

A. Proof of Proposition 2

Proof. For $X_u \in X_1$,

$$P(X_u | \setminus X_u) = \frac{P(X_u, X_{\setminus N(u)} \cap X_2 | X_1 \cup X_{N(u)} \setminus X_u)}{P(X_{\setminus N(u)} \cap X_2 | X_1 \cup X_{N(u)} \setminus X_u)}$$

Since $P(X_u | \setminus X_u) = P(X_u | X_1 \cup X_{N(u)} \setminus X_u)$ by the Markovian property of PMN, we have $X_u \perp\!\!\!\perp X_{\setminus N(u)} \cap X_2 | X_1 \cup X_{N(u)} \setminus X_u$.

$X_v \notin X_{N(u)}$ means $X_v \in X_{\setminus N(u)} \cap X_2$. Using the weak union rule for conditional independence (see e.g., (Koller & Friedman, 2009), 2.1.4.3), we obtain $X_u \perp\!\!\!\perp X_v | \setminus \{X_u, X_v\}$.

For $X_u \in X_2$, the proof is the same. \square

B. Proof of Theorem 3

Proof. We define that $\mathbf{B}(i)$ is the set of passages contains X_i . Here we only show the proof that Eq. (1) holds for GPR. Let's denote ϕ_B as short for $\phi_B(X_B)$.

$$\begin{aligned} & P(X_i | X_1 \cup X_{N(i)} \setminus X_i) \\ &= \frac{\frac{1}{Z} \int_{X_{\setminus N(i)} \cap X_2} P(X_1) P(X_2) \prod_{B \in \mathbf{B}(G)} \phi_B}{\frac{1}{Z} \int_{X_i} \int_{X_{\setminus N(i)} \cap X_2} P(X_1) P(X_2) \prod_{B \in \mathbf{B}(G)} \phi_B} \\ &= \left(\frac{P(X_1) \prod_{B \in \mathbf{B}(i)} \phi_B}{\int_{X_i} P(X_1) \prod_{B \in \mathbf{B}(i)} \phi_B} \right) \cdot \left(\frac{\int_{X_{\setminus N(i)} \cap X_2} P(X_2) \prod_{B \in \setminus \mathbf{B}(i)} \phi_B}{\int_{X_{\setminus N(i)} \cap X_2} P(X_2) \prod_{B \in \setminus \mathbf{B}(i)} \phi_B} \right) \\ &= \frac{P(X_1) \prod_{B \in \mathbf{B}(i)} \phi_B}{\int_{X_i} P(X_1) \prod_{B \in \mathbf{B}(i)} \phi_B} \\ &= \frac{P(X_1) \prod_{B \in \mathbf{B}(i)} \phi_B}{\int_{X_i} P(X_1) \prod_{B \in \mathbf{B}(i)} \phi_B} \cdot \frac{\frac{1}{Z} P(X_2) \prod_{B \in \setminus \mathbf{B}(i)} \phi_B}{\frac{1}{Z} P(X_2) \prod_{B \in \setminus \mathbf{B}(i)} \phi_B} \\ &= P(X_i | \setminus X_i), \end{aligned}$$

from which, we obtain the desired equality. Note that we used the fact that $X_{\mathbf{B}(i)} \cap (X_{\setminus N(i)} \cap X_2) = \emptyset$ from the second to the third and fourth line. \square

C. Proof of Theorem 4

Proof. This proof is constructive. Let's clarify some notations used in this proof. Lower-case bold letter \mathbf{a} is a vector-realization of a set of random variables A . $P(\mathbf{a}_K, \mathbf{c})$ means the probability of a realization where elements appearing on positions indexed by subgraph K are allowed to take random values, while other elements are fixed to value $\mathbf{c} \in \text{dom}(X)$. Note K might be \emptyset . We denote $P_1(X)$ as the equivalency of marginal $P(X_1)$.

First we define the following potential function:

$$\phi_S(X_S = \mathbf{x}_S) = \prod_{Z \subseteq S} \Delta_Z(X_Z = \mathbf{x}_Z)^{(-1)^{|S| - |Z|}},$$

where S is a subset of G , and

$$\Delta_Z(\mathbf{x}_Z) = \begin{cases} \frac{P(\mathbf{x}_Z, \mathbf{c})}{P_1(\mathbf{x}_Z, \mathbf{c}) P_2(\mathbf{x}_Z, \mathbf{c})}, & \exists B \in \mathbf{B}(G), B \subseteq Z, \\ 1 & \text{otherwise,} \end{cases} \quad (6)$$

First we show by construction, the multiplication of all potential functions over all subgraph structures, i.e., $\prod_{S \subseteq G} \phi_S$ will actually give us the **PR**.

