

高次元統計モデリングユニット 山田 誠 High-dimensional Statistical Modeling Unit Makoto Yamada



科学的発見のための機械学習技術の研究開発

問題意識: バイオデータのような非線形データからは既存法では 科学的発見ができない場合が多く, 新規の発見が見過ごされて いる可能性がある.

目標: 非線形データから自動的に新規の科学的発見をする枠組みの構築. 科学的発見を加速させる.

成果1:100万次元を超える超高次元非線形データから,中・小規模サーバを用いて自動で重要な特徴を選択する方法を開発.

成果2: 高次元非線形データから, 個人の特徴を選択する手法(個別化特徴選択手法)の枠組みを提案.

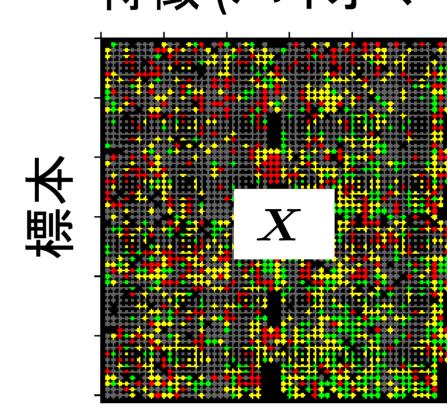
成果3: 非線形関係のあるデータから, 自動で重要なデータセットを選択し, そのp値を計算する方法を世界で初めて開発.

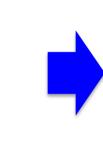
解釈性

既存 手法 研究対象

柔軟性(非線形性)

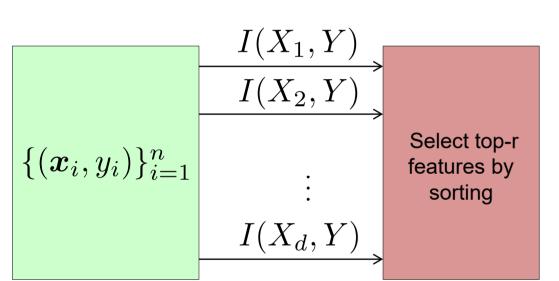
特徴 (バイオマーカー) ラベル

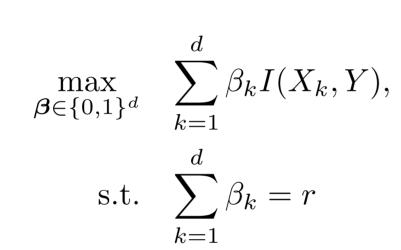




Block HSIC Lasso [biorXiv 2019]

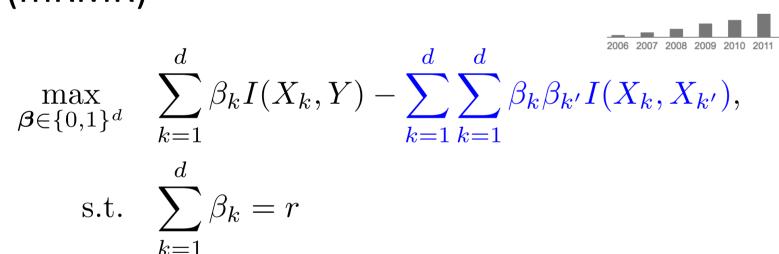
Sure Independence Screening (SIS法)





重複した特徴を選択してしまう⊗

Minimum Redundancy Maximum Relevance 法 (mRMR)



