

---

# **Bayesian Model for Prediction of Protein Residue-Residue Contacts**

---

**Susann Vorberg**

---

15.10.2017



Dissertation zur Erlangung des Doktorgrades der Fakultt fr  
Chemie und Pharmazie der Ludwig-Maximilians-Universitt  
Mnchen

---

# **Bayesian Model for Prediction of Protein Residue-Residue Contacts**

---

vorgelegt von  
Susann Vorberg  
geboren in Leipzig, Germany

Mnchen, den 15.10.2017



## **Erklärung**

Diese Dissertation wurde im Sinne von 7 der Promotionsordnung vom 28. November 2011 von Dr. Johannes Soeding betreut.

## **Eidesstattliche Versicherung**

Diese Dissertation wurde eigenständig und ohne unerlaubte Hilfe erarbeitet.

.....  
Ort, Datum

.....  
Susann Vorberg

Dissertation eingereicht am: 15.10.2017

Erstgutachter: Dr. Johannes Soeding .....

Zweitgutachter: Prof. Dr. Julien Gagneur .....

Tag der mündlichen Prfung: 15.12.2017



# Summary

Awesome contact prediction project abstract



# Acknowledgements

I thank the world.



# Table of Contents

|   |            |
|---|------------|
| <b>Summary</b>  | <b>i</b>   |
| <b>Acknowledgements</b>   | <b>iii</b> |
| <b>Table of Contents</b>  | <b>vi</b>  |
| <b>1 Introduction to Contact Prediction</b>                       | <b>1</b>   |
| 1.1 Local Statistical Models . . . . .                            | 1          |
| 1.2 Global Statistical Models . . . . .                           | 3          |
| 1.3 Machine Learning Methods and Meta-Predictors . . . . .        | 4          |
| 1.4 Modelling Protein Families with Potts Model . . . . .         | 5          |
| 1.5 Applications . . . . .  | 12         |
| 1.6 Evaluating Contact Prediction Methods . . . . .               | 16         |
| 1.7 Challenges for Coevolutionary Inference . . . . .             | 19         |
| <b>2 Interpretation of Coupling Matrices</b>                      | <b>25</b>  |
| 2.1 Single Coupling Values Carry Evidence of Contacts . . . . .   | 25         |
| 2.2 Physico-Chemical Fingerprints in Coupling Matrices . . . . .  | 27         |
| 2.3 Coupling Profiles Vary with Distance . . . . .                | 30         |
| 2.4 Higher Order Dependencies Between Couplings . . . . .         | 33         |
| <b>3 Contact Prior</b>  | <b>37</b>  |
| 3.1 Random Forest Classifiers . . . . .                           | 37         |
| 3.2 Hyperparameter Optimization for Random Forest . . . . .       | 39         |
| 3.3 Evaluating Random Forest Model as Contact Predictor . . . . . | 42         |
| <b>4 Methods</b>  | <b>49</b>  |
| 4.1 Dataset . . . . .   | 49         |
| 4.2 Computing Pseudo-Likelihood Couplings . . . . .               | 50         |
| 4.3 Sequence Reweighting . . . . .                                | 52         |

|                        |  |            |
|------------------------|--|------------|
| 4.4                    | Computing Amino Acid Frequencies . . . . .   | 53         |
| 4.5                    | Regularization . . . . .   | 53         |
| 4.6                    | The Potts Model . . . . .  | 54         |
| 4.7                    | Analysis of Coupling Matrices . . . . .  | 60         |
| 4.8                    | Optimizing Contrastive Divergence with Stochastic Gradient Descent                                       | 60         |
| 4.9                    | Tuning Regularization Coefficients for Contrastive Divergence . . .                                      | 72         |
| 4.10                   | Tuning the Gibbs Sampling Scheme for Contrastive Divergence . . .  | 73         |
| 4.11                   | Features used to train Random Forest Model . . . . .   | 75         |
| 4.12                   | Training Random Forest Contact Prior . . . . .   | 80         |
| <b>A</b>               | <b>Abbreviations</b>   | <b>83</b>  |
| <b>B</b>               | <b>Amino Acid Alphabet</b>   | <b>85</b>  |
| <b>C</b>               | <b>Dataset Properties</b>  | <b>87</b>  |
| <b>D</b>               | <b>Standard Deviation of Couplings for Noncontacts</b>   | <b>91</b>  |
| <b>E</b>               | <b>Amino Acid Interaction Preferences Reflected in Coupling Matrices</b>                                 | <b>95</b>  |
| E.1                    | Pi-Cation interactions . . . . .   | 95         |
| E.2                    | Disulfide Bonds . . . . .  | 95         |
| E.3                    | Aromatic-Proline Interactions . . . . .  | 96         |
| E.4                    | Network-like structure of aromatic residues . . . . .  | 96         |
| E.5                    | Aromatic Sidechains at small $C_b$ - $C_\beta$ distances . . . . .                                       | 96         |
| <b>F</b>               | <b>Training of the Random Forest Contact Prior</b>   | <b>99</b>  |
| F.1                    | Training Random Forest Model with pseudo-likelihood Feature . . .  | 99         |
| F.2                    | Evaluating window size with 5-fold Cross-validation . . . . .  | 99         |
| F.3                    | Evaluating non-contact threshold with 5-fold Cross-validation . . .                                      | 99         |
| F.4                    | Evaluating ratio of non-contacts and contacts in the training set with 5-fold Cross-validation . . . . . | 99         |
| <b>List of Figures</b> |  | <b>112</b> |
| <b>List of Tables</b>  |  | <b>113</b> |
| <b>References</b>      |  | <b>115</b> |

# 1

## Introduction to Contact Prediction

Contact prediction refers to the prediction of physical contacts between amino acid side chains in the 3D protein structure, given the protein sequence as input.

Historically, contact prediction was motivated by the idea that compensatory mutations between spatially neighboring residues can be traced down from evolutionary records [1]. As proteins evolve, they are under selective pressure to maintain their function and correspondingly their structure. Consequently, residues and interactions between residues constraining the fold, protein complex formation, or other aspects of function are under selective pressure. Highly constrained residues and interactions will be strongly conserved [2]. Another possibility to maintain structural integrity is the mutual compensation of unbeneficial mutations. For example, the unfavourable mutation of a small amino acid residue into a bulky residue in the densely packed protein core might have been compensated in the course of evolution by a particularly small side chain in a neighboring position. Other physico-chemical quantities such as amino acid charge or hydrogen bonding capacity can also induce compensatory effects[3]. The [MSA](#) of a protein family comprises homolog sequences that have descended from a common ancestor and are aligned relative to each other. According to the hypothesis, compensatory mutations show up as correlations between the amino acid types of pairs of [MSA](#) columns and can be used to infer spatial proximity of residue pairs (see Figure [1.1](#)).

The following sections will give an overview over important methods and developments in the field of contact prediction.

### 1.1 Local Statistical Models

Early contact prediction methods used local pairwise statistics to infer contacts that regard pairs of amino acids in a sequence as statistically independent from another.

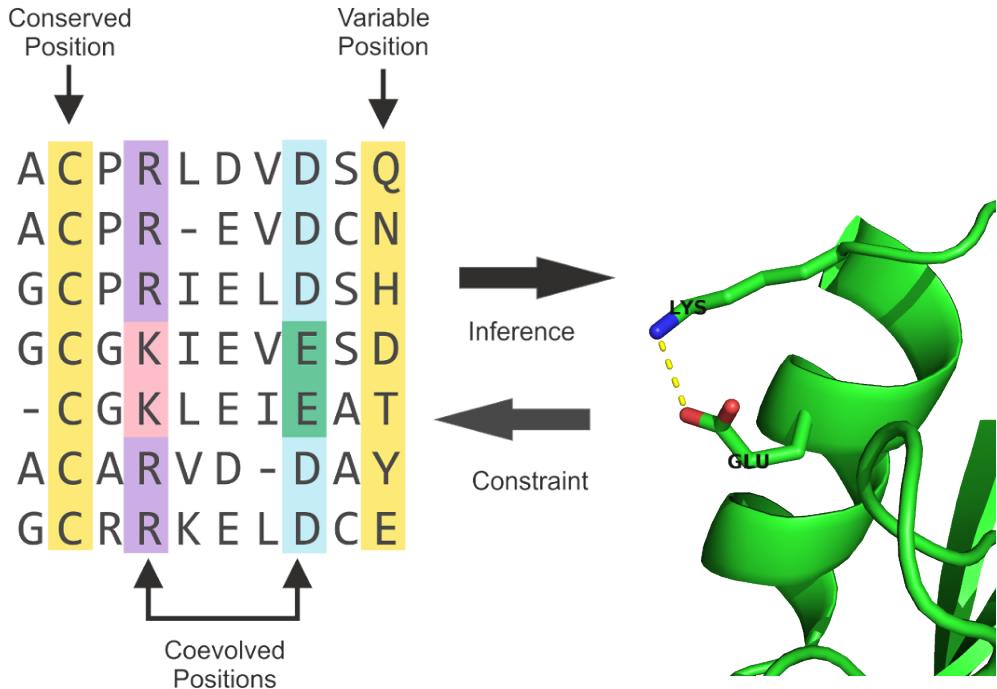


Figure 1.1: The evolutionary record of a protein family reveals evidence of compensatory mutations between spatially neighboring residues that are under selective pressure with respect to some physico-chemical constraints. Mining protein family sequence alignments for residue pairs with strong coevolutionary signals using statistical methods allows inference of spatial proximity for these residue pairs.

Several of these methods use correlation coefficient based measures, such as Pearson correlation between amino acid counts, properties associated with amino acids or mutational propensities at the sites of a MSA [1,3–6].

Many methods have been developed that are rooted in information theory and use MI measures to describe the dependencies between sites in the alignment [7–9]. Phylogenetic and entropic biases have been identified as strong sources of noise that confound the true coevolution signal [9–11]. Different variants of MI based approaches address these effects and improve on the signal-to-noise ratio [10,12,13]. The most prominent correction for background noises is APC that is still used by many modern methods and is discussed in section 1.4.4 [14]. Another popular method is OMES that essentially computes a chi-squared statistic to detect the differences between observed and expected pairwise amino acid frequencies for a pair of columns [15,16].

The traditional covariance approaches suffered from high false positive rates because of their inability to cope with transitive effects that arise from chains of correlations between multiple residue pairs [17–19]. The concept of transitve effects is illustrated in Figure 1.2. Considering three residues A, B and C, where A physically interacts with B and B with C. Strong statistical dependencies between pairs (A,B) and (B,C) can induce strong indirect signals for residues A and C, even though they are not physically interacting. These indirect correlations can become even larger than signals of other directly interacting pairs (D,E) and thus lead to false predictions [18].

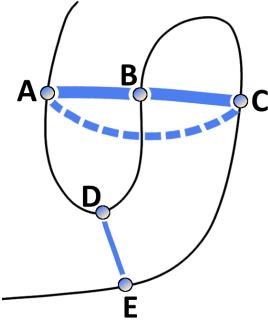


Figure 1.2: Effects of chained covariation obscure signals from true physical interactions. Consider residues A through E with physical interactions between the residue pairs A-B, B-C and D-E. The thickness of blue lines between residues reflects the strength of statistical dependencies between the corresponding alignment columns. Strong statistical dependencies between residue pairs (A,B) and (B,C) can induce a strong dependency between the spatially distant residues A and C. Covariation signals arising from transitive effects can become even stronger than other direct covariation signals and lead to false positive predictions.

Local statistical methods consider residue pairs independent of one another which is why they cannot distinguish between direct and indirect correlation signals. In contrast, global statistical models presented in the next section learn a joint probability distribution over all residues allowing to disentangle transitive effects [18,19]. Eventhough local statistical methods cannot compete with modern predictors, *OMES* and *MI* based scores often serve as a baseline in performance benchmarks for contact prediction [20,21].

## 1.2 Global Statistical Models

A huge leap forward was the development of sophisticated statistical models that make predictions for a single residue pair while considering all other pairs in the protein. These global models allow for the distinction between transitive and causal interactions which has been referred to in the literature as *DCA* [17,19].

In 1999 Lapedes et al. were the first to propose a global statistical approach for the prediction of residue-residue contacts in order to disentangle transitive effects [17]. They consider a Pott's model that can be derived under a maximum entropy assumption and use the model specific coupling parameters to infer interactions. At that time the wider implications of this advancement went unnoticed, but meanwhile the Pott's Model has become the most prominent statistical model for contact prediction. Section 1.4 deals extensively with the derivation and properties of the Pott's model, its application to contact prediction and its numerous realizations.

A global statistical model not motivated by the maximum entropy approach was proposed by Burger and Nijmegen in 2010 [18,22]. Their fast Bayesian network model incorporates additional prior information and phylogenetic correction via

[APC](#) but cannot compete with the pseudo-likelihood approaches presented in section [1.4.3.1](#).

## 1.3 Machine Learning Methods and Meta-Predictors

With the steady increase in protein sequence data, machine learning based methods have emerged that extract features from [MSAs](#) in order to learn associations between input features and residue-residue contacts. Sequence features typically include predicted solvent accessibility, predicted secondary structure, contact potentials, conservation scores, global protein features, pairwise coevolution statistics and averages of certain features over sequence windows. Numerous sequence-based methods have been developed using machine learning algorithms, such as support support vector machines (*SVMCon* [23], *SVM-SEQ* [24]), random forests (*ProC\_S3* [25], *TMhcp* [26], *PhyCMap* [27]), neural networks (*NETCSS* [28], *SAM* [29], [30], *SPINE-2D* [31], *NNCon* [32]) deep neural networks (*DNCon* [33], *CMAPro* [34]) and ensembles of genetic algorithm classifiers (*GaC* [35]).

Different contact predictors, especially when rooted in distinct principles like sequence-based and coevolution methods, provide orthogonal information on the likelihood that a pair of residues makes a contact [23,36]. The next logical step in method development therefore constitutes the combination of several base predictors and classical sequence-derived features in the form of meta-predictors.

The first published meta-predictor was *PconsC* in 2013, combining sequence features and predictions from the coevolution methods *PSICOV* and *plmDCA* [37]. In a follow-up version *PSICOV* has been replaced with *gaussianDCA* and the sequence-based method *PhyCMap* [38]. *EPC-MAP* was published in 2014 integrating *GREMLIN* as a coevolution feature with physicochemical information from predicted ab initio protein structures [39]. In 2015, *MetaPSICOV* was released combining predictions from *PSICOV*, *mfDCA* and *CCMpred* with other sequence derived feautures [40]. *RaptorX* uses *CCMpred* as coevolution feature and other standard contact prediction features within an ultra-deep neural network [41]. The newest developments *EPSILON-CP* and *NeBcon* both comprise the most comprehensive usage of contact prediction methods so far, combining five and eight state-of-the-art contact predictors, respectively [42,43].

Another conceptual advancement besides the combination of sources of information is based on the fact that contacts are not randomly or independently distributed. DiLena and colleagues found that over 98% of long-range contacts (sequence separation > 24 positions) are in close proximity of other contacts, compared to 30% for non-contacting pairs [34]. The distribution of contacts is governed by local structural elements, like interactions between helices or  $\beta$ -sheets, leading to characteristic patterns in the contact map that can be recognised [44]. Deep learning provides the means to model higher level abstractions of data and several methods apply multi-layered algorithms to refine predictions by learning patterns that reflect the local neighborhood of a contact [34,40,41,45].

Eventhough a benchmark comparing the recently developed meta-predictors is

yet to be made, it becomes clear from the recent CASP experiments, that meta-predictors outperform pure coevolution methods [46]. As coevolution scores comprise the most informative features among the set of input features, it is clear that meta-predictors will benefit from further improvements of pure coevolution methods [41,42].

## 1.4 Modelling Protein Families with Potts Model

Inferring contacts from a joint probability distribution over all residues in a protein sequence instead of using simple pairwise statistics has been proven to enable the distinction of direct statistical dependencies between residues from indirect dependencies mediated through other residues. The global statistical model that is commonly used to describe this joint probability distribution is the *Potts model*. It is a well-established model in statistical mechanics and can be derived from a maximum entropy assumption which is explained in the following.

The principle of maximum entropy, proposed by Jaynes in 1957 [47,48], states that the probability distribution which makes minimal assumptions and best represents observed data is the one that is in agreement with measured constraints (prior information) and has the largest entropy. In other words, from all distributions that are consistent with measured data, the distribution with maximal entropy should be chosen.

A protein family is represented by a MSA  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  of  $N$  protein sequences. Every protein sequence of the protein family represents a sample drawn from a target distribution  $p(\mathbf{x})$ , so that each protein sequence is associated with a probability. Each sequence  $\mathbf{x}_n = (\mathbf{x}_{n1}, \dots, \mathbf{x}_{nL})$  is of length  $L$  and every position constitutes a categorical variable  $x_i$  that can take values from an alphabet indexed by  $\{0, \dots, 20\}$ , where 0 stands for a gap and  $\{1, \dots, 20\}$  stand for the 20 types of amino acids. The measured constraints are given by the empirically observed single and pairwise amino acid frequencies that can be calculated as

$$f_i(a) = f(x_i=a) = \frac{1}{N} \sum_{n=1}^N I(x_{ni}=a) , \quad (1.1)$$

$$f_{ij}(a, b) = f(x_i=a, x_j=b) = \frac{1}{N} \sum_{n=1}^N I(x_{ni}=a, x_{nj}=b) . \quad (1.2)$$

According to the maximum entropy principle, the distribution  $p(\mathbf{x})$  should have maximal entropy and reproduce the empirically observed amino acid frequencies, so that

$$f(x_i = a) \equiv p(x_i = a) \\ = \sum_{\mathbf{x}'_1, \dots, \mathbf{x}'_L=1}^q p(\mathbf{x}') I(x'_i = a) \quad (1.3)$$

$$f(x_i = a, x_j = b) \equiv p(x_i = a, x_j = b) \\ = \sum_{\mathbf{x}'_1, \dots, \mathbf{x}'_L=1}^q p(\mathbf{x}') I(x'_i = a, x'_j = b). \quad (1.4)$$

Solving for the distribution  $p(\mathbf{x})$  that maximizes the Shannon entropy  $S = -\sum_{\mathbf{x}'} p(\mathbf{x}') \log p(\mathbf{x}')$  while satisfying the constraints given by the empirical amino acid frequencies in eq. (1.4) by introducing Lagrange multipliers  $\mathbf{w}_{ij}$  and  $v_i$ , results in the formulation of the *Potts model*,

$$p(\mathbf{x}|\mathbf{v}, \mathbf{w}) = \frac{1}{Z(\mathbf{v}, \mathbf{w})} \exp \left( \sum_{i=1}^L v_i(x_i) \sum_{1 \leq i < j \leq L} w_{ij}(x_i, x_j) \right). \quad (1.5)$$

The Lagrange multipliers  $\mathbf{w}_{ij}$  and  $v_i$  remain as model parameters to be fitted to data.  $Z$  is a normalization constant also known as *partition function* that ensures the total probability adds up to one by summing over all possible assignments to  $\mathbf{x}$ ,

$$Z(\mathbf{v}, \mathbf{w}) = \sum_{\mathbf{x}'_1, \dots, \mathbf{x}'_L=1}^q \exp \left( \sum_{i=1}^L v_i(x_i) \sum_{1 \leq i < j \leq L} w_{ij}(x_i, x_j) \right). \quad (1.6)$$

### 1.4.1 Model Properties

The Potts model is specified by singlet terms  $v_{ia}$  which describe the tendency for each amino acid  $a$  to appear at position  $i$ , and pair terms  $w_{ijab}$ , also called couplings, which describe the tendency of amino acid  $a$  at position  $i$  to co-occur with amino acid  $b$  at position  $j$ . In contrast to mere correlations, the couplings explain the causative dependence structure between positions by jointly modelling the distribution of all positions in a protein sequence and thus account for transitive effects. By doing so, a major source of noise in contact prediction methods is eliminated.

To get some intuition for the coupling coefficients, note that  $w_{ijab} = 1$  corresponds to a 2.7-fold higher probability for  $a$  and  $b$  to occur together than what is expected from the singlet frequencies if  $a$  and  $b$  were independent. Pairs of residues that are not in contact tend to have negligible couplings,  $\mathbf{w}_{ij} \approx 0$ , whereas pairs in contact tend to have vectors significantly different from 0. For contacting residues  $i$  and  $j$  in real world MSAs typical coupling strengths are on the order of  $\|\mathbf{w}_{ij}\| \approx 0.1$  (regularization dependent).

Maximum entropy models naturally give rise to exponential family distributions that express useful properties for statistical modelling, such as the convexity of the likelihood function which consequently has a unique, global minimum [49,50].

The Potts model is a discrete instance of what is referred to as a pairwise [Markov random field](#) in the statistics community. [MRFs](#) belong to the class of undirected graphical models, that represent the probability distribution in terms of a graph with nodes and edges characterizing the variables and the dependence structure between variables, respectively.

#### 1.4.1.1 Gauge Invariance

As every variable  $x_{ni}$  can take  $q = 21$  values, the model has  $L \times q + L(L - 1)/2 \times q^2$  parameters. But the parameters are not uniquely determined and multiple parametrizations yield identical probability distributions.

For example, adding a constant to all elements in  $v_i$  for any fixed position  $i$  or similarly adding a constant to  $v_{ia}$  for any fixed position  $i$  and amino acid  $a$  and subtracting the same constant from the  $qL$  coefficients  $w_{ijab}$  with  $b \in \{1, \dots, q\}$  and  $j \in \{1, \dots, L\}$  leaves the probabilities for all sequences under the model unchanged, since such a change will be compensated by a change of  $Z(\mathbf{v}, \mathbf{w})$  in eq. (1.6).

The overparametrization is referred to as *gauge invariance* in statistical physics literature and can be eliminated by removing parameters [19, 51]. An appropriate choice of which parameters to remove, referred to as *gauge choice*, reduces the number of parameters to  $L \times (q - 1) + L(L - 1)/2 \times (q - 1)^2$ . Popular gauge choices are the *zero-sum gauge* or *Ising-gauge* used by Weigt et al. [19] imposed by the restraints,

$$\sum_{a=1}^q v_{ia} = \sum_{a=1}^q w_{ijab} = \sum_{a=1}^q w_{ijba} = 0 \quad (1.7)$$

for all  $i, j, b$  or the *lattice-gas gauge* used by Morcos et al [51] and Marks et al [52] imposed by restraints

$$\mathbf{w}_{ij}(q, a) = \mathbf{w}_{ij}(a, q) = v_i(q) = 0 \quad (1.8)$$

for all  $i, j, a$  [53].

Alternatively, the indeterminacy can be fixed by including a regularization prior (see next section). The regularizer selects for a unique solution among all parametrizations of the optimal distribution and therefore eliminates the need to choose a gauge [54–56].

#### 1.4.2 Inferring Parameters for the Potts Model

Typically, parameter estimates are obtained by maximizing the log-likelihood function of the parameters over observed data. For the Potts model, the log-likelihood function is computed over sequences in the alignment  $\mathbf{X}$ :

$$\begin{aligned} \text{LL}(\mathbf{v}, \mathbf{w} | \mathbf{X}) &= \sum_{n=1}^N \log p(\mathbf{x}_n) \\ &= \sum_{n=1}^N \left[ \sum_{i=1}^L v_i(x_{ni}) + \sum_{1 \leq i < j \leq L} w_{ij}(x_{xn}, x_{nj}) - \log Z \right] \end{aligned} \quad (1.9)$$

The number of parameters in a Potts model is typically larger than the number of observations, i.e. the number of sequences in the [MSA](#). Considering a protein of length  $L = 100$ , there are approximately  $2 \times 10^6$  parameters in the model whereas the largest protein families comprise only around  $10^5$  sequences (see Figure 1.9). An underdetermined problem like this renders the use of regularizers necessary in order to prevent overfitting.

Typically, an L2-regularization is used that pushes the single and pairwise terms smoothly towards zero and is equivalent to the logarithm of a zero-centered Gaussian prior,

$$\begin{aligned} R(\mathbf{v}, \mathbf{w}) &= \log [\mathcal{N}(\mathbf{v} | \mathbf{0}, \lambda_v^{-1} I) \mathcal{N}(\mathbf{w} | \mathbf{0}, \lambda_w^{-1} I)] \\ &= -\frac{\lambda_v}{2} \|\mathbf{v}\|_2^2 - \frac{\lambda_w}{2} \|\mathbf{w}\|_2^2 + \text{const. ,} \end{aligned} \quad (1.10)$$

where the strength of regularization is tuned via the regularization coefficients  $\lambda_v$  and  $\lambda_w$  [57–59].

However, optimizing the log-likelihood requires computing the partition function  $Z$  given in eq. (1.6) that sums  $q^L$  terms. Computing this sum is intractable for realistic protein domains with more than 100 residues. Consequently, evaluating the likelihood function at each iteration of an optimization procedure is infeasible due to the exponential complexity of the partition function in protein length  $L$ .

Many approximate inference techniques have been developed to sidestep the infeasible computation of the partition function for the specific problem of predicting contacts that are briefly explained in the next section.

### 1.4.3 Solving the Inverse Potts Problem

In 1999 Lapedes et al. were the first to propose maximum entropy models for the prediction of residue-residue contacts in order to disentangle transitive effects [17]. In 2002 they applied their idea to 11 small proteins using an iterative Monte Carlo procedure to obtain estimates of the model parameters and achieved an increase in accuracy of 10-20% compared to the local statistical models [60]. As the calculations involved were very time-consuming and at that time required supercomputing resources, the wider implications were not noted yet.

Ten years later Weight et al proposed an iterative message-passing algorithm, here referred to as *mpDCA*, to approximate the partition function [19]. Eventhough their approach is computationally very expensive and in practice only applicable to

small proteins, they obtained remarkable results for the two-component signaling system in bacteria.

Balakrishnan et al were the first to apply pseudo-likelihood approximations to the full likelihood in 2011 [61]. The pseudo-likelihood optimizes a different objective and replaces the global partition function  $Z$  with local estimates. Balakrishnan and colleagues applied their method *GREMLIN* to learn sparse graphical models for 71 protein families. In a follow-up study in 2013, the authors proposed an improved version of *GREMLIN* that uses additional prior information [59].

Also in 2011, Morcos et al. introduced a naive mean-field inversion approximation to the partition function, named *mfDCA* [51]. This method allows for drastically shorter running times as the mean-field approach boils down to inverting the empirical covariance matrix calculated from observed amino acid frequencies for each residue pair  $i$  and  $j$  of the alignment. This study performed the first high-throughput analysis of intradomain contacts for 131 protein families and facilitated the prediction of protein structures from accurately predicted contacts in [52].

The initial work by Balakrishnan and colleagueas went almost unnoticed as it was not primarily targeted to the problem of contact prediction. Ekeberg and colleagueas independently developed the pseudo-likelihood method *plmDCA* in 2013 and showed its superior precision over *mfDCA* [55].

A related approach to mean-field approximation is sparse inverse covariance estimation, named *PSICOV*, developed by Jones et al. (2012) [21]. *PSICOV* uses an L1-regularization, known as graphical Lasso, to invert the correlation matrix and learn a sparse graphical model [62]. Both procedures, *mfDCA* and *PSICOV*, assume the model distribution to be a multivariate Gaussian. It has been shown by Banerjee et al. (2008)that this dual optimization solution also applies to binary data, as is the case in this application, where each position is encoded as a 20-dimensional binary vector [63].

Another related approach to *mfDCA* and *PSICOV* is *gaussianDCA*, proposed in 2014 by Baldassi et al. [64]. Similar to the other both approaches, they model the data as multivariate Gaussian but within a simple Bayesian formalism by using a suitable prior and estimating parameters over the posterior distribution.

So far, pseudo-likelihood has proven to be the most successful approximation of the likelihood with respect to contact prediction performance. Currently, there exist several implementations of pseudo-likelihood maximization that vary in slight details, perform similarly and thus are equally popular in the community, such as CCMpred [57], plmDCA[58] and GREMLIN [59].

#### 1.4.3.1 Maximum Likelihood Inference for Pseudo-Likelihood

The pseudo-likelihood is a rather old estimation principle that was suggested by Besag already in 1975 [65]. It represents a different objective function than the full likelihood and approximates the joint probability with the product over conditionals for each variable, i.e. the conditional probability of observing one variable given all the others:

$$\begin{aligned}
p(\mathbf{x}|\mathbf{v}, \mathbf{w}) &\approx \prod_{i=1}^L p(x_i|\mathbf{x}_{\setminus x_i}, \mathbf{v}, \mathbf{w}) \\
&= \prod_{i=1}^L \frac{1}{Z_i} \exp \left( v_i(x_i) \sum_{1 \leq i < j \leq L} w_{ij}(x_i, x_j) \right)
\end{aligned} \tag{1.11}$$

Here, the normalization term  $Z_i$  sums only over all assignments to one position  $i$  in sequence:

$$Z_i = \sum_{a=1}^q \exp \left( v_i(a) \sum_{1 \leq i < j \leq L} w_{ij}(a, x_j) \right) \tag{1.12}$$

Replacing the global partition function in the full likelihood with local estimates of lower complexity in the pseudo-likelihood objective resolves the computational intractability of the parameter optimization procedure. Hence, it is feasible to maximize the pseudo-log-likelihood function,

$$\begin{aligned}
\text{pLL}(\mathbf{v}, \mathbf{w} | \mathbf{X}) &= \sum_{n=1}^N \sum_{i=1}^L \log p(x_i | \mathbf{x}_{\setminus x_i}, \mathbf{v}, \mathbf{w}) \\
&= \sum_{n=1}^N \sum_{i=1}^L \left[ v_i(x_{ni}) + \sum_{j=i+1}^L w_{ij}(x_{ni}, x_{nj}) - \log Z_{ni} \right],
\end{aligned} \tag{1.13}$$

plus an additional regularization term in order to prevent overfitting and to fix the gauge to arrive at a **MAP** estimate of the parameters,

$$\hat{\mathbf{v}}, \hat{\mathbf{w}} = \underset{\mathbf{v}, \mathbf{w}}{\operatorname{argmax}} \text{pLL}(\mathbf{v}, \mathbf{w} | \mathbf{X}) + R(\mathbf{v}, \mathbf{w}). \tag{1.14}$$

Eventhough the pseudo-likelihood optimizes a different objective than the full-likelihood, it has been found to work well in practice for many problems, including contact prediction [50,54–56]. The pseudo-likelihood function retains the concavity of the likelihood and it has been proven to be a consistent estimator in the limit of infinite data for models of the exponential family [54,65,66]. That is, as the number of sequences in the alignment increases, pseudo-likelihood estimates converge towards the true full likelihood parameters.

#### 1.4.4 Computing Contact Maps

Model inference as described in the last section yields **MAP** estimates of the couplings  $\hat{\mathbf{w}}_{ij}$ . In order to obtain a scalar measure for the coupling strength between two residues  $i$  and  $j$ , all available methods presented in section 1.4.3 heuristically map the  $21 \times 21$  dimensional coupling matrix  $\mathbf{w}_{ij}$  to a single scalar quantity.

*mpDCA* [19] and *mfDCA* [51,52] employ a score called **DI**, that essentially computes the **MI** for two positions  $i$  and  $j$  using the couplings  $\mathbf{w}_{ij}$  instead of pairwise amino acid frequencies. Most pseudo-likelihood methods (*plmDCA* [55,58], *CCM-pred* [57], *GREMLIN* [59]) compute the *Frobenius norm* of the coupling matrix  $\mathbf{w}_{ij}$  to obtain a scalar contact score  $C_{ij}$ ,

$$C_{ij} = \|\mathbf{w}_{ij}\|_2 = \sqrt{\sum_{a,b=1}^q w_{ijab}^2}. \quad (1.15)$$

The Frobenius norm improves prediction performance over **DI** and further improvements can be obtained by computing the Frobenius norm only on the  $20 \times 20$  submatrix thus ignoring contributions from gaps [55,64,67]. *PSICOV* [21] uses an L1-norm on the  $20 \times 20$  submatrix instead of the Frobenius norm.

Furthermore it should be noted that the Frobenius norm is gauge dependent and is minimized by the *zero-sum gauge* [19]. Therefore, the coupling matrices should be transformed to *zero-sum gauge* before computing the Frobenius norm

$$\mathbf{w}'_{ij} = \mathbf{w}_{ij} - \mathbf{w}_{ij}(\cdot, b) - \mathbf{w}_{ij}(a, \cdot) + \mathbf{w}_{ij}(\cdot, \cdot), \quad (1.16)$$

where  $\cdot$  denotes average over the respective indices [55,57,58,64].

Another commonly applied heuristic known as **APC** has been introduced by Dunn et al. in order to reduce background noise arising from correlations between positions with high entropy or phylogenetic couplings [14]. **APC** is a correction term that is computed from the raw contact map as the product over average row and column contact scores  $\overline{C}_i$  divided by the average contact score over all pairs  $\overline{C}_{ij}$ . The corrected contact score  $C_{ij}^{APC}$  is obtained by subtracting the **APC** term from the raw contact score  $C_{ij}$ ,

$$C_{ij}^{APC} = C_{ij} - \frac{\overline{C}_i \overline{C}_j}{\overline{C}_{ij}}. \quad (1.17)$$

Visually, **APC** creates a *smoothing* effect on the contact maps that is illustrated in Figure 1.3 and it has been found to substantially boost contact prediction performance [14,59]. It was first adopted by *PSICOV* [21] but is now used by most methods to adjust raw contact scores.

It was long under debate why **APC** works so well and how it can be interpreted. Zhang et al. showed that **APC** essentially approximates the first principal component of the contact matrix and therefore removes the highest variability in the matrix that is assumed to arise from background biases [68]. Furthermore, they studied an advanced decomposition technique, called LRS matrix decomposition, that decomposes the contact matrix into a low-rank and a sparse component, representing background noise and true correlations, respectively.

Inferring contacts from the sparse component works astonishing well, improving precision further over **APC** independent of the underlying statistical model.

Dr Stefan Seemayer could show that the main component of background noise can be attributed to entropic effects and that a substantial part of **APC** amounts

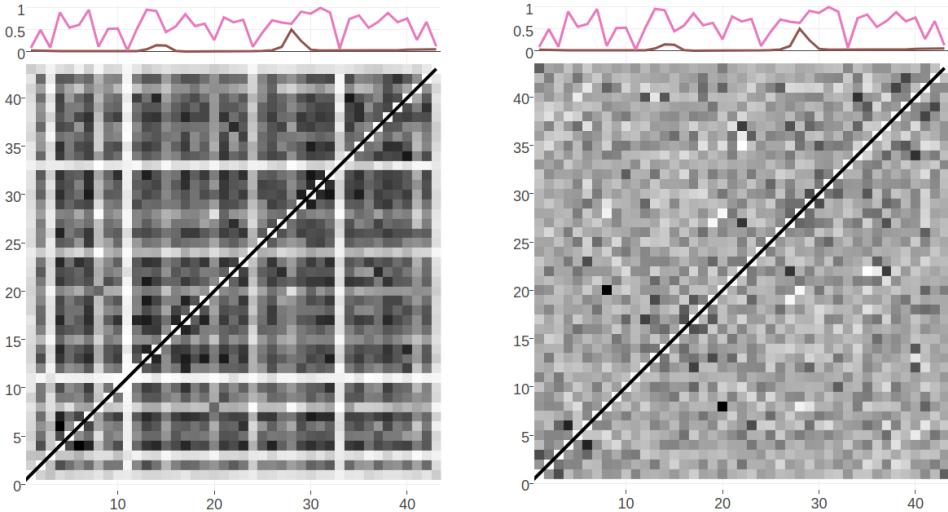


Figure 1.3: Contact maps computed from pseudo-likelihood couplings. Subplot on top of the contact maps illustrates the normalized Shannon entropy (pink line) and percentage of gaps for every position in the alignment (brown line). **a:** Contact map computed with Frobenius norm as in eq. (1.15). Overall coupling values are dominated by entropic effects, i.e. the amount of variation for a MSA position, leading to striped brightness patterns. For example, positions with high column entropy (e.g. positions 7, 12 or 31) have higher overall coupling values than positions with low column entropy (e.g. positions 11, 24 or 33). **b:** previous contact map but corrected for background noise with the APC as in eq. (1.17).

to correcting for these entropic biases (unpublished). In his doctoral thesis, he developed an entropy correction, computed as the geometric mean of per-column entropies, that correlates well with the APC correction term and yields similar precision for predicted contacts. The entropy correction has the advantage that it is computed from input statistics and therefore is independent of the statistical model used to infer the couplings. In contrast, APC and other denoising techniques such as LRS [68] discussed above, estimate a background model from the final contact matrix, thus depending on the statistical model used to infer the contact matrix.

## 1.5 Applications

The most popular and historically motivated application for contact prediction is contact-guided *de novo* structure prediction.

It has long been known that the native protein 3D structure can be reconstructed from an error-free contact map [69]. Also, protein fold reconstruction from sparse inter-residue proximity constraints obtained from experiments such as cross-linking/mass spectrometry, Foerster resonance energy transfer (FRET) or sparse nuclear Overhauser enhancement (NOE) distance data generated from NMR experiments has been demonstrated [70–75]. Predicted contacts, however, have long been regarded as being of little use for structure prediction because of

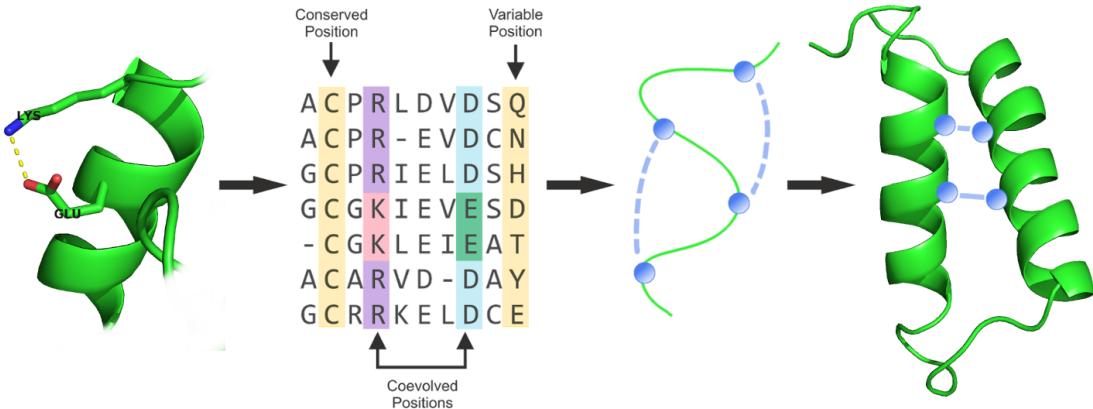


Figure 1.4: Generalized structure prediction pipeline integrating predicted contacts in form of distance constraints that guide conformational sampling.

their high false-positive rates [76,77]. Only with the emergence of global statistical models for contact prediction which drastically reduced false-positive rates there has been renewed interest in *de novo* structure prediction aided by predicted contacts. In 2011, Marks et al. showed that the top scoring contacts predicted with their mean-field approach *mfDCA* are sufficiently accurate to successfully deduce the native fold of the protein [52]. In the following years, methods to predict contacts have been improved and applied to model many more protein structures culminating in the high-throughput prediction of 614 protein structures out of which more than 100 represent novel folds by Ovchinnikov and colleagues in 2017 [78–86].

Many contact-guided protocols have been established since, that typically integrate predicted contacts in form of distance constraints into an energy function to guide the conformational sampling process: Unicon3D [87], RASREC [88], RBOAleph [89], GDFuzz3D [90], PconsFold [91], C2S\_Pipeline [92], FRAGFOLD + PSICOV [93], FILM3 [94], EVFold [52]. Figure 1.4 presents a generalized structure prediction pipeline using predicted contacts.

The optimal quality of inferred contacts and their effective utilization is still subject to discussion and further research. It has been demonstrated that only a small subset of native contacts is sufficient to produce accurate structural models [69,92,95–98]. Sathyapriya and colleagues developed a rational strategy to select important native contacts and successfully reconstructed the structure to near native resolution with only 8% of contacts [96]. Kim and colleagues formulated that only one correct contact for every 12 residues in the protein is sufficient to allow accurate topology level modeling given that the contacts are nonlocal and broadly distributed [95]. These studies emphasize that certain contacts are more important than others. Long-range contacts are rare and most informative for protein structure prediction because they define the overall fold and packing of tertiary structure whereas short-range contacts define local secondary structure [99]. It is a consistent finding that even though long-range contacts are of higher relevance than short-range contacts for structure reconstruction, their information alone is not sufficient [93,96,100]. Since a small number of correct residue-residue contacts is sufficient to improve protein structure prediction and many reconstruction pro-

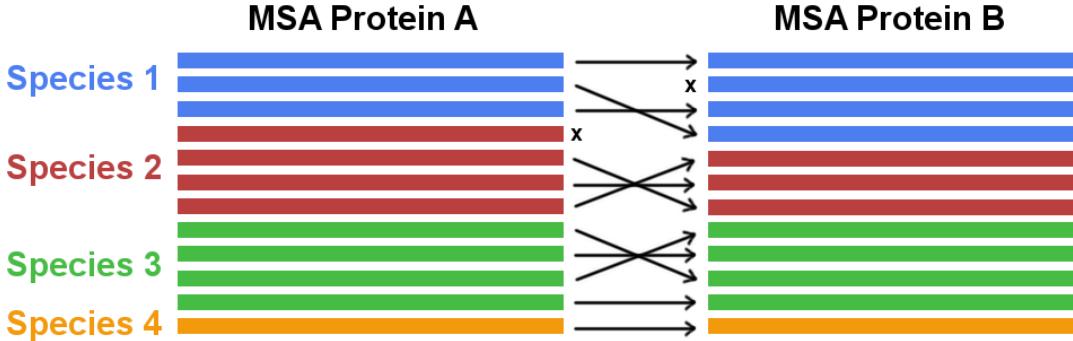


Figure 1.5: Concatenating two multiple sequence alignments. In case multiple paralogs exist for a gene in one species the correct interaction partner needs to be identified and matched (marked with arrows). Sequences that cannot be paired with a unique interaction partner need to be discarded (marked with x).

tocols can tolerate missing contact information much better than erroneous contact information, it has been stressed that methods development should focus on predicting a small number of high confident contacts [95,101]. Marks and colleagues observed that isolated false positives have a much stronger detrimental effect on structure prediction than false positives close to true contacts [52]. Zhang et al. found that their tool Touchstone II required an accuracy of long-range contact predictions of at least 22% to generate a positive effect to structure prediction [102]. Frequently, folding protocols employ a filtering step to eliminate unsatisfied or conflicting constraints possibly originating from false-positive contacts [103,104]. Generally it is assumed that higher precision of predicted long-range contacts results in improved structural models, albeit there is no strong correlation as model quality depends on many other factors such as the secondary structure composition of the protein, the domain size, the usage of additional sources of structural information, the type of distance constraint function and the particular structure reconstruction protocol [52,93,99,102,105].

Coevolution has not only been studied for residues pairs within a protein but also for residue pairs across protein–protein interfaces [79,80,85,106,107]. Eventhough the methodology of detecting coevolving amino acid pairs from the MSA is the same, a new challenge arises for the correct identification of orthologous interacting partners. Without the correct pairing of interacting partners for every species the detection of coevolutionary signals would be compromised. However, the generation of a MSA of paired sequences is complicated in the presence of multiple paralogs of a gene in a single genome. The problem of paralog matching is visualized in Figure 1.5. For prokaryotes, sequence paires are typically identified by exploiting the bacterial gene organisation in form of operons, i.e. co-localized genes will be co-expressed and are more likely to physically interact. Co-localisation of genes has also been applied to match genes from eukaryotes, assuming that Uniprot accession numbers can be used as a proxy for genomic distances [107]. New strategies have been developed based on the idea that an alignment with correctly matched paralogs will maximize the coevolution score [108,109].

