
ノンパラメトリックベイズ法による 教師なし形態素解析

持橋大地

NTTコミュニケーション科学基礎研究所

daichi@cslab.kecl.ntt.co.jp

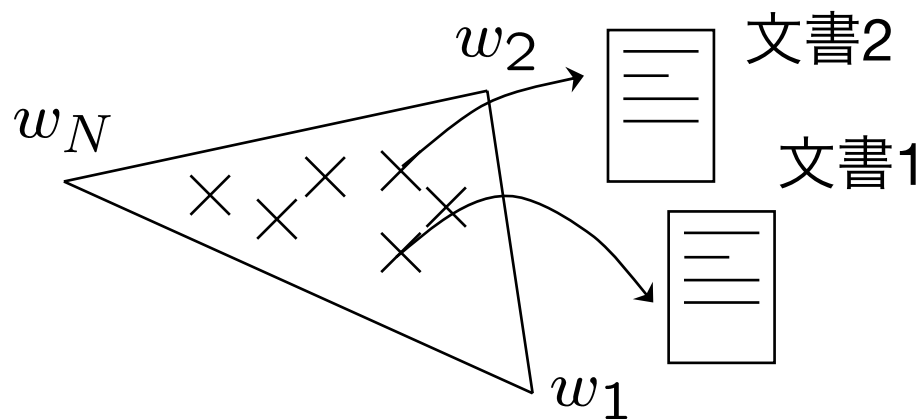
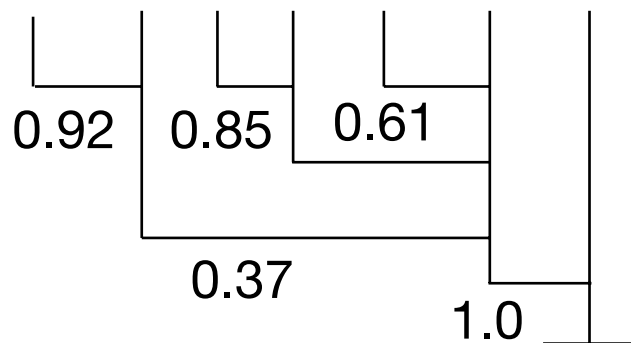
統計関連学会連合大会2009

「Bayes統計モデルのための計算技法とその応用」

自己紹介

- 研究分野：統計的自然言語処理
- 統計的自然言語処理とは
 - 大量のテキストデータの統計的な分析に基づく
 - 形態素解析（単語分割）
 - 構文解析・係り受け解析
 - 統計的意味解析
 - 文書の統計モデルと情報検索 etc, etc ...

彼女は花を買った。

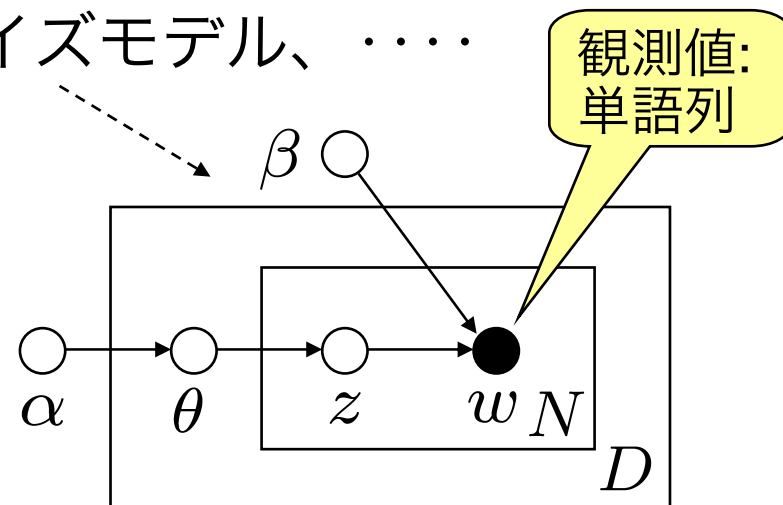


統計的自然言語処理

- 1990年代後半～からパラダイムシフト
 - 統計的機械学習の一部として重要な位置
- 論理式から、高度な統計モデルへ
 - チョムスキーの亡霊からの脱却
 - Webの登場と電子テキスト、計算資源の爆発的増大
 - 対数線形モデル、階層ベイズモデル、……

$$p(t|\mathbf{x}, \Lambda) = \frac{\exp(-\sum_i \lambda_i f_i(\mathbf{x}, t))}{\sum_{\mathbf{x}} \exp(-\sum_i \lambda_i f_i(\mathbf{x}, t))}$$

ある単語 \mathbf{x} の品詞
が形容詞である確率



形態素解析

- 日本語や中国語等は単語に分けられていない
……自然言語処理の非常に重要な課題

```
% echo “やあこんにちは, 同志社内はどうですか。”  
| mecab -O wakati  
やあ こんにちは, 同志社内はどうですか。  
(やあこんにちは, 同志社内はどうですか。×)
```

- Chasen, MeCab (NAIST)などが有名なツール
- これまで、教師あり学習 (supervised learning) によって学習されてきた
 - 人手で、単語分割の「正解例」を何万文も作成
 - 膨大な人手と手間のかかるデータ作成

形態素解析 (2)

S-ID:950117245-006 KNP:99/12/27

* 0 5D

一方 いっぽう * 接続詞 * * *

、 * 特殊 読点 * *

* 1 5D

震度 しんど * 名詞 普通名詞 * *

は は * 助詞 副助詞 * *

* 2 3D

揺れ ゆれ * 名詞 普通名詞 * *

の の * 助詞 接続助詞 * *

* 3 4D

強弱 きょうじゃく * 名詞 普通名詞 * *

毎日新聞
1995年度記事
から38,400文
(京大コーパス)
の例

- 膨大な人手で作成した教師(正解)データ
 - 対数線形モデルやその拡張を用いて識別器を学習
- 話し言葉の「正解」？ 古文？ 未知の言語？
 - |女御|更衣|あ|また|さ|ぶら|ひ|た|ま|ひける|中|に|、|...

