

## Introduction

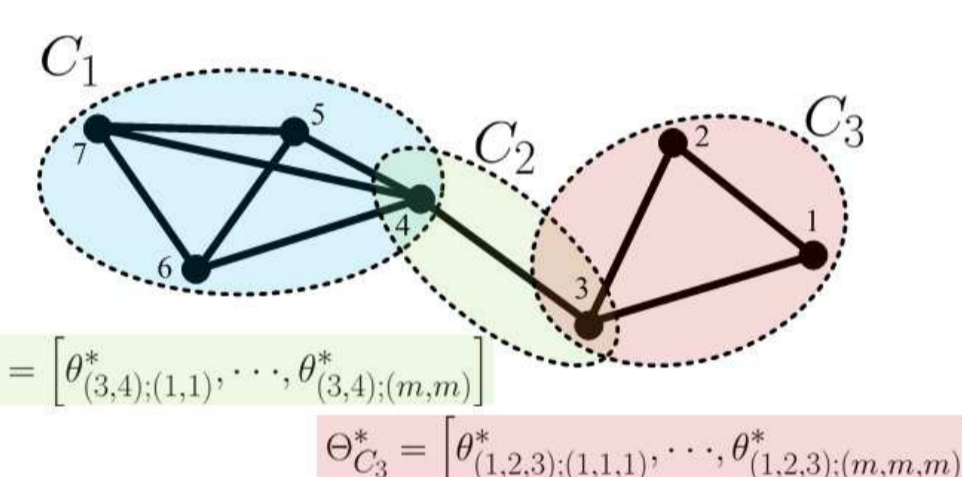
### ◇ Problem Setup

- Random variables  $X = (X_1, X_2, \dots, X_p)$  taking values in  $[m] = \{1, 2, \dots, m\}$
- Joint PDF  $\mathbb{P}(X) \propto \exp(\sum_{C \in \mathcal{C}} \phi_C(X_C))$  for some  $\mathcal{C} \subseteq 2^{[p]}$  and potentials  $\phi_C(\cdot)$
- Potential Functions  $\phi_C(X_C) = \sum_{v \in [m]^{|C|}} \theta_{C,v}^* \mathbb{I}[x_C = v]$

Given  $n$  iid samples  $X^{(i)}$ , learn  $\Theta_C^*$  for all  $C \in 2^{[p]}$

### ◇ Graphical Model

- Define the graph  $\mathcal{G}$  with  $p$  vertices corresponding to variables
- If  $\Theta_C^* \neq 0$  then make a complete graph on vertices  $X_C$



Example:

$$\mathbb{P}(X) \propto f_1(X_{C_1}; \Theta_{C_1}^*) f_2(X_{C_2}; \Theta_{C_2}^*) f_3(X_{C_3}; \Theta_{C_3}^*)$$

$$\Theta_{C_2}^* = [\theta_{(3,4),(1,1)}^*, \dots, \theta_{(3,4),(m,m)}^*]$$

$$\Theta_{C_3}^* = [\theta_{(1,2,3),(1,1,1)}^*, \dots, \theta_{(1,2,3),(m,m,m)}^*]$$

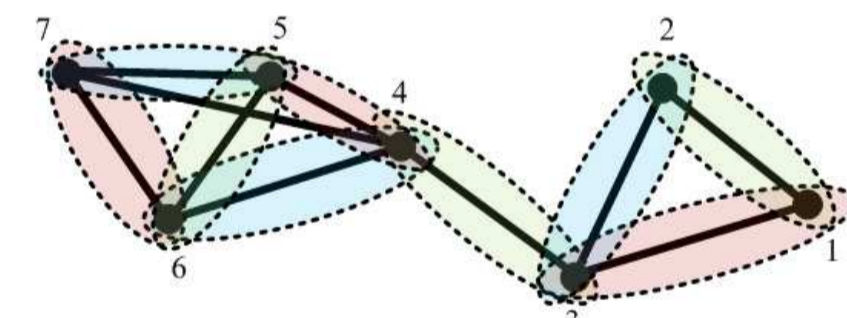
Given  $n$  ( $\propto \log(p)$ ) iid samples  $X^{(i)}$ , learn the *sparse* graph  $\mathcal{G}$