A related objective is the study of the oligomerization status of proteins. The study

of homo-oligomers is simplified in the sense that the identical protein sequence of both interaction partners renders the concatenation of two [MSAs](#) unnecessary and allows to work with one [MSA](#). A different challenge lies in the correct distinction between the physical contacts of the monomeric structure and the interprotein contacts. With the availability of monomeric structural data the idea is to filter out those high scoring contacts that form contacts in the monomeric structure or are located in the protein core. The remaining high scoring false positive contacts at the surface of the protein are potential contacts at the interface that can be incorporated into a docking protocol to drive complex formation [110,111].

Contacts are also used to analyse potential alternative conformations of proteins [112–117]. Coevolutionary analysis detects all evolutionarily significant residue–residue correlations, regardless of whether the interaction is formed in a transient state of the protein or its stable form. Therefore, predicted contact maps might capture multiple states of a protein, since they are of functional importance and thus under evolutionary pressure. Sfriso and colleagues developed an automated pipeline that introduces filtered predicted contacts as ensemble restraints into a molecular dynamics simulations and is able to detect alternative relevant conformational states [113].

Quality assessment of structural models, involving model selection and ranking, is a crucial task in structural biology. Predicted residue-residue contacts can indicate the best protein structure among a set of properly folded and misfolded structures by counting the number of satisfied contacts [77,118]. Besides ranking of models, predicted contacts have been used as features for training machine learning methods that predict the global quality of a structural model [119,120].

As mentioned before, methods for protein fold reconstruction from experimental distance constraints have been successfully applied for many years. Several integrative approaches have been developed that combine complementary sources of sparse structural constraints, including predicted contacts, to accurately determine protein structure [121,122].

Sadowski used predicted contacts to parse domain boundaries based on the simple idea that contacts are more abundant within domains than between domains [123].

Eventhough the coevolutionary methods have been developed for proteins, they have been successfully applied to analyse nucleotide coevolution and to predict RNA tertiary structures with the help of predicted nucleotide-nucleotide contacts [124–126]. Much less RNA sequences are required compared to protein sequences in order to extract statistically significant signals because of the reduced number of model parameters when working with a four letter alphabet (compared to a 20 letter alphabet with proteins). On the downside, alignment errors resulting from the complicated determination of RNA multiple sequence alignments limits the accuracy of coevolution analysis [126]. Despite the diminished accuracy, predicted nucleotide contacts have been demonstrated to improve RNA structure prediction over conventional methods [125].

The stastistical models used for coevolution analysis provide information about which residue pairs are important in evolution for folding or functional constraints. They can be used to assign probabilities to sequences that reflect the overal compliance of a sequence with the protein family under study and thereby provide

quantitative predictions of mutational effects [127–129]. Computational screening of mutational effects can support and complement the costly and time-consuming directed evolution or mutational screening experiments [127]. With a similar idea in mind, the coevolution models have been applied to sequences of human immune repertoires [130,131]. Antibody affinity maturation can be viewed as a Darwinian process with the affinity to the target antigen being the main fitness criterion. Therefore, given the model representing the antibody sequence family, the probability for a sequence reflects the binding affinity to the target antigen. Quantifying the effect of mutations is also helpful for protein design. Coevolving positions might be of particular interest as hotspots for engineering protein stability or functional specificity because they determine positions relevant to protein structure and function [132].

Skwark and colleagues applied the popular coevolution statistical models to genomes and developed a statistical method called *genomeDCA* [133]. They are able to identify coevolving polymorphic locus pairs based on the idea that the corresponding proteins form protein-protein interactions that are under strong evolutionary pressure. In a case study on two large human pathogen populations they found that three quarters of coevolving loci are located in genes that determine beta-lactam (antibiotic) resistance.

Fox and colleagues turn the idea of DCA upside down. They developed a benchmark for testing the accuracy of large MSAs by evaluating the agreement between the predicted and the native contacts [134]. Based on the assumption that better alignments provide more accurate contact predictions, the alignment quality is inferred from the precision of predicted contacts.

## 1.6 Evaluating Contact Prediction Methods

Choosing an appropriate benchmark for contact prediction is determined by the further utilization of the predictions. Most prominently, predicted contacts are used to assist structure prediction as outlined in the last section 1.5. Therefore, one could assess the quality of structural models computed with the help of predicted contacts. However, predicting structural models adds not only another layer of computational complexity but also raises questions about implementation details of the folding protocol.

It has been found that in general a small number of accurate contacts is sufficient to constrain the overall protein fold as already discussed. From these considerations emerged various standard benchmarks that have been established by the CASP community over many years [46,135,136]. CASP, the well-respected and independent competition for the structural bioinformatic's community introduced the contact prediction category in 1996. Taking place every two years, the progress in the field is assessed in a blind competition and the community discusses the outcome in a subsequent meeting. According to the CASP regulations, a pair of residues is defined to be in physical contact when the distance between their  $C_\beta$  atoms ( $C_\alpha$  in case of glycine) is less than  $8\text{\AA}$  in the reference protein structure.

The overall performance of a contact predictor is evaluated by the mean precision

over a testset of proteins with known high quality 3D structures against the top scoring predictions from every protein. The number of top scoring predictions per protein is typically normalized with respect to protein length  $L$  and precision is defined as the number of true contacts among the top scoring predicted contacts,

$$\text{precision} = \frac{TP}{TP + FP}, \quad (1.18)$$

where  $TP$  is a true positive contact and  $FP$  is false positive contact. A popular variant of this benchmark plot shows the mean precision of a certain fraction of top ranked predictions (e.g.  $L/5$  top ranked predictions) against specific properties of the test proteins such as protein length or alignment depth [137]. Another informative metric is mean error defined as:

$$\text{mean error} = \frac{\text{error}}{TP + FP} \begin{cases} \text{error} = \Delta C_\beta - T & \text{if } \Delta C_\beta > T \\ \text{error} = 0, & \text{otherwise} \end{cases} \quad (1.19)$$

where  $\Delta C_\beta$  is the actual distance of a residue pair in the native structure, and  $T$  is the distance threshold defining a true contact. The mean error helps to asses how wrong false positive predictions are. During CASP11 further evaluation metrics have been introduced, such as Matthews correlation coefficient, area under the precision-recall curve or F1 measure but they are rarely used in studies [46].

Currently best methods perform in the range XXX. Sequence feature based methods: Their performance is less dependent on the number of available sequence homologs compared to coevolution methods and therefore they can outperform pure coevolution methods in low data ranges [27,138]. TODOOOPLOT

### 1.6.1 Sequence Separation

Local residue pairs separated by only some positions in sequence (e.g  $|i - j| < 6$ ) are usually filtered out for evaluating contact prediction methods. They are trivial to predict as they typically correspond to contacts within secondary structure elements and reflect the local geometrical constraints. Figure 1.6 shows the distribution of  $C_\beta$  distances for various minimal sequence separation thresholds. Without filtering local residue pairs (sequence separation 1), there are several additional peaks in the distribution around  $5.5\text{\AA}$ ,  $7.4\text{\AA}$  and  $10.6\text{\AA}$  that can be attributed to local interactions in e.g. helices (see Figure 1.7).

Commonly, sequence separation bins are applied to distinguish short ( $6 < |i - j| \leq 12$ ), medium ( $12 < |i - j| \leq 24$ ) and long range ( $|i - j| > 24$ ) contacts [46,136]. Especially long range contacts are of importance for structure prediction as they are the most informative and able to constrain the overal fold of a protein [135].

### 1.6.2 Interpretation of Evaluation Results

There are certain subtleties to be considered when interpreting contact prediction evaluation results.

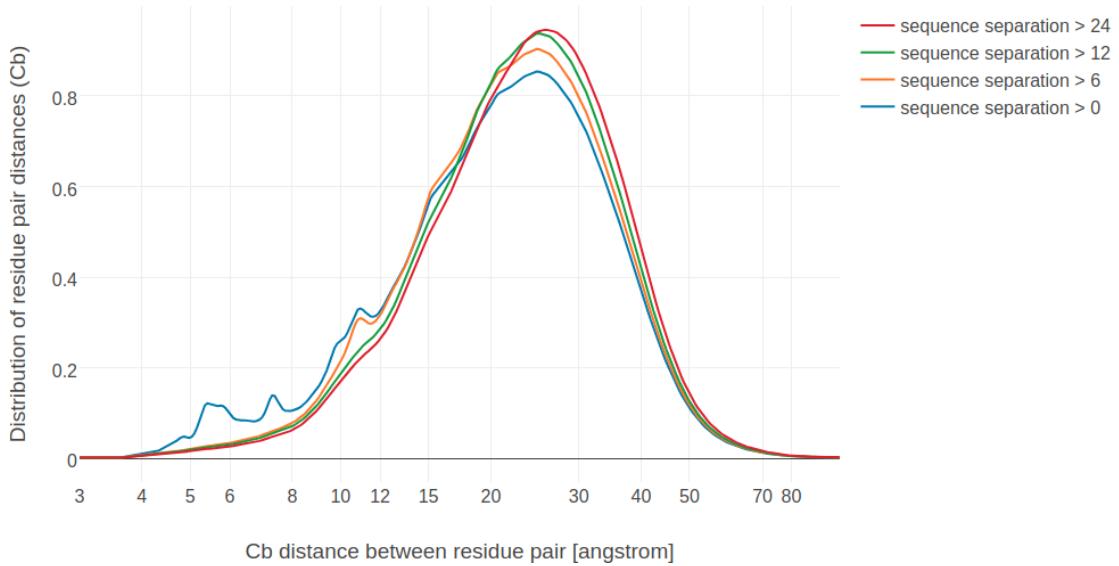


Figure 1.6: Distribution of residue pair  $C_\beta$  distances over 6741 proteins in the dataset (see Methods 4.1) at different minimal sequence separation thresholds.

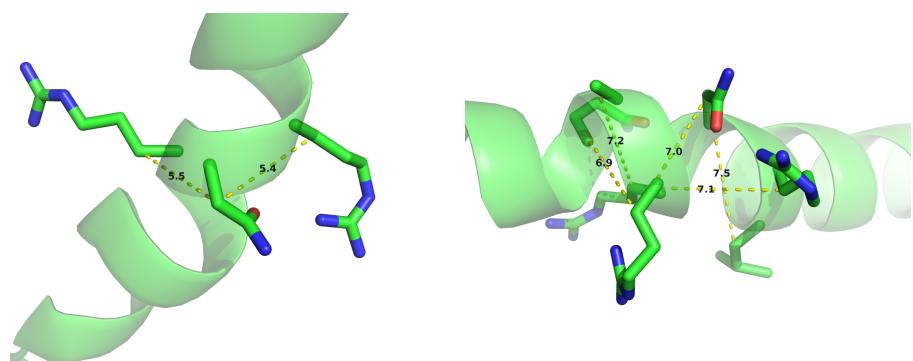


Figure 1.7:  $C_\beta$  distances between neighboring residues in  $\alpha$ -helices. Left: Direct neighbors in  $\alpha$ -helices have  $C_\beta$  distances around  $5.4\text{\AA}$  due to the geometrical constraints from  $\alpha$ -helical architecture. Right: Residues separated by two positions ( $|i-j| = 2$ ) are less geometrically restricted to  $C_\beta$  distances between  $7\text{\AA}$  and  $7.5\text{\AA}$ .

The rigid  $C_\beta$  distance definition of a contact is a very rough measure of true physical interactions between amino acid sidechains. More importantly, interactions between sidechains depend on their physico-chemical properties, on their orientation and different environments within proteins (see section ??) [139]. A simple  $C_\beta$  distance threshold not only misses to reflect biological interaction preferences of amino acids but also provides a questionable gold-standard for benchmarking. Other distance thresholds and definitions for physical contacts (e.g minimal atomic distances or distance between functional groups) have been studied as well. In fact, Duarte and colleagues found that using a  $C_\beta$  distance threshold between 9 Å and 11 Å yields optimal results when predicting the 3D structure from the respective contacts [101]. Anishchenko and colleagues analysed false positive predictions with respect to a minimal atom distance threshold < 5 Å, as they found that this cutoff optimally defines direct physical interactions of residue pairs [140].

Another issue concerns structural variation within a protein family. Evolutionary couplings are inferred from all family members in the [MSA](#) and therefore predicted contacts might be physical contacts in one family member but not in another. Anishchenko et al. could show that more than 80% of false positives at intermediate distances (minimal heavy atom distance 5-15 Å) are true contacts in at least one homolog structure [140]. Therefore, choosing the right trade-off between sensitivity and specificity when generating alignments is a crucial step as well as choosing the target protein structure for evaluation.

Finally, an important aspect not considered in the standard benchmarks is the spread of predicted contacts. It is perfectly possible to improve precision of predicted contacts without translating this improvement to better structural models. The reason being that structurally redundant contacts, that is contacts in the immediate sequence neighborhood of other contacts, do not give additional information to constrain the fold [36,52,95]. For example, given a contact between residues  $i$  and  $j$ , there is hardly an added value knowing that there is a contact between residues  $i+1$  and  $j+1$  when it comes to predicting the overall topology. This observation is highly relevant for deep learning methods due to their unique ability to abstract higher order interactions and recognize contact patterns. Several measures of the contact spread have been developed, like the mean euclidian distance between true and predicted contacts, but are not commonly evaluated yet [52,105].

## 1.7 Challenges for Coevolutionary Inference

Coevolution methods face several challenges when interpreting the covariation signals obtained from a [MSA](#). Some of these challenges have been successfully met (e.g. disentangling transitive effects with global statistical models), others are still open or open up new perspectives, such as dissecting different sources of coevolution signals.

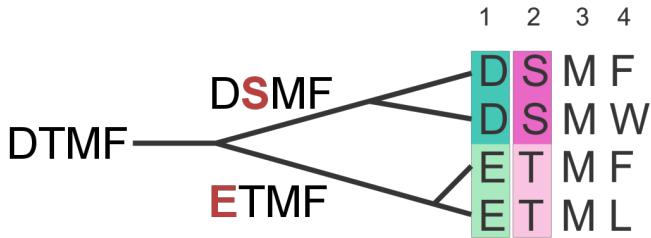


Figure 1.8: The phylogenetic dependence of closely related sequences can produce covariation signals. Here, two independent mutation events (highlighted in red) in two branches of the tree result in a perfect covariation signal for two positions.

### 1.7.1 Phylogenetic Effects as a Source of Noise

Sequences in [MSAs](#) do not represent independent samples of a protein family. In fact, there is selection bias from sequencing species of special interest (e.g human pathogens) or sequencing closely related species, e.g multiple strains. This uneven sampling of a protein family’s sequence space leaves certain regions unexplored whereas others are statistically overrepresented [51,53,141]. Furthermore, due to their evolutionary relationship, sequences of a protein family have a complicated dependence structure. Closely related sequences can cause spurious correlations between positions, as there was not sufficient time for the sequences to diverge from their common ancestor [13,17,18]. Figure 1.8 illustrates a simplified example, where dependence of sequences due to phylogeny leads to a covariation signal. To reduce the effects of redundant sequences, a popular sequence reweighting strategy has been found to improve contact prediction performance, where every sequence receives a weight that is the inverse of the number of similar sequences according to an identity threshold (see section 4.3) [21,51,53,142].

### 1.7.2 Entropic Effects as a Source of Noise

Another source for noise is entropy bias that is closely linked to phylogenetic effects. By nature, methods detecting signals from correlated mutations rely on a certain degree of covariation between sequence positions [18]. Highly conserved interactions pose a conceptual challenge, as changes from one amino acid to another cannot be detected if sequences do not vary. This results in generally higher co-evolution signals from positions with high entropy and underestimated signals for highly conserved interactions [11]. Several heuristics have been proposed to reduce entropy effects, such as Row-Column-Weighting (RCW) [13] or Average Product Correction (APC) [14] (see section 1.4.4).

### 1.7.3 Finite Sampling Effects

Spurious correlations can arise from random statistical noise and blur true co-evolution signals especially in low data scenarios. Consequently, false positive predictions attributable to random noise accumulate for protein families comprising low numbers of homologous sequences. This relationship was confirmed in many

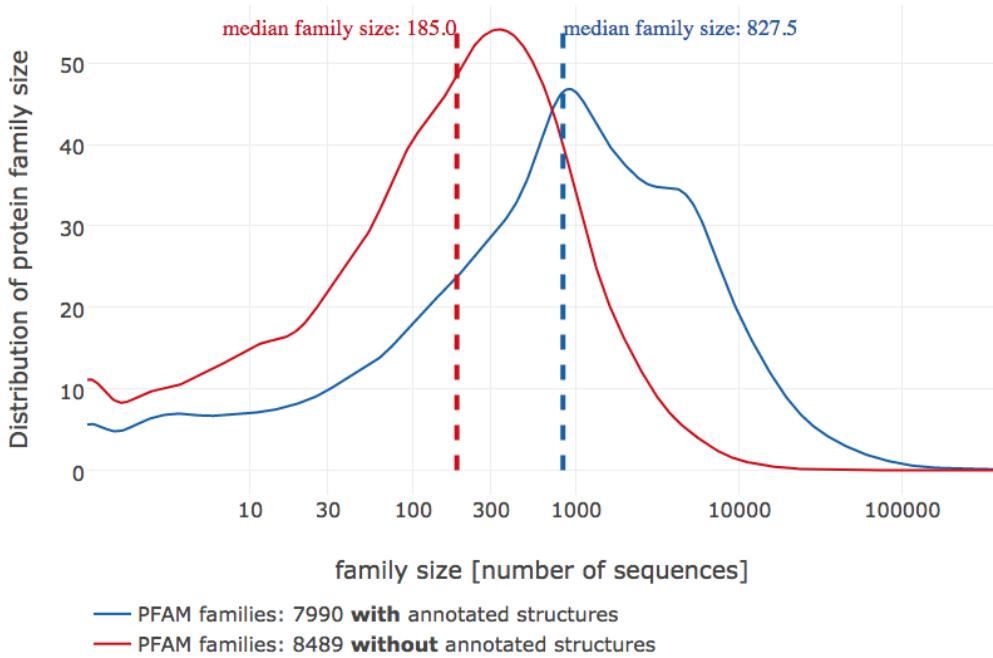


Figure 1.9: Distribution of PFAM family sizes. Less than half of the families in PFAM (7990 compared to 8489 families) do not have an annotated structure. The median family size in number of sequences for families with and without annotated structures is 185 and 827 respectively. Data taken from PFAM 31.0 (March 2017, 16712 entries) [144].

studies and as a rule of thumb it has been argued that proteins with  $L$  residues need at least  $5L$  sequences in order to obtain confident predictions that can be used for protein structure prediction [59,141]. Recently it was shown that precision of predicted contacts saturates for protein families with more than  $10^3$  diverse sequences and that precision is only dependent on protein length for families with small number of sequences [140].

Interesting targets for contact prediction are protein families without any associated structural information. As can be seen in Figure 1.9, those protein families generally comprise low numbers of homologous sequences with a median of 185 sequences per family and are thus susceptible to finite sampling effects.

With the rapidly increasing size of protein sequence databases (see section ??) the number of protein families with enough sequences for accurate contact predictions will increase steadily [59,143]. Nevertheless, because of the already mentioned sequencing biases, better and more sensitive statistical models are indispensable to extend the applicability domain of coevolutionary methods.

#### 1.7.4 Multiple Sequence Alignments

A correct MSA is the essential starting point for coevolution analysis as incorrectly aligned residues will confound the true signal. Highly sensitive and accurate alignment tools such as HHblits generate high quality alignments suitable for contact prediction [145]. However, there are certain subtleties to be kept in mind when generating alignments.

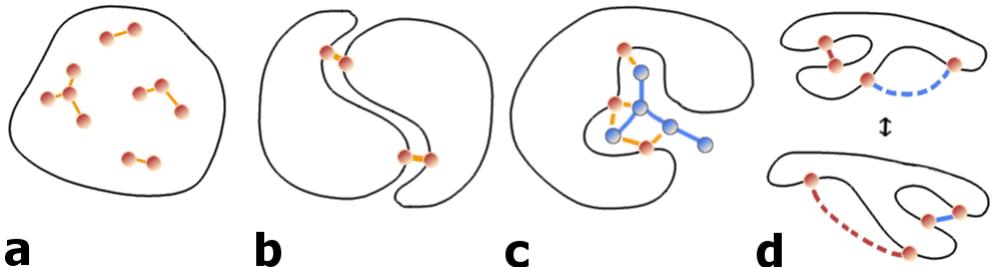


Figure 1.10: Possible sources of coevolutionary signals. **a)** Physical interactions between intra-domain residues. **b)** Interactions across the interface of predominantly homo-oligomeric complexes. **c)** Interactions mediated by ligands or metal atoms. **d)** Transient interactions due to conformational flexibility.

For example, proteins with repeated stretches of amino acids or with regions of low complexity are notoriously hard to align. Especially, repeat proteins have been found to produce many false positive contact predictions [140]. Therefore, MSAs need to be generated with great care and covariation methods need to be tailored to these specific types of proteins [146,147].

Furthermore, sensitivity of sequence search is critically dependent on the research question at hand and on the protein family under study. Many diverse sequences in general increase precision of predictions [137,148]. However, deep alignments can capture coevolutionary signals from different subfamilies [110]. If only a specific subfamily is of interest, many false predictions might arise from strong coevolutionary signals specific to another subfamily that constitutes a prominent subset in the alignment [132]. Therefore, a trade-off between specificity and diversity of the alignment is required to reach optimal results [78].

Another intrinsic characteristic of MSAs are repeated stretches of gaps that result from commonly utilized gap-penalty schemes assigning large penalties to insert a gap and lower penalties to gap extensions. Most statistical coevolution models for contact prediction treat gaps as the 21st amino acid. This introduces an imbalance as gaps and amino acids express different behaviours which can result in gap-induced artefacts [67].

### 1.7.5 Alternative Sources of Coevolution

Coevolutionary signals can not only arise from intra-domain contacts, but also from other sources, like homo-oligomeric contacts, alternative conformations, ligand-mediated interactions or even contacts over hetero-oligomeric interfaces (see Figure 1.10) [141]. With the objective to predict physical contacts it is therefore necessary to identify and filter these alternative sources of coevolutionary couplings.

Many proteins form homo-oligomers with evolutionary conserved interaction surfaces (Figure 1.10 b). Currently it is hard to reliably distinguish intra- and intermolecular contacts [110]. Anishchenko et al. found that approximately one third of strong co-evolutionary signals between residue pairs at long distances (minimal heavy atom distance  $>15\text{\AA}$ ) can be attributed to interactions across homo-oligomeric interfaces [140]. Several studies specifically analysed co-evolution across

homo-oligomeric interfaces for proteins of known structure by filtering for residue pairs with strong couplings at long distances [78,84,110,114,115,149] or used co-evolutionary signals to predict homo-dimeric complexes [111].

It has been proposed that co-evolutionary signals can also arise from ligand or atom mediated interactions between residues or from critical interactions in intermediate folding states (Figure 1.10 c) [142,150]. Confirming this hypothesis, a study showed that the cumulative strength of couplings for a particular residue can be used to predict functional sites [78,141].

Another important aspect is conformational flexibility (Figure 1.10 c). PDB structures used to evaluate coevolution methods represent only rigid snapshots taken in an unnatural crystalline environment. Yet proteins possess huge conformational plasticity and can adopt distinct alternative conformations or adapt shape when interacting with other proteins in an induced fit manner [151]. Several studies demonstrated successfully that coevolutionary signals can capture interactions specific to different distinct conformations [51,78,113,115].



# 2

## Interpretation of Coupling Matrices

Contact prediction methods learning a *Potts model* for the [MSA](#) of a protein family, map the inferred  $20 \times 20$  dimensional coupling matrices  $w_{ij}$  onto scalar values to obtain contact scores for each residue pair as outlined in section [1.4.4](#). As a result, the full information contained in coupling matrices is lost, such as the contribution of individual couplings  $w_{ijab}$ , whether a coupling is positive or negative, higher order dependencies between couplings or possibly biological meaningful signals. The following sections give some intuition for the information contained in coupling matrices.

### 2.1 Single Coupling Values Carry Evidence of Contacts

Given the success of [DCA](#) methods, it is clear that the inferred couplings  $\mathbf{w}_{ij}$  are good indicators of spatial proximity for residue pairs. As described in section [1.4.4](#), a contact score  $C_{i,j}$  for a residue pair  $(i, j)$  is commonly computed as the Frobenius norm over the coupling matrix,  $C_{i,j} = \|\mathbf{w}_{ij}\|_2 = \sqrt{\sum_{a,b=1}^{20} w_{ijab}^2}$ .

The plots in Figure [2.1](#) show the correlation of squared coupling values  $w_{ijab}^2$  with binary contact class (contact=1, non-contact=0) and the standard deviation of squared coupling values  $w_{ijab}^2$  for contacts computed on a dataset of 100.000 residue pairs per class (for details see methods section [4.7.1](#)). All couplings have a weak positive class correlation, meaning the stronger the squared coupling value, the more likely a contact can be inferred. Correlation is weak because most couplings  $w_{ijab}$  are close to zero since typically only few amino acid pairings per residue pair carry evidence and produce a signal. Generally, couplings that involve an aliphatic amino acid such as isoleucine (I), leucine (L), valine (V) or an alanine (A) express the strongest class correlation. In contrast, cysteine pairs (C-C) or pairs involving only the charged residus arginine (R), glutamic acid (E), lysine (K) or aspartic acid (D) correlate only weakly with contact class. Interestingly, for residue pairs being in physical contact, C-C and couplings involving charged

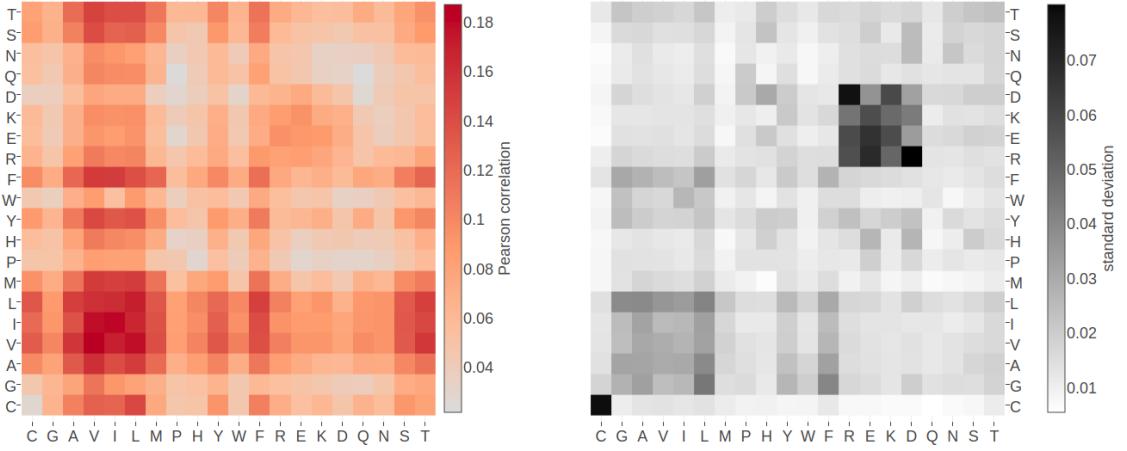


Figure 2.1: **Left** Pearson correlation of squared coupling values ( $w_{ijab}^2$ ) with contact class (contact=1, non-contact=0). **Right** Standard deviation of squared coupling values for residue pairs in contact. Dataset contains 100.000 residue pairs per class (for details see methods section 4.7.1). Amino acids are abbreviated with one-letter code and they are broadly grouped with respect to physico-chemical properties listed in Appendix B.

residues have the highest standard-deviation among all couplings as can be seen in the right plot in Figure 2.1. Standard deviation of squared coupling values from non-contacts shows no relevant patterns and is on average one magnitude smaller than for the contact class (see Appendix Figure D.1).

Different couplings are of varying importance for contact inference and have distinct characteristics. When looking at the raw coupling values (without squaring), these characteristics become even more pronounced. The plots in Figure 2.2 show the correlation of raw coupling values  $w_{ijab}$  with contact class and the standard deviation of coupling values for contacts. Standard deviation of coupling values for non-contacts shows no relevant patterns and is on average half as big as for the contact class (see Appendix Figure D.2). Interestingly, in contrast to the findings for squared coupling values, couplings for charged residue pairs, involving arginine (R), glutamic acid (E), lysine (K) and aspartic acid (D), have the strongest class correlation (positive and negative), whereas aliphatic coupling pairs correlate to a much lesser extent. This implies that squared coupling value is a better indicator of a contact than the raw signed coupling value for aliphatic couplings. On the contrary, the raw signed coupling values for charged residue pairs are much more indicative of a contact than the magnitude of their squared values. Raw couplings for cysteine (C-C) pairs, proline (P) and tryptophane (W) correlate only weakly with contact class. For these pairs neither a squared coupling value nor the raw coupling value seems to be a good indicator for a contact.

Looking only at correlations can be misleading if there are non-linear patterns in the data, for example higher order dependencies between couplings. For this reason it is advisable to take a more detailed view at coupling matrices and the distributions of their values.

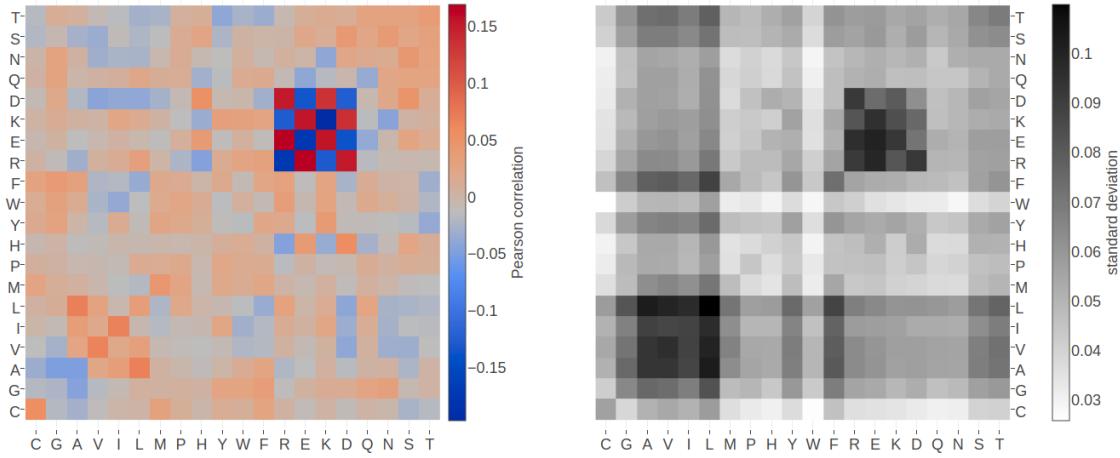


Figure 2.2: **Left** Pearson correlation of raw signed coupling values  $w_{ijab}$  with contact class (contact=1, non-contact=0). **Right** Standard deviation of coupling values for residue pairs in physical contact. Dataset contains 100.000 residue pairs per class (for details see section 4.7.1). Amino acids are abbreviated with one-letter code and they are broadly grouped with respect to physico-chemical properties listed in Appendix B.

## 2.2 Physico-Chemical Fingerprints in Coupling Matrices

The correlation analysis of coupling matrices in the last section revealed that certain couplings are more indicative of a contact than others. Individual coupling matrices for a residue pair that is in physical contact often display striking patterns that agree with the previous findings. These patterns allow a biological interpretation of the coupling values that reveal details of the physico-chemical interdependency between both residues.

Figure 2.3 visualizes the inferred coupling matrix and single potentials  $v_i$  and  $v_j$  for a residue pair  $(i, j)$  computed with the pseudo-likelihood method. The single potentials  $v_{ia}$  and  $v_{ja}$  describe the tendency for each amino acid  $a$  to appear at positions  $i$  and  $j$ , and the couplings  $w_{ijab}$  describe the tendency of amino acid  $a$  at position  $i$  to co-occur with amino acid  $b$  at position  $j$ . A cluster of strong coupling values can be observed for the couplings between the charged residues glutamic acid (E), aspartic acid (D), lysine (K) and arginine (R) and the polar residue glutamine (Q). Positive coupling values arise between positively charged residues (K, R) and negatively charged residues (E, D), whereas couplings between equally charged residues have negative values. These exemplary couplings (E-R, E-K, K-D) perfectly reflect the interaction preference for residues forming salt bridges. Indeed, in the protein structure the first residue (E) forms a salt bridge with the second residue (R) as can be seen in the left plot in Figure 2.5.

Figure 2.4 visualizes the coupling matrix for a pair of hydrophobic residues. Hydrophobic pairings, such as alanine (A) - isoleucine (I), or glycine (G) - isoleucine (I) have strong coupling values but the couplings also reflect a sterical constraint. Alanine is a small hydrophobic residue and it is favoured at both residue positions: it has strong positive single potentials  $v_i(A)$  and  $v_j(A)$  and strong positive

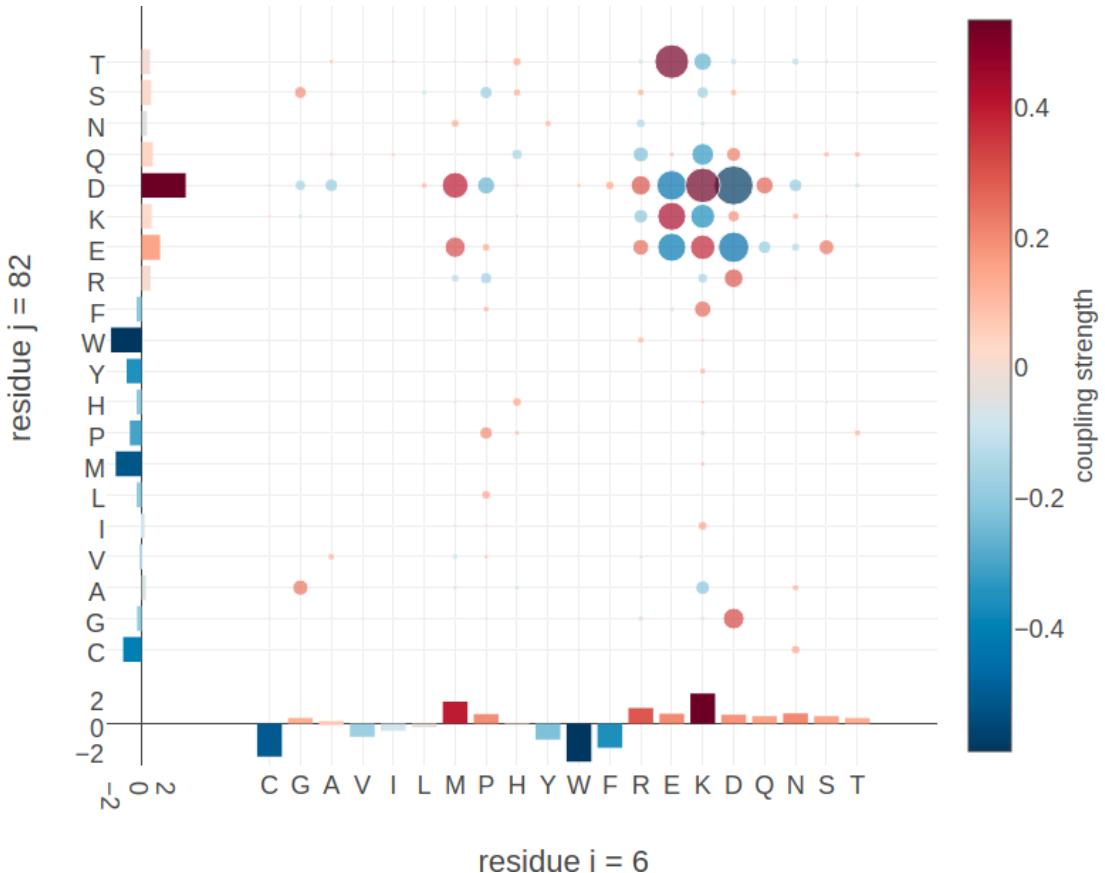


Figure 2.3: Couplings  $w_{ijab}$  and single potentials  $v_{ia}$  and  $v_{ja}$  computed with pseudo-likelihood for residues 6 and 82 in protein chain 1a9x\_A\_05. The matrix shows the  $20 \times 20$  couplings  $w_{ijab}$  with color representing coupling strength and direction (red = positive coupling value, blue = negative coupling value) and diameter of bubbles representing absolute coupling value  $|w_{ijab}|$ . Bars at the x-axis and y-axis correspond to the *Potts* model single potentials  $v_i$  and  $v_j$  respectively. Color reflects the value of single potentials. Amino acids are abbreviated with one-letter code and they are broadly grouped with respect to physico-chemical properties listed in Appendix B.

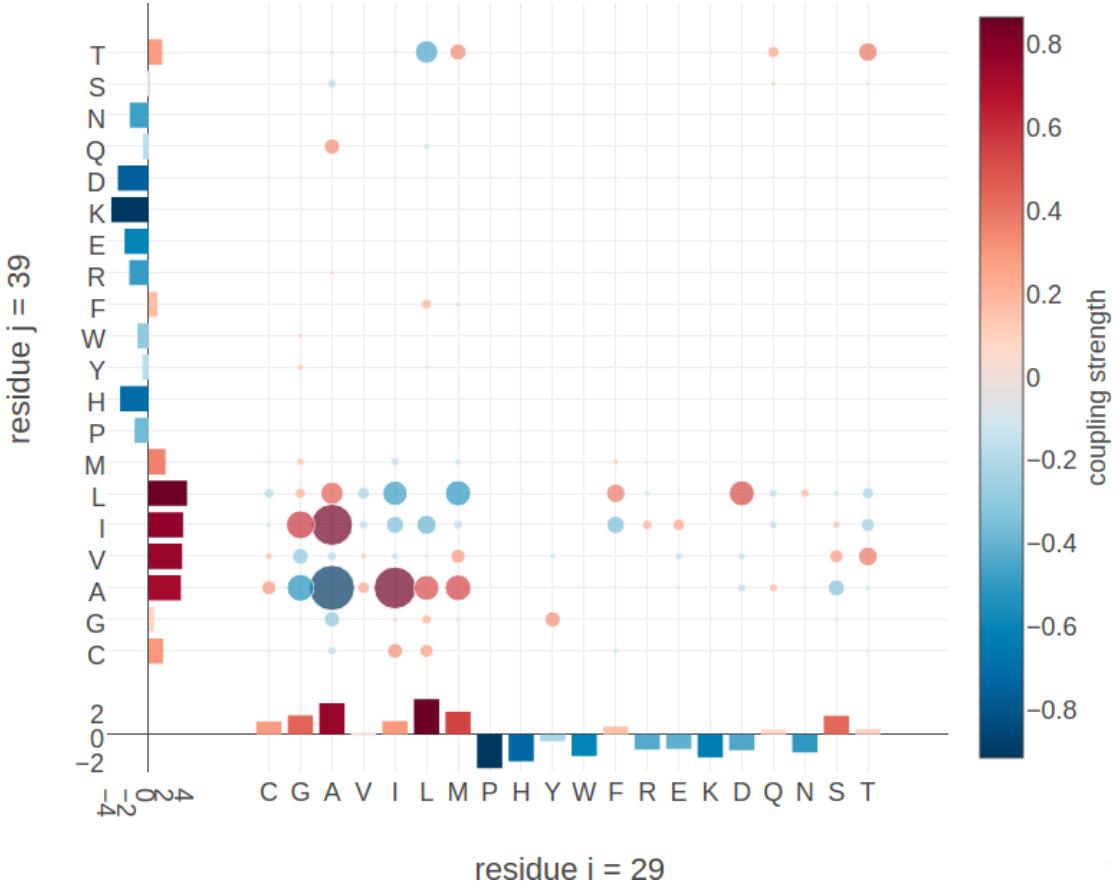


Figure 2.4: Couplings  $w_{ijab}$  and single potentials  $v_{ia}$  and  $v_{ja}$  computed with pseudo-likelihood for residues 29 and 39 in protein chain 1ae9\_A\_00. The matrix shows the 20x20 couplings  $w_{ijab}$  with color representing coupling strength and direction (red = positive coupling value, blue = negative coupling value) and diameter of bubbles representing absolute coupling value  $|w_{ijab}|$ . Bars at the x-axis and y-axis correspond to the *Potts* model single potentials  $v_i$  and  $v_j$  respectively. Color reflects the value of single potentials. Amino acids are abbreviated with one-letter code and they are broadly grouped with respect to physico-chemical properties listed in Appendix B.

couplings with isoleucine (I), leucine (L) and methionine (M). But alanine is disfavoured to appear at both positions at the same time since the A-A coupling is negative. Figure 2.5 illustrates the location of the two residues in the protein core. Here, hydrophobic residues are densely packed and the limited space allows for only small hydrophobic residues.

Many more biological interpretable signals can be identified from coupling matrices, including pi-cation interactions (see Appendix E.1), aromatic-proline interactions (see Appendix E.3), sulfur-aromatic interactions or disulphide bonds (see Appendix E.2).

Coucke and colleagues performed a thorough quantitative analysis of coupling matrices selected from confidently predicted residue pairs [152]. They showed that eigenmodes obtained from a spectral analysis of averaged coupling matrices are closely related to physico-chemical properties of amino acid interactions, like elec-

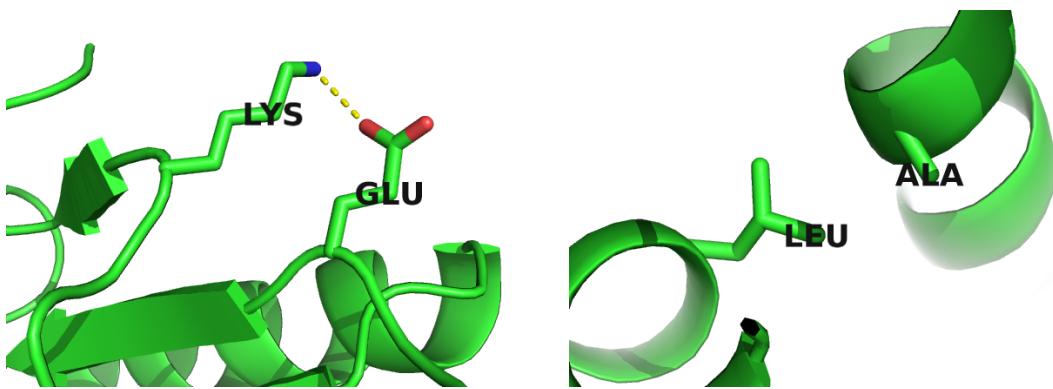


Figure 2.5: Interactions between protein side chains. **Left:** residue 6 (E) forms a salt bridge with residue 82 (R) in protein chain 1a9x\_A\_05. **Right:** residue 29 (A) and residue 39 (L) within the hydrophobic core of protein chain 1ae9\_A\_00.

trostaticity, hydrophobicity, steric interactions or disulphide bonds. By looking at specific populations of residues, like buried and exposed residues or residues from specific protein classes (small, mainly  $\alpha$ , etc), the eigenmodes of corresponding coupling matrices are found to capture very characteristic interactions for each class, e.g. rare disulfide contacts within small proteins and hydrophilic contacts between exposed residues. Their study confirms the qualitative observations presented above that amino acid interactions can leave characteristic physico-chemical fingerprints in coupling matrices.

## 2.3 Coupling Profiles Vary with Distance

Analyses in the previous sections showed that certain coupling values correlate more or less strong with contact class and that coupling matrices for contacts express biologically meaningful patterns.

More insights can be obtained by looking at the distribution of distinct coupling values for contacts, non-contacts and arbitrary populations of residue pairs. Figure 2.6 shows the distribution of selected couplings for filtered residue pairs with  $C_\beta - C_\beta$  distances  $< 5\text{\AA}$  (see methods section 4.7.2 for details). The distribution of R-E and E-E coupling values is shifted and skewed towards positive and negative values respectively. This is in accordance with attracting electrostatic interactions between the positively charged side chain of arginine and the negatively charged side chain of glutamic acid and also with repulsive interactions between the two negatively charged glutamic acid side chains.

