

# **Structure Learning for Generative Models of Protein Fold Families**

**Sivaraman Balakrishnan, Hetunandan Kamisetty,  
Jaime G. Carbonell, Christopher James Langmead\***

December 2009  
CMU-CS-09-177

School of Computer Science  
Carnegie Mellon University  
Pittsburgh, PA 15213

\*Department of Computer Science, Carnegie Mellon University, Pittsburgh, PA, USA 15213.  
E-mail: [cjl@cs.cmu.edu](mailto:cjl@cs.cmu.edu)

This research was supported by NSF IIS-0905193 and an award from Microsoft Research to CJL.

**Keywords:** Structure Learning, Generative Models, Probabilistic Graphical Models, Proteins

## Abstract

Statistical models of the amino acid composition of the proteins within a fold family are widely used in science and engineering. Existing techniques for learning probabilistic graphical models from multiple sequence alignments either make strong assumptions about the conditional independencies within the model (e.g., HMMs), or else use sub-optimal algorithms to learn the structure and parameters of the model. We introduce an approach to learning the topological structure *and* parameters of an undirected probabilistic graphical model. The learning algorithm uses block- $L_1$  regularization and solves a *convex* optimization problem, thus guaranteeing a globally *optimal* solution at convergence. The resulting model encodes both the position-specific conservation statistics *and* the correlated mutation statistics between sequential and long-range pairs of residues. Our model is generative, allowing for the design of new proteins that have corresponding statistical properties to those seen in nature. We apply our approach to two widely studied protein families: the WW and the PDZ folds. We demonstrate that our model is able to capture interactions that are important in folding and allostery. Our results additionally indicate that while the network of interactions within a protein is sparse, it is richer than previously believed.



# 1 Introduction

The patterns in the amino acid composition of the proteins within a fold family provide insights into the constraints that govern structure, function, and dynamics. While discriminative models of these patterns have been widely used to predict the structure and/or function of a given sequence, there is growing interest in the use of *generative* models to *design* proteins with a prescribed structure and/or function. This report introduces a new approach to learning the structure (i.e., topology) and the parameters of an undirected graphical model from a given multiple sequence alignment. Previous approaches to learning generative models have largely focused on the parameter estimation problem. The structure learning problem was first considered in [24], followed by some minor extensions [25], all of which are greedy algorithms that provide no guarantees as to the optimality of the resulting model. To address this deficiency, we introduce a principled approach to learning generative models from multiple sequence alignments. Our learning algorithm uses block- $L_1$  regularization and solves a *convex* optimization problem, thus guaranteeing a globally *optimal* solution at convergence. The resulting model encodes both the position-specific conservation statistics *and* the correlated mutation statistics between sequential and long-range pairs of residues. Our model is generative, allowing for the design of new proteins that have corresponding statistical properties to those seen in nature. We apply our approach to two widely studied protein families: the WW and the PDZ folds. We demonstrate that our model is able to capture interactions that are important in folding and allostery. Our results additionally indicate that while the network of interactions within a protein is sparse, it is richer than previously believed. While this report is limited to new more powerful models constructed from MSAs, we have previously shown in [15], that constraints learned from sequence can be effectively integrated with constraints learned from structure by using a probabilistic framework to model both sources of constraints.

## 2 Modeling Domain Families with Markov Random Fields

A protein is a polypeptide chain consisting of one or more mostly-independently evolving components called *domains*. A set of evolutionarily related domains from different proteins is called a family<sup>1</sup>. The domains within a family tend to have similar three dimensional structures and have similar biological functions. Thus, by examining the statistical patterns of sequence conservation and diversity within a domain family, we can gain insights into the constraints that determine structure and function. In what follows, we describe an approach to learn these statistical patterns from a given multiple sequence alignment. The resulting model is a probability distribution over amino acid sequences for a particular domain family.

Let  $X_i$  be the multinomial random variable representing the amino-acid composition at position  $i$  of the MSA of the domain family taking values in  $\{1\dots k\}$  where the number of states,  $k$ , is 21 (20 amino acids with one additional state corresponding to a gap). Let  $\mathbf{X} = \{X_1, X_2, \dots, X_p\}$  be the multi-variate random variable describing the amino acid composition of a MSA of length  $p$ . Our goal is to model  $P(\mathbf{X})$ , that is the amino-acid composition of the domain family.

---

<sup>1</sup>In this paper, the expression *domain family* is synonymous with *protein family*.

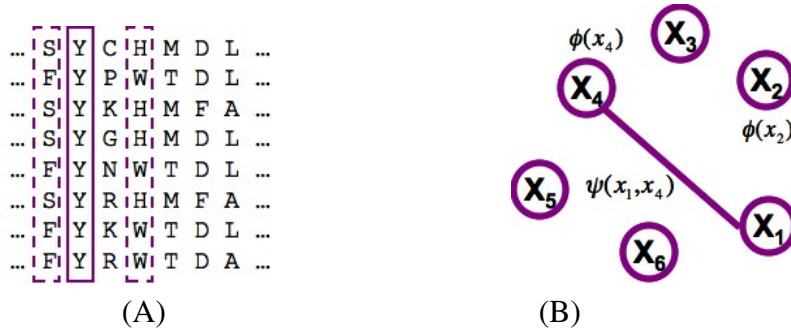


Figure 1: (A) A multiple sequence alignment (MSA) for a hypothetical domain family. (B) A portion of a Markov Random Field encoding the conservation in and the coupling in the MSA. The edge between random variables  $X_1$  and  $X_4$  reflects the coupling between positions 1 and 4 in the MSA.

Unfortunately,  $P(\mathbf{X})$  is a distribution over a space of size  $k^p$ , rendering the explicit modeling of the joint distribution computationally intractable for naturally occurring domains. However, by exploiting the properties of the distribution, one can significantly decrease the number of parameters required to represent this distribution.

