
ノンパラメトリックベイズ法による 言語モデル

持橋大地

統計数理研究所 モデリング研究系

daichi@ism.ac.jp

2012-3-15 (木), 統数研

Overview

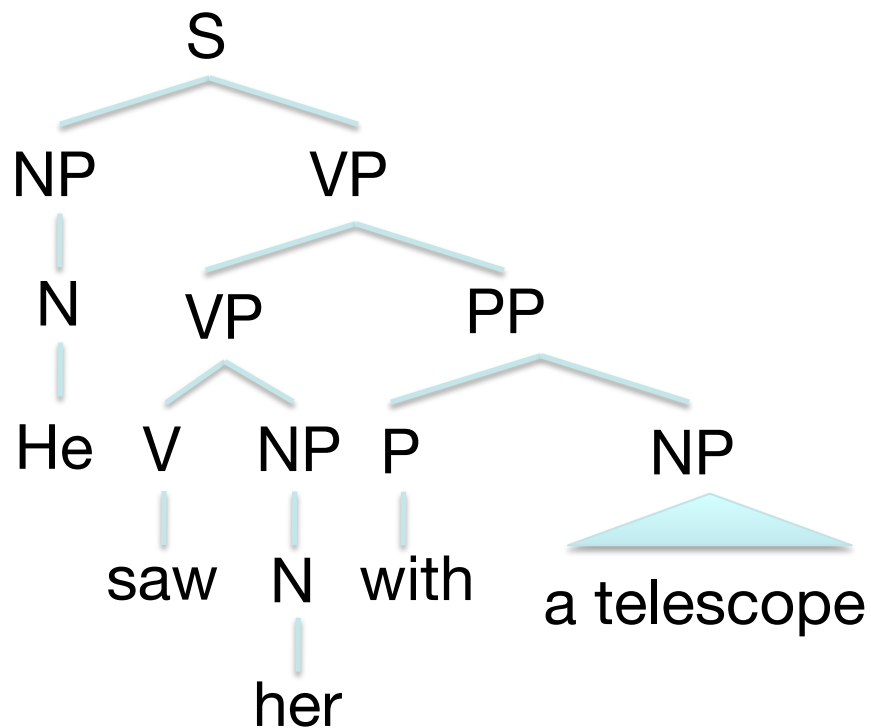
- 統計的自然言語処理・言語モデルとは
- ノンパラメトリック・ベイズ法とは (イントロ)
 - Dirichlet分布
 - Dirichlet過程
 - 階層Dirichlet過程
 - Chinese Restaurant Process (CRP)
- 階層Pitman-Yor過程に基づくnグラム言語モデル
- 階層言語モデルによる教師なし形態素解析

言語学と統計的自然言語処理

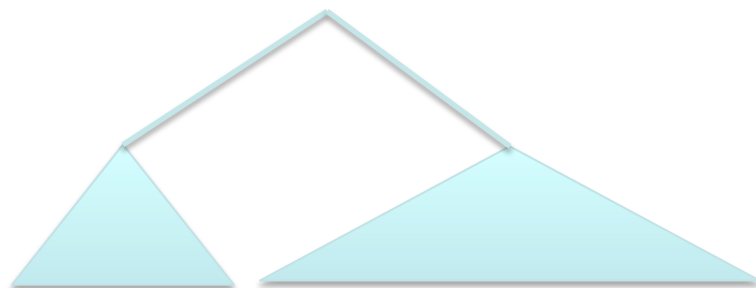
- 言語の研究 → 言語学科に行けばよいか？
- いわゆる「言語学」 = 手で書いたルールの固まり！
- 例：構文解析

S → NP VP
VP → V PP
VP → V NP
PP → P NP
NP → DET N
NP → N

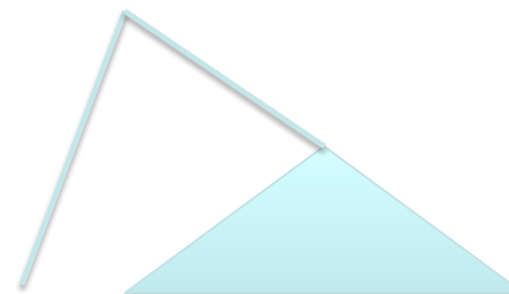
文法ルール



言語学と統計的自然言語処理 (2)



He saw her with a **telescope**

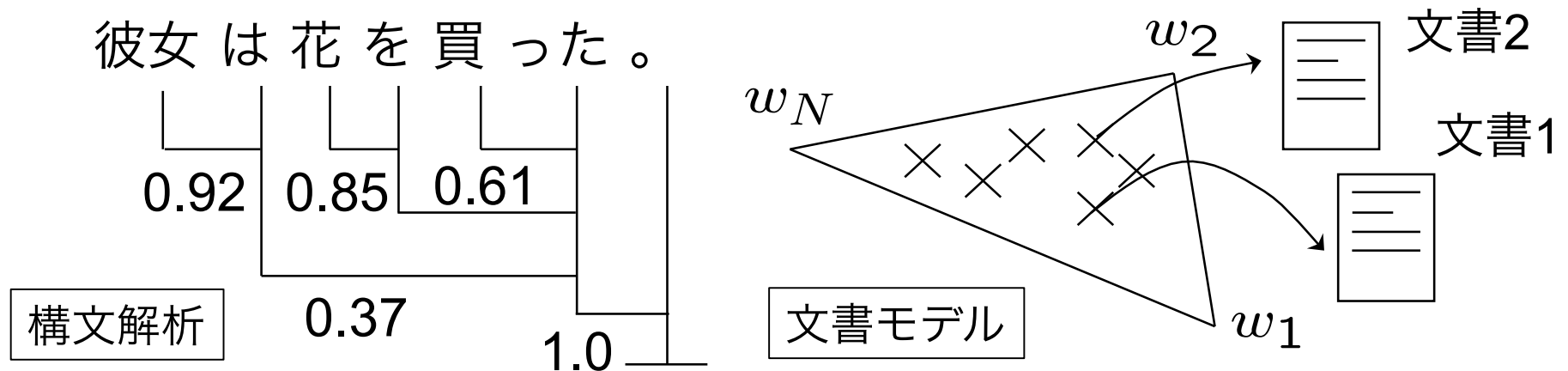


He saw her with a **hat**

- 解釈が名詞によってなぜ違うのか？
- 古典的な言語学では答えが出せない・
そもそも主観的
→ 確率モデル・統計学として数学的に
考え直す必要！
 - cf. 中世の天動説から地動説の数学理論へ

統計的自然言語処理

- 1990年代後半～: 大量の言語データから、言語の性質を統計的に学習
 - Webの出現、大量の電子テキスト
- 代表的な応用:
構文解析、形態素解析、文書モデル、意味極性分類、照応解析、言語進化モデル、……



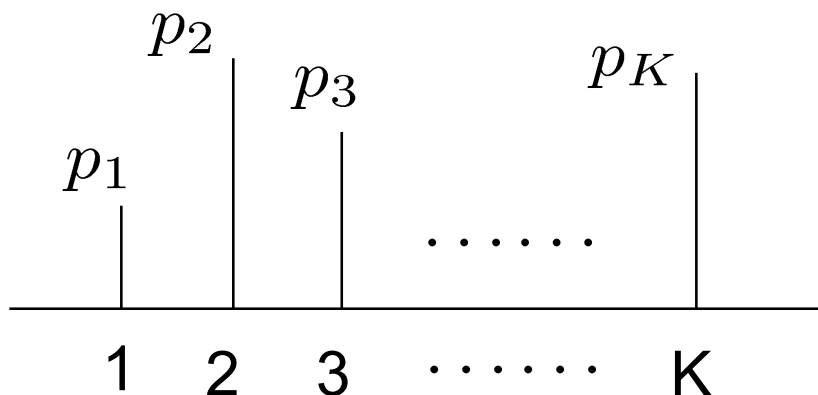
統計的言語モデル

- 統計的自然言語処理の最も基本的なモデル
- 単語列 $\mathbf{w} = w_1 w_2 \cdots w_n$ に対し、その確率 $p(\mathbf{w})$ 最大にする確率モデルを学習
 - 木構造やMarkovモデルなど
 - 情報理論と密接な関係 (良いモデル \equiv 良い符号化)
- 隠れ変数 \mathbf{z} があってよい

$$p(\mathbf{w}) = \sum_{\mathbf{z}} p(\mathbf{w}, \mathbf{z})$$

- \mathbf{z} は何でもよい!!
- 構文木、品詞列、感情極性、意味トピック、単語分割、etc,etc...
- 自然言語処理のほとんどすべての問題を含む

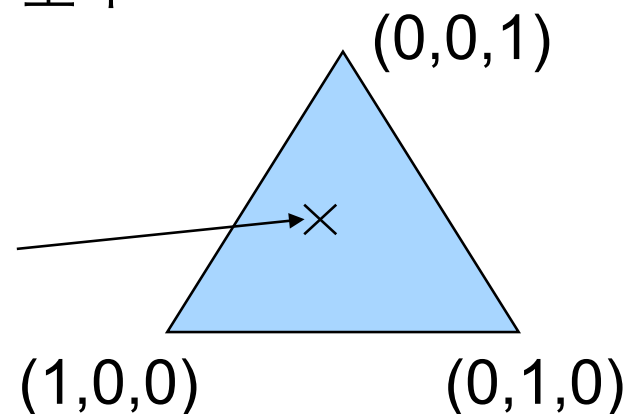
多項分布 (離散分布)



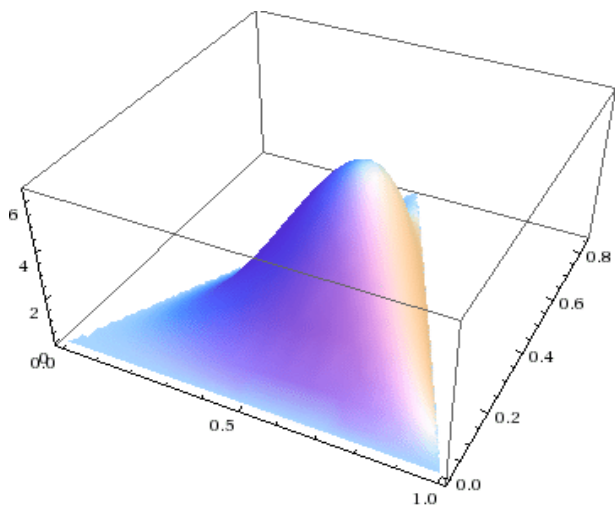
$$\sum_{k=1}^K p_k = 1,$$
$$\forall k, p_k \geq 0$$