Coupling values for cysteine pairs (C-C) have a broad distribution that is skewed towards positive values, reflecting the strong signals obtained from covalent disulphide bonds. The broad distribution for C-C, R-E and E-E agrees with the observation in section 2.1 that these specific coupling values have large standard deviations and that for charged residue pairings the signed coupling value is a strong indicator of a contact.

Hydrophobic pairs like V-I have an almost symmetric coupling distribution, confirming the finding that the direction of coupling is not indicative of a true contact

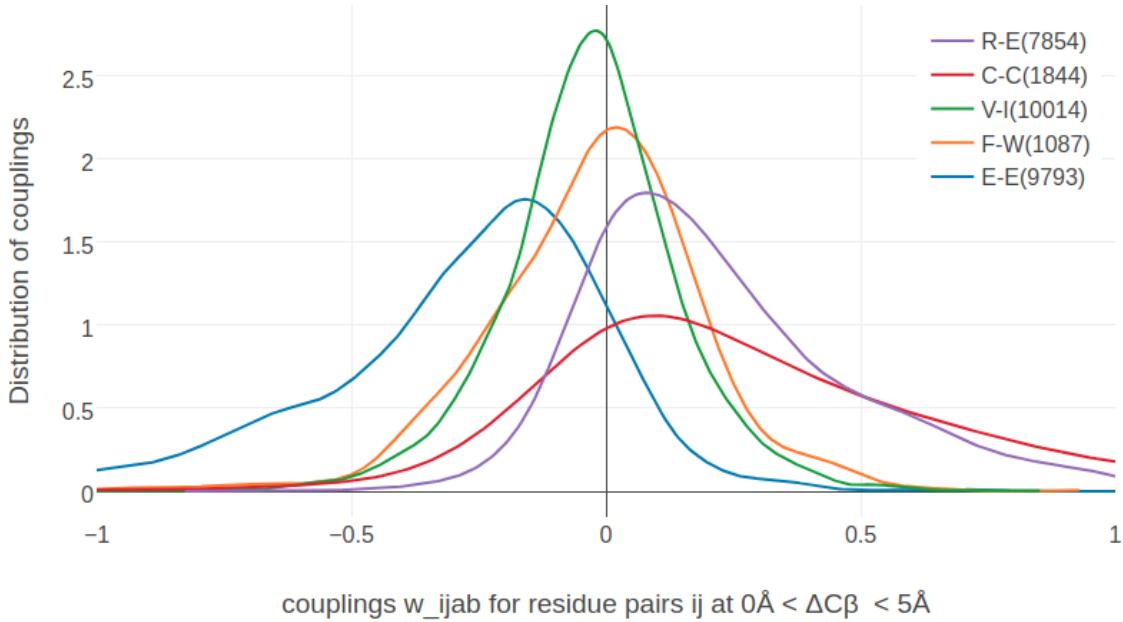


Figure 2.6: Distribution of selected couplings for filtered residue pairs with  $C_\beta - C_\beta$  distances  $< 5\text{\AA}$  (see methods section 4.7.2 for details). Number of coupling values used to determine the distribution is given in brackets in the legend. R-E = couplings for arginine and glutamic acid pairs, C-C = coupling for cystein residue pairs, V-I = coupling for valine and isoleucine pairs, F-W = coupling for phenylalanine and tryptophane pairs, E-E = coupling for glutamic acid residue pairs.

whereas the strength of the coupling is. The hydrophobic effect that determines hydrophobic interactions is not specific or directed. Therefore, hydrophobic interaction partners can commonly be substituted by other hydrophobic residues, which explains the not very pronounced positive coupling signal compared to more specific interactions, e.g. ionic interactions. It is not clear though, why hydrophobic pairs have an equally strong negative coupling signal at this distance range because this speaks against the hypothesis that hydrophobic pairs are commonly interchangeable. A vague explanation could be that a location in the tightly packed protein core calls for other very specific constraints, e.g. sterical fit or contact number, besides hydrophobic properties that are prohibitive for a particular hydrophobic residue at a certain position.

The distribution of aromatic coupling values like F-W is slightly skewed towards negative values, accounting for steric hindrance of their large sidechains at small distances. The yet very pronounced positive coupling signal for the bulky aromatic residues at this short distance range is not clear. The bulky planar aromatic rings of two aromatic residues often point away from each other when their  $C_\beta - C_\beta$  distances are small to avoid steric hindrance (see Appendix Figure E.5). A positive coupling signal might originate from other structural constraints from the local environment affecting both sidechains, similar to the scenario hypothetically explaining the negative coupling signal for hydrophobic residues.

In an intermediate  $C_\beta$  distance range between  $8\text{\AA}$  and  $12\text{\AA}$  the distributions for all coupling values are centered close to zero and are less broad. The distributions are still shifted and skewed, but less pronounced compared to the distributions at

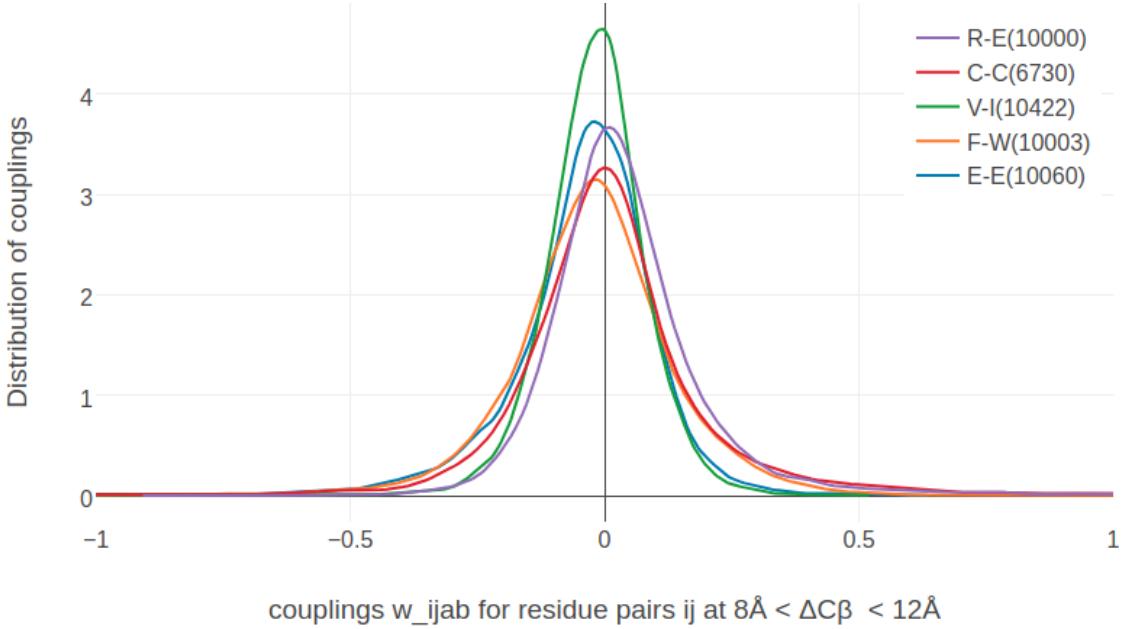


Figure 2.7: Distribution of selected couplings for filtered residue pairs with  $C_\beta - C_\beta$  distances between  $8\text{\AA}$  and  $12\text{\AA}$  (see methods section 4.7.2 for details). Number of coupling values used to determine the distribution is given in brackets in the legend. Couplings are the same as in Figure 2.6.

$C_\beta - C_\beta$  distances  $< 5\text{\AA}$ . For aromatic pairs like F-W, the distribution of coupling values has very long tails, suggesting rare but strong couplings for aromatic side chains at this distance.

Figure 2.8 shows the distribution of selected couplings for residue pairs far apart in the protein structure ( $C_\beta - C_\beta$  distances  $> 20\text{\AA}$ ).

The distribution for all couplings is centered at zero and has small variance. Only for C-C coupling values, the distribution has a long tail for positive values, presumably arising from the fact that the maximum entropy model cannot distinguish highly conserved signals of multiple disulphide bonds within a protein. This observation also agrees with the previous finding in section 2.1 that C-C coupling values, albeit having large standard-deviations, correlate only weakly with contact class. The same arguments apply to couplings of aromatic pairs that have a comparably broad distribution and do not correlate strongly with the contact class. The strong coevolution signals for aromatic pairs even at high distance ranges might result from some kind of cooperative effects. Aromatic residues are known to form network-like structures in the protein core that stabilize protein structure [153]. An example is given in Appendix Figure E.4. A possible explanation might be that the *Potts model* is limited to learning single positions and pairwise correlations. An extension to higher order couplings might resolve these cooperative effects observed between residues in the protein core.

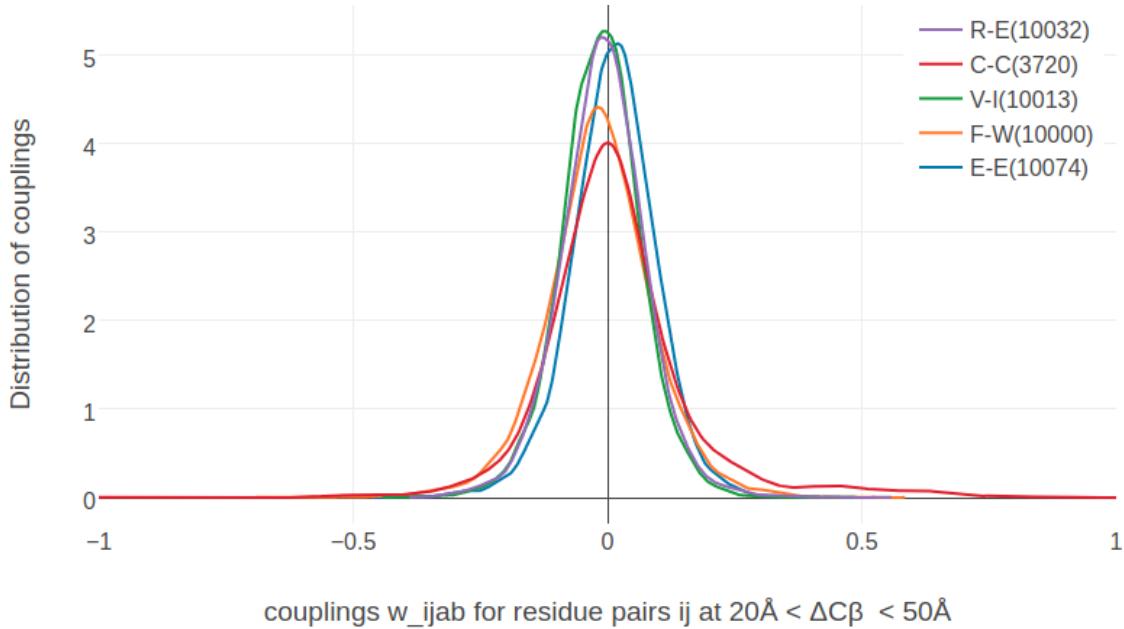


Figure 2.8: Distribution of selected couplings for filtered residue pairs with  $C_\beta - C_\beta$  distances between  $20\text{\AA}$  and  $50\text{\AA}$  (see methods section 4.7.2 for details). Number of coupling values used to determine the distribution is given in brackets in the legend. Couplings are the same as in Figure 2.6.

## 2.4 Higher Order Dependencies Between Couplings

The analyses in the previous sections focused on single coupling values picked from the  $20 \times 20$ -dimensional coupling matrices  $\mathbf{w}_{ij}$ . As mentioned before, analysing only single dimensions might be misleading when variables are dependent on each other and further insights might be concealed in higher order relationships. Unfortunately, it is not possible to reasonably visualize high dimensional coupling matrices.

Exploring two dimensional coupling scatter plots strengthens the observation that couplings matrices contain signals that reflect biological relevant amino acid interactions. The plots in the top row in Figure 2.9 show the distribution of couplings for filtered residue pairs with  $C_\beta - C_\beta$  distances  $< 8\text{\AA}$  between the ionic pairings of E-R and R-E and between the ionic pairing R-E and the equally charged residues E-E, respectively. Coupling values for R-E and E-R are positively correlated with predominantly positive values. This means when the amino acid pair R-E is frequently observed at two positions  $i$  and  $j$ , then it also likely that the amino acid pair E-R can be frequently observed. This situation indicates an important ionic interaction whereby the location of the positively and negatively charged residue at position  $i$  or  $j$  is irrelevant.

On the contrary, coupling values for R-E and E-E are negatively correlated, with positive values for R-E and negative values for E-E. This distribution can be interpreted with frequently occurring amino acid pairs R-E at two positions  $i$  and  $j$  while at the same time the amino acid pair E-E cannot be observed. Again, this

situation coincides with amino acid pairings that would be expected for an ionic interaction.

The bottom left plot in Figure 2.9 shows the distribution between couplings for the hydrophobic pairings I-L and V-I that is almost symmetric and broadly centered around zero. Coupling distributions for residue pairs that are not physically interacting ( $C_\beta \gg 8\text{\AA}$ ) resemble the distribution for hydrophobic pairings in that there is no correlation, but at high distance the distributions are much tighter centered around zero (bottom right plot in Figure 2.9).

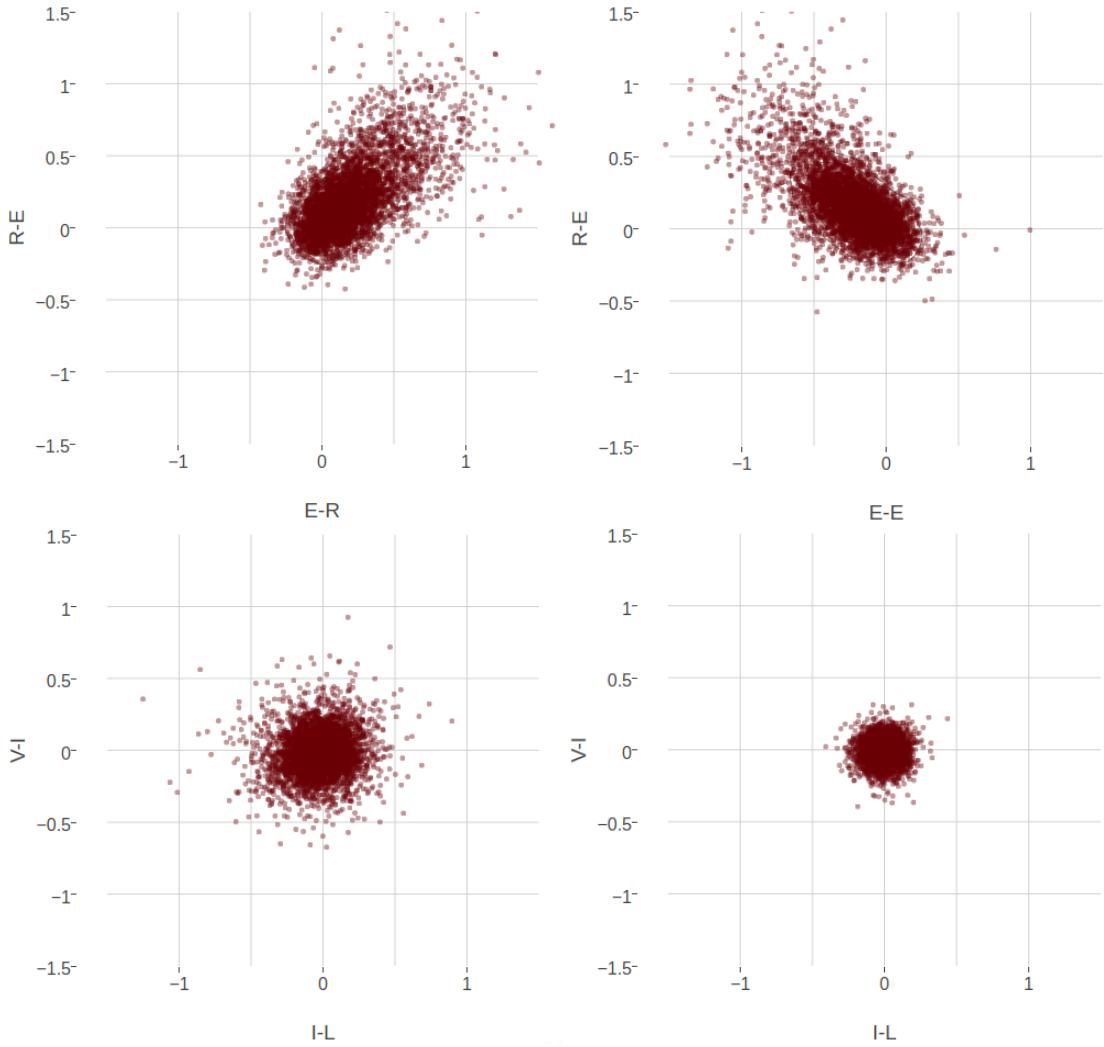


Figure 2.9: Two-dimensional distribution of approximately 10000 coupling values computed with pseudo-likelihood. **Top Left** The 2-dimensional distribution of couplings E-R and R-E for residue pairs with  $C_\beta - C_\beta$  distances  $< 8\text{\AA}$  is almost symmetric and the coupling values are positively correlated. **Top Right** The 2-dimensional distribution of couplings E-R and E-E for residue pairs with  $C_\beta - C_\beta$  distances  $< 8\text{\AA}$  is almost symmetric and the coupling values are negatively correlated. **Bottom Left** The 2-dimensional distribution of couplings I-L and V-I for residue pairs with  $C_\beta - C_\beta$  distances  $< 8\text{\AA}$  is symmetrically distributed around zero without visible correlation. **Bottom Right** The 2-dimensional distribution of couplings I-L and V-I for residue pairs with  $C_\beta - C_\beta$  distances  $> 20\text{\AA}$  is tightly distributed around zero.



# 3

## Contact Prior

The wealth of successful meta-predictors presented in section 1.3 highlights the importance to exploit other sources of information apart from coevolution statistics. Much information about residue interactions is typically contained in single position features that can be predicted from local sequence profiles, such as secondary structure, solvent accessibility or contact number, and in pairwise features such as the contact prediction scores for residue pairs  $(i, j)$  from a simple local statistical methods as presented in section 1.1.

For example, predictions of secondary structure elements and solvent accessibility are used by almost all modern machine learning predictors, such as MetaPsicov [40], NeBCon [43], EPSILON-CP [42], PconsC3 [38]. Other frequently used sequence derived features include pairwise contact potentials, sequence separation and conservation measures such as column entropy [40,43,154].

In the following sections I present a random forest classifier that uses sequence derived features to distinguish contacts from non-contacts. Methods section 4.11 lists all features used to train the classifier including the aforementioned standard features as well as some novel features.

The probabilistic predictions of the random forest model can be introduced directly as prior information into the Bayesian statistical model presented in the last section ?? to improve the overall prediction accuracy in terms of posterior probabilities. Furthermore, contact scores from coevolution methods can be added as additional feature to the random forest model in order to elucidate how much the combined information improves prediction accuracy over the single methods.

### 3.1 Random Forest Classifiers

Random Forests are supervised machine learning methods that belong to the class of ensemble methods [155–157]. They are easy to implement, fast to train and can handle large numbers of features due to implicit feature selection [158].

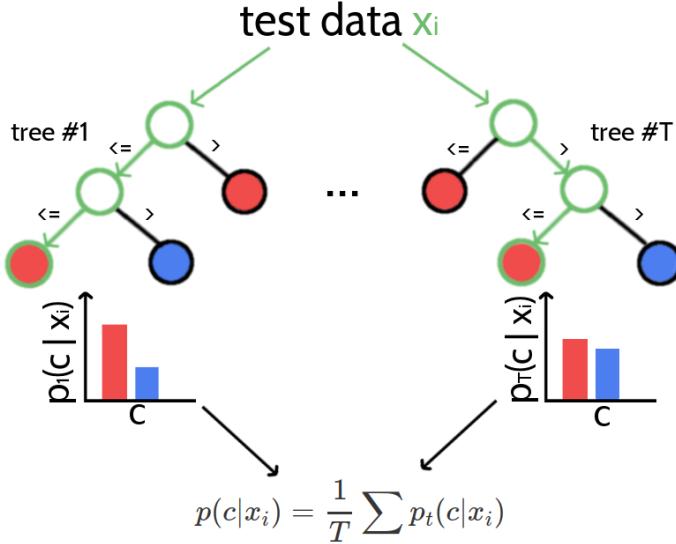


Figure 3.1: Classifying new data with random forests. A new data sample is run down every tree in the forest until it ends up in a leaf node. Every leaf node has associated class probabilities  $p(c)$  reflecting the fraction of training samples at this leaf node belonging to every class  $c$ . The color of the leaf nodes reflects the class with highest probability. The predictions from all trees in form of the class probabilities are averaged and yield the final prediction.

Ensemble methods combine the predictions of several independent base estimators with the goal to improve generalizability over a single estimator. Random forests are ensembles of decision trees where randomness is introduced in two ways:

1. every tree is build on a random sample that is drawn with replacement from the training set and has the same size as the training set (i.e., a bootstrap sample)
2. every split of a node is evaluated on a random subset of features

A single decision tree, especially when it is grown very deep is highly susceptible to noise in the training set and therefore prone to overfitting which results in poor generalization ability. As a consequence of randomness and averaging over many decision trees, the variance of a random forest predictor decreases and therefore the risk of overfitting [159]. It is still advisable to restrict the depth of single trees in a random forest, not only to counteract overfitting but also to reduce model complexity and to speedup the algorithm.

Random forests are capable of regression and classification tasks. For classification, predictions for new data are obtained by running each data sample down every tree in the forest and then either apply majority voting over single class votes or averaging the probabilistic class predictions. Probabilistic class predictions of single trees are computed as the fraction of training set samples of the same class in a leaf whereas the single class vote refers to the majority class in a leaf. Figure 3.1 visualizes the procedure of classifying a new data sample.

Typically, *Gini impurity*, which is a computationally efficient approximation to the entropy, is used as a split criterion to evaluate the quality of a split. It measures

the degree of purity in a data set regarding class labels as  $GI = (1 - \sum_{k=1}^K p_k^2)$ , where  $p_k$  is the proportion of class  $k$  in the data set. For every feature  $f$  in the random subset that is considered for splitting a particular node  $N$ , the *decrease in Gini impurity*  $\Delta GI_f$  will be computed as,

$$\Delta GI_f(N_{\text{parent}}) = GI_f(N_{\text{parent}}) - p_{\text{left}} GI_f(N_{\text{left}}) - p_{\text{right}} GI_f(N_{\text{right}})$$

where  $p_{\text{left}}$  and  $p_{\text{right}}$  refers to the fraction of samples ending up in the left and right child node respectively [158]. The feature  $f$  with highest  $\Delta GI_f$  over the two resulting child node subsets will be used to split the data set at the given node  $N$ .

Summing the *decrease in Gini impurity* for a feature  $f$  over all trees whenever  $f$  was used for a split yields the *Gini importance* measure, which can be used as an estimate of general feature relevance. Random forests therefore are popular methods for feature selection and it is common practice to remove the least important features from a data set to reduce the complexity of the model. However, feature importance measured with respect to *Gini importance* needs to be interpreted with care. The random forest model cannot distinguish between correlated features and it will choose any of the correlated features for a split, thereby reducing the importance of the other features and introducing bias. Furthermore, it has been found that feature selection based on *Gini importance* is biased towards selecting features with more categories as they will be chosen more often for splits and therefore tend to obtain higher scores [160].

## 3.2 Hyperparameter Optimization for Random Forest

There are several hyperparameters in a random forest model that need to be tuned to achieve best balance between predictive power and runtime. While more trees in the random forest generally improve performance of the model, they will slow down training and prediction. A crucial hyperparameter is the number of features that is randomly selected for a split at each node in a tree [161]. Stochasticity introduced by the random selection of features is a key characteristic of random forests as it reduces correlation between the trees and thus the variance of the predictor. Selecting many features typically increases performance as more options can be considered for each split, but at the same time increases risk of overfitting and decreases speed of the algorithm. In general, random forests are robust to overfitting, as long as there are enough trees in the ensemble and the selection of features for splitting a node introduces sufficient stochasticity. Overfitting can furthermore be prevented by restricting the depth of the trees, which is known as pruning or by enforcing a minimal leaf node size regarding the minimal number of data samples ending in a leaf node. Again, a positive side-effect of pruning and requiring minimal leaf node size is a speedup of the algorithm. [159]

In the following, I use 5-fold cross-validation to identify the optimal architecture of the random forest. Details about the training set and the cross-validation procedure can be found in method section 4.12. First I assessed performance of models for

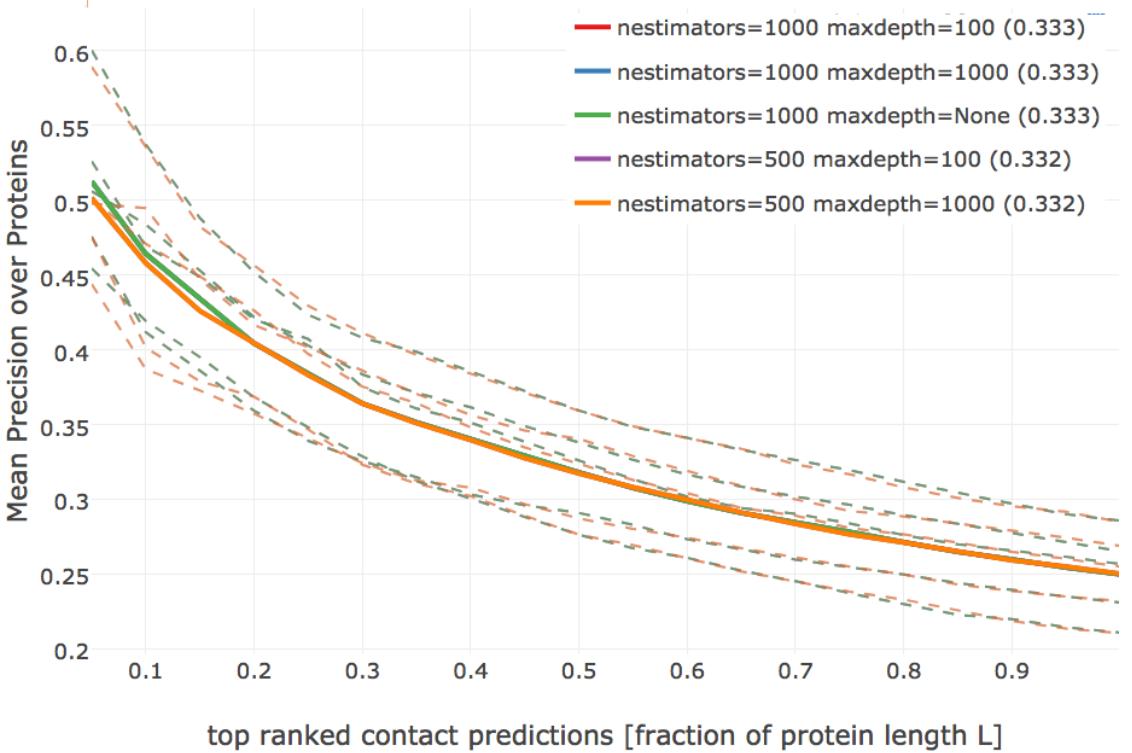


Figure 3.2: Mean precision over 200 proteins against highest scoring contact predictions from random forest models for different settings of  $n\_estimators$  and  $max\_depth$ . Dashed lines show the performance of models that have been learned on the five different subsets of training data. Solid lines give the mean precision over the five models. Only those models are shown that yielded the five highest mean precision values (given in parentheses in the legend). Random forest models with 1000 trees and maximum depth of trees of either 100, 1000 or unrestricted tree depth perform nearly identical (lines overlap). Random forest models with 500 trees and  $max\_depth=10$  or  $max\_depth=100$  perform slightly worse.

combinations of the parameter  $n\_estimators$ , defining the number of trees in the forest and the parameter  $max\_depth$  defining the maximum depth of the trees:

- $n\_estimators \in \{100, 500, 1000\}$
- $max\_depth \in \{10, 100, 1000, None\}$

Figure 3.2 shows that the top five parameter combinations perform nearly identical. Random forests with 1000 trees perform slightly better than models constituting 500 trees, irrespective of the depth of the trees. In order to keep model complexity small, I chose `n_estimators=1000` and `max_depth=100` for further analysis.

Next, I optimized the parameters  $min\_samples\_leaf$ , defining the minimum number of samples required at a leaf node and  $max\_features$ , defining the number of randomly selected features considered for each split using the following settings:

- $min\_samples\_leaf \in \{1, 10, 100\}$
- $max\_features \in \{8, 16, 38, 75\}$  representing  $\sqrt{N}$ ,  $\log 2N$ ,  $0.15N$  and  $0.3N$  respectively with  $N = 250$  being the number of features listed in method section 4.11.

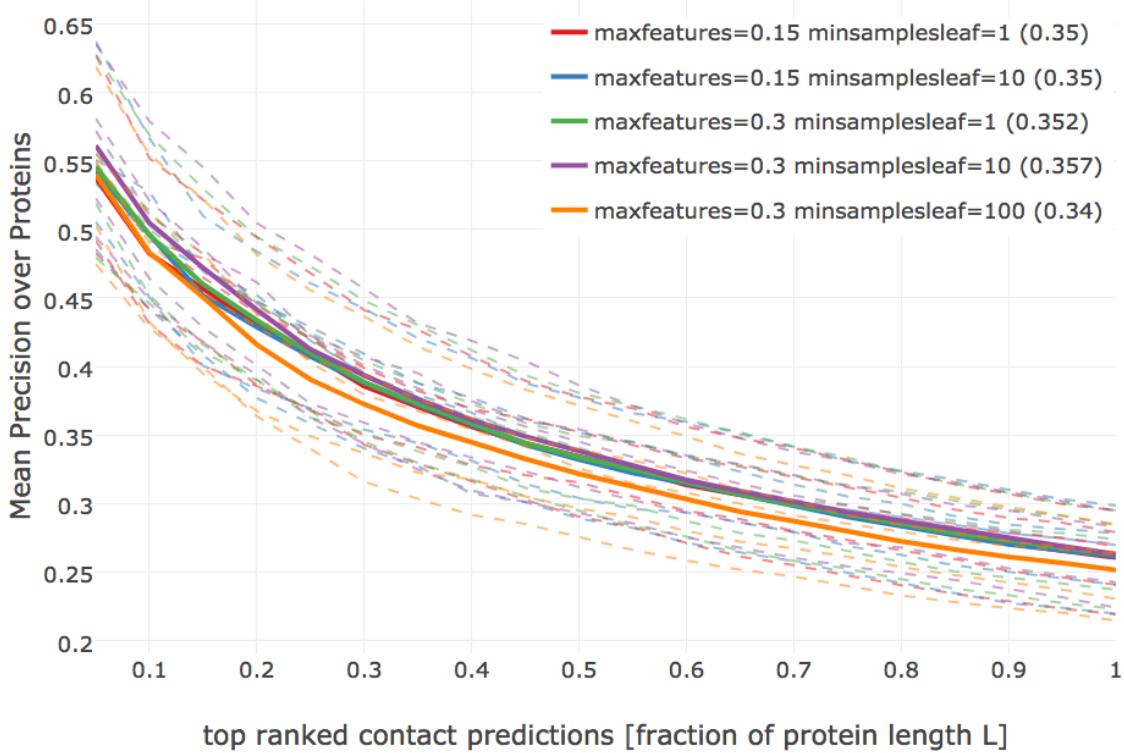


Figure 3.3: Mean precision over 200 proteins against highest scoring contact predictions from random forest models with different settings of *min\_samples\_leaf* and *max\_features*. Dashed lines show the performance of models that have been learned on the five different subsets of training data. Solid lines give the mean precision over the five models. Only those models are shown that yielded the five best mean precision values (given in parentheses in the legend).

Randomly selecting 30% of features (=75 features) and requiring at least 10 samples per leaf gives highest mean precision as can be seen in Figure 3.3. I chose `max_features=0.30` and `min_samples_leaf=10` for further analysis. Tuning the hyperparameters in a different order or on a larger dataset gives similar results.

In a next step I assessed dataset specific settings, such as the window size over which single positions features will be computed, the distance threshold to define non-contacts and the optimal proportions of contacts and non-contacts in the training set. I used the previously identified settings of random forest hyperparameters (`n_estimators=1000`, `min_samples_leaf=10`, `max_depth=100`, `max_features=0.30`).

- proportion of contacts/non-contacts  $\in \{1:2, 1:5, 1:10, 1:20\}$  while keeping total dataset size fixed at 300,000 residue pairs
- window size:  $\in \{5, 7, 9, 11\}$
- non-contact threshold  $\in \{8, 15, 20\}$

As can be seen in appendix F.2 and F.3, the default choice of using a window size of five positions and the non-contact threshold of  $8\text{\AA}$  proves to be the optimal setting. Furthermore, using five-times as many non-contacts as contacts in the

training set results in highest mean precision as can be seen in appendix F.4. These estimates might be biased in a way since the random forest hyperparameters have been optimized on a dataset using exactly these optimal settings.

### 3.3 Evaluating Random Forest Model as Contact Predictor

I trained a random forest classifier on the feature set described in methods section 4.11 and using the optimal hyperparameters identified with 5-fold cross-validation as described in the last section.

Figure 3.4 shows the ranking of the ten most important features according to *Gini importance*. Both local statistical contact scores, *OMES* [162] and *MI* (mutual information between amino acid counts), constitute the most important features besides the mean pair potentials according to Miyazawa & Jernigan [163] and Li&Fang[25]. Further important features include the relative solvent accessibility at both pair positions, the total percentage of gaps at both positions, the correlation between mean isoelectric point property at both positions, sequence separation and the beta-sheet propensity in a window of size five around position i.

Many features have low *Gini importance* scores which means they are rarely considered for splitting a node and can most likely be removed from the dataset. Removing irrelevant features from the dataset is a convenient procedure to reduce model complexity. It has been found, that prediction performance might even increase after removing the most irrelevant features [158]. For example, during the development of *EPSILON-CP*, a deep neural network method for contact prediction, the authors performed feature selection using boosted trees. By removing 75% of the most non-informative features (mostly features related to amino acid composition), the performance of their predictor increased slightly [42]. Other studies have also emphasized the importance of feature selection to improve performance and reduce model complexity [23,25].

As described in methods section 4.12.1, I performed feature selection by evaluating model performance on subsets of features of decreasing importance. Most models trained on subsets of the total feature space perform nearly identical compared to the model trained on all features, as can be seen in Figure 3.5. Performance of the random forest models drops noticeably when using only the 25 most important features. For the further analysis I am using the random forest model trained on the 75 most important features as this model constitutes the smallest set of features while performing nearly identical compared to the model trained on the complete feature set.

Figure 3.6 shows the mean precision for the random forest model trained on the 75 most important features. The random forest model has a mean precision of 0.33 for the top  $0.5 \cdot L$  contacts compared to a precision of 0.47 for pseudo-likelihood. Furthermore, the random forest model improves approximately ten percentage points in precision over the local statistical contact scores, *OMES* and mutual

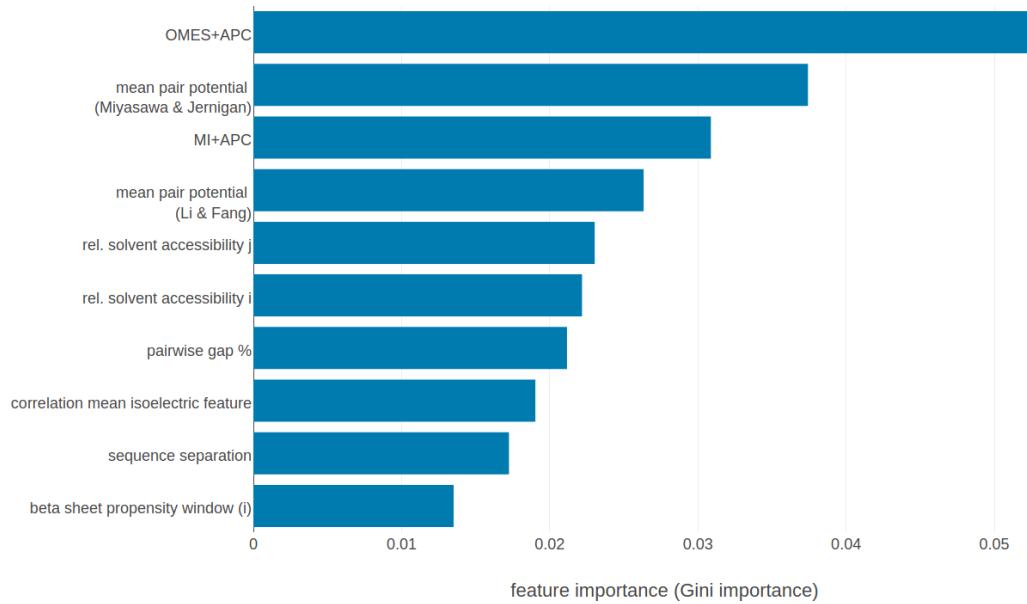


Figure 3.4: Top ten features ranked according to *Gini importance*. **OMES+APC**: APC corrected OMES score according to Fodor&Aldrich [162]. **mean pair potential (Miyasawa & Jernigan)**: average quasi-chemical energy of transfer of amino acids from water to the protein environment [163]. **MI+APC**: APC corrected mutual information between amino acid counts (using pseudo-counts). **mean pair potential (Li&Fang)**: average general contact potential by Li & Fang [25]. **rel. solvent accessibilty i(j)**: RSA score computed with Netsurfp (v1.0) [164] for position i(j). **pairwise gap%**: percentage of gapped sequences at either position i and j. **correlation mean isoelectric feature**: Pearson correlation between the mean isoelectric point feature (according to Zimmermann et al., 1968) for positions i and j. **sequence separation**:  $|j-i|$ . **beta sheet propensity window(i)**: beta-sheet propensity according to Psipred [165] computed within a window of five positions around i. Features are described in detail in methods section 4.11.

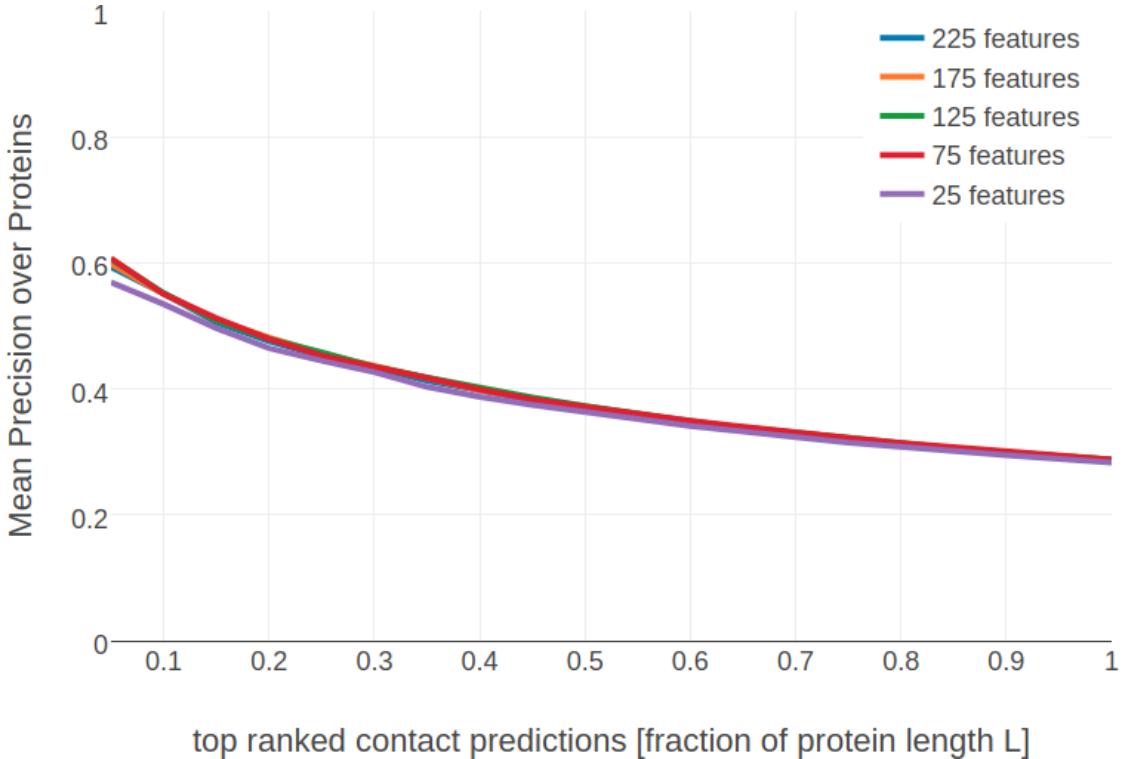


Figure 3.5: Mean precision of top ranked predictions over 200 proteins for random forest models trained on subsets of features of decreasing importance. Subsets of features have been selected as described in methods section 4.12.1.

information (MI). Both methods comprise important features of the random forest model as can be seen in Figure 3.4.

When analysing performance with respect to alignment size it can be found that the random forest model outperforms the pseudo-likelihood score for small alignments (see Figure 3.7).

Both, local statistical models *OMES* and *MI* also perform weak on small alignments, leading to the conclusion that the remaining sequence derived features are highly relevant when the alignment contains only few sequences. This finding is expected, as it is well known that models trained on simple sequence features perform almost independent of alignment size [42].

Figure 3.7 showed that the random forest predictor improves over the pseudo-likelihood coevolution method when the alignment consists of only few sequences. In order to assess this improvement in a more direct manner, it is possible to build a combined random forest predictor that is not only trained on the sequence derived features but also on the pseudo-likelihood contact score as an additional feature. As expected, the pseudo-likelihood score comprises the most important feature in the model (see Appendix Figure F.1) followed by the same sequence features that were found in the previous analysis in Figure 3.4. Models trained on subsets of features as described in method section 4.12.1 perform equally well as the model trained on the complete set of features (see Appendix Figure F.2). Only the model trained on the 26 most important features has slightly decreased precision for the top  $L/10$  ranked contacts. The model trained on 76 features was selected as the final model.

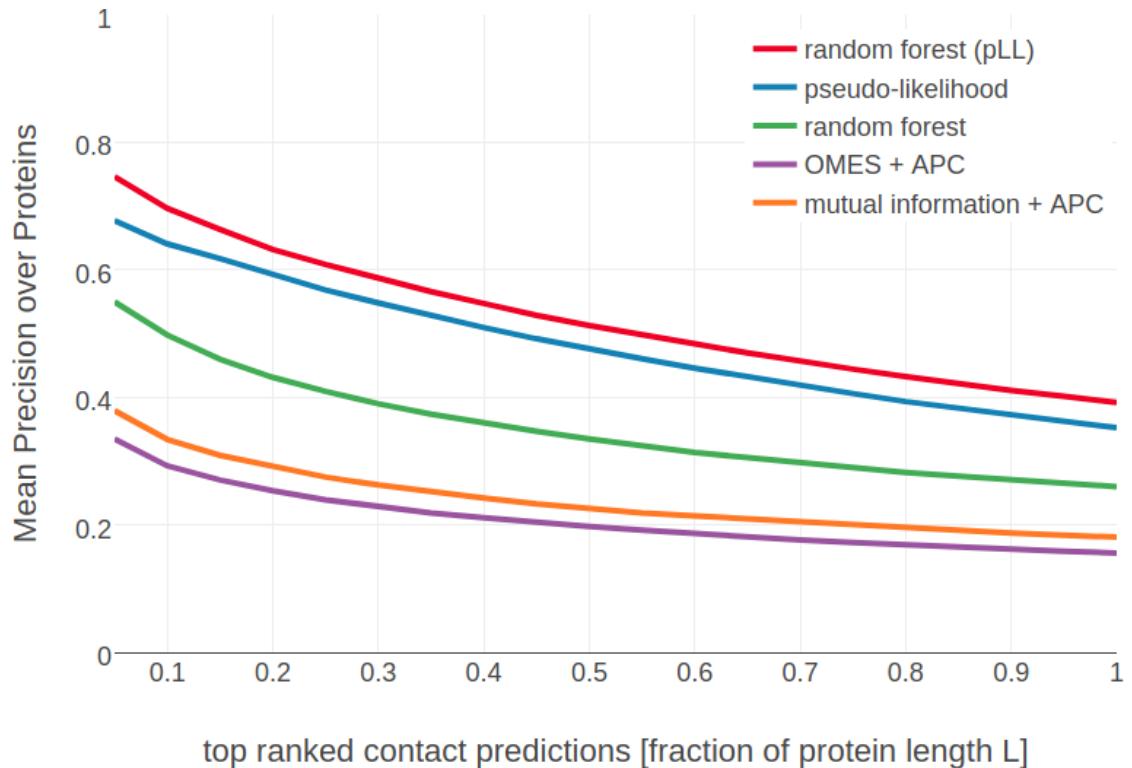


Figure 3.6: Mean precision for top ranked contacts on a test set of 774 proteins. **random forest (pLL)** = random forest model using sequence derived features and pseudo-likelihood contact score ([APC](#) corrected Frobenius norm of couplings). **pseudo-likelihood** = [APC](#) corrected Frobenius norm of couplings computed with pseudo-likelihood. **random forest** = random forest model trained on 75 sequence derived features. **OMES** = [APC](#) corrected *OMES* contact score according to Fodor&Aldrich [162]. **mutual information** = [APC](#) corrected mutual information between amino acid counts (using pseudo-counts).

