

相互カーネル行列補完によるタンパク質機能予測 Protein Function Prediction by Mutual Kernel Matrix Completion

リベロ レイシェル[†]
Rachelle Rivero

下山 愛祐美[†]
Ayumi Shimoyama

加藤 毅^{† § ¶}
Tsuyoshi Kato

1. まえがき

2000 年以降測定技術の急速な発達により、複数の種類のゲノムワイドのデータが得られるようになった。個々のデータタイプは細胞内の機構を一つの側面から見た観測値となっている。これらを統合して、細胞内でタンパク質の相互作用や細胞内の生物学的な機構の解明に用いられるようになった [1]。

カーネル行列を介した異種データの統合がしばしば用いられている [1]。データをカーネル行列で表現しておくことで、カーネル法と呼ばれる一連の方法論が利用可能になる。データの統合は単純で、各々のデータタイプから得られたカーネル行列を単純に平均をとるだけで統合することができる。この統合法の正当性は、カーネルの半正定値性を仮定してカーネル法が構成されており、半正定値カーネルは凸錐上にあることからいえる。

分子生物学において複数のデータを統合するとき一部のタンパク質に対して、一部のデータタイプが欠けているような状況がある。このような状況においてカーネルの線形結合によるデータ統合を行うには、2つのアプローチがある。その一つは、そのタンパク質を解析対象から除外するか、そのデータタイプを除外する方法である。すると、解析対象が少なくなる、もしくは、有用なデータタイプを一部のデータの欠落だけのために放棄することになってしまう。もう一つのアプローチは、欠落しているカーネルの値を推定する方法である [2, 3]。

従来のカーネル行列補完の方法は、1個のデータタイプだけ不完全で、そのほかのデータタイプはすべてのタンパク質に対して数値が観測できていることを仮定していた [2, 3] (図 1(a),(b) 参照)。しかし、現実的には、複数のデータタイプが不完全である場合 (図 1(c) 参照) がよくあり、従来法の想定とはギャップがある。本研究では、図 1(c) のように所与のカーネル行列各々が不完全でも相互に補完できる新しい方法論を開発した。

2. 問題設定

本研究で議論する相互カーネル行列補完 (MKMC) というタスクを述べる。K 個の不完全カーネル行列 $Q^{(k)}$ ($k = 1, \dots, K$) が図 1(c) のように与えられていたとする。 $Q^{(k)}$ の行列のサイズはいずれも $\ell \times \ell$ である。K 個の情報源のいずれかの対象に関する情報が欠損しているとする。すると、それぞれの情報源のカーネル行列 $Q^{(k)}$ における対応する行と列は欠損していることになる。 $Q_{vh,vh}^{(k)}$ は、 $Q^{(k)}$ の行と列の順番を次のように並び替えたものとする。最初の $n_k < \ell$ 個の対象に関

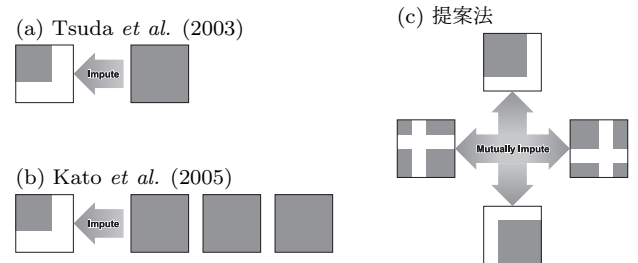


図 1: 問題設定.

する情報が得られているとし、残りの $(\ell - n_k)$ 個の対象に関する情報が欠損しているとする。このようにして $Q^{(k)}$ の行と列の順番を並び替えて得られる $Q_{vh,vh}^{(k)}$ を次のように分割する：

$$Q_{vh,vh}^{(k)} = \begin{bmatrix} Q_{v,v}^{(k)} & Q_{v,h}^{(k)} \\ Q_{h,v}^{(k)} & Q_{h,h}^{(k)} \end{bmatrix} \quad (1)$$

ただし、 $Q_{v,v}^{(k)} \in \mathbb{S}_{++}^{n_k}$ 、 $Q_{v,h}^{(k)} = (Q_{h,v}^{(k)})^\top \in \mathbb{R}^{n_k \times m_k}$ 、および $Q_{h,h}^{(k)} \in \mathbb{S}_{++}^{m_k}$ である。開発した算法では欠損した部分行列 $Q_{v,h}^{(k)}$ 、および $Q_{h,h}^{(k)}$ の値を推定する。 $\mathcal{H} := \{Q_{v,h}^{(k)}, Q_{h,h}^{(k)}\}_{k=1}^K$ とおく。