- K種類のアイテムのどれかが出る確率分布
 - 離散データの統計モデルの基本中の基本
- \mathbf{p} は $(K-1)$ 次元の単体(Simplex)の内部に存在

$$\mathbf{p} = (p_1, p_2, \dots, p_K)$$



ディリクレ分布



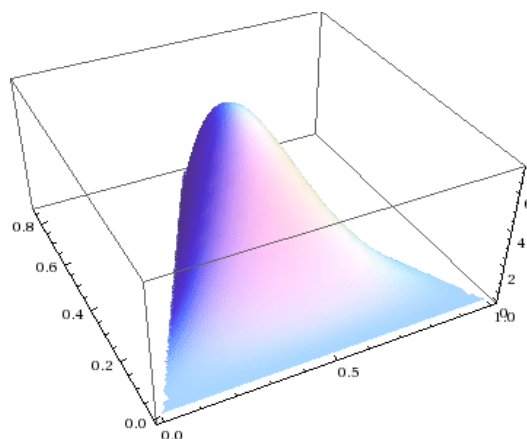
$$\text{Dir}(\mathbf{p}|\boldsymbol{\alpha}) \propto \prod_{k=1}^K p_k^{\alpha_k - 1}$$

パラメータ:

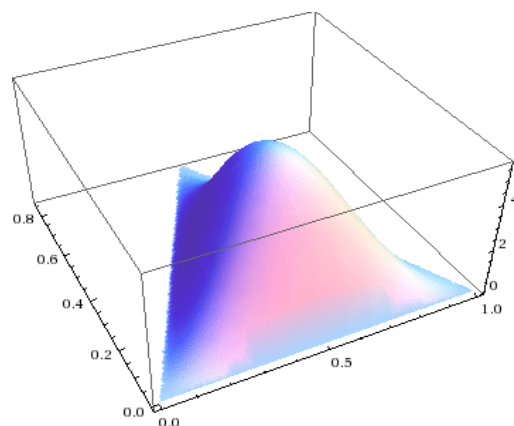
$$\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_K)$$

- ランダムな多項分布を生成する確率分布
- $\alpha_k \equiv 1$ のとき、単体上でUniformな分布
- 「期待値」 : $\bar{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_K) / \sum_k \alpha_k$
- 「分散」 : $\alpha = \sum_k \alpha_k$

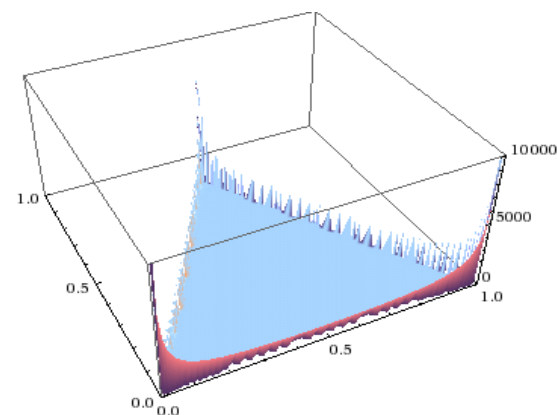
ディリクレ分布 (2)



$$\alpha = (2, 4, 2)$$



$$\alpha = (2, 2, 2)$$



$$\alpha = (0.5, 0.5, 0.5)$$

- $\alpha_k > 1$ のとき、上に凸
- $\alpha_k < 1$ のとき、下に凸
 - 統計的自然言語処理等では、多くの場合 $\alpha \ll 1$
($\alpha = 0.1 \sim 0.0001$ くらい)

ディリクレ分布に基づく予測

- ゆがんだ三面サイコロを振ったら、結果は $X = \{1, 2, 2, 3, 2, 3\}$ (1=1回, 2=3回, 3=2回) だった。次の目は?

- ベイズの定理: $p(\mathbf{p}|X) \propto p(X|\mathbf{p})p(\mathbf{p})$

$$\propto (p_1^1 \cdot p_2^3 \cdot p_3^2) \cdot \left(\prod_{k=1}^3 p_k^{\alpha_k - 1} \right)$$

$$= p_1^{\alpha_1 + 1 - 1} \cdot p_2^{\alpha_2 + 3 - 1} \cdot p_3^{\alpha_3 + 2 - 1}$$

$$= \text{Dir}(\alpha_1 + 1, \alpha_2 + 3, \alpha_3 + 2)$$

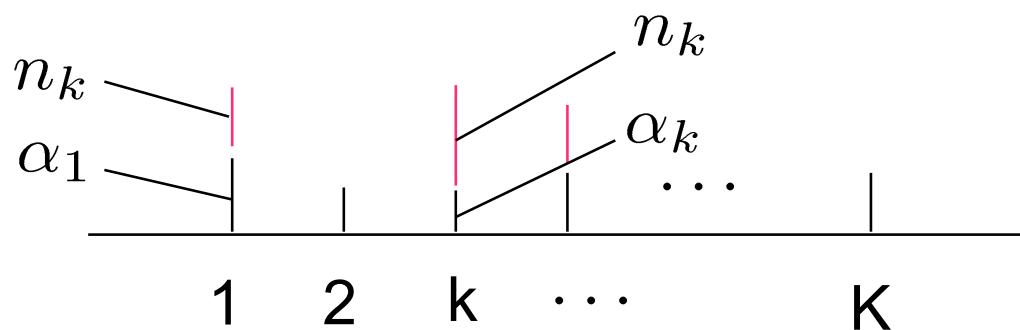
- \mathbf{p} の期待値は、

$$E[\mathbf{p}|X] = \left(\frac{\alpha_1 + 1}{\alpha + 6}, \frac{\alpha_2 + 3}{\alpha + 6}, \frac{\alpha_3 + 2}{\alpha + 6} \right) \quad (\alpha = \sum_k \alpha_k)$$

ディリクレ分布に基づく予測 (2)

- 一般に、 n 回の観測の中で k 番目のアイテムが出現したとすると、 $\boxed{n_k}$

$$p(k|X) = \frac{\alpha_k + n_k}{\alpha + n}.$$



注: $n_k = 0$ のとき、 $p(k|X) = \frac{\alpha_k}{\alpha + n}$

(出現しなかったアイテムにも正の確率)

ノンパラメトリック・ベイズ法とは

- モデルの複雑さを、**データの複雑さに応じて無限に伸縮**することのできるベイズ統計モデル
 - 「パラメータがない」という意味ではない
- 簡単な場合の例
 - GMMの混合数
 - HMMの隠れ状態数
 - 文書に存在する意味トピック数
 - 言語の文法的ルールの複雑さ
- 有名なもの: Dirichlet過程 (無限次元Dirichlet分布)

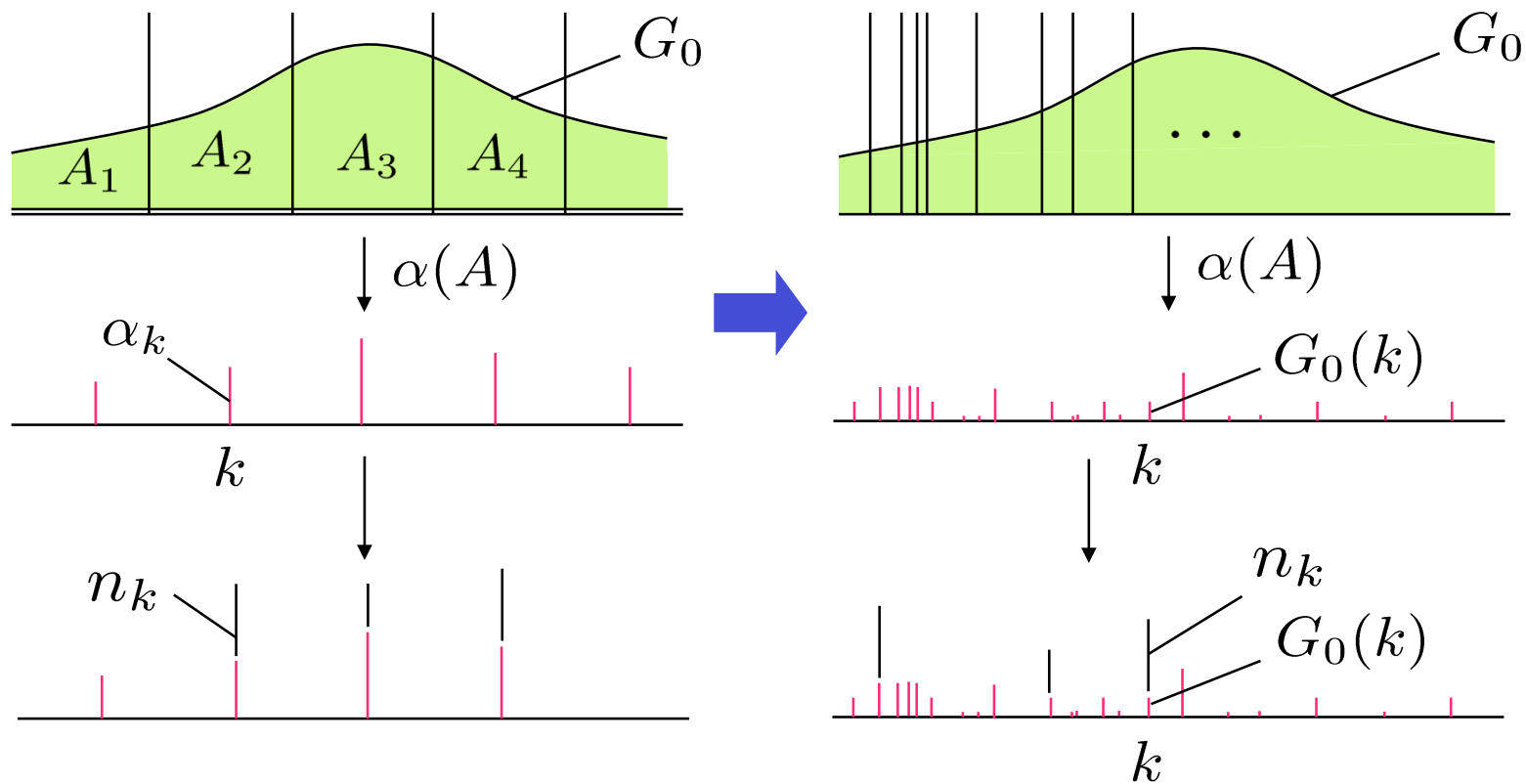
ディリクレ過程

- Dirichlet processとは要するに何?
→ 無限次元ディリクレ分布.
- DPの定義 (Ferguson 1973):

A stochastic process P is said to be a Dirichlet process on $(\mathcal{X}, \mathcal{A})$ with parameter $\bar{\alpha}$ for any measurable partition (A_1, \dots, A_k) of \mathcal{X} the random vector $(P(A_1), \dots, P(A_k))$ has a Dirichlet distribution with parameter $(\alpha(A_1), \dots, \alpha(A_k))$

- どういうこと??