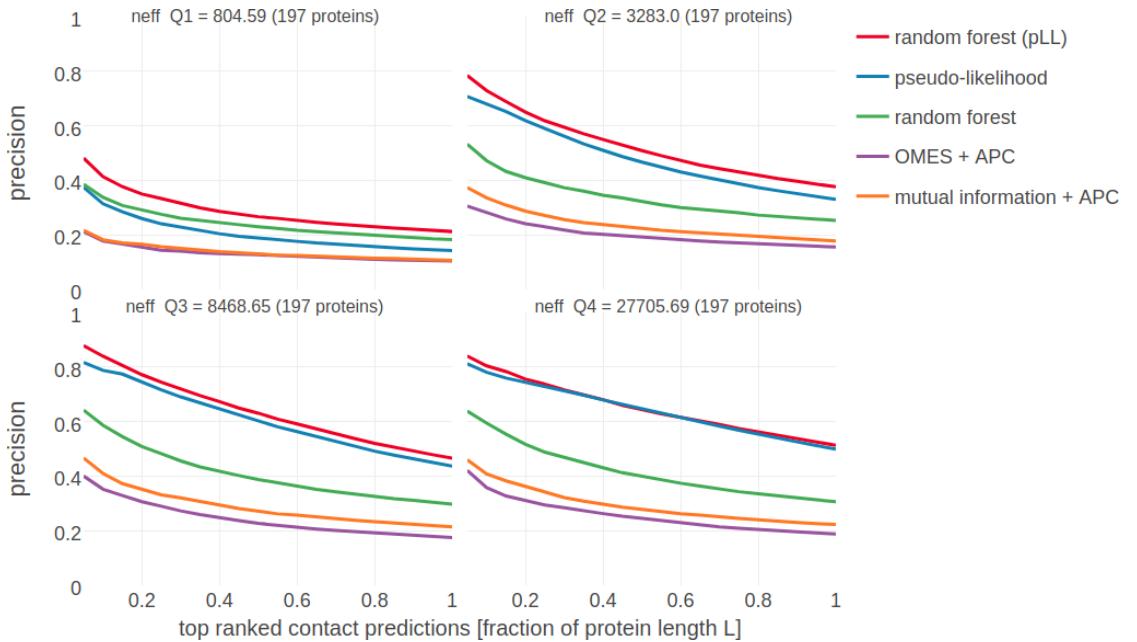


Figure 3.7: Mean precision for top ranked contacts on a test set of 774 proteins splitted into four equally sized subsets with respect to  $\text{Neff}$ . Subsets are defined according to quantiles of  $\text{Neff}$  values. Upper left: Subset of proteins with  $\text{Neff} < \text{Q1}$ . Upper right: Subset of proteins with  $\text{Q1} \leq \text{Neff} < \text{Q2}$ . Lower left: Subset of proteins with  $\text{Q2} \leq \text{Neff} < \text{Q3}$ . Lower right: Subset of proteins with  $\text{Q3} \leq \text{Neff} < \text{Q4}$ . **random forest (pLL)** = random forest model using sequence derived features and pseudo-likelihood contact score (**APC** corrected Frobenius norm of couplings). **pseudo-likelihood** = **APC** corrected Frobenius norm of couplings computed with pseudo-likelihood. **random forest** = random forest model trained on 75 sequence derived features. **OMES** = **APC** corrected *OMES* contact score according to Fodor&Aldrich [162]. **mutual information** = **APC** corrected mutual information between amino acid counts (using pseudo-counts).

Finally, comparing the random forest model trained on sequence features and pseudo-likelihood contact score to the pseudo-likelihood score in Figure 3.6 reveals that combining both types of information indeed improves predictive power over both single approaches. Especially for small alignments, the improvement is substantial as can be seen in the left upper plot in Figure 3.7. In contrast, the improvement on large alignments (right lower plot in Figure 3.7) is small, as the gain from simple sequence features compared to the much more powerful coevolution signals is neglectable.



# 4

## Methods

### 4.1 Dataset

A protein dataset has been constructed from the CATH (v4.1) [166] database for classification of protein domains. All CATH domains from classes 1(mainly  $\alpha$ ), 2(mainly  $\beta$ ), 3( $\alpha + \beta$ ) have been selected and filtered for internal redundancy at the sequence level using the `pdbfilter` script from the HH-suite[145] with an E-value cutoff=0.1. The dataset has been split into ten subsets aiming at the best possible balance between CATH classes 1,2,3 in the subsets. All domains from a given CATH topology (=fold) go into the same subsets, so that any two subsets are non-redundant at the fold level. Some overrepresented folds (e.g. Rossman Fold) have been subsampled ensuring that in every subset each class contains at max 50% domains of the same fold. Consequently, a fold is not allowed to dominate a subset or even a class in a subset. In total there are 6741 domains in the dataset.

Multiple sequence alignments were built from the CATH domain sequences (COMBS) using HHblits [145] with parameters to maximize the detection of homologous sequences:

```
hhblits -maxfilt 100000 -realign_max 100000 -B 100000 -Z 100000 -n 5  
-e 0.1 -all hhfilter -id 90 -neff 15 -qsc -30
```

The COMBS sequences are derived from the SEQRES records of the PDB file and sometimes contain extra residues that are not resolved in the structure. Therefore, residues in PDB files have been renumbered to match the COMBS sequences. The process of renumbering residues in PDB files yielded ambiguous solutions for 293 proteins, that were removed from the dataset. Another filtering step was applied to remove 80 proteins that do not hold the following properties:

- more than 10 sequences in the multiple sequence alignment ( $N > 10$ )
- protein length between 30 and 600 residues ( $30 \leq L \leq 600$ )
- less than 80% gaps in the multiple sequence alignment (percent gaps  $< 0.8$ )
- at least one residue-pair in contact at  $C_\beta < 8\text{\AA}$  and minimum sequence separation of 6 positions

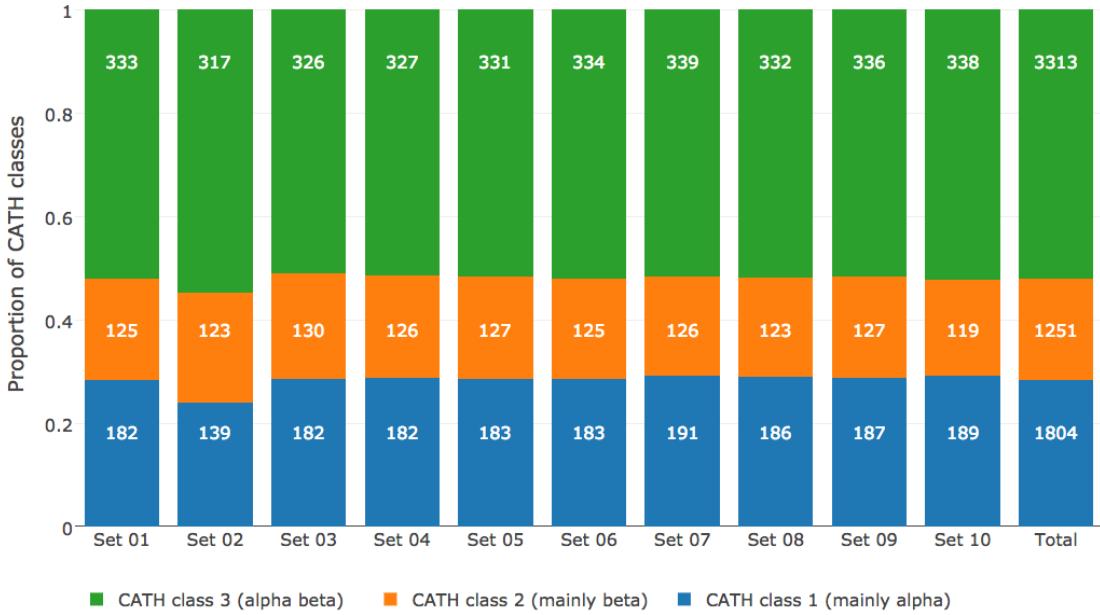


Figure 4.1: Distribution of CATH classes (1=mainly  $\alpha$ , 2=mainly  $\beta$ , 3= $\alpha - \beta$ ) in the dataset and the ten subsets.

The final dataset is comprised of **6368** proteins with almost evenly distributed CATH classes over the ten subsets (Figure 4.1).

## 4.2 Computing Pseudo-Likelihood Couplings

Dr Stefan Seemayer has reimplemented the open-source software CCMpred [57] in Python. CCMpred optimizes the regularized negative pseudo-log-likelihood using a conjugate gradients optimizer. Based on a fork of his private github repository I continued development and extended the software, which is now called CCMpredPy. It will soon be available at <https://github.com/soedinglab/CCMpredPy>. All computations in this thesis are performed with CCMpredPy unless stated otherwise.

### 4.2.1 Differences between CCMpred and CCMpredpy

CCMpredPy differs from CCMpred [57] which is available at <https://github.com/soedinglab/CCMpred> in several details:

Initialization of potentials  $\mathbf{v}$  and  $\mathbf{w}$ : - CCMpred initializes single potentials  $\mathbf{v}_i(a) = \log f_i(a) - \log f_i(a = "-")$  with  $f_i(a)$  being the frequency of amino acid a at position i and  $a = "-"$  representing a gap. A single pseudo-count has been added before computing the frequencies. Pair potentials  $\mathbf{w}$  are initialized at 0. - CCMpredPy initializes single potentials  $\mathbf{v}$  with the ML estimate of single potentials (see section 4.6.3) using amino acid frequencies computed as described in section 4.4. Pair potentials  $\mathbf{w}$  are initialized at 0.

Regularization:

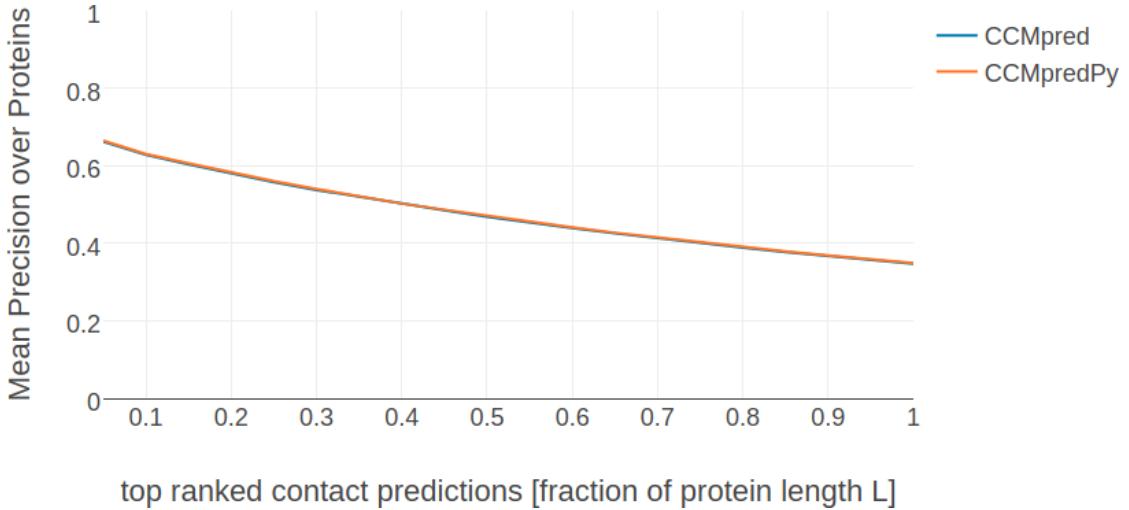


Figure 4.2: Mean precision over 3124 proteins of top ranked contacts computed as [APC](#) corrected Frobenius norm of couplings. Couplings have been computed with CCMpred [57] and CCMpredPy as specified in the legend. Specific flags that have been used to run both methods are described in detail in the text (see section [4.2.1](#)).

- CCMpred uses a Gaussian regularization prior centered at zero for both single and pair potentials. The regularization coefficient for single potentials  $\lambda_v = 0.01$  and for pair potentials  $\lambda_w = 0.2 * (L - 1)$  with  $L$  being protein length.
- CCMpredPy uses a Gaussian regularization prior centered at zero for the pair potentials. For the single potentials the Gaussian regularization prior is centered at the [ML](#) estimate of single potentials (see section [4.6.3](#)) using amino acid frequencies computed as described in section [4.4](#). The regularization coefficient for single potentials  $\lambda_v = 10$  and for pair potentials  $\lambda_w = 0.2 * (L - 1)$  with  $L$  being protein length.

Default settings for CCMpredPy have been chosen to best reproduce CCMpred results. A benchmark over a subset of approximately 3000 proteins confirms that performance measured as [PPV](#) for both methods is almost identical (see Figure [4.2](#)).

The benchmark in Figure [4.2](#) as well as all contacts predicted with CCMpred and CCMpredPy (using pseudo-likelihood) in my thesis have been computed using the following flags:

flags for computing pseudo-likelihood couplings with CCMpredPy:

```
--maxit 250          # Compute a maximum of MAXIT operations
--center-v           # Use a Gaussian prior for single potentials
                     # centered at ML estimate v*
--reg-l2-lambda-single 10    # regularization coefficient for
                           # single potentials
--reg-l2-lambda-pair-factor 0.2  # regularization coefficient for
                               # pairwise potentials computed as
```

```

# reg-12-lambda-pair-factor * (L-1)
--pc-uniform      # use uniform pseudocounts
                   # (1/21 for 20 amino acids + 1 gap state)
--pc-count 1       # defining pseudo count admixture coefficient
                   # rho = pc-count/( pc-count+ Neff)
--epsilon 1e-5     # convergence criterion for minimum decrease
                   # in the last K iterations
--ofn-pll          # using pseudo-likelihood as objective function
--alg-cg           # using conjugate gradient to optimize
                   # objective function

```

flags for computing pseudo-likelihood couplings with CCMpred:

```

-n 250    # NUMITER: Compute a maximum of NUMITER operations
-l 0.2    # LFACTOR: Set pairwise regularization coefficients
           # to LFACTOR * (L-1)
-w 0.8    # IDTHRES: Set sequence reweighting identity
           # threshold to IDTHRES
-e 1e-5   # EPSILON: Set convergence criterion for minimum
           # decrease in the last K iterations to EPSILON

```

## 4.3 Sequence Reweighting

As discussed in section 1.7, sequences in a MSA do not represent independent draws from a probabilistic model. To reduce the effects of overrepresented sequences, typically a simple weighting strategy is applied that assigns a weight to each sequence that is the inverse of the number of similar sequences according to an identity threshold [56]. It has been found that reweighting improves contact prediction performance [21,51,142] significantly but results are robust against the choice of the identity threshold in a range between 0.7 and 0.9 [51]. An identity threshold of 0.8 has been used for all analyses in this thesis.

Every sequence  $x_n$  of length  $L$  in an alignment with  $N$  sequences has an associated weight  $w_n = 1/m_n$ , where  $m_n$  represents the number of similar sequences:

$$w_n = \frac{1}{m_n}, m_n = \sum_{m=1}^N I(ID(x_n, x_m) \geq 0.8) ID(x_n, x_m) = \frac{1}{L} \sum_{i=1}^L I(x_n^i = x_m^i) \quad (4.1)$$

The number of effective sequences  $N_{\text{eff}}$  of an alignment is then the number of sequence clusters computed as:

$$N_{\text{eff}} = \sum_{n=1}^N w_n \quad (4.2)$$

TODO: Plot Performance for Seq weighting

## 4.4 Computing Amino Acid Frequencies

Single and pairwise amino acid frequencies are computed from amino acid counts of weighted sequences as described in the last section 4.3 and additional pseudocounts that are added to improve numerical stability.

Let  $a, b \in \{1, \dots, 20\}$  be amino acids and  $q_0(x_i = a), q_0(x_i = a, x_j = b)$  be the empirical single and pair frequencies without pseudocounts. The empirical single and pair frequencies with pseudocounts,  $q(x_i = a), q(x_i = a, x_j = b)$ , are defined

$$q(x_i = a) := (1 - \tau) q_0(x_i = a) + \tau \tilde{q}(x_i = a) \quad (4.3)$$

$$q(x_i = a, x_j = b) := (1 - \tau)^2 [q_0(x_i = a, x_j = b) - q_0(x_i = a)q_0(x_j = b)] + \tau \tilde{q}(x_i = a) q(x_j = b) \quad (4.4)$$

$$q(x_i = a) q(x_j = b) \quad (4.5)$$

with  $\tilde{q}(x_i = a) := f(a)$  being background amino acid frequencies and  $\tau \in [0, 1]$  is a pseudocount admixture coefficient, which is a function of the diversity of the multiple sequence alignment:

$$\tau = \frac{N_{pc}}{(N_{eff} + N_{pc})} \quad (4.6)$$

where  $N_{pc} > 0$ .

The formula for  $q(x_i = a, x_j = b)$  in the second line in eq (4.5) was chosen such that for  $\tau = 0$  we obtain  $q(x_i = a, x_j = b) = q_0(x_i = a, x_j = b)$ , and furthermore  $q(x_i = a, x_j = b) = q(x_i = a)q(x_j = b)$  exactly if  $q_0(x_i = a, x_j = b) = q_0(x_i = a)q_0(x_j = b)$ .

## 4.5 Regularization

*CCMpredPy* uses an L2-regularization per default that pushes the single and pairwise terms smoothly towards zero and is equivalent to the logarithm of a zero-centered Gaussian prior,

$$\begin{aligned} R(\mathbf{v}, \mathbf{w}) &= \log [\mathcal{N}(\mathbf{v} | \mathbf{v}^*, \lambda_v^{-1} I) \mathcal{N}(\mathbf{w} | \mathbf{w}^*, \lambda_w^{-1} I)] \\ &= -\frac{\lambda_v}{2} \|\mathbf{v} - \mathbf{v}^*\|_2^2 - \frac{\lambda_w}{2} \|\mathbf{w} - \mathbf{w}^*\|_2^2 + \text{const. ,} \end{aligned} \quad (4.7)$$

where the regularization coefficients  $\lambda_v$  and  $\lambda_w$  determine the strength of regularization.

The regularization coefficient  $\lambda_w$  for couplings  $\mathbf{w}$  is defined with respect to protein length  $L$  owing to the fact that the number of possible contacts in a protein increases quadratically with  $L$  whereas the number of observed contacts only increases linearly as can be seen in Figure 4.3.

Most previous pseudo-likelihood approaches using L2-regularization for pseudo-likelihood optimization set  $\mathbf{v}^* = \mathbf{w}^* = \mathbf{0}$  [57–59]. A different choice for  $v^*$  is discussed

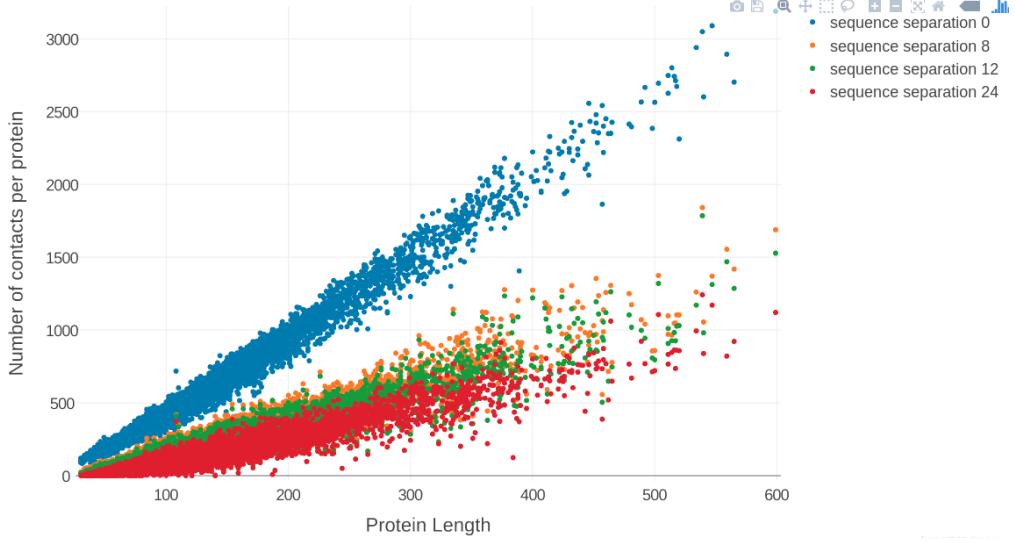


Figure 4.3: Number of contacts ( $C_\beta < 8\text{\AA}$ ) with respect to protein length and sequence separation has a linear relationship.

in section 4.6.3 that is used per default with *CCMpyp*. The single potentials will not be optimized with **CD** but will be fixed at  $v^*$  given in eq. (4.43). Furthermore, *CCMpyp* uses regularization coefficients  $\lambda_v = 10$  and  $\lambda_w = 0.2 \cdot (L - 1)$  for pseudo-likelihood optimization and the choice for  $\lambda_w$  used with **CD** is discussed in section 4.9.

## 4.6 The Potts Model

The  $N$  sequences of the **MSA**  $\mathbf{X}$  of a protein family are denoted as  $\mathbf{x}_1, \dots, \mathbf{x}_N$ . Each sequence  $\mathbf{x}_n = (\mathbf{x}_{n1}, \dots, \mathbf{x}_{nL})$  is a string of  $L$  letters from an alphabet indexed by  $\{0, \dots, 20\}$ , where 0 stands for a gap and  $\{1, \dots, 20\}$  stand for the 20 types of amino acids. The likelihood of the sequences in the **MSA** of the protein family is modelled with a *Potts Model*, as described in detail in section 1.4:

$$\begin{aligned} p(\mathbf{X}|\mathbf{v}, \mathbf{w}) &= \prod_{n=1}^N p(\mathbf{x}_n|\mathbf{v}, \mathbf{w}) \\ &= \prod_{n=1}^N \frac{1}{Z(\mathbf{v}, \mathbf{w})} \exp \left( \sum_{i=1}^L v_i(x_{ni}) \sum_{1 \leq i < j \leq L} w_{ij}(x_{ni}, x_{nj}) \right) \end{aligned} \quad (4.8)$$

The coefficients  $v_{ia}$  and  $w_{ijab}$  are referred to as single potentials and couplings, respectively that describe the tendency of an amino acid a (and b) to (co-)occur at the respective positions in the **MSA**.  $Z(\mathbf{v}, \mathbf{w})$  is the partition function that normalizes the probability distribution  $p(\mathbf{x}_n|\mathbf{v}, \mathbf{w})$ :

$$Z(\mathbf{v}, \mathbf{w}) = \sum_{y_1, \dots, y_L=1}^{20} \exp \left( \sum_{i=1}^L v_i(y_i) \sum_{1 \leq i < j \leq L} w_{ij}(y_i, y_j) \right) \quad (4.9)$$

The log likelihood is

$$LL(\mathbf{v}, \mathbf{w}) = \log p(\mathbf{X}|\mathbf{v}, \mathbf{w}) \quad (4.10)$$

$$= \sum_{n=1}^N \left[ \sum_{i=1}^L v_i(x_{ni}) \sum_{1 \leq i < j \leq L} w_{ij}(x_{ni}, x_{nj}) \right] - N \log Z(\mathbf{v}, \mathbf{w}). \quad (4.11)$$

The gradient of the log likelihood has single components

$$\frac{\partial LL(\mathbf{v}, \mathbf{w})}{\partial v_{ia}} = \sum_{n=1}^N I(x_{ni} = a) - N \frac{\partial}{\partial v_{ia}} \log Z(\mathbf{v}, \mathbf{w}) \quad (4.12)$$

$$= \sum_{n=1}^N I(x_{ni} = a) - N \sum_{y_1, \dots, y_L=1}^{20} \frac{\exp \left( \sum_{i=1}^L v_i(y_i) + \sum_{1 \leq i < j \leq L} w_{ij}(y_i, y_j) \right)}{Z(\mathbf{v}, \mathbf{w})} I(y_i = a) \quad (4.13)$$

$$= Nq(x_i = a) - Np(x_i = a | \mathbf{v}, \mathbf{w}) \quad (4.14)$$

and pair components

$$\frac{\partial LL(\mathbf{v}, \mathbf{w})}{\partial w_{ijab}} = \sum_{n=1}^N I(x_{ni} = a, x_{nj} = b) - N \frac{\partial}{\partial w_{ijab}} \log Z(\mathbf{v}, \mathbf{w}) \quad (4.15)$$

$$= \sum_{n=1}^N I(x_{ni} = a, x_{nj} = b) \quad (4.16)$$

$$- N \sum_{y_1, \dots, y_L=1}^{20} \frac{\exp \left( \sum_{i=1}^L v_i(y_i) + \sum_{1 \leq i < j \leq L} w_{ij}(y_i, y_j) \right)}{Z(\mathbf{v}, \mathbf{w})} I(y_i = a, y_j = b) \quad (4.17)$$

$$= Nq(x_i = a, x_j = b) - N \sum_{y_1, \dots, y_L=1}^{20} p(y_1, \dots, y_L | \mathbf{v}, \mathbf{w}) I(y_i = a, y_j = b) \quad (4.18)$$

$$= Nq(x_i = a, x_j = b) - Np(x_i = a, x_j = b | \mathbf{v}, \mathbf{w}) \quad (4.19)$$

#### 4.6.1 Treating Gaps as Missing Information

Treating gaps explicitly as 0'th letter of the alphabet will lead to couplings between columns that are not in physical contact. To see why, imagine a hypothetical alignment consisting of two sets of sequences as it is illustrated in Figure 4.4. The first set has sequences covering only the left half of columns in the MSA, while the second set has sequences covering only the right half of columns. The two blocks could correspond to protein domains that were aligned to a single query sequence.

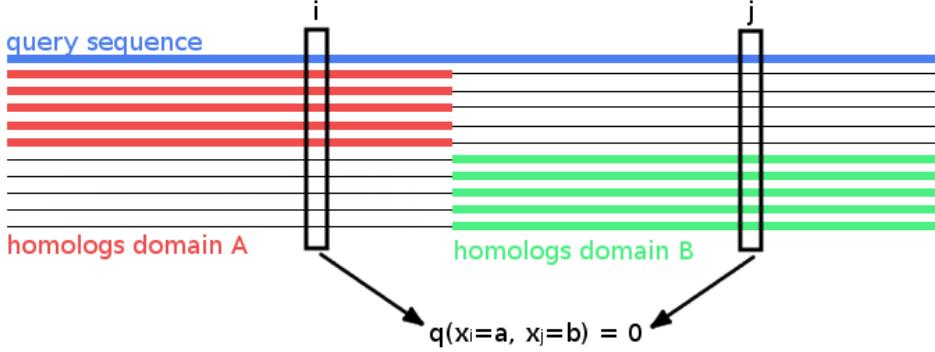


Figure 4.4: Hypothetical MSA consisting of two sets of sequences: the first set has sequences covering only the left half of columns, while the second set has sequences covering only the right half of columns. The two blocks could correspond to protein domains that were aligned to a single query sequence. Empirical amino acid pair frequencies  $q(x_i = a, x_j = b)$  will vanish for positions  $i$  from the left half and  $j$  from the right half of the alignment.

Now consider couplings between a pair of columns  $i, j$  with  $i$  from the left half and  $j$  from the right half. Since no sequence (except the single query sequence) overlaps both domains, the empirical amino acid pair frequencies  $q(x_i = a, x_j = b)$  will vanish for all  $a, b \in \{1, \dots, L\}$ .

According to the gradient of the log likelihood for couplings  $w_{ijab}$  given in eq (4.19), the empirical frequencies  $q(x_i = a, x_j = b)$  are equal to the model probabilities  $p(x_i = a, x_j = b | \mathbf{v}, \mathbf{w})$  at the maximum of the likelihood when the gradient vanishes. Therefore,  $p(x_i = a, x_j = b | \mathbf{v}, \mathbf{w})$  would have to be zero at the optimum when the empirical amino acid frequencies  $q(x_i = a, x_j = b)$  vanish for pairs of columns as described above. However,  $p(x_i = a, x_j = b | \mathbf{v}, \mathbf{w})$  can only become zero, when the exponential term is zero, which would only be possible if  $w_{ijab}$  goes to  $\infty$ . This is clearly undesirable, as physical contacts will be deduced from the size of the couplings.

The solution is to treat gaps as missing information. This means that the normalisation of  $p(\mathbf{x}_n | \mathbf{v}, \mathbf{w})$  should not run over all positions  $i \in \{1, \dots, L\}$  but only over those  $i$  that are not gaps in  $\mathbf{x}_n$ . Therefore, the set of sequences  $S_n$  used for normalization of  $p(\mathbf{x}_n | \mathbf{v}, \mathbf{w})$  in the partition function will be defined as:

$$S_n := \{(y_1, \dots, y_L) : 0 \leq y_i \leq 20 \wedge (y_i = 0 \text{ iff } x_{ni} = 0)\} \quad (4.20)$$

and the partition function becomes:

$$Z_n(\mathbf{v}, \mathbf{w}) = \sum_{\mathbf{y} \in S_n} \exp \left( \sum_{i=1}^L v_i(y_i) \sum_{1 \leq i < j \leq L} w_{ij}(y_i, y_j) \right) \quad (4.21)$$

To ensure that the gaps in  $\mathbf{y} \in S_n$  do not contribute anything to the sums, the parameters associated with a gap will be fixed to 0

$$v_i(0) = \mathbf{w}_{ij}(0, b) = \mathbf{w}_{ij}(a, 0) = 0 ,$$

for all  $i, j \in \{1, \dots, L\}$  and  $a, b \in \{0, \dots, 20\}$ .

Furthermore, the empirical amino acid frequencies  $q_{ia}$  and  $q_{ijab}$  need to be redefined such that they are normalised over  $\{1, \dots, 20\}$ ,

$$N_i := \sum_{n=1}^N w_n I(x_{ni} \neq 0) \quad q_{ia} = q(x_i = a) := \frac{1}{N_i} \sum_{n=1}^N w_n I(x_{ni} = a) \quad (4.22)$$

$$N_{ij} := \sum_{n=1}^N w_n I(x_{ni} \neq 0, x_{nj} \neq 0) \quad q_{ijab} = q(x_i = a, x_j = b) := \frac{1}{N_{ij}} \sum_{n=1}^N w_n I(x_{ni} = a, x_{nj} = b) \quad (4.23)$$

with  $w_n$  being sequence weights calculated as described in methods section 4.3. With this definition, empirical amino acid frequencies are normalized without gaps, so that

$$\sum_{a=1}^{20} q_{ia} = 1, \quad \sum_{a,b=1}^{20} q_{ijab} = 1. \quad (4.24)$$

#### 4.6.2 The Regularized Full Log Likelihood and its Gradient With Gap Treatment

In pseudo-likelihood based methods, a regularisation is commonly used that can be interpreted to arise from a prior probability. The same treatment will be applied to the full likelihood. Gaussian priors  $\mathcal{N}(\mathbf{v}|\mathbf{v}^*, \lambda_v^{-1}\mathbf{I})$  and  $\mathcal{N}(\mathbf{w}|\mathbf{0}, \lambda_w^{-1}\mathbf{I})$  will be used to constrain the parameters  $\mathbf{v}$  and  $\mathbf{w}$  and to fix the gauge. The choice of  $v^*$  is discussed in section 4.6.3. By including the logarithm of this prior into the log likelihood the regularised log likelihood is obtained,

$$LL_{\text{reg}}(\mathbf{v}, \mathbf{w}) = \log [p(\mathbf{X}|\mathbf{v}, \mathbf{w}) \mathcal{N}(\mathbf{v}|\mathbf{v}^*, \lambda_v^{-1}\mathbf{I}) \mathcal{N}(\mathbf{w}|\mathbf{0}, \lambda_w^{-1}\mathbf{I})] \quad (4.25)$$

or explicitly,

$$LL_{\text{reg}}(\mathbf{v}, \mathbf{w}) = \sum_{n=1}^N \left[ \sum_{i=1}^L v_i(x_{ni}) + \sum_{1 \leq i < j \leq L} w_{ij}(x_{ni}, x_{nj}) - \log Z_n(\mathbf{v}, \mathbf{w}) \right] \quad (4.26)$$

$$- \frac{\lambda_v}{2} \sum_{i=1}^L \sum_{a=1}^{20} (v_{ia} - v_{ia}^*)^2 - \frac{\lambda_w}{2} \sum_{1 \leq i < j \leq L} \sum_{a,b=1}^{20} w_{ijab}^2. \quad (4.27)$$

The gradient of the regularized log likelihood has single components

$$\frac{\partial LL_{\text{reg}}}{\partial v_{ia}} = \sum_{n=1}^N I(x_{ni} = a) - \sum_{n=1}^N \frac{\partial}{\partial v_{ia}} \log Z_n(\mathbf{v}, \mathbf{w}) - \lambda_v(v_{ia} - v_{ia}^*) \quad (4.28)$$

$$= N_i q(x_i = a) \quad (4.29)$$

$$- \sum_{n=1}^N \sum_{\mathbf{y} \in S_n} \frac{\exp \left( \sum_{i=1}^L v_i(y_i) + \sum_{1 \leq i < j \leq L} w_{ij}(y_i, y_j) \right)}{Z_n(\mathbf{v}, \mathbf{w})} I(y_i = a) \quad (4.30)$$

$$- \lambda_v(v_{ia} - v_{ia}^*) \quad (4.31)$$

and pair components

$$\frac{\partial LL_{\text{reg}}}{\partial w_{ijab}} = \sum_{n=1}^N I(x_{ni} = a, x_{nj} = b) - \sum_{n=1}^N \frac{\partial}{\partial w_{ijab}} \log Z_n(\mathbf{v}, \mathbf{w}) - \lambda_w w_{ijab} \quad (4.32)$$

$$= N_{ij} q(x_i = a, x_j = b) \quad (4.33)$$

$$- \sum_{n=1}^N \sum_{\mathbf{y} \in S_n} \frac{\exp \left( \sum_{i=1}^L v_i(y_i) + \sum_{1 \leq i < j \leq L} w_{ij}(y_i, y_j) \right)}{Z_n(\mathbf{v}, \mathbf{w})} I(y_i = a, y_j = b) \quad (4.34)$$

$$- \lambda_w w_{ijab} \quad (4.35)$$

Note that (without regularization  $\lambda_v = \lambda_w = 0$ ) the empirical frequencies  $q(x_i = a)$  and  $q(x_i = a, x_j = b)$  are equal to the model probabilities at the maximum of the likelihood when the gradient becomes zero.

If the proportion of gap positions in  $\mathbf{X}$  is small (e.g. < 5%, also compare percentage of gaps in dataset in Appendix Figure C.2), the sums over  $\mathbf{y} \in S_n$  in eqs. (4.31) and (4.35) can be approximated by  $p(x_i = a | \mathbf{v}, \mathbf{w})I(x_{ni} \neq 0)$  and  $p(x_i = a, x_j = b | \mathbf{v}, \mathbf{w})I(x_{ni} \neq 0, x_{nj} \neq 0)$ , respectively, and the partial derivatives become

$$\frac{\partial LL_{\text{reg}}}{\partial v_{ia}} = N_i q(x_i = a) - N_i p(x_i = a | \mathbf{v}, \mathbf{w}) - \lambda_v(v_{ia} - v_{ia}^*) \quad (4.36)$$

$$\frac{\partial LL_{\text{reg}}}{\partial w_{ijab}} = N_{ij} q(x_i = a, x_j = b) - N_{ij} p(x_i = a, x_j = b | \mathbf{v}, \mathbf{w}) - \lambda_w w_{ijab} \quad (4.37)$$

Note that the couplings between columns  $i$  and  $j$  in the hypothetical MSA presented in the last section 4.6.1 will now vanish since  $N_{ij} = 0$  and the gradient with respect to  $w_{ijab}$  is equal to  $-\lambda_w w_{ijab}$ .

### 4.6.3 The prior on $\mathbf{v}$

Most previous approaches chose a prior around the origin,  $p(\mathbf{v}) = \mathcal{N}(\mathbf{v} | \mathbf{0}, \lambda_v^{-1} \mathbf{I})$ , i.e.,  $v_{ia}^* = 0$ . It can be shown that the choice  $v_{ia}^* = 0$  leads to undesirable results.

Taking the sum over  $b = 1, \dots, 20$  at the optimum of the gradient of couplings in eq. (4.37), yields

$$0 = N_{ij} q(x_i = a, x_j \neq 0) - N_{ij} p(x_i = a | \mathbf{v}, \mathbf{w}) - \lambda_w \sum_{b=1}^{20} w_{ijab}, \quad (4.38)$$

for all  $i, j \in \{1, \dots, L\}$  and all  $a \in \{1, \dots, 20\}$ .

Note, that by taking the sum over  $a = 1, \dots, 20$  it follows that,

$$\sum_{a,b=1}^{20} w_{ijab} = 0. \quad (4.39)$$

At the optimum the gradient with respect to  $v_{ia}$  vanishes and according to eq. (4.36),  $p(x_i = a | \mathbf{v}, \mathbf{w}) = q(x_i = a) - \lambda_v(v_{ia} - v_{ia}^*)/N_i$ . This term can be substituted into equation (4.38), yielding

$$0 = N_{ij} q(x_i = a, x_j \neq 0) - N_{ij} q(x_i = a) + \frac{N_{ij}}{N_i} \lambda_v (v_{ia} - v_{ia}^*) - \lambda_w \sum_{b=1}^{20} w_{ijab}. \quad (4.40)$$

Considering a **MSA** without gaps, the terms  $N_{ij} q(x_i = a, x_j \neq 0) - N_{ij} q(x_i = a)$  cancel out, leaving

$$0 = \lambda_v (v_{ia} - v_{ia}^*) - \lambda_w \sum_{b=1}^{20} w_{ijab}. \quad (4.41)$$

Now, consider a column  $i$  that is not coupled to any other and assume that amino acid  $a$  was frequent in column  $i$  and therefore  $v_{ia}$  would be large and positive. Then according to eq. (4.41), for any other column  $j$  the 20 coefficients  $w_{ijab}$  for  $b \in \{1, \dots, 20\}$  would have to take up the bill and deviate from zero! This unwanted behaviour can be corrected by instead choosing a Gaussian prior centered around  $\mathbf{v}^*$  obeying

$$\frac{\exp(v_{ia}^*)}{\sum_{a'=1}^{20} \exp(v_{ia'}^*)} = q(x_i = a). \quad (4.42)$$

This choice ensures that if no columns are coupled, i.e.  $p(\mathbf{x} | \mathbf{v}, \mathbf{w}) = \prod_{i=1}^L p(x_i)$ ,  $\mathbf{v} = \mathbf{v}^*$  and  $\mathbf{w} = \mathbf{0}$  gives the correct probability model for the sequences in the **MSA**. Furthermore imposing the restraint  $\sum_{a=1}^{20} v_{ia} = 0$  to fix the gauge of the  $v_{ia}$  (i.e. to remove the indeterminacy), yields

$$v_{ia}^* = \log q(x_i = a) - \frac{1}{20} \sum_{a'=1}^{20} \log q(x_i = a'). \quad (4.43)$$

For this choice,  $v_{ia} - v_{ia}^*$  will be approximately zero and will certainly be much smaller than  $v_{ia}$ , hence the sum over coupling coefficients in eq. (4.41) will be close to zero, as it should be.

## 4.7 Analysis of Coupling Matrices

### 4.7.1 Correlation of Couplings with Contact Class

Approximately 100000 residue pairs have been filtered for contacts and non-contacts respectively according to the following criteria:

- sequence separation of residue pairs  $\geq 10$
- diversity ( $= \frac{\sqrt{N}}{L}$ ) of alignment  $\geq 0.3$
- number of non-gapped sequences  $\geq 1000$
- $C_\beta$  distance threshold for contact:  $< 8\text{\AA}$
- $C_\beta$  distance threshold for noncontact:  $> 25\text{\AA}$

### 4.7.2 Coupling Distribution Plots

For one-dimensional coupling distribution plots the residue pairs and respective pseudo-log-likelihood coupling values  $w_{ijab}$  have been selected as follows:

- sequence separation of residue pairs  $\geq 10$
- percentage of gaps per column  $\leq 30\%$
- evidence for a coupling  $w_{ijab}$  estimated from the alignment,  $N_{ij} \cdot q_i(a) \cdot q_j(b) \geq 100$  with:
  - $N_{ij}$ : number of sequences with no gaps at positions  $i$  or  $j$
  - $q_i(a)$ ,  $q_j(b)$ : frequencies of amino acids  $a$  and  $b$  at positions  $i$  and  $j$ , respectively (computed as described in section 4.4)

These criteria ensure that uninformative couplings are neglected, e.g. sequence neighbors albeit being contacts according to the  $C_\beta$  contact definition cannot be assumed to express biological meaningful coupling patterns, or couplings for amino acid pairings that do not have enough statistical power due to insufficient counts in the alignment.

The same criteria have been applied for selecting couplings for the two-dimensional distribution plots with the difference that evidence for a single coupling term has to be  $N_{ij} \cdot q_i(a) \cdot q_j(b) > 80$ .

## 4.8 Optimizing Contrastive Divergence with Stochastic Gradient Descent

This section describes hyperparameter tuning for the stochastic gradient descent optimization of CD.

The couplings  $w_{ijab}$  are initialized at 0 and single potentials  $v_i$  will not be optimized but rather kept fixed at their maximum-likelihood estimate  $v_i^*$  as described in methods section 4.6.3. The gradient of the full likelihood is approximated with

**CD** which involves Gibbs sampling of protein sequences according to the current model parametrization and is described in detail in methods section 4.10. Zero centered L2-regularization is used to constrain the coupling parameters  $\mathbf{w}$  using the regularization coefficient  $\lambda_w = 0.2L$  which is the default setting for optimizing the pseudo-likelihood with *CCMpredPy*. Performance will be evaluated by the mean precision of top ranked contact predictions over a benchmark set of 300 proteins, that is a subset of the data set described in methods section 4.1. Contact scores for couplings are computed as the **APC** corrected Frobenius norm as explained in section 1.4.4. Pseudo-likelihood couplings are computed with the tool *CCMpredPy* that is introduced in methods section 4.2.1 and the pseudo-likelihood contact score will serve as general reference method for tuning the hyperparameters.

#### 4.8.1 Convergence Criterion for Stochastic Gradient Descent

In theory the gradient descent algorithm has converged and the optimum of the objective function has been reached when the gradient becomes zero. In practice the gradients will never be exactly zero, especially due to the stochasticity of the gradient estimates when using stochastic gradient descent with **CD**. For this reason, it is crucial to define a suitable convergence criterion that can be tested during optimization and once the criterion is met, convergence is assumed and the algorithm is stopped. Typically, the objective function (or a related loss function) is periodically evaluated on a validation set and the optimizer is halted whenever the function value saturates or starts to increase. This technique is called early stopping and additionally prevents overfitting [167,168]. Unfortunately, we cannot compute the full likelihood function due to its complexity and need to define a different convergence criterion.

