パーシステント図に対する統計的機械学習

福水 健次 統計数理研究所

ENCOUNTER with MATHEMATICS 2017年12月22-23日 中央大学

1



1. イントロ

2. データ解析: カーネル法を中心として ・データ解析とは何か? ・カーネル法の基礎

3. パーシステント図のデータ解析

4. おわりに

TDA: 拡大する実応用 Big Dataの時代: 複雑な幾何構造を持つデータ → 特徴量/記述子の 導入が困難





Phylogenetic Tree of Life



Figure omitted

グラフィクス・画像

Figure omitted

遺伝子組み換え (Cámara et al. *PLOS Comp. Bio*. 2016) 葉の形のデータ Li et al. *Nature* 2017 形状の記述 (Freedman & Chen 2009)

音楽

Figure omitted

Tonnetzによる音楽表現 (Bergomi et al. *CTIC2016*)





統計的データへのトポロジーの応用?



パーシステントホモロジー

• すべての ε (スケール)を同時に考えられる. $X_{\varepsilon} \coloneqq \cup_{i=1}^{m} B_{\varepsilon}(x_{i})$



 ${ \longleftrightarrow }$

真の構造由来の2個のリング(1次元ホモロジー群の生成元) は 長いインターバルで存在するはず

データ解析とは何か?



Analysis of data is a process of inspecting, cleaning, transforming, and modeling data with the goal of highlighting useful information, suggesting conclusions, and supporting decision making. --*Wikipedia*



教科書的な(線形な)データ解析

• 数値のテーブル



行列による表示

$$\boldsymbol{X} = \begin{pmatrix} X_1^{(1)} \cdots & X_m^{(1)} \\ X_1^{(2)} \cdots & X_m^{(2)} \\ \vdots \ddots & \vdots \\ X_1^{(N)} \cdots & X_m^{(N)} \end{pmatrix}$$

m:次元 (列) N:データ数 (行)

- ・線形代数によるデータ解析
 - 相関
 - 線形回帰
 - 主成分分析 etc.



•例1:主成分分析(Principal component analysis, PCA)

PCA: データを低次元表現する方法 分散を最大にする部分空間に正射影する.

-1

-3

-3

1st direction =
$$\operatorname{argmax}_{\|a\|=1} \operatorname{Var}[a^T X]$$

 $\operatorname{Var}[a^T X] = \frac{1}{N} \sum_{i=1}^{N} \left\{ a^T \left(X^{(i)} - \frac{1}{N} \sum_{j=1}^{N} X^{(j)} \right) \right\}^2 = a^T V_{XX} a.$
 $V_{XX} = \frac{1}{N} \sum_{i=1}^{N} \left(X^{(i)} - \frac{1}{N} \sum_{j=1}^{N} X^{(j)} \right) \left(X^{(i)} - \frac{1}{N} \sum_{j=1}^{N} X^{(j)} \right)^T$

Xの分散共分散行列 (*m*×*m* matrix)

・第1主成分方向 = arg max $a^T V_{XX} a$ = argmax_{||a||=1} $a^T V_{XX} a$ = u_1 (最大固有値に対する固有ベクトル)

• 第*p*主成分方向 = 第*p*固有値に対する固有ベクトル

PCA \implies 分散共分散行列 V_{XX} の固有値問題

 具体例:Wine データ (taken from UCI Machine Learning Repository) 178種のイタリアワイン 13個の化学成分を計測 N = 178, m = 13.

→ PCA: 2次元の射影による可視化





- 例2:線形識別
 - 2クラス識別問題
 画像の識別:



Car? → Yes/No

• 訓練データ Input data Class label $\mathbf{X} = \begin{pmatrix} X_1^{(1)} \cdots X_m^{(1)} \\ X_1^{(2)} \cdots X_m^{(2)} \\ \vdots \ddots \vdots \\ X_1^{(N)} \cdots X_m^{(N)} \end{pmatrix} \qquad Y = \begin{pmatrix} Y^{(1)} \\ Y^{(2)} \\ \vdots \\ Y^{(N)} \end{pmatrix} \in \{\pm 1\}^N$



・ 学習: 線形識別関数を求める

 $h(x) = \operatorname{sgn}(a^T x + b)$ such that $h(X^{(i)}) = Y^{(i)}$ for all (or most) *i*.

