

A MEAN PERCENTILE RANK

We begin our definition of MPR by defining percentile rank (PR). First, given a set J , let $p_{i,J} = \Pr(J \cup \{i\} \mid J)$. The percentile rank of an item i given a set J is defined as

$$\text{PR}_{i,J} = \frac{\sum_{i' \notin J} \mathbf{1}(p_{i,J} \geq p_{i',J})}{|\mathcal{Y} \setminus J|} \times 100\%$$

where $\mathcal{Y} \setminus J$ indicates those elements in the ground set \mathcal{Y} that are not found in J .

For our evaluation, given a test set Y , we select a random element $i \in Y$ and compute $\text{PR}_{i,Y \setminus \{i\}}$. We then average over the set of all test instances \mathcal{T} to compute the mean percentile rank (MPR):

$$\text{MPR} = \frac{1}{|\mathcal{T}|} \sum_{Y \in \mathcal{T}} \text{PR}_{i,Y \setminus \{i\}}.$$

B HYPERPARAMETERS FOR EXPERIMENTS IN TABLE 2

Preventing numerical instabilities: The first term on the right side of Eq. 2 will be singular whenever $|Y_i| > K$, where Y_i is an observed subset. Therefore, to address this in practice we set K to the size of the largest subset observed in the data, K' , as in Gartrell et al. (2017). However, this does not entirely address the issue, as the first term on the right side of Eq. 2 may still be singular even when $|Y_i| \leq K$. In this case though, we know that we are not at a maximum, since the value of the objective function is $-\infty$. Numerically, to prevent such singularities, in our implementation we add a small ϵI correction to each L_{Y_i} when optimizing Eq. 2 (we set $\epsilon = 10^{-5}$ in our experiments).

We perform a grid search using a held-out validation set to select the best performing hyperparameters for each model and dataset. The hyperparameter settings used for each model and dataset are described below.

Symmetric low-rank DPP (Gartrell et al., 2016). For this model, we use K for the number of item feature dimensions for the symmetric component \mathbf{V} , and α for the regularization hyperparameter for \mathbf{V} . We use the following hyperparameter settings:

- Both Amazon datasets: $K = 30, \alpha = 0$.
- UK Retail dataset: $K = 100, \alpha = 1$.
- Instacart dataset: $K = 100, \alpha = 0.001$.
- Million Song dataset: $K = 150, \alpha = 0.0001$.

Baseline NDPP (Gartrell et al., 2019). For this model, to ensure consistency with the notation used in Gartrell et al. (2019), we use D to denote the number of item feature dimensions for the symmetric component \mathbf{V} , and D' to denote the number of item feature dimensions for the nonsymmetric components, \mathbf{B} and \mathbf{C} . As described in Gartrell et al. (2019), α is the regularization hyperparameter for the \mathbf{V} , while β and γ are the regularization hyperparameters for \mathbf{B} and \mathbf{C} , respectively. We use the following hyperparameter settings:

- Both Amazon datasets: $D = 30, \alpha = 0$.
- Amazon apparel dataset: $D' = 30$.
- Amazon three-category dataset: $D' = 100$.
- UK Retail dataset: $D = 100, D' = 20, \alpha = 1$.
- All datasets: $\beta = \gamma = 0$.

Scalable NDPP. As described in Section 3, we use K to denote the number of item feature dimensions for the symmetric component \mathbf{V} and the dimensionality of the nonsymmetric component \mathbf{C} . α is the regularization hyperparameter. We use the following hyperparameter settings:

- Amazon apparel dataset: $K = 30, \alpha = 0$.
- Amazon three-category dataset: $K = 100, \alpha = 1$.
- UK dataset: $K = 100, \alpha = 0.01$.

- Instacart dataset: $K = 100, \alpha = 0.001$.
- Million Song dataset: $K = 150, \alpha = 0.01$.

For all of the above model configurations we use a batch size of 200 during training, except for the scalable NDPPs trained on the Amazon apparel, Amazon three-category, Instacart, and Million Song datasets, where a batch size of 800 is used.

C TRAINING TIME

In Fig. 1, we report the wall-clock training time of the decomposition of Gartrell et al. (2019) (NDPP) and our scalable NDPP for the Amazon: 3-category (Fig. 1(a)) and UK Retail (Fig. 1(b)) datasets. For the Amazon: 3-category dataset, both approaches show comparable results, with the scalable NDPP converging 1.07 times faster than NDPP. But for the UK Retail dataset, which has a much larger ground set, our scalable NDPP achieves convergence about 8.31 times faster. We do not have a timing comparison for the Instacart dataset because the model with the decomposition of Gartrell et al. (2019) cannot be trained on this dataset due to prohibitive memory and computational costs.

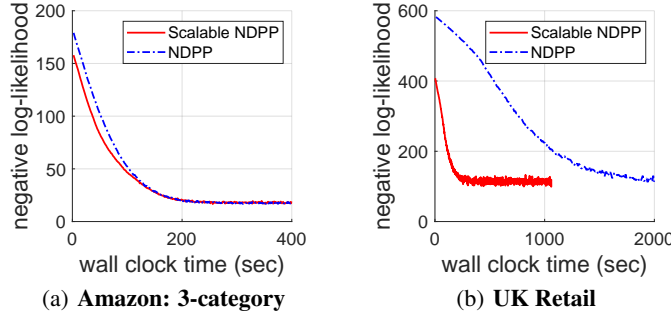


Figure 1: The negative log-likelihood of the training set for Gartrell et al. (2019)’s NDPP (blue, dashed) and our scalable NDPP (red, solid) versus wall-clock time for the (a) Amazon: 3-category and (b) UK Retail datasets.

