

数理システム論分野

統計学 機械学習 データサイエンス

教授 下平英寿 助教 劉言

教授プロフィール



下平英寿

@hshimodaira

ドラえもんを作るために統計学と機械学習の研究をしています。麻布→東大→統数研→東工大→阪大→京大情報 / 理研AIP

◎ 大阪府豊中市 ⌂ stat.sys.i.kyoto-u.ac.jp

■ 2009年4月からTwitterを利用しています

329 フォロー中 1,476 フォロワー

ツイート

ツイートと返信

メディア

いいね

スケジュール

- ・研究室のゼミ（全員参加）：毎週水曜 14:45 ~ 16時か17時
- ・テキストの輪読（どれか一つは参加）：週に 1 回。田中研・下平研の合同で開催
- ・研究打ち合わせは隨時（頻度は人によって全く違う）
- ・研究室の滞在時間（人によってぜんぜん違う）

輪読した本

- 2018年度 前期
 - (1) The Elements of Statistical Learning (Second Edition), Hastie, Tibshirani, Friedman (PDF無料, 機械学習の基本!)
 - (2) カーネル法入門—正定値カーネルによるデータ解析, 福水 健次 (数学的)
 - (3) Gaussian Process for Machine Learning, Rasmussen, Williams (PDF無料, ガウス過程と機械学習)
-
- 2018年度 後期
 - (1) Mathematical Foundations of Infinite-dimensional Statistical Models, Giné, Nickl (学内PDF無料, ノンパラ)
 - (2) 確率モデルによる画像処理技術入門, 田中和之 (基礎的な内容)
 - (3) Optimal transport: old and new , Villani (学内PDF無料, 最適輸送理論)