教師なし形態素解析

- 確率モデルに基づくアプローチ: 文字列 s について、それを分割した単語列 $p(\mathbf{w}|s)$ の確率

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmax}} p(\mathbf{w}|s)$$

を最大にする $\hat{\mathbf{w}}$ を探す

- 例: $p(\text{今日はもう見た}) > p(\text{今日はもう見た})$
 - 教師データを使わない; 辞書を使わない
 - 「言語として最も自然な分割」を学習する
- あらゆる単語分割の可能性を考える
 - たった50文字の文でも、
 $2^{50} = 1,125,899,906,842,624$ 通りの天文学的組み合わせ
(さらに無数の文が存在)

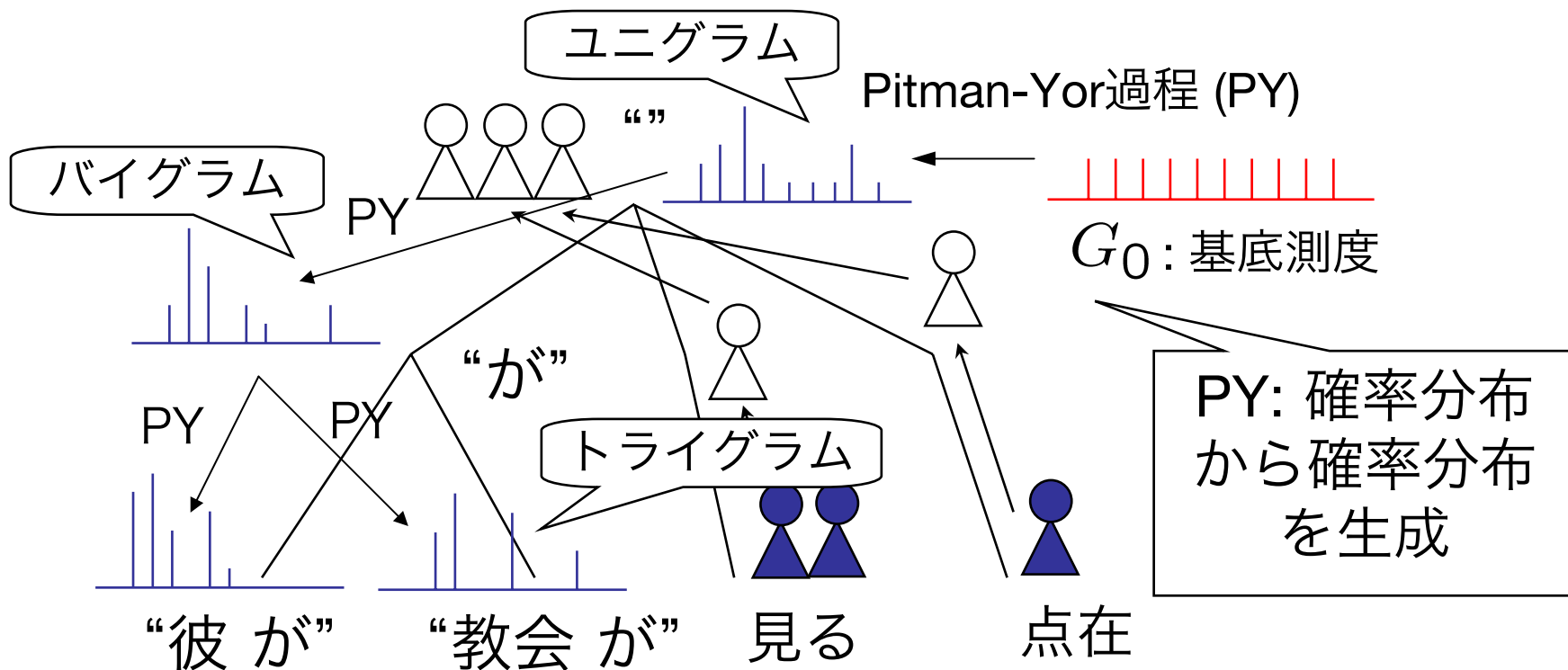
文の確率: nグラムモデル

$$p(\text{今日はもう見た}) \\ = p(\text{今日}|\wedge) \cdot p(\text{は}|\text{今日}) \cdot p(\text{もう}|\text{は}) \cdot p(\text{見た}|\text{もう})$$

