

## A Appendix

### A.1 Coordinate Descent Updates for Alternating Newton Coordinate Descent Method

In our alternating Newton coordinate descent algorithm, each element of  $\Delta_\Lambda$  is updated as follows:

$$(\Delta_\Lambda)_{ij} \leftarrow (\Delta_\Lambda)_{ij} - c_\Lambda + S_{\lambda_\Lambda/a_\Lambda}(c_\Lambda - \frac{b_\Lambda}{a_\Lambda}),$$

where  $S_r(w) = \text{sign}(w) \max(|w| - r, 0)$  is the soft-thresholding operator and

$$\begin{aligned} a_\Lambda &= \Sigma_{ij}^2 + \Sigma_{ii}\Sigma_{jj} + \Sigma_{ii}\Psi_{jj} + 2\Sigma_{ij}\Psi_{ij} + \Sigma_{jj}\Psi_{ii} \\ b_\Lambda &= (\mathbf{S}_{\mathbf{y}\mathbf{y}})_{ij} - \Sigma_{ij} - \Psi_{ij} + (\Sigma\Delta_\Lambda\Sigma)_{ij} + (\Psi\Delta_\Lambda\Sigma)_{ij} + (\Psi\Delta_\Lambda\Sigma)_{ji} \\ c_\Lambda &= \Lambda_{ij} + (\Delta_\Lambda)_{ij}. \end{aligned}$$

For  $\Theta$ , we perform coordinate-descent updates directly on  $\Theta$  without forming a second-order approximation of the log-likelihood to find a Newton direction, as follows:

$$\Theta_{ij} \leftarrow \Theta_{ij} - c_\Theta + S_{\lambda_\Theta/a_\Theta}(c_\Theta - \frac{b_\Theta}{a_\Theta}),$$

where

$$\begin{aligned} a_\Theta &= 2\Sigma_{jj}(\mathbf{S}_{\mathbf{xx}})_{ii} \\ b_\Theta &= 2(\mathbf{S}_{\mathbf{xy}})_{ij} + 2(\mathbf{S}_{\mathbf{xx}}\Theta\Sigma)_{ij} \\ c_\Theta &= \Theta_{ij}. \end{aligned}$$

### A.2 Time Complexity Analysis of Alternating Newton Block Coordinate Descent

In this section we describe the time complexity of the alternating Newton block coordinate descent method. The active set sizes for  $\Lambda$  and  $\Theta$  are  $m_\Lambda$  and  $m_\Theta$ , respectively. Also, because we use the iterative conjugate gradient method to compute columns of  $\Sigma$ , we assume that solving  $\Lambda\Sigma_i = \mathbf{e}_i$  takes at most  $K$  iterations.

#### A.2.1 Time Cost of Updating $\Lambda$

The time complexity of each  $\Lambda$  update is dominated by the cost of precomputing columns of  $\Sigma$  and  $\Psi$ . The cost of these precomputations is  $O\left(\left[1 + \frac{B_\Lambda}{q}\right][m_\Lambda Kq + nq^2]\right)$ , where  $B_\Lambda = \sum_{z \neq r} |\{j | i \in C_z, j \in C_r, (i, j) \in \mathcal{S}_\Lambda\}|$  is the number of cache misses and  $K$  is the number of conjugate gradient iterations. Although the worst-case of  $B_\Lambda = k_\Lambda q$  requires computing  $\Sigma$  and  $\Psi$  a total of  $k_\Lambda$  times, in practice, graph clustering dramatically reduces this additional cost of block-wise optimization. In the best case, when graph clustering identifies perfect block-diagonal structure in the active set, the number of cache misses  $B_\Lambda = 0$  and we incur no runtime penalty from limited memory.

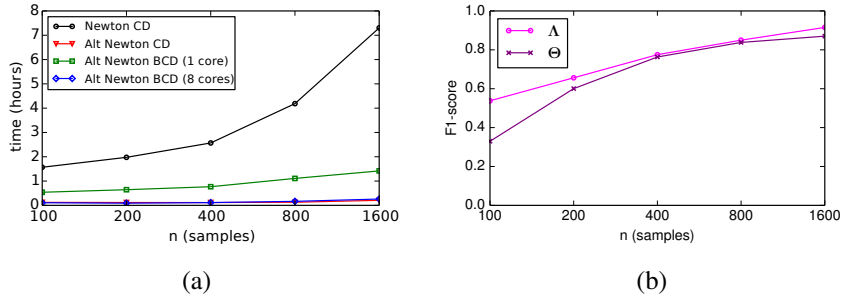


Figure 1: Results from varying sample size  $n$  on chain graph with  $p = q = 10,000$ . (a) Comparison of computation time of different methods. (b) Comparison of edge recovery accuracy as measured by  $F_1$ -score.

### A.2.2 Time Cost of Updating $\Theta$

The overall runtime is dominated by the cost of precomputing columns of  $\mathbf{S}_{xx}$  and  $\Sigma$ . The complexity of these operations is  $O(m_{\Lambda}Kq + m_{\Theta}q + n\tilde{p}B_{\Theta})$ , where  $\tilde{p}$  is the number of non-empty rows in  $\Theta$  and  $B_{\Theta} = \sum_{i,r} |\{i | j \in C_r, (i, j) \in \mathcal{S}_{\Theta}\}|$  is the number of cache misses. Without any row-wise sparsity we have  $\tilde{p} = p$  and  $B_{\Theta} = k_{\Theta}p$ , so the worst-case is that  $\mathbf{S}_{xx}$  is computed a total of  $k_{\Theta}$  times. The additional cost of computing  $\mathbf{S}_{xx}$  due to  $B_{\Theta}$  cache misses is substantially reduced in real datasets where most input variables influence few or none of the outputs. In the best case, if the active set of  $\Theta$  has a block structure, where each input influences only one group of outputs, overall  $\mathbf{S}_{xx}$  will be computed at most once per iteration and  $B_{\Theta} \leq \tilde{p}$ .

### A.3 Additional Results from Synthetic Data Experiments

We compare the performance of the different algorithms on synthetic datasets with different sample sizes  $n$ , using a chain graph structure with  $p = q = 10,000$ . Figure 1(a) shows that our methods run significantly faster than the previous method across all sample sizes. In Figure 1(b) we measure the accuracy in recovering the true chain graph structure in terms of  $F_1$ -score for different sample sizes  $n$ . At convergence,  $F_1$ -score was the same for all methods to three significant digits. As expected, the accuracy improves as the sample size increases.