To see the kinds of properties that we can exploit, let us consider a toy domain family represented by a MSA as shown in Fig. 1-(A). A close examination of the MSA reveals the following statistical properties of its composition: (i) the Tyrosine ('Y') at position 2 is conserved across the family; (ii) positions 1 and 4 are co-evolving – sequences with a (S) at position 1 have a Histidine (H) at position 4, while sequences with a Phenylalanine (F) at position 1 have a Tryptophan (W) at position 4; (iii) the remaining positions appear to evolve independent of each other. In probabilistic terms we say that  $X_1, X_3$  are co-varying, and that the remaining  $X_i$ 's are statistically independent. We can therefore encode the joint distribution over all positions in the MSA by storing one joint distribution  $P(X_1, X_4)$ , and the uni-variate distributions  $P(X_i)$ , for the remaining positions (since they are all statistically independent of every other variables). The ability to factor the full joint distribution,  $P(\mathbf{X})$ , in this fashion has an important consequence in terms of space complexity. Namely, we can reduce the space requirements from  $21^7$  to  $21^2 + 5 * 21$  parameters. This drastic reduction in space complexity translates to a corresponding reduction in time complexity for computations over the distribution. While this simple example utilizes independencies in the distribution; this kind of reduction is possible in the more general case of *conditional independencies*. A Probabilistic Graphical Model (PGM) exploits these (conditional) independence properties to store the joint probability distribution using a small number of parameters.

Intuitively, a PGM stores the joint distribution of a multivariate random variable over a graph; while any distribution can be modeled by a PGM with a complete graph, exploiting the conditional independencies in the distribution leads to a PGM with a (structurally) sparse graph. Following [24], we use a specific type of probabilistic graphical model called a Markov Random Field (MRF). In its commonly defined form with pair-wise log-linear potentials, a Markov Random Field (MRF) can be formally defined as a tuple  $\mathcal{M} = (\mathbf{X}, \mathcal{E}, \Phi, \Psi)$  where  $(\mathbf{X}, \mathcal{E})$  is an undirected graph over the

random variables, and  $\Phi, \Psi$  are a set of node and edge potentials, respectively, usually chosen to be log-linear functions of the form:

$$\phi_s = (e^{v_1^s} e^{v_2^s} \dots e^{v_k^s}); \quad \psi_{st} = \left\{ \begin{array}{l} e^{w_{11}^{st}} e^{w_{12}^{st}} \dots e^{w_{1k}^{st}} \\ e^{w_{21}^{st}} e^{w_{22}^{st}} \dots e^{w_{2k}^{st}} \\ \dots \\ e^{w_{k1}^{st}} e^{w_{k2}^{st}} \dots e^{w_{kk}^{st}} \end{array} \right\} \quad (1)$$

where  $\mathbf{v} = \{v_s | s = 1 \dots p\}$  and  $\mathbf{w} = \{w^{st} | (s, t) \in \mathcal{E}\}$  are node and edge “weights”, and  $k$  is the number of states that each random variable can take.

The probability of a particular sequence  $x = \{x_1, x_2, \dots, x_p\}$  according to  $\mathcal{M}$  is defined as:

$$P_{\mathcal{M}}(X) = \frac{1}{Z} \prod_{s \in V} \phi_s(X_s) \prod_{(s,t) \in E} \psi_{st}(X_s, X_t) \quad (2)$$

where  $Z$ , the so-called partition function, is a normalizing constant defined as a sum over all possible assignments to  $\mathbf{X}$ .

$$Z = \sum_{X \in \mathbf{X}} \prod_{s \in V} \phi_s(X_s) \prod_{(s,t) \in E} \psi_{st}(X_s, X_t) \quad (3)$$

The structure of the MRF for the MSA shown in Fig. 1(A) is shown in Fig. 1(B). For expositional clarity, only four nodes of the MRF are shown, corresponding to positions 1-4 in the MSA. The edge between variables  $X_1$  and  $X_4$  reflects the statistical coupling between those positions in the MSA.

### 3 Structure learning with $L_1$ Regularization

In the previous section we outlined how an MRF can parsimoniously model the probability distribution  $P(\mathbf{X})$ . In this section we consider the problem of *learning* the MRF from an MSA. This problem can be divided into two parts: (i) *structure learning* — learning the edges of the graph, and (ii) *parameter estimation* — learning  $\mathbf{v}, \mathbf{w}$  (since they completely define the potentials  $\Phi, \Psi$ ), given the structure of the graph.

Due to its importance and applicability in a broad spectrum of areas, the problem of structure learning for graphical models has received considerable attention from several communities. Broadly, the previously considered approaches to this problem are either constraint based [24, 23] or score based [7]. Constraint based methods estimate conditional independencies from data using hypothesis testing and then determine a graph that represents these independencies. Score based approaches combine a metric to measure goodness of fit with a metric to measure complexity of the graph to *score* each graph. This is combined with a (typically greedy) search procedure that generates candidate graphs. However, since the number of possible graphs is super exponential in the number of vertices the search problem is in general NP hard.

More recently several authors [29, 16, 14, 20] have considered convex approximations to the complexity metric and tractable (convex) approximations to the goodness of fit metric. Of these,

those based on  $L_1$  regularization are the most interesting because of their strong theoretical guarantees (consistency in both parameters and structure, i.e. as the number of samples increases we are guaranteed to find the true model, and high statistical efficiency, i.e. the number of samples needed to achieve this guarantee is small). See [28] for a recent review of L1-regularization. We use a similar convex optimization based approach for both structure learning and parameter estimation. To that end, we first describe a suitable objective function for the problem.

The log-likelihood of the parameters  $\Theta = (\mathcal{E}, \mathbf{v}, \mathbf{w})$ , given a set of sequences  $\mathcal{X} = \{\mathbf{X}^1, \mathbf{X}^2, \mathbf{X}^3, \dots, \mathbf{X}^n\}$ , is

$$ll(\Theta) = \frac{1}{n} \sum_{X^i \in \mathcal{X}} \left[ \sum_{s \in V} \log \phi_s(X_s^i) + \sum_{(s,t) \in E} \log \psi_{st}(X_s^i, X_t^i) \right] - \log Z \quad (4)$$

where the term in the braces is the unnormalized likelihood of each sequence, and  $Z$  is the global partition function. The problem of learning the structure *and* parameters of the MRF is now simply that of maximizing  $ll(\Theta)$ . To avoid over-fitting and learning densely connected structures, we need to regularize the log-likelihood. In what follows we describe a method to learn sparse structures by optimizing the pseudo-likelihood using block- $L_1$  regularizers.