### ◇ Pairwise Graphical Model

- Only cliques of size 2 exist
- Assign  $\Theta_{rt} \in \mathbb{R}^{m^2}$  to the edge  $(r, t)$  intersecting  $X_r$  and  $X_t$

Example:

$$\mathbb{P}(X) \propto \prod_{i=1}^{10} f_i(X_{C_i}; \Theta_{C_i}^*)$$



## Methodology

### ◇ Learning Neighborhoods

- For each vertex learn its neighborhood consistently
- Combine all neighborhoods to get the graph

### ◇ Conditional Likelihood

- Fixing a node  $X_r$ ;

$$\mathbb{P}[X_r = j | X_{\setminus r}] = \frac{\exp(\sum_{t \in [p] - \{r\}} \sum_{k=1}^m \theta_{rt,jk}^* \mathbb{I}[X_t = k])}{\sum_{j=1}^m \exp(\sum_{t \in [p] - \{r\}} \sum_{k=1}^m \theta_{rt,jk}^* \mathbb{I}[X_t = k])}$$

- Likelihood function

$$\mathcal{L}(\Theta_r) = -\frac{1}{n} \sum_{i=1}^n \log(\mathbb{P}[X_r = x_r^{(i)} | X_{\setminus r} = x_{\setminus r}^{(i)}])$$

## Our Algorithm

**Algorithm 1** (Discrete Graphical Model Learner). For a fixed vertex  $r$ , find  $\hat{\Theta}_r$ ; the solution of the following convex optimization program

$$\text{Minimize: } \mathcal{L}(\Theta_r) + \lambda \sum_{t \in [p] - \{r\}} \|\Theta_{rt}\|_2$$

Output the neighborhood  $\hat{\mathcal{N}}_r = \{t : \|\hat{\Theta}_{rt}\|_0 \neq 0\}$ .

## Pairwise Case

### ◇ Assumptions

Let  $Q_r^* = \mathbb{E}[\nabla^2 \log(\mathbb{P}[X_r | X_{\setminus r}])]$  be the Fisher information matrix.

- (A1) Minimum curvature:  $\Lambda_{\min}(Q_r^*) \geq C_{\min} > 0$ .
- (A2) Incoherence:  $\left\| Q_{\mathcal{N}_r^* \setminus \mathcal{N}_r^*}^* (Q_{\mathcal{N}_r^* \setminus \mathcal{N}_r^*}^*)^{-1} \right\|_{\infty, 2} \leq \frac{1-2\alpha}{\sqrt{|\mathcal{N}_r^*|}}$  for some  $\alpha \in (0, 1/2)$ .

**Theorem 1** (Pairwise model sparsistency). Under assumptions (A1)-(A2), suppose

$$n \geq Km^2 |\mathcal{N}_r^*|^2 \log(m^2 p)$$

and  $\lambda \geq \frac{8(2-\alpha)}{\alpha} \left( \sqrt{\frac{\log(p)}{n}} + \frac{m-1}{4\sqrt{n}} \right)$ . Then, Algorithm 1, with probability  $1 - c_1 \exp(-c_2 \lambda^2 n)$ , excludes all false neighbors and includes all neighbors  $t$  provided that  $\|\Theta_{rt}\|_2 \geq \frac{10}{C_{\min}} \lambda$ .

## General Higher-Order Case

### ◇ Further Assumption

- (A3) Bounded Mismatch: Non-Pairwise (high order) parameters  $\Theta_{P_c}^*$  satisfy  $\|\Theta_{P_c}^*\|_1 \leq K \frac{C_{\min}^2}{(m-1)^{|\mathcal{N}_r^*|}}$ .