下平英寿
@hshimodaira

研究室のサーバーです。日本で最初に稼働したTesla P100らしいです。みためではM40と区別つきません。最近のGPUはファンレスなんですね。



10:05 - 2016年10月15日

73件のリツイート 97件のいいね



1

73

97



別のツイートを追加



半額程度の鑫 @Chikeki6868 · 2016年10月16日

返信先: @hshimodaira さん

@NVIDIAJapan 研究室のサーバって、やはりスゲーな。見つけてだけで、満足しました。



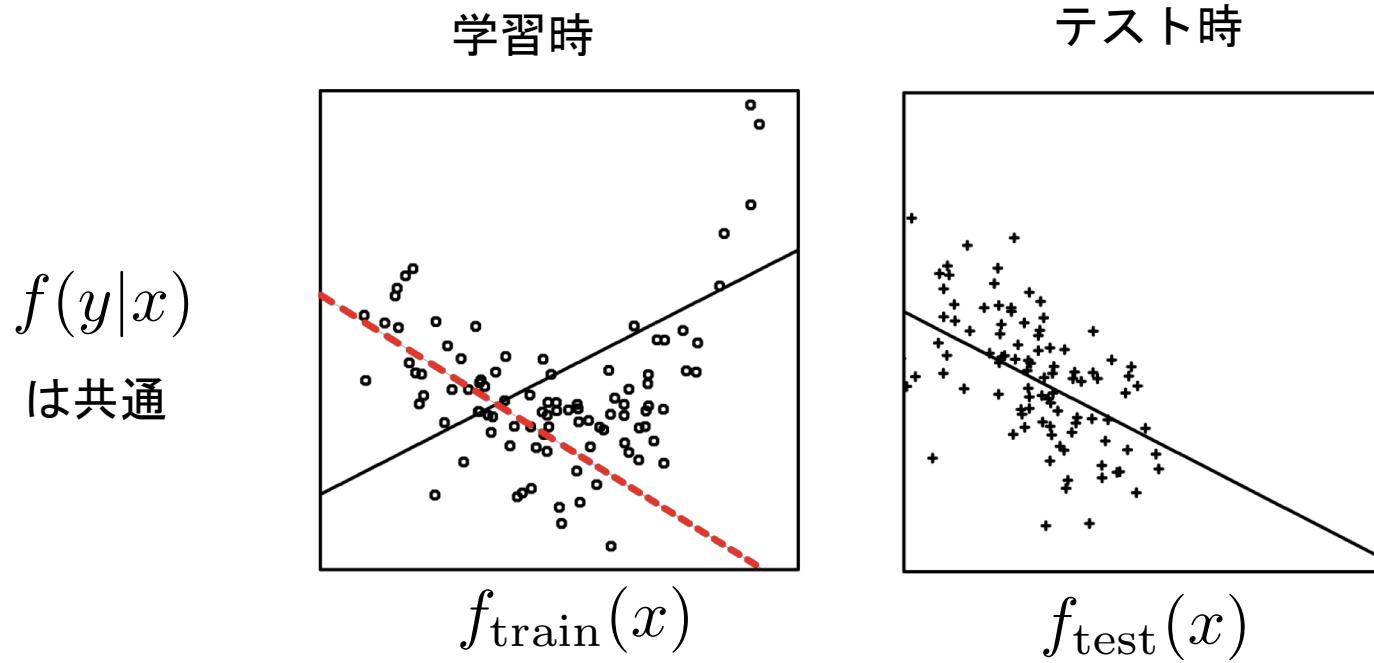
1



研究室のテーマ

統計学と機械学習

共変量シフト : x の分布 $f(x)$ が異なる



データの重み付けは確率密度比で行う : $w(x) = \frac{f_{\text{test}}(x)}{f_{\text{train}}(x)}$

「転移学習」として機械学習で使われている。因果推測にも関係

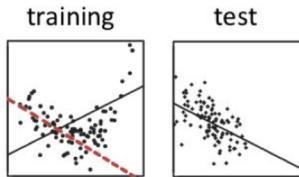
下平英寿
@hshimodaira

v

アニメでも普通に登場するようになって「ディープラーニングやりたい」っていう人があまりに多いので、あえてこんな研究室説明チラシつくってみたよ。研究がどう伸びていくかなんてわからないから、すこしロングスパンでみてもいいんじゃないかな。

stat.sys.i.kyoto-u.ac.jp/post-jp/445/

学習とテストでデータの分布が変わるべきの統計理論をひっそりと考える



確率密度比(density ratio)

$$w(x) = \frac{f_{\text{test}}(x)}{f_{\text{train}}(x)}$$

Shimodaira, Journal of Statistical Inference and Planning 2000



共変量シフトと命名！(covariate shift)

機械学習の分野でよく使われるようになる

Shimodaira (2000)の論文被引用数は700くらい



しらないうちに. . .

ディープラーニングを加速する手法
(Batch Normalization)に組み込まれる

Batch normalization: Accelerating deep network training by reducing internal covariate shift (Ioffe and Szegedy, ICML 2015)

彼らの論文被引用数はたった2年で4100くらい. . .

13:50 - 2018年3月31日

260件のリツイート 654件のいいね



260

654





下平英寿
@hshimodaira

数理統計の専門家としてマジレスしますと、有限のデータから普遍的な真実を導き出そうっていう帰納的推論は難しいんですよ。仮説検定において頻度論のp値もベイズの事後確率も問題があって、人類はまだ完璧な統計的方法には到達してないです。

Nature ダイジェスト、編集部 @NatureDigest

800人以上の科学者が「統計学的に有意」という概念そのものを放棄しようと提案。例えばp値という閾値の下なら効果や差が「あり」で上なら「なし」、などということではなく、そのように簡単に結論してはならないのだが、閾値で分けるやり方が蔓延し、多くの人が結果を誤解している現状...
[このスレッドを表示](#)

18:09 - 2019年3月23日

521件のリツイート 1,221件のいいね



1



521



1,221



下平英寿 @hshimodaira · 3月24日

苦労してあつめたデータを解析して判断する時、どのくらい信頼できるかを数値化したいわけですよ。p値も事後確率も適材適所で役立ってきたけど、本来想定された範囲をこえて使われるようになって問題が顕在化してきます。逆説的ですがデータ解析の普及が大成功したせいかも。



1



30



62



下平英寿
@hshimodaira

「カブトガニがクモの仲間」を統計的に検定するのに提案手法（マルチスケールブートストラップ法）が使われて嬉しい。数理の研究は地味なのでそれ自体が注目されることはない。

「生きた化石」カブトガニ なんとクモの仲間だった | ナショナルジオグラフィック日本版サイト



「生きた化石」カブトガニ なんとクモの仲間だった

原始的な海の生物カブトガニはクモ綱（こう）に属するらしいことが最新の遺伝子解析によって分かった。これまでの進化のストーリーが書き換えられるかもしれない...
[natgeo.nikkeibp.co.jp](#)

19:38 - 2019年3月3日

23件のリツイート 54件のいいね



1



23



54



昔つよかった自慢

2018-04-24 投稿者: SHIMO

Google Scholarの使い方（論文の被引用数）

論文調べで必須と思うけど意外と学生に知られてないのでメモしておきます。普通のGoogleはウェブサイトを検索するために使いますが、Google Scholarは論文や本を検索するために使います。使い方は普通のGoogleと同じで簡単。

Google Scholar



1. <https://scholar.google.co.jp> にアクセス（または“google scholar”を検索）

2. 検索窓に、キーワードとか著者名いれて検索する。たとえば shimodaira を検索すると. . .