3. 提案法：相互カーネル行列補完法

提案法では、モデル行列 $M \in \mathbb{S}_{++}^\ell$ を導入し、各 $Q^{(k)}$ が M との距離が最小になる \mathcal{H} を求めるものである。すなわち、目的関数を各 $Q^{(k)}$ と M とのカルバックライブラ (KL) 距離の和

$$J(\mathcal{H}, M) := \sum_{k=1}^K \text{KL}(Q^{(k)}, M) \quad (2)$$

で定義し、この目的関数を最小にする \mathcal{H} と M の組み合わせを見つける。KL 距離は、本来、確率分布間の擬距離であるが、本研究では、文献 [2, 3] に倣って、カーネル行列と正規分布とを関連付け、正規分布間の KL 距離

$$\text{KL}(Q^{(k)}, M) := \text{KL}(\mathcal{N}(\mathbf{0}, Q^{(k)}) || \mathcal{N}(\mathbf{0}, M)) \quad (3)$$

をカーネル行列間の KL 距離とした。この定義を用いると、最小解が正定値になることが保証されるので、カーネル行列たるために必要な正定値性に関する追加的制約が不要になる。しかし、 (\mathcal{H}, M) を同時最小化する解は閉形式では与えられないので、適当な正定値行列

[†]群馬大学大学院理工学府

[§]群馬大学次世代モビリティ社会実装センター (CRANTS)

[¶]早稲田大学規範科学総合研究所 (IIRS)

$M \in \mathbb{S}_{++}^{\ell}$ を初期値とし, E-step と M-step を繰り返す次の算法を用いることにする.

- 1: **repeat**
- 2: {E-step} $\mathcal{H} := \operatorname{argmin}_{\mathcal{H}} J(\mathcal{H}, M)$;
- 3: {M-step} $M := \operatorname{argmin}_{M} J(\mathcal{H}, M)$;
- 4: **until** convergence.

この算法の性質から, 目的関数は単調非増加することは明らかである. 加えて, 次のような望ましい性質を持つ:

Theorem 1. 上述の算法における E-step および M-step は閉形式で表される.

各ステップを, それぞれ, E-step および M-step と表記した所以は, それぞれ統計学でいう期待値最大化 (EM) 法の E-step および M-step で説明できるからである. その理由を次節で述べる. また, 紙面の制約から具体的な更新式は次節にのみ記述した.

4. EM 法との関係

EM 法再訪. EM 法は, 観測値 \mathbf{V} および潜在変数 \mathbf{H} の同時確率のモデル $p(\mathbf{V}, \mathbf{H} | \Theta)$ のパラメータ Θ を最尤推定するための枠組みである. すなわち, 経験分布 $q(\mathbf{V})$ に対して, EM 法は周辺分布 $p(\mathbf{V} | \Theta)$ の対数尤度関数

$$L(\Theta) := \mathbb{E}_{q(\mathbf{V})} [\log p(\mathbf{V} | \Theta)] \quad (4)$$

を最大化する Θ を見つけようとする. いま, 同時確率モデルの対数密度が

$$\log p(\mathbf{V}, \mathbf{H} | \Theta) = \langle \mathbf{S}(\mathbf{V}, \mathbf{H}), \mathbf{G}(\Theta) \rangle - A(\Theta) \quad (5)$$

のように指数分布族で表されているとする. ただし, $\mathbf{S}(\mathbf{V}, \mathbf{H}), \mathbf{G}(\Theta)$, および $A(\Theta)$ は十分統計量, 自然パラメータ, キュムラント関数である. すると, EM 法における Q 関数は

$$Q(\Theta; \Theta^{\text{old}}) := \langle \mathbb{E}_{\Theta^{\text{old}}} [\mathbf{S}(\mathbf{V}, \mathbf{H})], \mathbf{G}(\Theta) \rangle - A(\Theta) \quad (6)$$

と表される. E-step で期待値 $\mathbb{E}_{\Theta^{\text{old}}} [\mathbf{S}(\mathbf{V}, \mathbf{H})]$ を計算し, M-step では期待値を固定して, Q 関数を最大化するように Θ を更新する.

