

統計的自然言語処理と機械学習

持橋大地 数理・推論研究系 准教授 daichi@ism.ac.jp

“Linguista sum: linguistici nihil a me alienum puto.” / Roman Jakobson.

統計的自然言語処理を中心に、主に離散データの統計モデルの研究をしています。特に、離散データの裏に潜む連続性をいかに扱うかや、人間がカテゴリーをどうやって認識しているかに興味があります。必然的にモデルは教師なし学習を基にしたものになりますが、教師データとみなせる(自然な)データを統計モデルに取り入れることに躊躇はありません。

キーワード: ノンパラメトリックベイズ法, ガウス過程, ポアソン過程, MCMCなど

統計的自然言語処理

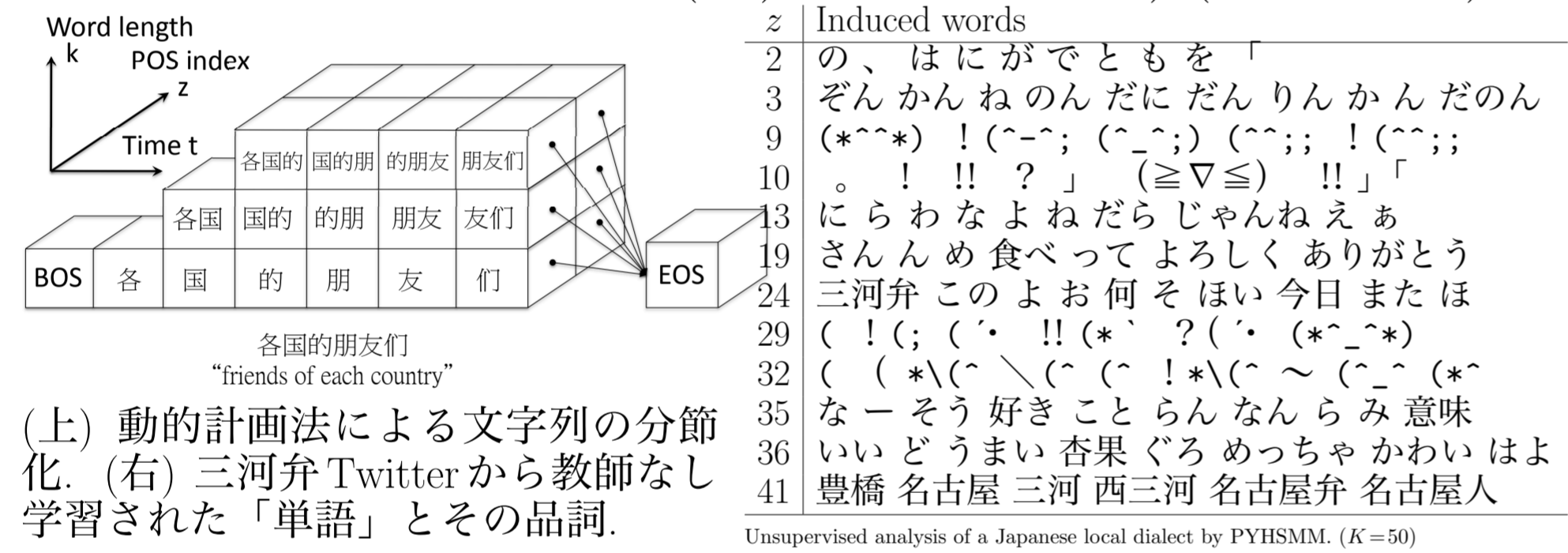
言語の生文字列のみから「単語」を学習する「教師なし形態素解析」(Mochihashi+, ACL 2009)の理論をさらに進め、条件付確率場(CRF)との統合モデルにより、少数の教師データを生かした**半教師あり形態素解析**の研究を行いました。(株)博報堂との共同研究 (Fujii+, TACL 2017)

$$p(c_l^k | s) = \sum_j p(c_l^k, c_{l-1}^j | s) \propto \sum_j \alpha[l-1][j] \cdot \beta[k][k-l+1] \cdot \exp[\lambda_0 \log p(c_l^k | c_{l-1}^j) + \gamma[l, k+1]] = \frac{\alpha[k][k-l+1] \cdot \beta[k][k-l+1]}{Z(s)}$$