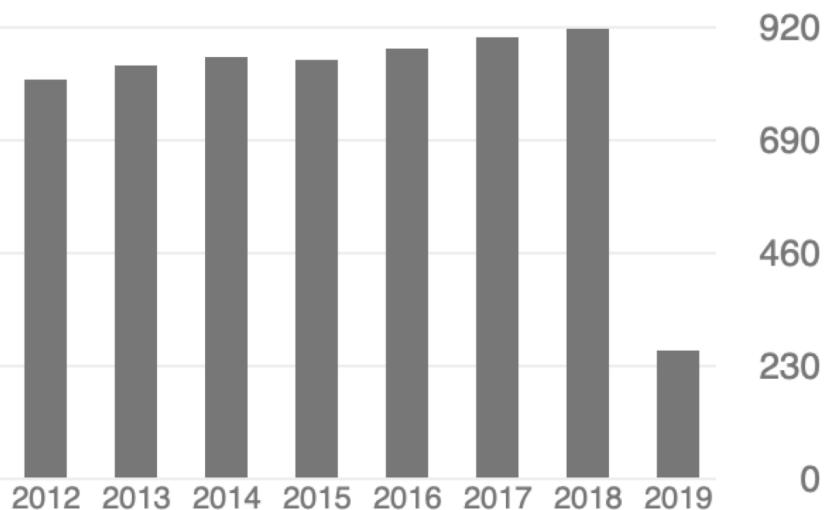
A screenshot of the Google Scholar search results for the query "shimodaira". The search bar at the top contains "shimodaira". Below it, there's a "記事" (Articles) section showing approximately 17,700 results found in 0.06 seconds. A sidebar on the left includes filters for "期間指定なし" (No date range), "2018年以降" (From 2018), "2017年以降" (From 2017), "2014年以降" (From 2014), and "期間を指定..." (Specify period). The main content area displays a list of articles, with the first one being a PDF titled "[PDF] Multiple comparisons of log-likelihoods with applications to phylogenetic inference" by H. Shimodaira, M. Hasegawa - Molecular biology and evolution, 1999 - researchgate.net. The page also includes a "関連性で並べ替え" (Sort by relevance) link and a "日付順に並べ替え" (Sort by date) link.

引用先

すべて表示

すべて 2014 年以来

	引用	12101	4665
	h 指標	23	18
	i10 指標	33	22



いまでもイケてる自慢



下平英寿

@hshimodaira

どうも研究室M1が若手奨励賞と最優秀ポスター賞らしい。研究室のslack画像みるとそう書いてある。ホントならスゴイな。「擬ユークリッド空間への単語埋め込み」言語処理学会第25回年次大会@名古屋 #NLP2019



19:25 - 2019年3月15日

2件のリツイート 22件のいいね



2

2

22



下平英寿

@hshimodaira

研究室メンバーのKIM, GEE WOOK (修士1年) が言語処理学会第25回年次大会 (NLP2019)において若手奨励賞と最優秀ポスター賞を受賞しました。

stat.sys.i.kyoto-u.ac.jp/post-ja/746/

擬ユークリッド空間への単語埋め込み

Kim Geewook, 奥野彰文, 下平英寿

京都大学大学院情報学研究科 / 理化学研究所革新知能総合研究センター



背景

- 表現学習において様々な類似度モデルが試されており、それらの性質が研究されている [OHS19, OKS19]
- 近年 Inner Product Difference Similarity (IPDS) [OKS19] という表現力の高い類似度モデルが提案されています。類似度閾値選択の問題の解決が期待されています

やったこと

- IPDSを単語の表現学習において試し、その有用性を調べた
- 実験としては。
 - WordNetの単語隣接グラフを埋め込みタスクを通じてIPDSの豊富な表現力の有用性を確認した
 - テキストコ-バスから単語を埋め込み、単語類似度タスクを通じて単語間の類似度をより精度高く学習できることを確認した
- IPDSはIPDSを構成する2つの内積に割り当てる次元選択の問題があるので、この点を改善する必要がある
- IPDSの汎化性能についてでは今後もっと調べて行く必要がある
- 上記の実験は最近提案された Weighted Inner-Product Similarity [KOF19] で改善されており、Inductiveな設定を含む豊富な実験によりその有用性も示されています
- 単語埋め込みでtargetとContextの分散表現を分離する/しないに関する問題は今後もっと調べていただきたい

