

方言解析のための 空間的インド料理店過程

持橋大地


統計数理研究所 / 国立国語研究所

daichi@ism.ac.jp

次世代言語科学研究センター 開所式
2025-3-3(月)

本研究は、菅澤翔之助氏(慶応大学経済学部)との共同研究です。

自己紹介

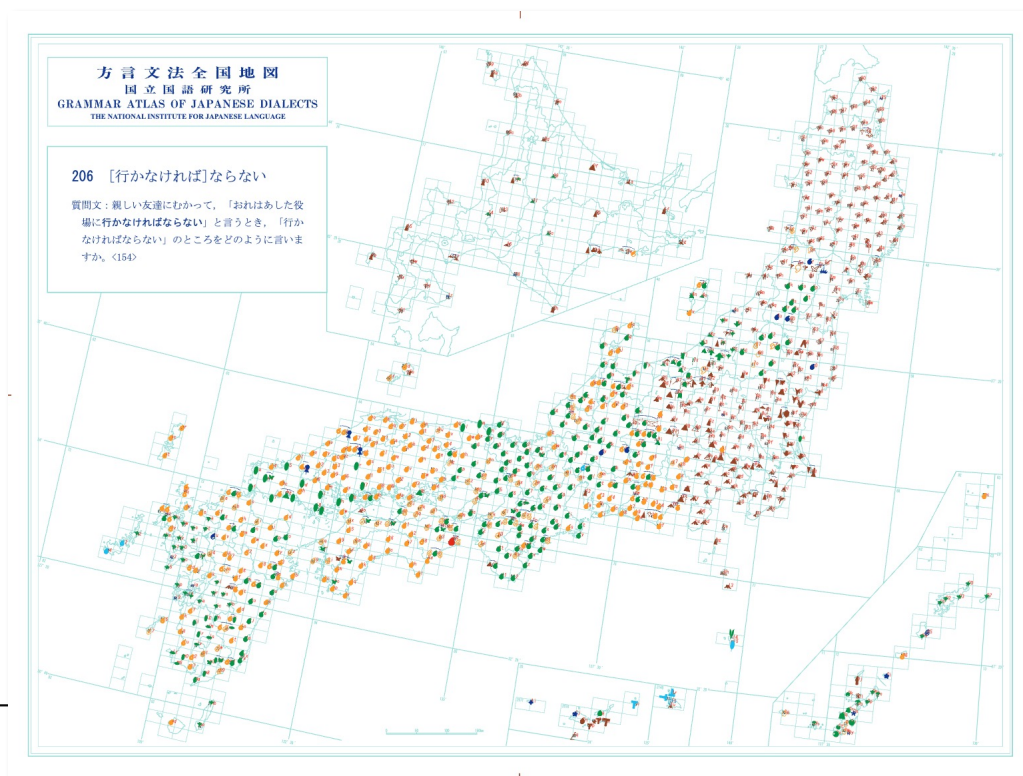
- 持橋大地
統計数理研究所 統計基盤数理研究系 教授
国立国語研究所 次世代言語科学研究センター 教授
(クロスアポイントメント)  国立国語研究所
次世代言語科学研究センター
- 専門: 自然言語処理、統計的機械学習



統数研:
唯一の統計学の
国立研究所
(立川市、国語研の隣)

方言のデータ化と分析

- 「日本言語地図」(1966年～), 「新日本言語地図」(2016年), 「方言文法全国地図」(1989年～)など、国立国語研究所が古くからデータ化と研究を行ってきた



言語地理学の研究を行うきっかけ

The screenshot shows a web browser window displaying the KAKEN project page. The browser address bar shows the URL: kaken.nii.ac.jp/ja/grant/KAKENHI-PROJECT-18KK0012/. The page header includes the KAKEN logo and navigation links for '研究課題をさがす' and '研究者をさがす'. The main content area displays the project title '時空間を融合する：GISと数理モデルを用いた新たな言語変化へのアプローチ' and a list of project details.

研究課題/領域番号	18KK0012
研究種目	国際共同研究加速基金(国際共同研究強化(B))
配分区分	基金
審査区分	中区分2:文学、言語学およびその関連分野
研究機関	国立民族学博物館
研究代表者	菊澤 律子 国立民族学博物館, 人類基礎理論研究部, 教授 (90272616)
研究分担者	村脇 有吾 京都大学, 情報学研究科, 講師 (70616606) 持橋 大地 統計数理研究所, 数理・推論研究系, 准教授 (80418508) 吉岡 乾 国立民族学博物館, 人類基礎理論研究部, 准教授 (20725345) 佐野 文哉 京都大学, アジア・アフリカ地域研究研究科, 客員研究員 (00965501)
研究期間 (年度)	2018-10-09 – 2024-03-31

- 国立民族学博物館 菊澤さんリーダーの国際科研費プロジェクト (2018～, 2024からも新しく採択)

