



## Note

# Liquid-theory analogy of direct-coupling analysis of multiple-sequence alignment and its implications for protein structure prediction

Akira R. Kinjo<sup>1</sup>

<sup>1</sup>Institute for Protein Research, Osaka University, Suita, Osaka 565-0871, Japan

Received October 23, 2015; accepted November 10, 2015

**The direct-coupling analysis is a powerful method for protein contact prediction, and enables us to extract “direct” correlations between distant sites that are latent in “indirect” correlations observed in a protein multiple-sequence alignment. I show that the direct correlation can be obtained by using a formulation analogous to the Ornstein-Zernike integral equation in liquid theory. This formulation intuitively illustrates how the indirect or apparent correlation arises from an infinite series of direct correlations, and provides interesting insights into protein structure prediction.**

**Key words:** Ornstein-Zernike equation, sequence profile, protein design, molecular evolution

Protein multiple-sequence alignments (MSA) are a useful means to extract various and valuable information about protein families [1]. It is well recognized that the frequency of amino acid residues at each alignment site is a useful measure of its functional importance. It has also been suggested that correlation between distant sites along the sequence is a rich source of information about the structure and function of the protein families [2]. In fact, recent years have seen a significant advance in our understanding of the site-site correlation observed in MSA. Of particular importance is the development of direct-coupling analysis (DCA) and related methods [3–5]. Although the basic idea has been already suggested in the last century [6], it is only by the recent

explosion of protein sequence data, in addition to theoretical development, that practical implementation of the idea was made possible. What DCA tells us is clear: The “apparent” correlation observed in a MSA is a result of “direct” correlations which are closely related to structural contacts. For example, if residues  $i$  and  $j$  are in physical contact (directly correlated), and so are residues  $j$  and  $k$ , then residues  $i$  and  $k$  may appear to be correlated even if they are not in contact.

There are many variants of DCA today. A major one is based on the principle of maximum entropy [3], others are based on the graphical Gaussian model [4] or phylogenetic analysis [5]. All of these methods are good predictors of physical contacts between residues in native protein structures. In this Note, I derive the direct correlation based on a formulation that is analogous to the integral equation theory of simple liquids [7]. This formulation has an advantage in that it intuitively shows how apparent correlations are realized by an infinite series of direct correlations. Based on the analogy with the liquid theory, it may be possible to elaborate the theory of direct correlations in MSA. More importantly, the intuitive picture that the present analysis provides helps us examine the mechanism of protein structure prediction from a new perspective, which may in turn lead to the development of new methods based on novel principles.

## Theory

A multiple-sequence alignment consisting of  $M$  ( $\gg 1$ ) amino acid sequences and  $N$  alignment sites may be regarded as an  $M \times N$  matrix of symbols. That is, each row represents an amino acid sequence including gap symbols and each column represents an alignment site. Let  $n_{k,i}(a) = 1$  if the residue type  $a$  appears at the site  $i$  of the sequence  $k$ , otherwise let

Corresponding author: Akira R. Kinjo, Institute for Protein Research, Osaka University, 3-2 Yamadaoka, Suita, Osaka 565-0871, Japan.  
e-mail: [akinjo@protein.osaka-u.ac.jp](mailto:akinjo@protein.osaka-u.ac.jp)

$n_{k,i}(a)=0$ . We first define the frequency  $n_i(a)$  of residue  $a$  at site  $i$  as

$$n_i(a) = \frac{1}{M} \sum_{k=1}^M n_{k,i}(a). \quad (1)$$

Next, the correlation (covariance) between residue  $a$  at site  $i$  and residue  $b$  at site  $j$  is defined as

$$C_{ij}(a,b) = \frac{1}{M} \sum_{k=1}^M [n_{k,i}(a) - n_i(a)][n_{k,j}(b) - n_j(b)]. \quad (2)$$

For simplicity, we assume that there are a sufficient number of sequences so that these statistics can be computed sufficiently accurately, and also ignore the effect of the phylogenetic bias in a family of sequences. Another caveat is required when there are completely conserved sites in the case of which the columns and rows corresponding to those conserved sites are zero. We assume this problem is properly taken care of, for example, by adding pseudo-counts. The correlations as a whole can be regarded as a  $21N \times 21N$  matrix by properly ordering residues and sites. Note that, since the equality  $\sum_{a=1}^{21} n_{k,i}(a) = 1$  holds for any sequence  $k$ , the matrix  $C$  is rank-deficient. Nevertheless, it can be made invertible by removing the rows and columns corresponding to the gap symbol, and hence the size of the matrix  $C$  is now  $20N \times 20N$ , which is assumed in the following.