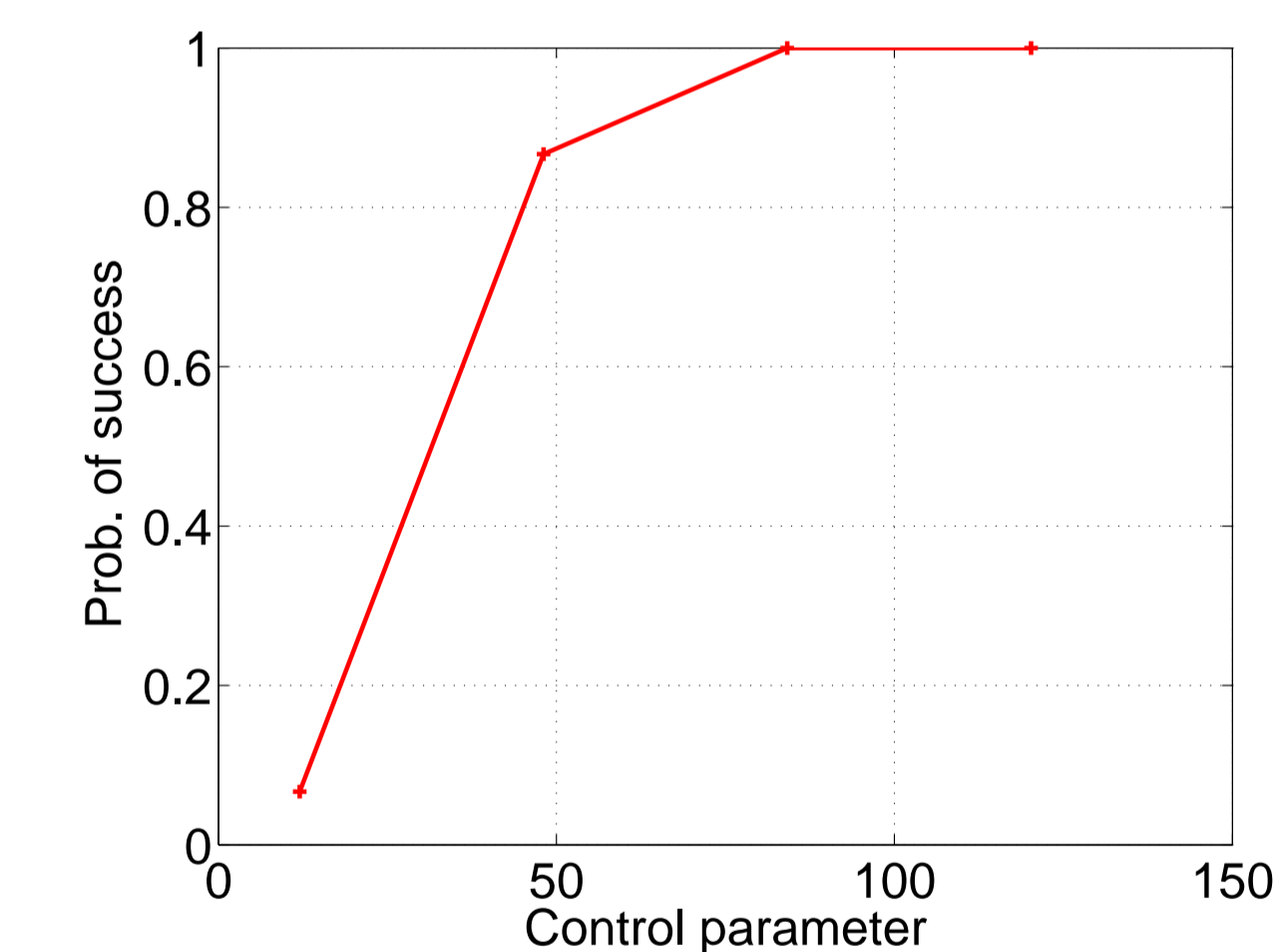
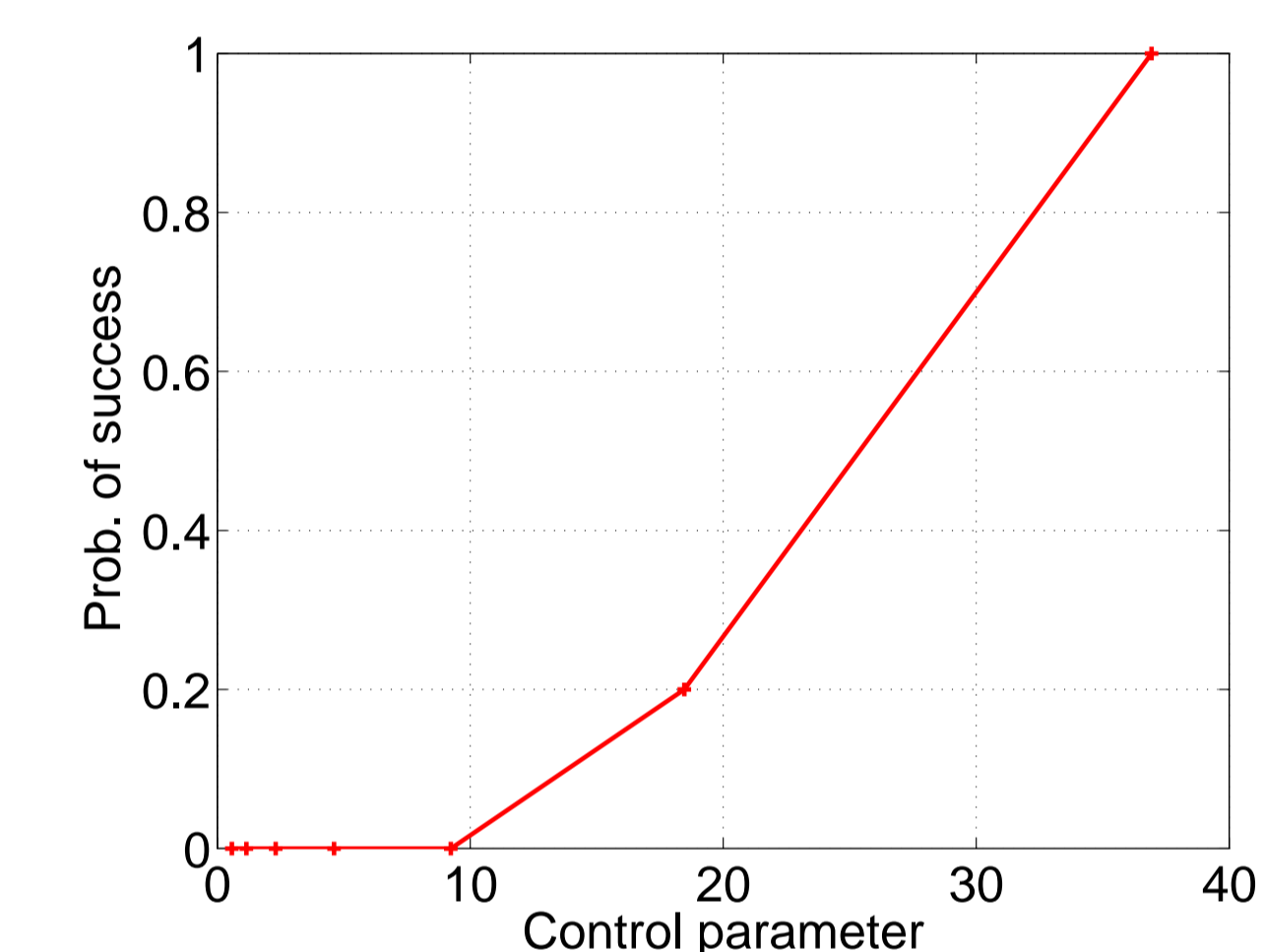
