

Researcher2Vec: ニューラル線形モデル による自然言語処理研究者の可視化と推薦

持橋大地

統計数理研究所 数理・推論研究系
日本学術振興会 学術情報分析センター
daichi@ism.ac.jp

言語処理学会年次大会 2021
2021-3-16 (火)

概要

- 論文の内容から研究者をベクトル化し、キーワードで検索できるサーバを公開
 - Word2vecと同様のニューラル文書ベクトルを、特異値分解によりDoc2Vecの20倍の速度で計算
 - Doc2Vec, LDAによる検索を超えて最高精度

安全ではありません — clml.ism.ac.jp/nlp2vec/

自然言語処理：研究内容検索

トピックモデル

名前	類似度(%)	最近の論文(2件)	研究者情報
真光九月	56.47	統計的トピックモデルの.. / 2ツイートを用いた対話..	[詳細]
津田裕亮	55.51	トピック変化点検出に基.. / LDAトピックモデルに基づ..	[詳細]
中村明	54.12	LDAの文脈長最適化による.. / トピックモデルを用いた..	[詳細]
速水悟	54.01	probabilistic Polynomia.. / ジャンル別LSIの結果統合..	[詳細]
小林一郎	50.24	潜在的意味を考慮した効.. / 潜在情報を考慮したグラ..	[詳細]
横本大輔	47.39	震災を題材としたニュー.. / 文書集合の話題俯瞰手法..	[詳細]
吉岡真治	47.35	時系列中国語ニュース・.. / 日中時系列ニュースにお..	[詳細]
持橋大地	41.80	条件付確率場とベイズ階.. / Gibbs Samplingによる確..	[詳細]
河田容英	38.36	「契約・解約」に関する.. / 震災を題材としたニュー..	[詳細]
田村哲嗣	36.77	probabilistic Polynomia.. / ジャンル別LSIの結果統合..	[詳細]
森信介	33.44	Combining Active Learni.. / 仮名漢字変換ログを用い..	[詳細]
福原知宏	32.52	震災を題材としたニュー.. / 文書集合の話題俯瞰手法..	[詳細]
中崎寛之	28.27	「犯罪」分野に関連する.. / Wikipediaエントリに関連..	[詳細]
山崎真理子	26.88	Wikipediaをトピック体系.. / 共起語集合の言語間差異	[詳細]

情報・使い方

- 言語処理学会年次大会1995-2013年に発表された約20年分の論文内容から、研究者を検索します。対象としているのは、この期間に5本以上の論文があった499人の研究者です。
- 空白で区切って調べたい単語を入力して下さい。単語ベクトル・文書ベクトルを基にしていますので、正確に一致していなくても検索できます。
- 「談話** 確率*」のように単語の最後に "*" をつけると, "*"の個数だけその単語を強調して検索します。
- 技術的な背景については、言語処理学会2021で発表予定の論文「[Researcher2Vec: ニューラル線形モデルによる自然言語処理研究者の可視化と推薦](#)」をお読み下さい。
- Web周りが得意で、ACL anthology, arXiv等英



概要

- 研究者ベクトルを2次元に可視化することで、研究内容の近い研究者が客観的に可視化できる
(既存のシステムはヒューリスティックが多い)



背景

- コンピュータサイエンス分野では、論文数が激増
 - 言語処理学会2020…338本, ACL 2020…779本, CVPR 2020…1470本, NeurIPS 2020…1,900本以上
- 人手による査読割り当ては、もはや限界
 - 上の採択数の4倍～10倍程度×3～5程度の査読数が必要 (CVPR 2020では19000個の査読が必要)
 - TPMSが使われているが、研究者別情報は非公開
 - 未だに研究者の専門性を人手で調べる必要がある

The screenshot shows the 'Meta-Reviewer Console' interface. At the top, there are navigation links for 'Submissions' and 'Reviewers', and a role selector set to 'Meta-Reviewer'. Below this is a 'Bidding' section with a 'Show:' dropdown set to '25' and buttons for 'Clear All Filters', 'Restore Columns', and 'Actions'. The main table has columns for Paper ID, Title, Subject Areas (Primary and Secondary), Suggestions, Meta-Review, Discussion & Feedback, Relevance, TPMS Rank, Your Bid, and Value Qu. The first row of data shows Paper ID 26, Title 'Research Paper Zero 1', Subject Areas 'MARINE VESSELS -> Hull' and 'AUTOMOBILES -> Engines', Relevance 0.32, and TPMS Rank 3 (highlighted with a red box). The 'Your Bid' column shows 'Not Entered'.

Paper ID	Title	Subject Areas		Suggestions	Meta-Review	Discussion & Feedback	Relevance	TPMS Rank	Your Bid	Value Qu		
		Primary	Secondary							Min	Max	Av
26	Research Paper Zero 1	MARINE VESSELS -> Hull	AUTOMOBILES -> Engines			Status: Awaiting Decision	0.32	3	Not Entered			

背景 (2)

- 研究の興味が多様化しており、自然言語処理なら誰に聞いても同じ、ではもはやない



Yoshihiro KANAMORI

@yshhrknmr

返信先: @yshhrknmrさん

...蛇足ながら、学生から見たら「この研究分野なら先生は専門家だからこの研究テーマもイケるだろう」と思って相談してみると、実は学生が思っているよりも研究分野が非常に広大で、指導教員が得意なのはそのうちの限られた範囲でしかなかった、ということで最初のつぶやきの現象が多発しています...

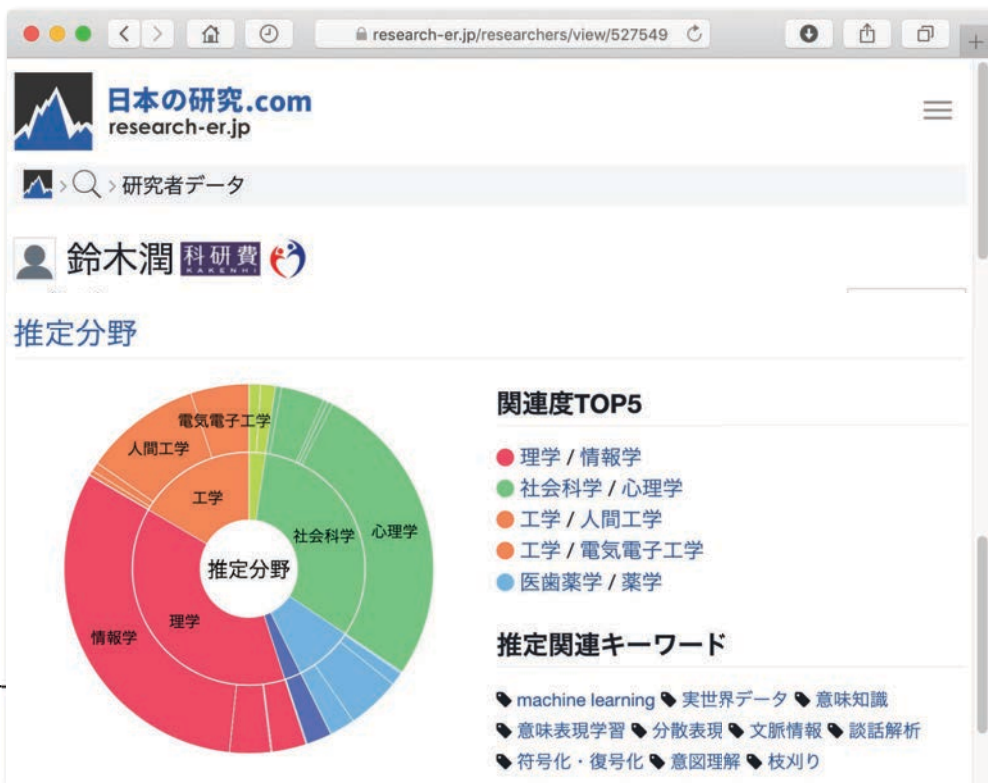
午前10:49 · 2020年12月23日 · Twitter Web App

背景 (2)