The general regularized structure learning problem can be formulated as:

$$\max_{\theta} ll(\theta) - R(\theta) \quad (5)$$

For the specific case of block- $L_1$  regularization,  $R(\theta)$  usually takes the form:

$$R(\theta) = \lambda_{node} \|\mathbf{v}\|_2 + \lambda_{edge} \sum_{1 \leq s < t \leq p} \|\mathbf{w}^{st}\|_q \quad (6)$$

where  $\lambda_{node}$  and  $\lambda_{edge}$  are regularization parameters that determine how strongly we penalize higher (absolute) weights. The value of  $\lambda_{node}$  and  $\lambda_{edge}$  control the trade-off between the log-likelihood term and the regularization term in our objective function.

The regularization described above groups all the parameters that describe an edge together in a *block*. The second term in Eq. 6 is the sum of the norms of each block. The choice of norm is usually selected from  $q \in \{1, 2, \infty\}$ . Since norms are always positive, this is exactly equivalent to penalizing the  $L_1$  norm of the vector of norms of each block with the penalty increasing with higher values of  $\lambda_{edge}$ .

The choice of norm affects the nature of the sparsity. Using  $q = 1$  is equivalent to penalizing the likelihood by the sum of the  $L_1$  norms of the individual parameters. That is, using  $q = 1$  encourages sparsity in the parameters. In contrast, using  $q = \{2, \infty\}$  encourages structural sparsity (sparsity in the edges).

In the following sections we present a method to tractably compute the objective, followed by a method to optimize it.

### 3.1 Pseudo Likelihood

The log-likelihood as defined in Eq. 4 is smooth, differentiable, and concave. However, maximizing the log-likelihood requires computing the global partition function  $Z$  and its derivatives,



which in general can take upto  $\mathcal{O}(k^p)$  time. While approximations to the partition function based on Loopy Belief Propagation [16] have been proposed as an alternative, such approximations can lead to inconsistent estimates.

Instead of approximating the true-likelihood using approximate inference techniques, we use a different approximation based on a pseudo-likelihood proposed by [5], and used in [29, 20]. The pseudo-likelihood is defined as:

$$\begin{aligned} pll(\Theta) &= \frac{1}{n} \sum_{X^i \in \mathcal{X}} \sum_{j=1}^p \log(P(X_j^i | X_{-j}^i)) \\ &= \frac{1}{n} \sum_{X^i \in \mathcal{X}} \sum_{j=1}^p \left[ \log \phi_j(X_j^i) + \sum_{k \in V_j'} \log \psi(X_j^i, X_k^i) - Z_j \right] \end{aligned}$$

where  $X_j^i$  is the residue at the  $j^{th}$  position in the  $i^{th}$  sequence of our MSA,  $X_{-j}^i$  denotes the ‘‘Markov blanket’’ of  $X_j^i$ , and  $Z_j$  is a local normalization constant for each node in the MRF. The set  $V_j'$  is the set of all vertices which connect to vertex  $j$  in the PGM. The replacement of a global partition function with local partition functions (which are sums over possible assignments to single nodes rather than a sum over all assignments to *all* nodes of the sequence) makes the pseudo-likelihood easier to compute than the true likelihood.

The pseudo-likelihood retains the concavity of the original problem, and so this approximation makes the problem tractable. Moreover, this approximation is known to yield a consistent estimate of the parameters. That is, as the number of samples increases, parameter estimates using pseudo-likelihood converge to the parameter values using true likelihood.

## 3.2 Optimizing L1-regularized Pseudo-Likelihood

In the previous two sections we described an objective function, and then a tractable and consistent approximation to it, given a set of weights (equivalently, potentials). However, to solve this problem we still need to be able to find the set of weights that maximizes the likelihood under the block-regularization form of Eq. 5. We note that the objective function associated with block- $L_1$  regularization is no longer smooth. In particular, its derivative with respect to any parameter is discontinuous at the point where the group containing the parameter is 0. We therefore consider an equivalent formulation where the non-differentiable part of the objective is converted into a constraint making the new objective function differentiable.

$$\begin{aligned} &\max_{\theta, \alpha} \ell(\theta) - \lambda_{node} \|\mathbf{v}\|_2 - \sum_{1 \leq s < t \leq p} \alpha_{st} \\ \text{subject to:} &\quad \forall (1 \leq s < t \leq p) : \alpha_{st} \geq \|w^{st}\|_q \end{aligned}$$

where the constraints hold with equality at the optimal  $(\theta, \alpha)$ .

One way to solve this reformulation is through a two stage procedure involving the use of projected gradients. In the first stage, we ignore the constraints, compute the gradient of the objective,

and then take a step in this direction. If the step results in any of the constraints being violated we solve an alternative (and simpler) Euclidean projection problem:

$$\begin{aligned} & \min_{\theta', \alpha'} \left\| \begin{bmatrix} \theta' \\ \alpha' \end{bmatrix} - \begin{bmatrix} \theta \\ \alpha \end{bmatrix} \right\|_2 \\ \text{subject to: } & \forall (1 \leq s < t \leq p) : \alpha_{st} \geq \|w^{st}\|_q \end{aligned}$$

which finds the closest parameter vector to the vector obtained by taking the gradient step (by minimizing the Euclidean distance), while satisfying the original constraints. This problem can be solved efficiently for block- $L_1$  norms using Spectral Projected Gradients (SPG), as shown in [20]. Thus, we used the algorithm from [20] to solve this problem in our experiments. Methods based on projected gradients are guaranteed to converge to a stationary point [6], and convexity ensures that this stationary point is globally optimal.

## 4 Related Work

The study of co-evolving residues in proteins has been a problem of key interest due to its wide utility. Much of the early work focused on detecting such pairs in order to predict contacts in a protein in the absence of a solved structure [1, 13] and to perform fold recognition. The pioneering work of [17] used an approach to determine probabilistic dependencies they call SCA and observed that analyzing such patterns could provide insights into the allosteric behavior of the proteins and be used to design new sequences[22]. Others have since developed similar methods [10, 11]. By focusing on co-variation or probabilistic *dependencies* between residues, such methods conflate direct and indirect influences and can lead to incorrect estimates. In contrast, [24] developed an algorithm that determine conditional independencies to learn a Markov Random Field over sequences. Their constraint-based algorithm proceeds by determining conditional independencies and adding edges in a greedy fashion. However, the algorithm can provide no guarantees on the correctness of the networks it learns. They then extended this approach to incorporate interaction data to learn models over pairs of interacting proteins [25] and also develop a sampling algorithm for protein design using such models[26]. More recently, [30] use a similar approach to determine residue contacts at a protein-protein interface. Their method uses a gradient descent approach using Loopy Belief Propagation to approximate likelihoods. Also, their algorithm does not regularize the model and can therefore be prone to over-fitting. In contrast, we use a Pseudo-Likelihood as our objective function thereby avoiding problems of convergence that Loopy BP based methods can face and regularize the model using block regularization to prevent over-fitting.