例) Fisher判別関数, Support Vector Machine (SVM), ロジスティック回帰 13

線形で十分か?



14

Body-Mass-Index vs Cancer Risk (death rate)

肥満とがんとの関係





Japan Public Health Center-based Prospective Study, National Cancer Center



データの非線形変換

・ 高次情報の抽出

 $(X, Y, Z) \rightarrow (X, Y, Z, X^2, Y^2, Z^2, XY, YZ, ZX, \cdots)$

• 元の空間の次元が高いと 計算は実現できない!

e.g. 10000 次元のデータ, 2 次までの特徴

 $_{10000}C_1 + _{10000}C_2 = 50,005,000$

・計算量爆発.
 より効率的な方法が必要 → カーネル法

非ベクトル的データ

- データ形式の多様化
 - 文字列(String) She keeps her room clean. Alice gave a present to Bob.
 - ・ヒストグラム
 - 画像のRGBヒストグラム表現

標準的なデータ解析手法は ベクトルデータを想定. 他の形式にはそのままでは 適用できない

カー



 $\Phi: \Omega \rightarrow H$

- ・内積計算が陽に可能。 $\langle \Phi(x), \Phi(y) \rangle = k(x, y)$ kernel trick
- ・多くの線形データ解析の計算は内積に依拠している.

正定値カーネル

<u>Def.</u>

Ω: set. $k: \Omega \times \Omega \rightarrow \mathbf{R}$ is a positive definite kernel if it satisfies

- 1) (symmetry) k(x, y) = k(y, x)
- 2) (positivity) for arbitrary points $x_1, ..., x_n$ in Ω , the Gram matrix $\begin{pmatrix}
 k(x_1, x_1) & \cdots & k(x_1, x_n) \\
 \vdots & \ddots & \vdots \\
 k(x_n, x_1) & \cdots & k(x_n, x_n)
 \end{pmatrix}$ is positive semidefinite,

i.e.,
$$\sum_{i,j=1}^{n} c_i c_j k(x_i, x_j) \ge 0$$
 for any $c_i \in \mathbf{R}$.

- R^m上の正定値カーネルの例
 - Euclid内積 $k(x, y) = x^T y$

• 多項式カーネル
$$k_P(x, y) = (c + x^T y)^d \quad (c > 0, d \in \mathbb{N})$$



• Gaussian カーネル

$$k_{G}(x, y) = \exp\left(-\left\|x - y\right\|^{2} / \sigma^{2}\right)$$
• Laplace カーネル

$$k_{L}(x, y) = \exp\left(-\alpha \sum_{i=1}^{m} |x_{i} - y_{i}|\right)$$

$$(\alpha > 0)$$



 カーネルトリック(内積の陽な計算)を成り立たせる関数は、正 定値カーネルである.

定理 (Moore-Aronszajn)
Ω上の正定値カーネル
$$k$$
 に対し, Ω上の関数からなるHilbert空間 H_k
が一意に存在して, 次が成り立つ.
1) $k(\cdot, x) \in H_k$ ($\forall x \in \Omega$).
2) span { $k(\cdot, x) \mid x \in \Omega$ } は H_k で稠密
3) (再生性)
 $\langle f, k(\cdot, x) \rangle = f(x)$ for any $f \in H_k, x \in \Omega$.

- 上の H_k を k が定める再生核ヒルベルト空間(reproducing kernel Hilbert space, RKHS)という. c.f. L^q 空間
- 再生性より特に, $\Phi(x) = k(\cdot, x)$ とおくと,

 $\langle \Phi(x), \Phi(y) \rangle = k(x, y).$

正定値カーネルによるデータ解析

- 正定値カーネル k を用意
- 特徴写像: $\Phi: \Omega \to H_k, x \mapsto k(\cdot, x)$

 $X_1, \dots, X_n \mapsto k(\cdot, X_1), \dots, k(\cdot, X_n)$

Space of original data Φ

非ベクトルデータの ベクトル化が可能

 $\langle \Phi(X_i), \Phi(X_j) \rangle = k(X_i, X_j)$

- 正定値カーネルを与えれば十分.
 - ・多くのデータ解析手法は、内積計算ができれば適用可能.
 - •特徴写像,特徴ベクトルを陽に知る必要はない.
 - カーネル法の計算は、グラム行列 $(k(X_i, X_j))_{ii}$ による計算となる.