ディリクレ過程 (2)



予測確率:

$$\frac{\alpha_k + n_k}{\alpha + n}$$

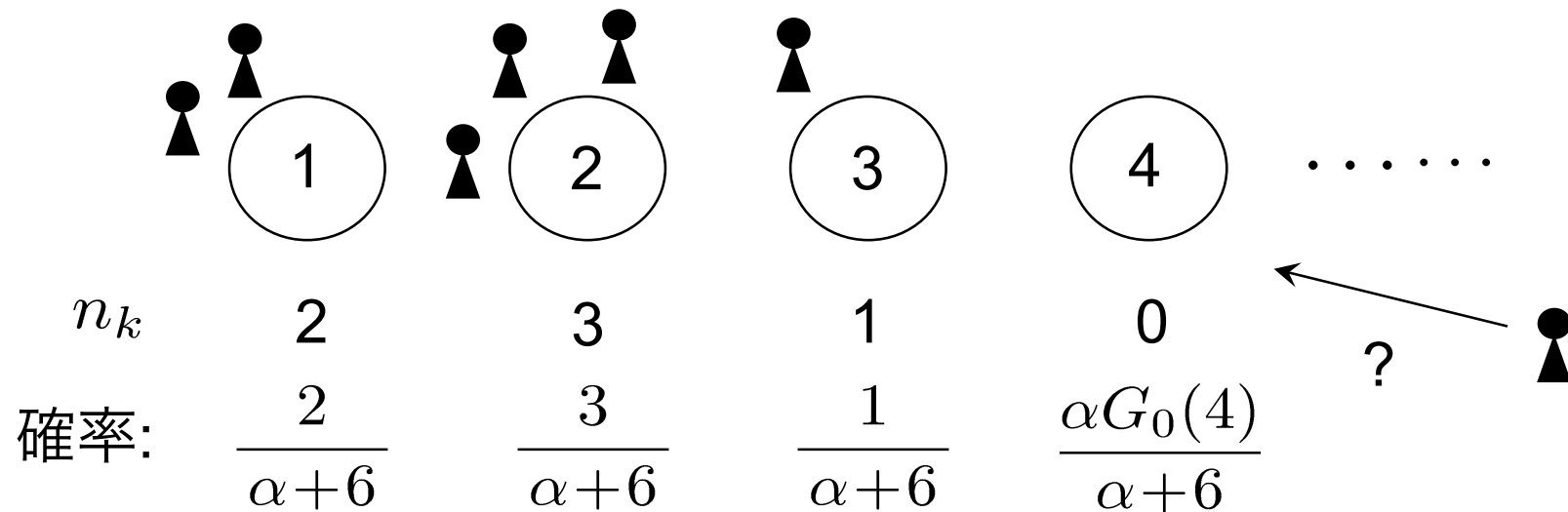
$$\frac{\alpha G_0(k) + n_k}{n + \alpha}$$

Chinese Restaurant Process (CRP)

- 予測確率

$$p(k|X) = \frac{\alpha_k + n_k}{\alpha + n} \text{ (Dirichlet), } \frac{\alpha G_0(k) + n_k}{\alpha + n} \text{ (DP)}$$

- ディリクレ分布/過程に従うと、頻度 n_k 高いものはさらに現れやすくなる (rich-gets-richer) → CRP



ディリクレ過程と言語モデル

- ディリクレ過程は、語彙が無限の場合の単語の確率分布ともみることができる

$$p(w|X) = \frac{c(w)}{\alpha + n} + \frac{\alpha}{\alpha + n} G_0(w)$$

- カウント $c(w)$ が 0 のどんな未知の単語 w でも、 $G_0(w) \cdot \alpha / (\alpha + n)$ の確率を持つ
- この確率分布は、 $p(w) \cdots$ ユニグラムモデル
 - 単語が独立に出現すると仮定している
 - 一般には、前の単語などに強く依存
 - “is going” → to, “united states of” → america など

nグラムモデルのベイズ学習

- nグラムモデル・・・古典的だが、音声認識や機械翻訳では未だ重要、基本的 (言葉のMarkovモデル)

$$p(\text{彼女が見る夢}) = p(\text{彼女}) \cdot p(\text{が} | \text{彼女}) \cdot p(\text{見る} | \text{が}) \cdot p(\text{夢} | \text{見る})$$

- nグラムモデルの問題: スムージング

$$\hat{p}(\text{yield} | \text{maximum likelihood will often}) = \frac{n(\text{maximum likelihood will often yield})}{n(\text{maximum likelihood will often})} = 0$$

現在のGoogle
カウント

- 頻度そのままではなく、何か値を足したりする必要!

ディリクレスムージング (MacKay 1994)

- nグラム確率分布 $p(w|h)$ はディリクレ事前分布 $p(\cdot|h) \sim \text{Dir}(\alpha)$ を仮定すると、結果は α_w を足すのと同じ

$$p(w|h) = \frac{n(w|h) + \alpha_w}{\sum_w (n(w|h) + \alpha_w)}$$

- α はバイグラムなら、Newton法で最適化できる

- 問題: 性能が意外と低い

- カウント $n(w|h)$ が0のとき、

$$p(w|h) = \frac{\alpha_w}{\sum_w (n(w|h) + \alpha_w)}$$

- $\alpha_w \ll 1$ なので、 $1 + \alpha_w \iff \alpha_w$ に物凄い差!!

大体0.1~0.001くらい

Kneser-Ney スムージング (Kneser, Ney 1995)

- 最高精度とよばれるスムージング法

$$p(w|h) = \frac{n(w|h) - D(n(w|h))}{n(h)} + \gamma(h)p(w|h')$$

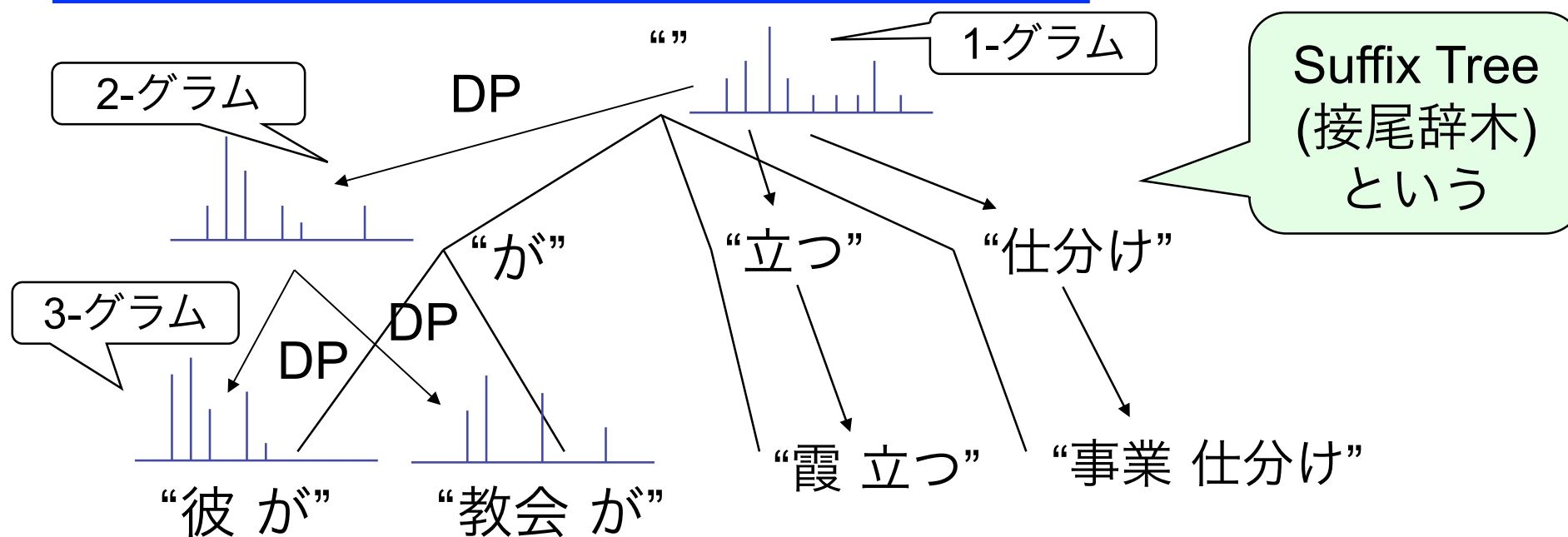
- 頻度 $n(w|h)$ から、一定数 D をディスカウント
- $\gamma(h)$ は h に後続する単語の種類数から決まる

- これは、下の階層 Pitman-Yor 過程による予測の近似であることが最近判明 (Goldwater+ 2006, Teh 2006)

$$p(w|h) = \frac{n(w|h) - d \cdot t_{hw}}{\theta + n(h)} + \frac{\theta + d \cdot t_h}{\theta + n(h)} p(w|h')$$

- 階層 Pitman-Yor 過程とは?