提案法は EM 法. 次のようにモデル分布と経験分布を定めると, 提案法は EM 法であることがわかる. 今, 完全データを $\mathbf{x}_1, \dots, \mathbf{x}_K \in \mathbb{R}^{\ell}$ とする. \mathbf{x}_k の各要素は $\mathbf{Q}^{(k)}$ のある行と列に関連付ける. $\mathbf{v}_k \in \mathbb{R}^{n_k}$ および $\mathbf{h}_k \in \mathbb{R}^{\ell - n_k}$ を \mathbf{x}_k の部分ベクトルとし, \mathbf{v}_k は $\mathbf{Q}^{(k)}$ の可視データに \mathbf{h}_k は $\mathbf{Q}^{(k)}$ の欠損データに対応する. EM 法における観測データを $\mathbf{V} := (\mathbf{v}_1, \dots, \mathbf{v}_K)$ で構成し, 潜在データを $\mathbf{H} := (\mathbf{h}_1, \dots, \mathbf{h}_K)$ で構成する. 経験分布 $q(\mathbf{V})$ は $\mathbb{E}_{q(\mathbf{V})} [\mathbf{v}_k \mathbf{v}_k^{\top}] = \mathbf{Q}_{v,v}^{(k)}$ を満たすものと

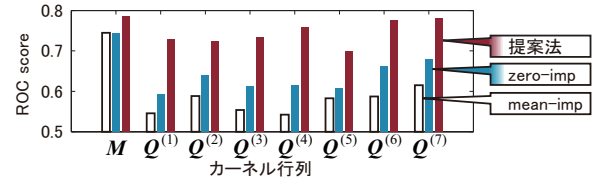


図 2: タンパク質機能予測の性能.

し, モデル分布は $\Theta = M$ とし, 完全データの密度は $p(\mathbf{V}, \mathbf{H} | \Theta) := \prod_{k=1}^K \mathcal{N}(\mathbf{x}_k | \mathbf{0}, M)$ で与えるとする. このモデルは,

$$\mathbf{S}(\mathbf{V}, \mathbf{H}) := \sum_{k=1}^K \mathbf{x}_k \mathbf{x}_k^{\top}, \text{ および } \mathbf{G}(M) = -\frac{1}{2} M^{-1}$$

とおくことにより, 指数分布族であることが示される. この設定を EM 法にあてはめる. $M_{h|v}^{(k)} := M_{h,h}^{(k)} - M_{h,v}^{(k)} (M_{v,v}^{(k)})^{-1} M_{v,h}^{(k)}$ とおくと,

$$\mathbb{E} [\mathbf{v}_k \mathbf{h}_k^{\top}] = \mathbf{Q}_{v,v}^{(k)} (M_{v,v}^{(k)})^{-1} M_{v,h}^{(k)}, \text{ および}$$

$$\mathbb{E} [\mathbf{h}_k \mathbf{h}_k^{\top}] = M_{h|v}^{(k)} + M_{h,v}^{(k)} (M_{v,v}^{(k)})^{-1} \mathbf{Q}_{v,v}^{(k)} (M_{v,v}^{(k)})^{-1} M_{v,h}^{(k)},$$

が成立し, それぞれ M を固定して \mathcal{H} を最適化したときの $\mathbf{Q}_{v,h}^{(k)}$ および $\mathbf{Q}_{h,h}^{(k)}$ に等しい. また, Q 関数を最大化する M は $\mathbf{Q}^{(k)}$ の平均で与えられるので, $\operatorname{argmin} J(\mathcal{H}, M)$ に等しくなる. 以上の導出により, 次の理論的結果を得た:

Theorem 2. 前節で述べた算法は EM 法の一つである.

5. 実験結果

文献 [1] で用いられていた 7 個のカーネル行列を使って性能の評価を行った. 0 で埋める方法 (zero-imp) および平均で埋める方法 (mean-imp) と提案法を比較した. 各データタイプで 50% を欠損させた. 無作為に選んだ 200 個を SVM の訓練に使い, 評価用タンパク質が膜タンパクか予測した. 各カーネル行列における予測性能を ROC スコアで評価したところ, 図 2 を得た. いずれのカーネル行列でも, 0 で埋めたり平均で埋めるような単純な補完法より提案法で補完したほうが顕著に高い汎化性能を得ることができた.

謝辞:本研究は JSPS 科研費 26249075, 40401236 の助成を受けたものである.

参考文献

- [1] G. R. Lanckriet et al. A statistical framework for genomic data fusion. *Bioinformatics*, 20(16):2626–35, Nov 2004.
- [2] Koji Tsuda et al. The em algorithm for kernel matrix completion with auxiliary data. *JMLR*, 4:67–81, 2003.
- [3] Tsuyoshi Kato, Koji Tsuda, and Kiyoshi Asai. Selective integration of multiple biological data for supervised network inference. *Bioinformatics*, 21:2488–2495, 2005.