今後の課題

- グラフ再構築タスク

- WordNetよりも複雑な単語の改変構造のグラフをそれぞれの内積で表現する
- 得られた分散表現を用いてそれぞれの類似度モデルに基づいたグラフを再構築する（リンクを予測する）
- リンク予測における ROC-AUC を評価尺度として用いた
- 一般語に割り当てるパラメータ数を変化しながら評価を行った

様々な類似度モデル

表現力



Inner-Product Similarity (IPS)

$$h_K^{IPS}(x, x') = \langle y, y' \rangle = \sum_{k=1}^K y_k y'_k$$

- ここで y は語 x の K 次元の分散表現
- 語語間のヨーリー度量で内積で表される。分散表現の間の内積が相似度を表すようにモデルizing
- 内積に基づいた広く用いられている類似度モデル
- 任意のPositive Definite (PD) カーネル関数を近似できることが知られている [OHS18], e.g., コサイン類似度関数

Shifted Inner-Product Similarity (SIPS)

$$h_K^{SIPS}(x, x') = \langle \tilde{y}, \tilde{y}' \rangle_{K-1} + b + b'$$

- 定数項付与の内積に基づいた類似度モデル
- 定数項を入れることで任意のPDカーネル関数はもちろん、任意

実験

グラフ再構築タスク

- WordNetよりも複雑な単語の改変構造のグラフをそれぞれの内積で表現する
- 得られた分散表現を用いてそれぞれの類似度モデルに基づいたグラフを再構築する（リンクを予測する）
- リンク予測における ROC-AUC を評価尺度として用いた
- 一般語に割り当てるパラメータ数を変化しながら評価を行った

表 1: 表現学習用語彙タスクの評価結果 (SIPSとIPDS)					
モデル	IPDS	SIPS	IPDS	SIPS	IPDS
L2Norm [Zhang et al., 2012]	80.00	81.71	81.80	81.80	81.80
IPDS [Ohsawa et al., 2019]	80.34	84.41	87.05	86.94	87.05
IPDS [Ohsawa et al., 2019] + Kotsopoulos and Kotsopoulos, 2017	80.34	84.41	87.05	86.94	87.05
基準	-	-	-	-	-

単語類似度タスク

- テキストコ-バスの中の単語間の共起情報をそれぞれの類似度モデルで学習し、単語の分散表現を学習する
- それらの類似度モデルによる類似度コストと様々なベンチマークデータセットとのスピアマン順位相関を計算し、性能を比較する

表 2: 単語類似度タスクの評価結果 (SIPSとIPDS)					
モデル	IPDS	SIPS	IPDS	SIPS	IPDS
IPDS [Zhang et al., 2012]	20.00	20.00	20.00	20.00	20.00
IPDS [Ohsawa et al., 2019]	20.00	20.00	20.00	20.00	20.00
IPDS [Ohsawa et al., 2019] + Kotsopoulos and Kotsopoulos, 2017	20.00	20.00	20.00	20.00	20.00
基準	-	-	-	-	-

参考文献

1:52 - 2019年3月23日

18件のリツイート 64件のいいね



2

2

64

2018年度の成果

数理統計学チーム@京大クラスタ



下平研究室(情報学研究科, 京都大学)

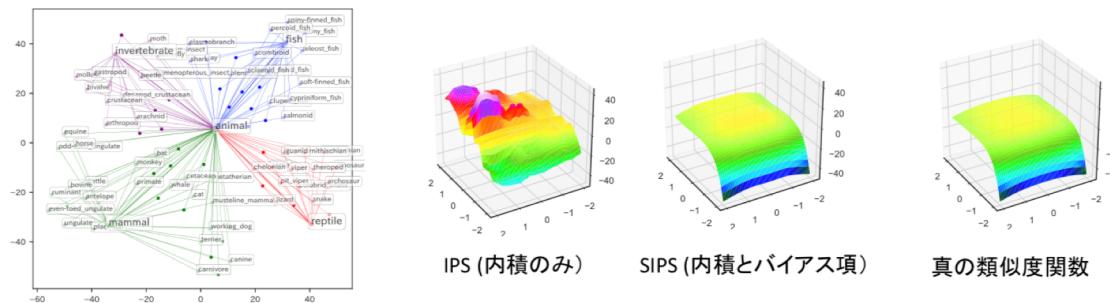
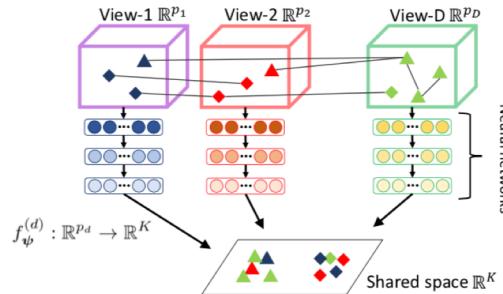
下平英寿(チームリーダー), Thong Pham(特別研究員), 福井一輝(パートタイマー),
奥野彰文(パートタイマー), 井上雅章(パートタイマー), KIM GEEWOOK(パートタイマー),
田中卓磨(パートタイマー), 岩田具治(客員研究員), 寺田吉壱(客員研究員), 伊森晋
平(客員研究員), 廣瀬慧(客員研究員), 白石友一(客員研究員), 劉言(客員研究員)