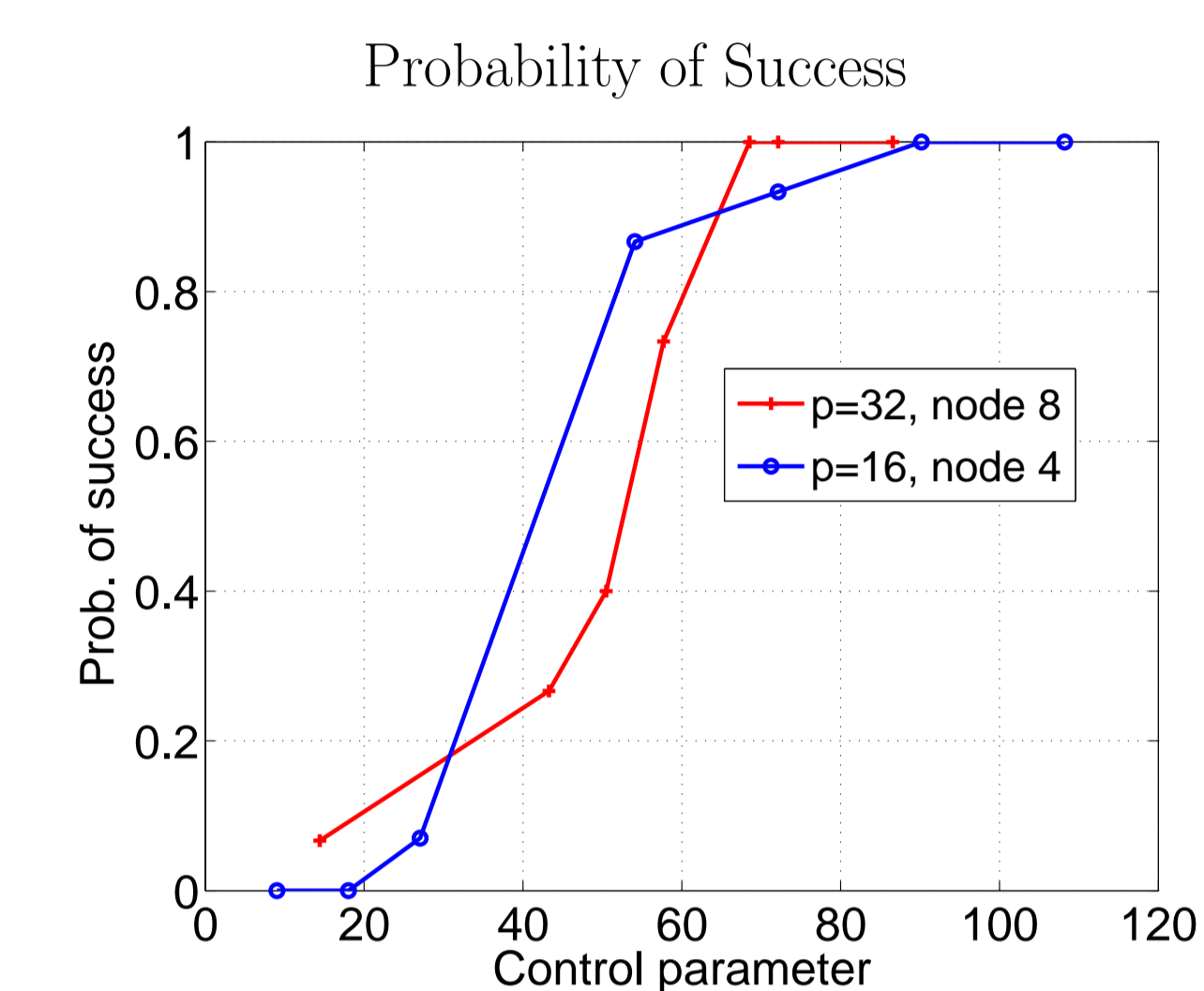
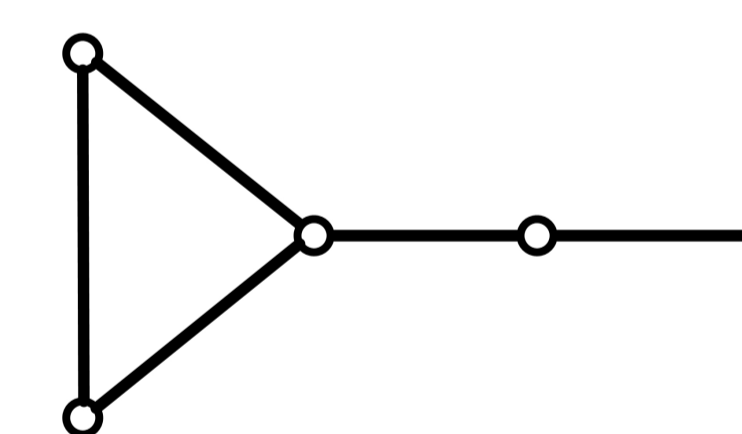
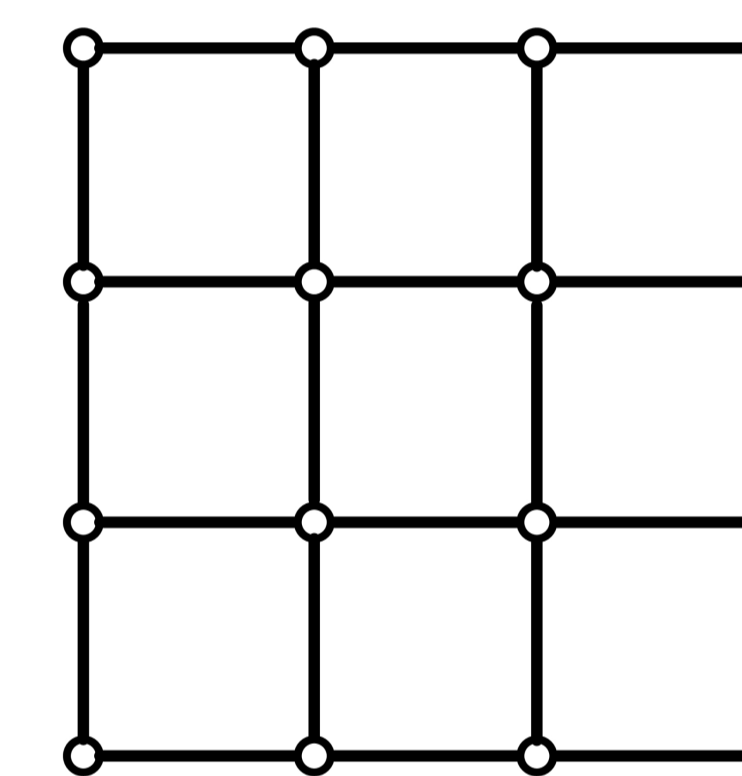
**Theorem 2** (Higher-Order model sparsistency). Under assumptions (A1)-(A3), suppose