One possibility is to stop learning when the norm of the gradients is close to zero [169]. As described earlier, only a subset of alignment sequences is used to estimate the gradient with **CD** which increases stochasticity of the gradient estimate and speeds up the algorithm. As a result of stochasticity, the norm of gradients does not converge to zero but saturates at a certain offset as it is described in section 4.10. Only by using large number of sequences to estimate the gradients with **CD**, the norm over gradients will converge towards a value close to zero. However, runtime increases linearly in the number of sequences used for estimating the gradient with **CD** by Gibbs sampling. Convergence can also be monitored as the relative change of the norm of gradients between a particular number of iterations and optimization can be stopped when the norm of gradients has reached a certain plateau. As gradient estimates are generally very noisy with stochastic gradient descent, gradient fluctuations complicate the proper assessment of this criterion.

Instead of the gradients, it is also possible to observe the relative change over the norm of parameter estimates  $\|\mathbf{w}\|_2$  and stop learning when it falls below a small threshold  $\epsilon$ ,

$$\frac{\|\mathbf{w}_{t-1}\|_2 - \|\mathbf{w}_t\|_2}{\|\mathbf{w}_{t-1}\|_2} < \epsilon. \quad (4.44)$$

This measure is less noisy than subsequent gradient estimates because the magnitude of parameter updates is bounded by the learning rate.

Another idea is to monitor the direction of gradients, since the assumption goes that gradients will start oscillating when approaching the optimum. However, this theoretical assumption is complicated by the fact that gradient oscillations are also typically observed when the parameter surface contains narrow valleys or when the learning rate is too big, as it is visualized in the right plot in Figure ???. Another interfering factor is momentum, as it is used in the *ADAM* optimizer. Parameters will be updated into the direction of a smoothed historical gradient and oscillations, regardless of which origin, will be damped. It is therefore hard to define a general convergence criteria that uses gradient directions and can distinguish these different scenarios.

Of course, the simplest strategy to assume convergence is to specify a maximum number of iterations for the optimization procedure, which also ensures that the algorithm will stop eventually if none of the other convergence criteria is met.

A necessary but not sufficient criterion for convergence for the full likelihood is given by  $\sum_{a,b=1}^{20} w_{ijab} = 0$ . This requirement is derived in section 4.6.3. When using plain stochastic gradient descent without momentum and without adaptive learning rates, this criterion is never violated when parameters are initialized uniformly. This is due to the fact that the 400 gradients  $w_{ijab}$  for  $a, b \in \{1, \dots, 20\}$  are not independent because the sum over the 400 pairwise amino acid counts at positions  $i$  and  $j$  is identical for the observed and the sampled alignment and amounts to,

$$\sum_{a,b=1}^{20} N_{ij} q(x_i=a, q_j=b) = N_{ij} . \quad (4.45)$$

Considering a residue pair  $(i,j)$  and assuming amino acid pair  $(a,b)$  has higher counts in the sampled alignment than in the observed input alignment, then this difference in counts must be compensated by other amino acid pairs  $(c,d)$  having less counts in the sampled alignment compared to the true alignment (see Figure 4.5). Therefore it holds,  $\sum_{a,b=1}^{20} \nabla_{w_{ijab}} LL_{\text{reg}}(\mathbf{v}, \mathbf{w}) = 0$ . This symmetry is translated into parameter updates as long as the same learning rate is used to update all parameters. However, when using adaptive learning rates, this symmetry is broken and the condition  $\sum_{a,b=1}^{20} w_{ijab} = 0$  can be violated during the optimization process. It is therefore interesting to monitor  $\sum_{1 \leq i < j \leq L} \sum_{a,b=1}^{20} w_{ijab}$ .

I set the maximum number of iterations to 5000 and stop the optimization algorithm when the relative change over the L2-norm of parameter estimates  $\|\mathbf{w}\|_2$  falls below the threshold of  $\epsilon = 1e - 8$ .

#### 4.8.2 Tuning Hyperparameters of Stochastic Gradient Descent Optimizer

The coupling parameters  $\mathbf{w}$  will be updated at each time step  $t$  by taking a step of size  $\alpha$  along the direction of the negative gradient of the regularized full log

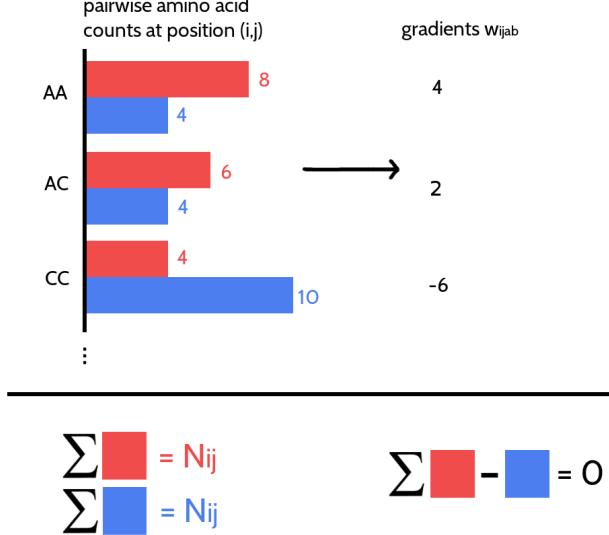


Figure 4.5: The 400 gradients  $\nabla_{w_{ijab}} LL_{\text{reg}}(\mathbf{v}, \mathbf{w})$  at position  $(i, j)$  for  $a, b \in \{1, \dots, 20\}$  are not independent. Red bars represent pairwise amino acid counts at position  $(i, j)$  for the sampled alignment. Blue bars represent pairwise amino acid counts at position  $(i, j)$  for the input alignment. The sum over pairwise amino acid counts at position  $(i, j)$  for both alignments is  $N_{ij}$ , which is the number of ungapped sequences. The gradient  $\nabla_{w_{ijab}} LL_{\text{reg}}(\mathbf{v}, \mathbf{w})$  is computed as the difference of pairwise amino acid counts for amino acids  $a$  and  $b$  at position  $(i, j)$ . The sum over gradients  $\nabla_{w_{ijab}} LL_{\text{reg}}(\mathbf{v}, \mathbf{w})$  at position  $(i, j)$  for all  $a, b \in \{1, \dots, 20\}$  is zero.

likelihood,  $-\nabla_w LL_{\text{reg}}(\mathbf{v}, \mathbf{w})$ , that has been approximated with [CD](#),

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \cdot \nabla_w LL_{\text{reg}}(\mathbf{v}, \mathbf{w}) . \quad (4.46)$$

In order to get a first intuition of the optimization problem, I tested initial learning rates  $\alpha_0 \in \{1e-4, 5e-4, 1e-3, 5e-3\}$  with a standard learning rate annealing schedule,  $\alpha = \frac{\alpha_0}{1+\gamma \cdot t}$  where  $t$  is the time step and  $\gamma$  is the decay rate that is set to 0.01 [[170](#)].

Figure 4.6 shows the mean precision for top ranked contacts computed from pseudo-likelihood couplings and from [CD](#) couplings optimized with stochastic gradient descent using the four different learning rates. Overall, mean precision for [CD](#) contacts is lower than for pseudo-likelihood contacts, especially when using the smallest ( $\alpha_0=1e-4$ ) and biggest ( $\alpha_0=5e-3$ ) learning rate.

Looking at individual proteins it turns out that the optimal learning rate depends on alignment size. The left plot in Figure 4.7 shows a convergence plot of [SGD](#) optimization using different learning rates for a protein with a small alignment. With a small initial learning rate  $\alpha_0=1e-4$  the optimization runs very slowly and does not reach convergence within 5000 iterations. Using a large initial learning rate  $\alpha_0=5e-3$  will result in slightly overshooting the optimum at the beginning of the optimization but with the learning rate decaying over time the parameter estimates converge. In contrast, for a protein with a big alignment (right plot in Figure 4.7) the choice of learning rate has a more pronounced effect. With

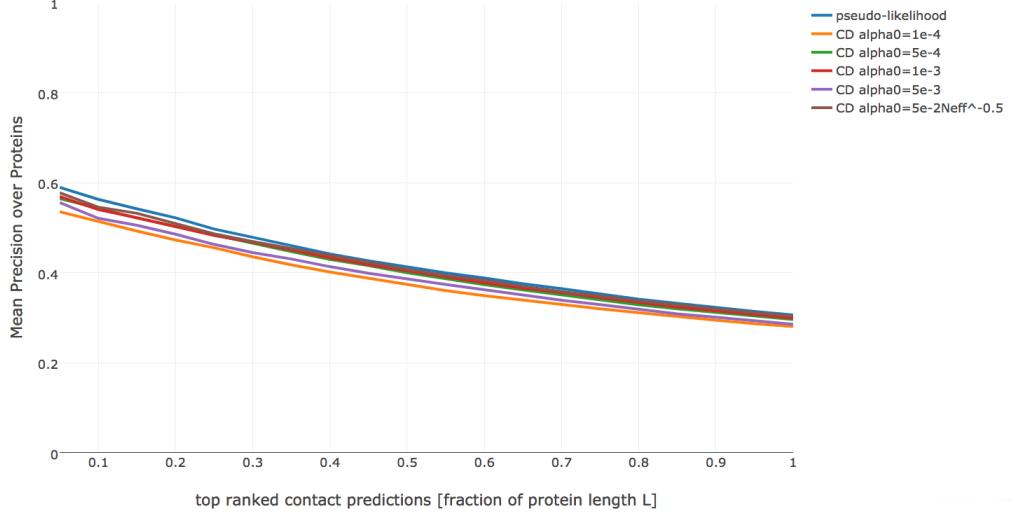


Figure 4.6: Mean precision for top ranked contact predictions over 286 proteins. Contact scores are computed as the APC corrected Frobenius norm of the couplings  $w_{ij}$ . **pseudo-likelihood**: couplings computed with pseudo-likelihood. **CD**: couplings computed with **CD** using stochastic gradient descent with different initial learning rates  $\alpha_0$  as specified in the legend.

a small initial learning rate  $\alpha_0 = 1e - 4$  the optimization runs slowly but almost converges within 5000 iterations. A large initial learning rate  $\alpha_0 = 5e - 3$  lets the parameters diverge quickly and the optimum cannot be recovered. With learning rates  $\alpha_0 = 5e - 4$  and  $\alpha_0 = 1e - 3$ , the optimum is well overshoot at the beginning of the optimization but the parameter estimates eventually converge as the learning rate decreases over time.

These observations can be explained by the fact that the magnitude of the gradient scales with the number of sequences in the alignment. The gradient is computed from amino acid counts as explained in section ???. Therefore, alignments with many sequences will generally produce larger gradients than alignments with few sequences, especially at the beginning of the optimization procedure when the difference in amino acid counts between sampled and observed sequences is largest. Following these observations, I defined the initial learning rate  $\alpha_0$  as a function of **Neff**, aiming at values for  $\alpha_0$  around 5e-3 for small **Neff** and values for  $\alpha_0$  around 1e-4 for large **Neff**,

$$\alpha_0 = \frac{5e-2}{\sqrt{N_{\text{eff}}}} . \quad (4.47)$$

For small **Neff**  $\approx 50$  this definition of the learning rate yields  $\alpha_0 \approx 7e-3$  and for big **Neff**  $\approx 20000$  this yields  $\alpha_0 \approx 3.5e-4$ . Using this learning rate defined as a function of **Neff**, precision improves over the previous fixed learning rates (see Figure 4.6). All following analyses are conducted using the **Neff**-dependent learning rate.

In a next step, I evaluated the following learning rate annealing schedules and decay rates using the **Neff**-dependent initial learning rate given in eq. (4.47):

- default linear learning rate schedule  $\alpha = \frac{\alpha_0}{1+\gamma t}$  with  $\gamma \in \{1e-3, 1e-2, 1e-1, 1\}$

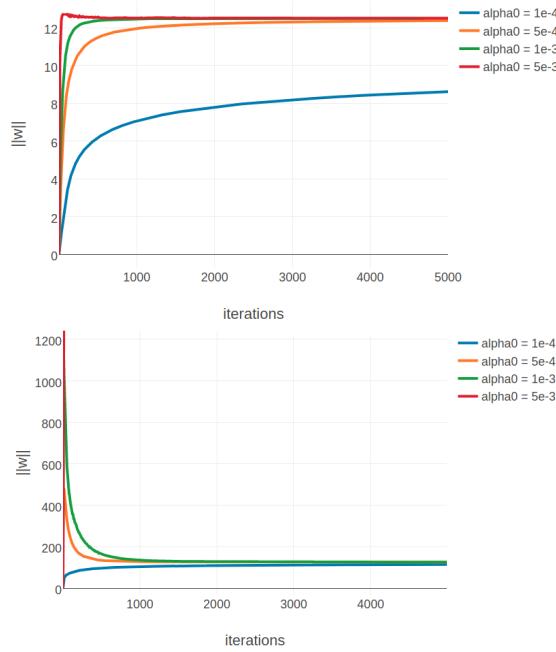


Figure 4.7: Convergence plots for two proteins during SGD optimization with different learning rates and convergence measured as L2-norm of the coupling parameters  $\|\mathbf{w}\|_2$ . Linear learning rate annealing schedule has been used with decay rate  $\gamma = 0.01$  and initial learning rates  $\alpha_0$  have been set as specified in the legend. **Left** Convergence plot for protein 1mkc\_A\_00 having protein length  $L=43$  and 142 sequences in the alignment ( $N_{eff}=96$ ). **Right** Convergence plot for protein 1c75\_A\_00 having protein length  $L=71$  and 28078 sequences in the alignment ( $N_{eff}=16808$ ). Figure is cut at the yaxis at  $\|\mathbf{w}\|_2 = 1500$ , but learning rate of  $5e-3$  reaches  $\|\mathbf{w}\|_2 \approx 13000$ .

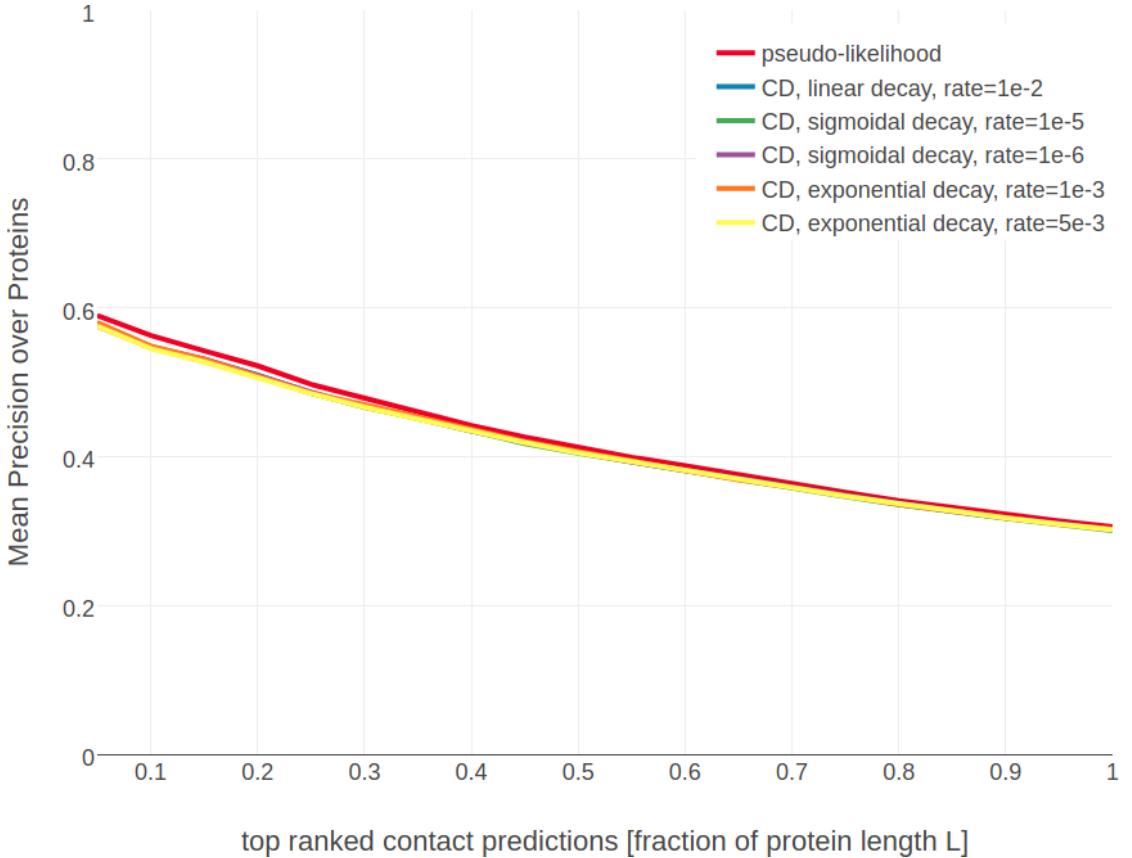


Figure 4.8: Mean precision for top ranked contact predictions over 288 proteins. Contact scores are computed as the [APC](#) corrected Frobenius norm of the couplings  $\mathbf{w}_{ij}$ . **pseudo-likelihood**: couplings computed with pseudo-likelihood. **CD**: couplings computed with [CD](#) using stochastic gradient descent with an initial learning rate defined with respect to [Neff](#). Learning rate annealing schedules and decay rates as specified in the legend.

- square root learning rate schedule  $\alpha = \frac{\alpha_0}{\sqrt{1+\gamma t}}$  with  $\gamma \in \{1e-2, 1e-1, 1\}$
- sigmoidal learning rate schedule  $\alpha_{t+1} = \frac{\alpha_t}{1+\gamma t}$  with  $\gamma \in \{1e-6, 1e-5, 1e-4, 1e-3\}$
- exponential learning rate schedule  $\alpha_{t+1} = \alpha_0 \cdot \exp(-\gamma t)$  with  $\gamma \in \{5e-4, 1e-4, 5e-3\}$

The learning rate annealing schedules are visualized for different decay rates in Appendix Figure ?? and the respective benchmark plots can be found in Appendix ???. Optimizing [CD](#) with [SGD](#) using any of the learning rate schedules listed above yields on average lower precision for the top ranked contacts than the pseudo-likelihood contact score. Several learning rate schedules perform almost equally well as can be seen in Figure 4.8. The highest precision, being one to two percentage points below the mean precision for the pseudo-likelihood contact score, is obtained with a linear learning rate schedule and decay rate  $\gamma = 1e-2$ , with a sigmoidal learning rate schedule and decay rates  $\gamma = 1e-5$  and  $\gamma = 1e-6$  and with an exponential learning rate schedule and decay rate  $\gamma = 1e-3$  and  $\gamma = 1e-5$ . The square root learning rate schedule gives ovarally bad results and does not lead to convergence because the learning rate decays slowly at later time steps.

In contrast to the findings regarding the initial learning rate earlier, an optimal

decay rate can be defined independent of the alignment size. Figure 4.9 shows convergence plots for the same two exemplary proteins as before. Proteins with low  $\text{Neff}$  are robust against the particular choice of learning rate schedule and decay rate (see left plot in Figure 4.9). The presumed optimum at  $\|w\|_2 \approx 12.5$  is almost always reached. Proteins with high  $\text{Neff}$  are stronger adversely affected by quickly decaying learning rates and the optimum at  $\|w\|_2 \approx 125$  cannot be reached before the learning rate diminishes which effectively prevents further optimization progress. Less quickly decaying learning rates, such as  $\gamma = 1e-2$  with a linear schedule or  $\gamma = 1e-6$  with a sigmoidal schedule, guide the parameter estimates close to the expected optimum and can be used with proteins having low  $\text{Neffs}$  as well as having high  $\text{Neffs}$ .

Several different learning rate annealing schedules yield almost identical mean precision for top ranked contacts, as was shown earlier (see Figure 4.8). But it can be found that they differ in convergence speed. Figure 4.10 shows the distribution over the number of iterations until convergence for SGD optimizations with five different learning rate schedules that yield similar performance. The optimization converges on average within less than 2000 iterations only when using either a sigmoidal learning rate annealing schedule with decay rate  $\gamma = 1e-5$  or an exponential learning rate annealing schedule with decay rate  $\gamma = 5e-3$ . On the contrary, the distribution of iterations until convergence has a median of 5000 when using a linear learning rate annealing schedule with  $\gamma = 1e-2$  or an exponential schedule with decay rate  $\gamma = 1e-3$ . Under these considerations, I chose a sigmoidal learning rate schedule with  $\gamma = 5e-6$  for all further analysis.

### 4.8.3 Tuning Hyperparameters of *ADAM* Optimizer

*ADAM* [171] stores an exponentially decaying average of past gradients and squared gradients,

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1)g \quad (4.48)$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2)g^2, \quad (4.49)$$

with  $g = \nabla_w LL_{\text{reg}}(\mathbf{v}, \mathbf{w})$  and the rate of decay being determined by hyperparameters  $\beta_1$  and  $\beta_2$ . Both terms  $m_t$  and  $v_t$  represent estimates of the first and second moments of the gradient, respectively. The following bias correction terms compensates for the fact that the vectors  $m_t$  and  $v_t$  are both initialized at zero and therefore are biased towards zero especially at the beginning of optimization,

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t} \quad (4.50)$$

$$\hat{v}_t = \frac{v_t}{1 - \beta_2^t}. \quad (4.51)$$

Parameters are then updated using step size  $\alpha$ , a small noise term  $\epsilon$  and the corrected moment estimates  $\hat{m}_t$ ,  $\hat{v}_t$ , according to

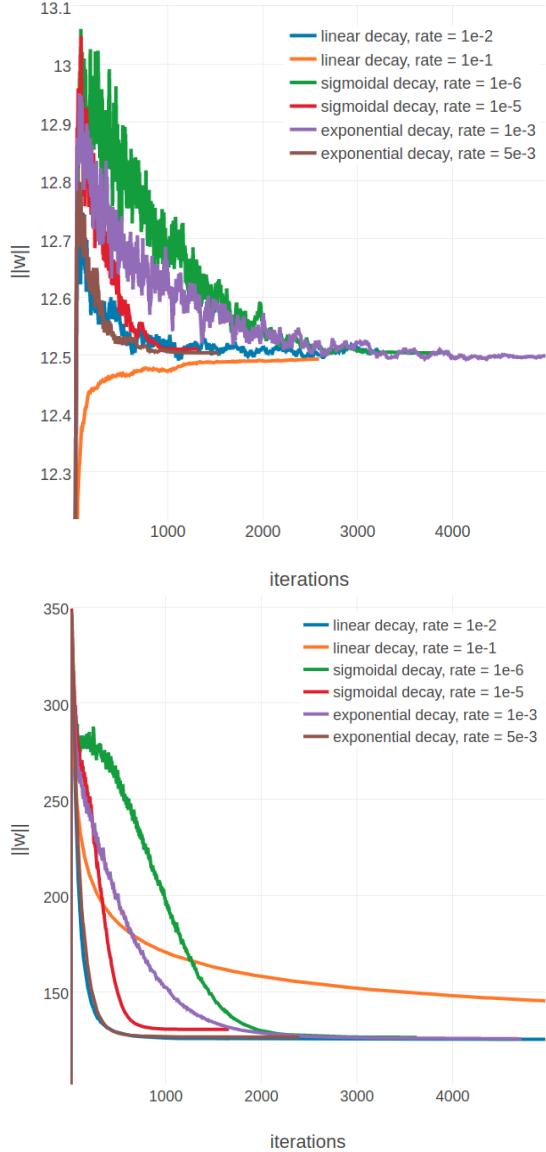


Figure 4.9: L2-norm of the coupling parameters  $\|\mathbf{w}\|_2$  during stochastic gradient descent optimization with different learning rates schedules. The initial learning rate  $\alpha_0$  is defined with respect to  $\text{Neff}$  as given in eq. (4.47). Learning rate schedules and decay rates are used according to the legend. **Left** Convergence plot for protein 1mkc\_A\_00 having protein length  $L=43$  and 142 sequences in the alignment ( $\text{Neff}=96$ ). **Right** Convergence plot for protein 1c75\_A\_00 having protein length  $L=71$  and 28078 sequences in the alignment ( $\text{Neff}=16808$ ).

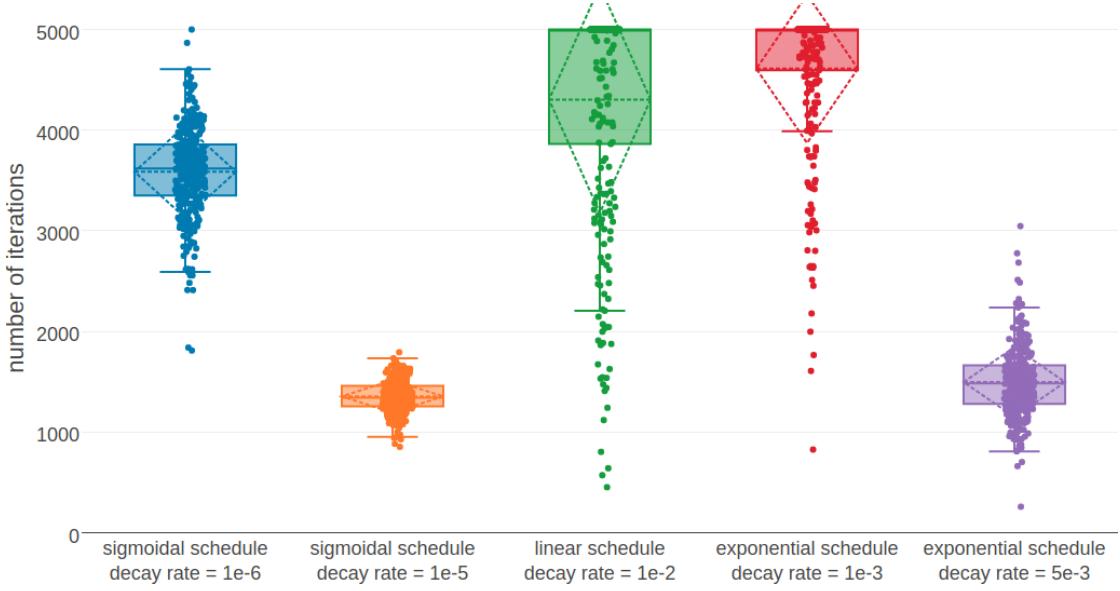


Figure 4.10: Distribution of the number of iterations until convergence for **SGD** optimizations of the full likelihood for different learning rate schedules. Convergence is reached when the relative difference of parameter norms  $\|\mathbf{w}\|_2$  falls below  $\epsilon = 1e - 8$ . Initial learning rate  $\alpha_0$  is defined with respect to **Neff** as given in eq. (4.47) and maximum number of iterations is set to 5000. Learning rate schedules and decay rates are used as specified in the legend.

$$x_{t+1} = x_t - \alpha \cdot \frac{\hat{m}_t}{\sqrt{\hat{v}_t + \epsilon}} \quad (4.52)$$

Kingma et al. proposed the default values  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$  and  $\epsilon = 1e8$  and a constant learning rate  $\alpha = 1e - 3$ .

For the two protein chains 1mkc\_A\_00 and 1c75\_A\_00, having 142 (**Neff**=96) and 28078 (**Neff**=16808) aligned sequences respectively, I analysed the convergence for **SGD** with different learning rates  $\alpha$  (see Figure 4.11). In contrast to plain stochastic gradient descent, with **ADAM** it is possible to use larger learning rates for proteins having big alignments, because the learning rate will be adapted to the magnitude of the gradient for every parameter individually. For protein 1mkc\_A\_00 having a small alignment, a learning rate of  $5e-3$  quickly leads to convergence whereas for protein 1c75\_A\_00 a larger learning rate can be chosen to obtain quick convergence. As a consequence, I defined the learning rate  $\alpha$  as a function of **Neff**,

$$\alpha = 2e-3 \log(N_{\text{eff}}) , \quad (4.53)$$

such that it will take values  $\sim 5e-3$  for proteins with small alignments and values  $\sim 1e-2$  for proteins with large alignments.

It is interesting to note in Figure 4.11, that the norm of the coupling parameters  $\|\mathbf{w}\|_2$  converges towards different values depending on the choice of the learning rate  $\alpha$ . This indicates that it is necessary to decrease the learning rate over time. By default, **ADAM** uses a constant learning rate, because the algorithm performs a kind of step size annealing by nature. However, popular implementations of

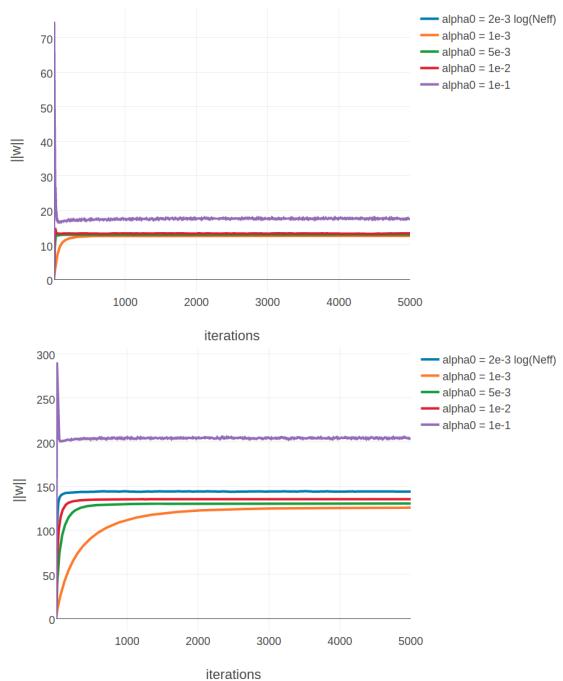


Figure 4.11: L2-norm of the coupling parameters  $\|\mathbf{w}\|_2$  during optimization with *ADAM* and different learning rates without annealing. The learning rate  $\alpha$  is specified in the legend. **Left** Convergence plot for protein 1mhc\_A\_00 having protein length  $L=43$  and 142 sequences in the alignment ( $N_{eff}=96$ ). **Right** Convergence plot for protein 1c75\_A\_00 having protein length  $L=71$  and 28078 sequences in the alignment ( $N_{eff}=16808$ ).

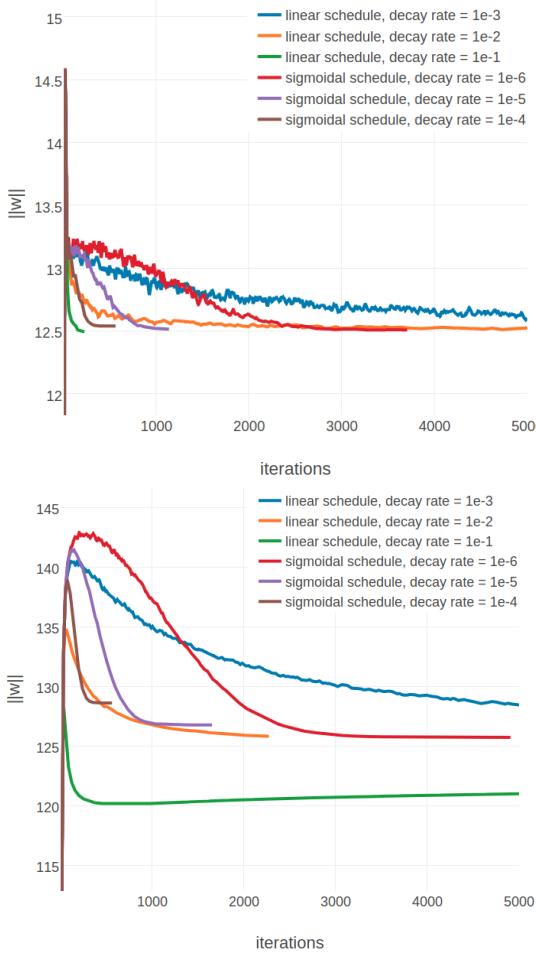


Figure 4.12: L2-norm of the coupling parameters  $\|\mathbf{w}\|_2$  during optimization with *ADAM* and different learning rate annealing schedules. The learning rate  $\alpha$  is specified with respect to  $N_{\text{eff}}$  as  $\alpha = 2e-3 \log(N_{\text{eff}})$ . The learning rate annealing schedule is specified in the legend. **Left** Convergence plot for protein 1mkc\_A\_00 having protein length  $L=43$  and 142 sequences in the alignment ( $N_{\text{eff}}=96$ ). **Right** Convergence plot for protein 1c75\_A\_00 having protein length  $L=71$  and 28078 sequences in the alignment ( $N_{\text{eff}}=16808$ ).

*ADAM* in the [Keras](#) [172] and [Lasagne](#) [173] packages allow the use of an annealing schedule. I therefore tested different learning rate annealing schedules for *ADAM* assuming that with decreasing learning rates the L2-norm of the coupling parameters  $\|\mathbf{w}\|_2$  will converge towards a consistent value. Indeed, as can be seen in Figure 4.12, when using a linear or sigmoidal learning rate annealing schedule with *ADAM*, the L2-norm of the coupling parameters  $\|\mathbf{w}\|_2$  converges roughly towards the same value that has been obtained with plain [SGD](#) shown in Figure 4.9.

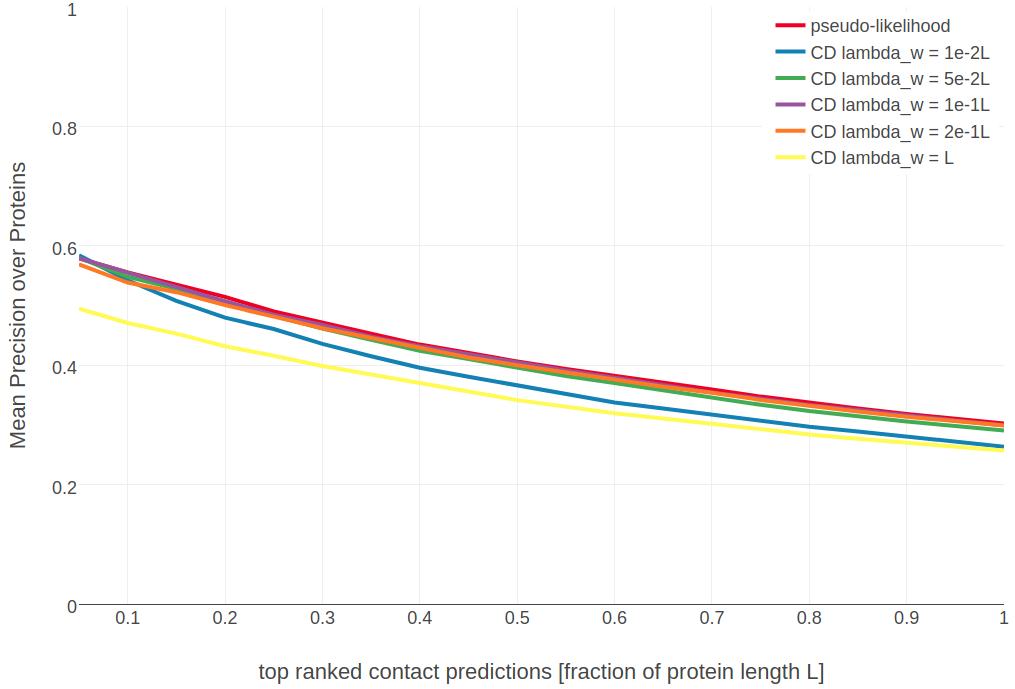


Figure 4.13: Mean precision for top ranked contact predictions over 288 proteins. Contact scores are computed as the [APC](#) corrected Frobenius norm of the couplings  $\mathbf{w}_{ij}$ . **pseudo-likelihood**: couplings computed with pseudo-likelihood. **CD lambda\_w = X**: couplings computed with [CD](#) using L2-regularization on the couplings  $\mathbf{w}$  with regularization coefficient  $\lambda_w$  chosen as specified in the legend and keeping the single potentials  $v_i$  fixed at their [MLE](#) optimum  $v_i^*$  given in eq. (4.43).

## 4.9 Tuning Regularization Coefficients for Contrastive Divergence

For tuning the hyperparameters of the stochastic gradient descent optimizer in the last section 4.8.2, the coupling parameters  $\mathbf{w}$  were constrained by a Gaussian prior  $\mathcal{N}(\mathbf{w}|0, \lambda_w^{-1} I)$  using the default pseudo-likelihood regularization coefficient  $\lambda_w = 1e-2L$  as described in methods section 4.5. It is conceivable that [CD](#) achieves optimal performance using stronger or weaker regularization than used for pseudo-likelihood optimization. Therefore, I evaluated performance of [CD](#) using the previously identified hyperparameters for [SGD](#) and different regularization coefficients  $\lambda_w \in \{1e-2L, 5e-2L, 1e-1L, 1e-2L, L\}$ . The single potentials  $\mathbf{v}$  are not subject to optimization and are kept fixed at their maximum-likelihood estimate  $v^*$  given in eq. (4.43).

As can be seen in Figure 4.13, using strong regularization for the couplings,  $\lambda_w = L$ , results in a drastic drop of mean precision. Using weaker regularization such as  $\lambda_w = 1e-2L$  or  $\lambda_w = 5e-2L$  improves precision for the top  $L/10$  and  $L/5$  predicted contacts but decreases precision when including lower ranked predictions. As a matter of fact, a slightly weaker regularization  $\lambda_w = 1e-1L$  than the default  $\lambda_w = 1e-2L$  improves mean precision especially for the top  $L/2$  contacts in such a way, that it is comparable to the pseudo-likelihood performance.

As mentioned before, a difference compared to pseudo-likelihood optimization is that the single potentials  $\mathbf{v}$  are not optimized with [CD](#) but rather set to their maximum-likelihood estimate as it is obtained in a single position model that is discussed in methods section [\(4.43\)](#). When the single potentials  $\mathbf{v}$  are optimized with [CD](#) using the same regularization coefficient  $\lambda_v = 10$  as it is used when optimizing the pseudo-likelihood, performance is almost indistinguishable compared to keeping the single potentials  $\mathbf{v}$  fixed as can be seen in appendix Figure ??.

## 4.10 Tuning the Gibbs Sampling Scheme for Contrastive Divergence

This section describes the default Gibbs sampling scheme that is used to approximate the gradients with [CD](#).

The gradient of the full log likelihood with respect to the couplings  $\mathbf{w}$  is computed as the difference of pairwise amino acid counts between the input alignment and a sampled alignment plus an additional regularization term as given in eq. (??). Pairwise amino acid counts for the input alignment are computed accounting for sequence weights (described in methods section [4.3](#)) and including pseudo counts (described in methods section [4.4](#)). Pairwise amino acid counts for the sampled alignment are computed in the same way using the same sequence weights that have been computed for the input alignment. A subset of sequences of size  $S = \min(10L, N)$ , with  $L$  being the length of sequences and  $N$  the number of sequences in the input alignment, is selected from the input alignment and used to initialize the Markov chains for the Gibbs sampling procedure. Because the input [MSA](#) is typically bigger than the sampled [MSA](#) (sampled alignment has  $\min(10L, N)$  sequences), the weighted pairwise amino acid counts of the sampled alignment need to be rescaled such that the total sample counts match the total counts from the input alignment.

The default implementation of the Gibbs sampler will sample new sequences by performing one full step of Gibbs sampling on each sequence as follows:

```
# Input: multiple sequence alignment X with N sequences of length L
# Input: model parameters v and w

N = dim(X)[0]      # number of sequences in alignment
L = dim(X)[1]      # length of sequences in alignment
S = min(10L, N)    # number of sequences that will be sampled
K = 1               # number of Gibbs steps

# randomly select S sequences from the input alignment X without replacement
sequences = random.select.rows(X, size=S, replace=False)
for seq in sequences:
    # perform K steps of Gibbs sampling
    for step in range(K):
        # iterate over permuted sequence positions i in {1, ..., L}
```

```

for i in shuffle(range(L)):
    # ignore gap positions
    if seq[i] == gap:
        continue
    # compute conditional probabilities for every amino acid a in {1, ...
    for a in range(20):
        p_cond[a] = p(seq[i]=a | (seq[1], ..., seq[i-1], seq[i+1], ...,
    # randomly select a new amino acid a in {1, ..., 20} for position i
    # according to conditional probabilities
    seq[i] = random.integer({1, ..., 20}, p_cond)

# sequences will now contain S newly sampled sequences
return sequences

```

#### 4.10.1 Varying the number of Gibbs Steps

The default [CD](#) algorithm as described by Hinton in 2002 applies only one full step of Gibbs sampling on each data sample to generate a sampled data set that will be used to approximate the gradient [174]. One full step of Gibbs sampling corresponds to sampling each position in a protein sequence according to the conditional probabilities computed from the current model probabilities as described in ???. The sampled sequences obtained after only one step of Gibbs sampling will be very similar to the input sequences. It has been shown that sampling with  $n > 1$  steps gives more precise results but increases computational cost per gradient evaluation [175,176].

In the following I analysed the impact on performance when Gibbs sampling each sequence with 1, 5 and 10 full steps. As can be seen, there is hardly an impact on precision while having much longer runtimes (by a factor of 5 and 10).

It is sometimes convenient to block some of the variables to improve mixing (Jensen, Kong, & Kjærulff, 1995; Wilkinson & Yeung, 2002). Mixing is slow if consecutive samples are highly dependent. Improve mixing by sampling sets of variables [177] Sample a percentage of all positions: - 100% - 90% - 65% - 50% - 1 = pLL

Another variant of [CD](#) that has been suggested by Tielemans in 2008 is [PCD](#)[176] that does not reset the Markov Chains at every iteration. The reason being that when using small learning rates, the model changes only slightly between iterations and the true data distribution can be better approximated. However, subsequent samples of [PCD](#) will be highly correlated creating a kind of momentum effect. Furthermore it has been found that [PCD](#) should be used with smaller learning rates and higher minibatch sizes.

As PCD might require smaller update steps and larger minibatches, I analysed the performance of PCD for the default settings of CD and additionally for smaller learning and decay rates and larger minibatches. Note that one Markov chain is kept for every sequence of the input alignment. At each iteration a subset  $N' < N$  of the Markov chains is randomly selected (without replacement) and used to for another round of Gibbs sampling at the current iteration.

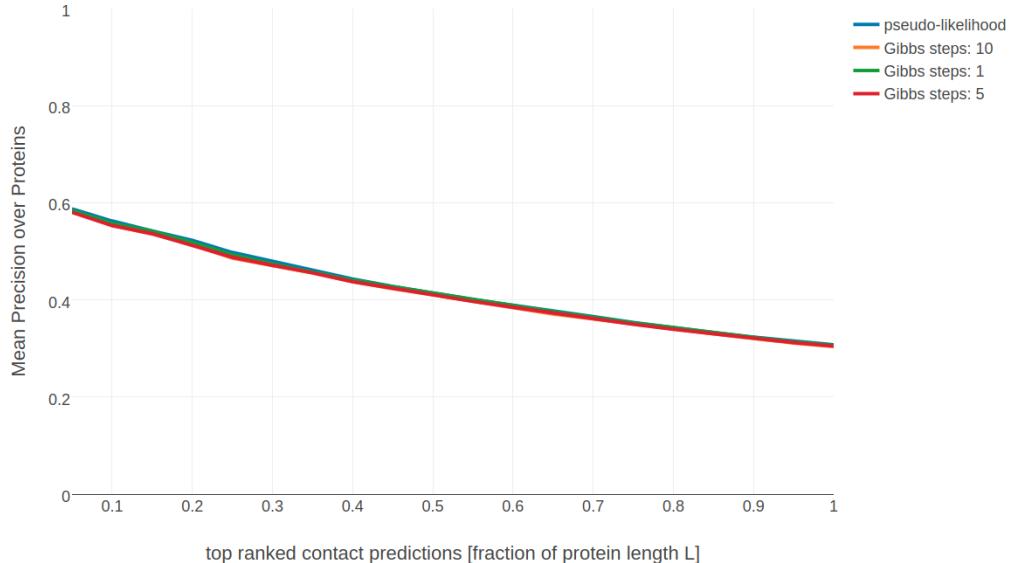


Figure 4.14: Performance of contrastive divergence optimization of the full likelihood with different number of Gibbs steps compared to pseudo-likelihood (blue) for 287 proteins. Contact scores are computed as the [APC](#) corrected Frobenius norm of the couplings  $\mathbf{w}_{ij}$ . pseudo-likelihood: contact scores computed from pseudo-likelihood. The other methods derive contact scores from couplings computed from [CD](#) with different number of Gibbs sampling steps.