Due to the *inclusion-exclusion* principle (see, e.g. [Koller & Friedman \(2009\)](#), 4.4.2.1), it can be shown that

$$\prod_{S \subseteq G} \phi_S(X_S = \mathbf{x}_S) = \Delta_G(\mathbf{x}).$$

If the graph G contains any passage, then by definition $\Delta_G(\mathbf{x}) = \frac{P(\mathbf{x})}{P_1(\mathbf{x})P_2(\mathbf{x})}$, which is exactly the PR. However, if G does not include any passage, meaning X_1 is completely independent of X_2 , then $\Delta_G(\mathbf{x}) = 1$ by definition, which is the exact value that a PR would take in such case.

Second, we show this construction under PMN condition is actually a **GPR**. Specifically, we show if S is not a passage, then $\phi_S(X_S = \mathbf{x}_S) = 1$, i.e. its potential function is nullified.

Obviously, for a ‘‘one-sided S ’’, $X_S \cap X_1 = \emptyset$ or $X_S \cap X_2 = \emptyset$, by definition, $\phi_S = 1$.

Otherwise, if S are ‘‘two-sided’’ but itself is not a passage, we should be able to find two nodes, indexed by $X_u \in X_1 \cap X_S$ and $X_v \in X_2 \cap X_S$, that are not connected by an edge. We may write the potential function for a subgraph S as

$$\phi_S(X_S = \mathbf{x}_S) = \prod_{W \subseteq S \setminus \{u,v\}} \left(\frac{\Delta_W(\mathbf{x}_W) \Delta_{W \cup \{u,v\}}(\mathbf{x}_{W \cup \{u,v\}})}{\Delta_{W \cup \{u\}}(\mathbf{x}_{W \cup \{u\}}) \Delta_{W \cup \{v\}}(\mathbf{x}_{W \cup \{v\}})} \right)^*,$$

where $*$ means we do not care the exact power which can be either -1 or 1, and

$$\begin{aligned} & \frac{\Delta_W(\mathbf{x}_W) \Delta_{W \cup \{u,v\}}(\mathbf{x}_W)}{\Delta_{W \cup \{u\}}(\mathbf{x}_{W \cup \{u\}}) \Delta_{W \cup \{v\}}(\mathbf{x}_{W \cup \{v\}})} = \\ & \frac{P_W P_{W \cup \{u,v\}}}{P_{W \cup \{u\}} P_{W \cup \{v\}}} \cdot \frac{P_{2W \cup \{v\}} P_{2W} P_{1W \cup \{u\}} P_{1W}}{P_{1W} P_{2W} P_{1W \cup \{u\}} P_{2W \cup \{v\}}}, \end{aligned} \quad (7)$$

where we have simplified the notation $P(\mathbf{x}_A, \mathbf{c})$ as P_A . The second factor in (7) is apparently 1. For the first factor in (7), we may divide both the numerator and denominator by $P_W \cdot P_W$. Then it yields $\frac{P(x_u, x_v | \mathbf{x}_W, \mathbf{c})}{P(x_u | \mathbf{x}_W, \mathbf{c}) P(x_v | \mathbf{x}_W, \mathbf{c})}$ which equals to one if and only if $X_u \perp\!\!\!\perp X_v | \{X_u, X_v\}$. This is guaranteed by PMN condition and Proposition 2. \square

D. Proof of Theorem 5

Since the PR is a density ratio between the joint density $p(\mathbf{x}_1, \mathbf{x}_2)$ and the product of two marginals $p(\mathbf{x}_1)p(\mathbf{x}_2)$, and the objective (5) is derived from the same sparsity inducing KLIEP criteria as it was discussed in [Liu et al. \(2015; 2016\)](#). The proof of Theorem 5 follows the primal-dual witness procedure ([Wainwright, 2009](#)).

First, the Assumptions 1, 2 and 3 we have made in Section 5 is essentially the same as those were imposed in Section 3.2 in [Liu et al. \(2016\)](#) (The Hessian of the negative log-likelihood is the sample Fisher information matrix). Then the proof follows the steps established in Section 4, [Liu et al. \(2016\)](#). However, the only thing we need to verify is that $\max_t \|\nabla_{\theta_t} \ell(\theta^*)\|$ is upper-bounded with high probability as $n \rightarrow \infty$. We formally state this in the following lemma:

Lemma 1. *If $\lambda_n \geq \frac{24(2-\alpha)}{\alpha} \cdot \sqrt{\frac{c \log(m^2+m)/2}{n}}$, then*

$$P \left(\max_t \|\nabla_{\theta_t} \ell(\theta^*)\| \geq \frac{\alpha \lambda_n}{4(2-\alpha)} \right) \leq 3 \exp(-c'' n),$$

where c and c'' are some constants.

