

Graph Cuts and Protein Design

Adi Mukund

Jennifer Zou

1 The Protein Design Problem

The protein design problem can be described as follows. Given a set of backbone coordinates $\mathbf{c} = (\vec{c}_1, \vec{c}_2, \dots, \vec{c}_n)$ and rotamer library R of length r , identify the optimal rotamer assignment sequence $\vec{r} = (r_1, r_2, \dots, r_n)$, $r_i \in R$, $1 \leq i \leq n$ according to an energy function $E(\vec{r})$. This assignment sequence is known as the global minimum energy conformation (GMEC).

2 Graph Labeling

The graph-label formulation of the protein design problem represents the set of all residues as a sparse interaction graph, a modification of the residue interaction graph $G = (V, \mathcal{E})$ in which each vertex represents a protein residue and each edge represents an interaction between two such residues. The sparse graph is created by omitting a set of interactions \mathcal{E}^* to create the graph $G^* = (V, \mathcal{E} - \mathcal{E}^*)$.

The label graph G_l is generated by adding a set of vertices V_r representing all of the rotamers in the rotamer library R to the vertex set V^* of G_* and connecting each vertex in V^* to each vertex in V_r to generate the vertex set V_l and corresponding edge set \mathcal{E}_l . The set of all edges connecting vertices in V_r and V^* is \mathcal{E}_r . Edges from \mathcal{E}_l are then pruned until every vertex in V^* is connected to exactly 1 vertex in V_r (that is, until every member of V^* is represented in \mathcal{E}_r exactly once) in order to generate a labeling for the G where each vertex in V^* is connected to a vertex in V_r representing its optimal rotameric assignment.

Given a function $\phi(v, r_v)$ that returns the cost of assigning label r_v to vertex v and a function $\psi(v, r_v, w, r_w)$ that returns the pairwise interaction between residues v and w when assigned labels r_v and r_w respectively, the energy of the graph (and associated rotamer assignment sequence) can be computed as follows:

The energy of the graph can then be computed by finding the sum over the rotamers $v, w \in V^*$ where $w > v$ and their associated labels r_v and r_w as follows:

$$E_G = \sum_v \phi(v, r_v) + \sum_v \sum_{w > v} \psi(v, r_v, w, r_w) \quad (1)$$

Thus, finding the GMEC is equivalent to minimizing equation 1 over the set of all possible labelings of G_l .

3 Approximate Labeling based on Linear Programming

Note: This section is based on the paper “Approximate Labeling via Graph Cuts Based on Linear Programming” by Komodakis and Tziritas.

3.1 Labeling

The Graph Labeling (GL) problem can be stated as follows: classify a set \mathcal{V} of n objects by assigning to each object a label from a given set \mathcal{L} of labels, given a weighted graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, w)$. For each $p \in \mathcal{V}$ there is a label cost $\mathbf{c}_p(a) \geq 0$ for assigning the label $a = f_p$ to p , and for every edge pq there is a pairwise cost $\mathbf{c}_{pq}(a, b) = w_{pq}d_{pq}(ab)$ where $d_{pq}(ab)$ is the distance between (or cost of) label a on vertex p and b on vertex q . Thus, the cost of a labeling f is as follows:

$$\text{COST}(f) = \sum_{p \in \mathcal{V}} \mathbf{c}_p(f_p) + \sum_{(p,q) \in \mathcal{E}} w_{pq}d_{pq}(f_p, f_q) \quad (2)$$

3.2 Conversion to a Metric Labeling Schema

The Metric Labeling problem requires the distance function $d(a, b)$ to be a metric; otherwise it is identical to the GL problem described above. This imposes three constraints:

1. $d(a, b) = 0 \iff a = b$
2. $d(a, b) = d(b, a) \geq 0$
3. $d(a, b) \leq d(a, c) + d(c, b)$

The first constraint can be achieved by letting the label set \mathcal{L} be the set of all residues over all rotamers, that is, each $l \in \mathcal{L}$ is not just a particular rotamer but a particular rotamer *at* a particular residue. Thus, allowing $d(a, b) = 0$ is admissible, as such a pairwise assignment can be avoided by setting the internal potential for assigning a rotamer to the wrong residue equal to infinity. The second is easily seen to be true from this definition.

The third constraint can be ignored based on the Komodakis/Tziritas paper, as it is not possible to ensure that it is true. For example, if a and b involve significant steric clash, but a and c and b and c do not then $d(a, b)$ will be greater than $d(a, c) + d(c, b)$.