D BENCHMARK ALGORITHMS FOR MAP INFERENCE

We test following approximate algorithms for MAP inference:

Greedy local search. This algorithm starts from the output of greedy, Y^G , and replaces $i \in Y^G$ with $j \notin Y^G$ that gives the maximum improvement of the determinant, if such i, j exist. Kathuria & Deshpande (2016) showed that running the search for such a swap $O(k^2 \log(k/\varepsilon))$ times with an accuracy parameter ε gives a tight approximation guarantee for MAP inference for symmetric DPPs. We set the number of swaps to $\lfloor k^2 \log(10k) \rfloor$ for $\varepsilon = 0.1$ and use greedy local search as a baseline, since it is strictly an improvement on the greedy solution. The proposed greedy conditioning can be used for fast greedy local search. Specifically, for each $i \in Y^G$, Algorithm 1 can compute marginal improvements conditioned by $Y^G \setminus \{i\}$ in time $O(MKk)$ thus its runtime can be $O(MKk^4 \log(k/\varepsilon))$. However, it is the slowest among all of our benchmark algorithms.

Stochastic greedy. This algorithm computes marginal gains of a few items chosen uniformly at random and selects the best among them. Mirzasoleiman et al. (2015) proved that $(M/k) \log(1/\varepsilon)$ samples are enough to guarantee an $(1 - 1/e - \varepsilon)$ -approximation ratio for submodular functions (i.e., symmetric DPPs). We choose $\varepsilon = 0.1$ and set the number of samples to $\lfloor (M/k) \log(10) \rfloor$. Under this setting, the time complexity of stochastic greedy is $O(MKk^2 \log(1/\varepsilon))$, which is better than the naïve exact greedy algorithm. However, we note that it is worse than that of our efficient greedy implement (Algorithm 1). This is because the stochastic greedy uses different random samples for every iteration and this does not take advantage of the amortized computations in Lemma 2. In our experiments, we simply modify line 10 in Algorithm 1 for stochastic greedy (argmax is operated

on a random subset of marginal gains), hence it can run in $O(MKk + (M/k) \log(1/\varepsilon))$ time. In practice, we observe that stochastic greedy is slightly slower than exact greedy due to the additional costs of random sampling process.

MCMC sampling. We also compare inference algorithms with sampling from a nonsymmetric DPP. To the best of our knowledge, exact sampling of a non-Hermitian DPP was studied in Poulson (2019), which requires the Cholesky decomposition with $O(M^3)$ complexity. This is infeasible for a large M . To resolve this, Markov Chain Monte-Carlo (MCMC) based sampling is preferred Li et al. (2016) for symmetric DPPs. In particular, we consider a Gibbs sampling for k -DPP, which begins with a random subset Y with size k , and picks $i \in Y$ and $j \notin Y$ uniformly at random. Then, it swaps them with probability

$$\frac{\det(\mathbf{L}_{Y \cup \{j\} \setminus \{i\}})}{\det(\mathbf{L}_{Y \cup \{j\} \setminus \{i\}}) + \det(\mathbf{L}_Y)} \quad (11)$$

and repeat this process for several steps. Li et al. (2016) showed that $O(Nk \log(k/\varepsilon))$ swaps are enough to approximate the ground-truth distribution under symmetric DPPs. However, for a fair runtime comparison to Algorithm 1, we set the number of swaps to $\lfloor 3N/K \rfloor$.

We provide the wall-clock time of the above algorithms for real-world datasets in Table 4. Observe that the greedy algorithm is the fastest method for all datasets except Million Song. For Million Song, MCMC sampling is faster than other approaches, but it has much larger relative errors in terms of log-determinant (see Table 3), which is not suitable for our purposes.

Table 4: Wall-clock time (in milliseconds) of MAP inference algorithms on NDPPs learned from real-world datasets.

| Algorithms | Amazon: Apparel | Amazon: 3-category | UK Retail | Instacart | Million Song |
|--|-----------------|--------------------|----------------|-----------------|------------------|
| Greedy local search (Kathuria & Deshpande, 2016) | 5.78 ms | 9.67 ms | 58.74 ms | 1.024 s | 7.277 s |
| Greedy (Algorithm 1) | 0.14 ms | 0.34 ms | 1.60 ms | 36.16 ms | 338.09 ms |
| Stochastic greedy (Mirzasoleiman et al., 2015) | 0.25 ms | 0.47 ms | 1.79 ms | 36.94 ms | 348.67 ms |
| MCMC sampling (Li et al., 2016) | 0.19 ms | 0.35 ms | 2.85 ms | 42.85 ms | 303.20 ms |

E COROLLARY OF THEOREM 2

Theorem 2 requires the technical condition $\sigma_{\min} > 1$ but in practice there is no particular evidence that this condition holds. While this condition can be achieved by multiplying \mathbf{L} by a constant, this leads to a (potentially large) additive term in Eq. 10. Here, we provide Corollary 1 which excludes the $\sigma_{\min} > 1$ assumption from Theorem 2, and quantifies this additive term.

Corollary 1. Consider a nonsymmetric low-rank DPP $\mathbf{L} = \mathbf{V}\mathbf{V}^\top + \mathbf{B}\mathbf{C}\mathbf{B}^\top$, where \mathbf{V}, \mathbf{B} are of rank K , and $\mathbf{C} \in \mathbb{R}^{K \times K}$. Given a cardinality budget k , let σ_{\min} and σ_{\max} denote the smallest and largest singular values of \mathbf{L}_Y for all $Y \subseteq [M]$ and $|Y| \leq 2k$. Let $\kappa := \sigma_{\max}/\sigma_{\min}$. Then,

$$\log \det(\mathbf{L}_{Y^G}) \geq \frac{4(1 - e^{-1/4})}{2 \log \kappa + 1} \log \det(\mathbf{L}_{Y^*}) - \left(1 - \frac{4(1 - e^{-1/4})}{2 \log \kappa + 1}\right) k (1 - \log \sigma_{\min}) \quad (12)$$

where Y^G is the output of Algorithm 1 and Y^* is the optimal solution of MAP inference in Eq. 4.