フィジーの場所



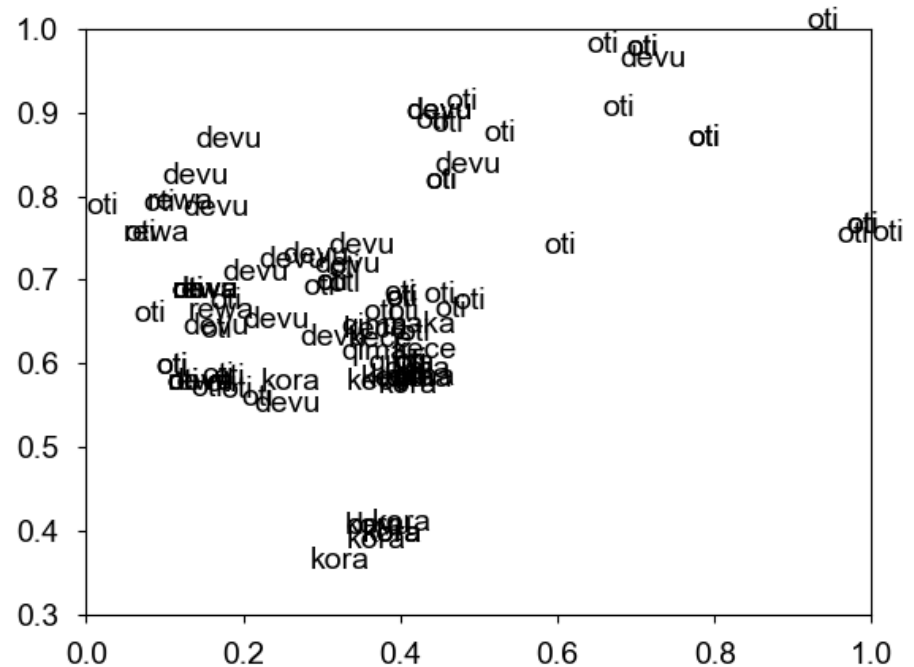
フィジー全体の地図



フィジー語の方言

180.58	39.84	oti devu
182.76	39.67	oti rewa
183.34	39.17	oti
183.90	39.86	oti
184.13	39.87	rewa devu
184.61	38.85	oti devu
185.01	40.03	devu
185.27	38.75	oti devu rewa
185.54	39.32	oti rewa devu
186.13	39.09	devu
186.16	39.83	devu
186.46	39.19	rewa oti
186.66	38.71	oti devu
186.94	40.25	devu
187.21	39.07	oti
187.24	38.79	oti
187.42	38.74	oti