$$n \geq Km^2 |\mathcal{N}_r^*|^{1.5c+0.5} \log(m^c p^{c-1})$$

and  $\lambda \geq \frac{8(2-\alpha)}{\alpha} \left( \sqrt{\frac{\log(p)}{n}} + \frac{m-1}{4\sqrt{n}} + \frac{1}{4} \|\Theta_{P_c}^*\|_1 \right)$ . Then, Algorithm 1, with probability  $1 - c_1 \exp(-c_2 (\lambda^2 - 2 \|\Theta_{P_c}^*\|_1) n)$ , excludes all false neighbors and includes all neighbors  $t$  provided that  $\|\Theta_{rt}\|_2 \geq \frac{10}{C_{\min}} \lambda$ .

## Numerical Experiments

True Graphical Model



Group-sparse regularization can efficiently learn general discrete graphical models.

## References

- [1] P. Ravikumar, M. J. Wainwright, and J. Lafferty. "High-dimensional ising model selection using  $\ell_1$ -regularized logistic regression", *Annals of Statistics*, 38(3), 1287:1319.
- [2] F. Bach. "Consistency of the group LASSO and multiple kernel learning", *Journal of Machine Learning Research*, 9:1179-1225, 2008.
- [3] N. Meinshausen and P. Buhlmann. "High dimensional graphs and variable selection with the lasso", *Annals of Statistics*, 34(3), 2006.