

MM アルゴリズムによる行列式点過程の学習

川島 貴大¹, 日野英逸^{2,3}

IBISML@OIST June 29, 2023

¹ 総合研究大学院大学 統計科学専攻

² 統計数理研究所

³ 理研 AIP



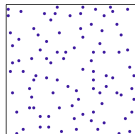


1. 背景
2. 提案手法
3. 既存法との関連
4. 実験
5. むすび

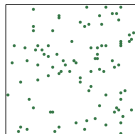
背景



全体集合からの多様なアイテムの生起を確率的にモデル化したい



DPP



Independent

格子点上の超一様サンプル

“porsche”

k=2



k=4



“philadelphia”

k=2



k=4



“cocker spaniel”

k=2



k=4



検索クエリからの多様な画像サジェスト

共起しやすい／しにくいアイテムの組合せを考慮できればよい



行列式点過程 (DPP; Determinantal Point Processes)

有限集合 $\mathcal{Y} = \{1, 2, \dots, N\}$ とその部分集合 $\mathcal{A} \subseteq \mathcal{Y}$ に対し、
行列式点過程が定める尤度関数は

$$P(\mathcal{A}|\mathbf{L}) = \frac{\det([\mathbf{L}]_{\mathcal{A}})}{\det(\mathbf{L} + \mathbf{I})}$$

と表される.

ここで $\mathbf{L} \in \mathbb{S}_+^N$ について, $[\mathbf{L}]_{\mathcal{A}} = (L_{ij})_{i,j \in \mathcal{A}} \in \mathbb{S}_+^{|\mathcal{A}|}$ は \mathcal{A} の要素によって定められる \mathbf{L} の主部分行列.

› 行列式によってアイテム間の共起・斥力が表現される



DPP の学習法はパラメータ L の構造により 3 種類に大別

- ・フルランク

- › EM [2], 不動点アルゴリズム [3]

- ・低ランク

- › 勾配ベース最適化 [4, 5]

- ・その他の構造

- › クロネッカー積に分解 [6], 対角 + 特殊な低ランク構造 [7]

計算量削減のため低ランク性や特殊な構造が L に仮定される



小中規模の問題はできれば L に余計な制約を与えず済ませたい

例：DPP の適用可能性を調べる予備実験

‣ フルランク DPP の学習

- ・ EM [Gillenwater et al., 2014]

- ・ Stiefel 多様体上での最適化が必要。複雑で不安定

- ・ 不動点アルゴリズム [Mariet and Sra, 2015]

- ・ シンプルだがハイパラを動かし加速させると収束が無保証

‣ 安定・高速・シンプルな DPP の学習則を考える

提案手法



平均対数尤度

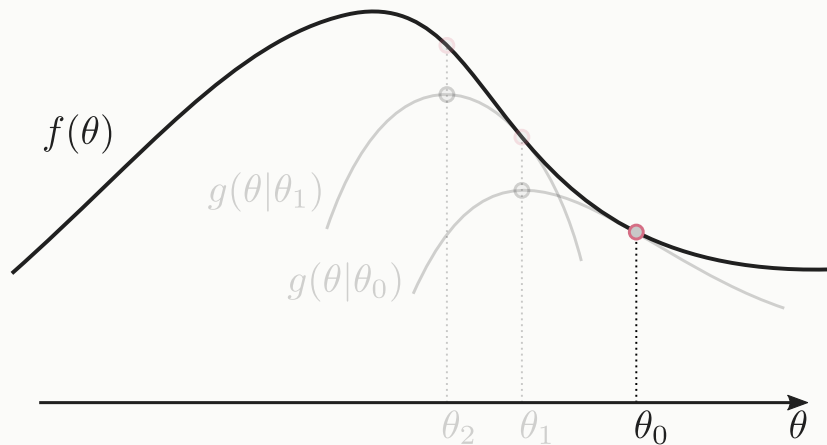
$$f(\mathbf{L}) = \frac{1}{M} \sum_{m=1}^M \log \det([\mathbf{L}]_{\mathcal{A}_m}) - \log \det(\mathbf{L} + \mathbf{I})$$