Block regularization is most similar in spirit to the group Lasso [31] and the multi-task Lasso [2]. Lasso [27] is the problem of finding a linear predictor, by minimizing the squared loss of the predictor with an  $L_1$  penalty. It is well known that the shrinkage properties of the  $L_1$  penalty lead to sparse predictors. The group Lasso extends this idea by grouping the weights of some features of the predictor using an  $L_2$  norm, [31] show that this leads to sparse selection of groups. The multi-task Lasso solves the problem of multiple separate (but similar) regression problems by grouping

the weight of a single feature across the multiple tasks. Intuitively, we solve a problem similar to a group Lasso, replacing the squared loss with an approximation to the negative log-likelihood, where we group all the feature weights of an edge in an undirected graphical model. Thus, sparse selection of groups gives our graphs the property of structural sparsity.

[16] consider learning with only an  $L_1$  penalty (and not a block- $L_1$  penalty), biasing the model to include fewer features, but not directly biasing the model towards sparse structures. They also use a different approximation to the likelihood term, using Loopy Belief Propagation. [20] apply block-regularized structure learning to the problem of detecting abnormalities in heart motion. Particularly they develop an efficient algorithm for tractably solving the convex structure learning problem, based on projected gradients. We use their algorithm in this paper.

## 5 Results

Given the probabilistic framework defined in Sec. 2, and the optimization objectives and algorithms defined in Sec. 3, we are now in a position to learn a graphical model given the sequence record of a protein family. The optimization framework has two major parameters that can be varied: the norm of the block-regularizer ( $L_1, L_2, L_\infty$ ) and the penalty parameters ( $\lambda_v, \lambda_e$ ). To understand the effects of these parameters, we first evaluated our method on artificial protein families whose sequence records were generated from known, randomly generated models. This let us evaluate the success of the various components of our framework in a controlled setting where the ground truth was known.

We then present our results on two real protein families, the PDZ and WW.

### 5.1 Simulations

We generated 32-node graphs, where each edge was including with probability  $\rho$ . Ten different values of  $\rho$  varying from 0.01 and 0.45 were used; for each value of  $\rho$ , twenty different graphs were generated resulting in a total of 200 graphs. For each edge that was included in a graph, edge and node weights were drawn from a Normal distribution (weights  $\sim \mathcal{N}(0,1)$ ). For each of these 200 graphical models, we then sampled 1000 sequences using a Gibbs sampler with a burn-in of 10,000 samples and discarding 1,000 samples between each accepted sequence. These 1000 sequences were then partitioned into two sets: a training set containing 500 sequences and a held-out set of 500 sequences used to test the model. The training set was then used to train a model using each of the three block regularization norms.

We first test our accuracy on structure learning. Since the structure of the model directly depends only on the regularization weight on the edges, the structures were learnt for each norm and each training set with different values of  $\lambda_e$  (between 1 and 500), keeping  $\lambda_v$  fixed at 1.

Fig. 2 shows our performance in predicting the true structure by using  $L_1-L_2$  (Fig. 2-A),  $L_1-L_\infty$  (Fig. 2-B), and  $L_1$  (Fig. 2-C). The accuracy is measured using the F-score (the harmonic mean of precision and recall) of the edge set. We observe that for all settings of  $\rho$  each of the block regularizers learn fairly accurate graphs at some value of  $\lambda_e$ . We see that  $L_1 - L_\infty$  requires higher regularization than the other norms to achieve comparable F-scores. This is because the  $L_\infty$  norm

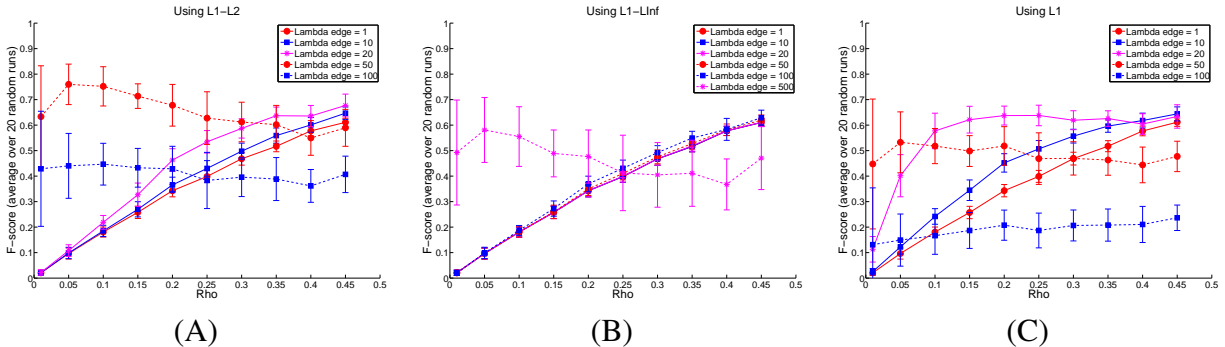


Figure 2: F-scores of structures learnt by using (A)  $L_1-L_2$  norm, (B)  $L_1-L_{\infty}$  norm and (C)  $L_1$  norm. Each figure shows the average and standard deviation of the F-score across 20 different graphs as a function of  $\rho$ , the probability of edge-occurrence.

of a vector is strictly less than its  $L_2$  and  $L_1$  norms. Fig. 3-A shows the global comparison across the three regularizers as a function of  $\rho$ , using in each case, the best model learnt across the different values of  $\lambda_e$ . Figure 3-A also compares our structure learning method with the algorithm in [24]. We evaluate their method over a wide range of parameter settings and select the best model. Figure 3-A shows that our methods significantly out-perform their method for *all* values of  $\rho$ . We see that over all settings our best model has an average F-score of *at least* 0.6. Thus, we are able to infer fairly accurate structures given the proper choice of settings.

Figure 3-B, shows the error in our parameter estimates given the true graph as a function of  $\rho$ . We also find that parameter estimation is reasonably robust to the choice of the regularization weights, as long as the regularization weights are non-zero.