Now we assume there exists a “direct correlation”  $D_{ij}(a,b)$  between residue  $a$  at site  $i$  and residue  $b$  at site  $j$ , and the correlation  $C$  is a result of an infinite series of the direct correlations:

$$\begin{aligned} C_{ij}(a,b) &= n_i(a)\delta_{i,j}\delta_{a,b} + n_i(a)D_{ij}(a,b)n_j(b) \\ &+ \sum_{k,c} n_i(a)D_{ik}(a,c)n_k(c)D_{kj}(c,b)n_j(b) \\ &+ \sum_{k_1,c_1,k_2,c_2} n_i(a)D_{ik_1}(a,c_1)n_{k_1}(c_1) \\ &\times D_{k_1k_2}(c_1,c_2)n_{k_2}(c_2)D_{k_2j}(c_2,b)n_j(b) \\ &+ \dots \end{aligned} \quad (3)$$

By defining the diagonal matrix  $\rho_{ij}(a,b) = n_i(a)\delta_{i,j}\delta_{a,b}$ , this equation is expressed as

$$C = \rho + \rho D \rho + \rho D \rho D \rho + \rho D \rho D \rho D \rho + \dots \quad (4)$$

$$= \rho + \rho D C. \quad (5)$$

This matrix equation is analogous to the Ornstein-Zernike integral equation in the theory of simple liquids [7] and can be expressed as a diagram in Figure 1 (where the left-hand side represents  $H = C - \rho$ ). By solving this equation for  $D$ , we have

$$D = \rho^{-1} - C^{-1} \quad (6)$$

which is essentially equivalent to the result of the mean-field DCA derived by Morcos *et al.* [3] based on the Plefka expansion [8].

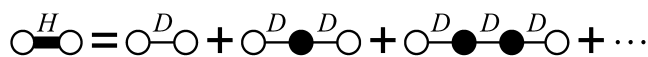


Figure 1 A diagrammatic representation of Eq. 4 with  $H = C - \rho$ .

## Discussion

While Morcos *et al.* [3] used direct correlations as *pair-wise interactions* between residues, direct correlations (in liquid theory) are generally different from interactions. In fact, the approach of Morcos *et al.* may be interpreted as the mean-spherical approximation [7] which is a particular closure condition for solving the Ornstein-Zernike equation. It may be interesting to investigate other choices of closure conditions such as those analogous to, for example, the Percus-Yevick (PY) or hypernetted-chain (HNC) approximations [7]. The HMSA closure [9] is another interesting possibility.

By rearranging Eq. (6), we have

$$\rho = (D + C^{-1})^{-1}. \quad (7)$$

This relation can be interpreted as a self-consistent condition (rather than a “definition”) for  $\rho$  when  $D$  is given, and shows how the position-specificity of residue frequencies depends on the entire context of a protein family and its structure. It is now widely accepted that sequence-based profile methods [10,11] are the best method for template-based structure prediction. Noting that the direct correlations well correspond to native contacts, Eq. (7) tells us that an infinite series of tertiary interactions are effectively convoluted into a sequence profile through the alignment of many evolutionarily related sequences. On the contrary, purely structure-based profile or threading methods [12], intuitively speaking, take into account only the first one or two terms in Eq. (4) where  $\rho$  in this case is *position-independent*. This may be a reason for the insufficient position-specificity, and hence the limited success, of purely structure-based profile methods.

The present analysis also has an implication for template-free or *de novo* structure prediction. All template-free methods are based on some empirical energy or scoring functions (whether physicochemical or statistical) and suffer from the problem of a rugged energy landscape that leads to many suboptimal non-native structures. In the mean time, studies on protein folding have shown that the energy landscape of natural proteins is minimally frustrated and funnel-like. This property can be readily modeled by the Go-like potentials in which only the native contacts are stabilizing [13,14]. It is conjectured that natural proteins have been naturally selected to satisfy such property in the course of molecular evolution [13]. This observation suggests a way to improve structure prediction by improving protein sequence design. That is, an empirical energy function that can reproduce the sequence profiles of (natural) protein families in the (re)designing process (i.e., generating sequences compatible with a given native structure) [15,16] may be expected to realize the “correct” direct correlation and development of such an energy

function may help improve structure prediction.