- この分野の研究をよく知っている先生は誰か？
→ 適切な指導者の発見 (学生、企業とも)
- 会議やジャーナルへの適切な研究者のリクルート
→ コネがなくてもコミュニティに加わられるようにしたい

既存システム

- 「日本の研究.com」, JDream Expert Finder, JSTサイエンスマップなど
 - 論文の内容ではなく、引用などメタ情報がベース
 - 本当の詳しい専門性は分からない、共著関係に引きずられる(コネ)



← 鈴木さんが
心理学??

既存研究

- 桂井ら (2016) : CiNIIの10万件の論文概要、300万語のテキストをトピック数 $K=500$ のLDAで解析
 - 次元圧縮しないベクトル空間モデルより高精度
- 持橋 (2019) : 学振内部の11万件の科研費申請書、3億語のテキストをトピック数 $K=4000$ の巨大なLDAで解析
 - LightLDA(WWW 2015)で高速化しても数日かかる
 - 科研費特別推進、基盤Sの審査に利用

ちなみに..



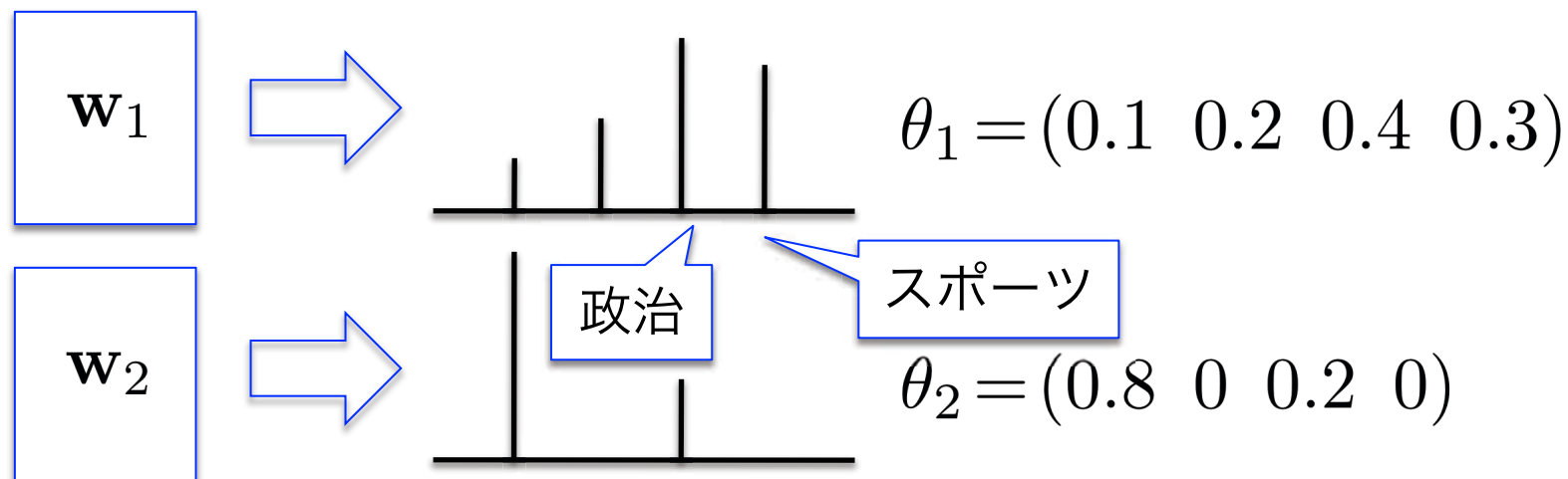
The screenshot shows a web browser displaying the KAKEN website. The URL is kaken.nii.ac.jp/ja/grant/KAKENHI-PROJECT-20H0448. The page title is 'Scholar2Vec: 研究者の多様な活動情報を埋め込める深層潜在空間の構築'. The page contains a table with the following information:

研究課題/領域番号	20H04484
研究種目	基盤研究(B)
配分区分	補助金
応募区分	一般
審査区分	小区分90020:図書館情報学および人文社会情報学関連
研究機関	同志社大学
研究代表者	桂井 麻里衣 同志社大学, 理工学部, 助教 (70744952)
研究分担者	大向 一輝 東京大学, 大学院人文社会系研究科(文学部), 准教授 (30413925) 梶原 智之 大阪大学, データビリティフロンティア機構, 特任助教(常勤) (70824960)
研究期間 (年度)	2020-04-01 - 2024-03-31

- 本研究とは、これとは完全に独立 (学振での必要性)
 - まだ論文が出ていない(はず)
 - 「様々な情報を埋め込む」 そうなので、今後に期待

LDAによる論文テキストの分析

- 確率的潜在意味解析 (LDA, Blei+ 2003) ...
各文書に対し、潜在的な話題 (トピック) 分布を推定



LDAの欠点

- 結果として、LDAは大変有効なモデルではあるが、
 - 負の相関を扱えない
 - 細かい意味的な違いに対処するのに限界がある
- 研究者や申請書の内容は和が1のトピック分布 θ

