

# Learning the structure of a Markov Random Field

Lee Richert (ECE), Raghav Kuppan (ECE), Yuanda Zhu (ECE)

## 1 Abstract

Learning the structure of a Markov Random Field from independent and identically distributed samples is an important problem. Unfortunately, calculation of the partition function is intractable for most graphs. Trace Lasso regularization of a pseudo-likelihood has been proposed as a means of learning the structure of probabilistic graphical models from samples. The Trace Lasso exhibits less volatility in response to correlation between random variables than Lasso while also promoting sparsity. We plan to adapt Trace Lasso for the purpose of learning the structure of a Markov Random field, using a maximum pseudo-likelihood approach.

## 2 Introduction

### 2.1 Defining the Ising Model

An Ising Model is a special case of a pairwise Markov Random Field where each vertex in the graph takes on values in  $\{-1, 1\}$ . The node potentials and edge potentials for an Ising model have very simple expressions thereby giving us a distribution of the form

$$\mathbb{P}_{\theta^*}(x) = \frac{1}{\mathbb{Z}(\theta^*)} \exp \sum_{(s,t) \in E} \theta_{st}^* x_s x_t$$

The Ising Model was proposed as a mathematical model for ferromagnetism in statistical mechanics but is also used in other applications such as Computer Vision and Neuroscience.

## 3 Prior Work on Model Selection

Learning the structure of Ising models is a challenging problem. For a general graph with  $p$  nodes of degree at most  $d$ , an exhaustive search across all possible edges takes around  $p^d$  computations. This is the time required to exhaustively search over all possible neighborhoods of a node and for each node test whether conditional independence assumptions hold. As  $d$  grows, the computational cost becomes untenable, so algorithms with lower

computational complexity are desired. In general, efficient algorithms for structural learning either restrict the graph structure or the nature of the interactions between the nodes. One possible assumption of the second kind is the Correlation Decay Property. A graphical model is said to have the correlation decay property if any two variables are asymptotically independent as the graph distance between them increases. This property holds for Ising models in certain real-world problems such as the ferromagnetic model in a high temperature regime. Alternative methods that do not explicitly require the Correlation Decay Property are usually based on Convex Optimization.

Ravikumar et. al [?] study the problem of signed edge recovery on Ising models and propose an  $l_1$ -regularized logistic regression approach to recover the signed edges in the graph. They establish sufficient conditions on the sample size  $n$ , dimension  $p$ , and maximum neighborhood size  $d$ , to get a consistent estimator. With this assumption, the structure of any bounded degree graph can be recovered with high probability once  $n/\log(p)$  is sufficiently large. This way, the signed neighborhood of every node can be estimated by optimizing a log pseudo-likelihood function with an  $\ell_1$  penalty. This way the entire set of signed edges of the graph can be recovered whereby the structure of the Ising model has been found. This technique is demonstrated on four-nearest neighbor lattices, eight-nearest neighbor lattices and a star graph as well as on a class of graphs with unbounded maximum neighborhood size, with the results being consistent with the theoretical conjectures

### 3.1 Learning Ising Models from Complete Data

### 3.2 Pseudolikelihood Functions

A Pseudolikelihood is an approximation to the likelihood function which provides a computationally simpler model for inference in a graphical model and is also exact provided certain assumptions are made about the model. The Ising model joint probability expression cannot be exactly computed since the partition function makes computations intractable. Our method involves computing a regularized-pseudolikelihood of node  $X_r$  conditioned on the other nodes. The pseudo-likelihood approach involves computing the conditional distribution of a node conditioned on all the other nodes, and then performing the optimization over the parameter set. As we are dealing with the Ising model, the conditional distribution has a very clean expression

given as

$$\mathbb{P}_{\theta^*}(x_r | x_{\setminus r})$$

Provided we have samples from the graphical model, this expression is easily computed.

Therefore, given a set of independent and identically distributed samples, the regularized pseudo-likelihood is of the form of a convex program

$$\min_{\theta_{\setminus r} \in \mathbb{R}^{p-1}} \{l(\theta; \mathbb{X}_1^n) + \lambda \|\theta_{\setminus r}\|\}$$

where  $\mathbb{X}_1^n$  is the set of i.i.d samples from the graph and  $\lambda$  is the regularization parameter.

Different Regularization techniques can be used according to the correlation structure of our data and the application where the Ising model is being used.

### 3.3 Regularization

Regularization is a method to avoid the overfitting problem in Machine Learning. In our case, it takes the form of a feature selection problem and is used to get a sparse approximation for our optimization. This makes our model cheaper and more interpretable. Among sparsity inducing norms, the  $\ell_1$  norm is the simplest and most widely used, leading to the Lasso when used in a least-squares framework. While the Lasso does well in high-dimensional settings, it is known to have stability problems in situations where the data exhibit strong correlation structures. Several solutions have been proposed which include the Elastic net, group Lasso and Sampling techniques. However these norms cannot just be plugged into the objective function as extra information is usually required, or in some cases, the problem is further complicated due to the addition of parameters to be estimated. The Trace norm takes into account the correlation structure of the data but does not require manual human intervention. It can be thought of as a way of turning the rank of a matrix into a norm. The Trace norm is adaptive and requires that only a single regularization parameter be chosen.