カーネルPCA

- PCA:線形な次元削減 → カーネルPCA:非線形な次元削減 (Schölkopf et al. 1998).

$$f_{\perp}$$
 f_{Φ}

(直交する方向は分散に効いてこない!)

(再生性を使うと) max $Var[\langle f, \Phi(X) \rangle] = \frac{1}{N} c^T \widetilde{K}_X^2 c$ subject to $||f||=1 \Leftrightarrow c^T \widetilde{K}_X c = 1$

$$(\widetilde{K}_{X})_{ij} = k(X^{(i)}, X^{(j)}) - \frac{1}{N} \sum_{b=1}^{N} k(X^{(i)}, X^{(b)}) - \frac{1}{N} \sum_{a=1}^{N} k(X^{(a)}, X^{(j)}) + \frac{1}{N^{2}} \sum_{a,b=1}^{N} k(X^{(a)}, X^{(b)})$$

(中心化Gram行列)

次の形の f_{N} を考えれば十分 $f = \sum_{i=1}^{N} c_{i} \left(\Phi(X^{(i)}) - \frac{1}{N} \sum_{j=1}^{N} \Phi(X^{(j)}) \right)$

・中心化Gram行列
$$\widetilde{K}_{X}$$
の計算
・ \widetilde{K}_{X} の固有分解 $\widetilde{K}_{X} = \sum_{i=1}^{N} \lambda_{i} u_{i} u_{i}^{T}$
 $\lambda_{1} \geq \lambda_{2} \geq \cdots \geq \lambda_{N} \geq 0$ eigenvalues
 $u_{1}, u_{2}, \dots, u_{N}$ unit eigenvectors
• 第p主成分方向 $f_{p} = \sum_{j} \frac{1}{\sqrt{\lambda_{p}}} u_{pj} \widetilde{\Phi}(X^{(j)}),$
 $\widetilde{\Phi}(X^{(j)}) = \Phi(X^{(j)}) - \frac{1}{N} \sum_{b=1}^{N} \Phi(X^{(b)})$:中心化特徴ベクトル
• $X^{(i)}$ の第p主成分 = $\langle f_{p}, \widetilde{\Phi}(X^{(i)}) \rangle = \sum_{j} \frac{1}{\sqrt{\lambda_{p}}} u_{pj} \widetilde{K}_{ji} = \sqrt{\lambda_{p}} u_{pi}$

カーネルPCAの例

Wine データ (UCI repository)
 クラスの情報はカーネルPCAには用いていない



Kernel PCA (Gaussian)



29

0.6

0.4

パーシステント図のデータ解析

PHによる統計的データ解析

•<u>単純な</u>位相的データ解析



<u>統計的な</u>位相的データ解析

(Kusano, Fukumizu, Hiraoka ICML2016; Reininghaus et al CVPR 2015; Kwitt et al NIPS2015; Fasy et al 2014)



パーシステント図のベクトル化

$$\Pi \coloneqq \{(b,d) \in \overline{\mathbf{R}}^2 | d > b\}, \Delta = \partial \Pi.$$

パーシステント図 $D = D_o \cup \Delta$
 $D_o = \{x_i\}$: multiset on Π. (重複度込み)
Γ: $|D_o| < \infty$ なるパーシステント図全体

• PD = 離散測度と同一視できる
$$D \leftrightarrow \sum_{x_i \in D_o} \delta_{x_i}$$
 δ : Diracのデルタ関数

・カーネル埋め込み $k: \Pi \bot の正定値カーネル$ $\varepsilon_k: \Gamma \rightarrow H_k,$ $\sum_i \delta_{x_i} \mapsto \sum_i k(\cdot, x_i) \in H_k$





カーネル埋め込み (Muandet, Fukumizu, Sriperumbudur, Schölkopf 2017)

(Ω, B): 可測空間, k: Ω上の可測で有界な正定値カーネル. *M*: (Ω, B)上の有限測度の族.以下の写像をカーネル埋め込みと呼ぶ.