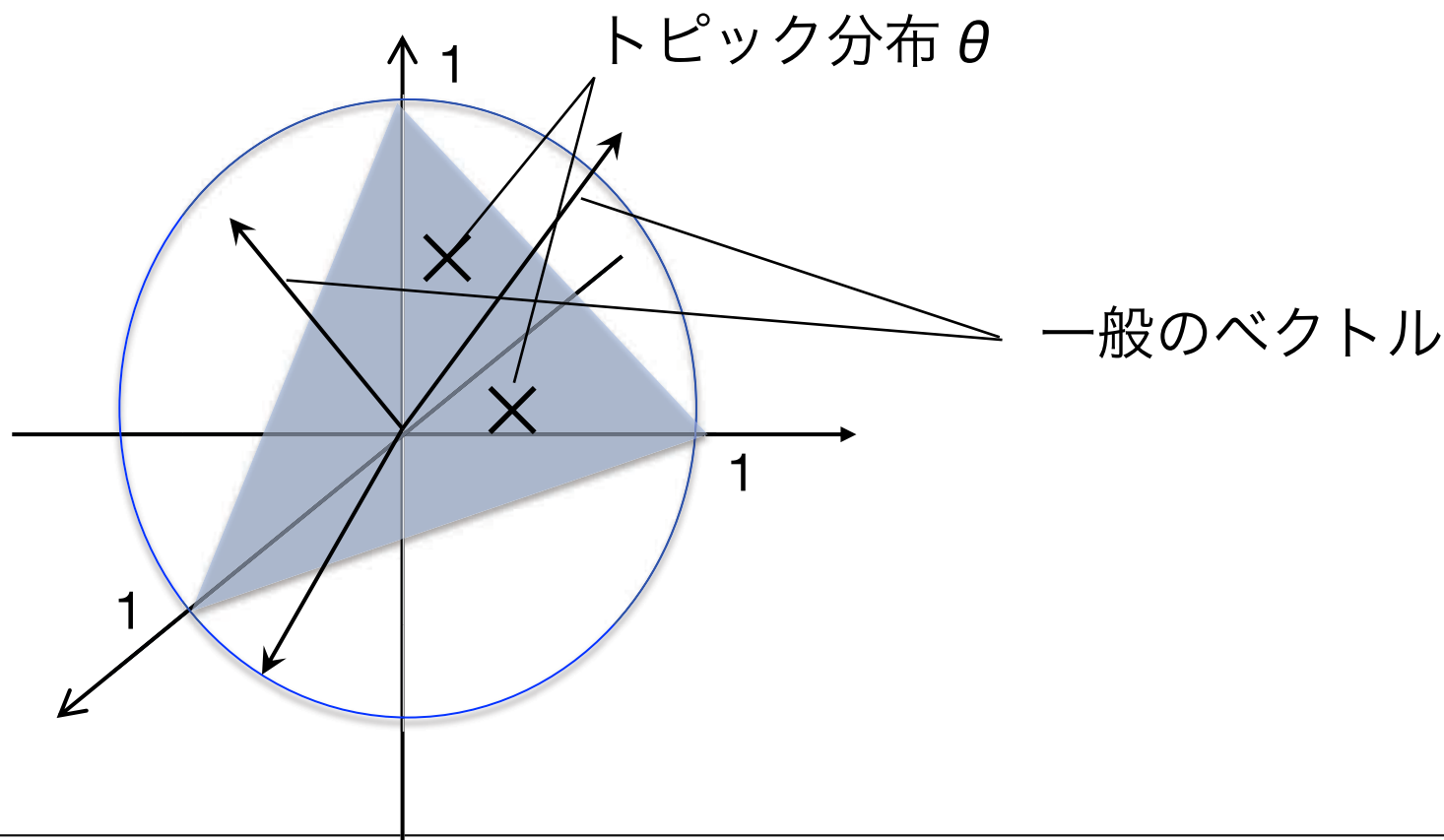
$$\theta = \begin{array}{c} | \\ | \\ | \\ | \\ \hline 1 \quad 2 \quad \dots \quad K \end{array}$$

で表現されるため、

- 「経済学者だが、数理的には解析学がベース」
「経済学者だが、数理的にはゲーム理論がベース」
などの細かい違いを θ で捉えるのは、非常に難しい
- 研究者の推薦精度にある程度限界がある

LDAの欠点 (2)

- 数学的には、LDAは単体の上でしか文書をモデル化していない
→ ベクトル空間のごく一部しか使っていない

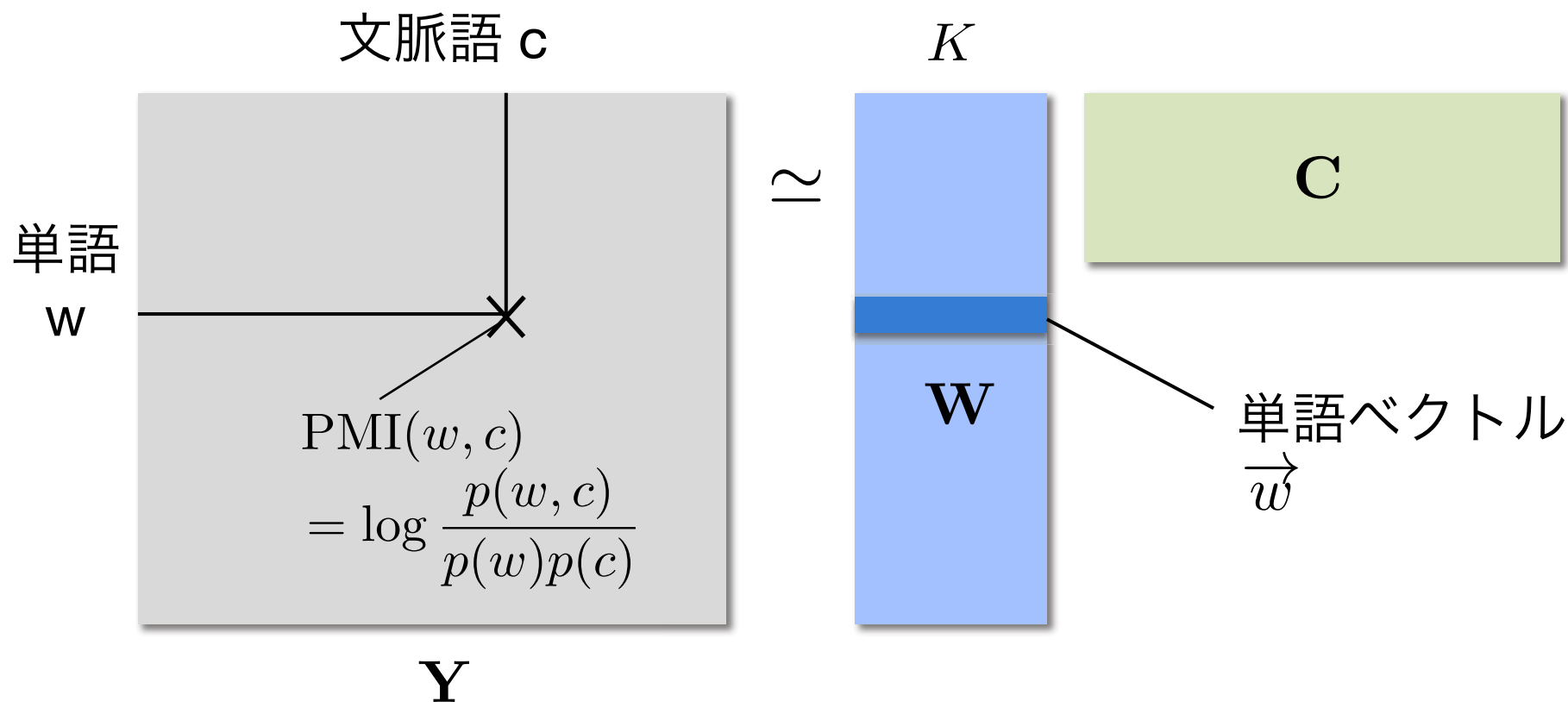


文書のベクトル表現

- 確率分布 θ で文書を表現するのをやめて、一般のベクトル \vec{d} で文書や研究者を表現すればよい?
→ RaP (Gehler+ 2006), RSM (Salakhutdinov 2009), Doc2Vec (Le and Mikolov 2014) , NVDM (Miao+ 2016) など多数あるが..
- ニューラル手法なので、一般に学習が難しい
- しかし...