#### 3.3.1 Lasso

Penalizing linear models by the number of variables used in the model is a method of variable selection in sparse high-dimensional settings. As this criterion is not convex, a convex relaxation for this norm is the  $\ell_1$  norm. But a solution to the likelihood function penalized with this penalty is not

equivariant to the rescaling of the predictors, so it is common to normalize the predictors. When normalizing the predictors and penalizing with the Lasso, we are implicitly using a regularization term that depends on the data matrix. We can write the normalized  $\ell_1$  norm as

$$\|\theta\|_1 = \sum_{i=1}^p \|\mathbf{X}^i\|_2 |\theta_i|$$

But the Lasso is known to perform poorly under three conditions:

1. Let  $p$  represent the number of features and  $n$  represent the number of observations. For  $p > n$ , Lasso selects at most  $n$  variables before it saturates; besides, Lasso is not well defined unless the bound on  $\ell_1$ -norm of the coefficients is smaller than a certain value.
2. For a group of variables whose pairwise correlations are very high, Lasso tends to select only one variable and does not care which one is selected.
3. For  $n > p$  case, when predictors have high correlations, the prediction performance of Lasso is dominated by ridge regression.

### 3.3.2 Group Lasso

Group Lasso divides the predictors into groups and penalize the sum of the ' $\ell_2$  norm of these groups. The Group Lasso is especially useful in case our model contains categorical variables. In case of categorical variables, the Lasso might leave out some levels in the category, as discussed in the previous section. The Group Lasso provides a method to incorporate prior information about the categories available. Given a partition set  $(S_i)$ , the group Lasso can be calculated as

$$\|\theta\|_{GL} = \sum_{i=1}^k \|\theta_i\|_2$$

The effect of this penalty function is to introduce sparsity at the group level: variables in a group are selected all together. The drawback here is that the Correlation structure of the data or information about the category partition must be known a priori.

### 3.3.3 Elastic Net

In order to address the third problem of Lasso, elastic net [?] was proposed in 2005 by adding the squared  $\ell_2$  norm, a strongly convex penalty term to the  $\ell_1$  norm in Lasso. Elastic net performs similarly to Lasso in scenario 1) and 2), but by encouraging grouping effect, has higher accuracy than Lasso in scenario 3). To be more specific, in scenario 3), some features are highly correlated with each other and are associated with response; thus elastic net aims to perform less shrinkage on those subsets of features. In addition, similar to Lasso, elastic net does both continuous shrinkage and automatic variable selection. Ridge regression has only continuous shrinkage but no automatic variable selection.

### 3.3.4 Trace Lasso

The Trace norm is a measure of the dimension of the subspace spanned by the selected predictors. The trace norm of a matrix is the sum of its singular values. The Trace lasso can be defined as  $\|\mathbf{X} \text{Diag}(\theta)\|_*$ . It has some interesting properties;

1. If all the predictors are orthogonal, then, it is equal to the  $\ell_1$  norm.
2. If all the predictors are equal, then, it is equal to the  $\ell_2$  norm

Thus when two predictors are strongly correlated, our norm will behave like Tikhonov regularization, while for almost uncorrelated predictors, it will behave like the Lasso.

Instead of blindly adding convexity in all directions, the trace lasso adds strong convexity only in the direction of highly correlated predictors. Thus, it always has a unique minimum and is much more stable than the lasso.

## 4 Approach

Our approach is to adopt the maximum psuedo-likelihood approach as in [?] but with the trace norm as our regularization parameter. The signed neighborhood of one node will be estimated by evaluating the conditional distribution of the node conditioned on all the other nodes and then maximizing this with a trace norm penalty. Computing the Pseudo-likelihood using the generated samples is not very hard for the Ising model. However there are

some difficulties in evaluating the trace norm. The usual method of optimization is by subgradient descent where the subgradient is the generalization of the derivative to functions which are not differentiable. Computing the subgradient of the trace norm is reported to be inefficient and the rate of convergence very slow, so a variational formulation for the trace norm [?] is to be used. The ensuing cost function can be optimized using a Conjugated Gradients method.

The Trace Norm of a matrix  $M$  can be expressed as

$$\|M\|_* = \frac{1}{2} \inf_{S \succeq 0} \text{tr}(M^T S^{-1} M) + \text{tr}(S)$$

and the infimum is attained for  $S = M^T M$ .

Therefore, our optimization problem can be reformulated as

$$\min_{\theta \in \mathbb{R}^{p-1}} l(\theta; \mathbb{X}_1^n) + \frac{\lambda}{2} \theta^T \text{Diag}(\text{diag}(X^T S^{-1} X)) \theta + \frac{\lambda}{2} \text{tr}(S)$$

where  $\text{Diag}(u)$  is the diagonal matrix whose diagonal elements are the vector  $u$  and  $\text{diag}(M)$  is the diagonal of the matrix  $M$ . This problem is jointly convex in  $(\theta, S)$  and can be solved by alternating the minimization over  $\theta$  and  $S$ . The matrix  $S$  can be computed with a closed form expression for the given problem which leaves us with a convex optimization program for the parameter  $\theta$ . This can be solved using a conjugated gradients method. The optimization is still not cheap as compared to the  $\ell_1$  norm and takes quite some time to run.

## 4.1 Applying Methods to Categorical Data

## 4.2 Learning Graphical Models for Classification

# 5 Datasets

# 6 Evaluating Performance

# 7 Logistics

Cleaning of the data and other basic data processing was done on the samples from the dataset using Python. This step also included the removal

of samples with missing data from the dataset. The Pseudo-likelihood optimization with regularization was done in Python using the NumPy and SciPy libraries. More specifically, the *minimize* procedure from the *optimize* library was used for the Conjugated Gradients method. The lasso, trace lasso and the elastic net were the regularization schemes implemented in the program.

## 8 Experiments

### 8.1 Synthetic Data

### 8.2 Congressional Voting Data

### 8.3 Mushroom Data

## 9 Conclusions

## 10 Acknowledgements