 $\varepsilon_k : \mathcal{M} \to H_k, \qquad \mu \mapsto \int k(\cdot, x) d\mu(x)$

<u>Def.</u> 可測で有界な正定値カーネルkが \mathcal{M} -<u>特性的(Characteristic)</u>である とは, ε_k が単射であることをいう.

Ω が R^m の部分集合で M がBorel 測度全体のときは、単に「特性的」ということにする.



- ベクトル化の実用的メリット
 - ガウスカーネルなどを用いると、単射なので測度の識別性は保たれる.
 - ベクトルデータに対する多くのデータ解析手法が適用できる
 - カーネルトリックにより、計算はグラム行列計算に還元される.

特性的なカーネルの特徴付け

<u>定理(Bochner)</u>

 $\psi \in \mathbf{R}^m$ 上の連続関数とする時、 ψ (**C**-値)正定値関数であることと、 \mathbf{R}^m 上の有限非負Borel測度 Λ が存在して

$$\psi(z) = \int \exp(\sqrt{-1}\omega^T z) d\Lambda(\omega)$$

と書けることは同値である.このときΛは一意に決まる.

<u>定理 (Sriperumbdur, Gretton, Fukumizu, Scholkopf, Lanckriet 2011)</u> \mathbf{R}^m 上の連続で平行移動不変な(**C**-値)正定値カーネル $k(x,y) = \psi(x - y)$ (ψ は正定値関数で上の積分表現を持つ)に対し, k が特性的であること と Supp(Λ) = \mathbf{R}^m とは同値である.

Proof idea: $\psi * \mu = 0 \iff \hat{\psi}\hat{\mu} = 0$. Example: ガウスカーネル

Persistence Weighted Gaussian Kernel: PD用のカーネル (Kusano, Fukumizu, Hiraoka ICML2016)

アイデア:対角線に近い生成元はノイズの可能性が高い → 重みを小さくする

$$k_{PWG}(x,y) = w(x)w(y)\exp\left(-\frac{\|y-x\|^2}{2\sigma^2}\right)$$

 $x = (x_1, x_2), y = (y_1, y_2) \in \mathbf{R}^2$

重み関数
$$w(x) = w_{C,p}(x) \coloneqq \arctan(C\operatorname{Pers}(x)^p)$$

(C, p > 0)
 $\operatorname{Pers}(x) \coloneqq d - b$ for $x \in \{(b,d) \in \mathbb{R}^2 | d \ge b\}$

PWGKによる埋め込みも特性的



Pers(x1

距離関数の安定性

• 安定性(stability)

d: PD上の距離 e.g. bottleneck distance

点集合上の距離として Hausdorff距離 D_H を考える

点集合の微小な変化はPDの 距離に微小な変化し与えない

PD ^ª

d が安定とは、ある定数 L が存在し、任意の点集合 X = $\{x_i\}, Y = \{y_j\}$ に対し $d(PD(X), PD(Y)) \le L \cdot D_H(X, Y)$

Point Cloud

が成り立つことをいう(Lipschitz連続性).

• Hausdorff距離

定義

X,*Y*:距離空間 (Ω,*d*)の集合,

D_H: Hausdorff距離

 $D_H(X,Y) = \max\{\sup_{x \in X} d(x,Y), \sup_{y \in Y} d(y,X)\}$



• <u>定理</u> (Bottleneck距離の安定性 Cohen-Steiner et al 2007; Chazal et al 2014) X, Y: 有限集合, $D_q(X), D_q(Y)$: 対応するq次パーシステント図(q任意) $d_B\left(D_q(X), D_q(Y)\right) \le D_H(X, Y).$

• PWGK の決める距離 $d_k(D_1, D_2) \coloneqq \|\mathcal{E}_k(D_1) - \mathcal{E}_k(D_2)\|_{H_k}$

<u>定理</u>(PWGK距離の安定性. Kusano, Hiraoka, Fukumizu ICML2016) $M: \mathbb{R}^d$ のコンパクト集合. $X \subset M, Y \subset \mathbb{R}^d$:有限集合. p > d + 1ならば, PWG カーネル $k_{p,C,\sigma}$ に対し $d_k(D_a(X), D_a(Y)) \leq A D_H(X, Y).$

ここでAは M, p, d, C, σ のみに依存する定数(X, Yには依存しない)

• ガウスカーネルでは安定性は知られていない.