Word2vecの数理

- 単語をベクトル化する、有名なWord2vec (Mikolov+ 2013) は、以下の自己相互情報量行列の行列分解と等価であることが示されている (Levy and Goldberg 2014)

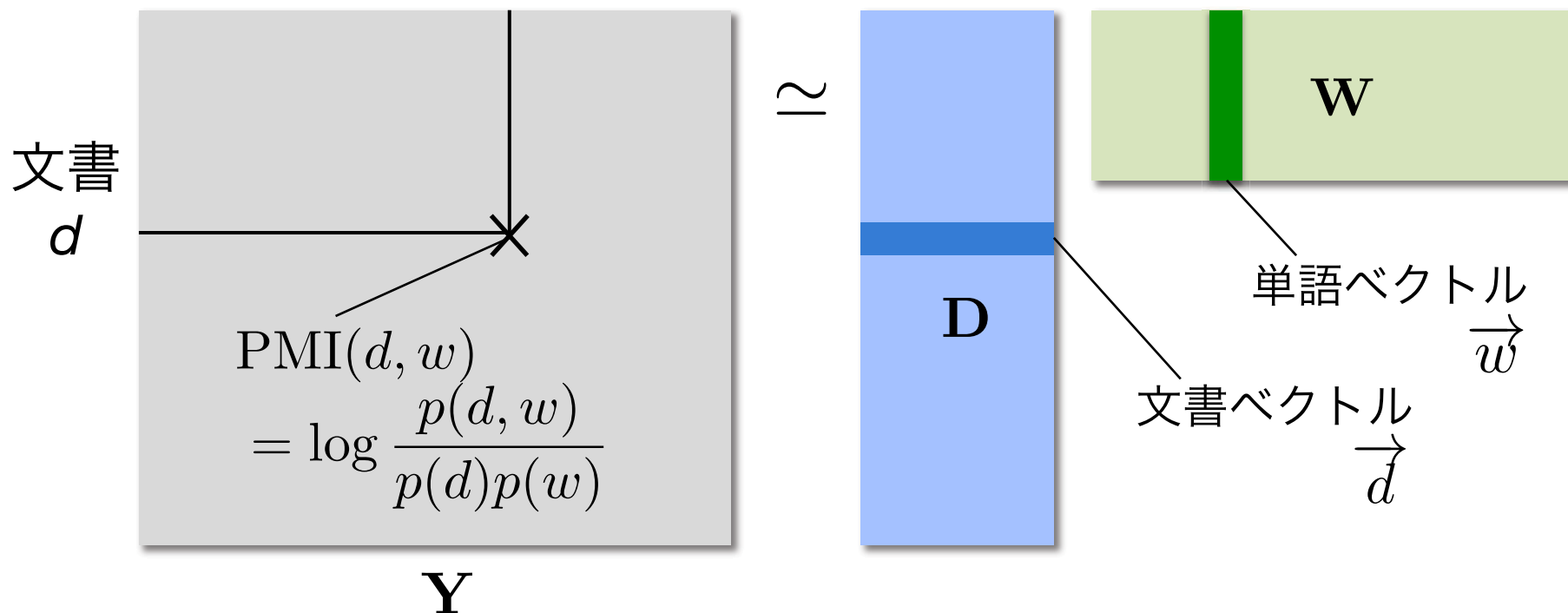


Word2vecから文書ベクトルへ

- 単語→文書、文脈語→含まれる単語 に置き換えれば、SVDで簡単に「文書ベクトル」と「単語ベクトル」を計算できる

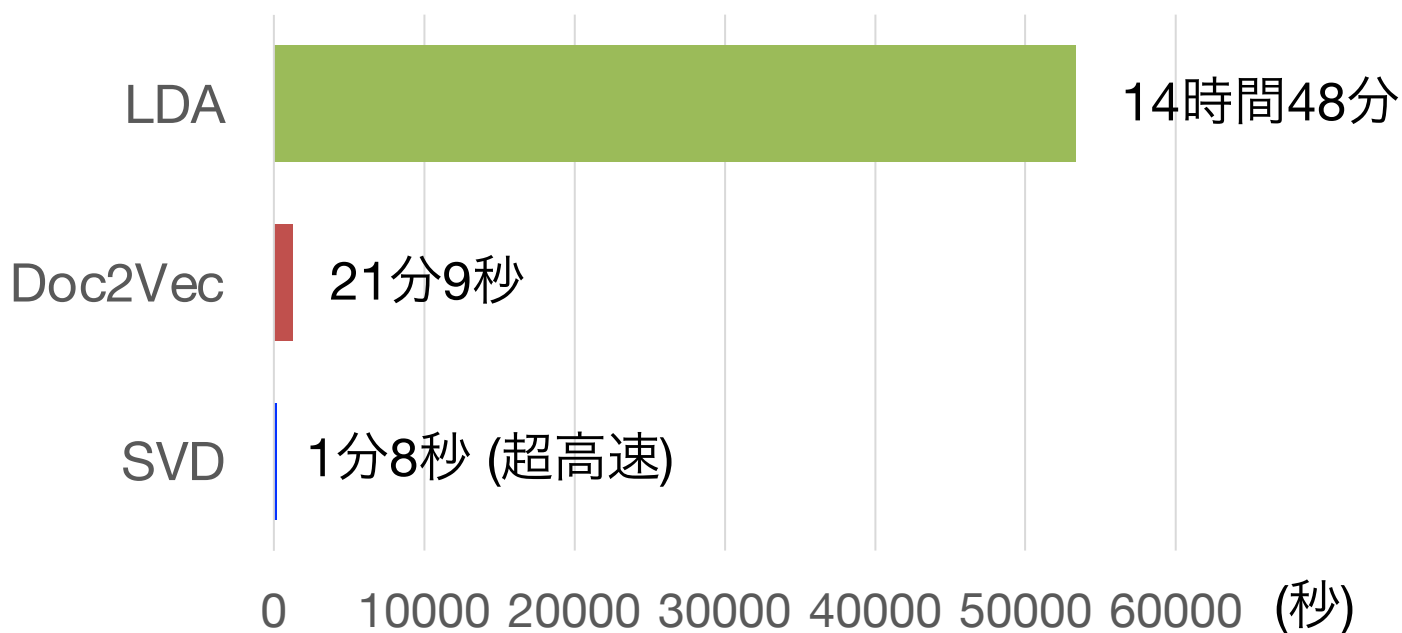
(注： $\log \frac{p(d, w)}{p(d)p(w)} = \log \frac{p(w|d)}{p(w)}$)

