3$. Another set of 63 charge features are generated by taking the absolute value of the partial charge on each atom before summations. The atomic partial charges are assigned using PDB2PQR software package.[1] An example usage of the software is `"pdb2pqr -ff=amber -ph-calc-method=propka -with-ph=7.0 input.pdb output.pqr"`. This will also generate the file that describes the pKa information.

**Coulomb interaction** The Coulomb interaction energy is computed on each atom against all other atoms. Sixty three features are generated by summing over combinations of element types $e_1, e_2, e_3, e_4, e_6$ and locations $l_1 - l_3$. Another set of 63 Coulomb interaction features are obtained by taking the absolute value on each atom before summations.

**Atomic electrostatic solvation free energy** Atomic electrostatic solvation free energy for each atom is modeled with the Poisson Boltzmann model. Sixty three features are generated by summing over the combinations of element types $e_1 - e_7$ and locations $l_1 - l_3$. The atomic electrostatic solvation free energy is modeled with Poisson Boltzmann equation solved using our in-house software MIBPB. The MIBPB software can be found at `http://weilab.math.msu.edu/MIBPB/`. An example usage is `"mibpb $pdbid h=0.5"`, which will carry out the Poisson Boltzmann computation for the input file `$pdbid.pqr` with grid size 0.5. Gird size of 0.5 Å is recommended while 0.8 Å is also acceptable.

The construction of the atomic features described above is further described in Algorithm 1.

**Neighborhood amino acid composition** Amino acid residues are grouped into hydrophobic, polar, positively charged, negatively charged, and special cases. The residues within 10 Å of the mutation site are examined. The count and percentage of nearby residues in each group are used as features describing the environment

of the mutation site providing 10 features. The accumulation, average, and variance of volume, surface area, weight, and hydropathy scores of the nearby residues are also calculated providing another 12 features.

**pKa shifts** The pKa values of ionizable groups are predicted. The maximum and summation of absolute value of pKa changes and summation of pKa changes between the wild type and the mutant as well as the largest pKa changes in both positive and negative directions are calculated resulting in 5 features. The pKa values and the differences between the wild type and the mutant at the mutation site, $C$ terminal, and $N$ terminal are also calculated resulting in another 9 features. Finally, the accumulated absolute and net pKa shifts for each of the seven ionizable groups make up an extra of 14 features. The pKa values of the ionizable groups are computed using PROPKA software package.[4]

**Secondary structure** The probability scores for the mutation site to be in a coil, helix, and strand as well as the torsion angle are predicted from the amino acid sequences for the wild type and the mutant. The secondary structure of the mutation site are also directly identified from the protein structure. The difference between the wild type and the mutant are also included resulting in 18 features. The secondary structure probability scores are predicted using SPIDER software package.[2] An example usage of SPIDER software is `"SPIDER2_local/misc/pred_pssm.py $PSSMFile"`

**Position-specific scoring matrix (PSSM)** The conservation score in the position-specific scoring matrix of the mutation site to be the amino acid in the wild type and the mutant as well as the difference of the two are used as features. Specifically, we collect the two PSSM scores for the mutation site to be in the wild residue type and mutant residue type and the difference between the two from the PSSMs constructed from wild type sequences. The PSSMs are computed using the psiblast command from BLAST+ software suite.[3] The specific version used is BLAST+/2.2.31 and an example usage is `"psiblast -query -db $PathToNRDatabase -num_iterations 5 -out $OutFileName -out_ascii_pssm $PSSMFileName"`. If the query sequence is too short and there is no hit found, the option `-evalue` can be set to a bigger number to loosen the criteria of including an entry as hit.

---

**Algorithm 1** The detailed steps of constructing atomic features (solvation energy, surface ares, partial charge, van der Waals, Coulomb).

---
 1: **function** DISTANCETOMUTATIONSITE(atom)
 2:     curdis $\leftarrow \infty$
 3:     **for** i $= 1 : $ m **do**                    ▷ Loop over all m atoms at the mutation site
 4:         dis $\leftarrow$ EUCLIDEAN(atom, atom$_i$)              ▷ The distance between two atoms
 5:         **if** dis $<$ curdis **then**
 6:             curdis $\leftarrow$ dis
 7:         **end if**
 8:     **end for**
 9:     **return** dis
10: **end function**

11: **function** ATOMLOCATION(atom)
12:     **if** atom is at mutation site **then**
13:         **return** 1
14:     **else if** DISTANCETOMUTATIONSITE(atom) $<$ 10 Å **then**
15:         **return** 2
16:     **else**
17:         **return** 3
18:     **end if**
19: **end function**