を MM アルゴリズムで最大化。つまり任意の $\mathbf{L}, \mathbf{L}^{(t)} \succeq 0$ に対し

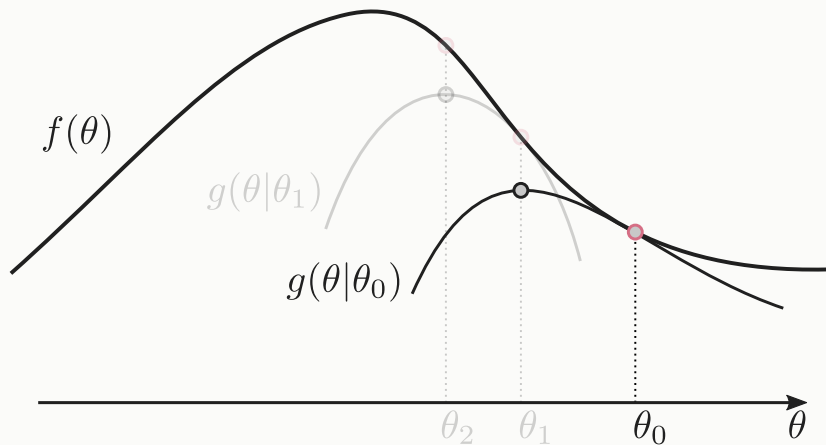
- $f(\mathbf{L}) \geq g(\mathbf{L}|\mathbf{L}^{(t)})$
- $f(\mathbf{L}^{(t)}) = g(\mathbf{L}^{(t)}|\mathbf{L}^{(t)})$

が成り立つ minorizer $g(\cdot|\mathbf{L}^{(t)})$ を設計して代理関数とし、繰り返し最大化を行う：

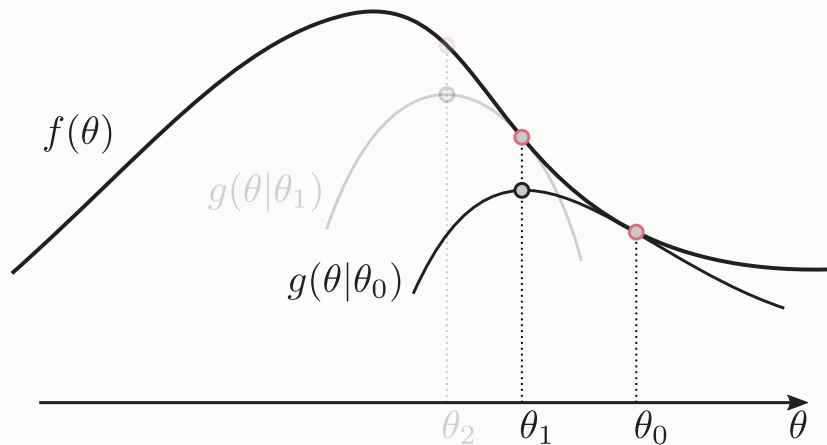
$$\mathbf{L}^{(t+1)} = \arg \max_{\mathbf{L} \succeq 0} g(\mathbf{L}|\mathbf{L}^{(t)})$$