単語 w



計算時間の比較

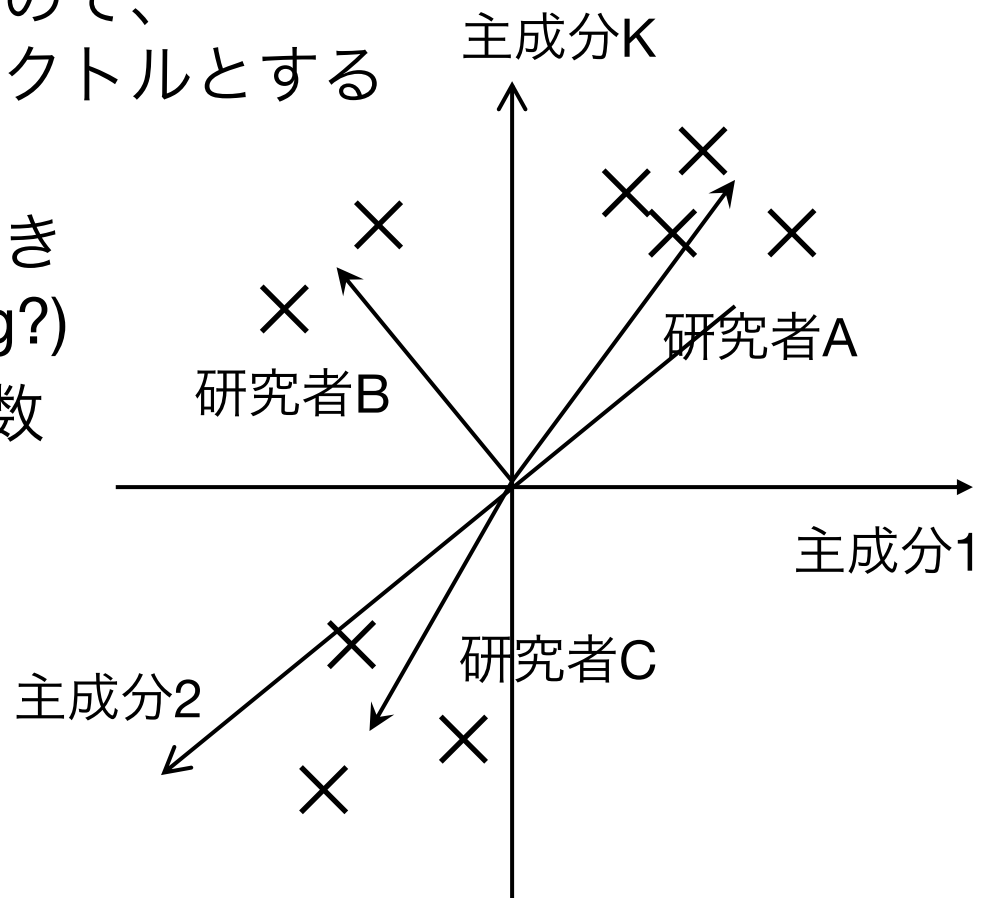
- K=1,000次元(トピック)で同じコーパスから学習した場合、



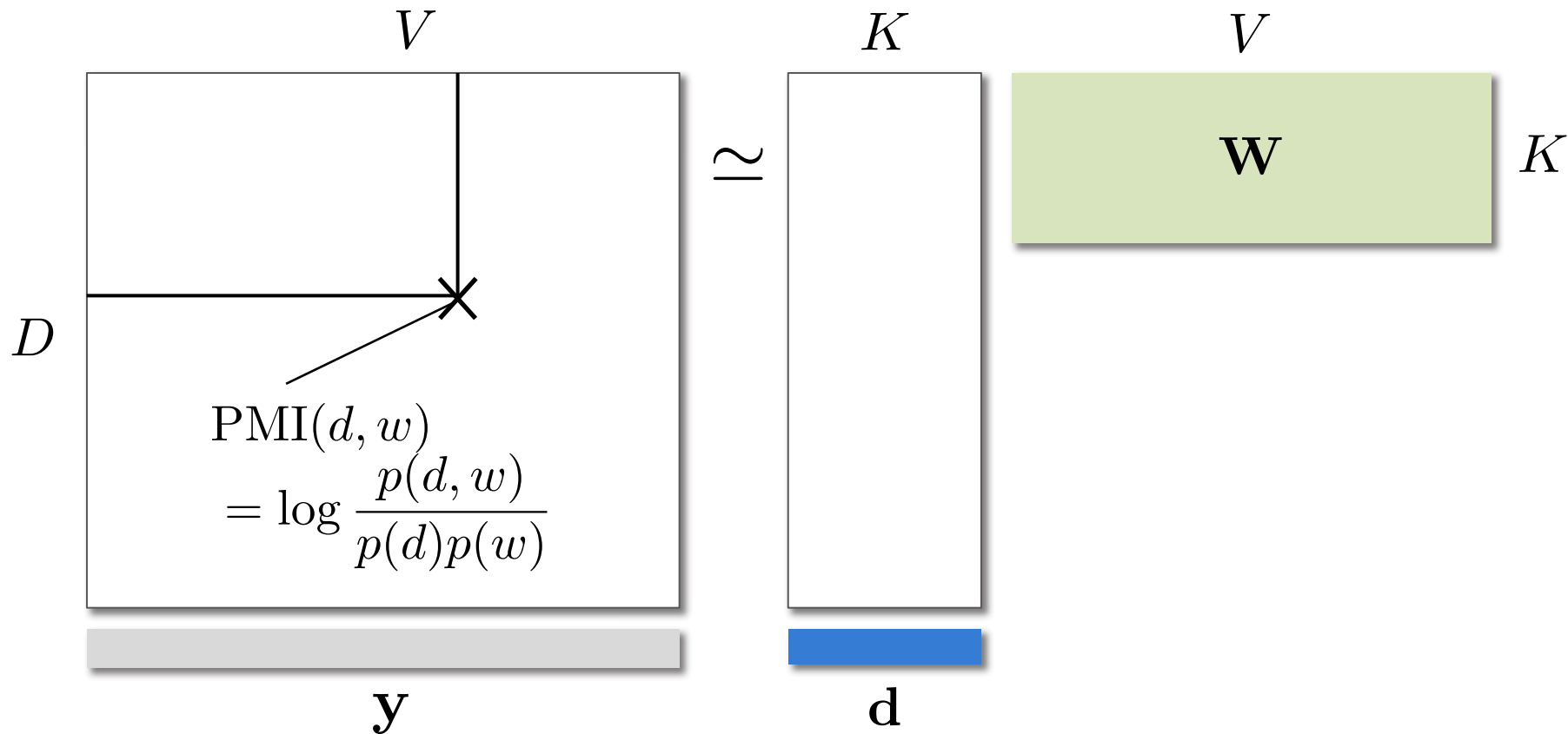
- LDAはGibbs 1,000 iteration, Doc2Vecは100 epochs
- データが巨大な場合、提案法はredsvdなども使える

文書ベクトルから研究者ベクトルへ

- 各研究者について、書いた論文/科研費申請書の文書ベクトルが得られるので、それらの平均を研究者ベクトルとする(最尤推定)
- 本来は分散も推定するべき(Kernel mean embedding?)
- 次元Kは、最大値は文書数(実験では3582)