• 再生核ヒルベルト空間上のガウスカーネル
$$K(\varphi_1, \varphi_2) = \exp\left(-\frac{\|\varphi_1 - \varphi_2\|_{H_k}^2}{2\tau^2}\right)$$

 $K(PD_i, PD_j) = \exp\left(-\frac{\|\varepsilon_k(PD_i) - \varepsilon_k(PD_j)\|_{H_k}^2}{2\tau^2}\right)$

 PD_i , PD_i : Persistence diagrams

PWGKの効率的近似計算

•計算量の問題

PDの点(PHの生成元)の数は多いこともある ($\geq 10^3, 10^4$) m個のPDの解析 $PD_i = \left\{ x_a^{(i)} \right\}_{i=1}^{N_a}$ $K(PD_i, PD_j) = \exp\left(-\frac{\|\varepsilon_k(PD_i) - \varepsilon_k(PD_j)\|_{H_k}^2}{2\tau^2}\right)$ の計算には $\left\| \mathcal{E}_k(PD_i) - \mathcal{E}_k(PD_j) \right\|_{H_1}^2$ $= \sum_{a=1}^{N_i} \sum_{b=1}^{N_i} k\left(x_a^{(i)}, x_b^{(i)}\right) + \sum_{a=1}^{N_j} \sum_{b=1}^{N_j} k\left(x_a^{(j)}, x_b^{(j)}\right) - 2\sum_{a=1}^{N_i} \sum_{b=1}^{N_j} k\left(x_a^{(i)}, x_b^{(j)}\right).$ $\exp\left(-\frac{\|x_a-x_b\|^2}{2\sigma^2}\right)$ の計算回数 = $O(m^2N^2)$ $\rightarrow N \approx 10^4$ では計算量莫大 $N = \max\{N_i | i = 1, ..., m\}$

• Random Fourier feature による効率的近似 (Rahimi & Recht 2008) Bochnerの定理 $\exp\left(-\frac{\|x_a - x_b\|^2}{2\sigma^2}\right) = C \int e^{\sqrt{-1}\omega^T (x_a - x_b)} \left(\frac{\sigma^2}{2\pi}\right) e^{-\frac{\sigma^2 \|\omega\|^2}{2}} d\omega$

積分(平均)をMonte Carlo サンプリングで近似: $\omega_1, \dots, \omega_L$: *i*.*i*.*d*. ~ Q_σ

$$\exp\left(-\frac{\|x_a - x_b\|^2}{2\sigma^2}\right) \approx C \frac{1}{L} \sum_{\ell=1}^{L} e^{\sqrt{-1}\omega_{\ell}^T x_a} \ \overline{e^{\sqrt{-1}\omega_{\ell}^T x_b}}$$

$$\sum_{a=1}^{N_{i}} \sum_{b=1}^{N_{j}} k\left(x_{a}^{(i)}, x_{b}^{(j)}\right) \approx \frac{c}{L} \sum_{\ell=1}^{L} \sum_{a=1}^{N_{i}} w\left(x_{a}^{(i)}\right) e^{\sqrt{-1}\omega_{\ell}^{T} x_{a}^{(i)}} \sum_{b=1}^{N_{j}} w\left(x_{b}^{(j)}\right) e^{\sqrt{-1}\omega_{\ell}^{T} x_{b}^{(j)}} L \operatorname{dim.}$$