統計学と機械学習

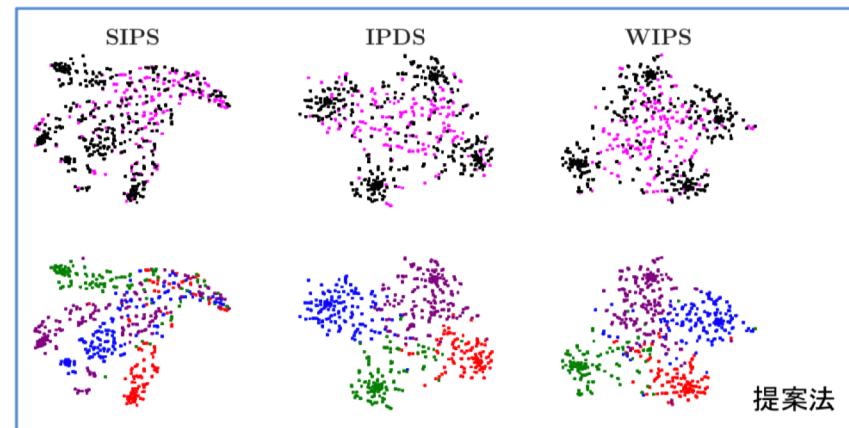
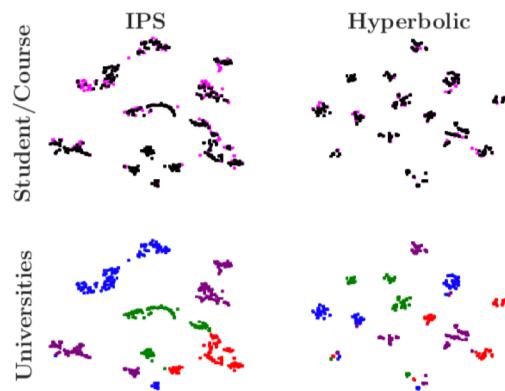
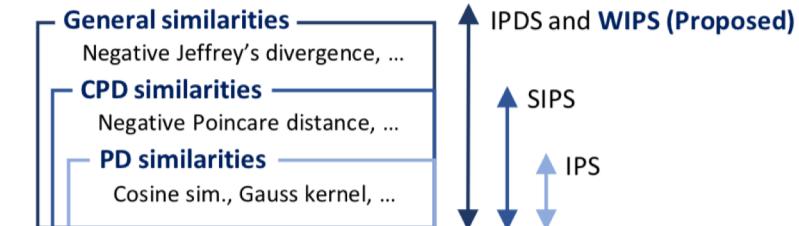
- 統計的仮説検定, 領域の検定, 選択的推測(selective inference)の理論
- マルチスケール・ブートストラップ法などリサンプリング法
- 汎化誤差を共変量シフト, 不完全観測などの設定で推定する情報量規準
- マルチドメインデータの多変量解析, ニューラルネットによるグラフ埋め込み
- 単語埋め込みなど自然言語処理, 画像と単語の埋め込みと相互検索
- 成長ネットワークの統計的推測
- 系統樹, 遺伝子発現データ, ゲノムデータの解析