20: **function** ELEMENTGROUP(groupindex)
21:     **if** groupindex==1 **then**
22:         **return** {C}
23:     **else if** groupindex==2 **then**
24:         **return** {N}
25:     **else if** groupindex==3 **then**
26:         **return** {O}

```
27:      else if groupindex==4 then
28:          return {S}
29:      else if groupindex==5 then
30:          return {H}
31:      else if groupindex==6 then
32:          return {C, N, O, S}
33:      else if groupindex==7 then
34:          return {C, N, O, S, H}
35:      end if
36: end function


37: function CONSTRUCTATOMICFEATURE(protein)
38:      featureatomic ← Empty vector
39:      featuresolveng ← $Zeros_{3 \times 7}$
40:      for i = 1 : n do                                    ▷ Loop over the n atoms of the protein
41:          for ie = 1 : 7 do
42:              if ELEMENT($atom_i$) in ELEMENTGROUP(ie) then
43:                  featuresolveng[ATOMLOCATION(atom),ie] += SOLVENG($atom_i$)
44:              end if
45:          end for
46:      end for
47:      Append FLATTEN(featuresolveng) to featureatomic
48:      featurearea ← $Zeros_{3 \times 7}$
49:      for i = 1 : n do
50:          for ie = 1 : 7 do
51:              if ELEMENT($atom_i$) in ELEMENTGROUP(ie) then
52:                  featurearea[ATOMLOCATION(atom),ie] += SUFACEAREA($atom_i$)
53:              end if
54:          end for
55:      end for
56:      Append FLATTEN(featurearea) to featureatomic
57:      featurecharge ← $Zeros_{3 \times 7 \times 2}$
58:      for i = 1 : n do
59:          for ie = 1 : 7 do
60:              if ELEMENT($atom_i$) in ELEMENTGROUP(ie) then
61:                  featurecharge[ATOMLOCATION(atom),ie,1] += PARTIALCHARGE($atom_i$)
62:                  featurecharge[ATOMLOCATION(atom),ie,2] += Abs(PARTIALCHARGE($atom_i$))
63:              end if
64:          end for
65:      end for
66:      Append FLATTEN(featurecharge) to featureatomic
67:      featurevdw ← $Zeros_{3 \times 5}$
68:      for i = 1 : n do
69:          for j = 1 : n do
70:              for ie in {1, 2, 3, 4, 6} do
71:                  if ELEMENT($atom_i$) in ELEMENTGROUP(ie) and ELEMENT($atom_j$) in ELEMENTGROUP(6) then
72:                      featurevdw[ATOMLOCATION(atom),ie] += VANDERWAALS($atom_i$, $atom_j$)
73:                  end if
74:              end for
75:          end for
76:      end for
77:      Append FLATTEN(featurevdw) to featureatomic
78:      featureclb ← $Zeros_{3 \times 5 \times 2}$
79:      for i = 1 : n do
```

```
80:        for j = 1 : n do
81:            for ie in {1, 2, 3, 4, 6} do
82:                if ELEMENT(atom_i) in ELEMENTGROUP(ie) and ELEMENT(atom_j) in ELEMENTGROUP(7) then
83:                    featurevdw[ATOMLOCATION(atom),ie,1] += COULOMB(atom_i, atom_j)
84:                    featurevdw[ATOMLOCATION(atom),ie,2] += Abs(COULOMB(atom_i, atom_j))
85:                end if
86:            end for
87:        end for
88:    end for
89:    Append FLATTEN(featureclb) to featureatomic
90:    return featureatomic
91: end function
```

**References**

[1] T. J. Dolinsky, P. Czodrowski, H. Li, J. E. Nielsen, J. H. Jensen, G. Klebe, and N. A. Baker. PDB2PQR: expanding and upgrading automated preparation of biomolecular structures for molecular simulations. *Nucleic Acids Research*, 35(S2):W522–5, 2007.

[2] Rhys Heffernan, Kuldip Paliwal, James Lyons, Abdollah Dehzangi, Alok Sharma, Jihua Wang, Abdul Sattar, Yuedong Yang, and Yaoqi Zhou. Improving prediction of secondary structure, local backbone angles, and solvent accessible surface area of proteins by iterative deep learning. *Scientific reports*, 5, 2015.

[3] Mark Johnson, Irena Zaretskaya, Yan Raytselis, Yuri Merezhuk, Scott McGinnis, and Thomas L Madden. Ncbi blast: a better web interface. *Nucleic Acids Research*, 36(suppl 2):W5–W9, 2008.

[4] H. Li, A. D. Robertson, and J. H. Jensen. Very fast empirical prediction and rationalization of protein pka values. *Proteins*, 61(4):704–21, 2005.

[5] Beibei Liu, Bao Wang, Rundong Zhao, Yiying Tong, and Guo Wei Wei. ESES: software for Eulerian solvent excluded surface. *Journal of Computational Chemistry*, 38:446–466, 2017.