

Markov Network Structure Learning

Sargur Srihari

srihari@cedar.buffalo.edu

Topics

- Structure Learning as model selection
- Structure Learning using Independence Tests
- Score based Learning: Hypothesis Spaces
- Objective Functions
 - Likelihood Functions
 - Bayesian Scores
- Optimization Task

Structure Learning

- Model Selection Problem:
 - Learning network structure from data
- Two types of solutions:
 - Constraint-based approaches
 - Search for graph structure satisfying independencies observed in empirical data
 - Score-based approach
 - Define objective functions for different models
 - Search for high-scoring model

Structure Learning using Independence Tests

- Constraint-based structure learning: use independence tests
- Assume
 - generating distribution P^* is positive
 - Can be represented as a Markov network H^* that is a perfect map of P^*
- Perform set of independence tests on P^* to recover H^*
- For tractability assume degree of nodes in H^* is at most d^*

Intractability of using local independencies

- Consider using Markov blanket independencies or Local independencies
 1. $(X \perp_{\chi - \{X\} - MB_P(X)} MB_P(X) | MB_P(X))$
 2. $(X \perp Y | \chi - \{X, Y\})$
- Neither is tractable
 - NP-hard (no polynomial time algorithm)
 - Need probabilities of exponential no. of events
 - Statistical issues
 - Independence tested on empirical data
 - Need exponential number of data points
- Constraint-based: useful only as starting point

Score-based Learning

- Structure learning formulated as an optimization problem
 - Hypothesis space consists of set of possible networks
 - Define an objective function
 - To score candidate networks
 - Formulate a search strategy

Score Based Learning: Hypothesis Spaces

- Structure learning is an optimization problem
- Several ways of formulating search space, varying in granularity
 - Coarse grained
 - Space of different structures
 - Model complexity measured in terms of clique size
 - Next level
 - Factor graphs
 - Measure complexity as Sizes of factors
 - Finest level
 - Individual features in a log-linear model
 - Measure sparsity at level of features included

Hypothesis Space Comparison

- More fine grained parameterization better matches properties of distribution
- Factor graph
 - Distinguishes between single large factor over k variables and a set of kC_2 pairwise factors requiring fewer parameters
- Feature-based approach
 - Distinguishes between full factor over k variables and a single log-linear feature over same variables

Fine-grained: obscures structure

- Sparsity in space of features selected does not correspond to sparsity in model structure
- Introducing a single feature $f(d)$ into the model has effect of introducing edges between all variables in d
- Models with few features can give rise to dense connectivity in networks
- Search algorithms are slower

Search space of log-linear models

- We have a set of features Ω that can have nonzero weight
- Select model structure M defined by some subset $\Phi[M] \subseteq \Omega$
- Let $\Theta[M]$ be the set of parameterizations θ compatible with model structure
 - i.e., those where $\theta_i \neq 0$ only if $f_i \in \Phi[M]$

Resulting log-linear distribution

- A structure and a compatible parameterization define a log-linear distribution

$$p(\chi|\mathbf{M},\theta) = \frac{1}{Z} \exp \left\{ \sum_{i \in \Phi[\mathbf{M}]} \theta_i f_i(\xi) \right\} = \frac{1}{Z} \exp \{ f^T \theta \}$$

- Where because of compatibility of θ with \mathbf{M} ,
- A feature not in $\Phi[\mathbf{M}]$ does not influence the final vector product
- Since it is multiplied by a parameter that is 0¹¹

Objective Function: Likelihood

- What to optimize in score-based approach
- Most straightforward is the likelihood of training data

$$score_L(\mathbf{M} : D) = \max_{\theta \in \Theta[\mathbf{M}]} \ln P(D | \mathbf{M}, \theta) = \ell\left(\left(\mathbf{M}, \hat{\theta}_M\right) : D\right)$$

where $\hat{\theta}_M$ are the maximum likelihood parameters compatible with \mathbf{M}

- Score measures fitness of model to data
- It prefers more complex models

Objective Function: Bayesian Scores

- Bayesian score: marginal likelihood that integrates over all possible parameterizations

$$\int P(D | M, \theta) P(\theta | M) d\theta$$

- Difficult to evaluate for MNs (unlike BNs)
- Simplest asymptotic approximation- BIC score

$$score_{BIC}(M : D) = \ell\left(\left(M, \hat{\theta}_M\right) : D\right) - \frac{\dim(M)}{2} \ln M$$