东软集团	19	\$K\$ \$s\$	6
游景玉	17	%/%U\$1	6
任亮森	17	%U/%m\$col<	5
南昌铁路	16	https	4
东方红三号"卫星	13	%U/%m\$	4
刘积仁	13	%U/%m\$	3
internet	11	%U/%m\$P	3
东宝	11	\$s\$ \$s\$	3
张肇群	10	4F`	3
彭云	10	%W	3
玲英	10	%U/%m\$col<	2
杭州	10	%U/%m\$ml<	2
亚仿	10	?!\$GJ'\$s\$	2
南丁格尔	9	%U/\$1' \$H	2
中远香港集团	7	%/%U\$ \$	2
		% (%%TR\$	2

(左) ノンパラメトリックベイズ法との統合モデルによる前向き-後向き確率。(右) 中国語微博および日本語 Twitter から学習された単語。

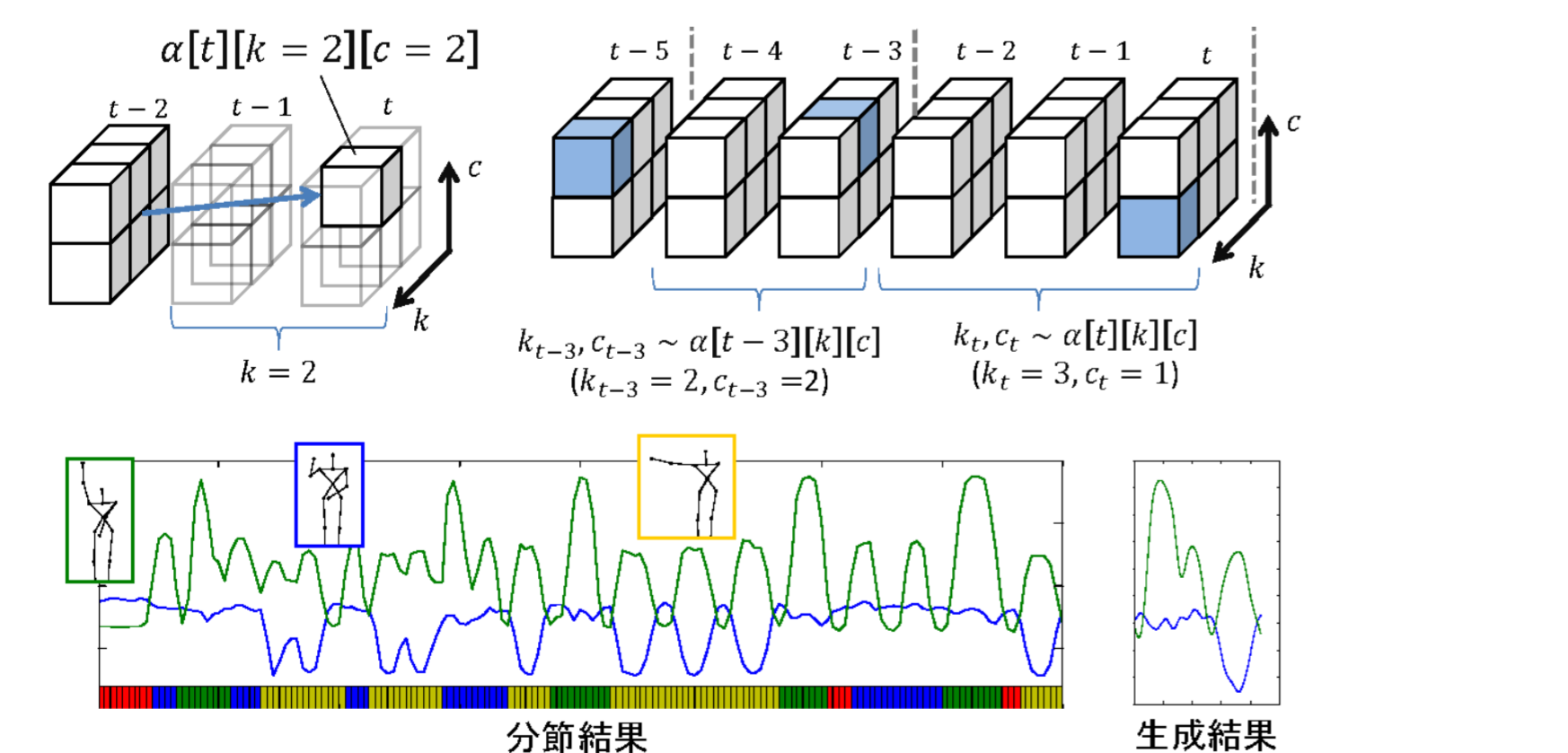
また、単語への分割だけでなく、各単語の**品詞**も教師なし学習する**完全教師なし形態素解析**の研究を行いました。(デンソーアイティラボラトリ(株)との共同研究 (ACL 2015))