文頭を表す特殊文字

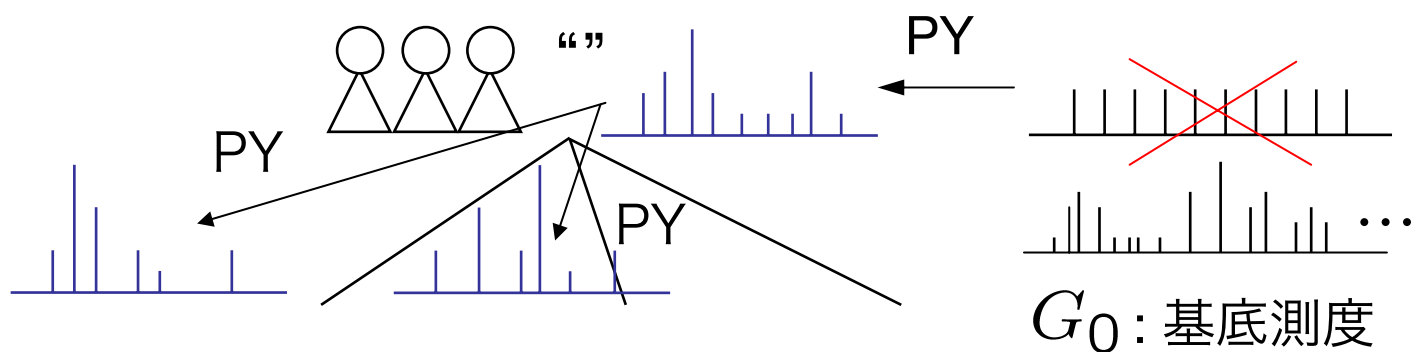
- 条件付き確率の積で文の確率を計算
 - 自然言語処理では、きわめて強力 (Shannon 1948)
 - 確率のテーブルは、ほとんどが0
 - 階層的なスムージングが不可欠
 - あらゆる部分文字列が「単語」になりうる
- ➡ 階層ベイズモデル: 階層Pitman-Yor過程言語モデル (HPYLM) (Teh 2006; Goldwater+ 2005)
- Pitman-Yor過程: ディリクレ過程 (GEM分布) の一般化

準備: HPYLM n-gram



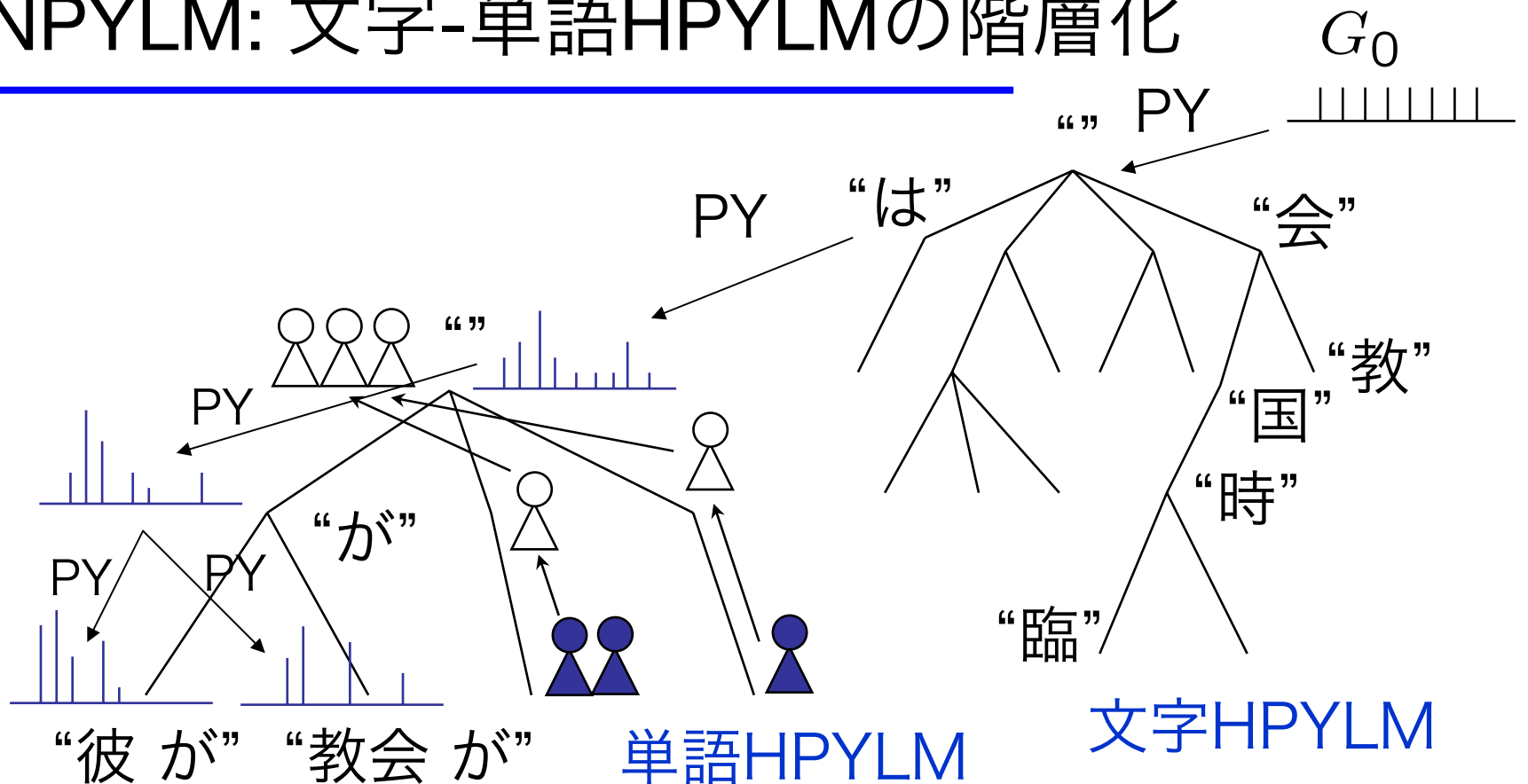
- カウントが0でも、より低いオーダーのMarkovモデルを用いて階層ベイズでスムージング
 - 注目している単語がそもそも存在しなかったら？