- Another is Laplace Approximation

$$score_{Laplace}(M : D) = \ell\left(\left\langle M, \hat{\theta}_M \right\rangle : D\right) - \frac{\dim(M)}{2} \ln M$$

where $\hat{\theta}_M$ are parameters obtained from MAP estimation:

$$\hat{\theta}_M = \arg \max_{\theta} P(D | \theta, M) P(\theta | M) \text{ and } A \text{ is the negative Hessian matrix evaluated at } \hat{\theta}_M$$

$$A_{i,j} = - \frac{\partial}{\partial \theta_i} \frac{\partial}{\partial \theta_j} \ell\left(\left\langle M, \hat{\theta}_M \right\rangle : D\right) + \ln P(\theta | M)$$

Parameter Penalty Scores

- Alternative to marginal likelihood
- Evaluate maximum posterior probability

$$score_{MAP}(\mathbf{M} : D) = \max_{\theta \in \Theta[\mathbf{M}]} \ell\left(\left(\mathbf{M}, \hat{\theta}_M\right) : D\right) + \ln P\left(\hat{\theta}_M \mid \mathbf{M}\right)$$

- Prior regularizes the likelihood, moving it away from the maximum likelihood values
- L2 regularized Map leads to a fully-connected structure
- L1 regularization leads to a sparser structure

Optimization Task

- Greedy structure search
- External search over structures is implemented as a local search

Greedy score-based structure search for log-linear models

```

Procedure Greedy-MN-Structure-Search (
     $\Omega$ ,    // All possible features
     $\mathcal{F}_0$ ,  // initial set of features
     $\text{score}(\cdot : \mathcal{D})$ ,  // Score
)
1   $\mathcal{F}' \leftarrow \mathcal{F}_0$     // New feature set
2   $\theta \leftarrow 0$ 
3  do
4       $\mathcal{F} \leftarrow \mathcal{F}'$ 
5       $\theta \leftarrow \text{Parameter-Optimize}(\mathcal{F}, \theta, \text{score}(\cdot : \mathcal{D}))$ 
6          // Find parameters that optimize the score objective, relative to
           // current feature set, initializing from the current parameters
7      for each  $f_k \in \mathcal{F}$  such that  $\theta_k = 0$ 
8           $\mathcal{F} \leftarrow \mathcal{F} - f_k$ 
9          // Remove inactive features
10     for each operator  $o$  applicable to  $\mathcal{F}$ 
11         Let  $\hat{\Delta}_o$  be the approximate improvement for  $o$ 
12         Choose some subset  $\mathcal{O}$  of operators based on  $\hat{\Delta}$ 
13          $\mathcal{F}' \leftarrow \mathcal{O}(\mathcal{F})$     // Apply selected operators to  $\mathcal{F}$ 
14     while termination condition not reached
15 return  $(\mathcal{F}, \theta)$ 

```


Successor Evaluation

- Considerably more expensive than for BNs
- At each stage need to evaluate the score for all candidates we wish to examine
- Requires estimating parameters for the structure
- Use heuristic that a single change to the structure does not result in drastic changes to model

Choice of Scoring Function

- Greedy algorithm template can be used with any objective function
- Log-likelihood or L2-regularized log-likelihood give rise to non-zero parameters
- L1-regularized Likelihood can be optimized in a way that guarantees convergence to globally optimal solution

L_1 -Regularization for Structure Learning

- L_1 -regularized likelihood

$$score_{L_1}(\theta : D) = \ell(\langle M, \theta \rangle : D) - \|\theta\|_1$$

- Implement optimization as a double loop algorithm where we separately consider structure and parameters

Evaluating Changes in model

- Candidate evaluation step 11 of *Greedy-MN-Structure Search*
- Simply avoid computing the exact score of successors and use simpler heuristics
- A more refined approach is called *Grafting*
 - Estimate benefit of introducing a feature f_k by computing the gradient of the likelihood relative to θ_k for score X

$$\Delta_X^{grad}(\theta_k : \theta^0, D) = \frac{\partial}{\partial \theta_k} score_X(\theta : D)$$