Proof. For conveniences, let’s denote the approximated PR model $\exp \left(\sum_{u \leq v} \theta_{u,v}^\top \psi(\mathbf{x}_{u,v}) \right) / \hat{N}(\theta)$ as $\hat{g}(\mathbf{x}; \theta)$. Since $\hat{g}(\mathbf{x}; \theta) = \frac{N(\theta)}{\hat{N}(\theta)} g(\mathbf{x}; \theta)$, and $\frac{\hat{N}(\theta)}{N(\theta)} = \frac{1}{\binom{n}{2}} \sum_{j \neq k} g(\mathbf{x}^{[j,k]}; \theta)$ is always bounded by $[C_{\min}, C_{\max}]$, we can see $\hat{g}(\mathbf{x}; \theta)$ is also bounded. For simplicity, we write

$$0 < C'_{\min} \leq \hat{g}(\mathbf{x}; \theta) \leq C'_{\max} < \infty.$$

We have

$$\nabla_{\boldsymbol{\theta}_t} \ell(\boldsymbol{\theta}^*) = - \left[\frac{1}{n} \sum_{i=1}^n \mathbf{f}_t(\mathbf{x}^{(i)}) \right] + \left[\frac{1}{\binom{n}{2}} \sum_{j \leq k} \hat{g}(\mathbf{x}^{[j,k]}; \boldsymbol{\theta}^*) \mathbf{f}_t(\mathbf{x}^{[j,k]}) \right].$$

First we show that $\|\nabla_{\boldsymbol{\theta}_t} \ell(\boldsymbol{\theta}^*)\|$ can be upper-bounded as:

$$\begin{aligned} \|\nabla_{\boldsymbol{\theta}_t} \ell(\boldsymbol{\theta}^*)\| &\leq \underbrace{\left\| -\frac{1}{n} \sum_{i=1}^n \mathbf{f}_t(\mathbf{x}^{(i)}) + \mathbb{E}_p[\mathbf{f}_t(\mathbf{x})] \right\|}_{a_n} + \underbrace{\left\| \frac{1}{\binom{n}{2}} \sum_{j \neq k} \hat{g}(\mathbf{x}^{[j,k]}; \boldsymbol{\theta}^*) \mathbf{f}_t(\mathbf{x}^{[j,k]}) - \frac{1}{\binom{n}{2}} \sum_{j \neq k} g(\mathbf{x}^{[j,k]}; \boldsymbol{\theta}^*) \mathbf{f}_t(\mathbf{x}^{[j,k]}) \right\|}_{b_n} \\ &\quad + \underbrace{\left\| \frac{1}{\binom{n}{2}} \sum_{j \neq k} g(\mathbf{x}^{[j,k]}; \boldsymbol{\theta}^*) \mathbf{f}_t(\mathbf{x}^{[j,k]}) - \mathbb{E}_{p1,p2}[g(\mathbf{x}; \boldsymbol{\theta}^*) \mathbf{f}_t(\mathbf{x})] \right\|}_{\|\mathbf{w}_n\|}, \end{aligned}$$

We now need Hoeffding inequality (Hoeffding, 1963) for bounded-norm vector random variables which has appeared in previous literatures such as (Steinwart & Christmann, 2008): For a set of bounded zero-mean vector-valued random variable $\{\mathbf{y}_i\}_{i=1}^n$, $\|\mathbf{y}\| \leq c$, we have

$$P\left(\left\| \sum_{i=1}^n \mathbf{y}_i \right\| \geq n\epsilon\right) \leq \exp\left(\frac{-n\epsilon^2}{2c^2}\right),$$

for all $\epsilon \geq \frac{2c}{\sqrt{n}}$. Now it is easy to see

$$P(a_n \geq \epsilon) \leq \exp\left(-\frac{2n\epsilon^2}{C_{\mathbf{f}_t, \max}^{\prime 2}}\right) \quad (8)$$