Unsupervised analysis of a Japanese local dialect by PYHSM. (K=50)

ロボティクス

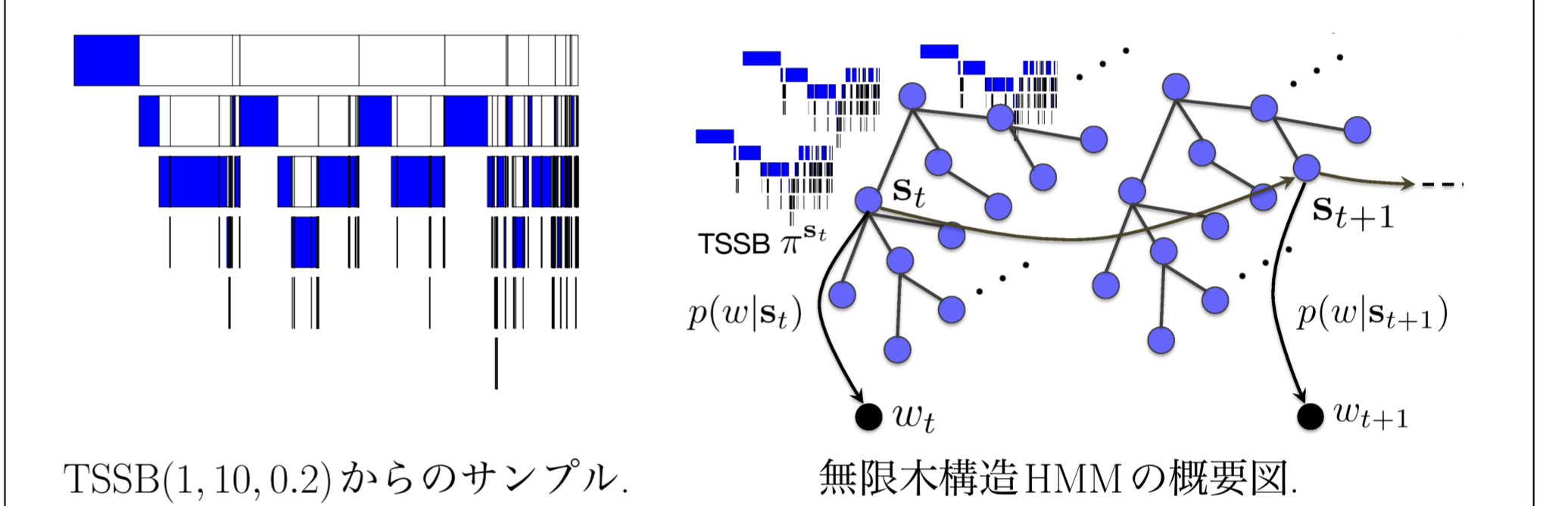
上の教師なし形態素解析を応用して、各分節で文字列ではなく、**ガウス過程**による軌道を生成することで、ロボティクスにおいて関節の角度時系列ベクトルから「動作」を統計的に学習する研究を行っています。これは**隠れセミマルコフモデル**であるため各分節は潜在状態を持ち、さらに**階層ディリクレ過程**を用いることで潜在状態の数も自動的に推定可能としました。この一連の研究は、ロボティクスのトップ国際会議IROSにおいて本会議論文として発表されています(Nakamura+2017, Nagano+ IROS 2018, Nagano+ IROS 2019)。