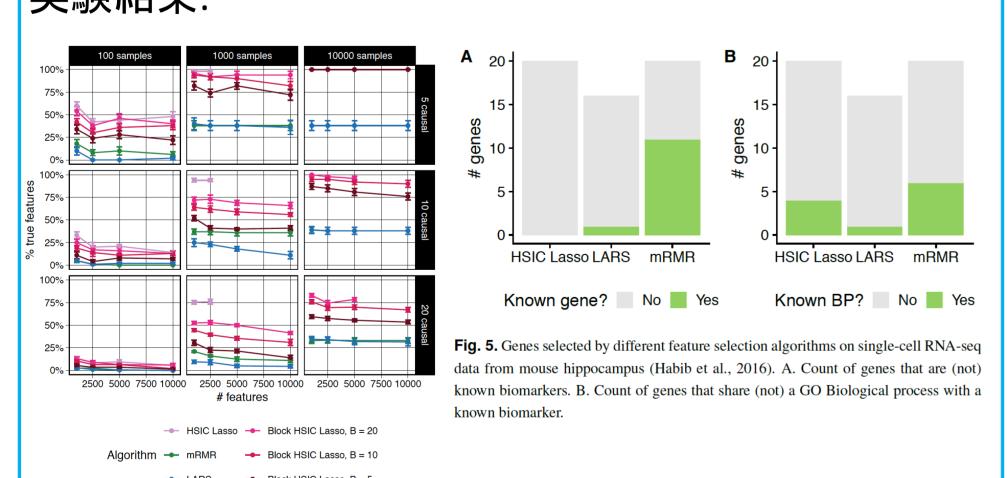
Block HSIC Lasso, mRMR法の凸最適化版

$$\max_{\boldsymbol{\alpha} \geq 0} \sum_{k=1}^{d} \alpha_k \mathrm{HSIC}_b(\boldsymbol{f}_k, \boldsymbol{y}) - \frac{1}{2} \sum_{k,k'=1}^{d} \alpha_k \alpha_{k'} \mathrm{HSIC}_b(\boldsymbol{f}_k, \boldsymbol{f}_{k'}) - \lambda \|\boldsymbol{\alpha}\|_1,$$

$$\mathrm{HSIC}_b(\boldsymbol{f}_k, \boldsymbol{y}) = \frac{B}{n} \sum_{\ell=1}^{n/B} \mathrm{HSIC}_v(\boldsymbol{f}_k^{(\ell)}, \boldsymbol{y}^{(\ell)}).$$

Github: https://github.com/riken-aip/pyHSICLasso

実験結果:



理論: 大域的最適解がもとまる.

アルゴリズム: 数百万次元×数万サンプルからの特徴選択が可能. 少数特徴で高い予測精度.

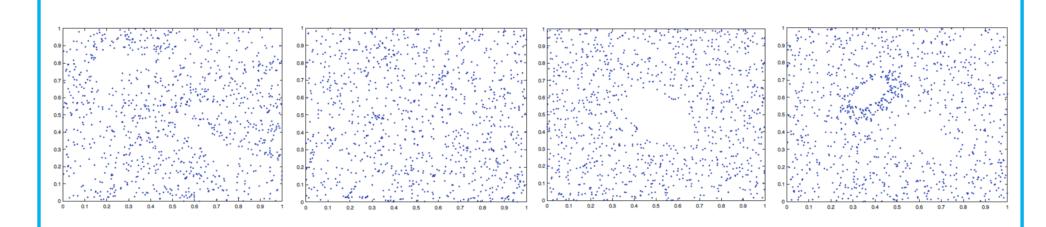
発表実績: BioRxiv 2019. ISMB 2019投稿.

応用展開: 骨髄性急性白血病データの解析(ヒト疾患モデル研究チーム), アトピー性皮膚炎データの解析(健康医療データAI予測推論開発ユニット).

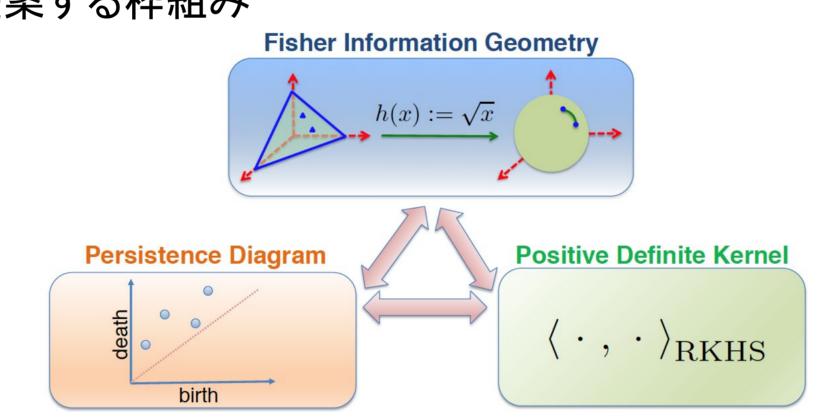
理論展開: 1. 理論解析. 2. グラフデータに対する特徴 選択の問題に拡張.