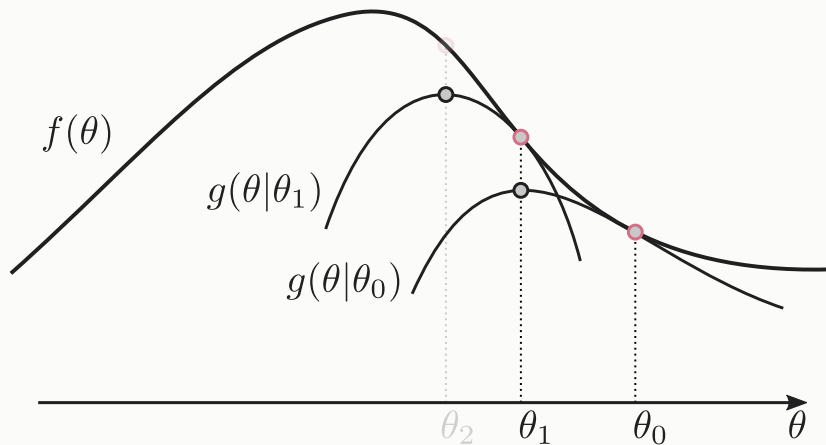
MM アルゴリズムの概要図



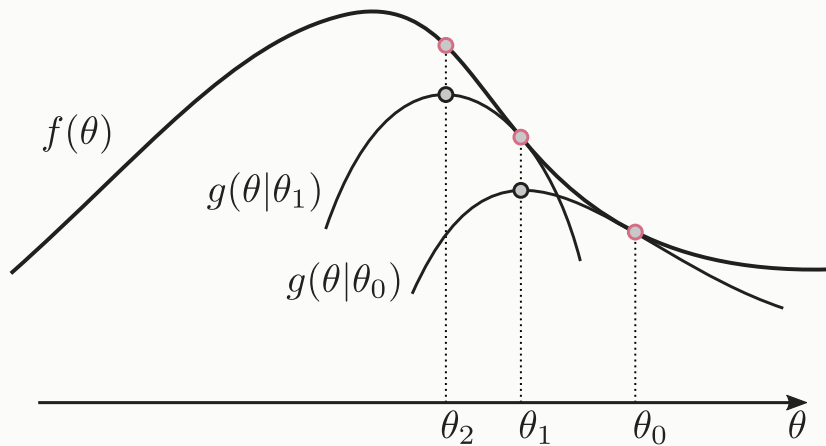
MM アルゴリズムの概要図



MM アルゴリズムの概要図



MM アルゴリズムの概要図



MM アルゴリズムの概要図



結果として、平均対数尤度 $f(\mathbf{L})$ の minorizer として次を得る：

行列式点過程のための minorizer

$$g(\mathbf{L}|\mathbf{L}^{(t)}) = -\frac{1}{M} \sum_{m=1}^M \text{tr}\{\mathbf{L}^{(t)} \mathbf{U}_{\mathcal{A}_m}^\top [\mathbf{L}^{(t)}]_{\mathcal{A}_m}^{-1} \mathbf{U}_{\mathcal{A}_m} \mathbf{L}^{(t)} \mathbf{L}^{-1}\} \\ - \text{tr}\{(\mathbf{L}^{(t)} + \mathbf{I})^{-1} \mathbf{L}\} + \text{const.}$$

ここで $\mathbf{U}_{\mathcal{A}_m}$ は N 次単位行列から \mathcal{A}_m の要素に対応する行のみを残した $|\mathcal{A}_m| \times N$ バイナリ行列.

> $[\mathbf{L}]_{\mathcal{A}} = \mathbf{U}_{\mathcal{A}} \mathbf{L} \mathbf{U}_{\mathcal{A}}^\top$ と書ける

👍 $g(\mathbf{L}|\mathbf{L}^{(t)})$ は凹関数



minorizer $g(\mathbf{L}|\mathbf{L}^{(t)})$ の最大化は、一次の最適性条件から

$$\underline{-\mathbf{L}(\mathbf{L}^{(t)} + \mathbf{I})^{-1}\mathbf{L} + \mathbf{Q}_M^{(t)} = \mathbf{O}}$$

$$\mathbf{Q}_M^{(t)} = \mathbf{L}^{(t)} \left(\frac{1}{M} \sum_{m=1}^M \mathbf{U}_{\mathcal{A}_m}^\top [\mathbf{L}^{(t)}]_{\mathcal{A}_m}^{-1} \mathbf{U}_{\mathcal{A}_m} \right) \mathbf{L}^{(t)}$$

を満たす \mathbf{L} を求めるという問題に帰着

これは CARE¹ という二次の行列方程式の特殊形

‣ 🍀 Schur method などにより $\mathcal{O}(N^3)$ で 数値的に解ける

提案法ではこの CARE を繰り返し解いて \mathbf{L} の最尤推定値を得る

¹continuous algebraic Riccati equation.