キーワード検索の方法

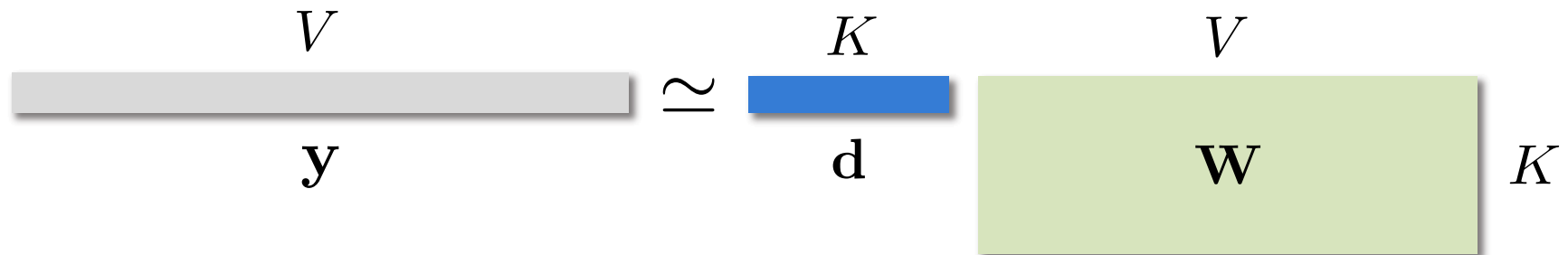


- クエリを仮想的な「文書」 y と思うと、

$$y \simeq dW$$

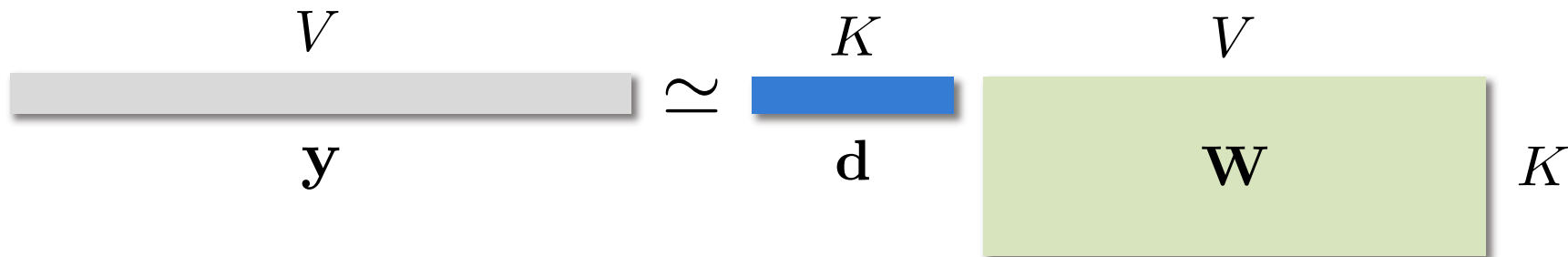
が成り立っている (d は対応する文書ベクトル)

キーワード検索の方法 (2)

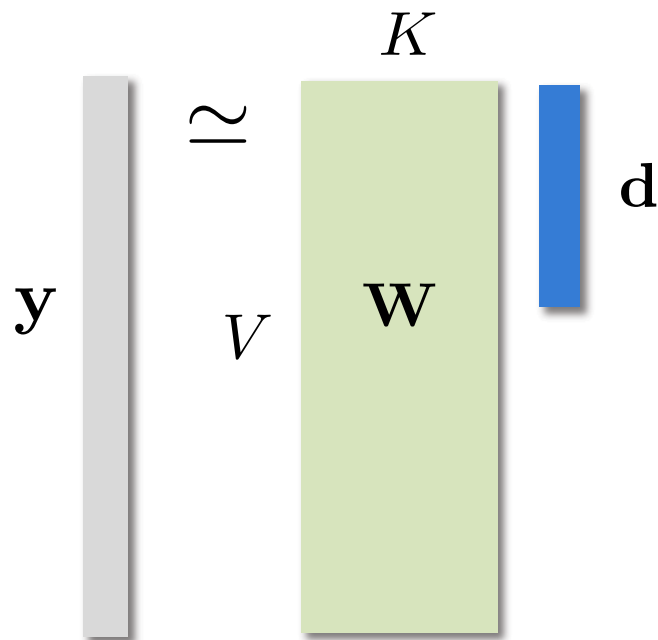


- y の要素はPMI $\log p(w|y)/p(w)$ なので、クエリの部分に1を立て、残りは $\log(1)=0$ のベクトル
 - クエリ単語の最後に * を付けると、1ではなく2,3,... にしてその単語を強調できる (例: “neural** model”)
- $y \sim dW$ の方程式は等式ではなく近似
- 得られた「文書ベクトル」 d を研究者ベクトルと比べればよい

キーワード検索の方法 (3)



を書き直すと、



- これは線形回帰モデル！

$$y \simeq Wd$$

- よって、 d の最適解は通常のOLSで、

$$d = (W^T W)^{-1} W^T y$$

キーワード検索の方法 (4)

- 線形回帰の基本ですが、二乗誤差を最小化したいので

$$\begin{aligned} E &= |\mathbf{y} - \mathbf{W}\mathbf{d}|^2 \\ &= (\mathbf{y} - \mathbf{W}\mathbf{d})^T (\mathbf{y} - \mathbf{W}\mathbf{d}) \\ &= \mathbf{y}^T \mathbf{y} - 2\mathbf{d}^T \mathbf{W}^T \mathbf{y} + \mathbf{d}^T \mathbf{W}^T \mathbf{W} \mathbf{d} \end{aligned}$$

- よって

$$\begin{aligned} \frac{\partial E}{\partial \mathbf{d}} &= -2\mathbf{W}^T \mathbf{y} + 2\mathbf{W}^T \mathbf{W} \mathbf{d} = 0 \\ \therefore \mathbf{d} &= (\mathbf{W}^T \mathbf{W})^{-1} \mathbf{W}^T \mathbf{y} \end{aligned}$$

- 事前に $\mathbf{R} = (\mathbf{W}^T \mathbf{W})^{-1} \mathbf{W}^T$ を計算しておけば、

$$\mathbf{d} = \mathbf{R}\mathbf{y}$$

で一瞬で求まる

キーワード検索の方法 (5)

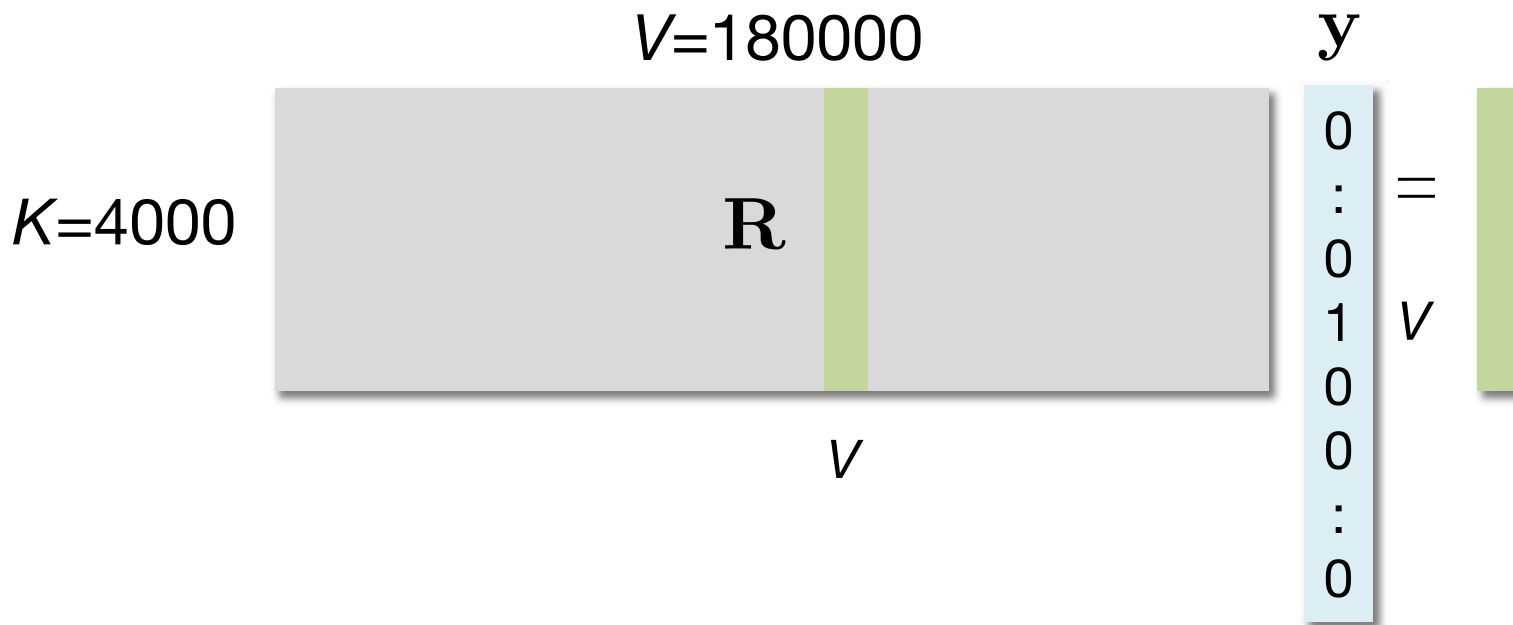
- $\mathbf{R} = (\mathbf{W}^T \mathbf{W})^{-1} \mathbf{W}^T$ は、`numpy.linalg.solve (dot(W.T,W), W.T)` で求められる
- W は $V \times K$ なので、 $W^T W$ は $K \times K$, $R = (W^T W)^{-1} W^T$ は $(K \times K) \times (K \times V) = K \times V$ の行列
 - 言語処理学会では $V=18000$, $K=2000 \rightarrow R$ は 277MB
 - 学振の場合は $V=180000$, $K=4000$ なので R (および W) は 要素数 7億2000万個の巨大な行列、 $\sim 4\text{GB}$