HPYLM: 無限語彙モデル



- 基底測度 G_0 は、単語の事前確率を表す
 - 語彙 V が有限なら、 $G_0(w \in V) = 1/|V|$
- G_0 は可算無限でもよい！ → 無限語彙
 - PYに従って、必要に応じて「単語」が生成される
 - 「単語」の確率は、文字n-gram=もう一つのHPYLM
 - 他の方法で与えてもよい (が、再学習が面倒)

NPYLM: 文字-単語HPYLMの階層化



- HPYLM-HPYLMの埋め込み言語モデル
 - つまり、階層Markovモデル
- 文字HPYLMの G_0 は, 文字数分の1 (日本語なら1/6879)

NPYLMの学習問題の定式化

- データ: $\mathbf{X} = \{s_1, s_2, \dots, s_X\}$ (文の集合)
 - 文: $s = c_1 c_2 \dots c_N$ (文字列)
 - 隠れ変数: $\mathbf{z} = z_1 z_2 \dots z_N$ ($z_i = 1$ のとき単語境界)
 - 隠れ変数の組み合わせは指数的に爆発

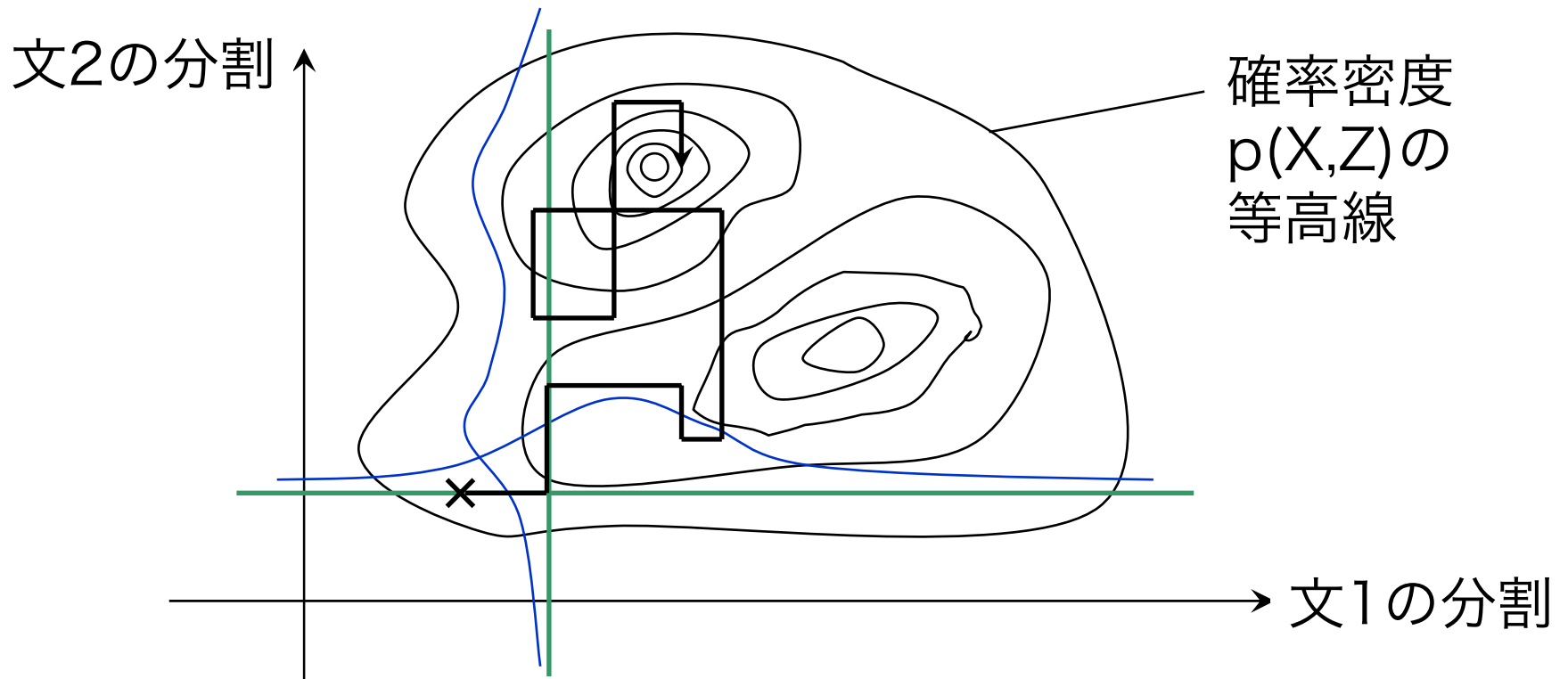
- 文がそれぞれ独立だと仮定すると、

$$p(\mathbf{X}) = \prod_{n=1}^X p(s_n) \quad (1)$$

$$p(s_n) = \sum_{\mathbf{z}_n} p(s_n, \mathbf{z}_n) \quad (2)$$

- 各文 s_n の分割 \mathbf{z}_n を、どうやって推定するか?
→ ブロック化ギブスサンプリング、MCMC.

Blocked Gibbs Sampling



- 確率 $p(X,Z)$ を最大にする単語分割を求める
- 単語境界は、前後の「単語」に強い依存関係
→ 文ごとに、可能な単語分割をまとめてサンプル (Blocked Gibbs sampler)

Blocked Gibbs Sampler for NPYLM

- 各文の単語分割を確率的にサンプリング
→ 言語モデル更新
→ 別の文をサンプリング
...を繰り返す.

- アルゴリズム:

0. For $s = s_1 \dots s_X$ do

$\text{parse_trivial}(s, \Theta)$.