The proof of Corollary 1 is provided in Appendix G.5. Note that instead of $\log(\sigma_{\max})/\log(\sigma_{\min})$, Corollary 1 has a $\log(\sigma_{\max}/\sigma_{\min})$ term in the denominator.

F PERFORMANCE GUARANTEE FOR GREEDY MAP INFERENCE

The matrices learned on real datasets are too large to compute the exact MAP solution, but we can compute exact MAP for small matrices. In this section, we explore the performance of the greedy algorithm studied in Theorem 2 for 5×5 synthetic kernel matrices. More formally, we first pick $K = 3$ singular values s_1, s_2, s_3 from a kernel learned for the “Amazon: 3-category” dataset (a plot of these singular values can be seen in Fig. 2(c)) and generate $\mathbf{L} = \mathbf{V}_1 \text{diag}([s_1, s_2, s_3]) \mathbf{V}_2^\top$, where $\mathbf{V}_1, \mathbf{V}_2 \in \mathbb{R}^{5 \times 3}$ are random orthonormal matrices. To ensure that \mathbf{L} is a P_0 matrix, we repeatedly

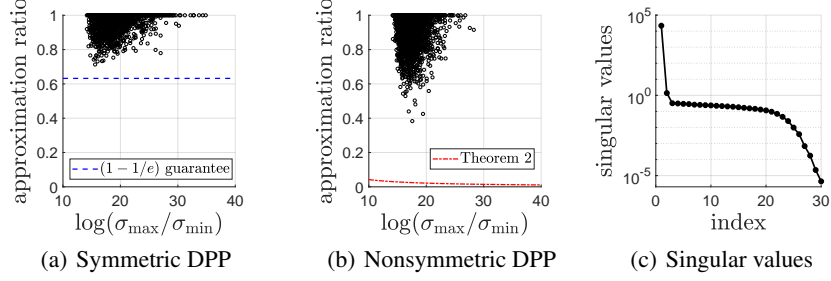


Figure 2: Approximation ratios of greedy with respect to different values of $\log(\sigma_{\max}/\sigma_{\min})$ from Corollary 1 under (a) symmetric DPP and (b) nonsymmetric DPP. (c) The singular values of the kernels learned for the “Amazon: 3-category” dataset. We construct 10,000 random P_0 matrices $L \in \mathbb{R}^{5 \times 5}$, with rank $K = 3$, whose singular values are from the learned kernels.

sample V_1, V_2 until all principal minors of L are nonnegative. We also evaluate the performance of the symmetric DPP, where the kernel matrices are generated similarly to the NDPP, except we set $V_1 = V_2$. We set $k = 3$ and generate 10,000 random kernels for both symmetric DPPs and NDPPs.

The results for symmetric and nonsymmetric DPPs are shown in Fig. 2(a) and Fig. 2(b), respectively. We plot the approximation ratio of Algorithm 1, i.e., $\log \det(L_{YG}) / \log \det(L_{Y*})$, with respect to $\log(\sigma_{\max}/\sigma_{\min})$, from Corollary 1. We observe that the greedy algorithm for both often shows approximation ratios close to 1. However, the worst-case ratio for NDPPs is worse than that of symmetric DPPs; $\log \det(L_Y)$ for $L \in P_0^+$ is non-submodular, and the greedy algorithm with a nonsubmodular function does not have as tight of a worst-case bound as in the symmetric case.

G PROOFS

G.1 PROOF OF LEMMA 1

Lemma 1. *Let $\ell \leq M$ be an even integer and let $\mathbf{A} \in \mathbb{R}^{M \times M}$ be a skew-symmetric matrix with rank ℓ . Then, there exist $\mathbf{B} \in \mathbb{R}^{M \times \ell}$ and positive numbers $\lambda_1, \dots, \lambda_{\ell/2}$, such that $\mathbf{A} = \mathbf{B}\mathbf{C}\mathbf{B}^\top$, where $\mathbf{C} \in \mathbb{R}^{\ell \times \ell}$ is the block-diagonal matrix with $(\ell/2)$ diagonal blocks of size 2 given by Σ_i , $i = 1, \dots, \ell/2$.*

Proof. First, $\mathbf{A} = \mathbf{P}\mathbf{\Sigma}\mathbf{P}^\top$ for some orthogonal matrix $\mathbf{P} \in \mathbb{R}^{M \times M}$ and

$$\mathbf{\Sigma} = \begin{pmatrix} 0 & \lambda_1 & & & & \\ -\lambda_1 & 0 & & & & \\ & & 0 & \lambda_2 & & \\ & & -\lambda_2 & 0 & & \\ & & & & \ddots & \\ & & & & & 0 & \lambda_{\ell/2} \\ & & & & & -\lambda_{\ell/2} & 0 \\ & & & & & & & 0 \\ & & & & & & & & \ddots \\ & & & & & & & & & 0 \end{pmatrix} \quad (13)$$

(see, e.g., (Thompson, 1988, Proposition 2.1), which is easily extended to the case when M is odd).

Let \mathbf{C} be the $\ell \times \ell$ submatrix of $\mathbf{\Sigma}$ obtained by keeping its first ℓ rows and columns and let $\mathbf{Q} = \begin{pmatrix} \mathbf{I}_\ell \\ 0 \end{pmatrix}$, where \mathbf{I}_ℓ is the $\ell \times \ell$ identity matrix. Then, $\mathbf{\Sigma} = \mathbf{Q}\mathbf{C}\mathbf{Q}^\top$ and one can write $\mathbf{A} = \mathbf{P}\mathbf{Q}\mathbf{C}\mathbf{Q}^\top\mathbf{P}^\top$. Setting $\mathbf{B} = \mathbf{P}\mathbf{Q}$ proves the lemma. \square