#### PLOT PCD for different LEARNIGN RATEWS and SAMPLE SIZES

Discussion: - as could be seen: improved gradients and different solutions do not translate into improved precision of top ranked contacts - APC corrected l2norm might not be an appropriate measure for CD couplings: look at correlation plot of couplings (pll vs cd) and l2norm (pll vs cd) and apc l2norm (pll vs cd) -> differences vanish - ranking of residues might not be influenced by subtle changes in parameters when crude l2norm is computed -> rank plot: merged list of top ranked contacts from both methods what can we see: - generally pll has stronger scores (see also boxplot over all proteins? statistic?) - ranking is very similar, especially for top ranked contacts (thats why benchmark plots so similar)

## 4.11 Features used to train Random Forest Model

Given a multiple sequence alignment of a protein family, various sequence features can be derived that have been found to be informative of a residue-residue contact.

In total there are 250 features that can be divided into global, single position and pairwise features and are described in the following sections. If not stated otherwise, *weighted* features have been computed using amino acid counts or amino acid frequencies based on weighted sequences as described in section [4.3](#).

### 4.11.1 Global Features

These features describe alignment characteristics. Every pair of residues ( $i, j$ ) from the same protein will be attributed the same feature.

Table 4.1: Features characterizing the total alignment

| Feature                      | Description  | No. Features per residue pair ( $i, j$ ) |
|------------------------------|--|--|
| L                            | log of protein length  | 1  |
| N                            | number of sequences  | 1  |
| Neff                         | number of effective sequences computed as the sum over sequence weights (see section 4.3)  | 1  |
| gaps                         | average percentage of gaps over all positions  | 1  |
| diversity                    | $\frac{\sqrt{N}}{L}$ , N=number of sequences, L=protein length   | 1  |
| amino acid composition       | weighted amino acid frequencies in alignment   | 20                                       |
| Psipred                      | secondary structure prediction by PSIPRED (v4.0)[165] given as average three state propensities                                  | 3  |
| NetsurfP                     | secondary structure prediction by NetsurfP (v1.0)[164] given as average three state propensities                                 | 3  |
| contact prior protein length | simple contact predictor based on expected number of contacts per protein with respect to protein length (see description below) | 1  |

There are in total 32 global alignment features per residue pair.

The last feature listed in table 4.1 (“contact prior protein length”) stands for a simple contact predictor based on expected number of contacts per protein with respect to protein length. The average number of contacts per residue, computed as the observed number of contacts divided by protein length  $L$ , has a non-linear relationship with protein length  $L$  as can be seen in Figure 4.15.

In log space, the average number of contacts per residue can be fitted with a linear regression (see Figure 4.16) and yields the following functions:

- $f(L) = 1.556 + 0.596 \log(L)$  for sequence separation of 0 positions
- $f(L) = -1.273 + 0.59 \log(L)$  for sequence separation of 8 positions
- $f(L) = -1.567 + 0.615 \log(L)$  for sequence separation of 12 positions
- $f(L) = -2.0 + 0.624 \log(L)$  for sequence separation of 24 positions

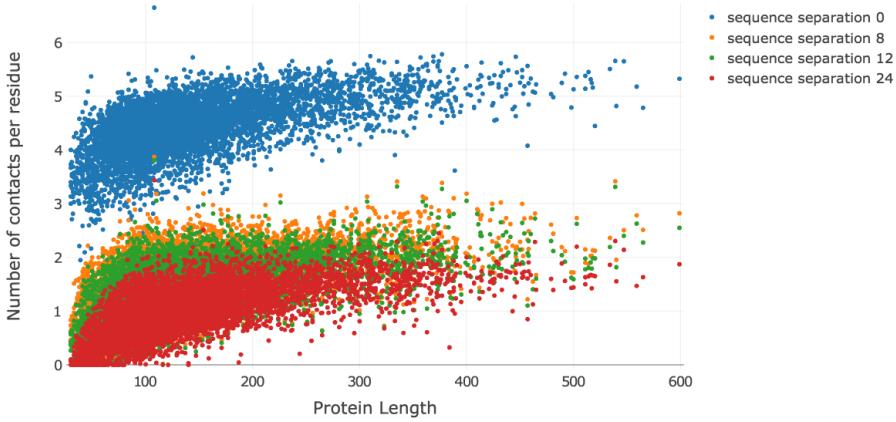


Figure 4.15: Observed number of contacts per residue has a non-linear relationship with protein length. Distribution is shown for several thresholds of sequence separation  $|j-i|$ .

A simple contact predictor can be formulated as the ratio of the expected number of contacts per residue, given by  $f(L)$ , and the possible number of contacts per residue which is  $L - 1$ ,

$$p(r_{ij} = 1 | L) = \frac{f(L)}{L - 1},$$

with  $r_{ij} = 1$  representing a contact between residue  $i$  and  $j$ .

#### 4.11.2 Single Position Features

These features describe characteristics of a single alignment column. Every residue pair  $(i, j)$  will be described by two features, once for each position.

Table 4.2: Caption hereSingle Position Sequence Features

| Feature                          | Description   | No. Features per residue pair $(i, j)$ |
|----------------------------------|---|--|
| shannon entropy (excluding gaps) | $-\sum_{a=1}^{20} p_a \log p_a$   | 2                                      |
| shannon entropy (including gaps) | $-\sum_{a=1}^{21} p_a \log p_a$   | 2                                      |
| kullback leibler divergence      | between weighted observed and background amino acid frequencies [178]           | 2                                      |
| jennson shannon divergence       | between weighted observed and background amino acid frequencies [178]           | 2                                      |
| PSSM                             | log odds ratio of weighted observed and background amino acid frequencies [178] | 40                                     |

| Feature                                | Description  | No. Features per residue pair $(i, j)$ |
|--|--|--|
| secondary structure prediction         | three state propensities PSIPRED (v4.0) [165]  | 6                                      |
| secondary structure prediction         | three state propensities Netsurfp (v1.0) [164]   | 6                                      |
| solvent accessibility prediction       | RSA and RSA Z-score Netsurfp (v1.0) [164]  | 4                                      |
| relative position in sequence          | $\frac{i}{L}$ for a protein of length $L$  | 2                                      |
| number of ungapped sequences           | $\sum_n w_n I(x_{ni} \neq 20)$ for sequences $x_n$ and sequence weights $w_n$                          | 2                                      |
| percentage of gaps                     | $\frac{\sum_n w_n I(x_{ni}=20)}{N_{\text{eff}}}$ for sequences $x_n$ and sequence weights $w_n$        | 2                                      |
| Average Atchley Factor                 | Atchley Factors 1-5 [179]  | 10                                     |
| Average polarity (Grantham)            | Polarity according to Grantham [180]. Data taken from AAindex Database [181].                          | 2                                      |
| Average polarity (Zimmermann)          | Polarity according to Zimmermann et al. [182]. Data taken from AAindex Database [181].                 | 2                                      |
| Average isoelectricity                 | Isoelectric point according to Zimmermann et al. [182]. Data taken from AAindex Database [181].        | 2                                      |
| Average hydrophobicity (Wimley&White)  | Hydrophobicity scale according to Wimley & White [183]. Data taken from UCSF Chimera [183].            | 2                                      |
| Average hydrophobicity (Kyte&Dolittle) | Hydrophobicity index according to Kyte & Doolittle [184]. Data taken from AAindex Database [181].      | 2                                      |
| Average hydrophobicity (Cornette)      | Hydrophobicity according to Cornette [185].  | 2                                      |
| Average bulkiness                      | Bulkiness according to Zimmerman et al. [182]. Data taken from AAindex Database [181].                 | 2                                      |
| Average volume                         | Average volumes of residues according to Pontius et al. [186]. Data taken from AAindex Database [181]. | 2                                      |

There are 48 single sequence features per residue and consequently 96 single se-

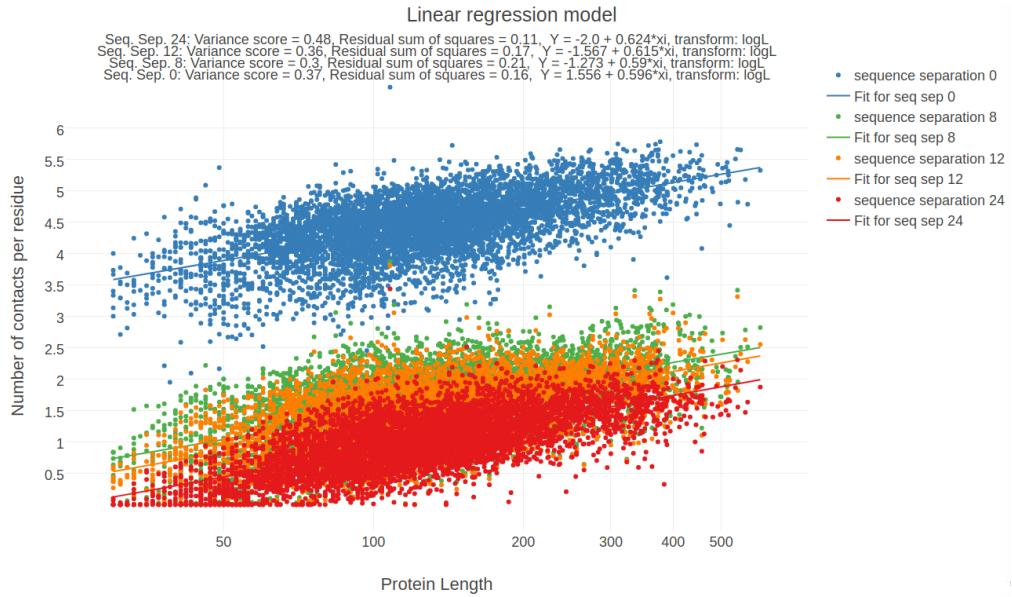


Figure 4.16: Distribution of average number of contacts per residue against protein length and corresponding linear regression fits. Protein length is on logarithmic scale. Distribution and linear regression fits are shown for different sequence separation thresholds  $|j-i|$ .

quence features per residue pair.

Additionally, all single features will be computed within a window of size 5. The window feature for center residue  $i$  will be computed as the mean feature over residues  $[i-2, \dots, i, \dots, i+2]$ . Whenever the window extends the range of the sequence (for  $i < 2$  and  $i > (L-2)$ ), the window feature will be computed only for valid sequence positions. This results in additional 96 window features per residue pair.

### 4.11.3 Pairwise Features

These features are computed for every pair of columns  $(i, j)$  in the alignment with  $i < j$ .

Table 4.3: Pairwise Sequence Features

| Feature                      | Description   | No. Features<br>per residue<br>pair $(i, j)$ |
|------------------------------|---|--|
| sequence separation          | $j - i$   | 1  |
| gaps                         | pairwise percentage of gaps using weighted sequences  | 1  |
| number of ungapped sequences | $\sum_n w_n I(x_{ni} \neq 20, x_{nj} \neq 20)$ for sequences $x_n$ and sequence weights $w_n$ | 1  |

| Feature                                | Description  | No. Features per residue pair $(i, j)$ |
|--|--|--|
| correlation physico-chemical features  | pairwise correlation of all physico-chemical properties listed in table 4.2  | 13                                     |
| pairwise potential (buried)            | Average quasi-chemical energy of interactions in an average buried environment according to Miyazawa&Jernigan [163]. Data taken from AAindex Database [181].                 | 1                                      |
| pairwise potential (water)             | Average quasi-chemical energy of transfer of amino acids from water to the protein environment according to Miyazawa&Jernigan [163]. Data taken from AAindex Database [181]. | 1                                      |
| pairwise potential (Li&Fang)           | Average general contact potential by Li&Fang [25]  | 1                                      |
| pairwise potential (Zhu&Braun)         | Average statistical potential from residue pairs in beta-sheets by Zhu&Braun [187]   | 1                                      |
| joint shannon entropy (excluding gaps) | $-\sum_{a=1}^{20} \sum_{b=1}^{20} p(a, b) \log p(a, b)$  | 1                                      |
| joint shannon entropy (including gaps) | $-\sum_{a=1}^{21} \sum_{b=1}^{21} p(a, b) \log p(a, b)$  | 1                                      |
| normalized MI                          | normalized mutual information of amino acid counts at two positions  | 1                                      |
| MI (+pseudo-counts)                    | mutual information of amino acid counts at two positions, including uniform pseudo-counts  | 1                                      |
| MI (+pseudo-counts + APC)              | mutual information of amino acid counts at two positions; including pseudo-counts and average product correction   | 1                                      |
| OMES coevolution score                 | according to Fodor&Aldrich [162] with and without APC  | 2                                      |

There are in total 26 pairwise sequence features.

## 4.12 Training Random Forest Contact Prior

Proteins constitute highly imbalanced datasets with respect to the number of residue pairs that form and do not form physical contacts. As can be seen in

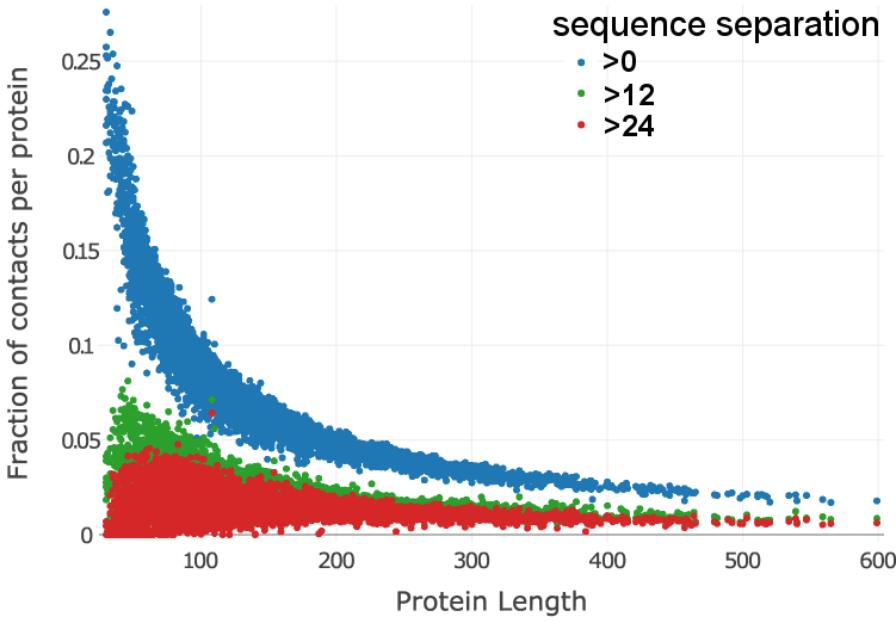


Figure 4.17: Fraction of contacts among all possible contacts in a protein against protein length L. The distribution has a non-linear relationship. At a sequence separation  $>8$  positions the fraction of contacts for intermediate size proteins with length  $>100$  is approximately 2%. Data set contains 6368 proteins and is explained in methods section 4.1.

Figure 4.17, depending on the enforced sequence separation threshold and protein length the percentage of contacts per protein varies between 25% and 0%. Most studies applying machine learning algorithms for predicting residue-residue contacts rebalanced the data set by undersampling of the majority class. Table 4.4 lists choices for the proportion of contacts to non-contacts used to train some machine learning contact predictors. I followed the same strategy and undersampled residue pairs that are not physical contacts with a proportion of contacts to non-contacts of 1:5.

Table 4.4: Important machine learning contact prediction approaches and their choices for rebalancing the data set.

| Study                     | Machine Learning Algorithm | Proportion of Contacts : Non-contacts         |
|---------------------------|----------------------------|---|
| Wu et al. (2008) [24]     | SVM                        | 1:4   |
| Li et al. (2011) [25]     | Random Forest              | 1:1, 1:2                                      |
| Wang et al. (2011) [26]   | Random Forest              | 1:4   |
| DiLena et al. (2012) [34] | deep neural network        | 1: $\approx$ 4 (sampling 20% of non-contacts) |
| Wang et al. (2013) [27]   | Random Forest              | 1: $\approx$ 4 (sampling 20% of non-contacts) |

The total training set is comprised of 50,000 residue pairs  $< 8\text{\AA}$  (“contacts”)

and 250,000 residue pairs  $> 8\text{\AA}$  (“non-contacts”). I filtered residue pairs using a sequence separation of 12 positions and selected at maximum 100 contacts and 500 non-contacts per protein. The data is collected in equal parts from data subsets 1-5 (see methods section 4), so that the training set consists of five subsets that are non-redundant at the fold level. Each of the five models for cross-validation will be trained on 40,000 contacts and 200,000 non-contacts originating from four of the five subsets. As the training set has been undersampled for non-contacts, it is not representative of real world proteins and the models need to be validated on a more realistic validation set. Therefore, each of the five trained models is not validated on the hold-out set but on separate validation sets containing 40 proteins at a time. The proteins of the validation sets are randomly selected from the respective fifth data subset and consequently are non-redundant at the fold level with training data. Performance is assessed by means of the standard contact prediction benchmark (mean precision against top ranked contacts).

I used the module `RandomForestClassifier` in the Python package `sklearn` (v. 0.19) [188] and trained the models on features extracted from `MSAs` which are listed in methods section 4.11.

### 4.12.1 Feature Selection

A random forest model is trained on the total set of features. Given the distribution of *Gini importance* values of features from the model, subsets of features are defined by features having *Gini importance* values larger than the {10, 30, 50, 70, 90}-percentile of the distribution. Performance of the models trained on these subsets of features is evaluated on the same validation set.

# A

## Abbreviations

**APC** Avarage Product Correction

**CASP** Critical Assessment of protein Structure Prediction

**CD** Contrastive Divergence

**DCA** Direct Coupling Analysis

**DI** Direct Information

**EM** electron microscopy

**IDP** intrinsically disordered proteins

**MAP** Maximum a posteriori

**MCMC** Markov Chain Monte Carlo

**MI** mutual information

**ML** Maximum-Likelihood

**MLE** Maximum-Likelihood Estimate

**MRF** Markov-Random Field

**MSA** Multiple Sequence Alignment

**Neff** Number of effective sequences

**PCD** Persistent Contrastive Divergence

**PDB** protein data bank

**SGD** stochastic gradient descent



# B

## Amino Acid Alphabet

Table B.1: Amino acid abbreviations and physico-chemical properties according to Livingstone et al., 1993 [189]

| One letter Code | Three letter Code | Amino Acid    | Physico-chemical properties                        |
|-----------------|-------------------|---------------|--|
| A               | Ala               | Alanine       | tiny, hydrophobic                                  |
| C               | Cys               | Cysteine      | small, hydrophobic, polar ( $C_{S-H}$ )            |
| D               | Asp               | Aspartic Acid | small, negatively charged, polar                   |
| E               | Glu               | Glutamic Acid | negatively charged, polar                          |
| F               | Phe               | Phenylalanine | aromatic, hydrophobic                              |
| G               | Gly               | Glycine       | tiny, hydrophobic                                  |
| H               | His               | Histidine     | hydrophobic, aromatic, polar, (positively charged) |
| I               | Ile               | Isoleucine    | aliphatic, hydrophobic                             |
| K               | Lys               | Lysine        | positively charged, polar                          |
| L               | Leu               | Leucine       | aliphatic, hydrophobic                             |
| M               | Met               | Methionine    | hydrophobic  |
| N               | Asn               | Asparagine    | small, polar                                       |
| P               | Pro               | Proline       | small  |
| Q               | Gln               | Glutamine     | tiny, hydrophobic                                  |
| R               | Arg               | Arginine      | positively charged, polar                          |
| S               | Ser               | Serine        | tiny, polar  |
| T               | Thr               | Threonine     | hydrophobic, polar                                 |
| V               | Val               | Valine        | small, aliphatic                                   |
| W               | Trp               | Tryptophan    | aromatic, hydrophobic, polar                       |
| Y               | Tyr               | TYrosine      | aromatic, hydrophobic, polar                       |

# C

## Dataset Properties

The following figures display various statistics about the dataset used throughout this thesis. See section [4.1](#) for information on how this dataset has been generated.

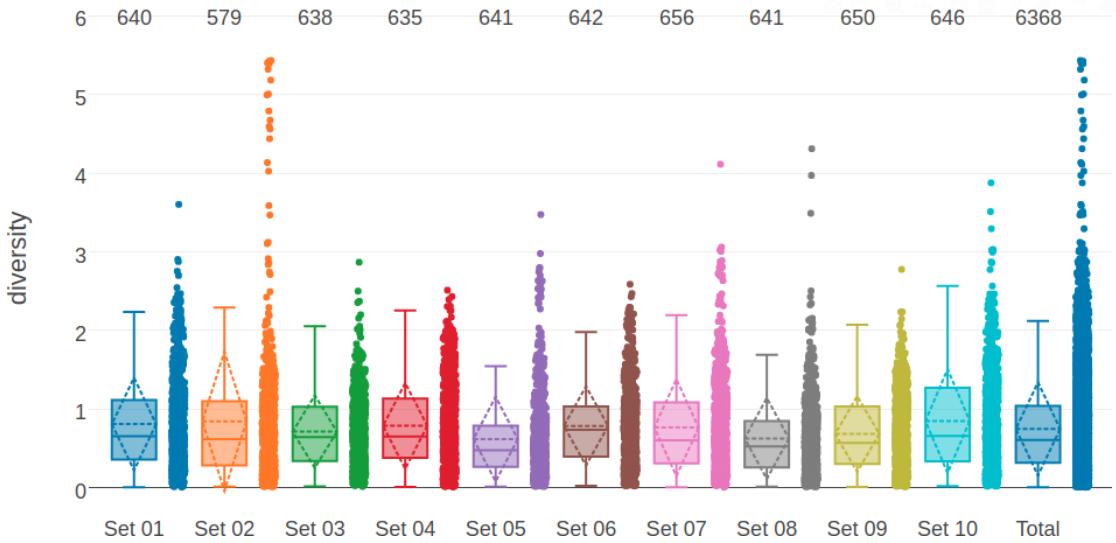


Figure C.1: Distribution of alignment diversity ( $= \sqrt{\frac{N}{L}}$ ) in the dataset and its ten subsets.

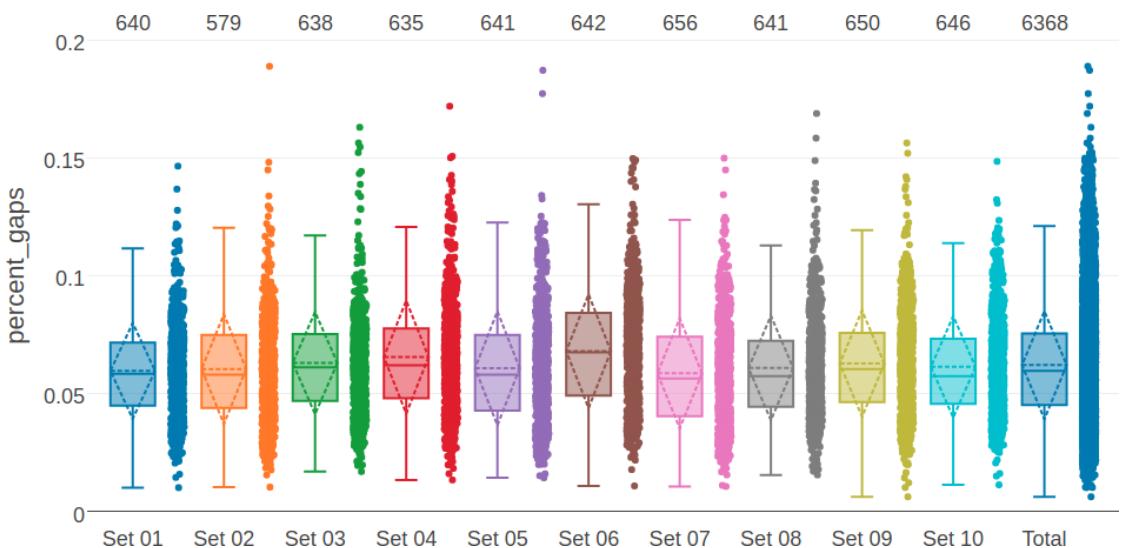


Figure C.2: Distribution of gap percentage of alignments in the dataset and its ten subsets.

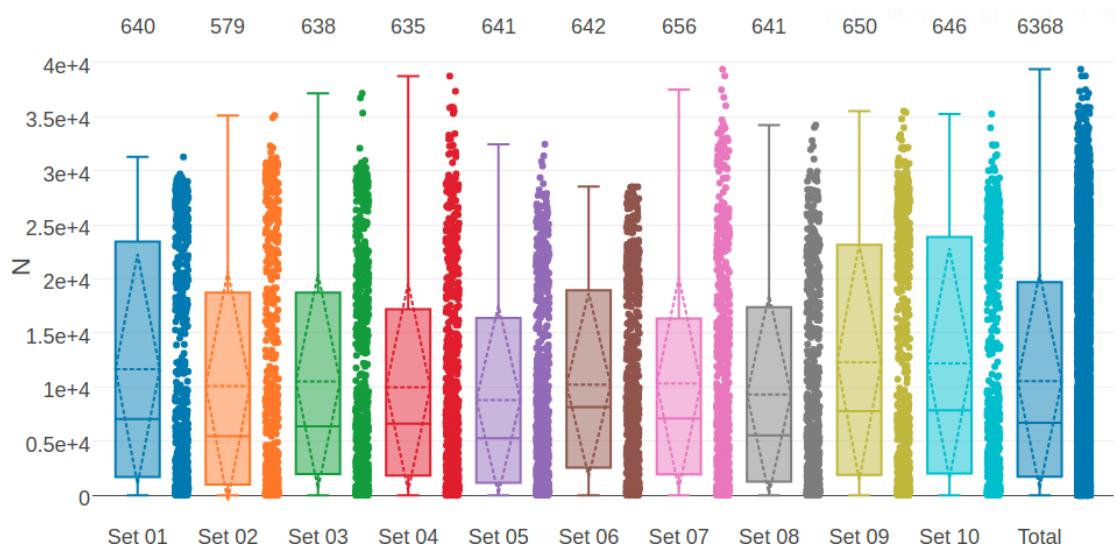


Figure C.3: Distribution of alignment size (number of sequences  $N$ ) in the dataset and its ten subsets.

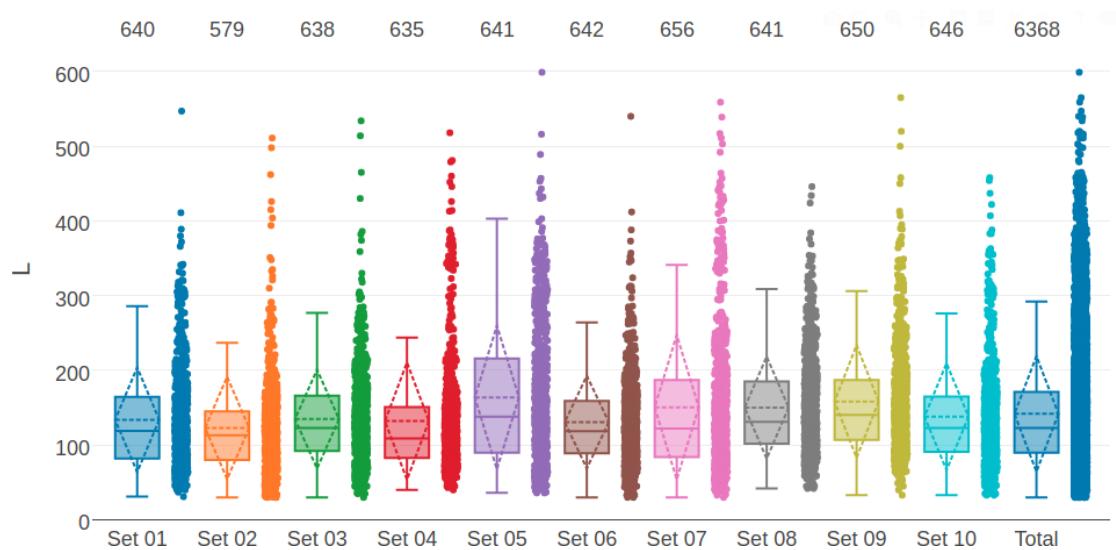


Figure C.4: Distribution of protein length  $L$  in the dataset and its ten subsets.



# D

## Standard Deviation of Couplings for Noncontacts

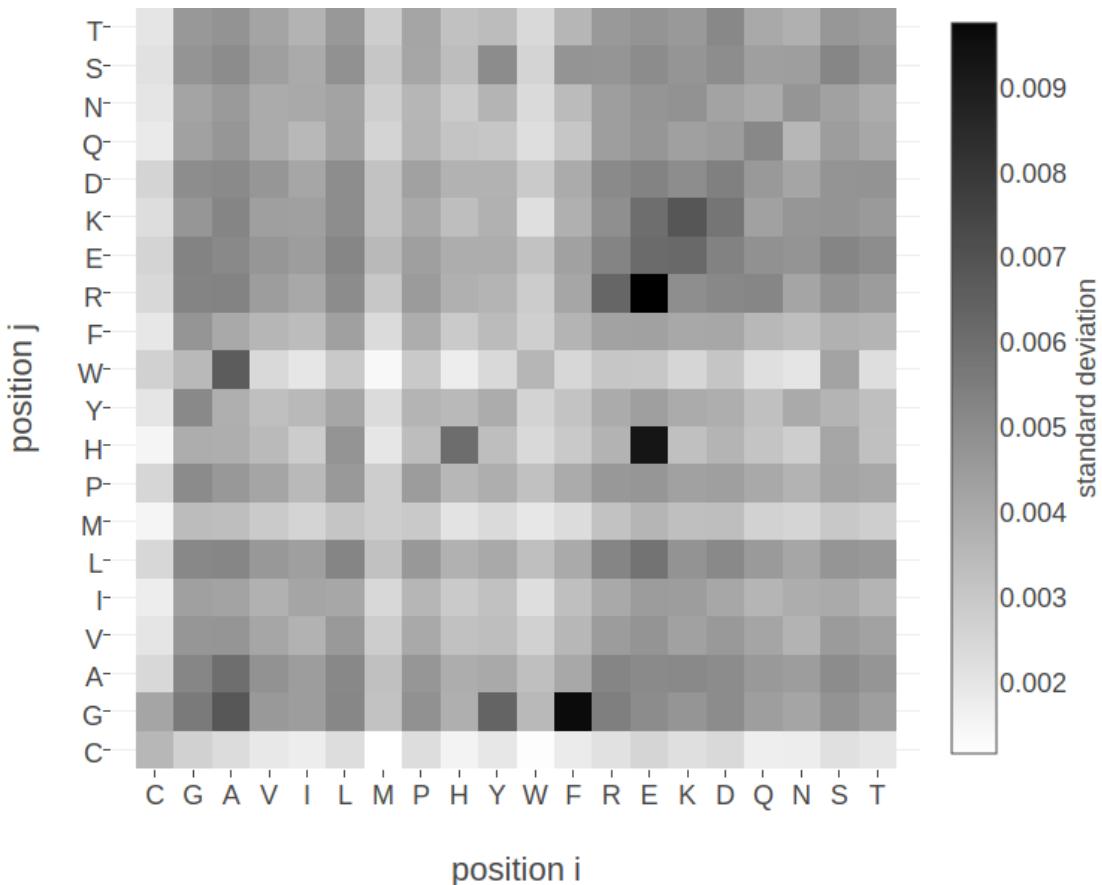


Figure D.1: Standard deviation of squared coupling values  $w_{ijab}^2$  for residue pairs not in physical contact ( $\Delta C_\beta > 25\text{\AA}$ ). Dataset contains 100.000 residue pairs per class (for details see methods section 4.7.1). Amino acids are abbreviated with one-letter code and they are broadly grouped with respect to physico-chemical properties listed in Appendix B

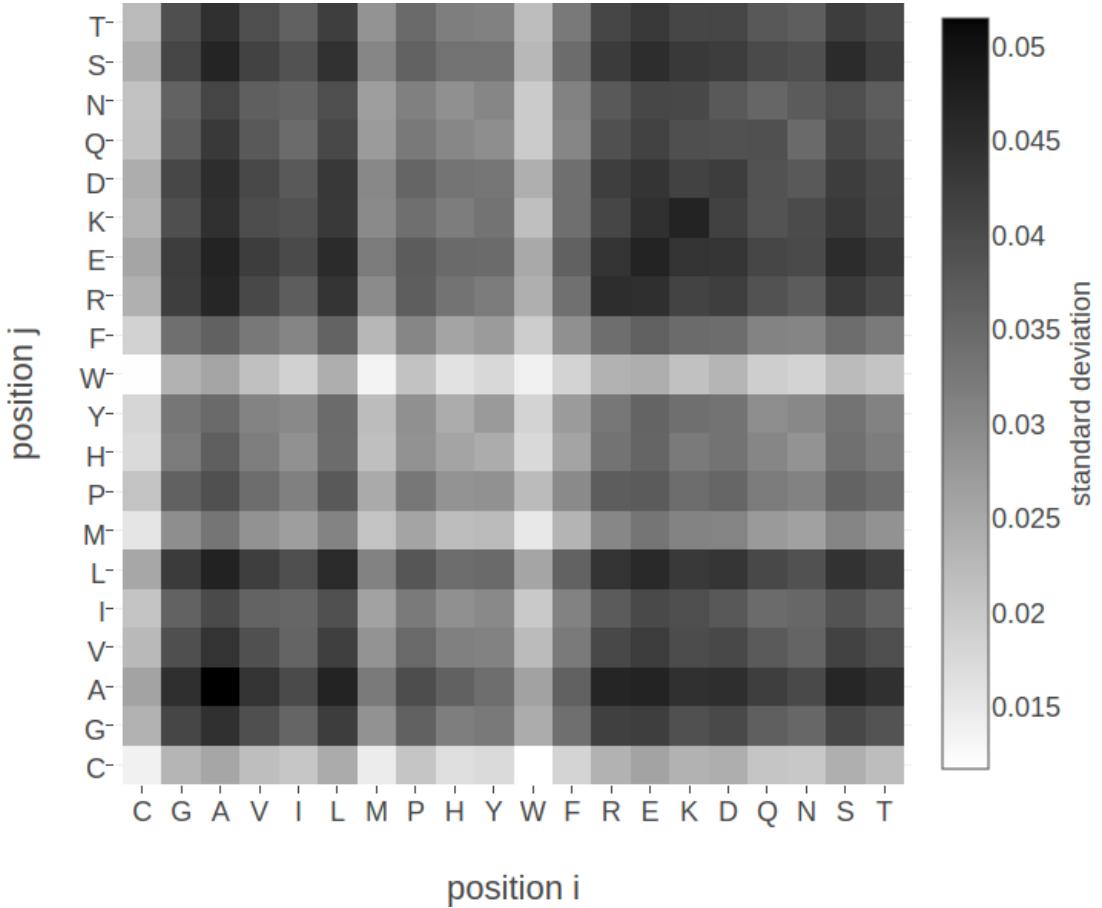


Figure D.2: Standard deviation of coupling values  $w_{ijab}$  for residue pairs not in physical contact ( $\Delta C_\beta > 25\text{\AA}$ ). Dataset contains 100.000 residue pairs per class (for details see section 4.7.1). Amino acids are abbreviated with one-letter code and they are broadly grouped with respect to physico-chemical properties listed in Appendix B.



# E

## Amino Acid Interaction Preferences Reflected in Coupling Matrices

### E.1 Pi-Cation interactions

Figure E.1 shows a Tyrosine and a Lysine residue forming a cation- $\pi$  interaction in protein 2ayd. The corresponding coupling matrix in figure ?? reflects the strong interaction preference.

### E.2 Disulfide Bonds

Figure E.2 shows two cysteine residues forming a covalent disulfide bond in protein 1alu. The corresponding coupling matrix in figure ?? reflects the strong interaction preference of cysteines.

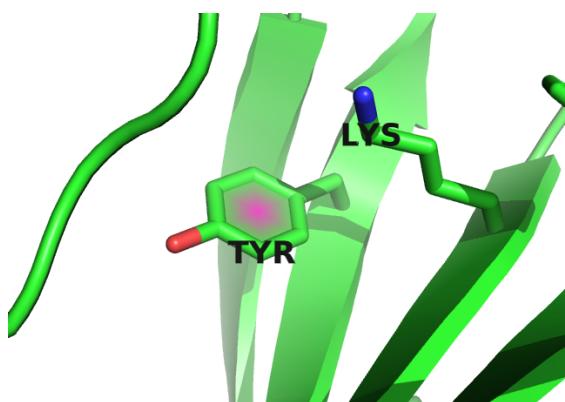


Figure E.1: Tyrosine (residue 37) and Lysine (residue 48) forming a cation- $\pi$  interaction in protein 2ayd.

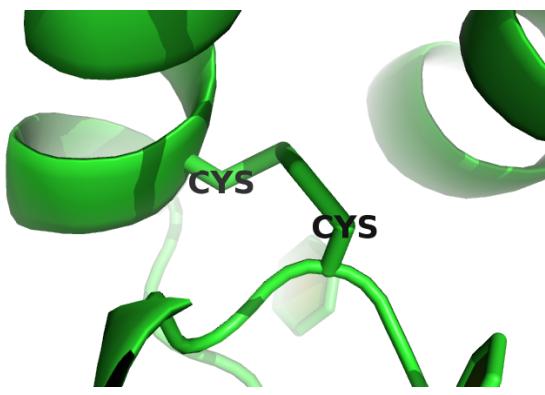


Figure E.2: Two cystein residues (residues 54 and 64) forming a covalent disulfide bond in protein 1alu.

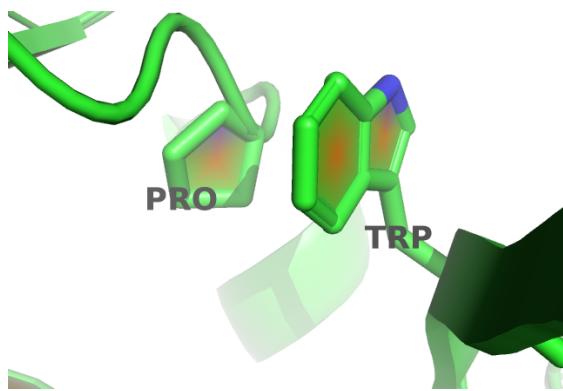


Figure E.3: Proline and tryptophan (residues 17 and 34) stacked on top of each otherengaging in a CH/π interaction in protein chain 1aol\_A\_00.

### E.3 Aromatic-Proline Interactions

Figure E.3 shows a proline and a tryptophan residue forming such a CH/π interaction in protein 1aol. The corresponding coupling matrix in figure ?? reflects this interaction with strong positive coupling between proline and tryptophan.

### E.4 Network-like structure of aromatic residues

### E.5 Aromatic Sidechains at small $C_b-C_\beta$ distances

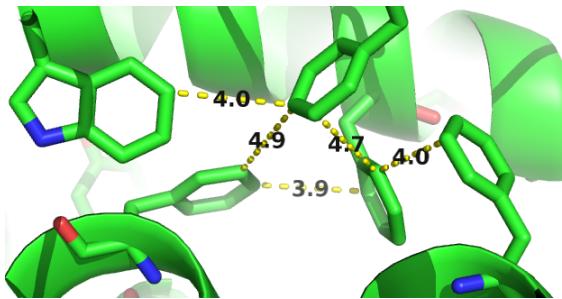


Figure E.4: Network-like structure of aromatic residues in the protein core. 80% of aromatic residues are involved in such networks that are important for protein stability [153].

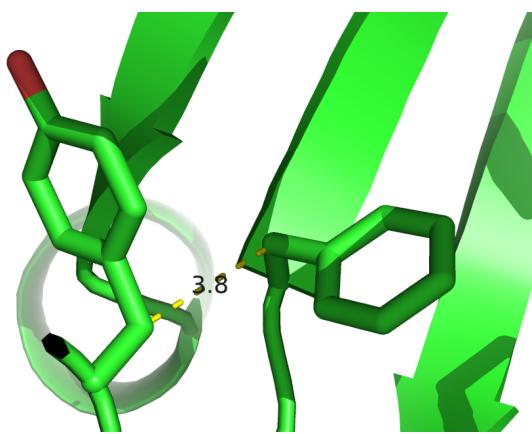


Figure E.5: The planar ring system of aromatic sidechains at short  $C_\beta$ - $C_\beta$  distances (e.g.  $\Delta C_\beta < 5\text{\AA}$ ) often points away from each other to avoid steric hindrance.



# F

## Training of the Random Forest Contact Prior

- F.1 Training Random Forest Model with pseudo-likelihood Feature**
- F.2 Evaluating window size with 5-fold Cross-validation**
- F.3 Evaluating non-contact threshold with 5-fold Cross-validation**
- F.4 Evaluating ratio of non-contacts and contacts in the training set with 5-fold Cross-validation**

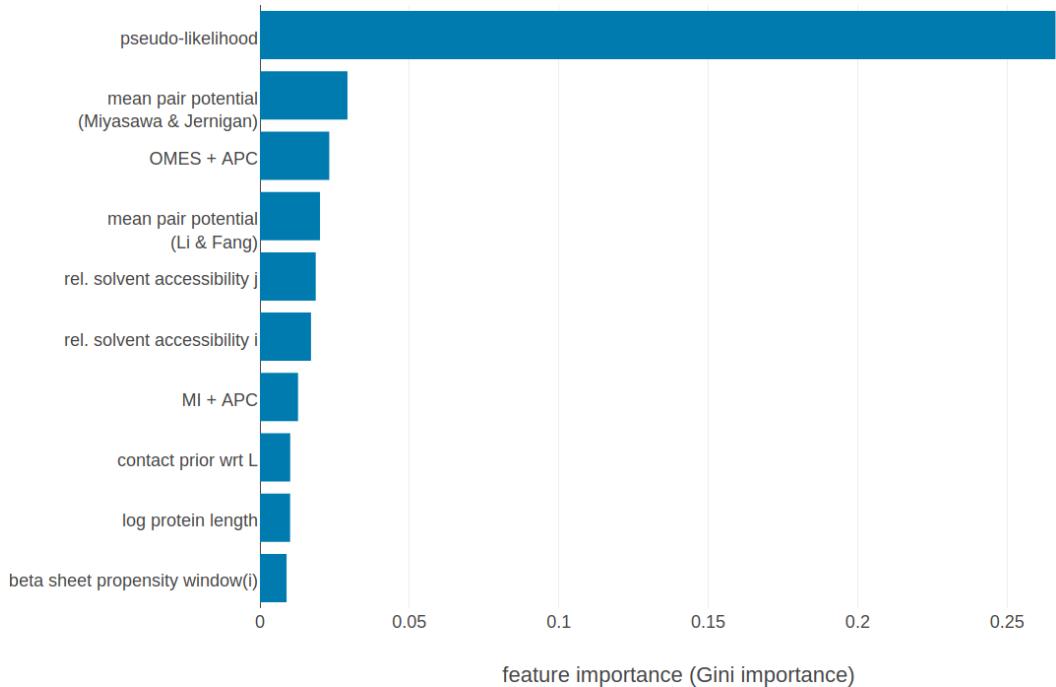


Figure F.1: Top ten features ranked according to *Gini importance*. **pseudo-likelihood**: [APC](#) corrected Frobenius norm of couplings computed with pseudo-likelihood. **mean pair potential (Miyasawa & Jernigan)**: average quasi-chemical energy of transfer of amino acids from water to the protein environment [163]. **OMES+APC**: [APC](#) corrected OMES score according to Fodor&Aldrich [162]. **mean pair potential (Li&Fang)**: average general contact potential by Li & Fang [25]. **rel. solvent accessibilty i(j)**: RSA score computed with Net-surfpr (v1.0) [164] for position i(j). **MI+APC**: [APC](#) corrected mutual information between amino acid counts (using pseudo-counts). **contact prior wrt L**: simple contact prior based on expected number of contacts wrt protein length (see methods section ??). **log protein length**: logarithm of protein length. **beta sheet propensity window(i)**: beta-sheet propensity according to Psipred [165] computed within a window of five positions around i. Features are described in detail in methods section 4.11.