階層ディリクレ過程 (HDP)

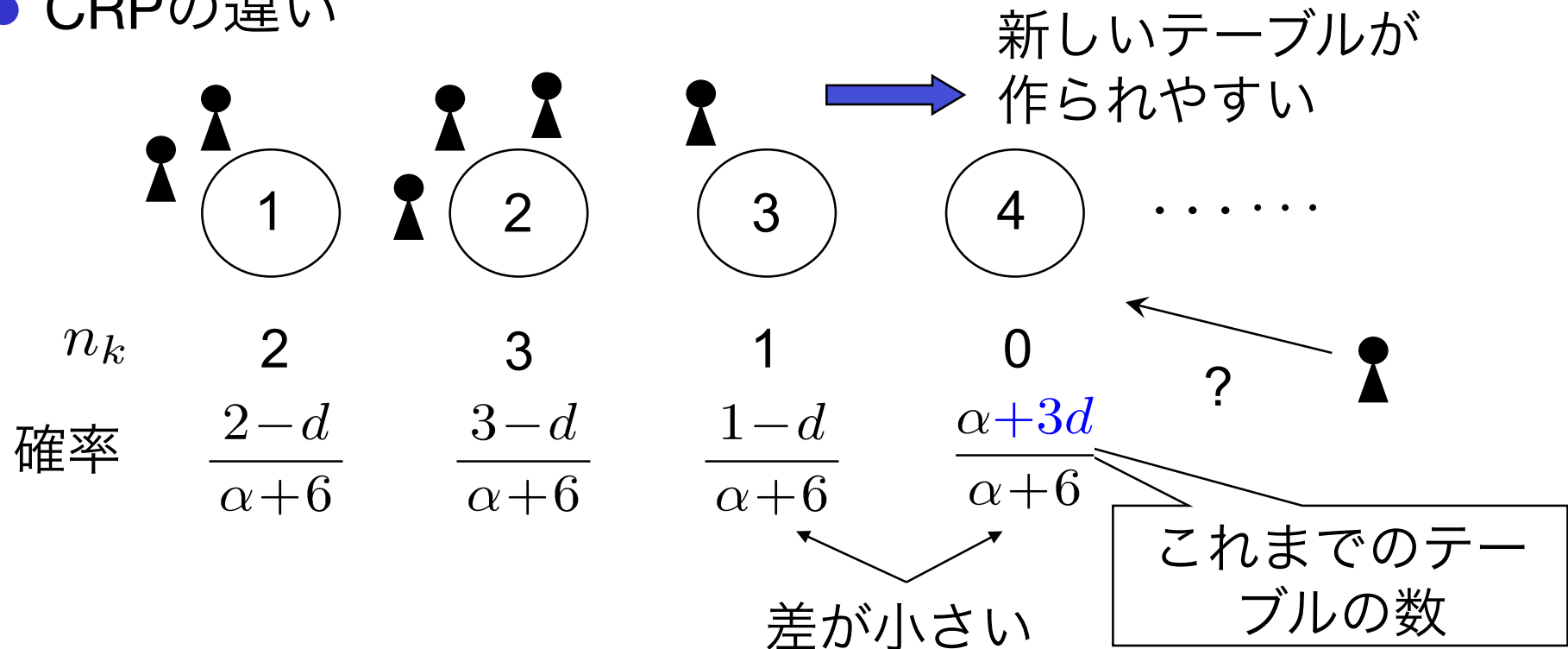


- 統計的自然言語処理の広い範囲で、 n グラムモデル (=言葉のMarkovモデル)が重要
 - n グラム... 前の $(n-1)$ 語に依存して次の語が出現
- n グラム分布を基底測度として、DPで $(n+1)$ グラム分布を生成する

Pitman-Yor過程

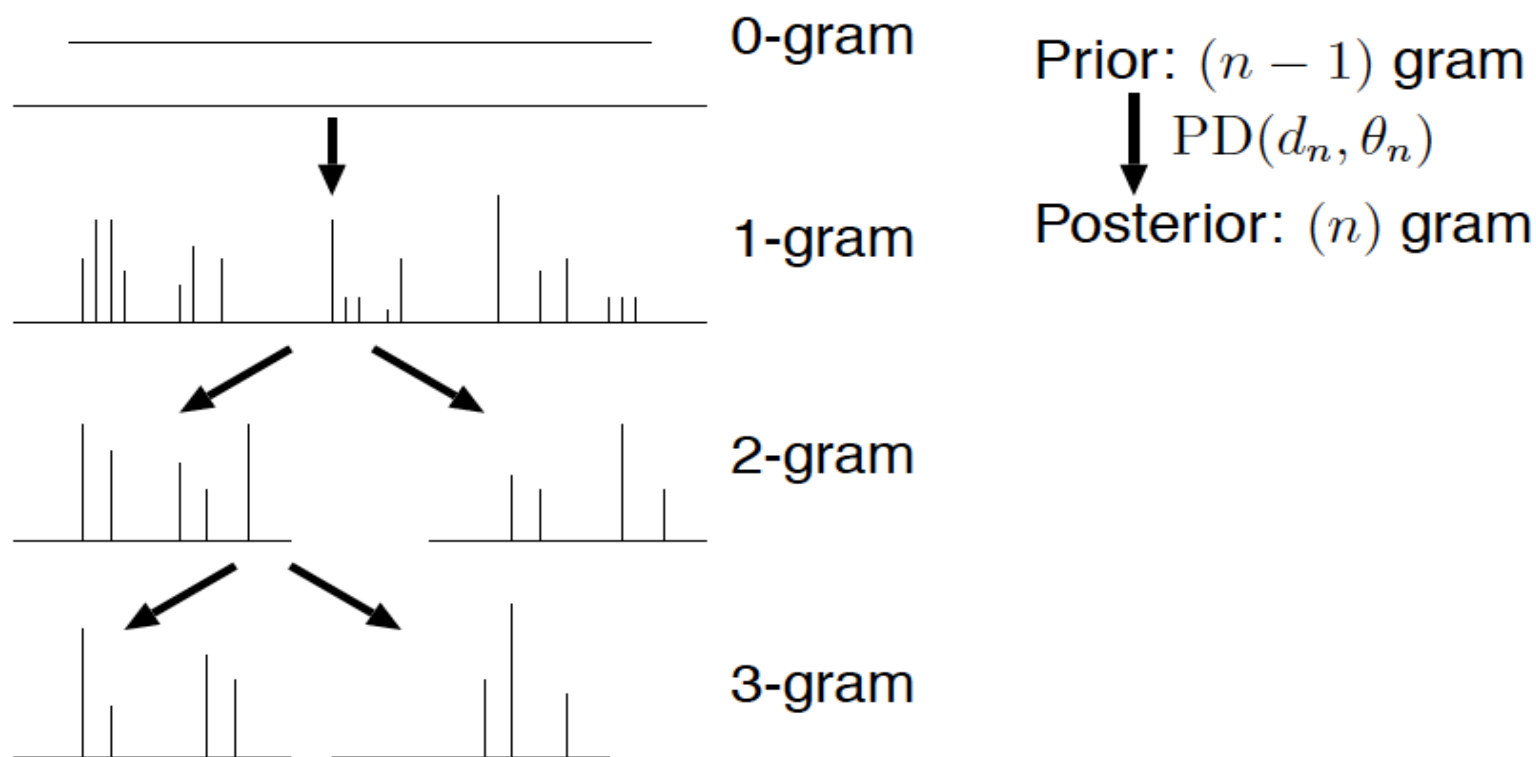
- Pitman-Yor過程 (Pitman and Yor 1997): $PY(\alpha, d)$
 - ディリクレ過程の拡張, Poisson-Dirichlet
 - 新たにディスカウント係数 d を持つ

- CRPの違い

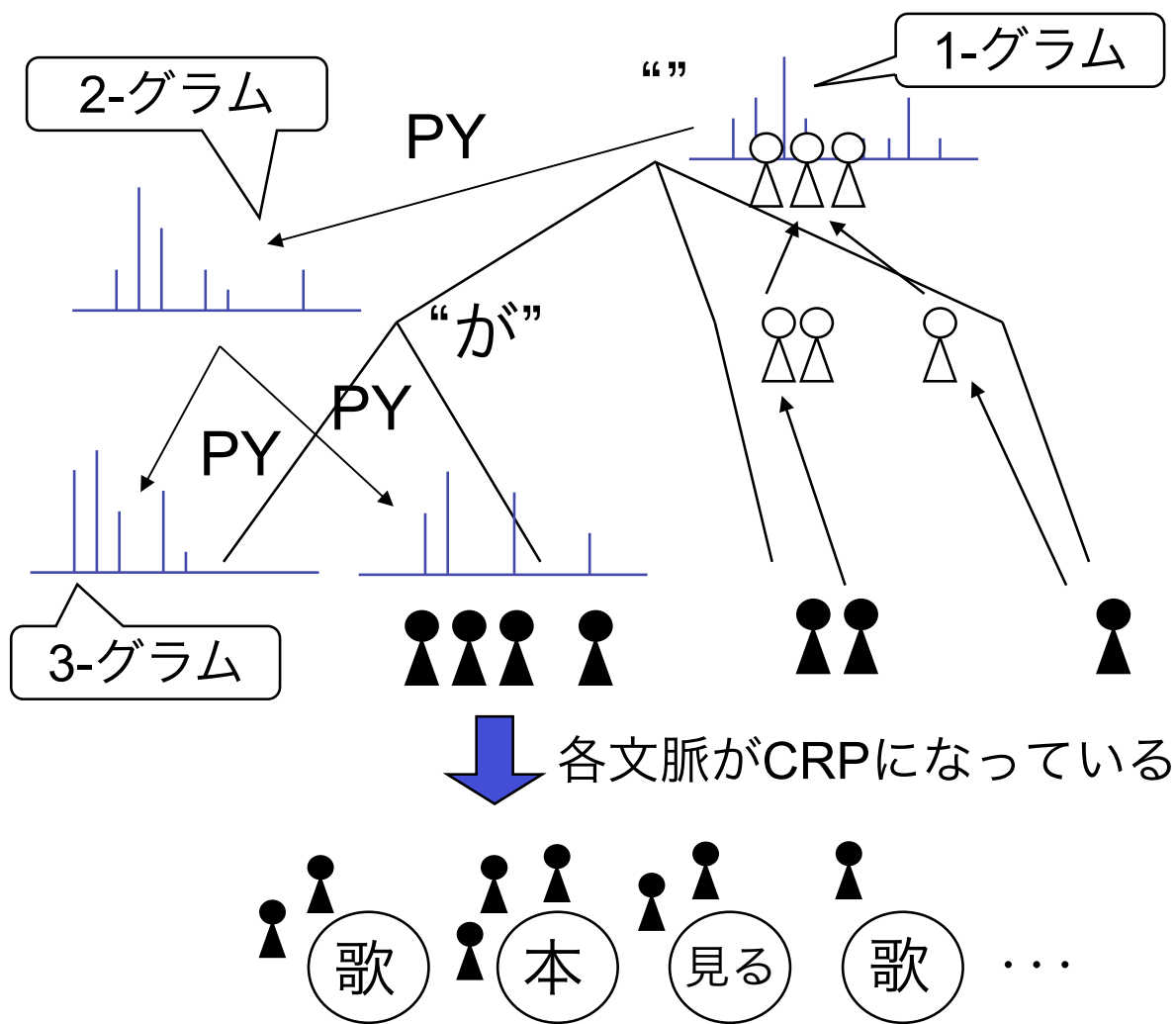


階層Pitman-Yor過程 (1)