G.2 PROOF OF THEOREM 1

Theorem 1. *Given an NDPP with kernel $\mathbf{L} = \mathbf{V}\mathbf{V}^\top + \mathbf{B}\mathbf{C}\mathbf{B}^\top$, parameterized by \mathbf{V} of rank K , \mathbf{B} of rank K , and a $K \times K$ matrix \mathbf{C} , we can compute the regularized log-likelihood (Eq. 2) and its gradient in $O(MK^2 + K^3 + nK'^3)$ time, where K' is the size of the largest of the n training subsets.*

Proof. We first show that the log-likelihood can be computed in time linear in M . Using the matrix determinant lemma, one can easily verify that the DPP normalization term can be computed as

$$\det(\mathbf{I} + \mathbf{L}) = \det\left(\mathbf{I} + (\mathbf{V} \quad \mathbf{B}\mathbf{C}) \begin{pmatrix} \mathbf{V}^\top \\ \mathbf{B}^\top \end{pmatrix}\right) = \det\left(\mathbf{I}_{2K} + \begin{pmatrix} \mathbf{V}^\top \\ \mathbf{B}^\top \end{pmatrix} (\mathbf{V} \quad \mathbf{B}\mathbf{C})\right) \quad (14)$$

where \mathbf{I}_{2K} is the identity matrix with dimension $2K$. As Eq. 14 requires a matrix-multiplication between $(2K) \times M$ matrices and the determinant of $(2K) \times (2K)$ matrices, this allows us to transform a $O(M^3)$ operation into an $O(MK^2 + K^3)$ one.

Having established that the normalization term in the likelihood can be computed in $O(MK^2 + K^3)$ time, we proceed with characterizing the complexity of the other terms in the likelihood. The first term in Eq. 2 consists of determinants of size $|Y_i|$. Assuming that these never exceed size K' , each can be computed in at most $O(K'^3)$ time. The regularization term is a simple sum of norms that can be computed in $O(MK)$ time. Therefore, the full regularized log-likelihood can be computed in $O(MK^2 + K^3 + nK'^3)$ time.

To prove that the gradient of the log-likelihood can be computed in time linear in M , we begin by showing that the logarithm of DPP normalization term can be factorized as follows:

$$Z = \log \det(\mathbf{I} + \mathbf{L}) \quad (15)$$

$$= \log \det \left(\mathbf{I}_{2K} + \begin{pmatrix} \mathbf{V}^\top \\ \mathbf{B}^\top \end{pmatrix} (\mathbf{V} \ \mathbf{B}) \begin{pmatrix} \mathbf{I}_K & \mathbf{0} \\ \mathbf{0} & \mathbf{C} \end{pmatrix} \right) \quad (16)$$

$$= \log \det \left(\begin{pmatrix} \mathbf{I}_K & \mathbf{0} \\ \mathbf{0} & \mathbf{C}^{-1} \end{pmatrix} + \begin{pmatrix} \mathbf{V}^\top \\ \mathbf{B}^\top \end{pmatrix} (\mathbf{V} \ \mathbf{B}) \right) + \log \det \begin{pmatrix} \mathbf{I}_K & \mathbf{0} \\ \mathbf{0} & \mathbf{C} \end{pmatrix} \quad (17)$$

$$= \log \det \begin{pmatrix} \mathbf{I}_K + \mathbf{V}^\top \mathbf{V} & \mathbf{V}^\top \mathbf{B} \\ \mathbf{B}^\top \mathbf{V} & \mathbf{C}^{-1} + \mathbf{B}^\top \mathbf{B} \end{pmatrix} + \log \det(\mathbf{C}) \quad (18)$$

$$= \log \det(\mathbf{I}_K + \mathbf{V}^\top \mathbf{V}) + \log \det(\mathbf{C}^{-1} + \mathbf{B}^\top (\mathbf{I} - \mathbf{V}(\mathbf{I}_K + \mathbf{V}^\top \mathbf{V})^{-1} \mathbf{V}^\top) \mathbf{B}) + \log \det(\mathbf{C}) \quad (19)$$

where Eq. 17 follows from the determinant commutativity (i.e., $\det(\mathbf{AB}) = \det(\mathbf{A})\det(\mathbf{B})$) and Eq. 18 and Eq. 19 come from the Schur's determinant identity³. For simplicity, we write $\mathbf{X} = \mathbf{I} - \mathbf{V}(\mathbf{I}_K + \mathbf{V}^\top \mathbf{V})^{-1} \mathbf{V}^\top$ and $(\mathbf{C}^{-1})^\top = \mathbf{C}^{-\top}$ and note that \mathbf{X} depends only on \mathbf{V} .

The gradient of Z has three parts: $\nabla Z = (\nabla_{\mathbf{V}} Z, \nabla_{\mathbf{B}} Z, \nabla_{\mathbf{C}} Z)$ where each can be computed as

$$\begin{aligned} \nabla_{\mathbf{V}} Z &= \nabla_{\mathbf{V}} \log \det(\mathbf{I}_K + \mathbf{V}^\top \mathbf{V}) + \nabla_{\mathbf{V}} \log \det(\mathbf{C}^{-1} + \mathbf{B}^\top \mathbf{X} \mathbf{B}) \\ &= 2\mathbf{V}(\mathbf{I}_K + \mathbf{V}^\top \mathbf{V})^{-1} \end{aligned} \quad (20)$$

$$- \mathbf{X} \mathbf{B}((\mathbf{C}^{-1} + \mathbf{B}^\top \mathbf{X} \mathbf{B})^{-1} + (\mathbf{C}^{-\top} + \mathbf{B}^\top \mathbf{X} \mathbf{B})^{-1}) \mathbf{B}^\top \mathbf{X} \mathbf{V} \quad (21)$$

$$\nabla_{\mathbf{B}} Z = \nabla_{\mathbf{B}} \log \det(\mathbf{C}^{-1} + \mathbf{B}^\top \mathbf{X} \mathbf{B}) \quad (22)$$

$$= \mathbf{X} \mathbf{B}((\mathbf{C}^{-1} + \mathbf{B}^\top \mathbf{X} \mathbf{B})^{-1} + (\mathbf{C}^{-\top} + \mathbf{B}^\top \mathbf{X} \mathbf{B})^{-1}) \quad (23)$$

$$\nabla_{\mathbf{C}} Z = \nabla_{\mathbf{C}} \log \det(\mathbf{C}) + \nabla_{\mathbf{C}} \log \det(\mathbf{C}^{-1} + \mathbf{X}) \quad (24)$$

$$= \mathbf{C}^{-\top} - \mathbf{C}^{-\top}(\mathbf{C}^{-1} + \mathbf{B}^\top \mathbf{X} \mathbf{B})^{-\top} \mathbf{C}^{-\top} \quad (25)$$