Persistent Fisher Kernel [NeurIPS 2018]

位相データ解析 (TDA)



提案する枠組み



Persistence Fisher Kernel

$$k_{\mathrm{PF}}(\mathrm{Dg}_i,\mathrm{Dg}_j) := \exp\left(-td_{\mathrm{FIM}}(\mathrm{Dg}_i,\mathrm{Dg}_j)\right)$$

Algorithm 1 Compute d_{FIM} for persistence diagrams

Input: Persistence diagrams Dg_i , Dg_j , and a bandwith $\sigma > 0$ for smoothing **Output:** d_{FIM}

1: Let $\Theta \leftarrow \mathrm{Dg}_i \cup \mathrm{Dg}_{j\Delta} \cup \mathrm{Dg}_j \cup \mathrm{Dg}_{i\Delta}$ (a set for smoothed and normalized measures)

2: Compute $\bar{\rho_i} = \rho_{\left(\mathrm{Dg}_i \cup \mathrm{Dg}_{j\Delta}\right)} \leftarrow \left[\frac{1}{Z} \sum_{u \in \mathrm{Dg}_i \cup \mathrm{Dg}_{j\Delta}} \mathbb{N}(x; u, \sigma I)\right]_{x \in \Theta}$ where $Z \leftarrow \sum_{x \in \Theta} \sum_{u \in \mathrm{Dg}_i \cup \mathrm{Dg}_{j\Delta}} \mathbb{N}(x; u, \sigma I)$

MPEG7

- 3: Compute $\bar{\rho_j} = \rho_{\left(\mathrm{Dg}_j \cup \mathrm{Dg}_{i\Delta}\right)}$ similarly as $\bar{\rho_i}$.
- 4: Compute $d_{\text{FIM}} \leftarrow \arccos\left(\left\langle\sqrt{\bar{\rho_i}}, \sqrt{\bar{\rho_j}}\right\rangle\right)$ where $\langle\cdot,\cdot\rangle$ is a dot product and $\sqrt{\cdot}$ is element-wise.

 $73.33 \pm 4.17 \mid 72.38 \pm 2.41$

 74.83 ± 4.36 76.63 ± 0.66

Orbit

実験結果: 分類, 変化点検出

 k_{PSS}

	k_{SW} Prob+ k_{G} Tang+ k_{G}	$76.83 \pm 3.$ $55.83 \pm 5.$ $66.17 \pm 4.$ $80.00 \pm 4.$	45 72.89 01 77.32		
k _{PSS}	k _{PWG}	k _{sw}	Prob + k _G	Tang + k _G	k _{PF}
Granular packing system					
0 10 20 30 (id = 23)		0 10 20 30 (id = 23)	0 10 20 30 (id = 23)		

理論:パーシステント図間の類似度を精度よく測ることが可能.

発表実績: 難関国際会議NeurIPS 2018発表.

応用展開: 材料科学分野への応用.

理論展開: 1. 手法の高速化. 2. ベイズ最適化.

mmdInf [ICLR 2019]

Post Selection Inference: 選択したデータセットが統計的有意かどうかを検定する.

- H_0 : $\sum_{s=1}^d \eta_s \widehat{D}(\boldsymbol{X}^{(s)}, \boldsymbol{Y}^{(s)}) = 0 \mid \mathcal{S} \text{ was selected},$
- H_1 : $\sum_{s=1}^d \eta_s \widehat{D}(\boldsymbol{X}^{(s)}, \boldsymbol{Y}^{(s)}) \neq 0 \mid \mathcal{S} \text{ was selected},$

特徴

- 1. 二標本検定を複数データに拡張。
- 2. ノンパラメトリック.
- 3. 不完全U統計量Maximum Mean Discrepancy (MMD).