as long as

$$\epsilon \geq \frac{C_{\mathbf{f}_t, \max}'}{2\sqrt{n}}. \quad (9)$$

As to b_n , it can be upper-bounded by

$$\begin{aligned} b_n &= \left\| \frac{1}{\binom{n}{2}} \sum_{j \neq k} \hat{g}(\mathbf{x}^{[j,k]}; \boldsymbol{\theta}^*) \mathbf{f}_t(\mathbf{x}^{[j,k]}) - \frac{1}{\binom{n}{2}} \sum_{j \neq k} \hat{g}(\mathbf{x}^{[j,k]}; \boldsymbol{\theta}^*) \mathbf{f}_t(\mathbf{x}^{[j,k]}) \right\| \\ &= \left\| \frac{\hat{N}(\boldsymbol{\theta}^*)}{N(\boldsymbol{\theta}^*)} \frac{1}{\binom{n}{2}} \sum_{j \neq k} \hat{g}(\mathbf{x}^{[j,k]}; \boldsymbol{\theta}^*) \mathbf{f}_t(\mathbf{x}^{[j,k]}) - \frac{1}{\binom{n}{2}} \sum_{j \neq k} \hat{g}(\mathbf{x}^{[j,k]}; \boldsymbol{\theta}^*) \mathbf{f}_t(\mathbf{x}^{[j,k]}) \right\| \\ &\leq \left\| \frac{1}{\binom{n}{2}} \sum_{j \neq k} \hat{g}(\mathbf{x}^{[j,k]}; \boldsymbol{\theta}^*) \mathbf{f}_t(\mathbf{x}^{[j,k]}) \right\| \cdot \left\| \frac{\hat{N}(\boldsymbol{\theta}^*)}{N(\boldsymbol{\theta}^*)} - 1 \right\| \\ &\leq C_{\max}' C_{\mathbf{f}_t, \max}' \left| \frac{1}{\binom{n}{2}} \sum_{j \neq k} g(\mathbf{x}^{[j,k]}; \boldsymbol{\theta}^*) - 1 \right|, \end{aligned}$$

and due to Hoeffding inequality of the U-statistics (see (Hoeffding, 1963), 5b) we may obtain:

$$P(b_n > \epsilon) < 2 \exp\left(-\frac{2n\epsilon^2}{C_{\max}^{\prime 2} C_{\max}^{\prime 2} C_{\mathbf{f}_t, \max}^{\prime 2}}\right). \quad (10)$$

As to \mathbf{w}_n , we first bound its i -th element $w_{i,n}$ using Hoeffding inequality for U-statistics,

$$P(|w_{i,n}| \geq \epsilon) \leq 2 \exp\left(-\frac{2nb\epsilon^2}{C_{\max}^2 C_{\mathbf{f}_t, \max}^2}\right),$$

thus by using the union bound, we have

$$P(\|\mathbf{w}_n\|_{\infty} \geq \epsilon) \leq 2b \exp\left(-\frac{2nb\epsilon^2}{C_{\max}^2 C_{\mathbf{f}_t, \max}^2}\right),$$

and since $\|\mathbf{w}_n\| \leq \sqrt{b}\|\mathbf{w}_n\|_{\infty}$, we have

$$P(\|\mathbf{w}_n\| \geq \epsilon) \leq P(\sqrt{b}\|\mathbf{w}_n\|_{\infty} \geq \epsilon) \leq 2b \exp\left(-\frac{2n\epsilon^2}{C_{\max}^2 C_{\mathbf{f}_t, \max}^2}\right). \quad (11)$$

Therefore, combining (8), (10) and (11):

$$P(\|\nabla_{\boldsymbol{\theta}_t} \ell(\boldsymbol{\theta}^*)\| \geq 3\epsilon) \leq P(a_n + b_n + c_n \geq 3\epsilon) \leq c'' \exp\left(-\frac{n\epsilon^2}{c'}\right),$$

where c' is a constant defined as $c' = \max\left(\frac{1}{2}C_{\max}^2 C_{\max}^2 C_{\mathbf{f}_t, \max}^{\prime 2}, \frac{1}{2}C_{\max}^2 C_{\mathbf{f}_t, \max}^2, \frac{1}{2}C_{\mathbf{f}_t, \max}^{\prime 2}\right)$, and $c'' = 2b + 3$, given $\epsilon \geq \frac{2C_{\mathbf{f}_t, \max}^{\prime}}{\sqrt{n}}$. Applying the union-bound for all $t \in S \cup S^c$,