Observe that \mathbf{X} combines a $M \times M$ identity matrix with $M \times K$ matrices, hence multiplying it with a $M \times K$ matrix (e.g., $\mathbf{X} \mathbf{V}$ or $\mathbf{X} \mathbf{B}$) can be computed in $O(MK^2)$ time. Since each of the remaining matrix inverses in Eq. 21, Eq. 23, and Eq. 25 involve a $K \times K$ matrix inverse, with a cost of $O(K^3)$ operations, we have a net computational cost of $O(MK^2 + K^3)$ for computing $\nabla \log \det(\mathbf{I} + \mathbf{L})$.

The gradient of the first term in Eq. 2 involves computing gradients of determinants of size at most K' , which results in size K' matrix inverses, since for a matrix \mathbf{A} , $\frac{\partial}{\partial A_{ij}}(\log \det(\mathbf{A})) = (\mathbf{A}^{-1})_{ij}^\top$. Each of these inverses can be computed in $O(K'^3)$ time. The gradient of the simple sum-of-norms regularization term can be computed in $O(MK)$ time. Therefore, combining these results with the results above for the complexity of the gradient of the normalization term, we have the following overall complexity of the gradient for the full log-likelihood: $O(MK^2 + K^3 + nK'^3)$. \square

G.3 PROOF OF LEMMA 2

Lemma 2. Given $\mathbf{B} \in \mathbb{R}^{M \times K}$, $\mathbf{C} \in \mathbb{R}^{K \times K}$, and $Y = \{a_1, \dots, a_k\} \subseteq [M]$, let $\mathbf{b}_i \in \mathbb{R}^{1 \times K}$ be the i -th row in \mathbf{B} and $\mathbf{B}_Y \in \mathbb{R}^{|Y| \times K}$ be a matrix containing rows in \mathbf{B} indexed by Y . Then, it holds that

$$\mathbf{B}_Y^\top (\mathbf{B}_Y \mathbf{C} \mathbf{B}_Y^\top)^{-1} \mathbf{B}_Y = \sum_{j=1}^k \mathbf{p}_j^\top \mathbf{q}_j, \quad (6)$$

where row vectors $\mathbf{p}_j, \mathbf{q}_j \in \mathbb{R}^{1 \times K}$ for $j = 1, \dots, k$ satisfy $\mathbf{p}_1 = \mathbf{b}_{a_1} / (\mathbf{b}_{a_1} \mathbf{C} \mathbf{b}_{a_1}^\top)$, $\mathbf{q}_1 = \mathbf{b}_{a_1}$, and

$$\mathbf{p}_{j+1} = \frac{\mathbf{b}_{a_j} - \mathbf{b}_{a_j} \mathbf{C}^\top \sum_{i=1}^j \mathbf{q}_i^\top \mathbf{p}_i}{\mathbf{b}_{a_j} \mathbf{C} (\mathbf{b}_{a_j} - \mathbf{b}_{a_j} \mathbf{C}^\top \sum_{i=1}^j \mathbf{q}_i^\top \mathbf{p}_i)^\top}, \quad \mathbf{q}_{j+1} = \mathbf{b}_{a_j} - \mathbf{b}_{a_j} \mathbf{C} \sum_{i=1}^j \mathbf{p}_i^\top \mathbf{q}_i. \quad (7)$$

³ $\det \begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{pmatrix} = \det(\mathbf{A}) \det(\mathbf{D} - \mathbf{C} \mathbf{A}^{-1} \mathbf{B}).$

Proof. We prove by induction on k . When $k = 1$, the result is trivial because

$$B_Y^\top (B_Y C B_Y^\top)^{-1} B_Y = \mathbf{b}_{a_1}^\top (\mathbf{b}_{a_1} C \mathbf{b}_{a_1}^\top)^{-1} \mathbf{b}_{a_1} = \mathbf{p}_1^\top \mathbf{q}_1. \quad (26)$$

Now we assume that the statement holds for $k - 1$. Let $Y := \{a_1, \dots, a_{k-1}\}$ and $a := a_k$. From the inductive hypothesis, it holds

$$B_Y^\top (B_Y C B_Y^\top)^{-1} B_Y = \sum_{j=1}^{k-1} \mathbf{p}_j^\top \mathbf{q}_j. \quad (27)$$

Now we write

$$B_{Y \cup \{a\}}^\top \left(B_{Y \cup \{a\}} C B_{Y \cup \{a\}}^\top \right)^{-1} B_{Y \cup \{a\}} \quad (28)$$

$$= B_{Y \cup \{a\}}^\top \left(\begin{pmatrix} B_Y \\ \mathbf{b}_a \end{pmatrix} C \begin{pmatrix} B_Y^\top & \mathbf{b}_a^\top \end{pmatrix} \right)^{-1} B_{Y \cup \{a\}} \quad (29)$$

$$= \begin{pmatrix} B_Y^\top & \mathbf{b}_a^\top \end{pmatrix} \begin{pmatrix} B_Y C B_Y^\top & B_Y C \mathbf{b}_a^\top \\ \mathbf{b}_a C B_Y^\top & \mathbf{b}_a C \mathbf{b}_a^\top \end{pmatrix}^{-1} \begin{pmatrix} B_Y \\ \mathbf{b}_a \end{pmatrix}. \quad (30)$$