高い表現力をもつ類似度関数の学習理論

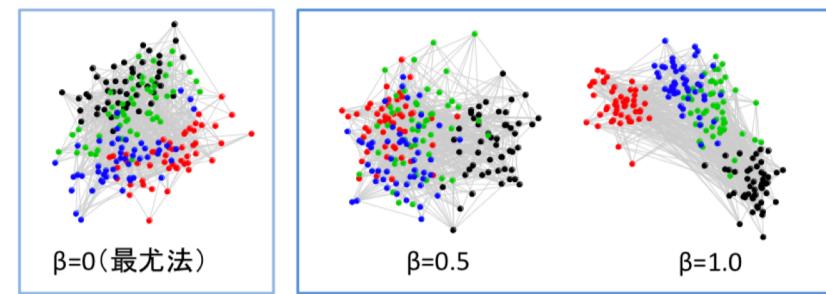
- Okuno, Hada, Shimodaira ([ICML 2018](#)) A probabilistic framework for multi-view feature learning with many-to-many associations via neural networks
- 内積だけでも十分に大きなニューラルネットワークを併用することで広いクラスの類似度関数を学習できることを数学的に証明 (Universal Approximation定理とMercer定理, 推定量の漸近的一致性の証明)
- グラフ埋め込みでの学習アルゴリズム (mini-batch SGD)
- 画像, 単語などマルチドメインの多変量解析への拡張
- 自然言語処理 (単語埋め込み), 画像検索などの応用
- Okuno, Shimodaira ([ICML workshop Theoretical Foundations and Applications of Deep Generative Models 2018](#)) On representation power of neural network-based graph embedding and beyond
- Okuno, Kim, Shimodaira (to appear [AISTATS 2019](#)) Graph Embedding with Shifted Inner Product Similarity and Its Improved Approximation Capability
- 内積類似度 (IPS) が表現できるのは 正定値関数 (positive definite, PD)だけである (例) コサイン類似度
- 内積だけでなくバイアス項も学習するモデル Shifted Inner Product Similarity (SIPS) の提案
- SIPSは任意の条件付き正定値関数 (conditionally PD, CPD)を学習できることを数学的に証明
- CPDは十分に広いクラスである (例) ユークリッド空間の距離, Poincare埋め込み, 双曲空間の距離, (一部の)Wasserstein 距離など, 機械学習でよく利用される距離がたいてい表現できる
- さらに, 2つの内積の差を学習するモデル Inner Product Difference Similarity (IPDS)も提案
- IPDSは不定値 (indefinite)カーネルを含む類似度関数を学習できることを数学的に証明



- Kim, Okuno, Fukui, Shimodaira ([arXiv:1902.10409](#)) Representation Learning with Weighted Inner Product for Universal Approximation of General Similarities
- 重み付き内積を学習するモデル Weighted Inner Product Similarity (WIPS)の提案と高速でスケーラブルな学習法の実装
- 重みの値として負も許して学習することにより、不定値 (indefinite)カーネルを含む類似度関数を学習できることを数学的および実験で証明



- Okuno, Shimodaira (to appear [AISTATS 2019](#)) Robust Graph Embedding with Noisy Link Weights
- 類似度関数学習やグラフ埋め込みがノイズに強くなるようにロバスト化する手法の提案
- 提案手法 (β -グラフ埋め込み)は新規に考案した損失関数(経験積率 β -スコア)の最小化を実行
- ロバスト統計学におけるdensity power divergenceの理論を発展させたもの

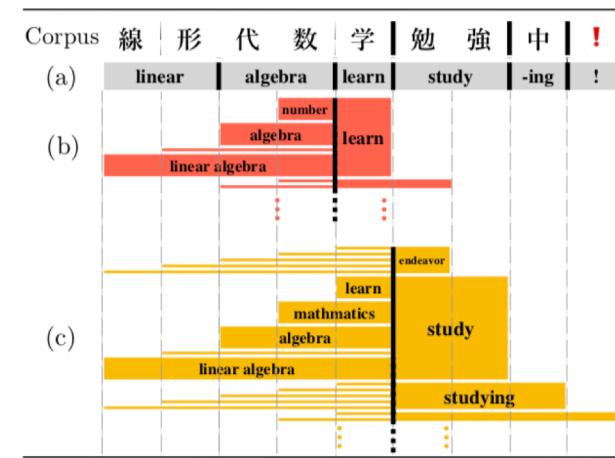
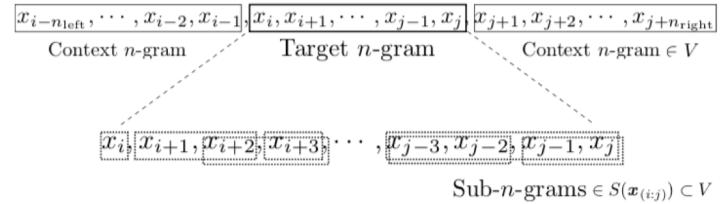
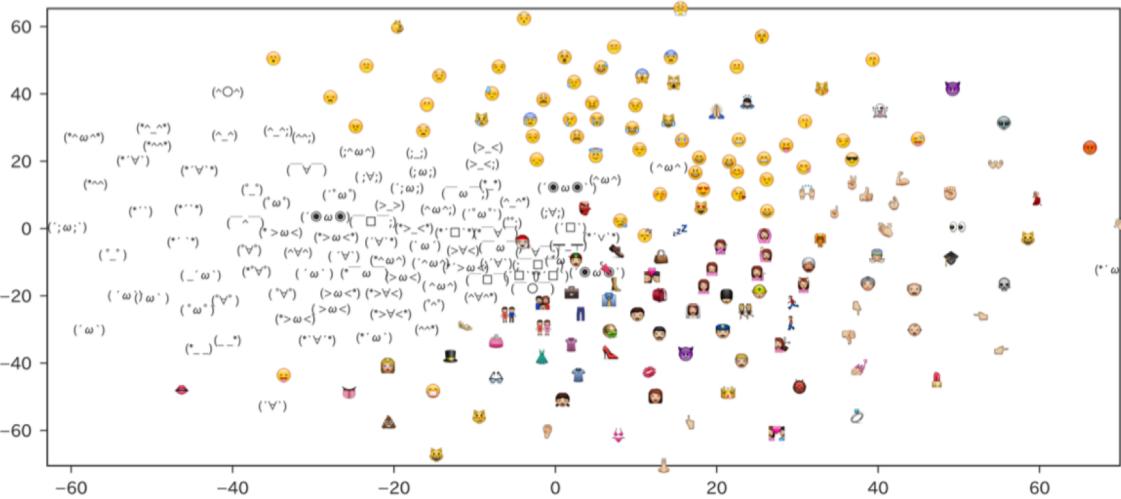