← 文字列全体が一つの「単語」

1. For $j = 1 \dots M$ do

 For $s = \text{randperm}(s_1 \dots s_X)$ do

 言語モデルから $\text{words}(s)$ を削除

$\text{words}(s) \sim p(w|s, \Theta)$ をサンプリング

 言語モデルに $\text{words}(s)$ を追加して更新

← Θ : 言語モデルのパラメータ

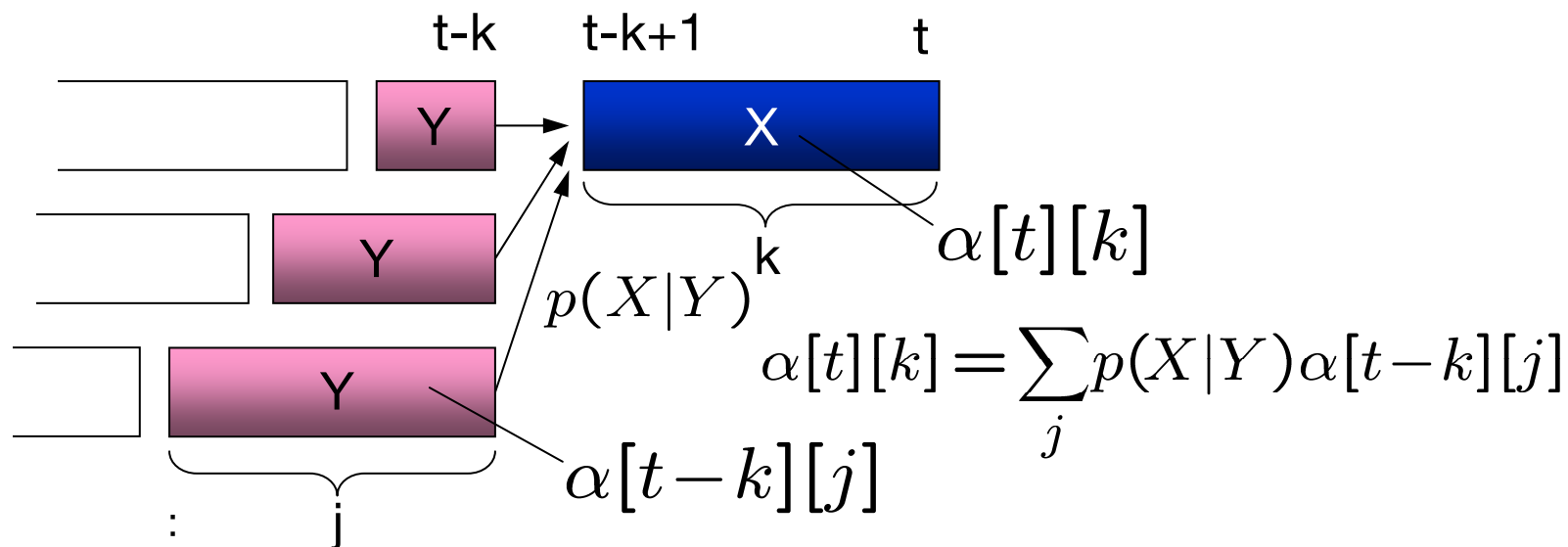
done.

Gibbs Samplingと単語分割

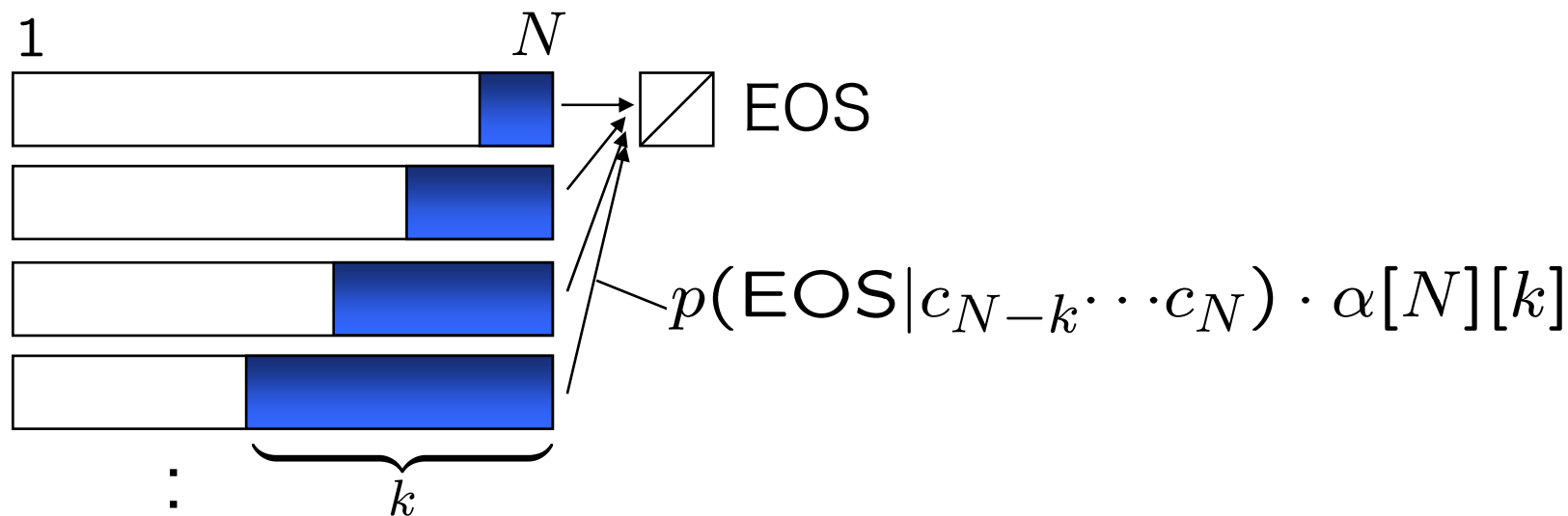
- 1 神戸では異人館 街の 二十棟 が破損した。
 - 2 神戸 では 異人館 街の 二十棟 が破損した。
 - 10 神戸 では 異人館 街の 二十棟 が破損した。
 - 50 神戸 では異人 館 街 の 二十棟 が破損した。
 - 100 神戸 では 異 人館 街 の 二十棟 が破損した。
 - 200 神戸 では 異人館 街 の 二十棟 が破損した。
- ギブスサンプリングを繰り返すごとに、単語分割とそれに基づく言語モデルを交互に改善していく。