To handle the inverse matrix we employ the Schur complement, which yields

$$\begin{pmatrix} X & \mathbf{y} \\ \mathbf{z} & w \end{pmatrix}^{-1} = \begin{pmatrix} X^{-1} & \mathbf{0} \\ \mathbf{0} & 0 \end{pmatrix} + \frac{1}{(w - \mathbf{z} X^{-1} \mathbf{y})^{-1}} \begin{pmatrix} X^{-1} \mathbf{y} \mathbf{z} X^{-1} & -X^{-1} \mathbf{y} \\ -\mathbf{z} X^{-1} & 1 \end{pmatrix} \quad (31)$$

for any non-singular square matrix $X \in \mathbb{R}^{k \times k}$, column vector $\mathbf{y} \in \mathbb{R}^k$ and row vector $\mathbf{z} \in \mathbb{R}^{1 \times k}$, unless $(w - \mathbf{z} X^{-1} \mathbf{y})^{-1} = 0$. Applying this, we have

$$\begin{pmatrix} B_Y C B_Y^\top & B_Y C \mathbf{b}_a^\top \\ \mathbf{b}_a C B_Y^\top & \mathbf{b}_a C \mathbf{b}_a^\top \end{pmatrix}^{-1} = \begin{pmatrix} (B_Y C B_Y^\top)^{-1} & \mathbf{0} \\ \mathbf{0} & 0 \end{pmatrix} + \frac{1}{\mathbf{b}_a C \mathbf{b}_a^\top - \mathbf{b}_a C B_Y^\top (B_Y C B_Y^\top)^{-1} B_Y C \mathbf{b}_a^\top} \begin{pmatrix} (B_Y C B_Y^\top)^{-1} B_Y C \mathbf{b}_a^\top \mathbf{b}_a C B_Y^\top (B_Y C B_Y^\top)^{-1} & -(B_Y C B_Y^\top)^{-1} B_Y C \mathbf{b}_a^\top \\ -\mathbf{b}_a C B_Y^\top (B_Y C B_Y^\top)^{-1} & 1 \end{pmatrix} \quad (32)$$

Substituting Eq. 32 into Eq. 30, we obtain

$$B_{Y \cup \{a\}}^\top \left(B_{Y \cup \{a\}} C B_{Y \cup \{a\}}^\top \right)^{-1} B_{Y \cup \{a\}} \quad (33)$$

$$= B_Y^\top (B_Y C B_Y^\top)^{-1} B_Y + \frac{(\mathbf{b}_a^\top - B_Y^\top (B_Y C B_Y^\top)^{-1} B_Y C \mathbf{b}_a^\top) (\mathbf{b}_a - \mathbf{b}_a C B_Y^\top (B_Y C B_Y^\top)^{-1} B_Y)}{\mathbf{b}_a C (\mathbf{b}_a^\top - B_Y^\top (B_Y C B_Y^\top)^{-1} B_Y C \mathbf{b}_a^\top)} \quad (34)$$

$$= \sum_{j=1}^{k-1} \mathbf{p}_j^\top \mathbf{q}_j + \frac{(\mathbf{b}_a^\top - \sum_{j=1}^{k-1} \mathbf{p}_j^\top \mathbf{q}_j C \mathbf{b}_a^\top) (\mathbf{b}_a - \mathbf{b}_a C \sum_{j=1}^{k-1} \mathbf{p}_j^\top \mathbf{q}_j)}{\mathbf{b}_a C (\mathbf{b}_a^\top - \sum_{j=1}^{k-1} \mathbf{p}_j^\top \mathbf{q}_j C \mathbf{b}_a^\top)} \quad (35)$$

$$= \sum_{j=1}^{k-1} \mathbf{p}_j^\top \mathbf{q}_j + \mathbf{p}_k^\top \mathbf{q}_k \quad (36)$$

where the third line holds from the inductive hypothesis Eq. 27 and the last line holds from the definition of $\mathbf{p}_k, \mathbf{q}_k \in \mathbb{R}^{1 \times K}$. \square

G.4 PROOF OF THEOREM 2

Theorem 2. Consider a nonsymmetric low-rank DPP $\mathbf{L} = \mathbf{V} \mathbf{V}^\top + \mathbf{B} \mathbf{C} \mathbf{B}^\top$, where \mathbf{V}, \mathbf{B} are of rank K , and $\mathbf{C} \in \mathbb{R}^{K \times K}$. Given a cardinality budget k , let σ_{\min} and σ_{\max} denote the smallest and largest singular values of \mathbf{L}_Y for all $Y \subseteq \llbracket M \rrbracket$ and $|Y| \leq 2k$. Assume that $\sigma_{\min} > 1$. Then,

$$\log \det(\mathbf{L}_{Y^G}) \geq \frac{4(1 - e^{-1/4})}{2(\log \sigma_{\max} / \log \sigma_{\min}) - 1} \log \det(\mathbf{L}_{Y^*}) \quad (10)$$

where Y^G is the output of Algorithm 1 and Y^* is the optimal solution of MAP inference in Eq. 4.

Proof. The proof of Theorem 2 relies on an approximation guarantee of nonsubmodular greedy maximization (Bian et al., 2017, Theorem 1). We introduce their result in below.