- n グラム分布が、階層的に $(n-1)$ グラム分布からのPitman-Yor過程によって生成されたと仮定
 - 最初はUniform, だんだん急峻になる



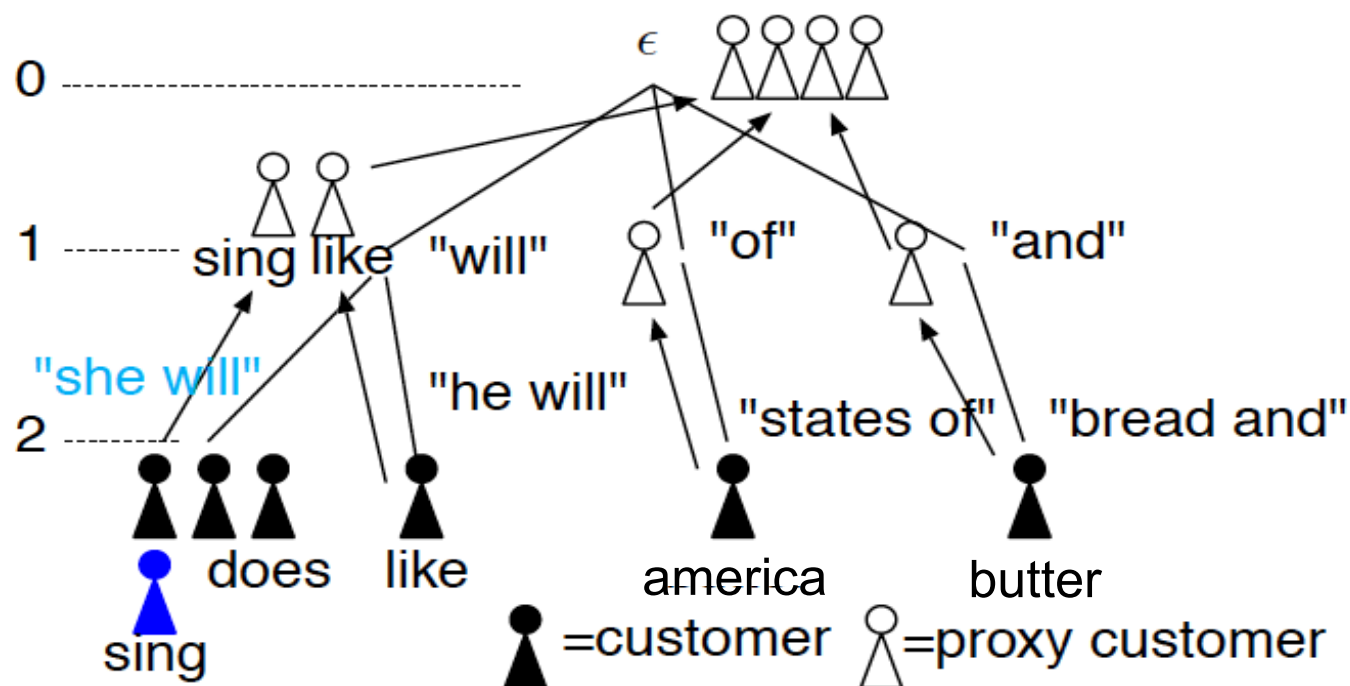
階層Pitman-Yor過程 (2)



- 各nグラム文脈hでの単語出現がCRPに従う
- カウントを複数のテーブルで表現
→ テーブルが増えたとき、客のコピー (代理客) を(n-1)グラム文脈に送る
 - sing a|song →
a|song →
song

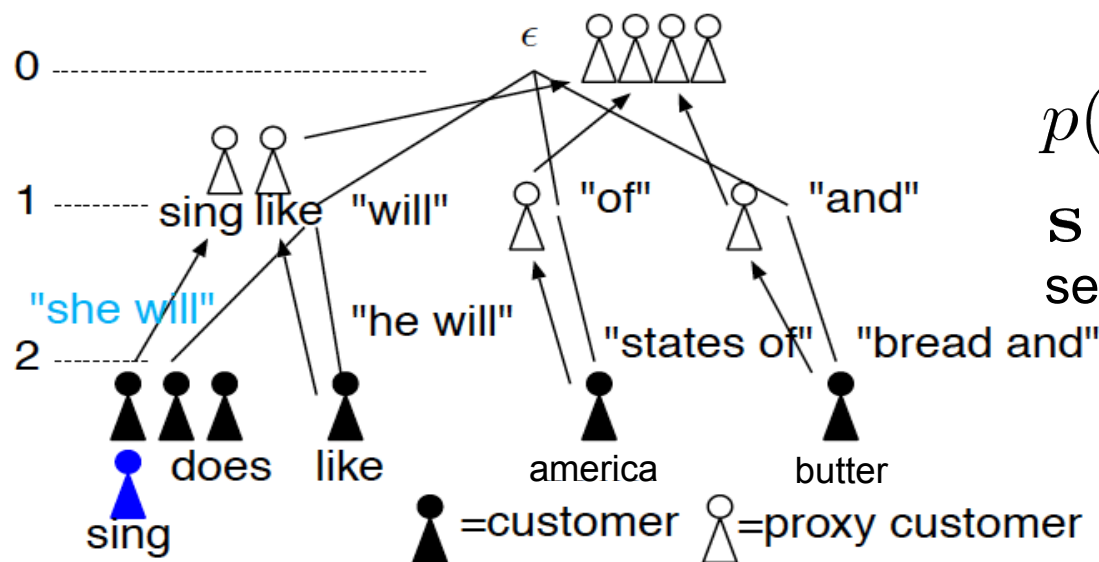
階層CRP表現

- 実際のカウント = 黒い客は、常に深さ(n-1)へ追加
 - 下のトライグラムの場合は、深さ2
- 適宜、確率的に、スムージングのためのカウントを親に再帰的に送る (白い客、代理客)



HPYLMの学習

- **HPYLM** (hierarchical Pitman-Yor language model) の学習 = 潜在的な代理客の最適配置
- Gibbs sampling: 客を一人削除して再追加、を繰り返す
 - For each $w = \text{randperm}(\text{all counts in the corpus})$,
 - 客 w と関連する代理客をモデルから削除
 - 客 w をモデルに追加 = 代理客を再サンプル



$$p(\mathbf{w}) = \sum_{\mathbf{s}} p(\mathbf{w}, \mathbf{s})$$

\mathbf{s} : 白い代理客の
seating arrangements

HPYLMの予測確率 (再掲)

$$p(w|h) = \frac{n(w|h) - d \cdot t_{hw}}{\theta + n(h)} + \frac{\theta + d \cdot t_h}{\theta + n(h)} p(w|h')$$

- 文脈 h の下での単語 w の予測確率
 - 一つ短い文脈 h' での同様な予測確率との混合
- $t_{hw} \equiv 1$ のとき、Kneser-Ney スムージングと一致
 - 実際には t_{hw} は $O(\log(n))$ で増えることが示されている