計算量 O(LN) → 2層目のGram 行列計算 $O(mLN + m^2L)$. c.f. $O(m^2N^2)$ $L,m \ll N$ なら大きなゲイン

他のベクトル化: Persistent Landscape

• Persistent Landscape (Bubenik, JMLR 2015)

 $\lambda_D(k,t) = \underset{(b,d)\in D}{\max} \Lambda_{(b,d)}(t)$ $\Lambda_{(b,d)}(t) \coloneqq \begin{cases} t-b, & t \in [b,(b+d)/2] \\ d-t, & t \in [(b+d)/2,d] \\ 0, & \text{otherwise} \end{cases}$



Persistence landscape に基づく内積・距離
 D₁, D₂: (q次元) パーシステント図

$$(D_1, D_2)_{PL} \coloneqq \sum_{k=1}^{\infty} \int \lambda_{D_1}(k, t) \lambda_{D_2}(k, t) dt \qquad (L2内積)$$

$$d_{PL,p}(D_1, D_2) \coloneqq \left(\sum_{k=1}^{\infty} \left\|\lambda_{D_1}(k, t) - \lambda_{D_2}(k, t)\right\|_p^p\right)^{1/p} \qquad (Lp \mathbb{E})$$

Persistence landscape 一つの計算量 = O(N²) N:生成元の数
 (Bubenik & Dłotko, 2017)

Persistence Scale-Space Kernel Reininghaus, S. Huber, U. Bauer, and R. Kwitt, CVPR2015

• PSSカーネル: 熱方程式から着想 $k_{PSS}(D_1, D_2) = \frac{1}{8\pi t} \sum_{x_i \in D_1} \sum_{y_j \in D_2} \exp\left(-\frac{\|x_i - y_j\|^2}{8t}\right) - \exp\left(-\frac{\|x_i - \overline{y_j}\|^2}{8t}\right)$ y = (b, d)に対して $\overline{y} = (d, b)$

- *x_i* または *y_j* が対角線に近づくと和の成分は0に近づく.
 PWGKと同様の効果
- 実は、PWGKでweightを特別に取ったものに一致.
- ・ガウスのバンド幅と寿命による重みが1つの
 パラメータ t で制御されている → 自由度低い (59p参照)



さまざまなPDの距離

- Bottleneck distance d_B
- カーネル埋め込みによる距離 d_k
- Persistence Landscapeによる距離 d_{PL} (計算量 $O(N^2)$) N: 生成元の数
- Wasserstein distance of degree p (一般の非負測度に対する距離)

$$W_p(D_1, D_2) \coloneqq \inf_{\gamma: D_1 \to D_2} \left(\sum_{x \in D_1} \|x - \gamma(x)\|_{\infty}^p \right)^{\frac{1}{p}}$$

γ:対角線および重複度を込めた全単射

(計算量 *O*(*N*^{2.5}) (d'Amico 2006))

(計算量 O(N²), 近似によりO(N))

Fact:
$$W_{\infty} = d_B$$

計算量: $O(N^3)$ (N:生成元の数, Kuhn-Munkres アルゴリズム,

 安定性: W_p, d_{PL}に対しても安定性は知られているが、弱い形(厳密な Lipschitz連続ではない.詳細は Cohen-Steiner et al 2010; Bubenik 2015).





 リサンプリングによる信頼集合の構成 P: 確率測度, 連続な密度関数 p(x). M = Supp(P). S_n : $X_1, \dots, X_n \sim P$, i.i.d. Goal: $PD(S_n)$ は確率変数 $\rightarrow d_k(PD(S_n), PD(M))$ を評価せよ. Idea: 安定性 $d_k(PD(S_n), PD(M)) \leq A d_H(S_n, M)$ を用いる. $S_{h,n}^{(\ell)} \coloneqq \{X_{i_1}, \dots, X_{i_h}\}$: $S_n,$ からのサブサンプル. $\ell = 1, \dots, \binom{n}{h} =: N (サイズb サブ)$ サンプルの総数) nが大きいとき, 高い確率で $d_H(S_n, M) \leq d_H(S_{b,n}^{(\ell)}, S_n)$.

確率 ≥ $1 - \alpha + o(1)$ で, 各生成元 $z \in PD(S_n)$ に対し $x \in PD(M)$ が 存在する.