## 5.2 Evaluating Structure and Parameters Jointly

In a simulated setting, structure and parameter estimates can be compared against known ground truth. However, for real domain families we need other evaluation methods. We evaluate the structure and parameters for real domain families by measuring the imputation error of the learnt models. Informally, the imputation error measures the probability of *not* being able to “generate” a complete sequence, given an incomplete one. The imputation error of a column is measured by erasing it in the test MSA, and then computing the probability that the true (known) residues would be predicted by the learnt model. This probability is calculated by performing inference on the erased columns, conditioned on the rest of the MSA. The imputation error of a model is the average of its imputation error over columns.

Using imputation error directly for model selection generally gives us models that are too dense. Intuitively, once we have identified the true model, adding extra edges decreases the imputation error by a very small amount, probably a reflection of the finite-sample bias. On the other hand, we note that there is a distinct “knee” in the graphs of the number of edges versus the imputation error (see Fig. 5 and Fig. 6). We believe the true model would be in the vicinity of this knee.

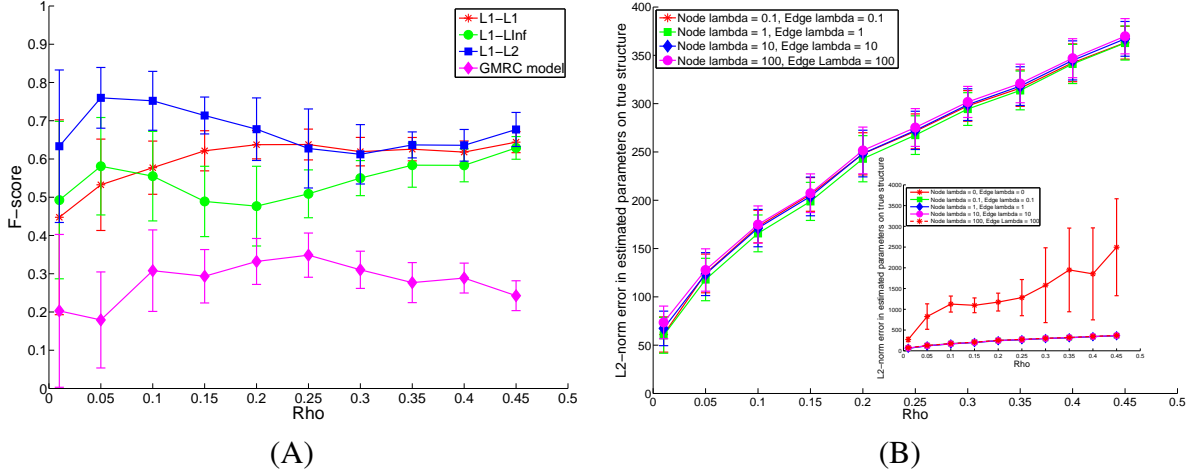


Figure 3: (A) Edge occurrence probability  $\rho$  versus F-score for the structure learning methods we propose, and the method proposed in [24]. (B)  $L_2$  norm of the error in the estimated parameters as a function of the weight of the regularization in stage two.

### 5.3 A generative model for the WW domain

The WW domain family (Pfam id: PF00397 [3]) is a small protein interaction module with two highly conserved tryptophans that adopts a curved three-stranded  $\beta$ -sheet structure with a binding site for proline-containing peptides. In [22] and [19], the authors determine, using Statistical Coupling Analysis (SCA), that the residues can be divided into two clusters: the first cluster contains a set of 8 strongly coupled residues (highlighted in yellow in Fig. 4), and the second cluster contains everything else. Based on this finding, the authors then designed 44 sequences that satisfy co-evolution constraints of the first cluster, of which 12 actually fold *in vivo*. An alternative set of control sequences, which did not satisfy the constraints, failed to fold.

We first constructed a MSA by starting with the PFAM alignment and removing sequences to construct a non-redundant alignment (no pair of sequences was greater than 90% similar). This resulted in an MSA with 700 sequences of which two thirds were used as a training set and the rest were used as a test set. The training set was used to learn models for each of the three norms, using multiple settings of  $\lambda_e$ . Given the structure of the graph, parameters were learned using  $\lambda_v = 1, \lambda_e = 1$ .

Fig. 5-A shows the imputation error of each of the learnt models on the test set. As can be seen, the first few edges contribute to a significant decrease (the first 20 edges contribute to a third of the total decrease in imputation error). This is consistent with [22, 24] who find a small set of vertices and edges to be important. However, the knee of the curve is much further down, at 122 edges. This indicates that the actual pattern of interactions in a domain family, while sparse, is richer than previously thought. In the figure we compare our results to the GMRC method of [24], and to a method that adds edges in the order of their statistical coupling ( $\Delta\Delta G^{stat}$ ) that is also used by the SCA method. We notice that our imputation errors are considerably lower to the methods we compare to.

To see which residues are affected by these edges, we performed the following experiment. We constructed a shuffled MSA by taking the natural MSA and randomly permuting the amino acids within the same column for each column. The new MSA now contains no co-evolving residues but has the same conservation profile as the original MSA. We then computed a coupling profile for the domain by calculating the difference in the imputation errors between sequences in the test set and the shuffled MSA. Intuitively, having a high imputation error difference means that the position was strongly constrained in the original MSA. Fig. 4 shows the results of this analysis; we identify 15 positions in the MSA including all 8 positions previously identified by [19].

In addition to this analysis we also performed a retrospective analysis of the artificial sequences designed by [19]. We attempt to distinguish sequences that folded from those that didn't. To make a fair comparison we select a model of comparable sparsity to that in [19] (shown with a black circle in Fig. 5). Although this is a discriminative (folded or not) test of a generative model we achieve a high AUC of 0.843 (the ROC curve is shown in the Fig. 7). We therefore postulate that the additional constraints that we identified are indeed critical to the stability of the WW fold. The AUC is comparable to the published results of [24]. However, their model has a much higher imputation error as shown in Fig. 5.

## 5.4 Allosteric regulation in the PDZ domain

The PDZ domain is a family of small, evolutionarily well represented protein binding motifs. The domain is most commonly found in signaling proteins and helps to anchor trans-membrane proteins to the cytoskeleton and hold together signaling complexes. The PDZ domain is also interesting because it is considered an *allosteric* protein. The domain, and its members have been studied extensively, in multiple studies, using a wide range of techniques ranging from computational approaches based on statistical coupling ([17]) and Molecular Dynamics simulations [8], to NMR based experimental studies ([12]).