動的計画法による推論

- $\text{words}(s) \sim p(w|s, \Theta)$: 文 s の単語分割のサンプリング
- 確率的Forward-Backward (Viterbiだとすぐ局所解)
 - Forwardテーブル $\alpha[t][k]$ を用いる
 - $\alpha[t][k]$: 文字列 $c_1 c_2 \dots c_t$ が、時刻 t から k 文字前までを単語として生成された確率
 - それ以前の分割について周辺化...動的計画法で再帰

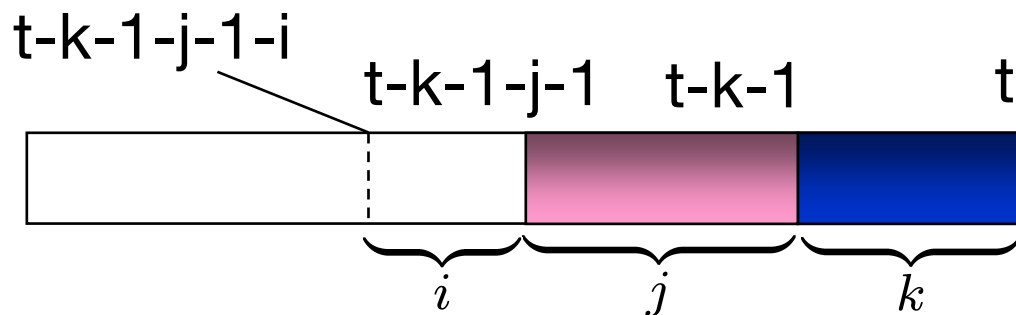


動的計画法によるデコード



- $\alpha[N][k]$ = 文字列の最後の k 文字が単語となる文字列確率なので、EOS に接続する確率に従って後ろから k をサンプル
- $c_{N-k} \dots c_N$ が最後の単語だとわかったので、 $\alpha[N-k-1][k']$ を使ってもう一つ前の単語をサンプル
- 以下文頭まで繰り返す

動的計画法による推論 (トライグラムの場合)



- トライグラムの場合は、Forward 変数として $\alpha[t][k][j]$ を用いる
 - $\alpha[t][k][j]$: 時刻 t までの文字列の k 文字前までが単語、さらにその j 文字前までが単語である確率
 - 動的計画法により、 $\alpha[t-k-1][j][i]$ ($i = 0 \dots L$) を使って再帰
 - プログラミングが超絶ややこしい ;_;
 - (文字列は有限なので前が存在しないことがある)

実験: 日本語 & 中国語コーパス

- 京大コーパス & SIGHAN Bakeoff 2005 中国語単語分割公開データセット
- 京大コーパスバージョン4
 - 学習: 37,400文、評価: 1000文(ランダムに選択)
- 日本語話し言葉コーパス: 国立国語研究所
- 中国語
 - 簡体中国語: MSRセット, 繁体中国語: CITYUセット
 - 学習: ランダム50,000文、評価: 同梱テストセット
- 学習データをそれぞれ2倍にした場合も同時に実験

京大コーパスの教師なし形態素解析結果

一方、村山富市首相の周囲にも韓国の状況や立場を知る高官はいない。

日産自動車は、小型乗用車「ブルーバード」の新モデル・S Vシリーズ5車種を12日から発売した。

季刊誌で、今月三十日発行の第一号は「車いすテニス新世代チャンピオン誕生－斎田悟司 ジャパンカップ 松本、平和カップ 広島連覇」「フェスピック北京大会－日本健闘メダル獲得総数88個」「ジャパンパラリンピック－日本の頂点を目指す熱い闘い」などの内容。