Physicochemically, it is the sequence that determines the structure. Evolutionarily, however, it is the structure that molds the pattern of a family of sequences. The DCA sheds new light especially on the latter aspect of proteins by explicitly providing the relation between the observed correlation  $C$  (i.e., the pattern of sequences) and the direct correlation  $D$  ( $\approx$  physical contacts). I hope the present analysis help further clarify the meaning of this intricate relationship between protein sequences and structures.

### Acknowledgments

I thank Mr. Iseo Nose whose lecture on Goethe's morphology motivated me to write this note.

### Conflict of Interest

None declared.

### Author Contributions

ARK did everything.

### References

- [1] Durbin, R., Eddy, S., Krogh, A. & Mitchison, G. *Biological Sequence Analysis: Probabilistic Models of Proteins and Nucleic Acids* (Cambridge Univ. Press, Cambridge, U.K., 1999).
- [2] Toh, H. *Bioinformatics for functional analyses of proteins* (Kodan-sha, Tokyo, Japan, 2004. In Japanese).
- [3] Morcos, F., Pagnani, A., Lunt, B., Bertolino, A., Marks, D. S., Sander, C., *et al.* Direct-coupling analysis of residue coevolution captures native contacts across many protein families. *Proc. Natl. Acad. Sci. USA* **108**, E1293–E1301 (2011).
- [4] Jones, D. T., Buchan, D. W. A., Cozzetto, D. & Pontil, M. PSICOV: precise structural contact prediction using sparse inverse covariance estimation on large multiple sequence alignments. *Bioinformatics* **28**, 184–190 (2012).
- [5] Miyazawa, S. Prediction of contact residue pairs based on co-substitution between sites in protein structures. *PLoS One* **8**, e54252 (2013).
- [6] Lapedes, A. S., Giraud, B., Liu, L. & Stormo, G. D. Correlated mutations in models of protein sequences: phylogenetic and structural effects. In *Statistics in Molecular Biology and Genetics*, 236–256 (Institute of Mathematical Statistics, Hayward, CA, USA, 1999).
- [7] Hansen, J.-P. & McDonald, I. R. *Theory of Simple Liquids* (Academic Press, London, 3rd edition, 2006).
- [8] Plefka, T. Convergence condition of the TAP equation for the infinite-ranged Ising spin glass model. *J. Phys. A. Math. Gen.* **15**, 1971–1978 (1982).
- [9] Zerah, G. & Hansen, J.-P. Self-consistent integral equations for fluid pair distribution functions: Another attempt. *J. Chem. Phys.* **84**, 2336–2343 (1986).
- [10] Altschul, S. F., Madden, T. L., Schäffer, A. A., Zhang, J., Zhang, Z., Miller, W., *et al.* Gapped blast and PSI-blast: a new generation of protein database search programs. *Nucleic Acids Res.* **25**, 3389–3402 (1997).
- [11] Krogh, A., Brown, M., Mian, I. S., Sjölander, K. & Haussler, D. Hidden Markov models in computational biology: Applications to protein modeling. *J. Met. Biol.* **235**, 1501–1531 (1994).
- [12] Eidhammer, I., Jonassen, I. & Taylor, W. R. *Protein bioinformatics* (Wiley & Sons, Chichester, England, 2004).
- [13] Gō, N. Theoretical studies of protein folding. *Annu. Rev. Biophys. Bioeng.* **12**, 183–210 (1983).
- [14] Onuchic, J. N. & Wolynes, P. G. Theory of protein folding. *Curr. Opin. Struct. Biol.* **14**, 70–75 (2004).
- [15] Kuhlman, B. & Baker, D. Native protein sequences are close to optimal for their structures. *Proc. Natl. Acad. Sci. USA* **97**, 10383–10388 (2000).
- [16] Ollikainen, N. & Kortemme, T. Computational protein design quantifies structural constraints on amino acid covariation. *PLoS Comput. Biol.* **9**, e1003313 (2013).