We use the MSA from [17]. We chose a random sub-sample with two-thirds of the sequences as the training set and use the rest as a test set. Using this training set, we learnt generative models for each of the block regularizers, with multiple settings of  $\lambda_e$  in each case and then computed the imputation error on the test set (shown in Fig. 6-A). We selected the model that formed the knee on the curve with the best imputation error. This model had around 700 edges and is shown in full in Fig. 8. Fig. 6-B shows a subset of the edges colored according to their strength (red strongest; blue weakest).

The SCA based approach of [17] identified a set of residues that were coupled to a residue near the active site (HIS-70) including a residue at a distal site on the other end of the protein (GLY-49 in this case). Since the SCA approach can only determine the presence of a dependence but cannot distinguish between direct and indirect couplings, only a cluster of residues was identified. Our model also identifies this interaction, but more importantly, it determines that this interaction is mediated by ALA-74 with position 74 *directly* interacting with both these positions. By providing such a list of sparse interactions our model can provide a small list of hypotheses to an experimentalist looking for possible mechanisms of such allosteric behavior.

In addition to the pathway between HIS-70 and GLY-49, we also identify residues not on the pathway that are connected to other parts of the protein including, for example ASN-61 of the

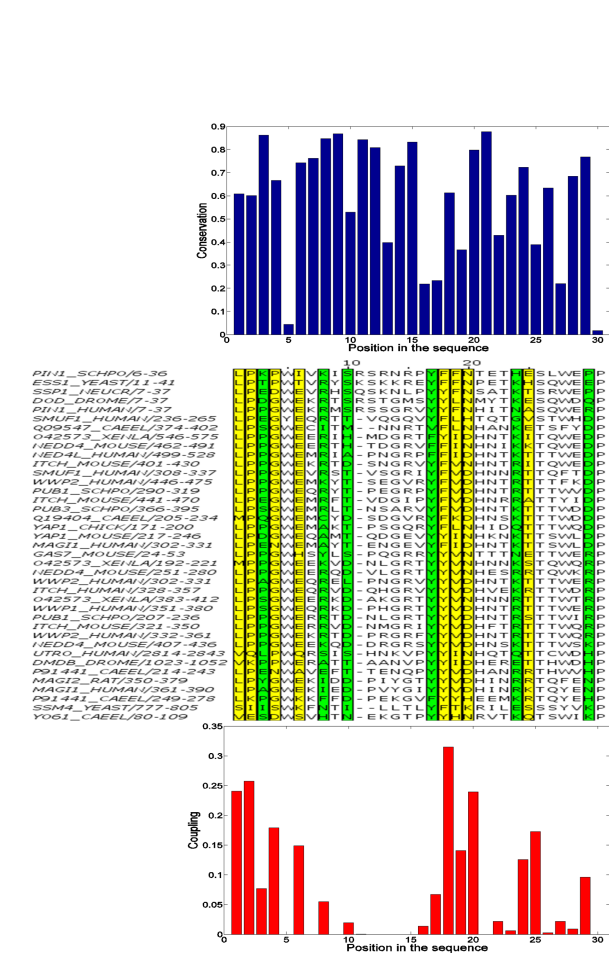


Figure 4: Part of the WW MSA we used. In yellow are positions identified by [19] as being critical to folding. Positions we additionally identify are in green. The conservation profile (top) shows the entropy (scaled) at each position. The coupling profile is shown below the MSA.

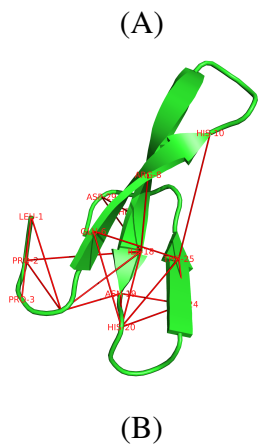
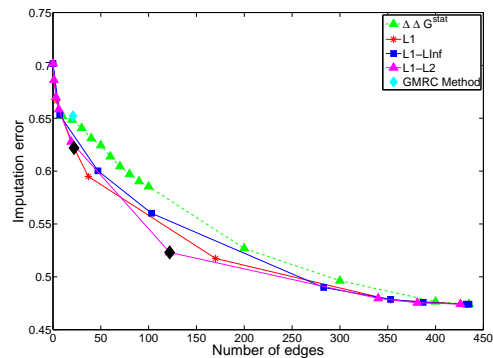


Figure 5: (A) Number of edges versus imputation error for the WW domain. The model at the knee was a model that minimized the  $L_1-L_2$  norm and had 122 edges (shown with a black diamond). Comparisons to the GMRC method and a method based on statistical coupling ( $\Delta\Delta G^{stat}$ ) are shown. (B) The edges of the model used in the discriminative task, overlaid on the structure of the WW domain of a ubiquitin protein ligase (PDB id: 1I5H)[4]

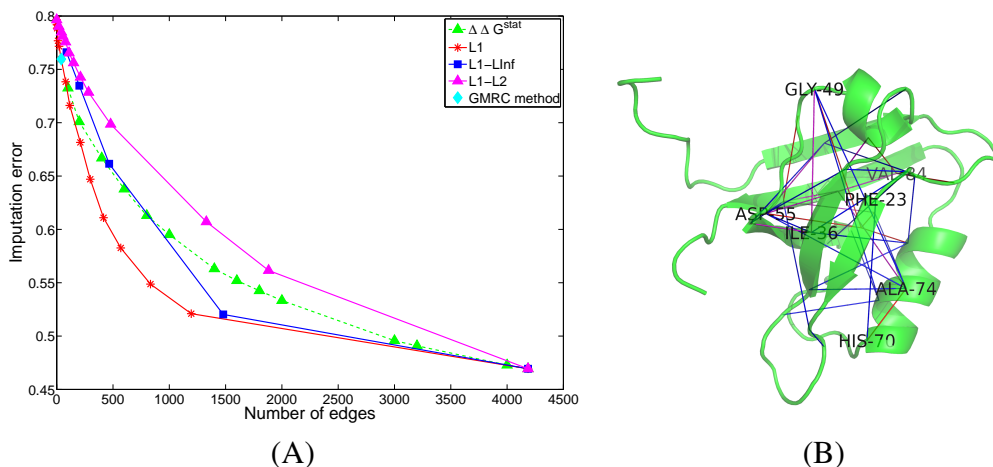


Figure 6: (A) Number of edges versus imputation error for the PDZ domain. Comparisons to the GMRC method and a method based on statistical coupling ( $\Delta \Delta G^{stat}$ ) are shown. (B) Edges learnt from our models overlaid on the structure of PDZ domain of PSD-95(PDB id:1BE9). Edge colors indicate the strength of the coupling (red being the strongest, and blue being the weakest)

protein. This position is connected to ALA-88 and VAL-60 in our model, and does not appear in the network suggested by [17], but has been implicated by NMR experiments [12] as being dynamically linked to the active site. Thus, our method appears to capture a richer set of interactions than that are possible using SCA.