階層言語モデルによる教師なし形態素解析

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。)
```

X

- 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|、|...

形態素解析 (3)

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、…

קפיטליזם היא שיטה כלכלית וחברתית שהתפתחה באירופה בין המאה הששעשרה והמאה התשעעשרה, המבוססת בעיקרה על הזכות של פרטיסוקבות לבעלות פרטית על הרכוש לשימוש בובאופן חופשי, תוך הסתמכות על אכיפת זכויות הקניין באמצעות הרשות השופטת.

- 古語や、未知の言語の文に関しては
そもそも何が単語なのかわからない！
- 世界の他の言語でも同様の問題

中国語:

Mozilla Firefox

ファイル(E) 編集(E) 表示(V) 履歴(S) ブックマーク(B) ツール(T) ヘルプ(H)

http://www.tsinghua.edu.cn/qhdwzy/board_xxqk2.jsp?board=12&bid2=1201&pagen☆ Google

清华大学
Tsinghua University

学校概况 | 院系设置 | 管理机构 | 科学研究 | 教师队伍 | 人才培养 | 招生就业 | 虚拟校园

学校沿革
历任校长
统计资料

校长致辞



喜迎百年华诞 再铸新的辉煌

——2011年新年献辞

● 校长 顾秉林

一世纪沧桑砥砺，一百年春华秋实。此时此刻，2010年的余晖散去，2011年的曙光在前！此时此刻，清华正在送走她的第一个百年，迈向新百年的征程！在辞旧迎新、钟声回荡之际，请允许我代表学校，向清华全体同学、教职员工、离退休人员和海内外广大校友，向长期关心支持清华发展的各界人士，致以崇高的敬意和新年的祝福！

一百年来，清华大学的发展始终与国家民族的命运休戚与共，形成了优良的精神传统和鲜明的办学特色。一代代清华人“自强不息、厚德载物”，涌现出众多学术大师、兴业英才和治国栋梁，为中国社会进步和世界文明发展作出了重要贡献。特别是近年来，在国家的大力支持下，学校致力于世界一流大学建设，积极探索中国特色的“大学之道”，各项事业不断取得新的进展，正在跻身世界一流大学的行列。

大学之道，育人为本。一年来，以“清华新百年人才培养的使命与战略”为主题的第23次教育工作讨论会顺利举行，全校师生就推动办学优势转化、培养拔尖创新人才进一步取得共识。“清华学堂人才培养计划”以及多项教育教学改革措施相继实施。招生工作大力推进多元评价、兼顾拔尖与公平，生源质量进一步提高。就业毕业生中，超过80%选择到国家重要行业和领域建功立业。外国留学生规模不断扩大，结构进一步优化，外国研究生在学规模居全国高校首位。

大学之道，学术为魂。一年来，我校积极面向国际学术研究前沿和国家重大战略需求开展高

タイ語:

iGoogle - Mozilla Firefox

ファイル(E) 編集(E) 表示(V) 履歴(S) ブックマーク(B) ツール(T) ヘルプ(H)

http://www.chula.ac.th/

Google

ขอเชิญชาวจุฬาฯ ร่วมนับถอยหลังสู่การเป็นเจ้าภาพ “จามจุรีเกมส์”
กีฬามหาวิทยาลัยแห่งประเทศไทย ครั้งที่ 38 วันที่ 15-22 มกราคม 2554

จามจุรีเกมส์ 38
Research and Development

Chulalongkorn University
จุฬาลงกรณ์มหาวิทยาลัย
เสาหลักของแผ่นดิน

มหาวิทยาลัย
อันดับหนึ่ง
ของประเทศไทย
ผ่านสื่อหลัก

- ข่าวด่วน
- ค้นหาด้วยเสิร์จเอนจินจฬาร
- หน้าหลัก
- ข่าสดปัดจ, ปัน
- บุคลากร
- ศิษย์เก่า

- สภามหาวิทยาลัย
- ติดต่อเรา
- แผนผังเว็บไซต์
- ธงชาติ

== เมนูด่วน ==

ศูนย์ประสานงานสื่อมวลชน
แข่งขันกีฬามหาวิทยาลัยแห่งประเทศไทยครั้งที่ 38

จุฬาฯ เปิดศูนย์ประสานงานสื่อมวลชน (Press Center)
การแข่งขันกีฬามหาวิทยาลัยฯ ครั้งที่ ๓๘ “จามจุรีเกมส์”

จามจุรีเกมส์ 38
สนับสนุนการเกิดโดยสภา
การแข่งขันกีฬามหาวิทยาลัยแห่งประเทศไทย ครั้งที่ 38
ระหว่างวันที่ 15-22 มกราคม 2554
ณ กรุงเทพมหานคร 10110, ประเทศไทย

ศาสตราจารย์
ดร.วิจิตร ศรีสอ้าน

ศูนย์ประสานงานสื่อมวลชน
แข่งขันกีฬามหาวิทยาลัยแห่งประเทศไทยครั้งที่ 38

การรับสมัคร คณะ สำนักวิชา และสถาบัน ผู้มาเยือน แนะนำจุฬาฯ สื่อสารองค์กร จุฬากับนานาชาติ การค้นคว้าวิจัย กิจกรรมและการบริการสังคม

ペルシャ語:

The screenshot shows the homepage of the Isfahan University of Technology (IUT) website. The browser is Mozilla Firefox, and the URL is http://www.iut.ac.ir/. The page features a blue header with the university's name in Persian calligraphy and its logo. Below the header, there are navigation links for various departments and services. The main content area includes a large banner image of a park and a list of news items. On the right side, there are sections for 'دانشجویان' (Students), 'استاد و کارکنان' (Faculty and Staff), 'سرویس ها و خدمات' (Services and Facilities), and 'فناوری اطلاعات' (Information Technology).