整備新幹線へ投入する予算があるのなら、在来線を改良するなどして、高速化を推進し輸送力増強を図ればよい。

国連による対イラク制裁解除に向け、関係の深い仏に一層の協力を求めるのが狙いとみられる。

この日、検査されたのはワシントン州から輸出された「レッドデリシャス」、五二トン。

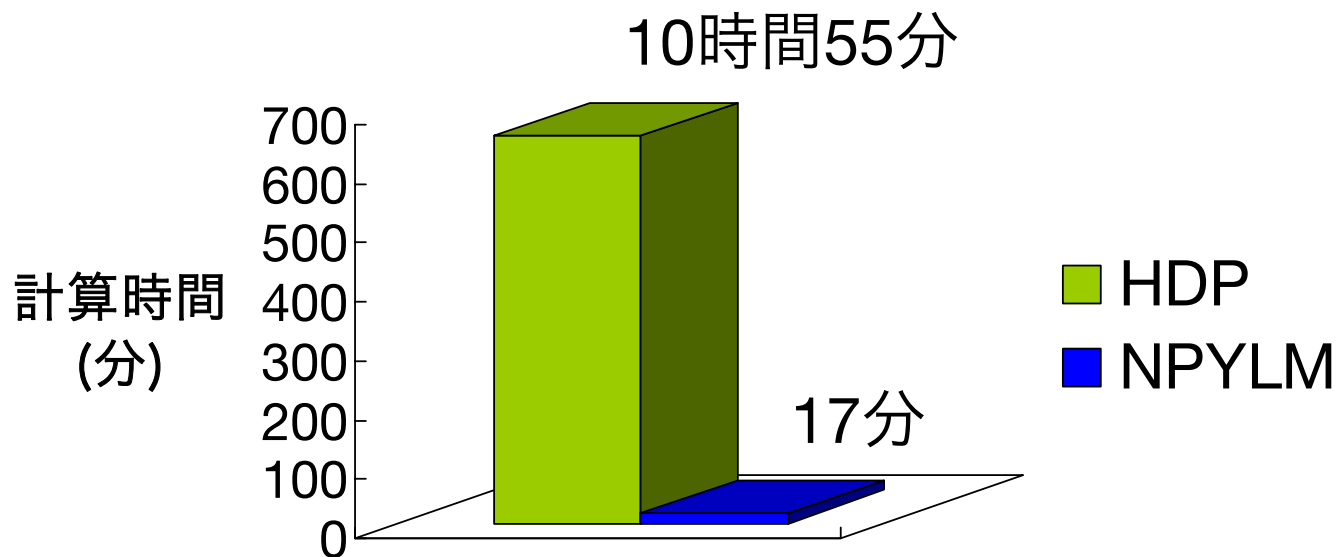
ビタビアルゴリズムで効率的に計算可能
(先行研究では不可能)

“正解”との一致率 (F値)

モデル	MSR	CITYU	京大
NPY(2)	0.802 (51.9)	0.824 (126.5)	0.621 (23.1)
NPY(3)	0.807 (48.8)	0.817 (128.3)	0.666 (20.6)
NPY(+)	0.804 (38.8)	0.823 (126.0)	0.682 (19.1)
ZK08	0.667 (—)	0.692 (—)	—

- NPY(2), NPY(3) = NPYLM 単語バイグラム or トライグラム + 文字 ∞ グラム
 - NPY(+)はNPY(3)でデータを2倍にしたもの
- 中国語: ZK08 = (Zhao&Kit 2008)での最高値と比べ、大きく改善
 - ZK08はヒューリスティックな手法をさらに混合したもの

計算時間の比較



- HDP(Goldwater+ ACL 2006): 学習データのすべての文字について1文字ずつサンプリング
 - モデルは単語2グラムのみ (文字モデルなし)
- NPYLM: 文毎に動的計画法により効率的にサンプリング
 - 単語3グラム-文字 ∞ グラムの階層ベイズモデル

日本語話し言葉コーパス (国立国語研究所)

うーんうんになってしまおうところでしょうねへーあーでもいいいいこと
ですよねうーん

うーん自分にも凄くプラスになりますものねそうですねふーん羨ましい
です何かうーん精神的にもう子供達に何かこう支えられるようないーも
のってやっぱりあるんですよやっているとうーんうーんうーん

うーん長くやってればそんなものがうんうんそうでしょうねたくさんやっ
ぱりありますねうんうーんなるほど…



うーん うん になってしまおう ところ でしょうね へー あー でも いい いい
こと ですよねうーん

うーん 自分 にも 凄く プラス になります ものね そう ですね ふーん
羨ましい です 何か うーん 精神的 にもう 子供達 に何か こう 支えられる
ようないー もの って やっぱり ある んですよ やっていると うーん

うーん うーん うーん 長く やって れば そんな ものが うん うん そう
でしょうね たくさん やっぱり あります ね うん うーん なる ほど…

「源氏物語」の教師なし形態素解析

しばしは夢かとのみたどられしを、やうやう思ひしづまるにしも、さむべき方なくたへがたきは、いかにすべきわざにかとも、問ひあはすべき人だになきを、忍びては参りたまひなんや。若宮の、いとおぼつかなく、露けき中に過ぐしたまふも、心苦しう思さるるを、とく参りたまへ』など、はかばかしうも、のたまはせやらず、むせかへらせたまひつつ、かつは人も心弱く見たてまつるらむと、思しつつまぬにしもあらぬ御気色の……



しばしは夢かとのみたどられしを、やうやう思ひしづまるにしも、さむべき方なくたへがたきは、いかにすべきわざにかとも、問ひあはすべき人だになきを、忍びては参りたまひなんや。若宮の、いとおぼつかなく、露けき中に過ぐしたまふも、心苦しう思さるるを、とく参りたまへ』など、はかばかしうも、のたまはせやらず、むせかへらせたまひつつ、かつは人も心弱く見たてまつるらむと、思しつつまぬにしもあらぬ御気色の……

アラビア語教師なし形態素解析

- Arabic Gigawords から40,000文 (Arabic AFP news)

الفلستيني بسبب تظاهرة لانصار حركة المقاومة الاسلامية حماس
و اذا تحقق ذلك فان كيسلو فسكيه قد حاز ثلاثه جري في ابرز ثلاثة

صحية
+ قائد
الا يقل

Google translate:

“Filstinebsbptazahrplansarhrkpalmquaompalaslami
phamas.”

وقالت دانيل تومسون التي كتبت السيناريو. وقد استغرق اعداد خمسة اعوام. "تاريخي

↓ NPYLM

الفلستيني بسبب تظاهرة لانصار حركة المقاومة الاسلامية حماس
و اذا تحقق ذلك ف ان كيسلو فسكي يكون قد حاز ثلاثه جري في ابرز ثلاثة

صحية
سطينية
مالا يقل

Google translate:

“Palestinian supporters of the event because of
the Islamic Resistance Movement, Hamas.”