$$P\left(\max_{t \in S \cup S^c} \|\nabla_{\boldsymbol{\theta}_t} \ell(\boldsymbol{\theta}^*)\| \geq 3\epsilon\right) \leq \frac{c''(m^2 + m)}{2} \exp\left(-\frac{n\epsilon^2}{c'}\right),$$

$$P\left(\max_{t \in S \cup S^c} \|\nabla_{\boldsymbol{\theta}_t} \ell(\boldsymbol{\theta}^*)\| \geq \frac{\alpha\lambda_n}{4(2-\alpha)}\right) \leq \frac{c''(m^2 + m)}{2} \exp\left(-\left(\frac{\alpha\lambda_n}{12(2-\alpha)}\right)^2 \frac{n}{c'}\right),$$

and when $\lambda_n \geq \frac{24(2-\alpha)}{\alpha} \sqrt{\frac{c' \log(m^2 + m)/2}{n}}$,

$$P\left(\max_{t \in S \cup S^c} \|\nabla_{\boldsymbol{\theta}_t} \ell(\boldsymbol{\theta}^*)\| \geq \frac{\alpha\lambda_n}{4(2-\alpha)}\right) \leq c''' \exp(-c'''n),$$

where c''' is a constant. Assume that $\log \frac{m^2 + m}{2} > 1$ and we set λ_n as

$$\lambda_n \geq \frac{24(2-\alpha)}{\alpha} \sqrt{\frac{(c' + C_{\mathbf{f}_t, \max}^2) \log(m^2 + m)/2}{n}},$$

then (9), the condition of using vector Hoeffding-inequality is satisfied. \square

Given Lemma 1, we may obtain other technical results, such as the estimation error bound, using the same proof as it was demonstrated in Section 4, Liu et al. (2016).

E. Experimental Settings

We measure the performance of three methods using True Positive Rate (TPR) and True Negative Rate (TNR) that are used in Zhao et al. (2014). The TPR and TFR are defined as:

$$\text{TPR} = \frac{\sum_{t' \in S} \delta(\hat{\boldsymbol{\theta}}_{t'} \neq \mathbf{0})}{\sum_{t' \in S} \delta(\boldsymbol{\theta}_{t'}^* \neq \mathbf{0})}, \quad \text{TNR} = \frac{\sum_{t'' \in S^c} \delta(\hat{\boldsymbol{\theta}}_{t''} = \mathbf{0})}{\sum_{t'' \in S^c} \delta(\boldsymbol{\theta}_{t''}^* = \mathbf{0})},$$

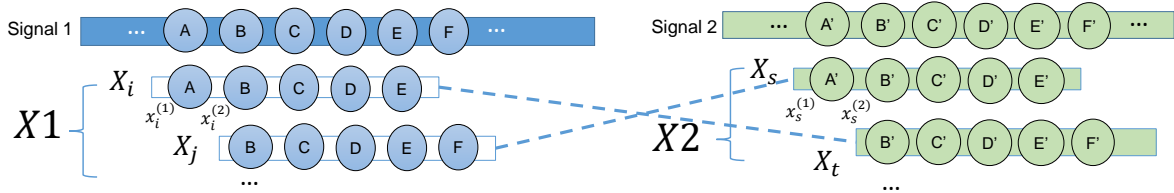


Figure 7. The illustration of sequence matching problem formulation.

where δ is the indicator function.

The differential learning method (Zhao et al., 2014) used in Section 6.1 learns the difference between two precision matrices. In our setting, if one can learn the difference between the precision matrices of $p(\mathbf{x})$ and $p(\mathbf{x}_1)p(\mathbf{x}_2)$, one can figure out all edges that go across two groups (\mathbf{x}_1 and \mathbf{x}_2).

This method requires sample covariance matrices of $p(\mathbf{x})$ and $p(\mathbf{x}_1)p(\mathbf{x}_2)$ respectively. The sample covariance of $p(\mathbf{x})$ is easy to compute given joint samples. However, to obtain the sample covariance of $p(\mathbf{x}_1)p(\mathbf{x}_2)$, we would again need the U-statistics (Hoeffding, 1963) introduced in line Section 4. We may approximate the u, v -th element of the covariance matrix of $p(\mathbf{x}_1)p(\mathbf{x}_2)$ using the formula: $\Sigma_{u,v} = \frac{1}{\binom{n}{2}} \sum_{j \neq k} x_v^{[j,k]} x_u^{[j,k]}$, assuming the joint distribution has zero mean.

F. Illustration of Sequence Matching

We plot the illustrations of our sequence matching problem formulation from two sequences in Figure 7.