PWGKでも同様の議論が可能

カーネル法によるPD解析の応用例

人工データでの基礎実験

- •2値識別問題
 - 点集合:大きな円S1上の点 小円 S0 上の点は、ある場合と ない場合がある.
 - 正解 $Y = XOR(Z_1, Z_2)$
 - Z₁: SOが存在するか Yes/No

小さなスケール

SO

5 **S1**

 PD_i

Y = 1

noise

- Z₂: S1に対応する生成元が ((b(*S*1)<1 && d(*S*1)>4)? Yes/NO 大きなスケール
- ノイズを付加.
- ・訓練データ100セット, テストデータ100セット
- S V M で 識別 機構 成





• テストデータに対する正解率

-					
	埋込み用のカーネル		2層目のカーネル		
	Kernel	Weight	Linear	Gauss	
-		PWGK: $p = 1$	51.5	83.8	
		PWGK: $p = 5$	50.5	84.8	
$k_{PWG}(x, y) =$ $w(x)w(y)\exp^{-\frac{1}{2}}$ $k_{Lin}(x, y) =$ $w(x)w(y)x^{7}$	Weighted Gauss	PWGK: $p = 10$	49.4	84.8	
	$O\left(-\frac{\ y-x\ ^2}{2\sigma^2}\right)$	W _{pers}	56.5	57.5	$w_{pers} \coloneqq Pers(x)$
		w = 1 (Gauss)	51.2	51.5	
		atan $p = 1$	49.4	52.5	
		atan $p=5$	50.5	50.5	
	Weighted Linear	atan $p = 10$	50.5	51.5	
	ν	W _{pers}	49.4	51.5	
	, ,	w = 1 (Linear)	48.4	57.5	
	PSSK (K _{PSS})		50.5	58.5	
	Persistence La	ndscape (K_{PL})	52.5	54.5	

応用例1:タンパク質の識別

- ・タンパク質の幾何的形状 → 機能に大きな役割
- PDをデータとして識別器を構成. サポートベクターマシン(SVM) 識別問題の defacto standard.

データA: Protein-drug 結合

Figure omitted

 A型インフルエンザウイルスのM2 channel 抗インフルエンザ薬のターゲット

> Cang, Mu, Wu, Opron, Xia, Wei, *Molecular Based Mathematical Biology* (2015) Fig. 3

- Rimantadine が結合しているか否かを, M2 channel のNMR(核磁気共鳴)による立体構造データから識別.
 - 結合/非結合それぞれ15データ.
 - ランダムに選んだ20データ(各10)でSVMを訓練.残りでテスト. 100ランダムセット行う.

データB: ヘモグロビンの2状態の識別

- X線回折による立体構造データを用いて、2 状態 Relaxed (R) / Taut (T) を識別.
- R:9データ, T:10データ.

Figure omitted

各クラス1データずつをテストに残し、
 残りで訓練、2データでテスト。
 全ての組み合わせに対して行う。

Relaxed (R) Taut (T)

Cang, Mu, Wu, Opron, Xia, Wei, *Molecular Based Mathematical Biology* (2015) Fig. 4

比較:

Cang et al (2015)では、PHに対して、Molecular Topological Fingerprint (MTF)と呼ぶ経験的な13個の記述子を定義.

#	Dim	Description
1	0	2番目に長い生成元の寿命
2	0	3番目に長い生成元の寿命
3	0	寿命の総和
4	0	寿命の平均
5	1	最長生成元の生成点
6	1	最長生成元の寿命
7	1	長さ1.5Å以上の中の最短生成元の生成点
8	1	長さ1.5Å以上の生成元の中間点の平均
9	1	[4.5,5.5]Åに存在する生成元の数/総原子数
10	1	[3.5, 4.5)Åと(5.5, 6.5]Åに存在する生成元の数/総原子数
11	1	寿命の総和
12	1	寿命の平均
13	2	最初の生成元の生成点

- PWGK + ガウスカーネルSVM
 PWGK: p = 5, C とカーネルバンド幅はクロスバリデーションで選択
- 識別結果

クロスバリデーションによる平均正答率

	A. Protein-Drug	B. Hemoglobin
PWGK	100	88.90
MTF*	(nbd) 93.91 / (bd) 98.31	84.50

- MTF の結果は, Cang et al. Molecular Based Mathematical Biology (2015) より抜粋.
- PWGKは, 1次元PHのみ使用.