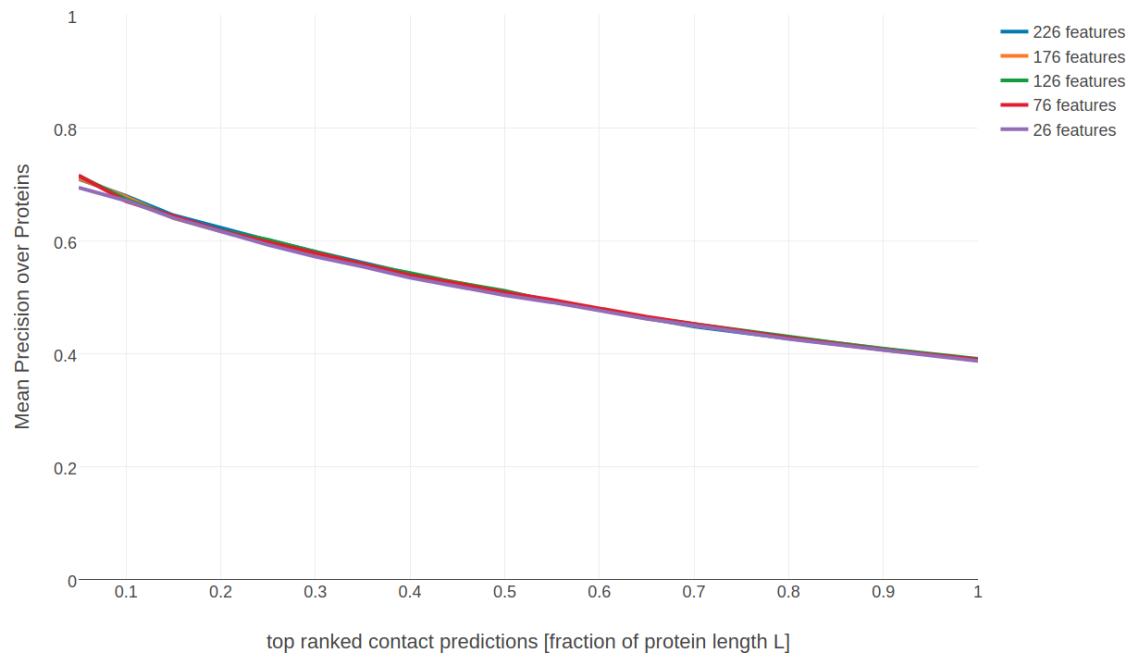


Figure F.2: Mean precision for top ranked contacts over 200 proteins for various random forest models trained on subsets of features. Subsets of features have been selected as described in section 4.12.1.

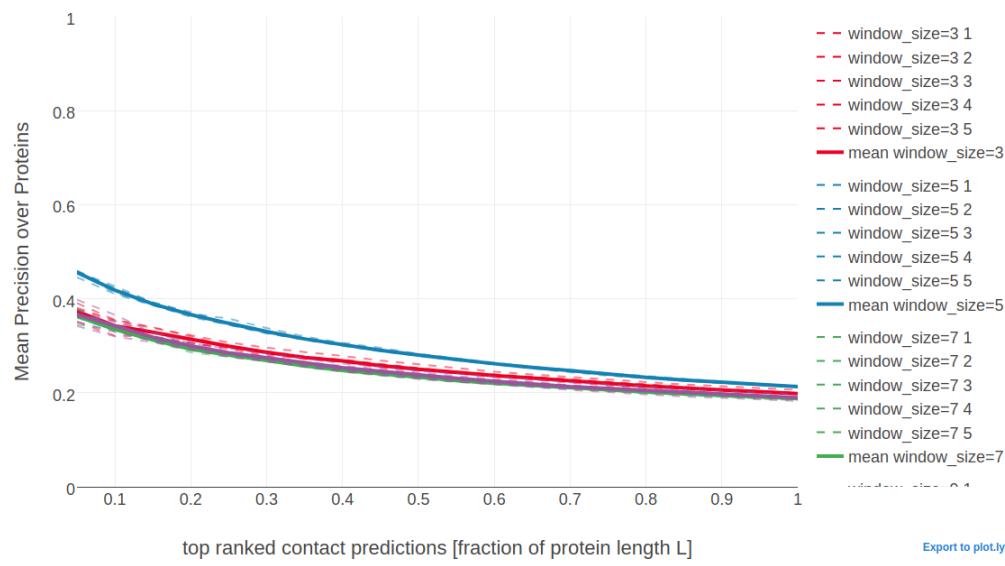


Figure F.3: Mean precision over validation set of 200 proteins for top ranked contact predictions for different choices of window size for single position features. Dashed lines represent the models trained on four subsets of the training data according to the 5-fold cross-validation scheme. Solid lines represent the mean over the five cross-validation models.

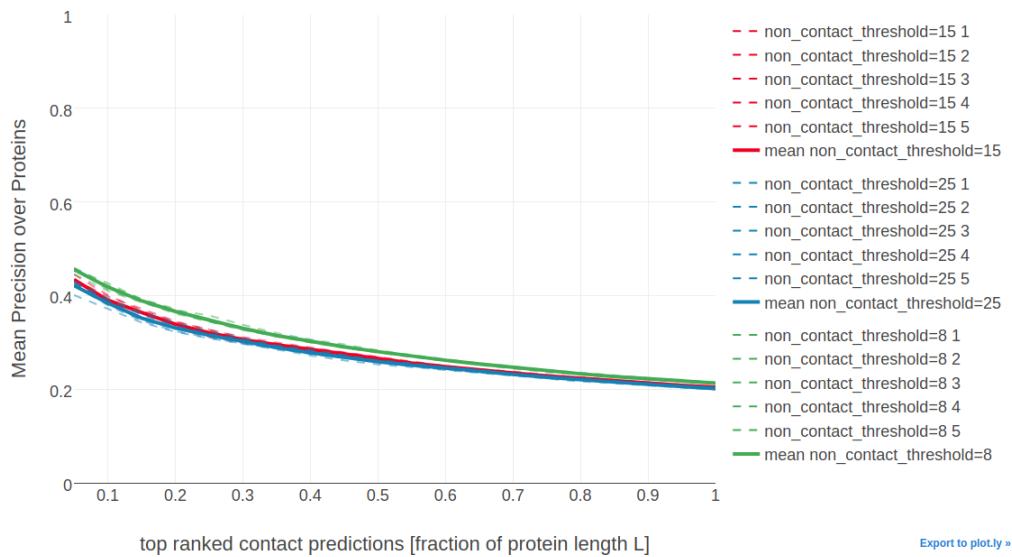


Figure F.4: Mean precision over validation set of 200 proteins for top ranked contact predictions for different choices of the non-contact threshold to define non-contacts. Dashed lines represent the models trained on four subsets of the training data according to the 5-fold cross-validation scheme. Solid lines represent the mean over the five cross-validation models.

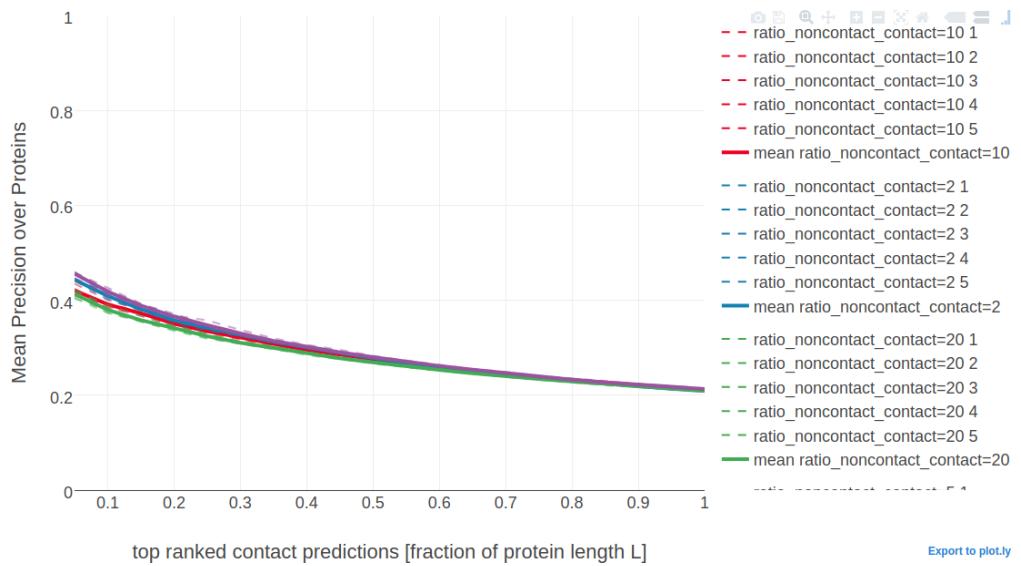


Figure F.5: Mean precision over validation set of 200 proteins for top ranked contact predictions for different choices of dataset composition with respect to the ratio of contacts and non-contacts. Dashed lines represent the models trained on four subsets of the training data according to the 5-fold cross-validation scheme. Solid lines represent the mean over the five cross-validation models.

# List of Figures

- 1.1 The evolutionary record of a protein family reveals evidence of compensatory mutations between spatially neighboring residues that are under selective pressure with respect to some physico-chemical constraints. Mining protein family sequence alignments for residue pairs with strong coevolutionary signals using statistical methods allows inference of spatial proximity for these residue pairs. . . . . 2
- 1.2 Effects of chained covariation obscure signals from true physical interactions. Consider residues A through E with physical interactions between the residue pairs A-B, B-C and D-E. The thickness of blue lines between residues reflects the strength of statistical dependencies between the corresponding alignment columns. Strong statistical dependencies between residue pairs (A,B) and (B,C) can induce a strong dependency between the spatially distant residues A and C. Covariation signals arising from transitive effects can become even stronger than other direct covariation signals and lead to false positive predictions. . . . . 3
- 1.3 Contact maps computed from pseudo-likelihood couplings. Subplot on top of the contact maps illustrates the normalized Shannon entropy (pink line) and percentage of gaps for every position in the alignment (brown line). **Left:** Contact map computed with Frobenius norm as in eq. (1.15). Overall coupling values are dominated by entropic effects, i.e. the amount of variation for a MSA position, leading to striped brightness patterns. For example, positions with high column entropy (e.g. positions 7, 12 or 31) have higher overall coupling values than positions with low column entropy (e.g. positions 11, 24 or 33). **b:** previous contact map but corrected for background noise with the APC as in eq. (1.17). . . . . 12
- 1.4 Generalized structure prediction pipeline integrating predicted contacts in form of distance constraints that guide conformational sampling. . . . . 13
- 1.5 Concatenating two multiple sequence alignments. In case multiple paralogs exist for a gene in one species the correct interaction partner needs to be identified and matched (marked with arrows). Sequences that cannot be paired with a unique interaction partner need to be discarded (marked with x). . . . . 14

|      |   |    |
|------|---|----|
| 1.6  | Distribution of residue pair $C_\beta$ distances over 6741 proteins in the dataset (see Methods 4.1) at different minimal sequence separation thresholds. . . . .   | 18 |
| 1.7  | $C_\beta$ distances between neighboring residues in $\alpha$ -helices. Left: Direct neighbors in $\alpha$ -helices have $C_\beta$ distances around $5.4\text{\AA}$ due to the geometrical constraints from $\alpha$ -helical architecture. Right: Residues separated by two positions ( $ i - j  = 2$ ) are less geometrically restricted to $C_\beta$ distances between $7\text{\AA}$ and $7.5\text{\AA}$ . . . . .  | 18 |
| 1.8  | The phylogenetic dependence of closely related sequences can produce covariation signals. Here, two independent mutation events (highlighted in red) in two branches of the tree result in a perfect covariation signal for two positions. . . . .  | 20 |
| 1.9  | Distribution of PFAM family sizes. Less than half of the families in PFAM (7990 compared to 8489 families) do not have an annotated structure. The median family size in number of sequences for families with and without annotated structures is 185 and 827 respectively. Data taken from PFAM 31.0 (March 2017, 16712 entries) [144]. . . . .   | 21 |
| 1.10 | Possible sources of coevolutionary signals. <b>a)</b> Physical interactions between intra-domain residues. <b>b)</b> Interactions across the interface of predominantly homo-oligomeric complexes. <b>c)</b> Interactions mediated by ligands or metal atoms. <b>d)</b> Transient interactions due to conformational flexibility. . . . .   | 22 |
| 2.1  | <b>Left</b> Pearson correlation of squared coupling values $(w_{ijab})^2$ with contact class (contact=1, non-contact=0). <b>Right</b> Standard deviation of squared coupling values for residue pairs in contact. Dataset contains 100.000 residue pairs per class (for details see methods section 4.7.1). Amino acids are abbreviated with one-letter code and they are broadly grouped with respect to physico-chemical properties listed in Appendix B. . . . . | 26 |
| 2.2  | <b>Left</b> Pearson correlation of raw signed coupling values $w_{ijab}$ with contact class (contact=1, non-contact=0). <b>Right</b> Standard deviation of coupling values for residue pairs in physical contact. Dataset contains 100.000 residue pairs per class (for details see section 4.7.1). Amino acids are abbreviated with one-letter code and they are broadly grouped with respect to physico-chemical properties listed in Appendix B. . . . .         | 27 |

|     |   |    |
|-----|---|----|
| 2.3 | Couplings $w_{ijab}$ and single potentials $v_{ia}$ and $v_{ja}$ computed with pseudo-likelihood for residues 6 and 82 in protein chain 1a9x_A_05. The matrix shows the 20x20 couplings $w_{ijab}$ with color representing coupling strength and direction (red = positive coupling value, blue = negative coupling value) and diameter of bubbles representing absolute coupling value $ w_{ijab} $ . Bars at the x-axis and y-axis correspond to the <i>Potts</i> model single potentials $v_i$ and $v_j$ respectively. Color reflects the value of single potentials. Amino acids are abbreviated with one-letter code and they are broadly grouped with respect to physico-chemical properties listed in Appendix B. . . . .  | 28 |
| 2.4 | Couplings $w_{ijab}$ and single potentials $v_{ia}$ and $v_{ja}$ computed with pseudo-likelihood for residues 29 and 39 in protein chain 1ae9_A_00. The matrix shows the 20x20 couplings $w_{ijab}$ with color representing coupling strength and direction (red = positive coupling value, blue = negative coupling value) and diameter of bubbles representing absolute coupling value $ w_{ijab} $ . Bars at the x-axis and y-axis correspond to the <i>Potts</i> model single potentials $v_i$ and $v_j$ respectively. Color reflects the value of single potentials. Amino acids are abbreviated with one-letter code and they are broadly grouped with respect to physico-chemical properties listed in Appendix B. . . . . | 29 |
| 2.5 | Interactions between protein side chains. <b>Left:</b> residue 6 (E) forms a salt bridge with residue 82 (R) in protein chain 1a9x_A_05. <b>Right:</b> residue 29 (A) and residue 39 (L) within the hydrophobic core of protein chain 1ae9_A_00. . . . .  | 30 |
| 2.6 | Distribution of selected couplings for filtered residue pairs with $C_\beta - C_\beta$ distances $< 5\text{\AA}$ (see methods section 4.7.2 for details). Number of coupling values used to determine the distribution is given in brackets in the legend. R-E = couplings for arginine and glutamic acid pairs, C-C = coupling for cystein residue pairs, V-I = coupling for valine and isoleucine pairs, F-W = coupling for phenylalanine and tryptophane pairs, E-E = coupling for glutamic acid residue pairs. . . . .  | 31 |
| 2.7 | Distribution of selected couplings for filtered residue pairs with $C_\beta - C_\beta$ distances between $8\text{\AA}$ and $12\text{\AA}$ (see methods section 4.7.2 for details). Number of coupling values used to determine the distribution is given in brackets in the legend. Couplings are the same as in Figure 2.6. . . . .  | 32 |
| 2.8 | Distribution of selected couplings for filtered residue pairs with $C_\beta - C_\beta$ distances between $20\text{\AA}$ and $50\text{\AA}$ (see methods section 4.7.2 for details). Number of coupling values used to determine the distribution is given in brackets in the legend. Couplings are the same as in Figure 2.6. . . . .   | 33 |

2.9 Two-dimensional distribution of approximately 10000 coupling values computed with pseudo-likelihood. **Top Left** The 2-dimensional distribution of couplings E-R and R-E for residue pairs with  $C_\beta - C_\beta$  distances  $< 8\text{\AA}$  is almost symmetric and the coupling values are positively correlated. **Top Right** The 2-dimensional distribution of couplings E-R and E-E for residue pairs with  $C_\beta - C_\beta$  distances  $< 8\text{\AA}$  is almost symmetric and the coupling values are negatively correlated. **Bottom Left** The 2-dimensional distribution of couplings I-L and V-I for residue pairs with  $C_\beta - C_\beta$  distances  $< 8\text{\AA}$  is symmetrically distributed around zero without visible correlation. **Bottom Right** The 2-dimensional distribution of couplings I-L and V-I for residue pairs with  $C_\beta - C_\beta$  distances  $> 20\text{\AA}$  is tightly distributed around zero. . . . . 35

3.1 Classifying new data with random forests. A new data sample is run down every tree in the forest until it ends up in a leaf node. Every leaf node has associated class probabilities  $p(c)$  reflecting the fraction of training samples at this leaf node belonging to every class  $c$ . The color of the leaf nodes reflects the class with highest probability. The predictions from all trees in form of the class probabilities are averaged and yield the final prediction. . . . . 38

3.2 Mean precision over 200 proteins against highest scoring contact predictions from random forest models for different settings of  $n\_estimators$  and  $max\_depth$ . Dashed lines show the performance of models that have been learned on the five different subsets of training data. Solid lines give the mean precision over the five models. Only those models are shown that yielded the five highest mean precision values (given in parentheses in the legend). Random forest models with 1000 trees and maximum depth of trees of either 100, 1000 or unrestricted tree depth perform nearly identical (lines overlap). Random forest models with 500 trees and  $max\_depth=10$  or  $max\_depth=100$  perform slightly worse. . . . . 40

3.3 Mean precision over 200 proteins against highest scoring contact predictions from random forest models with different settings of  $min\_samples\_leaf$  and  $max\_features$ . Dashed lines show the performance of models that have been learned on the five different subsets of training data. Solid lines give the mean precision over the five models. Only those models are shown that yielded the five best mean precision values (given in parentheses in the legend). . . . . 41

|     |  |    |
|-----|--|----|
| 3.4 | Top ten features ranked according to <i>Gini importance</i> . <b>OMES+APC</b> : APC corrected OMES score according to Fodor&Aldrich [162]. <b>mean pair potential (Miyasawa &amp; Jernigan)</b> : average quasi-chemical energy of transfer of amino acids from water to the protein environment [163]. <b>MI+APC</b> : APC corrected mutual information between amino acid counts (using pseudo-counts). <b>mean pair potential (Li&amp;Fang)</b> : average general contact potential by Li & Fang [25]. <b>rel. solvent accessibility i(j)</b> : RSA score computed with Netsurfp (v1.0) [164] for position i(j). <b>pairwise gap%</b> : percentage of gapped sequences at either position i and j. <b>correlation mean isoelectric feature</b> : Pearson correlation between the mean isoelectric point feature (according to Zimmermann et al., 1968) for positions i and j. <b>sequence separation</b> : $ j-i $ . <b>beta sheet propensity window(i)</b> : beta-sheet propensity according to Psipred [165] computed within a window of five positions around i. Features are described in detail in methods section 4.11. . . . . | 43 |
| 3.5 | Mean precision of top ranked predictions over 200 proteins for random forest models trained on subsets of features of decreasing importance. Subsets of features have been selected as described in methods section 4.12.1. . . . .  | 44 |
| 3.6 | Mean precision for top ranked contacts on a test set of 774 proteins. <b>random forest (pLL)</b> = random forest model using sequence derived features and pseudo-likelihood contact score (APC corrected Frobenius norm of couplings). <b>pseudo-likelihood</b> = APC corrected Frobenius norm of couplings computed with pseudo-likelihood. <b>random forest</b> = random forest model trained on 75 sequence derived features. <b>OMES</b> = APC corrected <i>OMES</i> contact score according to Fodor&Aldrich [162]. <b>mutual information</b> = APC corrected mutual information between amino acid counts (using pseudo-counts). . . . .  | 45 |
| 3.7 | Mean precision for top ranked contacts on a test set of 774 proteins splitted into four equally sized subsets with respect to Neff. Subsets are defined according to quantiles of Neff values. Upper left: Subset of proteins with $\text{Neff} < Q1$ . Upper right: Subset of proteins with $Q1 \leq \text{Neff} < Q2$ . Lower left: Subset of proteins with $Q2 \leq \text{Neff} < Q3$ . Lower right: Subset of proteins with $Q3 \leq \text{Neff} < Q4$ . <b>random forest (pLL)</b> = random forest model using sequence derived features and pseudo-likelihood contact score (APC corrected Frobenius norm of couplings). <b>pseudo-likelihood</b> = APC corrected Frobenius norm of couplings computed with pseudo-likelihood. <b>random forest</b> = random forest model trained on 75 sequence derived features. <b>OMES</b> = APC corrected <i>OMES</i> contact score according to Fodor&Aldrich [162]. <b>mutual information</b> = APC corrected mutual information between amino acid counts (using pseudo-counts). . . . .   | 46 |

|     |   |    |
|-----|---|----|
| 4.1 | Distribution of CATH classes (1=mainly $\alpha$ , 2=mainly $\beta$ , 3= $\alpha - \beta$ ) in the dataset and the ten subsets. . . . .  | 50 |
| 4.2 | Mean precision over 3124 proteins of top ranked contacts computed as APC corrected Frobenius norm of couplings. Couplings have been computed with CCMpred [57] and CCMpredPy as specified in the legend. Specific flags that have been used to run both methods are described in detail in the text (see section 4.2.1). . . . .  | 51 |
| 4.3 | Number of contacts ( $C_\beta < 8\text{\AA}$ ) with respect to protein length and sequence separation has a linear relationship. . . . .  | 54 |
| 4.4 | Hypothetical MSA consisting of two sets of sequences: the first set has sequences covering only the left half of columns, while the second set has sequences covering only the right half of columns. The two blocks could correspond to protein domains that were aligned to a single query sequence. Empirical amino acid pair frequencies $q(x_i = a, x_j = b)$ will vanish for positions $i$ from the left half and $j$ from the right half of the alignment. . . . .   | 56 |
| 4.5 | The 400 gradients $\nabla_{w_{ijab}} LL_{\text{reg}}(\mathbf{v}, \mathbf{w})$ at position $(i, j)$ for $a, b \in \{1, \dots, 20\}$ are not independent. Red bars represent pairwise amino acid counts at position $(i, j)$ for the sampled alignment. Blue bars represent pairwise amino acid counts at position $(i, j)$ for the input alignment. The sum over pairwise amino acid counts at position $(i, j)$ for both alignments is $N_{ij}$ , which is the number of ungapped sequences. The gradient $\nabla_{w_{ijab}} LL_{\text{reg}}(\mathbf{v}, \mathbf{w})$ is computed as the difference of pairwise amino acid counts for amino acids $a$ and $b$ at position $(i, j)$ . The sum over gradients $\nabla_{w_{ijab}} LL_{\text{reg}}(\mathbf{v}, \mathbf{w})$ at position $(i, j)$ for all $a, b \in \{1, \dots, 20\}$ is zero. . . . . | 63 |
| 4.6 | Mean precision for top ranked contact predictions over 286 proteins. Contact scores are computed as the APC corrected Frobenius norm of the couplings $\mathbf{w}_{ij}$ . <b>pseudo-likelihood</b> : couplings computed with pseudo-likelihood. <b>CD</b> : couplings computed with CD using stochastic gradient descent with different initial learning rates $\alpha_0$ as specified in the legend. . . . .   | 64 |
| 4.7 | Convergence plots for two proteins during SGD optimization with different learning rates and convergence measured as L2-norm of the coupling parameters $\ \mathbf{w}\ _2$ . Linear learning rate annealing schedule has been used with decay rate $\gamma = 0.01$ and initial learning rates $\alpha_0$ have been set as specified in the legend. <b>Left</b> Convergence plot for protein 1mkc_A_00 having protein length $L=43$ and 142 sequences in the alignment (Neff=96). <b>Right</b> Convergence plot for protein 1c75_A_00 having protein length $L=71$ and 28078 sequences in the alignment (Neff=16808). Figure is cut at the yaxis at $\ \mathbf{w}\ _2 = 1500$ , but learning rate of $5e-3$ reaches $\ \mathbf{w}\ _2 \approx 13000$ . . . . .   | 65 |

|      |  |    |
|------|--|----|
| 4.8  | Mean precision for top ranked contact predictions over 288 proteins. Contact scores are computed as the APC corrected Frobenius norm of the couplings $\mathbf{w}_{ij}$ . <b>pseudo-likelihood</b> : couplings computed with pseudo-likelihood. <b>CD</b> : couplings computed with CD using stochastic gradient descent with an initial learning rate defined with respect to $\text{Neff}$ . Learning rate annealing schedules and decay rates as specified in the legend. . . . .   | 66 |
| 4.9  | L2-norm of the coupling parameters $\ \mathbf{w}\ _2$ during stochastic gradient descent optimization with different learning rates schedules. The initial learning rate $\alpha_0$ is defined with respect to $\text{Neff}$ as given in eq. (4.47). Learning rate schedules and decay rates are used according to the legend. <b>Left</b> Convergence plot for protein 1mkc_A_00 having protein length $L=43$ and 142 sequences in the alignment ( $\text{Neff}=96$ ). <b>Right</b> Convergence plot for protein 1c75_A_00 having protein length $L=71$ and 28078 sequences in the alignment ( $\text{Neff}=16808$ ). . . . . | 68 |
| 4.10 | Distribution of the number of iterations until convergence for SGD optimizations of the full likelihood for different learning rate schedules. Convergence is reached when the relative difference of parameter norms $\ \mathbf{w}\ _2$ falls below $\epsilon = 1e - 8$ . Initial learning rate $\alpha_0$ is defined with respect to $\text{Neff}$ as given in eq. (4.47) and maximum number of iterations is set to 5000. Learning rate schedules and decay rates are used as specified in the legend. . . . .  | 69 |
| 4.11 | L2-norm of the coupling parameters $\ \mathbf{w}\ _2$ during optimization with <i>ADAM</i> and different learning rates without annealing. The learning rate $\alpha$ is specified in the legend. <b>Left</b> Convergence plot for protein 1mkc_A_00 having protein length $L=43$ and 142 sequences in the alignment ( $\text{Neff}=96$ ). <b>Right</b> Convergence plot for protein 1c75_A_00 having protein length $L=71$ and 28078 sequences in the alignment ( $\text{Neff}=16808$ ). . . . .  | 70 |
| 4.12 | L2-norm of the coupling parameters $\ \mathbf{w}\ _2$ during optimization with <i>ADAM</i> and different learning rate annealing schedules. The learning rate $\alpha$ is specified with respect to $\text{Neff}$ as $\alpha = 2e-3 \log(\text{Neff})$ . The learning rate annealing schedule is specified in the legend. <b>Left</b> Convergence plot for protein 1mkc_A_00 having protein length $L=43$ and 142 sequences in the alignment ( $\text{Neff}=96$ ). <b>Right</b> Convergence plot for protein 1c75_A_00 having protein length $L=71$ and 28078 sequences in the alignment ( $\text{Neff}=16808$ ). . . . .      | 71 |
| 4.13 | Mean precision for top ranked contact predictions over 288 proteins. Contact scores are computed as the APC corrected Frobenius norm of the couplings $\mathbf{w}_{ij}$ . <b>pseudo-likelihood</b> : couplings computed with pseudo-likelihood. <b>CD lambda_w = X</b> : couplings computed with CD using L2-regularization on the couplings $\mathbf{w}$ with regularization coefficient $\lambda_w$ chosen as specified in the legend and keeping the single potentials $v_i$ fixed at their MLE optimum $v_i^*$ given in eq. (4.43). . . . .  | 72 |

|  |    |
|--|----|
| 4.14 Performance of contrastive divergence optimization of the full likelihood with different number of Gibbs steps compared to pseudo-likelihood (blue) for 287 proteins. Contact scores are computed as the APC corrected Frobenius norm of the couplings $w_{ij}$ . pseudo-likelihood: contact scores computed from pseudo-likelihood. The other methods derive contact scores from couplings computed from CD with different number of Gibbs sampling steps. . . . . | 75 |
| 4.15 Observed number of contacts per residue has a non-linear relationship with protein length. Distribution is shown for several thresholds of sequence separation $ j-i $ . . . . .  | 77 |
| 4.16 Distribution of average number of contacts per residue against protein length and corresponding linear regression fits. Protein length is on logarithmic scale. Distribution and linear regression fits are shown for different sequence separation thresholds $ j-i $ . . . . .  | 79 |
| 4.17 Fraction of contacts among all possible contacts in a protein against protein length L. The distribution has a non-linear relationship. At a sequence separation $>8$ positions the fraction of contacts for intermediate size proteins with length $>100$ is approximately 2%. Data set contains 6368 proteins and is explained in methods section 4.1. . . . .  | 81 |
| C.1 Distribution of alignment diversity ( $= \sqrt{(\frac{N}{L})}$ ) in the dataset and its ten subsets. . . . .   | 88 |
| C.2 Distribution of gap percentage of alignments in the dataset and its ten subsets. . . . .   | 88 |
| C.3 Distribution of alignment size (number of sequences N) in the dataset and its ten subsets. . . . .   | 89 |
| C.4 Distribution of protein length L in the dataset and its ten subsets. . . . .   | 89 |
| D.1 Standard deviation of squared coupling values $w_{ijab}^2$ for residue pairs not in physical contact ( $\Delta C_\beta > 25\text{\AA}$ ). Dataset contains 100.000 residue pairs per class (for details see methods section 4.7.1). Amino acids are abbreviated with one-letter code and they are broadly grouped with respect to physico-chemical properties listed in Appendix B . . . . .   | 92 |
| D.2 Standard deviation of coupling values $w_{ijab}$ for residue pairs not in physical contact ( $\Delta C_\beta > 25\text{\AA}$ ). Dataset contains 100.000 residue pairs per class (for details see section 4.7.1). Amino acids are abbreviated with one-letter code and they are broadly grouped with respect to physico-chemical properties listed in Appendix B. . . . .  | 93 |
| E.1 Tyrosine (residue 37) and Lysine (residue 48) forming a cation- $\pi$ interaction in protein 2ayd. . . . .   | 95 |
| E.2 Two cystein residues (residues 54 and 64) forming a covalent disulfide bond in protein 1alu. . . . .   | 96 |

|     |  |     |
|-----|--|-----|
| E.3 | Proline and tryptophan (residues 17 and 34) stacked on top of each otherengaging in a CH/π interaction in protein chain 1aol_A_00. . . . .   | 96  |
| E.4 | Network-like structure of aromatic residues in the protein core. 80% of aromatic residues are involved in such networks that are important for protein stability [153]. . . . .  | 97  |
| E.5 | The planar ring system of aromatic sidechains at short $C_\beta$ - $C_\beta$ distances (e.g. $\Delta C_\beta < 5\text{\AA}$ ) often points away from each other to avoid steric hindrance. . . . .   | 97  |
| F.1 | Top ten features ranked according to <i>Gini importance</i> . <b>pseudo-likelihood</b> : APC corrected Frobenius norm of couplings computed with pseudo-likelihood. <b>mean pair potential (Miyasawa &amp; Jernigan)</b> : average quasi-chemical energy of transfer of amino acids from water to the protein environment [163]. <b>OMES+APC</b> : APC corrected OMES score according to Fodor&Aldrich [162]. <b>mean pair potential (Li&amp;Fang)</b> : average general contact potential by Li & Fang [25]. <b>rel. solvent accessibility i(j)</b> : RSA score computed with Netsurfpp (v1.0) [164] for position i(j). <b>MI+APC</b> : APC corrected mutual information between amino acid counts (using pseudo-counts). <b>contact prior wrt L</b> : simple contact prior based on expected number of contacts wrt protein length (see methods section ??). <b>log protein length</b> : logarithm of protein length. <b>beta sheet propensity window(i)</b> : beta-sheet propensity according to Psipred [165] computed within a window of five positions around i. Features are described in detail in methods section 4.11. . . . . | 100 |
| F.2 | Mean precision for top ranked contacts over 200 proteins for variaous random forest models trained on subsets of features. Subsets of features have been selected as described in section 4.12.1. . . . .  | 101 |
| F.3 | Mean precision over validation set of 200 proteins for top ranked contact predictions for different choices of window size for single position features. Dashed lines represent the models trained on four subsets of the training data according to the 5-fold cross-validation scheme. Solid lines represent the mean over the five cross-validation models. . . . .   | 101 |
| F.4 | Mean precision over validation set of 200 proteins for top ranked contact predictions for different choices of the non-contact threshold to define non-contacts. Dashed lines represent the models trained on four subsets of the training data according to the 5-fold cross-validation scheme. Solid lines represent the mean over the five cross-validation models. . . . .   | 102 |

F.5 Mean precision over validation set of 200 proteins for top ranked contact predictions for different choices of dataset composition with respect to the ratio of contacts and non-contacts. Dashed lines represent the models trained on four subsets of the training data according to the 5-fold cross-validation scheme. Solid lines represent the mean over the five cross-validation models. . . . . 102

## List of Tables

|     |   |    |
|-----|---|----|
| 4.1 | Features characterizing the total alignment . . . . .   | 76 |
| 4.2 | Caption hereSingle Position Sequence Features . . . . .   | 77 |
| 4.3 | Pairwise Sequence Features . . . . .  | 79 |
| 4.4 | Important machine learning contact prediction approaches and<br>their choices for rebalancing the data set. . . . . | 81 |
| B.1 | Amino acid abbreviations and physico-chemical properties accord-<br>ing to Livingstone et al., 1993 [189] . . . . . | 86 |



## References

1. Gbel, U., Sander, C., Schneider, R., and Valencia, A. (1994). Correlated mutations and residue contacts in proteins. *Proteins* *18*, 309–17. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/8208723>.
2. Godzik, A., and Sander, C. (1989). Conservation of residue interactions in a family of Ca-binding proteins. "Protein Eng. Des. Sel." *2*, 589–596. Available at: <https://academic.oup.com/peds/article-lookup/doi/10.1093/protein/2.8.589>.
3. Neher, E. (1994). How frequent are correlated changes in families of protein sequences? *Proc. Natl. Acad. Sci. U. S. A.* *91*, 98–102. Available at: <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=42893{\&}tool=pmcentrez{\&}rendertype=fulltext>.
4. Taylor, W.R., and Hatrik, K. (1994). Compensating changes in protein multiple sequence alignments. "Protein Eng. Des. Sel." *7*, 341–348. Available at: <http://peds.oxfordjournals.org/content/7/3/341.abstract>.
5. Oliveira, L., Paiva, A.C.M., and Vriend, G. (2002). Correlated mutation analyses on very large sequence families. *Chembiochem* *3*, 1010–7. Available at: [http://onlinelibrary.wiley.com/doi/10.1002/1439-7633\(20021004\)3:10{\&}3C1010::AID-CBIC1010{\&}3E3.0.CO;2-T/full](http://onlinelibrary.wiley.com/doi/10.1002/1439-7633(20021004)3:10{\&}3C1010::AID-CBIC1010{\&}3E3.0.CO;2-T/full).
6. Shindyalov, I., Kolchanov, N., and Sander, C. (1994). Can three-dimensional contacts in protein structures be predicted by analysis of correlated mutations? "Protein Eng. Des. Sel." *7*, 349–358. Available at: <http://peds.oxfordjournals.org/content/7/3/349>.
7. Clarke, N.D. (1995). Covariation of residues in the homeodomain sequence family. *Protein Sci.* *4*, 2269–78. Available at: <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=42893{\&}tool=pmcentrez{\&}rendertype=fulltext>.
8. Korber, B. (1993). Covariation of Mutations in the V3 Loop of Human Immunodeficiency Virus Type 1 Envelope Protein: An Information Theoretic Analysis. *Proc. Natl. Acad. Sci.* *90*, 7176–7180. Available at: <http://www.pnas.org/content/90/15/7176.abstract?ijkey=6e1c9cbdef66bd0beefedc88ea07591b837f2213>.
9. Martin, L.C., Gloor, G.B., Dunn, S.D., and Wahl, L.M. (2005). Using information theory to search for co-evolving residues in proteins. *Bioinformatics* *21*, 4116–24. Available at: <http://bioinformatics.oxfordjournals.org/content/21/22/4116.full>.
10. Atchley, W.R., Wollenberg, K.R., Fitch, W.M., Terhalle, W., and Dress, A.W. (2000). Correlations Among Amino Acid Sites in bHLH Protein Domains: An Information Theoretic Analysis. *Mol. Biol. Evol.* *17*, 164–178. Available at:

- <http://mbe.oxfordjournals.org/content/17/1/164.abstract?ijkey=2a1f0a044a8fd2213955e4f2c17f>
11. Fodor, A.A., and Aldrich, R.W. (2004). Influence of conservation on calculations of amino acid covariance in multiple sequence alignments. *Proteins* *56*, 211–21. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/15211506>.
  12. Tillier, E.R., and Lui, T.W. (2003). Using multiple interdependency to separate functional from phylogenetic correlations in protein alignments. *Bioinformatics* *19*, 750–755. Available at: <http://bioinformatics.oxfordjournals.org/content/19/6/750.abstract>
  13. Gouveia-Oliveira, R., and Pedersen, A.G. (2007). Finding coevolving amino acid residues using row and column weighting of mutual information and multi-dimensional amino acid representation. *Algorithms Mol. Biol.* *2*, 12. Available at: <http://www.almob.org/content/2/1/12>.
  14. Dunn, S.D., Wahl, L.M., and Gloor, G.B. (2008). Mutual information without the influence of phylogeny or entropy dramatically improves residue contact prediction. *Bioinformatics* *24*, 333–40. Available at: <http://bioinformatics.oxfordjournals.org/content/24/3/333>.
  15. Kass, I., and Horovitz, A. (2002). Mapping pathways of allosteric communication in GroEL by analysis of correlated mutations. *Proteins* *48*, 611–7. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/12211028>.
  16. Noivirt, O., Eisenstein, M., and Horovitz, A. (2005). Detection and reduction of evolutionary noise in correlated mutation analysis. *Protein Eng. Des. Sel.* *18*, 247–53. Available at: <http://peds.oxfordjournals.org/content/18/5/247.full>.
  17. Lapedes, A., Giraud, B., Liu, L., and Stormo, G. (1999). Correlated mutations in models of protein sequences: phylogenetic and structural effects. *93*, 236–256. Available at: <http://www.citeulike.org/user/qluo/article/5092214>.
  18. Burger, L., and Nimwegen, E. van (2010). Disentangling direct from indirect co-evolution of residues in protein alignments. *PLoS Comput. Biol.* *6*, e1000633. Available at: <http://dx.plos.org/10.1371/journal.pcbi.1000633>.
  19. Weigt, M., White, R.A., Szurmant, H., Hoch, J.A., and Hwa, T. (2009). Identification of direct residue contacts in protein-protein interaction by message passing. *Proc. Natl. Acad. Sci. U. S. A.* *106*, 67–72. Available at: <http://www.pnas.org/content/106/1/67.abstract>.
  20. Juan, D. de, Pazos, F., and Valencia, A. (2013). Emerging methods in protein co-evolution. *Nat. Rev. Genet.* *14*, 249–61. Available at: <http://www.readcube.com/articles/10.1038/nrg3414?locale=en>.
  21. Jones, D.T., Buchan, D.W.A., Cozzetto, D., and Pontil, M. (2012). PSICOV: precise structural contact prediction using sparse inverse covariance estimation on large multiple sequence alignments. *Bioinformatics* *28*, 184–90. Available at: <http://bioinformatics.oxfordjournals.org/content/28/2/184.full>.
  22. Burger, L., and Nimwegen, E. van (2008). Accurate prediction of protein-protein interactions from sequence alignments using a Bayesian method. *Mol. Syst. Biol.* *4*, 165. Available at: <http://msb.embopress.org/content/4/1/>

**165.abstract.**

23. Cheng, J., and Baldi, P. (2007). Improved residue contact prediction using support vector machines and a large feature set. *BMC Bioinformatics* *8*, 113. Available at: <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC1852326/>{\&}tool=pmcentrez{\&}db=bioassess
24. Wu, S., and Zhang, Y. (2008). A comprehensive assessment of sequence-based and template-based methods for protein contact prediction. *Bioinformatics* *24*, 924–31. Available at: <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC2648832/>{\&}tool=pmcentrez
25. Li, Y., Fang, Y., and Fang, J. (2011). Predicting residue-residue contacts using random forest models. *Bioinformatics* *27*, 3379–84. Available at: <http://bioinformatics.oxfordjournals.org/content/27/24/3379.long>.
26. Wang, X.-F., Chen, Z., Wang, C., Yan, R.-X., Zhang, Z., and Song, J. (2011). Predicting residue-residue contacts and helix-helix interactions in transmembrane proteins using an integrative feature-based random forest approach. *PLoS One* *6*, e26767. Available at: <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3203928/>{\&}tool=pmcentrez
27. Wang, Z., and Xu, J. (2013). Predicting protein contact map using evolutionary and physical constraints by integer programming. *Bioinformatics* *29*, i266–73. Available at: <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3694661/>{\&}tool=pmcentrez
28. Fariselli, P., Olmea, O., Valencia, A., and Casadio, R. (2001). Prediction of contact maps with neural networks and correlated mutations. *Protein Eng. Des. Sel.* *14*, 835–843. Available at: <http://peds.oxfordjournals.org/content/14/11/835.long>.
29. Shackelford, G., and Karplus, K. (2007). Contact prediction using mutual information and neural nets. *Proteins* *69 Suppl 8*, 159–64. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/17932918>.
30. Hamilton, N., Burrage, K., Ragan, M.A., and Huber, T. (2004). Protein contact prediction using patterns of correlation. *Proteins Struct. Funct. Bioinforma.* *56*, 679–684. Available at: <http://onlinelibrary.wiley.com/doi/10.1002/prot.20160/abstract>.
31. Xue, B., Faraggi, E., and Zhou, Y. (2009). Predicting residue-residue contact maps by a two-layer, integrated neural-network method. *Proteins* *76*, 176–83. Available at: <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC2716487/>{\&}tool=pmcentrez
32. Tegge, A.N., Wang, Z., Eickholt, J., and Cheng, J. (2009). NNcon: improved protein contact map prediction using 2D-recursive neural networks. *Nucleic Acids Res.* *37*, W515–8. Available at: <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC2703959/>{\&}tool=pmcentrez
33. Eickholt, J., and Cheng, J. (2012). Predicting protein residue–residue contacts using deep networks and boosting. *28*, 3066–3072.
34. Di Lena, P., Nagata, K., and Baldi, P. (2012). Deep architectures for protein contact map prediction. *Bioinformatics* *28*, 2449–57. Available at: <http://bioinformatics.oxfordjournals.org/content/28/19/2449.full{\#}sec-14>.
35. Chen, P., and Li, J. (2010). Prediction of protein long-range contacts using an ensemble of genetic algorithm classifiers with sequence profile centers. *BMC Struct. Biol.* *10 Suppl 1*, S2. Available at: <http://bmcbioinformatics.biomedcentral.com/articles/10.1186/1471-2105-10-S1-S2>.