従来法はリンクのノイズに影響される

提案法はリンクにノイズがあっても β を大きくすることでクラスタが分離できている

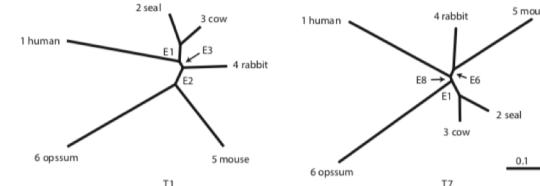
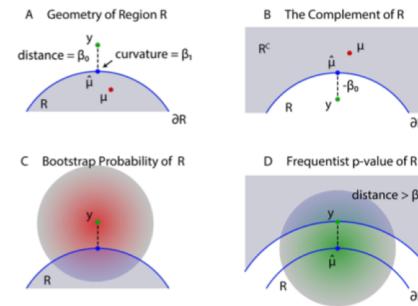
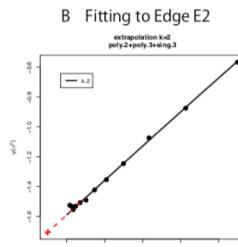
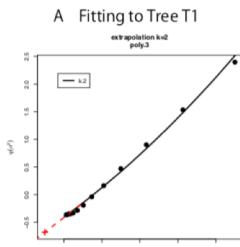
単語分割を必要としない単語埋め込み

- Kim, Fukui, Shimodaira ([EMNLP Workshop on Noisy User-generated Text W-NUT 2018](#)), Word-like character n -gram embedding
 - Kim, Fukui, Shimodaira (to appear [NAACL-HLT 2019](#)), Segmentation-free compositional n -gram embedding
 - 既存の単語埋め込み法では、文書データを前処理によって分割して、単語の列を生成する必要があった。英語などスペースで区切られた言語では問題ないが、日本語、中国語、韓国語では形態素解析による前処理が必要であり、これが障害となることがあった
 - とくにソーシャルメディアでは表記ゆれ、新語の問題がある
 - そこで「単語分割を全く行わない単語埋め込み法」の提案
 - 単語、フレーズ、文に対して連続的に適用できる
 - コーパスのすべての部分文字列に対して埋め込みを効率的に計算

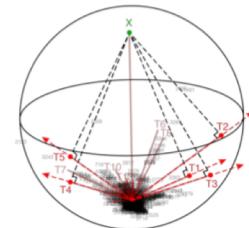


選択的推測の理論とその系統樹推定への応用

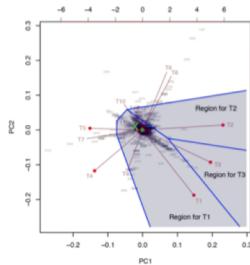
- Shimodaira, Terada ([arXiv:1902.04964](https://arxiv.org/abs/1902.04964)), Selective Inference for Testing Trees and Edges in Phylogenetics
- 従来の統計的仮説検定では、事前に仮説を定める必要がある。ところが医学、科学全般で必ずしもこれが守られず、実際にはデータを見てから仮説を設定し、同じデータを再び用いて仮説検定を実行するため、偽の発見となりやすくその危険性が指摘されてきた。頻度論のp値の代わりにベイズ事後確率を用いてもこの問題は解決しない
- 「データを見てから仮説を選択すること」を設定に組み込んだ選択的推測 (selective inference) の理論が構築されつつある
- 弊チームの先行研究 (Terada, Shimodaira [arXiv:1711.00949](https://arxiv.org/abs/1711.00949)) ではマルチスケール・ブートストラップ法 (Shimodaira 2002, 2004, 2008)によってブートストラップ確率のスケーリング則から選択的推測のp-値を計算する数理統計理論を与えていた
- 本研究ではその手法をもとに分子進化系統樹を推定する問題へ実際に応用し、クラスタリングや系統樹のクレードのp-値において選択的推測による調整の重要性を示した



