

## 20.320 — Problem Set # 7

November 5<sup>th</sup>, 2010

*Due on November 12<sup>th</sup>, 2010 at 11:59am. No extensions will be granted.*

General Instructions:

1. You are expected to state all your assumptions and provide step-by-step solutions to the numerical problems. Unless indicated otherwise, the computational problems may be solved using Python/MATLAB or hand-solved showing all calculations. Both the results of any calculations and the corresponding code must be printed and attached to the solutions. For ease of grading (and in order to receive partial credit), your code must be well organized and thoroughly commented, with meaningful variable names.
2. You will need to submit the solutions to each problem to a separate mail box, so please prepare your answers appropriately. Staple the pages for each question separately and make sure your name appears on each set of pages. (The problems will be sent to different graders, which should allow us to get the graded problem set back to you more quickly.)
3. Submit your completed problem set to the marked box mounted on the wall of the fourth floor hallway between buildings 8 and 16.
4. The problem sets are due at noon on Friday the week after they were issued. There will be no extensions of deadlines for any problem sets in 20.320. Late submissions will not be accepted.
5. Please review the information about acceptable forms of collaboration, which was provided on the first day of class and follow the guidelines carefully.

# 1 Designing helical peptides

In this problem, you will attempt to rationally redesign a peptide to adopt a desired structure based on your knowledge of the principles of secondary structure formation and aided by a computer algorithm. You have been provided with a Python implementation of the Chou-Fasman algorithm.

- a) Start with the peptide  $(\text{Ala}_5\text{Gly}_4\text{Ser})_2$ . What do you predict its secondary structure to be? In solution, do you think that a large or a small fraction of the peptide will actually adopt the predicted structure? Justify your answer.
- b) Find a variant that differs from the starting peptide by at most 3 amino acids and is not predicted to be helical.
- c) Find 3 variants that each differ from the original peptide by at most 5 amino acids and are strongly predicted to be helical through their length. *Make sure the 3 variants are all different and non-trivial; the graders will use common sense to judge whether a variant is a genuinely distinct peptide or a minor variation on another of your solutions.*
- d) Based on your knowledge of the C-F algorithm, explain why each of your variants in the above questions produced the changes in expected helicity.