$V=180000$

$K=4000$

\mathbf{R}

キーワード検索の方法 (6)



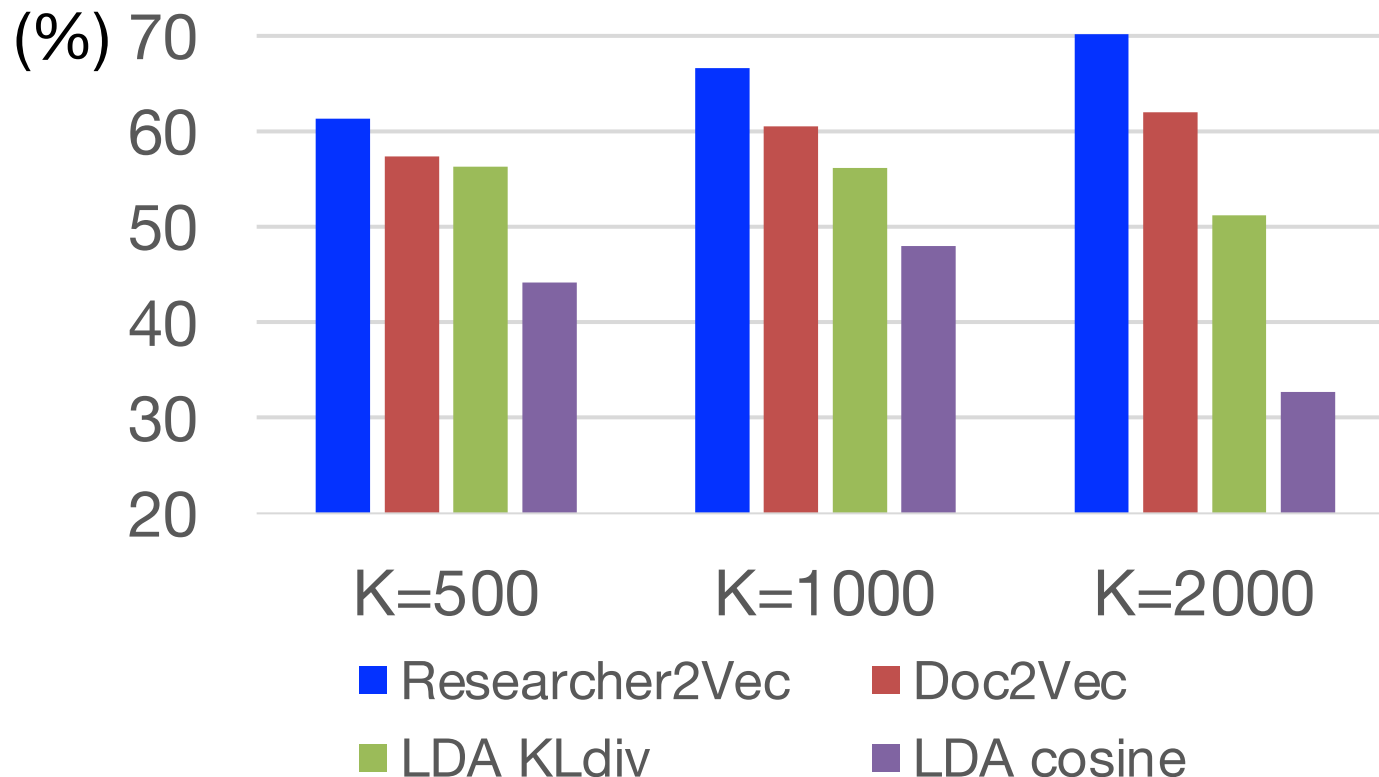
- y の要素はほとんど0なので、 Ry の掛け算は結局、 R の対応する列を取り出してくればよい
- R はディスクにmmap()しておけば、メモリ使用も最小(対応済み)

実験とデータ

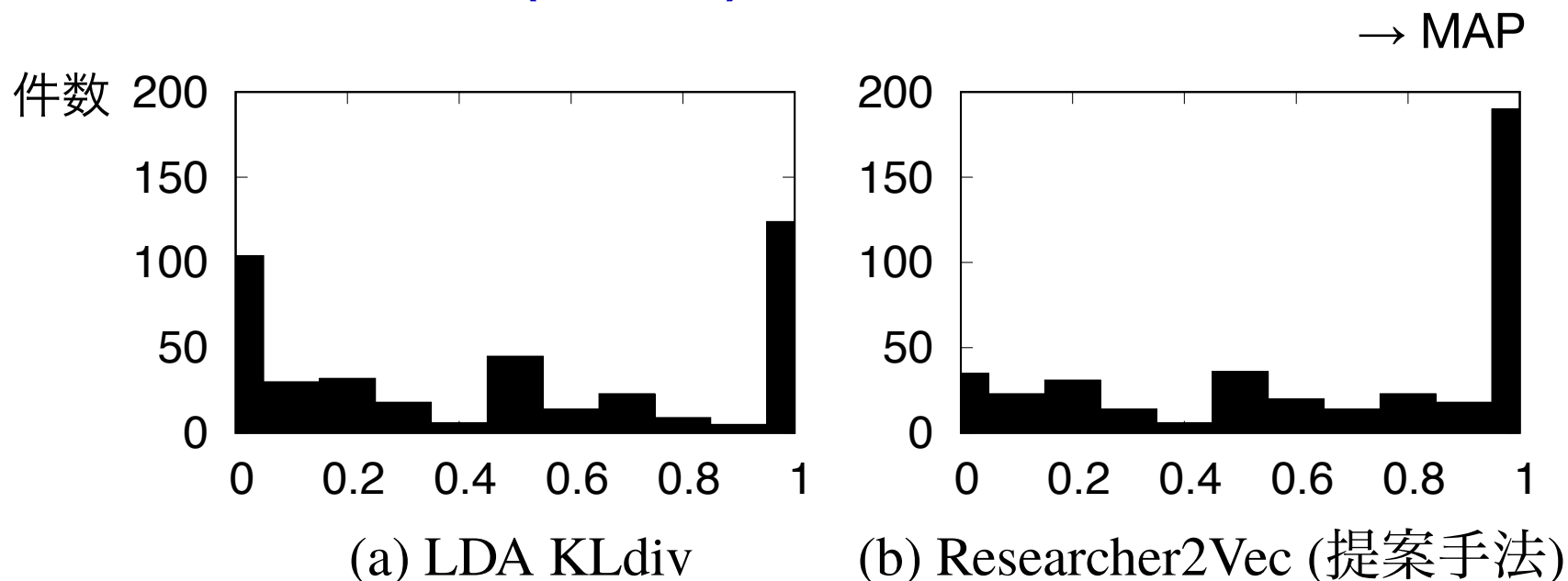
- 言語処理学会年次大会の20年分の論文データ (1995-2013) を使って実験
 - 20周年記念コーパスなので、ニューラル以前なのに注意
- MeCabで形態素解析し、語彙18,135語、論文4,082本で13,654,061語のデータ (1300万単語)
- この期間で5本以上の論文がある著者499人/3660人を実験の対象
- テストデータ500文書の著者を推定し、スコア順に並べた際の平均適合率 (Mean Average Precision, MAP)を計算
 - $MAP=1$: スコア最上位がすべて真の著者
 - Doc2Vec, LDAと比較