## 6 Discussion and Future Work

In this paper we have proposed a statistical sequence-based approach to modeling the evolutionary pressures on a protein family. Overall, we find that by employing sound probabilistic modeling and convex structure (and parameter) learning, we are able to find a good balance between structural sparsity (simplicity) and goodness of fit. We demonstrate the utility of our method in identifying constraints useful both in protein design and in furthering our understanding of protein function and regulation.

One limitation associated with a sequence-only approach to learning a statistical model for a domain family is that the correlations observed in the MSA can be inflated due to phylogeny [18, 9]. There are a number of ways to incorporate phylogenetic information into our model. For example, given a phylogenetic clustering of sequences, we can incorporate a single additional node in the graphical model reflecting the cluster to which the sequence belongs. This would allow us to distinguish functional coupling from coupling caused due to phylogenetic variations.

Designing proteins from a generative sequence based model such as ours could be greatly enhanced by incorporating structure based information which explicitly models the physical constraints of the protein. Such information could easily be incorporated either through the use of informative priors (e.g., interaction energies, etc), or by the addition of edge features.



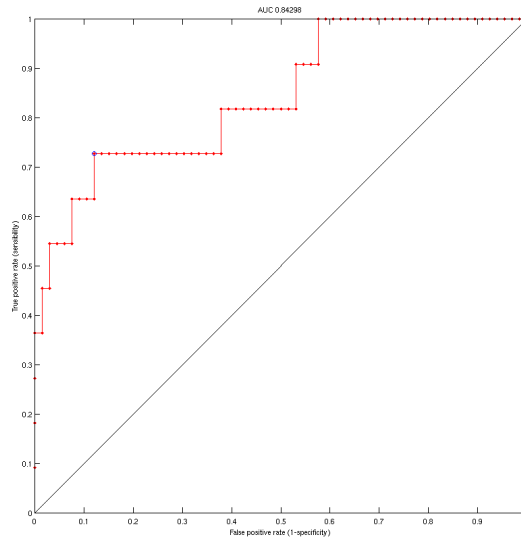


Figure 7: ROC curve of our model for the task of distinguishing artificial WW sequences that fold from those that don't. All sequences and their labels (folded in vivo or not) are from [19]

Recently, [21] proposed a new method specifically to optimize costly functions, where the projection step is cheap, by using a quasi-Newton algorithm which uses local curvature of the objective to approximate its second derivative. Typically, this leads to much faster convergence. We expect this to be applicable to our method.

## Acknowledgments

This research was supported by NSF IIS-0905193 and an award from Microsoft Research to CJL.

## References

- [1] D. Altschuh, T. Vernet, P. Berti, D. Moras, and K. Nagai. Coordinated amino acid changes in homologous protein families. *Protein Eng.*, 2(3):193–199, September 1988.
- [2] Andreas Argyriou, Theodoros Evgeniou, and Massimiliano Pontil. Convex multi-task feature learning. In *Machine Learning*. press, 2007.
- [3] A. Bateman, E. Birney, L. Cerruti, R. Durbin, L. Etwiller, S.R. Eddy, S. Griffiths-Jones, K.L. Howe, M. Marshall, and E.L.L. Sonnhammer. The Pfam protein families database. *Nucleic acids research*, 30(1):276, 2002.