وقد استغرق اعداد ه خمسة اعوام . و قال ت دانيل تومسون التي " تاريخي

“Alice in Wonderland”の解析

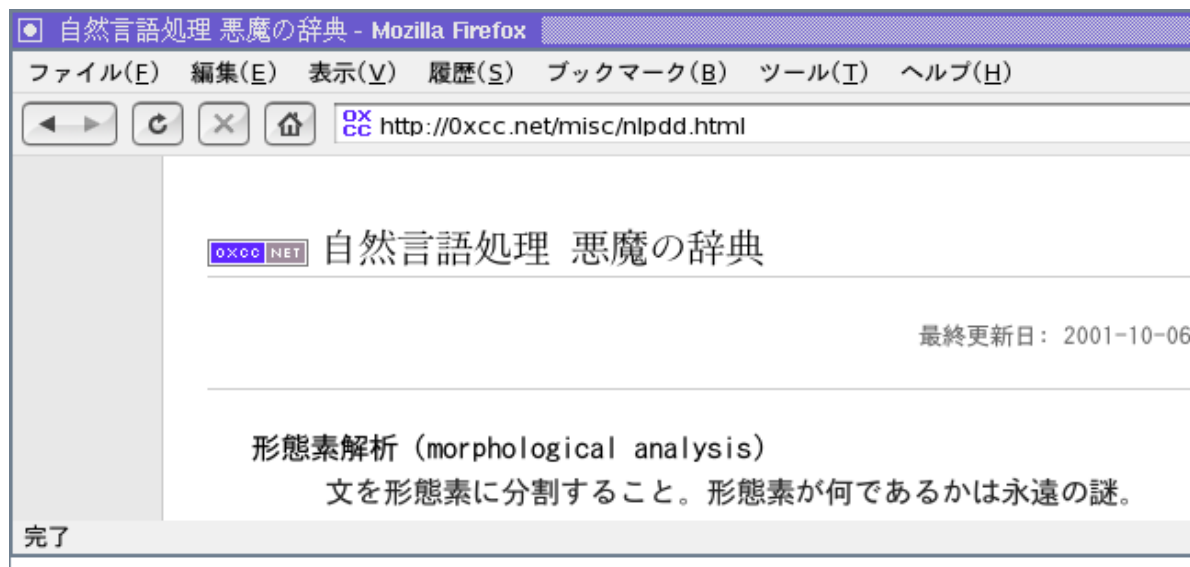


first, she dream ed of little alic e herself , and once again the tiny hand s were clasped up on her knee , and the bright eager eyes were looking up into hers -- she could hear the very tone s of her voice , and see that queer little toss of her head to keep back the wandering hair that would always get into hereyes -- and still as she listened , or seemed to listen , the whole place a round her became alive the strange creatures of her little sister 's dream. the long grass rustled at her feet as the white rabbit hurried by -- the frightened mouse splashed his way through the neighbour ing pool -- she could hear the rattle of the tea cups as the march hare and his friends shared their never -ending meal , and the shrill voice of the queen...



first, she dream ed of little alic e herself , and once again the tiny hand s were clasped up on her knee , and the bright eager eyes were looking up into hers -- she could hear the very tone s of her voice , and see that queer little toss of her head to keep back the wandering hair that would always get into hereyes -- and still as she listened , or seemed to listen , the whole place a round her became alive the strange creatures of her little sister 's dream. the long grass rustled at her feet as the white rabbit hurried by -- the frightened mouse splashed his way through the neighbour ing pool -- she could hear the rattle of the tea cups as the march hare and his friends shared their never -ending meal , and the ...

“形態素”の再定義



- “自然言語処理 悪魔の辞典”: 高林哲氏
 - 「形態素が何であるかは永遠の謎」

教師あり学習では、
確かに謎



- 今や謎ではない！
 - “形態素”とは、文字列の生成確率を最大にするような統計的な単位として導くことができる。

まとめ

- ベイズ単語nグラム-文字nグラムを階層的に統合した言語モデルによる、教師なし形態素解析
 - 動的計画法+MCMCによる効率的な学習
- あらゆる自然言語に適用できる
 - データに自動的に適応、「未知語」問題がない
 - 識別学習と違い、学習データをいくらでも増やせる
 - 話し言葉、ブログ、未知の言語、古文、...
- あらゆる言語の文字列から直接、「単語」を推定しながら言葉のモデルを学習する方法ともみなせる

実装

- 数万～数十万文 (数百万～数千万文字)の学習テキストに対してGibbsサンプリングを繰り返すため、高速な実装が不可欠
 - MATLABやRでは計算が追いつかない
- C++&Cで実装, 6000行程度
 - 解析速度は100～200文/秒 (10ms/文以下)
 - 1つの文を解析するのに、nグラム確率を40000回程度計算する必要
 - 階層的データ構造の動的なアップデート
 - 学習時間: 10～20時間程度

おわり

ご清聴ありがとうございました。

展望

- 教師あり学習と異なり、学習データをいくらでも増やせる → 学習の高速化、並列化
 - HDP-LDAのGibbsの並列化 (Welling+, NIPS 2007-2008) が適用可能
- 識別学習との融合による半教師あり学習
 - Loglinearの枠組で統合するにも、生成モデルが必要
 - これまで、生成モデルが存在しなかった
 - 提案法は、CRFのForward-Backwardの教師なし版のようなもの
 - POS Tagging: CRF+HMM (鈴木,藤野+ 2007)で提案