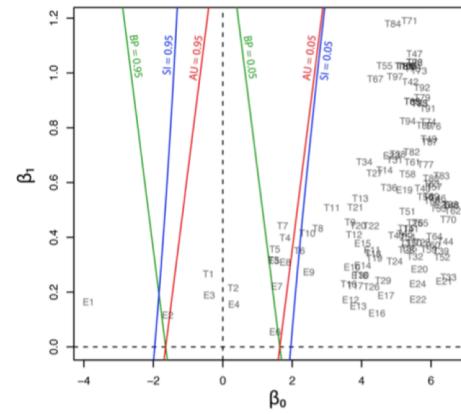
A Projection of data point X to tree models



B Regions for trees with data point X

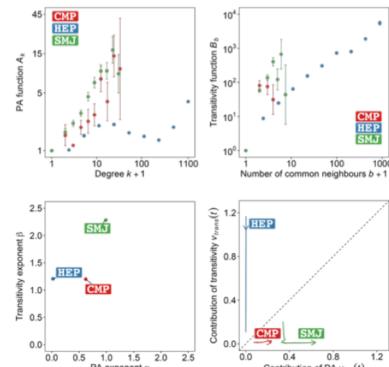


Estimated β_0 and β_1 for Mammal Dataset



複雑ネットワーク成長モデルの統計推測

- Pham, Sheridan, Shimodaira (to appear Journal of Statistical Software 2019), PAFit: an R Package for Estimating Preferential Attachment and Node Fitness in Temporal Complex Networks
- Inoue, Pham, Shimodaira (Unifying Themes in Complex Systems IX. ICCS 2018. Springer Proceedings in Complexity. Springer), Transitivity vs Preferential Attachment: Determining the Driving Force Behind the Evolution of Scientific Co-Authorship Networks
- 複雑ネットワークの成長モデルをネットワーク時系列データから統計推測する手法とソフトウェア(PAFit)の開発
- 優先的選択(知名度)と適応度(能力)の効果または優先的選択(知名度)と推移度(共通の友人)の効果のノンパラメトリック同時推定
- 共通する論文共著者が一人でもいれば、次に共著する確率は急増

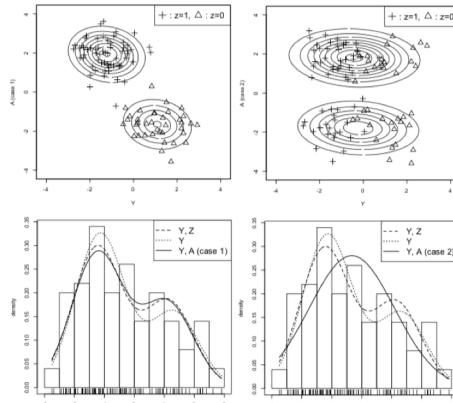


共著者ネットワークデータの分野

- CMP:複雑ネットワーク
- HEP:高エネルギー物理
- SMJ:経営戦略論

補助変数選択の情報量規準

- Imori, Shimodaira ([arXiv:1902.07954](https://arxiv.org/abs/1902.07954)), An information criterion for auxiliary variable selection in incomplete data analysis
- 主要変数に加えて学習データで補助変数が観測できるとき、補助変数も利用したほうが主要変数のモデル推定精度が高くなる事がある
- 主要変数に無関係な補助変数はノイズとなり有害
- 主要変数の一部は観測できない
- このような設定で汎化誤差の推定量である情報量規準を初めて導出
- 補助変数や主要変数の選択、もっと一般にモデル選択を行える
- 既存の赤池情報量規準(AIC)などを特殊な場合として含む一般化



良い補助変数
を利用するすると推定精度が向上

悪い補助変数
を利用すると推定精度が悪化

つまり、学生次第というのがホンネ

数理工学コース分野配属のための研究室見学について

工学部総合校舎 下平 または 学生室・劉助教

見学期間（4月5日(金)～12日(金)）に隨時行います。なるべく事前にメール連絡をおねがいします。何人かまとまって来ると対応しやすいです。