論文著者の推薦精度 (平均適合率)

- Doc2Vecを超えて、提案法が常に最高精度
 - LDAは、桂井ら(2016)の確率分布のコサイン類似度よりKLダイバージェンスで測った方がよい



平均適合率(MAP)の分布



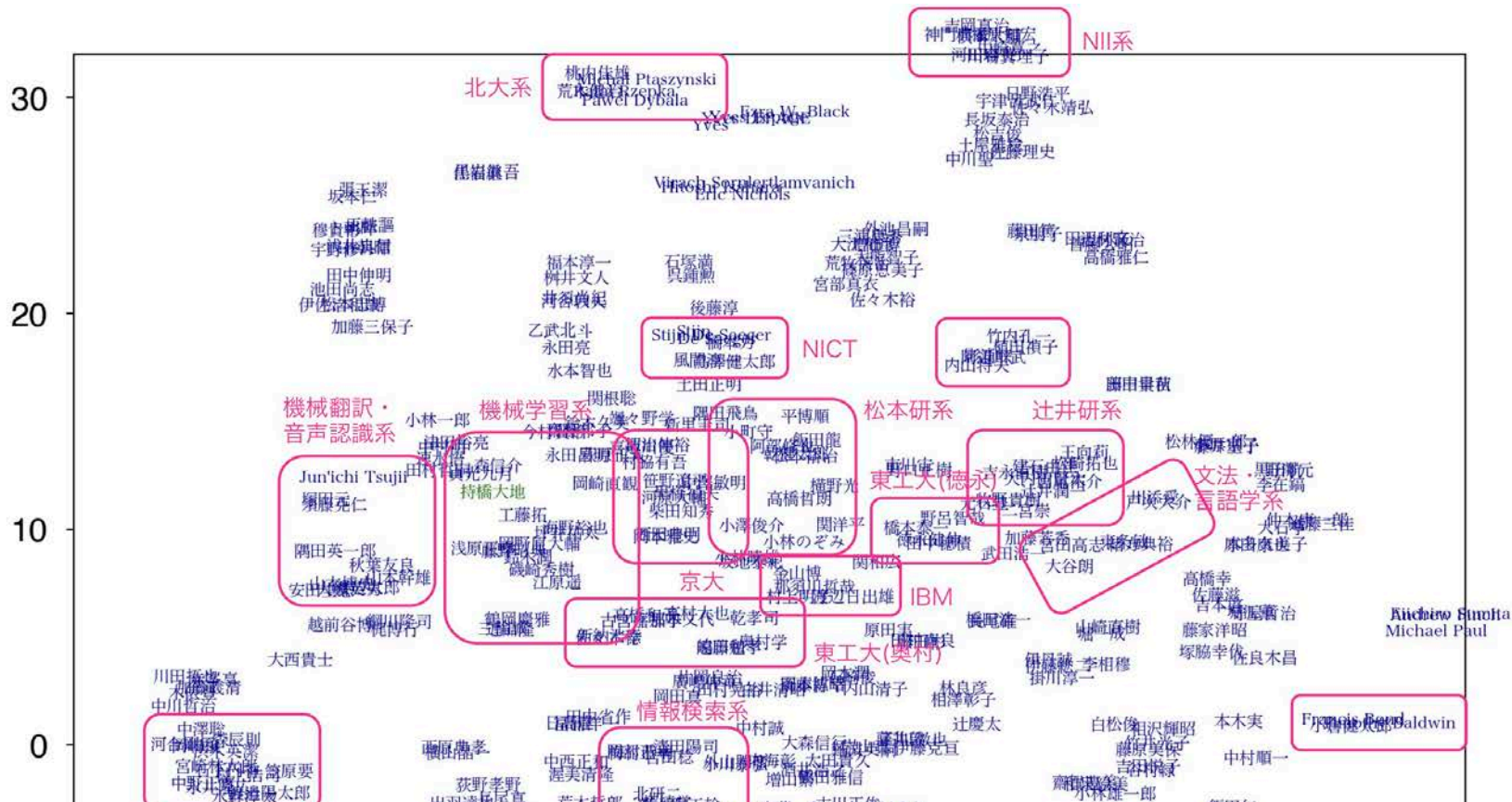
- 提案手法では、ほとんどの場合に平均適合率=1
→ 著者をスコア順に並び替えたとき、**真の著者が最上位を占める**
- それ以外は、先生と興味が違う学生の論文や、英語論文

デモ

<http://clml.ism.ac.jp/nlp2vec/>

- 公開サーバですので、誰でも検索を試すことができます
 - 言語処理学会の過去の論文のビューアにもなります
(著者別)

研究者の可視化



- t-SNEによる可視化 (可視化は1通りではありません)
- 詳しくは論文のAppendixを参照



将来の課題

- 今回は直感の働く自然言語処理分野の論文で試しただけで、他の分野へも適用可能
- arXivやACL anthologyの英語論文への適用
 - 査読者や講演者を探すのに実際困っている!
- Web周りが得意で(非常に)役に立つシステムを作りたい学生さんなどを募集しています
- 埋め込み空間での分布全体を考慮した研究者推薦(kernel mean embedding?)
- 「ある分野全体をカバー」する研究者集合の求め方

まとめ

- 論文からWord2vecと同等のニューラル文書ベクトルをSVDで高速に求め、それを書いた研究者の“研究者ベクトル”を計算する手法を提案
 - Doc2Vecの20～40倍高速、解析的な検索解、省メモリ
 - Doc2VecおよびLDAを超えて最高精度
- 実際に言語処理学会の年次大会コーパスから、論文に含まれる単語で研究者を検索できるシステムを公開(nlp2vec)
- 今後は、日本学術振興会で同様のシステムを開発したい