## Python code for Problem 2

choufasman.py:

```
1
2
3 pepseq = ['ALA', 'ALA', 'ALA', 'ALA', 'ALA',
4           'GLY', 'GLY', 'GLY', 'GLY', 'SER',
5           'ALA', 'ALA', 'ALA', 'ALA', 'ALA',
6           'GLY', 'GLY', 'GLY', 'GLY', 'SER']
7
8 PA = {"ALA":1.42,\
9       "ARG":0.98,\
10      "ASN":1.01,\
11      "ASP":0.67,\
12      "CYS":0.70,\
13      "GLN":1.51,\
14      "GLU":1.11,\
15      "GLY":0.57,\
16      "HIS":1.00,\
17      "ILE":1.08,\
18      "LEU":1.21,\
19      "LYS":1.14,\
20      "MET":1.45,\
21      "PHE":1.13,\
22      "PRO":0.57,\
23      "SER":0.77,\
24      "THR":0.83,\
25      "TRP":1.08,\
26      "TYR":0.69,\
27      "VAL":1.06}
28
29 PA2 = {"ALA":"H",\
30        "ARG":"i",\
31        "ASN":"b",\
32        "ASP":"I",\
33        "CYS":"i",\
34        "GLN":"h",\
35        "GLU":"H",\
36        "GLY":"B",\
37        "HIS":"I",\
38        "ILE":"h",\
39        "LEU":"H",\
40        "LYS":"h",\
41        "MET":"H",\
42        "PHE":"h",\
43        "PRO":"B",\
44        "SER":"i",\
45        "THR":"i",\
46        "TRP":"h",\
47        "TYR":"b",\
48        "VAL":"h"}
49
50
51 def P_average(window):
52     total = 0.0
53     for residue in window:
54         total += PA[residue]
```

```

55     return (total/float(len(window)))
56
57 def findAlpha(seq,PA):
58     """
59     Uses Chou-Fasman criteria to suggest alpha helical regions
60     but does not take beta sheets into account
61     Inputs:
62     seq == (list) the amino acids sequence of the protein
63     PA == dictionary whose keys are amino acids and values are the
64     CF <Palph> parameters from the table in your problem set
65     PA2 == dictionary of CF a-helix Classification for each amino acid
66     Outputs:
67     AHindices == (list) contains the residue indices of seq that are
68     predicted to form helices
69     """
70     AHindices=[]
71     #Search for helix nucleation region
72     for i in range(len(seq)-5):
73         window = seq[i:i+6]
74         if not 'PRO' in window:
75             helix_propensity = 0.0
76             breakers = 0
77             for aa in window:
78                 if PA2[aa] == 'H' or PA2[aa] == 'h':
79                     helix_propensity += 1.0
80                 if PA2[aa] == 'I':
81                     helix_propensity += 0.5
82                 if PA2[aa] == 'b' or PA2[aa] == 'B':
83                     breakers += 1
84             if helix_propensity >= 4.0 and breakers < 2:
85                 begin = i
86                 end = i+5
87                 helix = (begin,end)
88                 #Extend nucleation region
89                 while (begin-4) >= 0:
90                     score = P.average(seq[begin-4:begin])
91                     if ('PRO' in seq[begin-4:begin]) or (score < 1.0): break
92                     else:
93                         begin -= 1 #Extend nucl. region in the N-term direction
94                         helix = (begin,end)
95                 while (end+4) < len(seq):
96                     score = P.average(seq[end+1:end+5])
97                     if ('PRO' in seq[end+1:end+5]) or (score < 1.0): break
98                     else:
99                         end += 1 #Extend nucl. region in the C-term direction
100                        helix = (begin,end)
101                #Store residues in AH indices
102                for n in range(begin,end+1):
103                    if not n in AHindices:
104                        AHindices.append(n)
105            return AHindices
106
107 myindices = findAlpha(pepseq,PA)
108
109 print myindices

```

## 2 Multiple sequence alignment

Receptor tyrosine kinases of the Epidermal Growth Factor Receptor (EGFR) family are essential to numerous physiological and pathological processes. In human, 12 EGFR family ligands have been identified and a significantly conserved section of the multiple sequence alignment (MSA) of some members of this family is shown below. We have also included the extracellular matrix protein Tenascin-C which contains EGF-like domains known to activate EGF receptors. Some gaps have been omitted to make to simplify the problem. In the MSA, the amino acids are represented by their one-letter amino acid code. Capital letters indicate a significant alignment while lowercase letters indicate no significant alignment was found.

### MSA Alignment

```
AREG_HUMAN/142-182      KKNPCNaefqNFCIH-GECKYIEH---LEAVTCKCQQEYFGERCG
BTC_HUMAN/65-105       HFSRCPkqykHYCIK-GRCRFVVA---EQTPSCVCDEGYIGARCE
EGF_HUMAN/972-1013    SDSECP1shdGYCLHDGVCMYIEA---LDKYACNCVVG YIGERCQ
EREG_HUMAN/64-104     SITKCSsdmngYCLH-GQCIYLVD---MSQNYCRCEVGYTGVRCE
HBEGF_HUMAN/104-144   KRDPCLrkykDFCIH-GECKYVKE---LRAPSCICHPGYHGERCH
NRG1_HUMAN/178-222    HLVKCAekekTFCVNGGECFMVKD1snPSRYLCKCQPGFTGARCT
NRG2_HUMAN/341-382    HARKCNetakSYCVNGGVCY YIEG---INQLSCKCPNGFFGQRCL
NRG3_HUMAN/286-329    HFKPCRdkd1AYCLNDGECFVIET1-tGSHKHCRCKEGYQGVRCD
NRG4_HUMAN/5-46       HEEPCGpshkSFCLNGGLCYVIPT---IPSPFCRCVENYTGARCE
TGFA_HUMAN/43-83      HFNDCPdshtQFCFH-GTCRFLVQ---EDKPACVCHSGYVGARCE
TENA_HUMAN/559-590    KEQRCP----SDCHGQGRCVDG-----QCICHEGFTGLDCG
```