isfahan university of technology - Google 検索 - Mozilla Firefox

ファイル(E) 編集(E) 表示(V) 履歴(S) ブックマーク(B) ツール(T) ヘルプ(H)

http://www.iut.ac.ir/

جستجوی سایت

وب سایت دانشگاه

EN فارسی

کتابخانه آموزشهای الکترونیکی سامانه الکترونیکی دروس سامانه اتوماسیون اداری سامانه گلستان پست الکترونیکی

[دانشجویان]
خدمات دانشجویی امور فرهنگی مرکز مشاوره

[استاد و کارکنان]
اداره رفاه اداره کارگزینی بهداشت و درمان استاد

[سرویس ها و خدمات]
سرویسهای چند رسانه ای ارتباط با مسئولین

[فناوری اطلاعات]
مرکز فناوری اطلاعات هسته محتوای دیجیتال پورتال

اطلاعیه پذیرش و ثبت نام موفق دانشجویان مقطع دکتری ورودی نیمسال دوم 1389-90

تقدیر از استاد دانشگاه صنعتی اصفهان در ششمین جشنواره تحلیلی از پژوهشگران استان اصفهان

قطب علمی مطالعات کودکی آب و خاک دانشگاه صنعتی اصفهان در نمایشگاه دستاوردهای پژوهش و فناوری کشور

بیانیه دانشگاه صنعتی اصفهان به مناسبت "سالگرد نهم دی ماه روز حملانه و بصیرت"

سومین شماره هفته نامه الکترونیک دانشگاه صنعتی اصفهان منتشر شد

سومین گردهمایی روسای مراکز کار آفرینی دانشگاه های منطقه مرکزی کشور به میزبانی دانشگاه صنعتی اصفهان برگزار شد

حضور پژوهشگرده شهبود اعتباری بسیج دانشجویی دانشگاه صنعتی اصفهان در نمایشگاه پژوهش و فناوری کشور

شیوه نامه اجرائی آیین نامه میهمانی و انتقال دانشجویان دانشگاه صنعتی اصفهان

دانشگاه صنعتی اصفهان برترین دانشگاه کشور در تبدیل علم به ثروت ترشید اخبار

درباره ی دانشگاه افتخارات رویدادها اخبار

(Isfahan university of technology, Iran)

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 w の $p(w|s)$ 確率

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

を最大にする \hat{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)

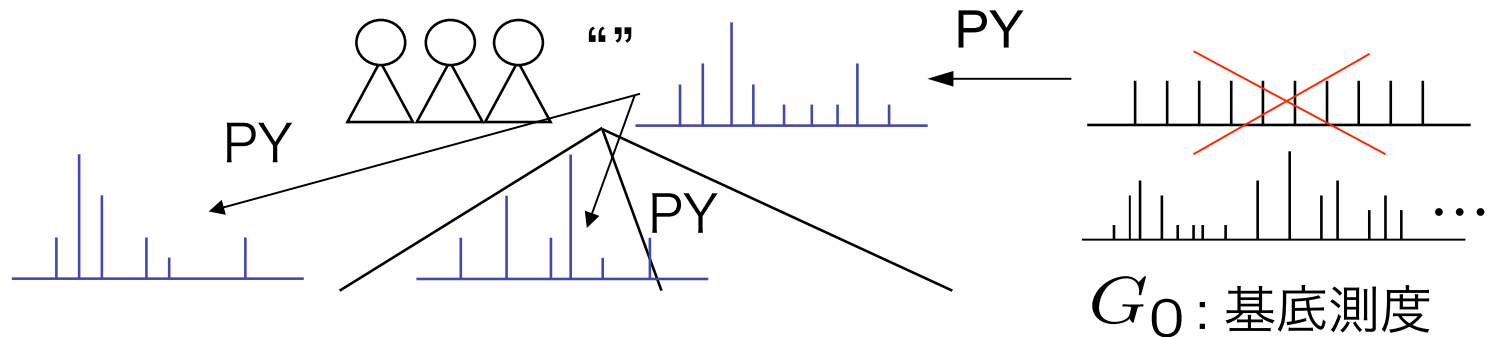
文の確率: nグラムモデル

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

文頭を表す特殊文字

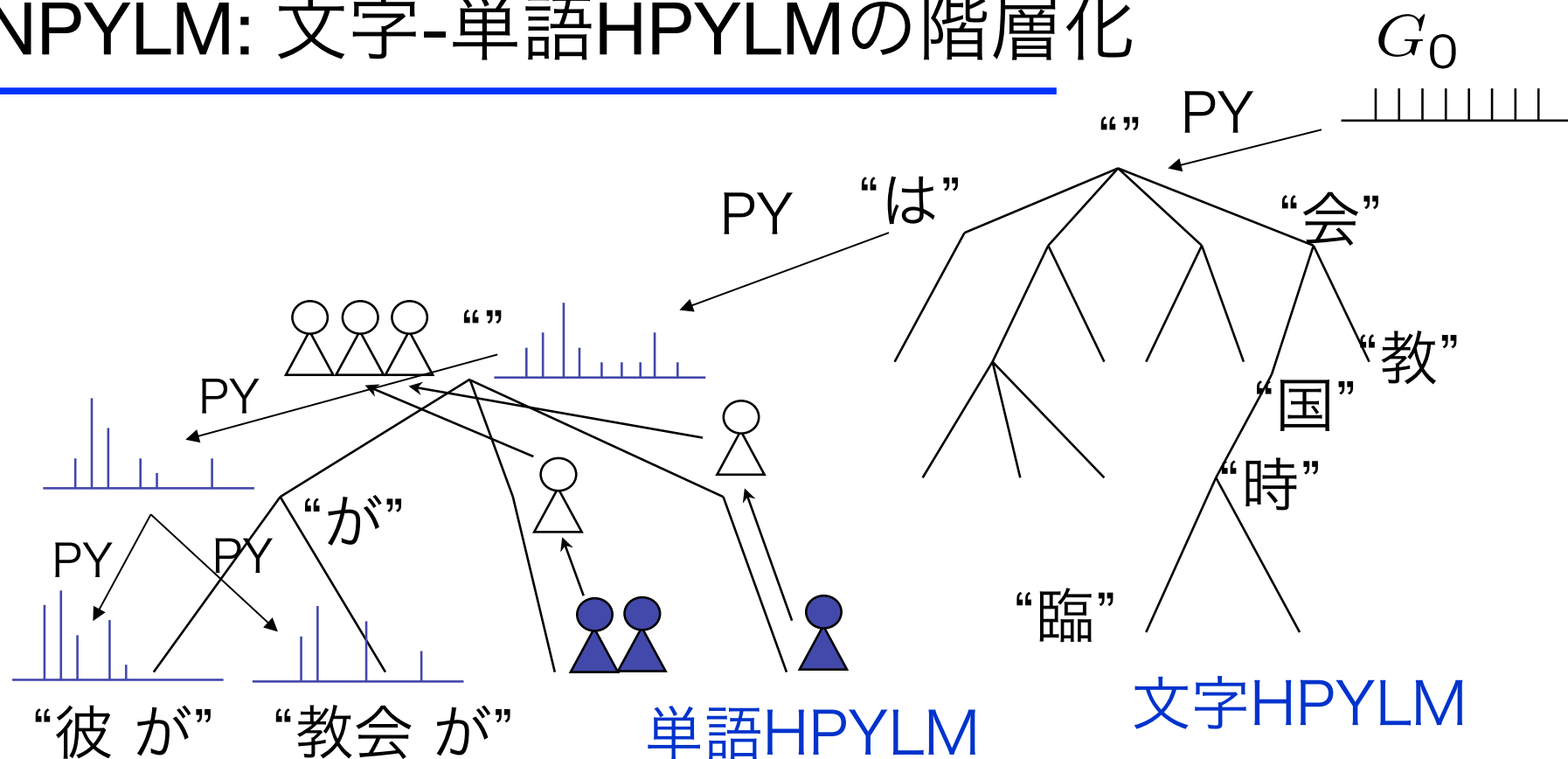
- 条件付き確率の積で文の確率を計算
 - 自然言語処理では、きわめて強力 (Shannon 1948)
 - 確率のテーブルは、ほとんどが0
 - 階層的なスムージングが不可欠
 - あらゆる部分文字列が「単語」になりうる
- ➡ 階層ベイズモデル: 階層Pitman-Yor過程言語モデル (HPYLM) (Teh 2006; Goldwater+ 2005)
- Pitman-Yor過程: ディリクレ過程 (GEM分布) の一般化

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

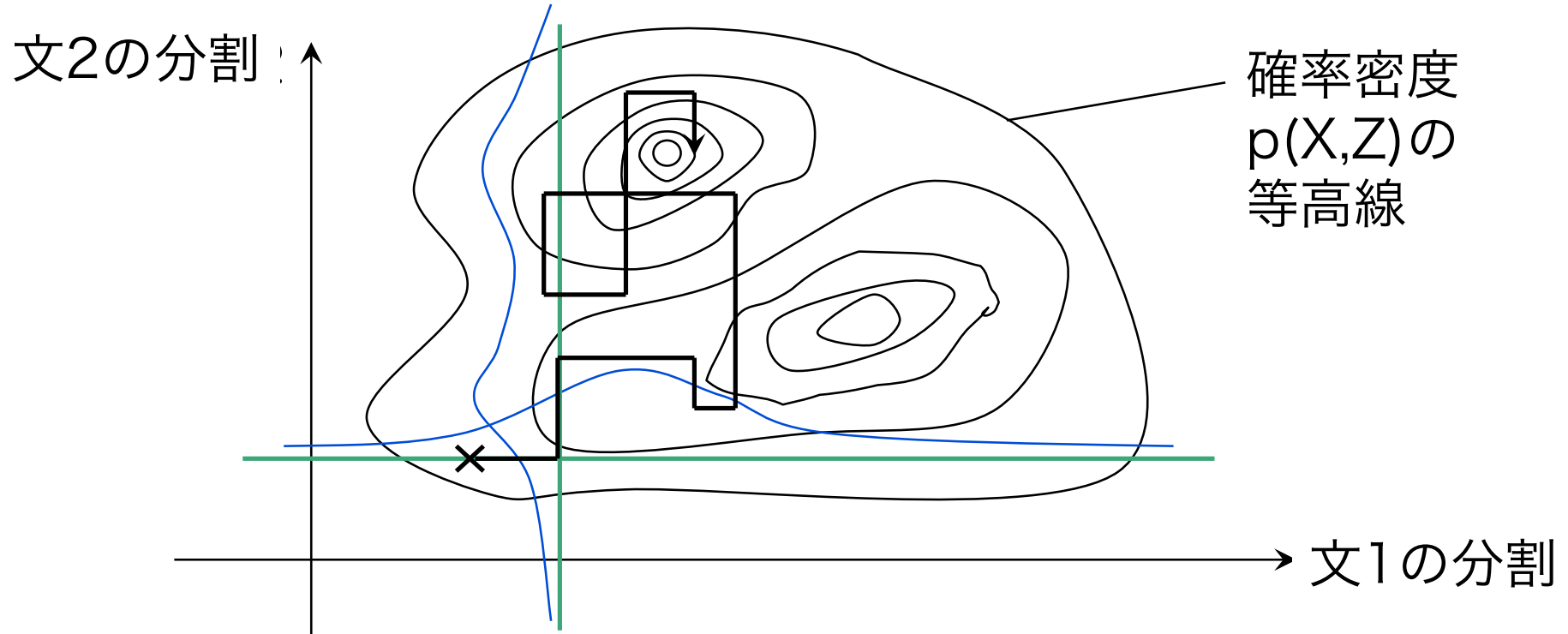
- 文がそれぞれ独立だと仮定すると、

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

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

- 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

done.

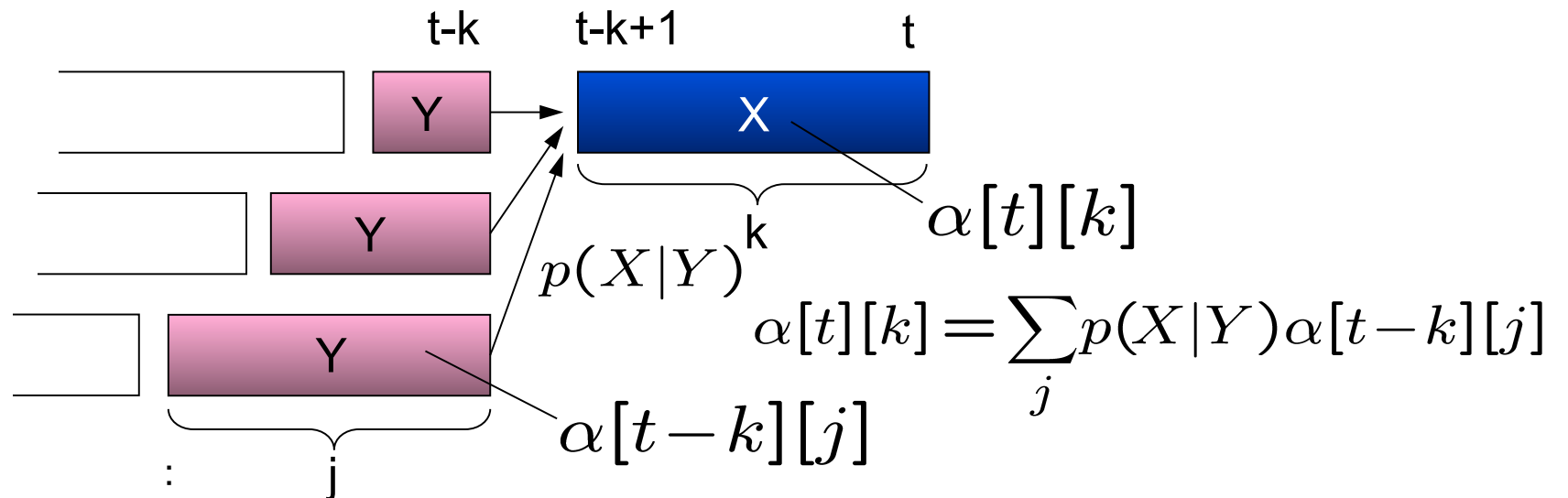
← Θ : 言語モデル
のパラメータ

Gibbs Samplingと単語分割

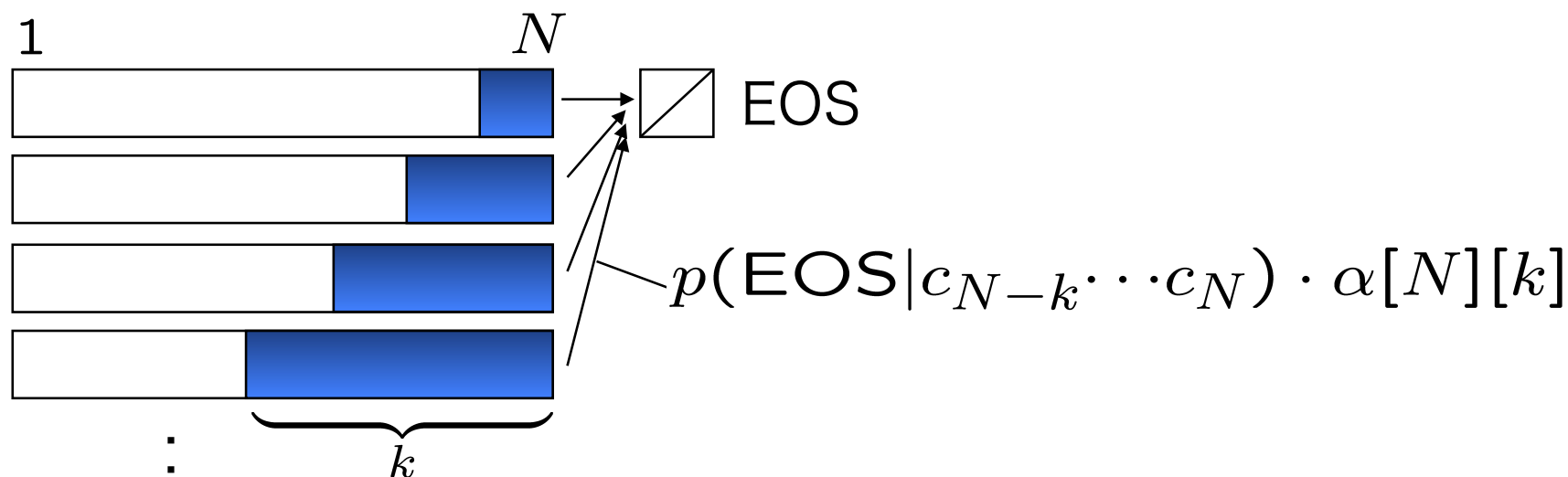
- 1 神戸では異人館 街の 二十棟 が破損した。
 - 2 神戸 では 異人館 街の 二十棟 が破損した。
 - 10 神戸 では 異人館 街の 二十棟 が破損した。
 - 50 神戸 では異人館 街 の 二十棟 が破損した。
 - 100 神戸 では 異人館 街 の 二十棟 が破損した。
 - 200 神戸 では 異人館 街 の 二十棟 が破損した。
- ギブスサンプリングを繰り返すごとに、単語分割とそれに基づく言語モデルを交互に改善していく。