Theorem 1 (Lee et al., 2016) Consider a stochastic data generating process $\mathcal{Z} \sim N(\mu, \Sigma)$. If a feature selection event is characterized by $A\mathcal{Z} \leq b$ for a matrix A and a vector b that do not depend on \mathcal{Z} , then, for any fixed vector $\eta \in \mathbb{R}^d$,
切断正規分布の累積密度関数

切断正規分布の累積密度関数
$$F_{\boldsymbol{\eta}^{\top}\boldsymbol{\mu},\boldsymbol{\eta}^{\top}\boldsymbol{\Sigma}\boldsymbol{\eta}}^{[V^{-}(\boldsymbol{A},\boldsymbol{b}),V^{+}(\boldsymbol{A},\boldsymbol{b})]}(\boldsymbol{\eta}^{\top}\boldsymbol{\mathcal{Z}}) \mid \boldsymbol{A}\boldsymbol{\mathcal{Z}} \leq \boldsymbol{b} \sim \mathrm{Uni}(0,1),$$

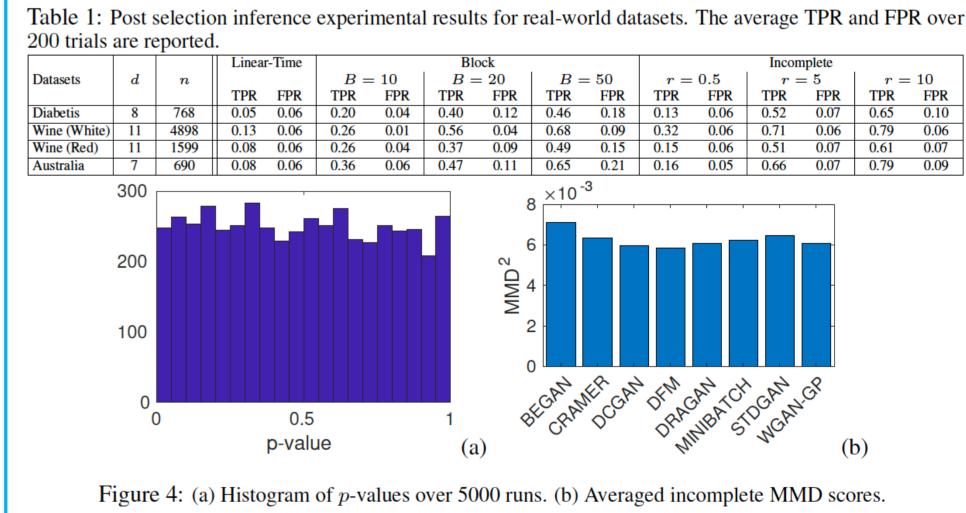
where $F_{t,u}^{[v,w]}(\cdot)$ is the cumulative distribution function of the uni-variate truncated normal distribution with the mean t, the variance u, and the lower and the upper truncation points v and w, respectively. Furthermore, using $\mathbf{c} := \frac{\Sigma \eta}{\eta^{\top} \Sigma \eta}$, the lower and the upper truncation points are given as

$$V^{-}(\boldsymbol{A}, \boldsymbol{b}) := \max_{j:(\boldsymbol{A}\boldsymbol{c})_{j} < 0} \left\{ \frac{b_{j} - (\boldsymbol{A}\boldsymbol{z})_{j}}{(\boldsymbol{A}\boldsymbol{c})_{j}} \right\} + \boldsymbol{\eta}^{\top} \boldsymbol{z},$$

$$V^{+}(\boldsymbol{A}, \boldsymbol{b}) := \min_{j:(\boldsymbol{A}\boldsymbol{c})_{j} > 0} \left\{ \frac{b_{j} - (\boldsymbol{A}\boldsymbol{z})_{j}}{(\boldsymbol{A}\boldsymbol{c})_{j}} \right\} + \boldsymbol{\eta}^{\top} \boldsymbol{z}.$$

$$(2)$$

実験結果:



理論:標本数が十分に大きければ, Type I errorを理論的にコントロールできる.

アルゴリズム: シンプルであり,標本数が大きい場合にも高速に計算可能.

発表実績: 難関国際会議ICLR 2019採択.

応用展開:薬剤誘発による細胞の形態変化解析に利用するため,実証実験を開始. Incomplete MMDを用いたアルゴリズムの開発.

理論展開: 高次元小標本の場合への拡張.