既存法との関連



Mariet and Sra (2015) は不動点アルゴリズムとして更新則

$$\mathbf{L}^{(t+1)} = \mathbf{L}^{(t)} + \mathbf{L}^{(t)} \left(\frac{1}{M} \sum_m U_{\mathcal{A}_m}^\top [\mathbf{L}^{(t)}]_{\mathcal{A}_m}^{-1} U_{\mathcal{A}_m} - (\mathbf{L}^{(t)} + \mathbf{I})^{-1} \right) \mathbf{L}^{(t)}$$

を導いているが、実はこれは minorizer を

$$\begin{aligned} h(\mathbf{L}|\mathbf{L}^{(t)}) &= -\frac{1}{M} \sum_{m=1}^M \text{tr}\{\mathbf{L}^{(t)} U_{\mathcal{A}_m}^\top [\mathbf{L}^{(t)}]_{\mathcal{A}_m}^{-1} U_{\mathcal{A}_m} \mathbf{L}^{(t)} \mathbf{L}^{-1}\} \\ &\quad - \log \det(\mathbf{L}) - \text{tr}\{(\mathbf{L}^{(t)} + \mathbf{I})^{-1} \mathbf{L}^{-1} \mathbf{L}^{(t)}\} + \text{const.} \end{aligned}$$

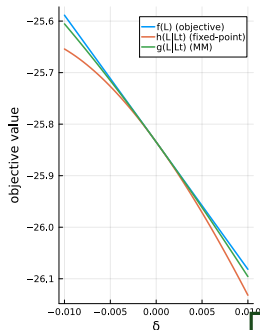
なる 非凹関数 とおいた MM アルゴリズムとみなせる。

➤ minorizer の最適性は必ずしも保証されない

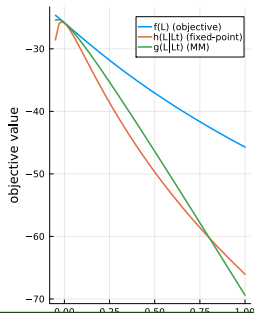


提案法と既存法の関連

$L^{(t)}$ 近傍の L に対して, $g(L|L^{(t)}) \geq h(L|L^{(t)})$.