Theorem 3 ((Bian et al., 2017, Theorem 1)). *Consider a set function f defined on all subsets of $\{1, \dots, M\} = \llbracket M \rrbracket$ is monotone nondecreasing and nonnegative, i.e., $0 \leq f(Y) \leq f(X)$ for $\forall Y \subseteq X \subseteq \llbracket M \rrbracket$. Given a cardinality budget $k \geq 1$, let Y^* be the optimal solution of $\max_{|Y|=k} f(Y)$ and $Y^0 = \emptyset$, $Y^t := \{a_1, \dots, a_t\}$, $t = 1, \dots, k$ be the successive chosen by the greedy algorithm with budget k . Denote γ be the largest scalar such that*

$$\sum_{i \in X \setminus Y^t} (f(Y^t \cup \{i\}) - f(Y^t)) \geq \gamma(f(X \cup Y^t) - f(Y^t)), \quad (37)$$

for $\forall X \subseteq \llbracket M \rrbracket$, $|X| = k$ and $t = 0, \dots, k-1$, and α be the smallest scalar such that

$$f(Y^{t-1} \cup \{i\} \cup X) - f(Y^{t-1} \cup X) \geq (1 - \alpha)(f(Y^{t-1} \cup \{i\}) - f(Y^{t-1})). \quad (38)$$

for $\forall X \subseteq \llbracket M \rrbracket$, $|X| = k$ and $i \in Y^{k-1} \setminus X$. Then, it holds that

$$f(Y^k) \geq \frac{1}{\alpha} (1 - e^{-\alpha\gamma}) f(Y^*). \quad (39)$$

In order to apply this result for MAP inference of NDPPs, the objective should be monotone nondecreasing and nonnegative. We first show that $\sigma_{\min} > 1$ is a sufficient condition for both monotonicity and nonnegativity.

Lemma 3. *Given a P_0 matrix $L \in \mathbb{R}^{M \times M}$ and the budget $k \geq 0$, a set function $f(Y) = \log \det(L_Y)$ defined on $Y \subseteq \llbracket M \rrbracket$ is monotone nondecreasing and nonnegative for $|Y| \leq k$ when $\sigma_{\min} > 1$.*

The proof of Lemma 3 is provided in Appendix G.6. Next, we aim to find proper bounds on α and γ . To resolve this, we provide the following upper and lower bounds of the marginal gain for $f(Y) = \log \det(L_Y)$.

Lemma 4. *Let $f(Y) = \log \det(L_Y)$ and assume that $\sigma_{\min} > 1$. Then, for $Y \subseteq \llbracket M \rrbracket$, $|Y| < 2k$ and $i \notin Y$, it holds that*

$$f(Y \cup \{i\}) - f(Y) \geq \log \sigma_{\min}, \quad (40)$$

$$f(Y \cup \{i\}) - f(Y) \leq 2 \log \sigma_{\max} - \log \sigma_{\min} \quad (41)$$

where σ_{\min} and σ_{\max} are the smallest and largest singular values of L_Y for all $Y \subseteq \llbracket M \rrbracket$, $|Y| \leq 2k$.

The proof of Lemma 4 is provided in Appendix G.7. To bound γ , we consider $X \subseteq \llbracket M \rrbracket$, $|X| = k$ and denote $X \setminus Y^t = \{x_1, \dots, x_r\} \neq \emptyset$. Then,

$$\sum_{i \in X \setminus Y^t} (f(Y^t \cup \{i\}) - f(Y^t)) = \sum_{j=1}^r f(Y^t \cup \{x_j\}) - f(Y^t) \geq r \log \sigma_{\min} \quad (42)$$

where the last inequality comes from Eq. 40. Similarly, we get

$$f(X \cup Y^t) - f(Y^t) = \sum_{j=1}^r f(\{x_1, \dots, x_j\} \cup Y^t) - f(\{x_1, \dots, x_{j-1}\} \cup Y^t) \quad (43)$$

$$\leq r(2 \log \sigma_{\max} - \log \sigma_{\min}) \quad (44)$$

where the last inequality comes from Eq. 41. Combining Eq. 42 to Eq. 44, we obtain that

$$\frac{\sum_{i \in X \setminus Y^t} f(Y^t \cup \{i\}) - f(Y^t)}{f(X \cup Y^t) - f(Y^t)} \geq \frac{\log \sigma_{\min}}{2 \log \sigma_{\max} - \log \sigma_{\min}} \quad (45)$$

which allows us to choose $\gamma = \left(2 \frac{\log \sigma_{\max}}{\log \sigma_{\min}} - 1\right)^{-1}$.

To bound α , we similarly use Lemma 4 to obtain

$$\frac{f(X \cup Y^{t-1} \cup \{i\}) - f(X \cup Y^{t-1})}{f(Y^{t-1} \cup \{i\}) - f(Y^{t-1})} \geq \frac{\log \sigma_{\min}}{2 \log \sigma_{\max} - \log \sigma_{\min}} \quad (46)$$

and we can choose $\alpha = 1 - \frac{\log \sigma_{\min}}{2 \log \sigma_{\max} - \log \sigma_{\min}} = \frac{2(\log \sigma_{\max} - \log \sigma_{\min})}{2 \log \sigma_{\max} - \log \sigma_{\min}}$.

Now let $\kappa = \frac{\log \sigma_{\max}}{\log \sigma_{\min}}$. Then $\gamma = \frac{1}{2\kappa-1}$ and $\alpha = \frac{2(\kappa-1)}{2\kappa-1}$. Putting γ and α into Eq. 39, we have

$$\frac{1}{\alpha}(1 - e^{-\alpha\gamma}) \geq \frac{2\kappa-1}{2(\kappa-1)} \left(1 - e^{-\frac{2(\kappa-1)}{(2\kappa-1)^2}}\right) \quad (47)$$

$$\geq \frac{2\kappa-1}{2(\kappa-1)} 4 \exp(-1/4) \frac{2(\kappa-1)}{(2\kappa-1)^2} \quad (48)$$

$$= \frac{4 \exp(-1/4)}{2\kappa-1} \quad (49)$$

where the second inequality holds from the fact that $\max_{\kappa \geq 1} \frac{2(\kappa-1)}{(2\kappa-1)^2} = \frac{1}{4}$ and $1 - e^{-x} \geq 4 \exp(-1/4)x$ for $x \in [0, 1/4]$. \square