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

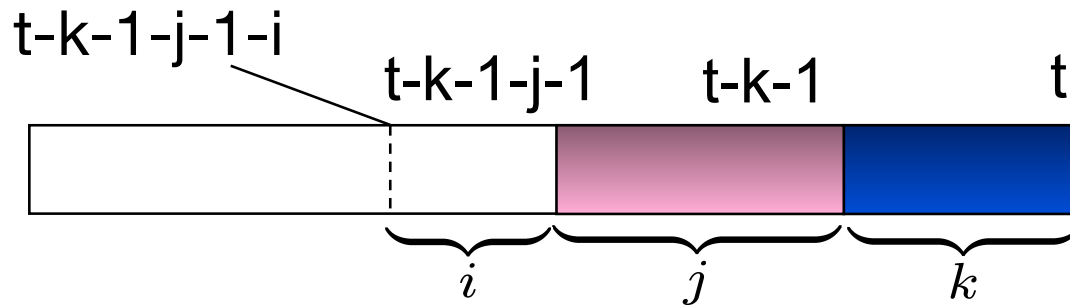


動的計画法によるデコード



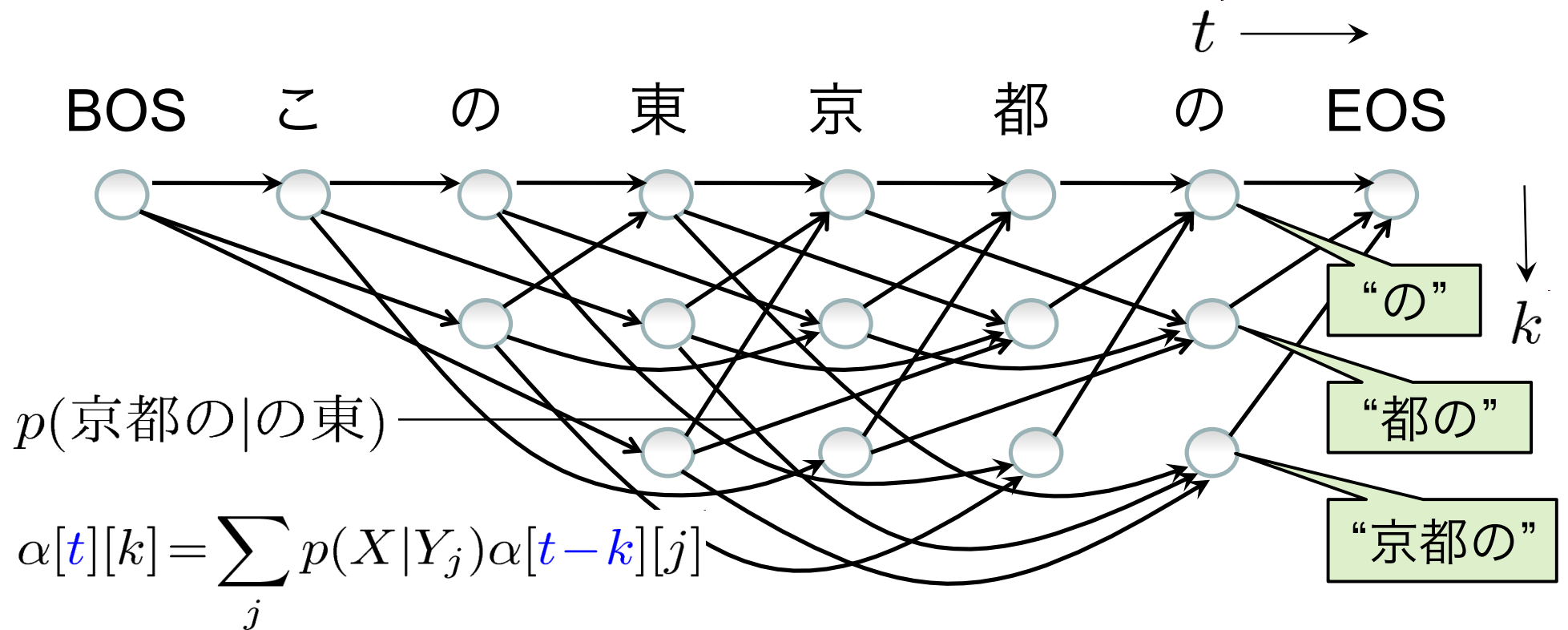
- $\alpha[N][k]$ = 文字列の最後の k 文字が単語となる文字列確率なので、EOSに接続する確率に従って後ろから k をサンプル
- $c_{N-k} \cdots c_N$ が最後の単語だとわかったので、 $\alpha[N-k-1][k']$ 使ってもう一つ前の単語をサンプル
- 以下文頭まで繰り返す

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)

NPYLM as a Semi-Markov model



- Semi-Markov HMM (Murphy 02, Ostendorf 96)の教師なし学習+MCMC法
- 状態遷移確率(nグラム)を超精密にスムージング

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。

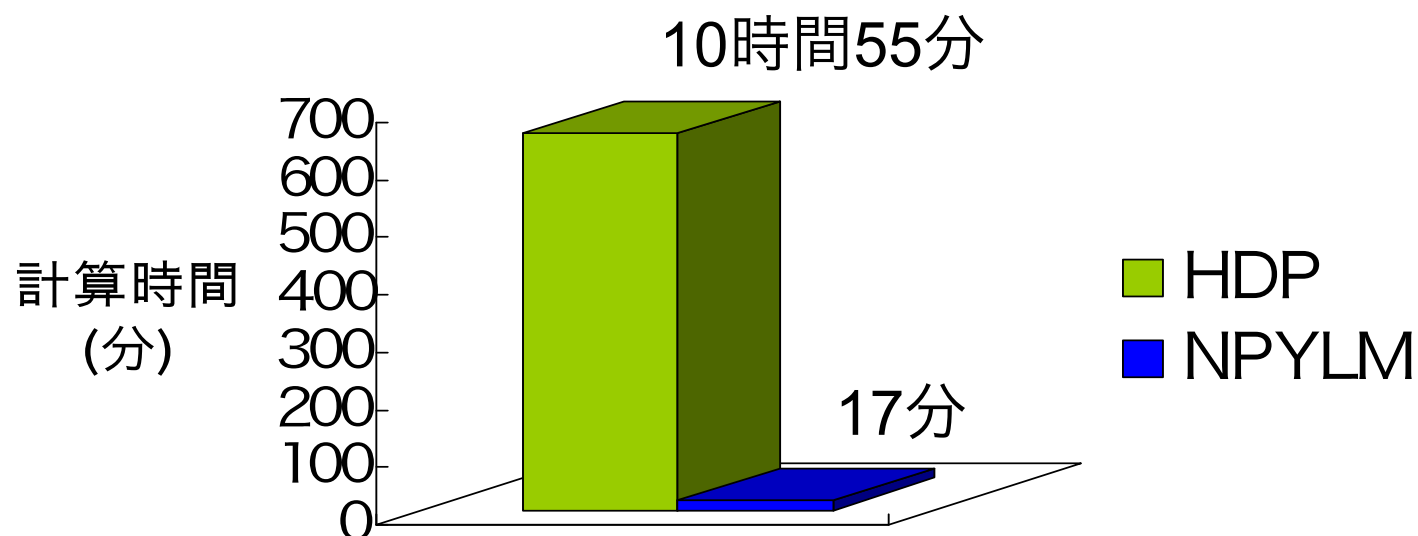
ビタビアルゴリズムで効率的に計算可能
(先行研究では不可能)

“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 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 ですよ 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 うー もの っ て や っ ぱ り ある ん だ す よ や っ て 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 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 ですよ 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 うー もの っ て や っ ぱ り ある ん だ す よ や っ て 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 ل انصار حركة المقاومة الاسلامية حماس .
واذا تحقق ذلك فان كى سل وفس كى ي ك ون قد حاز ثلث حواشي على ابرز ثلث اثة

Google translate:
"Filstinebsbptazahrplansarhrkpalmquaompalaslami phamas"
اي قل .

يخي "وقد استغرق اعداده خمسة اعوام . وقال ت دان ييل تومسون التي كتبت ال س ي ن ا ر ي و



الفلستيني بسبب تظاهرة ل انصار حركة المقاومة الاسلامية حماس .
تحقق ذلك ف ان كى سل وفس كى ي ك ون قد حاز ثلث حواشي على ابرز ثلث اثة

Google translate:
"Palestinian supporters of the event because of the Islamic
Resistance Movement, Hamas."
اي قل .

وقد استغرق اعداده خمسة اعوام . وقال ت دان ييل تومسون التي " تاريخي

“Alice in Wonderland”の解析



first, she dreamed of little Alice herself, and once again the tiny hands were clasped up on her knee, and the bright eager eyes were looking up into hers -- she could hear the very tones of her voice, and see that queer little toss of her head to keep back the wandering hair that would always get into her eyes -- and still as she listened, or seemed to listen, the whole place around 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 neighbouring 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 Alice herself ,and once again the tiny hand s were clasped upon 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 whiterabbit 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 -endingme a l ,and the ...

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

本研究のまとめ

- 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

全体のまとめ

- ノンパラメトリック・ベイズ法
 - ... 複雑なデータから、真に本質的な構造を取り出すための統計モデル
 - モデル選択や頻度での足切りと異なる精密なモデル
 - 潜在パラメータ数の組み合わせ爆発
- この他にも、非常に高度なモデルが存在
- 様々な分野に適用が進んでいる
 - 自然言語処理 (文法学習, 統計的機械翻訳, ...)
 - 画像処理 (画像分割, 画像認識, ...)
 - 機械学習全般 (IRM, Mondrian process, ...)

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