$L^{(t)}$ 近傍



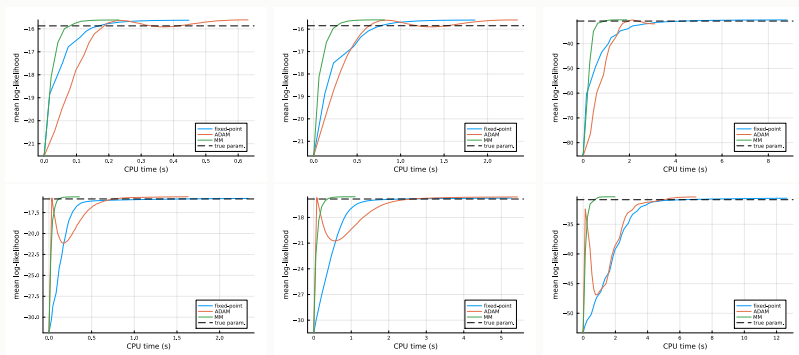
$L^{(t)}$ 非近傍

👍 提案 minorizer は局所的にタイト

実験



$L^* = V^*V^{*\top}$, $v_{ij}^* \sim \mathcal{U}(0, 10/N)$ から M 個のサンプルを生成



$N = 32, M = 2,500$

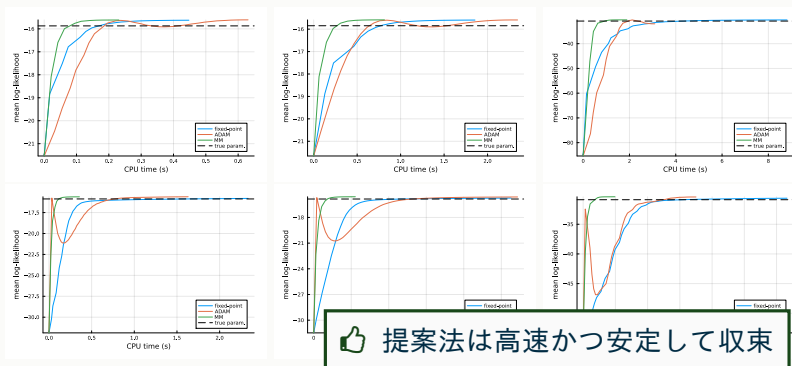
$N = 32, M = 10,000$

$N = 128, M = 2,500$

(上段) Wishart 分布から初期化, (下段) 一様分布から初期化.
 青線 : 不動点アルゴリズム, 橙線 : Adam, 緑線 : 提案法.



$L^* = V^*V^{*\top}$, $v_{ij}^* \sim \mathcal{U}(0, 10/N)$ から M 個のサンプルを生成



提案法は高速かつ安定して収束

$N = 32, M = 2,500$

$N = 32, M = 10,000$

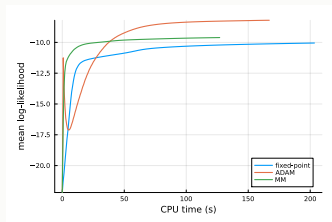
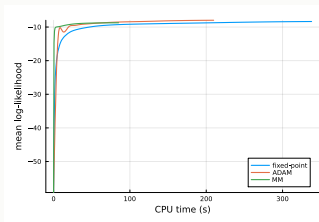
$N = 128, M = 2,500$

(上段) Wishart 分布から初期化, (下段) 一様分布から初期化.
 青線 : 不動点アルゴリズム, 橙線 : Adam, 緑線 : 提案法.



フォーク音楽における各鍵盤の打鍵データ ($N = 88$, $\bar{M} = 6,182$)

Method	WISHART		BASIC	
	Log-likelihood	Runtime (s)	Log-likelihood	Runtime (s)
FP	-8.30 ± 0.22	46.58 ± 3.66	-10.14 ± 0.28	37.26 ± 7.47
Adam	-7.98 ± 0.76	29.36 ± 8.28	-8.89 ± 2.96	29.00 ± 11.91
MM	-9.51 ± 0.25	24.27 ± 6.66	-9.59 ± 0.22	21.56 ± 4.24

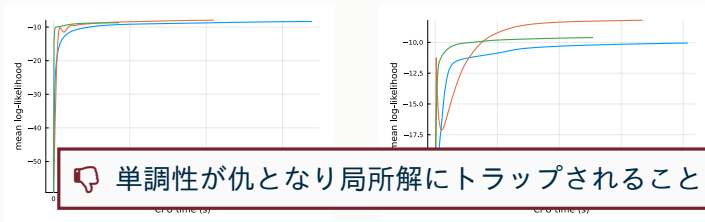


(左) Wishart 分布で初期化, (右) 一様分布で初期化. 緑線が提案法.