G.5 PROOF OF COROLLARY 1

Corollary 1. Consider a nonsymmetric low-rank DPP $\mathbf{L} = \mathbf{V}\mathbf{V}^\top + \mathbf{B}\mathbf{C}\mathbf{B}^\top$, where \mathbf{V}, \mathbf{B} are of rank K , and $\mathbf{C} \in \mathbb{R}^{K \times K}$. Given a cardinality budget k , let σ_{\min} and σ_{\max} denote the smallest and largest singular values of \mathbf{L}_Y for all $Y \subseteq [M]$ and $|Y| \leq 2k$. Let $\kappa := \sigma_{\max}/\sigma_{\min}$. Then,

$$\log \det(\mathbf{L}_{Y^G}) \geq \frac{4(1 - e^{-1/4})}{2 \log \kappa + 1} \log \det(\mathbf{L}_{Y^*}) - \left(1 - \frac{4(1 - e^{-1/4})}{2 \log \kappa + 1}\right) k (1 - \log \sigma_{\min}) \quad (12)$$

where Y^G is the output of Algorithm 1 and Y^* is the optimal solution of MAP inference in Eq. 4.

Proof. Now consider $\mathbf{L}' = (\frac{e}{\sigma_{\min}})\mathbf{L}$ where e is the exponential constant. Then, $\sigma'_{\min} = \sigma_{\min}(\frac{e}{\sigma_{\min}}) = e$ and $\sigma'_{\max} = \sigma_{\max}(\frac{e}{\sigma_{\min}})$. Using the fact that $\log \det(\mathbf{L}'_Y) = \log \det(\mathbf{L}_Y) - |Y| \log \sigma_{\min}$, we obtain the result. \square

G.6 PROOF OF LEMMA 3

Before stating the proof, we introduce interlacing properties of singular values.

Theorem 4 (Interlacing Inequality for Singular Values, (Thompson, 1972, Theorem 1)). Consider a real matrix $\mathbf{A} \in \mathbb{R}^{M \times N}$ with singular values $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_{\min(M,N)}$ and its supmatrix $\mathbf{B} \in \mathbb{R}^{P \times Q}$ with singular values $\beta_1 \geq \beta_2 \geq \dots \geq \beta_{\min(P,Q)}$. Then, the singular values have the following interlacing properties:

$$\sigma_i \geq \beta_i, \quad \text{for } i = 1, \dots, \min(P, Q), \quad (50)$$

$$\beta_i \geq \sigma_{i+M-P+N-Q}, \quad \text{for } i = 1, \dots, \min(P+Q-M, P+Q-N). \quad (51)$$

Note that when $M = N$ and $P = Q = N - 1$, it holds that $\beta_i \geq \sigma_{i+2}$ for $i = 1, \dots, N - 2$.

We are now ready to prove Lemma 3.

Lemma 3. Given a P_0 matrix $\mathbf{L} \in \mathbb{R}^{M \times M}$ and the budget $k \geq 0$, a set function $f(Y) = \log \det(\mathbf{L}_Y)$ defined on $Y \subseteq [M]$ is monotone nondecreasing and nonnegative for $|Y| \leq k$ when $\sigma_{\min} > 1$.

Proof. Since $\mathbf{L} \in P_0$, its all principal submatrices are also in P_0 . By the definition of a P_0 matrix, it holds that

$$|\det(\mathbf{L}_Y)| = \det(\mathbf{L}_Y) = \prod_i \sigma_i(\mathbf{L}_Y) \quad (52)$$

where $\sigma_i(\mathbf{L}_Y)$ for $i \in [Y]$ are singular values of \mathbf{L}_Y . Then, $F(Y) = \sum_i \log(\sigma_i(\mathbf{L}_Y))$ is nonnegative for all Y such that $|Y| \leq K$ due to $\sigma_i(\mathbf{L}_Y) \geq \sigma_{\min} > 1$. Similarly, we have $F(Y \cup \{a\}) - F(Y) = \sum_{i=1}^{|Y|+1} \log \sigma_i(\mathbf{L}_{Y \cup \{a\}}) - \sum_{i=1}^{|Y|} \log \sigma_i(\mathbf{L}_Y) \geq \log \sigma_{\min} > 0$ from Eq. 50. \square

G.7 PROOF OF LEMMA 4

Lemma 4. *Let $f(Y) = \log \det(\mathbf{L}_Y)$ and assume that $\sigma_{\min} > 1$. Then, for $Y \subseteq \llbracket M \rrbracket$, $|Y| < 2k$ and $i \notin Y$, it holds that*

$$f(Y \cup \{i\}) - f(Y) \geq \log \sigma_{\min}, \quad (40)$$

$$f(Y \cup \{i\}) - f(Y) \leq 2 \log \sigma_{\max} - \log \sigma_{\min} \quad (41)$$

where σ_{\min} and σ_{\max} are the smallest and largest singular values of \mathbf{L}_Y for all $Y \subseteq \llbracket M \rrbracket$, $|Y| \leq 2k$.

Proof. For a P_0 matrix, we remark that its determinant is equivalent to the product of all singular values. For $Y \subseteq \llbracket M \rrbracket$ and $i \notin Y$, from the interlacing inequality of Eq. 50 we have that

$$F(Y \cup \{i\}) - F(Y) = \sum_{j=1}^{|Y|+1} \log \sigma'_j - \sum_{j=1}^{|Y|} \log \sigma_j \geq \log \sigma'_{|Y|+1} \geq \log \sigma_{\min} \quad (53)$$

where σ'_j and σ_j are the j -th largest singular value of $\mathbf{L}_{Y \cup \{i\}}$ and \mathbf{L}_Y , respectively. Similarly, using Eq. 51, we get

$$F(Y \cup \{i\}) - F(Y) \leq \log(\sigma'_1 \sigma'_2) - \log \sigma_{|Y|} \leq 2 \log \sigma_{\max} - \log \sigma_{\min}. \quad (54)$$

□