応用例2:シリカ(SiO₂)の液相-ガラス相

SiO₂を液体から急冷すると,ガラス状態になる. 目的:液相からガラス相に転移する温度を特定したい.

アモルファスの記述子: 標準的なものがない

データ: SiO₂分子動力学(MD)シミュレーション (Nakamura et al 2015 Nanotechnology; Hiraoka et al 2016 PNAS)

- •80個の温度で、3次元原子配置データを取得(スナップショット)
- ・原子の3次元配置データから、PDを計算.
- SiとO原子では、異なる半径の球を用いている.
- 物理学的方法:エンタルピー曲線を描いて、微分の推定を行い、その不連続 点を推定する。正確な推定は難しい。
- ・提案法:PD図のカーネル埋め込みに対する、変化点検出問題として転移点を推定する. (Kusano et al ICML2016)



パーシステント図



ガラス相





- 変化点検出(一般的な話)
 - パラメータtに沿ったデータが変化するt_{*}を 見つける。 $Y_t, t = 1, ..., T.$



- カーネル変化点検出法 (Harchoui et al 2009) Fisher判別スコアを用いる
 - 各 t に対し前後で平均を計算: $\hat{m}_{1:t} = \frac{1}{t} \sum_{i=1}^{t} \Phi(X_i), \hat{m}_{t+1:T} = \frac{1}{T-t} \sum_{i=t+1}^{T} \Phi(X_i).$ • $\Delta_t \coloneqq \left\| (V_{1:t} + V_{t+1:T} + \gamma I)^{-\frac{1}{2}} (\widehat{m}_{1:t} - \widehat{m}_{t+1:T}) \right\|_{H_t}^2$. $(V_{1:t}, V_{t+1:T} : 前後の分散)$

 - Find $\max_{t} \Delta_t$.
- ガラス-液相転移点問題の場合は、 $Y_t = \mathcal{E}_k(D_{\varepsilon_t})$ (t = 1, ..., 80). 無限次元

•SiO₂の転移点検出

 $\operatorname{KFDR}_{n,\ell,\gamma}(\mathcal{D})$ Δ_t PWGK PSSK $\overline{20}$ 3040 5060

検出された変化点 = 3100K Enthalpyによる方法: [2000K, 3500K]

・低次元表現: カーネル主成分分析

カーネル埋め込み(ベクトル)に対して、さらにカーネル主成分分析 を行い、3次元で表示.



液相-ガラス相は、変化点検出の 結果に基づいて色付けした.

まとめと展望

まとめ

- カーネル法によるPDのベクトル化 → さまざまなデータ解析
 - 多数のパーシステント図に対する統計的データ解析
 - PWGK: 生成元のlifetimeを考慮した柔軟なカーネル
 - さまざまな標準的データ解析手法が利用可能
 - 主成分分析,正準相関分析,判別分析,識別,クラスタリング, etc …



- PDの信頼度
 - ・現状の信頼集合の議論はPDの距離に基づく
 → 短い寿命だが統計的に安定な生成元の信頼度は計れない.



安定性の不等式を超えた解析(ランダムPDの分布)が可能か?
 Signal/noiseの分離

- 時系列
 - 代数的: ●──▶●◀──●──▶●◀──●
 - ・統計的: PD_t, PD_{t+1}, PD_{t+2}, … → 統計的な時系列解析

Reference

 Persistence weighted Gaussian kernel for topological data analysis. Genki Kusano, Yasuaki Hiraoka, Kenji Fukumizu.

[Conference paper] Proceedings of The 33rd International Conference on Machine Learning, PMLR 48:2004-2013, 2016. <u>http://proceedings.mlr.press/v48/kusano16.html</u> [Full journal version] <u>arXiv:1601.01741</u> [math.AT]