フォーク音楽における各鍵盤の打鍵データ ($N = 88$, $\bar{M} = 6,182$)

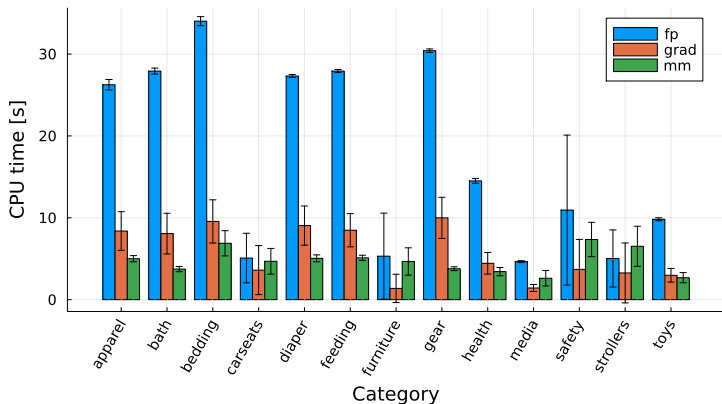
Method	WISHART		BASIC	
	Log-likelihood	Runtime (s)	Log-likelihood	Runtime (s)
FP	-8.30 ± 0.22	46.58 ± 3.66	-10.14 ± 0.28	37.26 ± 7.47
Adam	-7.98 ± 0.76	29.36 ± 8.28	-8.89 ± 2.96	29.00 ± 11.91
MM	-9.51 ± 0.25	24.27 ± 6.66	-9.59 ± 0.22	21.56 ± 4.24



(左) Wishart 分布で初期化, (右) 一様分布で初期化. 緑線が提案法.



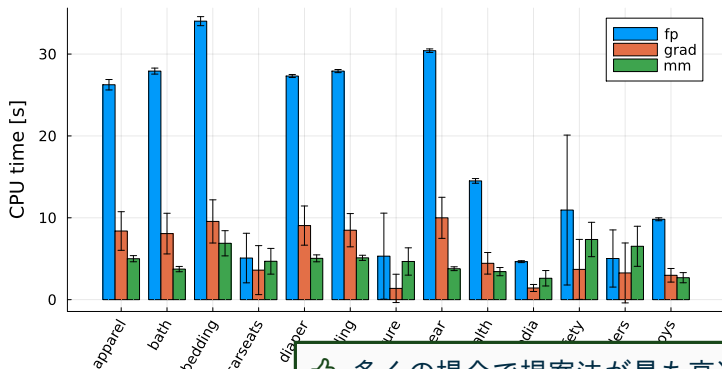
育児関連商品についての欲しいものリスト内の商品を 13 のカテゴリに細分化 ($\bar{N} = 71$, $\bar{M} = 8,585$)



各カテゴリに対する計算時間. 緑が提案法.



育児関連商品についての欲しいものリスト内の商品を 13 のカテゴリに細分化 ($\bar{N} = 71$, $\bar{M} = 8,585$)



👍 多くの場合で提案法が最も高速

各カテゴリに対する計算時間. 緑が提案法.

むすび



- [1] Alex Kulesza and Ben Taskar. Determinantal point processes for machine learning. *Foundations and Trends® in Machine Learning*, 5(2-3):123–286, 2012.
- [2] Jennifer A Gillenwater, Alex Kulesza, Emily Fox, and Ben Taskar. Expectation-Maximization for Learning Determinantal Point Processes. In *Advances in Neural Information Processing Systems*, volume 27, 2014.
- [3] Zelda Mariet and Suvrit Sra. Fixed-point algorithms for learning determinantal point processes. In *Proceedings of the 32nd International Conference on Machine Learning*, pages 2389–2397, 2015.
- [4] Mike Gartrell, Ulrich Paquet, and Noam Koenigstein. Low-Rank Factorization of Determinantal Point Processes. *Proceedings of the AAAI Conference on Artificial Intelligence*, 31(1), 2017.



- [5] Takayuki Osogami, Rudy Raymond, Akshay Goel, Tomoyuki Shirai, and Takanori Maehara. Dynamic Determinantal Point Processes. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1), 2018.
- [6] Zelda E. Mariet and Suvrit Sra. Kronecker Determinantal Point Processes. In *Advances in Neural Information Processing Systems*, volume 29, 2016.
- [7] Christophe Dupuy and Francis Bach. Learning Determinantal Point Processes in Sublinear Time. In *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*, pages 244–257, 2018.