[biomedcentral.com/articles/10.1186/1472-6807-10-S1-S2](http://biomedcentral.com/articles/10.1186/1472-6807-10-S1-S2).

36. Jones, D.T., Singh, T., Kosciolak, T., and Tetchner, S. (2015). MetaPSICOV: combining coevolution methods for accurate prediction of contacts and long range hydrogen bonding in proteins. *Bioinformatics* *31*, 999–1006. Available at: <http://bioinformatics.oxfordjournals.org/content/31/7/999.short>.
37. Skwark, M.J., Abdel-Rehim, A., and Elofsson, A. (2013). PconsC: combination of direct information methods and alignments improves contact prediction. *Bioinformatics* *29*, 1815–6. Available at: <http://bioinformatics.oxfordjournals.org/content/29/14/1815.long>.
38. Skwark, M.J., Michel, M., Menendez Hurtado, D., Ekeberg, M., and Elofsson, A. (2016). Accurate contact predictions for thousands of protein families using PconsC3. *bioRxiv*.
39. Schneider, M., and Brock, O. (2014). Combining Physicochemical and Evolutionary Information for Protein Contact Prediction. *PLoS One* *9*, e108438. Available at: <http://www.plosone.org/article/info{\%}3Adoi{\%}2F10.1371{\%}2Fjournal.pone.0108438>.
40. Jones, D.T., Singh, T., Kosciolak, T., and Tetchner, S. (2015). MetaPSICOV: combining coevolution methods for accurate prediction of contacts and long range hydrogen bonding in proteins. *Bioinformatics* *31*, 999–1006. Available at: <http://bioinformatics.oxfordjournals.org/content/31/7/999.short>.
41. Wang, S., Sun, S., Li, Z., Zhang, R., and Xu, J. (2016). Accurate De Novo Prediction of Protein Contact Map by Ultra-Deep Learning Model. *PLoS Comput. Biol.* *13*, e1005324. Available at: <http://arxiv.org/abs/1609.00680> <http://www.ncbi.nlm.nih.gov/pubmed/28056090> <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC5000000/>
42. Stahl, K., Schneider, M., and Brock, O. (2017). EPSILON-CP: using deep learning to combine information from multiple sources for protein contact prediction. *BMC Bioinformatics* *18*, 303. Available at: <http://bmcbioinformatics.biomedcentral.com/articles/10.1186/s12859-017-1713-x>.
43. He, B., Mortuza, S.M., Wang, Y., Shen, H.-B., and Zhang, Y. (2017). NeBcon: Protein contact map prediction using neural network training coupled with nave Bayes classifiers. *Bioinformatics*. Available at: <https://academic.oup.com/bioinformatics/article-lookup/doi/10.1093/bioinformatics/btx164>.
44. Andreani, J., and Sding, J. (2015). Bbcontacts: Prediction of  $\beta$ -strand pairing from direct coupling patterns. *Bioinformatics* *31*, 1729–37. Available at: <http://bioinformatics.oxfordjournals.org/content/31/11/1729>.
45. Skwark, M.J., Raimondi, D., Michel, M., and Elofsson, A. (2014). Improved contact predictions using the recognition of protein like contact patterns. *PLoS Comput. Biol.* *10*, e1003889. Available at: <http://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1003889>.
46. Monastyrskyy, B., D'Andrea, D., Fidelis, K., Tramontano, A., and Kryshtafovych, A. (2015). New encouraging developments in contact prediction: Assessment of the CASP11 results. *Proteins*. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/26474083>.
47. Jaynes, E.T. (1957). Information Theory and Statistical Mechanics I. *Phys. Rev.* *106*, 620–630. Available at: <https://link.aps.org/doi/10.1103/1103>

*PhysRev.* **106**. 620.

48. Jaynes, E.T. (1957). Information Theory and Statistical Mechanics. II. *Phys. Rev.* **108**, 171–190. Available at: <https://link.aps.org/doi/10.1103/PhysRev.108.171>.
49. Wainwright, M.J., and Jordan, M.I. (2007). Graphical Models, Exponential Families, and Variational Inference. *Found. Trends Mach. Learn.* **1**, 1–305. Available at: <http://www.nowpublishers.com/article/Details/MAL-001>.
50. Murphy, K.P. (2012). Machine Learning: A Probabilistic Perspective (MIT Press).
51. Morcos, F., Pagnani, A., Lunt, B., Bertolino, A., Marks, D.S., Sander, C., Zecchina, R., Onuchic, J.N., Hwa, T., and Weigt, M. (2011). Direct-coupling analysis of residue coevolution captures native contacts across many protein families. *Proc. Natl. Acad. Sci. U. S. A.* **108**, E1293–301. Available at: <http://www.pnas.org/content/108/49/E1293.full>.
52. Marks, D.S., Colwell, L.J., Sheridan, R., Hopf, T.A., Pagnani, A., Zecchina, R., and Sander, C. (2011). Protein 3D structure computed from evolutionary sequence variation. *PLoS One* **6**, e28766. Available at: <http://dx.plos.org/10.1371/journal.pone.0028766>.
53. Cocco, S., Feinauer, C., Figliuzzi, M., Monasson, R., and Weigt, M. (2017). Inverse Statistical Physics of Protein Sequences: A Key Issues Review. arXiv. Available at: <https://arxiv.org/pdf/1703.01222.pdf>.
54. Koller, D., and Friedman, N.I.R. (2009). Probabilistic graphical models: Principles and Techniques (MIT Press).
55. Ekeberg, M., Lvkvist, C., Lan, Y., Weigt, M., and Aurell, E. (2013). Improved contact prediction in proteins: Using pseudolikelihoods to infer Potts models. *Phys. Rev. E* **87**, 012707. Available at: <https://link.aps.org/doi/10.1103/PhysRevE.87.012707>.
56. Stein, R.R., Marks, D.S., and Sander, C. (2015). Inferring Pairwise Interactions from Biological Data Using Maximum-Entropy Probability Models. *PLOS Comput. Biol.* **11**, e1004182. Available at: <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC4520494/>
57. Seemayer, S., Gruber, M., and Sding, J. (2014). CCMpred-fast and precise prediction of protein residue-residue contacts from correlated mutations. *Bioinformatics, btu500*. Available at: <http://bioinformatics.oxfordjournals.org/content/early/2014/08/12/bioinformatics.btu500>.
58. Ekeberg, M., Hartonen, T., and Aurell, E. (2014). Fast pseudolikelihood maximization for direct-coupling analysis of protein structure from many homologous amino-acid sequences. *J. Comput. Phys.* **276**, 341–356. Available at: <http://www.sciencedirect.com/science/article/pii/S0021999114005178>.
59. Kamisetty, H., Ovchinnikov, S., and Baker, D. (2013). Assessing the utility of coevolution-based residue-residue contact predictions in a sequence-and structure-rich era. *Proc. Natl. Acad. Sci. U. S. A.* **110**, 15674–9. Available at: <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3785744/>
60. Lapedes, A., Giraud, B., and Jarzynski, C. (2012). Using Sequence Alignments

- to Predict Protein Structure and Stability With High Accuracy. Available at: <http://arxiv.org/abs/1207.2484>.
61. Balakrishnan, S., Kamisetty, H., Carbonell, J.G., Lee, S.-I., and Langmead, C.J. (2011). Learning generative models for protein fold families. *Proteins* *79*, 1061–78. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/21268112>.
  62. Friedman, J., Hastie, T., and Tibshirani, R. (2008). Sparse inverse covariance estimation with the graphical lasso. *Biostatistics* *9*, 432–41. Available at: <http://biostatistics.oxfordjournals.org/content/9/3/432.abstract>.
  63. Banerjee, O., El Ghaoui, L., and D'Aspremont, A. (2008). Model Selection Through Sparse Maximum Likelihood Estimation for Multivariate Gaussian or Binary Data. *J. Mach. Learn. Res.* *9*, 485–516. Available at: <http://dl.acm.org/citation.cfm?id=1390681.1390696>.
  64. Baldassi, C., Zamparo, M., Feinauer, C., Procaccini, A., Zecchina, R., Weigt, M., and Pagnani, A. (2014). Fast and accurate multivariate gaussian modeling of protein families: predicting residue contacts and protein-interaction partners. *PLoS One* *9*, e92721. Available at: <http://dx.plos.org/10.1371/journal.pone.0092721>.
  65. Besag, J. (1975). Statistical Analysis of Non-Lattice Data. *Source Stat.* *24*, 179–195. Available at: <http://www.jstor.org/stable/2987782>.
  66. Gidas, B. (1988). Consistency of maximum likelihood and pseudo-likelihood estimators for Gibbs Distributions. *Stoch. Differ. Syst. Stoch. Control Theory Appl.* Available at: <http://www.researchgate.net/publication/244456377/Consistency-of-likelihood-estimators-for-Gibbs-Distributions>.
  67. Feinauer, C., Skwark, M.J., Pagnani, A., and Aurell, E. (2014). Improving contact prediction along three dimensions. *19*. Available at: <http://arxiv.org/abs/1403.0379>.
  68. Zhang, H., Gao, Y., Deng, M., Wang, C., Zhu, J., Li, S.C., Zheng, W.-M., and Bu, D. (2016). Improving residue-residue contact prediction via low-rank and sparse decomposition of residue correlation matrix. *Biochem. Biophys. Res. Commun.* Available at: <http://www.sciencedirect.com/science/article/pii/S0006291X16301838>.
  69. Vendruscolo, M., Kussell, E., and Domany, E. (1997). Recovery of protein structure from contact maps. *Fold. Des.* *2*, 295–306. Available at: [http://dx.doi.org/10.1016/S1359-0278\(97\)00041-2](http://dx.doi.org/10.1016/S1359-0278(97)00041-2).
  70. Li, W., Zhang, Y., and Skolnick, J. (2004). Application of sparse NMR restraints to large-scale protein structure prediction. *Biophys. J.* *87*, 1241–8. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/15298926> <http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Search&term=PMC1304462>.
  71. Yu, X., Wu, X., Bermejo, G.A., Brooks, B.R., and Taraska, J.W. (2013). Accurate high-throughput structure mapping and prediction with transition metal ion FRET. *Structure* *21*, 9–19. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/23273426> <http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Search&term=PMC3700372>.
  72. Kalinin, S., Peulen, T., Sindbert, S., Rothwell, P.J., Berger, S., Restle, T., Goody, R.S., Gohlke, H., and Seidel, C.A.M. (2012). A toolkit and bench-

- mark study for FRET-restrained high-precision structural modeling. *Nat. Methods* *9*, 1218–1225. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/23142871> <http://www.nature.com/doifinder/10.1038/nmeth.2222>.
73. Bowers, P.M., Strauss, C.E., and Baker, D. (2000). De novo protein structure determination using sparse NMR data. *J. Biomol. NMR* *18*, 311–8. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/11200525>.
74. Kolinski, A., and Skolnick, J. (1998). Assembly of protein structure from sparse experimental data: An efficient Monte Carlo model. *Proteins Struct. Funct. Genet.* *32*, 475–494. Available at: <http://doi.wiley.com/10.1002/10.1002/28SICI29109701342819980901> <http://doi.wiley.com/10.1002/29323A43C4753A3AAID-PROT63E3.0.CO>
75. Aszdi, A., Taylor, W.R., and Gradwell, M.J. (1995). Global Fold Determination from a Small Number of Distance Restraints. *J. Mol. Biol.* *251*, 308–326. Available at: <http://www.sciencedirect.com/science/article/pii/S0022283685704366?via%3Dhub>.
76. Wu, S., Szilagyi, A., and Zhang, Y. (2011). Improving protein structure prediction using multiple sequence-based contact predictions. *Structure* *19*, 1182–91. Available at: <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3154634/> [&tool=pmcentrez](http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3154634/)
77. Tress, M.L., and Valencia, A. (2010). Predicted residue-residue contacts can help the scoring of 3D models. *Proteins Struct. Funct. Bioinforma.* *78*, NA–NA. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/20408174> <http://doi.wiley.com/10.1002/prot.22714>.
78. Hopf, T.A., Colwell, L.J., Sheridan, R., Rost, B., Sander, C., and Marks, D.S. (2012). Three-dimensional structures of membrane proteins from genomic sequencing. *Cell* *149*, 1607–21. Available at: <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3641149/>
79. Ovchinnikov, S., Kamisetty, H., and Baker, D. (2014). Robust and accurate prediction of residue-residue interactions across protein interfaces using evolutionary information. *Elife* *3*, e02030. Available at: <http://elifesciences.org/content/3/e02030.abstract>.
80. Hopf, T.A., Schrfe, C.P.I., Rodrigues, J.P.G.L.M., Green, A.G., Sander, C., Bonvin, A.M.J.J., and Marks, D.S. (2014). Sequence co-evolution gives 3D contacts and structures of protein complexes. Available at: <http://arxiv.org/abs/1405.0929>.
81. Hayat, S., Sander, C., Marks, D.S., and Elofsson, A. (2015). All-atom 3D structure prediction of transmembrane  $\beta$ -barrel proteins from sequences. *Proc. Natl. Acad. Sci. U. S. A.* *112*, 5413–8. Available at: <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC4364406/> [&tool=pmcentrez](http://www.ncbi.nlm.nih.gov/pmc/articles/PMC4364406/)
82. Hopf, T.A., Morinaga, S., Ihara, S., Touhara, K., Marks, D.S., and Benton, R. (2015). Amino acid coevolution reveals three-dimensional structure and functional domains of insect odorant receptors. *Nat. Commun.* *6*, 6077. Available at: <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC4364406/> [&tool=pmcentrez](http://www.ncbi.nlm.nih.gov/pmc/articles/PMC4364406/)
83. Raval, A., Piana, S., Eastwood, M.P., and Shaw, D.E. (2015). Assessment of the utility of contact-based restraints in accelerating the prediction of protein structure using molecular dynamics simulations. *Protein Sci.* Available at: <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC4364406/> [&tool=pmcentrez](http://www.ncbi.nlm.nih.gov/pmc/articles/PMC4364406/)

//www.ncbi.nlm.nih.gov/pubmed/26266489.

84. Wang, Y., and Barth, P. (2015). Evolutionary-guided de novo structure prediction of self-associated transmembrane helical proteins with near-atomic accuracy. *Nat. Commun.* *6*, 7196. Available at: <http://www.nature.com/ncomms/2015/150521/ncomms8196/abs/ncomms8196.html>.
85. Ovchinnikov, S., Kinch, L., Park, H., Liao, Y., Pei, J., Kim, D.E., Kamisetty, H., Grishin, N.V., and Baker, D. (2015). Large scale determination of previously unsolved protein structures using evolutionary information. *Elife* *4*, e09248. Available at: <http://elifesciences.org/content/early/2015/09/03/eLife.09248.abstract>.
86. Ovchinnikov, S., Park, H., Varghese, N., Huang, P.-S., Pavlopoulos, G.A., Kim, D.E., Kamisetty, H., Kyripides, N.C., and Baker, D. (2017). Protein structure determination using metagenome sequence data. *Science* *355*, 294–298. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/28104891> <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC5493203/>.
87. Bhattacharya, D., Cao, R., and Cheng, J. (2016). UniCon3D: de novo protein structure prediction using united-residue conformational search via stepwise, probabilistic sampling. *Bioinformatics*, btw316. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/27259540>.
88. Braun, T., Koehler Leman, J., and Lange, O.F. (2015). Combining Evolutionary Information and an Iterative Sampling Strategy for Accurate Protein Structure Prediction. *PLoS Comput. Biol.* *11*, e1004661. Available at: <http://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1004661>.
89. Mabrouk, M., Putz, I., Werner, T., Schneider, M., Neeb, M., Bartels, P., and Brock, O. (2015). RBO Aleph: leveraging novel information sources for protein structure prediction. *Nucleic Acids Res.* *43*, W343–8. Available at: <http://nar.oxfordjournals.org/content/43/W1/W343>.
90. Pietal, M.J., Bujnicki, J.M., and Kozlowski, L.P. (2015). GDFuzz3D: a method for protein 3D structure reconstruction from contact maps, based on a non-Euclidean distance function. *Bioinformatics*, btv390. Available at: <http://bioinformatics.oxfordjournals.org/content/early/2015/07/23/bioinformatics.btv390>.
91. Michel, M., Hayat, S., Skwark, M.J., Sander, C., Marks, D.S., and Elofsson, A. (2014). PconsFold: improved contact predictions improve protein models. *Bioinformatics* *30*, i482–i488. Available at: <http://bioinformatics.oxfordjournals.org/content/30/17/i482.long>.
92. Konopka, B.M., Ciombor, M., Kurczynska, M., and Kotulska, M. (2014). Automated Procedure for Contact-Map-Based Protein Structure Reconstruction. *J. Membr. Biol.* Available at: <http://www.ncbi.nlm.nih.gov/pubmed/24682239>.
93. Kosciolak, T., and Jones, D.T. (2014). De novo structure prediction of globular proteins aided by sequence variation-derived contacts. *PLoS One* *9*, e92197. Available at: <http://dx.plos.org/10.1371/journal.pone.0092197>.
94. Nugent, T., and Jones, D.T. (2012). Accurate de novo structure prediction of large transmembrane protein domains using fragment-assembly and correlated

- mutation analysis. Proc. Natl. Acad. Sci. U. S. A. *109*, E1540–7. Available at: <http://www.pnas.org/content/early/2012/05/21/1120036109>.
95. Kim, D.E., Dimaio, F., Yu-Ruei Wang, R., Song, Y., and Baker, D. (2014). One contact for every twelve residues allows robust and accurate topology-level protein structure modeling. Proteins *82 Suppl 2*, 208–18. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/23900763>.
96. Sathyapriya, R., Duarte, J.M., Stehr, H., Filippis, I., and Lappe, M. (2009). Defining an essence of structure determining residue contacts in proteins. PLoS Comput. Biol. *5*, e1000584. Available at: <http://dx.plos.org/10.1371/journal.pcbi.1000584>.
97. Chen, Y., Ding, F., and Dokholyan, N.V. (2007). Fidelity of the Protein Structure Reconstruction from Inter-Residue Proximity Constraints. J. Phys. Chem. B *111*, 7432–7438. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/17542631> <http://pubs.acs.org/doi/abs/10.1021/jp068963t>.
98. Vassura, M., Margara, L., Di Lena, P., Medri, F., Fariselli, P., and Casadio, R. (2007). Reconstruction of 3D structures from protein contact maps. IEEE/ACM Trans. Comput. Biol. Bioinform. *5*, 357–67. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/18670040>.
99. Adhikari, B., Bhattacharya, D., Cao, R., and Cheng, J. (2017). Assessing Predicted Contacts for Building Protein Three-Dimensional Models. Methods Mol. Biol. *1484*, 115–126. Available at: [http://link.springer.com/10.1007/978-1-4939-6406-2\\\_\\\_9](http://link.springer.com/10.1007/978-1-4939-6406-2\_\_9).
100. Di Lena, P., Vassura, M., Margara, L., Fariselli, P., and Casadio, R. (2009). On the Reconstruction of Three-dimensional Protein Structures from Contact Maps. Algorithms *2*, 76–92. Available at: <http://www.mdpi.com/1999-4893/2/1/76>.
101. Duarte, J.M., Sathyapriya, R., Stehr, H., Filippis, I., and Lappe, M. (2010). Optimal contact definition for reconstruction of contact maps. BMC Bioinformatics *11*, 283. Available at: <http://www.biomedcentral.com/1471-2105/11/283>.
102. Zhang, Y., Kolinski, A., and Skolnick, J. (2003). TOUCHSTONE II: a new approach to ab initio protein structure prediction. Biophys. J. *85*, 1145–64. Available at: <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=1303233\&tool=pmcentrez\&>
103. Wang, S., Li, W., Zhang, R., Liu, S., and Xu, J. (2016). CoinFold: a web server for protein contact prediction and contact-assisted protein folding. Nucleic Acids Res., gkw307. Available at: <http://nar.oxfordjournals.org/content/early/2016/04/25/nar.gkw307.long>.
104. Adhikari, B., Bhattacharya, D., Cao, R., and Cheng, J. (2015). CONFOLD: Residue-residue contact-guided ab initio protein folding. Proteins *83*, 1436–49. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/25974172>.
105. Oliveira, S.H.P. de, Shi, J., and Deane, C.M. (2016). Comparing co-evolution methods and their application to template-free protein structure prediction. Bioinformatics, btw618. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/27682254>.
106. Rodriguez-Rivas, J., Marsili, S., Juan, D., and Valencia, A. (2016).

- Conservation of coevolving protein interfaces bridges prokaryote-eukaryote homologies in the twilight zone. Proc. Natl. Acad. Sci. U. S. A. *113*, 15018–15023. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/27965389> <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC5206536/>.
107. Feinauer, C., Szurmant, H., Weigt, M., and Pagnani, A. (2016). Inter-Protein Sequence Co-Evolution Predicts Known Physical Interactions in Bacterial Ribosomes and the Trp Operon. PLoS One *11*, e0149166. Available at: <http://journals.plos.org/plosone/article?id=10.1371/journal.pone.0149166>.
108. Gueudr, T., Baldassi, C., Zamparo, M., Weigt, M., and Pagnani, A. (2016). Simultaneous identification of specifically interacting paralogs and inter-protein contacts by Direct-Coupling Analysis. 19. Available at: <http://arxiv.org/abs/1605.03745>.
109. Bitbol, A.-F., Dwyer, R.S., Colwell, L.J., and Wingreen, N.S. (2016). Inferring interaction partners from protein sequences. Proc. Natl. Acad. Sci. *113*, 12180–12185. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/27663738> <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC5087060/> <http://www.pnas.org/lookup/doi/10.1073/pnas.1606762113>.
110. Uguzzoni, G., John Lovis, S., Oteri, F., Schug, A., Szurmant, H., and Weigt, M. (2017). Large-scale identification of coevolution signals across homo-oligomeric protein interfaces by direct coupling analysis. Proc. Natl. Acad. Sci. *114*, E2662–E2671. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/28289198> <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC5380090/> <http://www.pnas.org/lookup/doi/10.1073/pnas.1615068114>.
111. Dos Santos, R.N., Morcos, F., Jana, B., Andricopulo, A.D., and Onuchic, J.N. (2015). Dimeric interactions and complex formation using direct coevolutionary couplings. Sci. Rep. *5*, 13652. Available at: <http://www.nature.com/srep/2015/150904/srep13652/full/srep13652.html>.
112. Parisi, G., Zea, D.J., Monzon, A.M., and Marino-Buslje, C. (2015). Conformational diversity and the emergence of sequence signatures during evolution. Curr. Opin. Struct. Biol. *32C*, 58–65. Available at: <http://www.sciencedirect.com/science/article/pii/S0959440X15000147>.
113. Sfriso, P., Duran-Frigola, M., Mosca, R., Emperador, A., Aloy, P., and Orozco, M. (2016). Residues Coevolution Guides the Systematic Identification of Alternative Functional Conformations in Proteins. Structure *24*, 116–126. Available at: <http://linkinghub.elsevier.com/retrieve/pii/S0969212615004657>.
114. Sutto, L., Marsili, S., Valencia, A., and Gervasio, F.L. (2015). From residue coevolution to protein conformational ensembles and functional dynamics. Proc. Natl. Acad. Sci. U. S. A., 1508584112. Available at: <http://www.pnas.org/content/early/2015/10/20/1508584112.abstract>.
115. Jana, B., Morcos, F., and Onuchic, J.N. (2014). From structure to function: the convergence of structure based models and co-evolutionary information. Phys. Chem. Chem. Phys. *16*, 6496. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/24603809> <http://xlink.rsc.org/?DOI=c3cp55275f>.
116. Morcos, F., Jana, B., Hwa, T., and Onuchic, J.N. (2013). Coevolutionary

- signals across protein lineages help capture multiple protein conformations. Proc. Natl. Acad. Sci. U. S. A. 110, 20533–8. Available at: <http://www.pnas.org/content/110/51/20533.short>.
117. Jeon, J., Nam, H.-J., Choi, Y.S., Yang, J.-S., Hwang, J., and Kim, S. (2011). Molecular evolution of protein conformational changes revealed by a network of evolutionarily coupled residues. Mol. Biol. Evol. 28, 2675–85. Available at: <http://mbe.oxfordjournals.org/content/28/9/2675.full>.
118. Woniak, P., Kotulska, M., and Vriend, G. (2017). Correlated mutations distinguish misfolded and properly folded proteins. Bioinformatics 33, 1497–1504.
119. Cao, R., Adhikari, B., Bhattacharya, D., Sun, M., Hou, J., and Cheng, J. (2016). QACon: single model quality assessment using protein structural and contact information with machine learning techniques. Bioinformatics 14, btw694. Available at: <https://academic.oup.com/bioinformatics/article-lookup/doi/10.1093/bioinformatics/btw694>.
120. Terashi, G., Nakamura, Y., Shimoyama, H., and Takeda-Shitaka, M. (2014). Quality Assessment Methods for 3D Protein Structure Models Based on a Residue–Residue Distance Matrix Prediction. Chem. Pharm. Bull. 62, 744–753. Available at: [https://www.jstage.jst.go.jp/article/cpb/62/8/62{\\\_}c13-00973/{\\\_}article](https://www.jstage.jst.go.jp/article/cpb/62/8/62{\_}c13-00973/{\_}article).
121. Tang, Y., Huang, Y.J., Hopf, T.A., Sander, C., Marks, D.S., and Montelione, G.T. (2015). Protein structure determination by combining sparse NMR data with evolutionary couplings. Nat. Methods advance on. Available at: <http://dx.doi.org/10.1038/nmeth.3455>.
122. Ward, A.B., Sali, A., and Wilson, I.A. (2013). Biochemistry. Integrative structural biology. Science 339, 913–5. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/23430643> <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3633482/>.
123. Sadowski, M.I. (2013). Prediction of protein domain boundaries from inverse covariances. Proteins 81, 253–60. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/22987736> <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3563215/>.
124. Nawy, T. (2016). Structural biology: RNA structure from sequence. Nat. Methods 13, 465–465. Available at: <http://dx.doi.org/10.1038/nmeth.3892>.
125. Weinreb, C., Gross, T., Sander, C., and Marks, D.S. (2015). 3D RNA from evolutionary couplings. Available at: <http://arxiv.org/abs/1510.01420>.
126. De Leonidis, E., Lutz, B., Ratz, S., Cocco, S., Monasson, R., Schug, A., and Weigt, M. (2015). Direct-Coupling Analysis of nucleotide coevolution facilitates RNA secondary and tertiary structure prediction. Nucleic Acids Res., gkv932. Available at: <http://nar.oxfordjournals.org/content/early/2015/09/29/nar.gkv932.full>.
127. Hopf, T.A., Ingraham, J.B., Poelwijk, F.J., Schrf, C.P.I., Springer, M., Sander, C., and Marks, D.S. (2017). Mutation effects predicted from sequence co-variation. Nat. Biotechnol. 35, 128–135. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/28092658> <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC553769/> <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC553769/>
128. Wu, N.C., Du, Y., Le, S., Young, A.P., Zhang, T.-H., Wang, Y., Zhou, J., Yoshizawa, J.M., Dong, L., and Li, X. et al. (2016). Coupling high-throughput

- genetics with phylogenetic information reveals an epistatic interaction on the influenza A virus M segment. *BMC Genomics* *17*, 46. Available at: <http://bmcgenomics.biomedcentral.com/articles/10.1186/s12864-015-2358-7>.
129. Figliuzzi, M., Jacquier, H., Schug, A., Tenaillon, O., and Weigt, M. (2015). Coevolutionary landscape inference and the context-dependence of mutations in beta-lactamase TEM-1. *Mol. Biol. Evol.*, msv211. Available at: <http://mbe.oxfordjournals.org/content/early/2015/10/06/molbev.msv211.abstract>.
130. Asti, L., Uguzzoni, G., Marcatili, P., and Pagnani, A. (2016). Maximum-Entropy Models of Sequenced Immune Repertoires Predict Antigen-Antibody Affinity. *PLoS Comput. Biol.* *12*, e1004870. Available at: <http://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1004870>.
131. Elhanati, Y., Murugan, A., Callan, C.G., Mora, T., and Walczak, A.M. (2014). Quantifying selection in immune receptor repertoires. *Proc. Natl. Acad. Sci. U. S. A.* *111*, 9875–80. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/24941953> <http://www.ncbi.nlm.nih.gov/articlerender.fcgi?artid=PMC4103359>.
132. Franceus, J., Verhaeghe, T., and Desmet, T. (2016). Correlated positions in protein evolution and engineering. *J. Ind. Microbiol. Biotechnol.*, 1–9. Available at: <http://link.springer.com/10.1007/s10295-016-1811-1>.
133. Skwark, M.J., Croucher, N.J., Puranen, S., Chewapreecha, C., Pesonen, M., Xu, Y.Y., Turner, P., Harris, S.R., Beres, S.B., and Musser, J.M. *et al.* (2017). Interacting networks of resistance, virulence and core machinery genes identified by genome-wide epistasis analysis. *PLOS Genet.* *13*, e1006508. Available at: <http://dx.plos.org/10.1371/journal.pgen.1006508>.
134. Fox, G., Sievers, F., and Higgins, D.G. (2016). Using de novo protein structure predictions to measure the quality of very large multiple sequence alignments. *Bioinformatics* *32*, 814–20. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/26568625>.
135. Monastyrskyy, B., Fidelis, K., Tramontano, A., and Kryshtafovych, A. (2011). Evaluation of residue-residue contact predictions in CASP9. *Proteins Suppl 1*, 119–25. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/21928322> <http://www.ncbi.nlm.nih.gov/articlerender.fcgi?artid=PMC3226919>.
136. Monastyrskyy, B., D'Andrea, D., Fidelis, K., Tramontano, A., and Kryshtafovych, A. (2014). Evaluation of residue-residue contact prediction in CASP10. *Proteins 82 Suppl 2*, 138–53. Available at: <http://www.ncbi.nlm.nih.gov/articlerender.fcgi?artid=PMC3226919>.
137. Ashkenazy, H., Unger, R., and Kliger, Y. (2009). Optimal data collection for correlated mutation analysis. *Proteins* *74*, 545–55. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/18655065>.
138. Kosciol, T., and Jones, D.T. (2015). Accurate contact predictions using co-evolution techniques and machine learning. *Proteins Struct. Funct. Bioinforma.*, n/a–n/a. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/26205532>.
139. Betts, M.J., and Russell, R.B. Amino Acid Properties and Consequences of Substitutions. In *Bioinforma. genet.* (Chichester, UK: John Wiley & Sons, Ltd), pp. 289–316. Available at: <http://doi.wiley.com/10.1002/0470867302.ch14>.
140. Anishchenko, I., Ovchinnikov, S., Kamisetty, H., and Baker, D. (2017). Ori-

- gins of coevolution between residues distant in protein 3D structures. Proc. Natl. Acad. Sci., 201702664. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/28784799> <http://www.pnas.org/lookup/doi/10.1073/pnas.1702664114>.
141. Marks, D.S., Hopf, T.A., and Sander, C. (2012). Protein structure prediction from sequence variation. Nat. Biotechnol. 30, 1072–1080. Available at: <http://www.nature.com/nbt/journal/v30/n11/full/nbt.2419.html>.
142. Buslje, C.M., Santos, J., Delfino, J.M., and Nielsen, M. (2009). Correction for phylogeny, small number of observations and data redundancy improves the identification of coevolving amino acid pairs using mutual information. Bioinformatics 25, 1125–31. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/19276150> <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC2672635>.
143. The UniProt Consortium (2013). Update on activities at the Universal Protein Resource (UniProt) in 2013. Nucleic Acids Res. 41, D43–7. Available at: <http://nar.oxfordjournals.org/content/41/D1/D43>.
144. Finn, R.D., Coggill, P., Eberhardt, R.Y., Eddy, S.R., Mistry, J., Mitchell, A.L., Potter, S.C., Punta, M., Qureshi, M., and Sangrador-Vegas, A. *et al.* (2016). The Pfam protein families database: towards a more sustainable future. Nucleic Acids Res. 44, D279–D285. Available at: <https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/gkv1344>.
145. Remmert, M., Biegert, A., Hauser, A., and Sding, J. (2012). HHblits: lightning-fast iterative protein sequence searching by HMM-HMM alignment. Nat. Methods 9, 173–5. Available at: <http://dx.doi.org/10.1038/nmeth.1818>.
146. Espada, R., Parra, R.G., Mora, T., Walczak, A.M., and Ferreiro, D. (2015). Capturing coevolutionary signals in repeat proteins. BMC Bioinformatics 16, 207. Available at: <http://arxiv.org/abs/1407.6903>.
147. Toth-Petroczy, A., Palmedo, P., Ingraham, J., Hopf, T.A., Berger, B., Sander, C., Marks, D.S., Alexander, P., He, Y., and Chen, Y. *et al.* (2016). Structured States of Disordered Proteins from Genomic Sequences. Cell 167, 158–170.e12. Available at: <http://linkinghub.elsevier.com/retrieve/pii/S0092867416312430>.
148. Avila-Herrera, A., and Pollard, K.S. (2015). Coevolutionary analyses require phylogenetically deep alignments and better null models to accurately detect inter-protein contacts within and between species. BMC Bioinformatics 16, 268. Available at: <http://www.biomedcentral.com/1471-2105/16/268>.
149. Lee, B.-C., and Kim, D. (2009). A new method for revealing correlated mutations under the structural and functional constraints in proteins. Bioinformatics 25, 2506–13. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/19628501>.
150. Ovchinnikov, S., Kim, D.E., Wang, R.Y.-R., Liu, Y., DiMaio, F., and Baker, D. (2015). Improved de novo structure prediction in CASP11 by incorporating Co-evolution information into rosetta. Proteins. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/26677056>.
151. Noel, J.K., Morcos, F., and Onuchic, J.N. (2016). Sequence co-evolutionary information is a natural partner to minimally-frustrated models of biomolecular dy-

- namics. F1000Research 5. Available at: <http://f1000research.com/articles/5-106/v1>.
152. Coucke, A., Uguzzoni, G., Oteri, F., Cocco, S., Monasson, R., and Weigt, M. (2016). Direct coevolutionary couplings reflect biophysical residue interactions in proteins. *J. Chem. Phys.* 145, 174102. Available at: <http://scitation.aip.org/content/aip/journal/jcp/145/17/10.1063/1.4966156>.
153. Burley, S., and Petsko, G. (1985). Aromatic-aromatic interaction: a mechanism of protein structure stabilization. *Science* (80-. ). 229, 23–28. Available at: <http://www.sciencemag.org/cgi/doi/10.1126/science.3892686>.
154. Ma, J., Wang, S., Wang, Z., and Xu, J. (2015). Protein contact prediction by integrating joint evolutionary coupling analysis and supervised learning. *Bioinformatics*, btv472. Available at: <http://bioinformatics.oxfordjournals.org/content/early/2015/09/04/bioinformatics.btv472>.
155. Ho, T.K. (1998). The random subspace method for constructing decision forests. *IEEE Trans. Pattern Anal. Mach. Intell.* 20, 832–844. Available at: <http://ieeexplore.ieee.org/document/709601/>.
156. Tin Kam Ho (1995). Random decision forests. In Proc. 3rd int. conf. doc. anal. recognit. (IEEE Comput. Soc. Press), pp. 278–282. Available at: <http://ieeexplore.ieee.org/document/598994/>.
157. Breiman, L. (2001). Random Forests. *Mach. Learn.* 45, 5–32. Available at: <http://link.springer.com/10.1023/A:1010933404324>.
158. Menze, B.H., Kelm, B.M., Masuch, R., Himmelreich, U., Bachert, P., Petrich, W., and Hamprecht, F.A. (2009). A comparison of random forest and its Gini importance with standard chemometric methods for the feature selection and classification of spectral data. *BMC Bioinformatics* 10, 213. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/19591666> <http://www.ncbi.nlm.nih.gov/entrez/fetch?artid=PMC2724423>.
159. Louppe, G. (2014). Understanding Random Forests: From Theory to Practice. Available at: <http://arxiv.org/abs/1407.7502>.
160. Strobl, C., Boulesteix, A.-L., Zeileis, A., and Hothorn, T. (2007). Bias in random forest variable importance measures: Illustrations, sources and a solution. *BMC Bioinformatics* 8, 25. Available at: <http://bmcbioinformatics.biomedcentral.com/articles/10.1186/1471-2105-8-25>.
161. Bernard, S., Heutte, L., and Adam, S. (2009). Influence of Hyperparameters on Random Forest Accuracy. In (Springer, Berlin, Heidelberg), pp. 171–180. Available at: [http://link.springer.com/10.1007/978-3-642-02326-2{\\\_}18](http://link.springer.com/10.1007/978-3-642-02326-2{\_}18).
162. Fodor, A.A., and Aldrich, R.W. (2004). Influence of conservation on calculations of amino acid covariance in multiple sequence alignments. *Proteins* 56, 211–21. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/15211506>.
163. Miyazawa, S., and Jernigan, R.L. (1999). Self-consistent estimation of inter-residue protein contact energies based on an equilibrium mixture approximation of residues. *Proteins* 34, 49–68. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/10336383>.
164. Petersen, B., Petersen, T.N., Andersen, P., Nielsen, M., and Lundegaard, C.

- (2009). BMC Structural Biology A generic method for assignment of reliability scores applied to solvent accessibility predictions. *BMC Struct. Biol.* *9*. Available at: <http://www.biomedcentral.com/1472-6807/9/51>.
165. Jones, D.T. (1999). Protein secondary structure prediction based on position-specific scoring matrices 1 1Edited by G. Von Heijne. *J. Mol. Biol.* *292*, 195–202. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/10493868> <http://linkinghub.elsevier.com/retrieve/pii/S0022283699930917>.
166. Sillitoe, I., Lewis, T.E., Cuff, A., Das, S., Ashford, P., Dawson, N.L., Furnham, N., Laskowski, R.A., Lee, D., and Lees, J.G. *et al.* (2015). CATH: comprehensive structural and functional annotations for genome sequences. *Nucleic Acids Res.* *43*, D376–D381. Available at: <https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/gku947>.
167. Bengio, Y. (2012). Practical Recommendations for Gradient-Based Training of Deep Architectures. In *Neural networks: Tricks of the trade* (Springer Berlin Heidelberg), pp. 437–478. Available at: <https://arxiv.org/pdf/1206.5533v2.pdf>.
168. Mahsereci, M., Balles, L., Lassner, C., and Hennig, P. (2017). Early Stopping without a Validation Set. *arXiv*. Available at: <http://arxiv.org/abs/1703.09580>.
169. Carreira-Perpin, M. a, and Hinton, G.E. (2005). On Contrastive Divergence Learning. *Artif. Intell. Stat.* *0*, 17. Available at: <http://learning.cs.toronto.edu/{~}hinton/absps/cdmis.pdf>
170. Bottou, L. (2012). Stochastic Gradient Descent Tricks. In *Neural networks: Tricks of the trade* (Springer, Berlin, Heidelberg), pp. 421–436. Available at: [http://link.springer.com/10.1007/978-3-642-35289-8{\\\_}25](http://link.springer.com/10.1007/978-3-642-35289-8{\_}25).
171. Kingma, D., and Ba, J. (2014). Adam: A Method for Stochastic Optimization. Available at: <http://arxiv.org/abs/1412.6980>.
172. Chollet, F. and others (2015). Keras. Available at: <https://github.com/fchollet/keras>.
173. Dieleman, S., Schlter, J., Raffel, C., Olson, E., Snderby, S.K., Nouri, D., Maturana, D., Thoma, M., Battenberg, E., and Kelly, J. *et al.* (2015). Lasagne: First release. Available at: <https://zenodo.org/record/27878>.
174. Hinton, G.E. (2002). Training Products of Experts by Minimizing Contrastive Divergence. *Neural Comput.* *14*, 1771–1800. Available at: <http://www.gatsby.ucl.ac.uk/publications/tr/tr00-004.pdf>.
175. Bengio, Y., and Delalleau, O. (2009). Justifying and Generalizing Contrastive Divergence. *Neural Comput.* *21*, 1601–21. Available at: <http://www.iro.umontreal.ca/{~}lisa/publications2/index.php/attachments/single/105>.
176. Tieleman, T. (2008). Training Restricted Boltzmann Machines using Approximations to the Likelihood Gradient. *Proc. 25th Int. Conf. Mach. Learn.* *307*, 7.
177. Jensen, C.S., Kjærulff, U., and Kong, A. (1995). Blocking Gibbs sampling in very large probabilistic expert systems. *Int. J. Hum. Comput. Stud.* *42*,

- 647–666. Available at: <http://www.sciencedirect.com/science/article/pii/S1071581985710294>.
178. Robinson, A.B., and Robinson, L.R. (1991). Distribution of glutamine and asparagine residues and their near neighbors in peptides and proteins. *Proc. Natl. Acad. Sci. U. S. A.* *88*, 8880–4. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/1924347> <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC373232/>
179. Atchley, W.R., Zhao, J., Fernandes, A.D., and Drke, T. (2005). Solving the protein sequence metric problem. *Proc. Natl. Acad. Sci. U. S. A.* *102*, 6395–400. Available at: <http://www.pnas.org/content/102/18/6395.abstract>.
180. Grantham, R. (1974). Amino acid difference formula to help explain protein evolution. *Science* *185*, 862–4. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/4843792>
181. Kawashima, S., Pokarowski, P., Pokarowska, M., Kolinski, A., Katayama, T., and Kanehisa, M. (2008). AAindex: amino acid index database, progress report 2008. *Nucleic Acids Res.* *36*, D202–5. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/17998252> <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC2270073/>
182. Zimmerman, J.M., Eliezer, N., and Simha, R. (1968). The characterization of amino acid sequences in proteins by statistical methods. *J. Theor. Biol.* *21*, 170–201. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/5700434>.
183. Wimley, W.C., and White, S.H. (1996). Experimentally determined hydrophobicity scale for proteins at membrane interfaces. *Nat. Struct. Biol.* *3*, 842–8. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/8836100>.
184. Kyte, J., and Doolittle, R.F. (1982). A simple method for displaying the hydrophilicity character of a protein. *J. Mol. Biol.* *157*, 105–132. Available at: <http://www.sciencedirect.com/science/article/pii/0022283682905150>
185. Cornette, J.L., Cease, K.B., Margalit, H., Spouge, J.L., Berzofsky, J.A., and DeLisi, C. (1987). Hydrophobicity scales and computational techniques for detecting amphipathic structures in proteins. *J. Mol. Biol.* *195*, 659–685. Available at: <http://linkinghub.elsevier.com/retrieve/pii/0022283687901896>.
186. Pontius, J., Richelle, J., and Wodak, S.J. (1996). Deviations from Standard Atomic Volumes as a Quality Measure for Protein Crystal Structures. *J. Mol. Biol.* *264*, 121–136. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/8950272> <http://linkinghub.elsevier.com/retrieve/pii/S0022283696906282>.
187. Zhu, H., and Braun, W. (1999). Sequence specificity, statistical potentials, and three-dimensional structure prediction with self-correcting distance geometry calculations of beta-sheet formation in proteins. *Protein Sci.* *8*, 326–42. Available at: <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC129133/> <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC129133/tools=pmcentrez>
188. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., and Dubourg, V. *et al.* (2011). Scikit-learn: Machine Learning in Python. *J. Mach. Learn. Res.* *12*, 2825–2830. Available at: <http://jmlr.csail.mit.edu/papers/v12/pedregosa11a.html>.
189. Livingstone, C.D., and Barton, G.J. (1993). Protein sequence alignments: a strategy for the hierarchical analysis of residue conservation. *Bioinformat-*

ics 9, 745–756. Available at: [http://bioinformatics.oxfordjournals.org/  
content/9/6/745](http://bioinformatics.oxfordjournals.org/content/9/6/745).