分節化の結果から得られた「動作」。これは、電通大・お茶の水女子大・阪大との共同研究です。

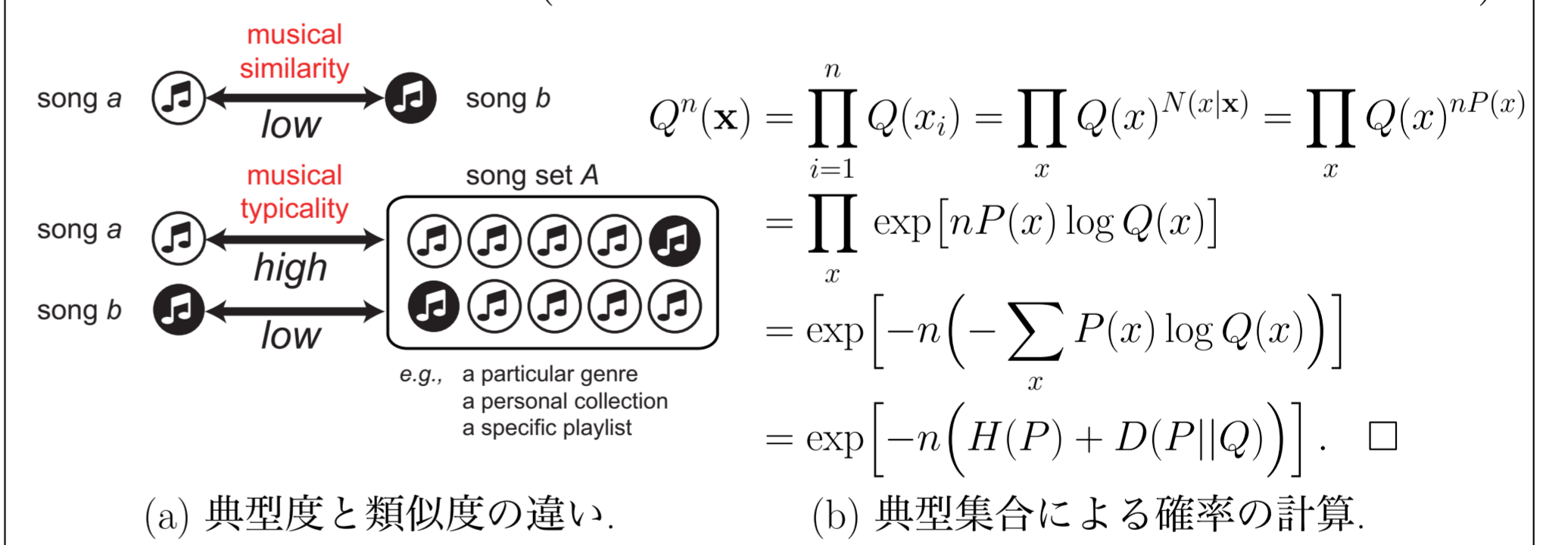
統計的機械学習

離散的な隠れマルコフモデル(HMM)は自然言語処理だけでなく、経済学やバイオインフォマティクスなど様々な分野に使われる基本的な統計モデルですが、状態が“2”“19”のように単純で、その間の相関が表せないという問題がありました。これに対し、木構造Stick-breaking過程(Adams+2010)をさらに階層化することで、無限の深さと無限個の分岐を持つ潜在的な木構造を状態空間に持つ**無限木構造隠れマルコフモデル**を定義し、特別なMCMC法により学習が行えることを示しました。(持橋&能地2016)本研究は2017年の情報処理学会山下研究賞を受賞しています。



音楽情報処理, 音声認識

音声認識においても、音から教師データを介さず、直接「単語」を認識する**教師なし音声認識**の研究を行っています(東工大・篠崎研究室との共同研究)。また、音楽情報処理において曲がどれほど「ありがち」かを定量化することは様々な意味で重要ですが、これを情報理論における**典型集合**としてとらえることで、「ありがち度」を情報理論的尺度として計算する研究を行いました。この研究は、音楽情報処理のトップ国際会議ISMIR 2016に採録されています(産総研 後藤グループとの共同研究)。



その他の分野

この他にも、教育工学, 計量政治学, 動物学, 言語学, 計量国語学, バイオインフォマティクスなどの分野で共同研究を行っています。

現在の学生: 博士後期課程3名 (NEC 中央研究所, (株)AGC, (株)メルペイ), 特別共同利用研究員1名 (慶応大学大学院政治学専攻)。
その他: 日本学術振興会 学術情報分析センター研究員, 日本医療研究開発機構(AMED) 科学技術調査員