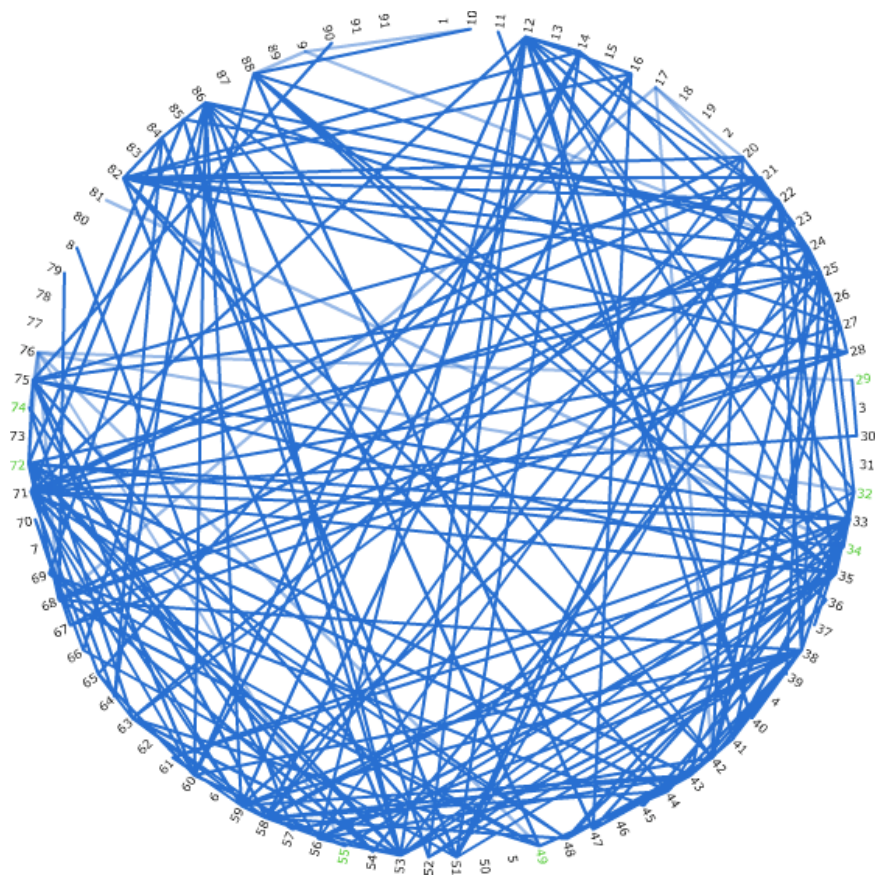


Figure 8: Graph showing all edges identified for the PDZ domain

- [4] H.M. Berman, J. Westbrook, Z. Feng, G. Gilliland, T.N. Bhat, H. Weissig, I.N. Shindyalov, and P.E. Bourne. The protein data bank. *Nucl. Acids Res.*, 28:235–242, 2000.
- [5] J. Besag. Efficiency of pseudolikelihood estimation for simple Gaussian fields. *Biometrika*, 64(3):616–618, 1977.
- [6] Stephen Boyd and Lieven Vandenberghe. *Convex Optimization*. Cambridge University Press, March 2004.
- [7] David Maxwell Chickering and Craig Boutilier. Optimal structure identification with greedy search. *Journal of Machine Learning Research*, 3:507–554, 2002.
- [8] Anne Dhulesia, Joerg Gsponer, and Michele Vendruscolo. Mapping of two networks of residues that exhibit structural and dynamical changes upon binding in a pdz domain protein. *Journal of the American Chemical Society*, 130(28):8931–8939, July 2008.
- [9] Joseph Felsenstein. *Inferring Phylogenies*. Sinauer Associates, September 2003.
- [10] Anthony A. Fodor and Richard W. Aldrich. On evolutionary conservation of thermodynamic coupling in proteins. *Journal of Biological Chemistry*, 279(18):19046–19050, April 2004.
- [11] Angelika Fuchs, Antonio J. Martin-Galiano, Matan Kalman, Sarel Fleishman, Nir Ben-Tal, and Dmitrij Frishman. Co-evolving residues in membrane proteins. *Bioinformatics*, 23(24):3312–3319, December 2007.
- [12] E.J. Fuentes, C.J. Der, and A.L. Lee. Ligand-dependent dynamics and intramolecular signaling in a PDZ domain. *Journal of molecular biology*, 335(4):1105–1115, 2004.
- [13] Ulrike Göbel, Chris Sander, Reinhard Schneider, and Alfonso Valencia. Correlated mutations and residue contacts in proteins. *Proteins: Structure, Function, and Genetics*, 18(4):309–317, April 1994.
- [14] Holger Hofling and Robert Tibshirani. Estimation of sparse binary pairwise markov networks using pseudo-likelihoods. *Journal of Machine Learning Research*, 10:883–906, April 2009.
- [15] H. Kamisetty, B. Ghosh, C. Bailey-Kellogg, and C.J. Langmead. Modeling and Inference of Sequence-Structure Specificity. In *Proc. of the 8th International Conference on Computational Systems Bioinformatics (CSB)*, pages 91–101, 2009.
- [16] Su-In Lee, Varun Ganapathi, and Daphne Koller. Efficient structure learning of markov networks using  $l_1$ -regularization. In B. Schölkopf, J. Platt, and T. Hoffman, editors, *Advances in Neural Information Processing Systems 19*, pages 817–824. MIT Press, Cambridge, MA, 2007.
- [17] S. W. Lockless and R. Ranganathan. Evolutionarily conserved pathways of energetic connectivity in protein families. *Science*, 286:295–299, Oct 1999.

- [18] D. D. Pollock and W. R. Taylor. Effectiveness of correlation analysis in identifying protein residues undergoing correlated evolution. *Protein Eng.*, 10(6):647–657, June 1997.
- [19] W. P. Russ, D. M. Lowery, P. Mishra, M. B. Yaffe, and R. Ranganathan. Natural-like function in artificial WW domains. *Nature*, 437:579–583, Sep 2005.
- [20] Mark Schmidt, Kevin Murphy, Glenn Fung, and Rmer Rosales. Structure learning in random fields for heart motion abnormality detection. In *CVPR*. IEEE Computer Society, 2008.
- [21] Mark Schmidt, Ewout van den Berg, Michael P. Friedlander, and Kevin Murphy. Optimizing costly functions with simple constraints: A limited-memory projected quasi-newton algorithm. *Proc. of Conf. on Artificial Intelligence and Statistics*, pages 456–463, 2009.
- [22] M. Socolich, S. W. Lockless, W. P. Russ, H. Lee, K. H. Gardner, and R. Ranganathan. Evolutionary information for specifying a protein fold. *Nature*, 437:512–518, Sep 2005.
- [23] Peter Spirtes, Clark Glymour, and Richard Scheines. *Causation, Prediction, and Search (Lecture Notes in Statistics)*. Springer-Verlag, 1993.
- [24] J. Thomas, N. Ramakrishnan, and C. Bailey-Kellogg. Graphical models of residue coupling in protein families. *IEEE/ACM Transactions on Computational Biology and Bioinformatics (TCBB)*, 5(2):183–197, 2008.
- [25] J. Thomas, N. Ramakrishnan, and C. Bailey-Kellogg. Graphical models of protein-protein interaction specificity from correlated mutations and interaction data. *Proteins: Structure, Function, and Bioinformatics*, 76(4):911–29, 2009.
- [26] J. Thomas, N. Ramakrishnan, and C. Bailey-Kellogg. Protein Design by Sampling an Undirected Graphical Model of Residue Constraints. *IEEE/ACM Transactions on Computational Biology and Bioinformatics (TCBB)*, 6(3):506–516, 2009.
- [27] Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society, Series B*, 58:267–288, 1994.
- [28] JA Tropp. Just relax: Convex programming methods for identifying sparse signals in noise. *IEEE Transactions on Information Theory*, 52(3):1030–1051, 2006.
- [29] Martin J. Wainwright, Pradeep Ravikumar, and John D. Lafferty. High-dimensional graphical model selection using  $\ell_1$ -regularized logistic regression. In B. Schölkopf, J. Platt, and T. Hoffman, editors, *Advances in Neural Information Processing Systems 19*, pages 1465–1472. MIT Press, Cambridge, MA, 2007.
- [30] M. Weigt, R. A. White, H. Szurmant, J. A. Hoch, and T. Hwa. Identification of direct residue contacts in protein-protein interaction by message passing. *Proc. Natl. Acad. Sci. U.S.A.*, 106:67–72, Jan 2009.
- [31] Ming Yuan and Yi Lin. Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society, Series B*, 68:49–67, 2006.