- Deduce the PROSITE consensus pattern for the above alignment by **manual** pattern recognition of the MSA. If you are not familiar with PROSITE notation:  
[http://en.wikipedia.org/wiki/Sequence\\_motif](http://en.wikipedia.org/wiki/Sequence_motif).
- What amino acids were absolutely preserved throughout the evolution of this family? Give a rationale why each was preserved.
- What amino acids were somewhat preserved? (can be mutated to another amino acid with similar properties)
- Compute the log-odds matrix for the first five positions of this alignment. Assume that all amino acids are equally probable in the background and add a pseudocount of 0.1%.
- From a structural perspective, why would a particular amino acid be conserved while another would not at a specific position in a protein?
- Why are proline and glycine likely to disrupt helical secondary structures?

### 3 Bone morphogenic protein-2

Bone morphogenic protein-2 (BMP-2) is a 116 amino acid protein morphogen and part of a larger family of BMPs. Bone morphogenetic proteins regulate many developmental processes during embryogenesis as well as tissue homeostasis in the adult. Signaling of bone morphogenetic proteins is accomplished by binding to two types of serine/threonine kinase transmembrane receptors termed type I and type II. Because a large number of ligands signal through a limited number of receptors, ligand-receptor interaction in the BMP superfamily is highly promiscuous, with a ligand binding to various receptors and a receptor binding many different BMP ligands. You can obtain structural data for BMP2 from the Protein Databank (PDB) using accession [2QJB](#). The co-ordinates of the  $\alpha$ -carbons of the amino acids can be extracted from the PDB file using PyRosetta or any other method of your choice. For the remainder of this problem, consider only the 116 amino acids of chain A of the BMP2 dimer bound to its receptor. Be mindful that there is not fully characterized structural information for all residues in this protein chain. Indeed only residues 12-114 have been resolved. PyRosetta will read those residues with index 1-103.

- Determine the 103\*103 'distance matrix' for BMP2 chain A by computing the distances between the alpha carbons of all amino acid pairs. Represent this distance matrix as a 'binary contact matrix' (BCM). The contents of the BCM must be either 0 (for no contact) or 1 (for contact). For computing the BCM, you may assume that any two amino acids with  $C_\alpha$  atoms that are less than 6 angstroms apart make contact with each other. The use of MATLAB or Python to derive the BCM is permitted.
- Based on the BCM, determine the sum of sequence separations  $\Delta S_{ij}$  in residues (between contacting residues  $i$  and  $j$ ), for all the residues in BMP2.

$$S_{\text{TOT}, r} = \sum_{j=1}^n \Delta S_{rj}$$

- For which amino acid is  $S_{\text{TOT}}$  the greatest (only one amino acid)? What types of interactions does this amino acid make with its contacts? *Hint: find which residue it is most likely to contact with given the BMP results. The PDB file can help you.* How will these interactions influence the stability and the kinetics of folding?
- Calculate the total number of amino acid contacts ( $N$ ) within BMP (*i.e.* the number of pairwise interaction between residues of the protein) and use this to estimate the contact order (CO) for chain A of the insulin using:

$$\text{CO} = \sum_{i=1}^N \sum_{j=i+1}^N \frac{\Delta S_{ij}}{L \cdot N}$$

- Assuming that the logarithmic rate constant of protein folding  $\ln(k_{\text{eff}})$  is proportional to the contact order (CO) as shown by [Plaxco et al. J Mol Biol. 1998](#), estimate the rate constant of folding ( $k_{\text{eff}}$ ) using a constant of proportionality of 300, *i.e.* using the expression:

$$\ln(k_{\text{eff}}) = -300 \cdot \text{CO}$$

f) Do you suppose that the relationship between folding rate and contact order provided here holds for all proteins - why or why not? Does this relationship hold for BMP2 - why or why not? If not, will the actual rate constant of folding for BMP2 be lesser/greater than the calculated rate constant.

MIT OpenCourseWare  
<http://ocw.mit.edu>

20.320 Analysis of Biomolecular and Cellular Systems  
Fall 2012

For information about citing these materials or our Terms of Use, visit: <http://ocw.mit.edu/terms>.