- どうやって、こうした空間的変異を扱うか？



プロジェクトで作成したデータ

- 100 Word Lists : 100個の意味 x 村の方言の行列
 - 空間的な離散データ、列ごとに語彙が異なる

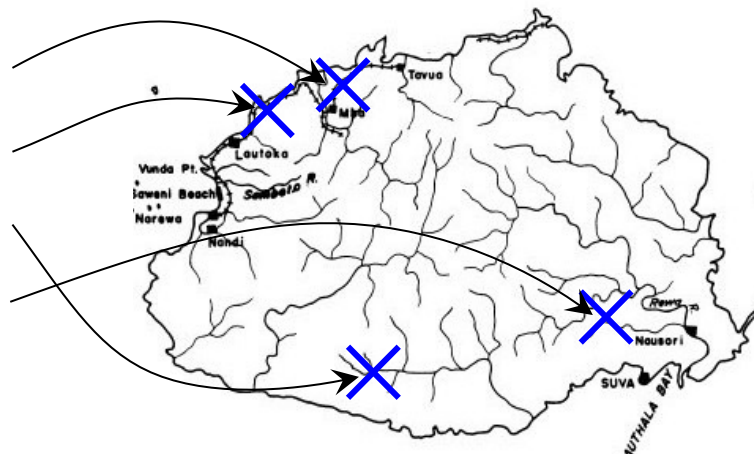
	A	B	C	D	E	F	AI	AJ	AK	AL	AM	AN	AO	AP	AQ	AR											
1	OBJECTID	VillagesComCodes	VillagesComCor	VillagesCor	ComCodeGrp	ComGrou	e	at	in	na	the	nō	gen	N	oqō	tl	oqori	th	oyā	that	iko	you	koya	hin	keda	us	keimami
2	1	Tuvana village			Ono	Lau	i	a	we	kaiīiei	xxx	xxx	iko	koikoyaiko	kīketaketa	keimami											
3	2	Matavualevu Settlement			Ono	Lau	i	a	we	kaiīiei	xxx	xxx	iko	koikoyaiko	kīketaketa	keimami											
4	3	Matokana village			Ono	Lau	i	a	we	kaiīiei	xxx	xxx	iko	koikoyaiko	kīketaketa	keimami											
5	4	Nukuni village			Ono	Lau	i	a	we	kaiīiei	xxx	xxx	iko	koikoyaiko	kīketaketa	keimami											
6	5	Lovoni village			Ono	Lau	i	a	we	kaiīiei	xxx	xxx	iko	koikoyaiko	kīketaketa	keimami											
7	6	Doi village			Ono	Lau	i	a	we	kaiīiei	xxx	xxx	iko	koikoyaiko	kīketaketa	keimami											
8	7	Vatoo village		LAU	Laucake	Lau	i	a	o-	iei, ī	iqore, iqoya	maiei, maīi	iko	koikoya	keiketa, ke	keimami											
9	8	Ogea Village		LAU	Laucake	Lau	i	a	o-	iei, ī	iqore, iqoya	maiei, maīi	iko	koikoya	keiketa, ke	keimami											
10	9	Muanaicake Village			Laucake	Lau	i	a	o-	iei, ī	iqore, iqoya	maiei, maīi	iko	koikoya	keiketa, ke	keimami											
11	10	Muanaicake Village			Laucake	Lau	i	a	o-	iei, ī	iqore, iqoya	maiei, maīi	iko	koikoya	keiketa, ke	keimami											
12	11	Matanuku Village	RAVITAKI	KADAVU	Ravitaki	Kadavu	i	na	nō-	kā	kari	kacei	iko	kia	keda	kēmī											
13	12	Burelevu Village	RAVITAKI	KADAVU	Ravitaki	Kadavu	i	na	nō-	kā	kari	kacei	iko	kia	keda	kēmī											
14	13	Muani Village		KADAVU	Ravitaki	Kadavu	i	na	nō-	kā	kari	kacei	iko	kia	keda	kēmī											
15	14	Levuka Village	NAKUKOLEVU	KADAVU	Nabukelevu	Kadavu				kea	meri	mē	iko	kaia	kēdā												
16	15	Kabariki Village	NAKUKOLEVU	KADAVU	Nabukelevu	Kadavu				kea	meri	mē	iko	kaia	kēdā												
17	16	Nasau Village	NAKUKOLEVU	KADAVU	Nabukelevu	Kadavu				kea	meri	mē	iko	kaia	kēdā												
18	17	Qaliira Village	NAKUKOLEVU	KADAVU	Nabukelevu	Kadavu				kea	meri	mē	iko	kaia	kēdā												
19	18	Nadaviqelevu Village	NAKUKOLEVU	KADAVU	Nabukelevu	Kadavu				kea	meri	mē	iko	kaia	kēdā												
20	19	Naividamu Village			Laucake	Lau	i	a	o-	iei, ī	iqore, iqoya	maiei, maīi	iko	koikoya	keiketa, ke	keimami											
21	20	Cevai Village	TAVUKI	KADAVU	Tavuki	Kadavu	i	na	nō-	kā	kari	kacei	iko	kia	keda	kēmī											
22	21	Namanusa Village	RAVITAKI	KADAVU	Ravitaki	Kadavu	i	na	nō-	kā	kari	kacei	iko	kia	keda	kēmī											
23	22	Nabukelevuirua Village	NAKUKOLEVU	KADAVU	Nabukelevu	Kadavu				kea	meri	mē	iko	kaia	kēdā												
24	23	Lomati Village	NAKUKOLEVU	KADAVU	Nabukelevu	Kadavu				kea	meri	mē	iko	kaia	kēdā												
25	24	Talaulia Village		KADAVU	Nabukelevu	Kadavu				kea	meri	mē	iko	kaia	kēdā												
26	25	Mokoisa Village		KADAVU	Ravitaki	Kadavu	i	na	nō-	kā	kari	kacei	iko	kia	keda	kēmī											
27	26	Dagai Village		KADAVU	Nabukelevu	Kadavu				kea	meri	mē	iko	kaia	kēdā												
28	27	Wailevu Village	RAVITAKI	KADAVU	Ravitaki	Kadavu	i	na	nō-	kā	kari	kacei	iko	kia	keda	kēmī											
29	28	Baidamudamu Village	TAVUKI	KADAVU	Tavuki	Kadavu	i	na	nō-	kā	kari	kacei	iko	kia	keda	kēmī											
30	29	Waisomo Village		KADAVU	Tavuki	Kadavu	i	na	nō-	kā	kari	kacei	iko	kia	keda	kēmī											
31	30	Natumua Village	TAVUKI	KADAVU	Tavuki	Kadavu	i	na	nō-	kā	kari	kacei	iko	kia	keda	kēmī											

問題

- 村ごとに100個の概念に対する方言のデータから、**潜在的なパターン**を見つけたい
- 全体を単にグループ化するのではなく、各村は、**幾つかの未知の潜在的な因子**を持っていると仮定する

因子

	1	2	3	4	5	6	...
村 1	■	□	■	□	□	■	□
村 2	■	■	□	□	□	■	□
村 3	□	□	■	■	□	□	□
村 4	□	■	□	□	□	■	■

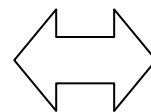


潜在特徴モデルとクラスタリング

- こうしたモデルは、潜在特徴モデルとよばれる
 - クラスタリングでは、各村は1つのグループにだけ所属
 - 潜在特徴モデルでは、各村は**複数の因子が「オン」**になっている (組み合わせは 2^K 通り)

潜在特徴モデル

村 1	■	□	■	□	□	■	□
村 2	■	■	□	□	□	■	□
村 3	□	□	■	■	□	□	□
村 4	□	■	□	□	□	■	■

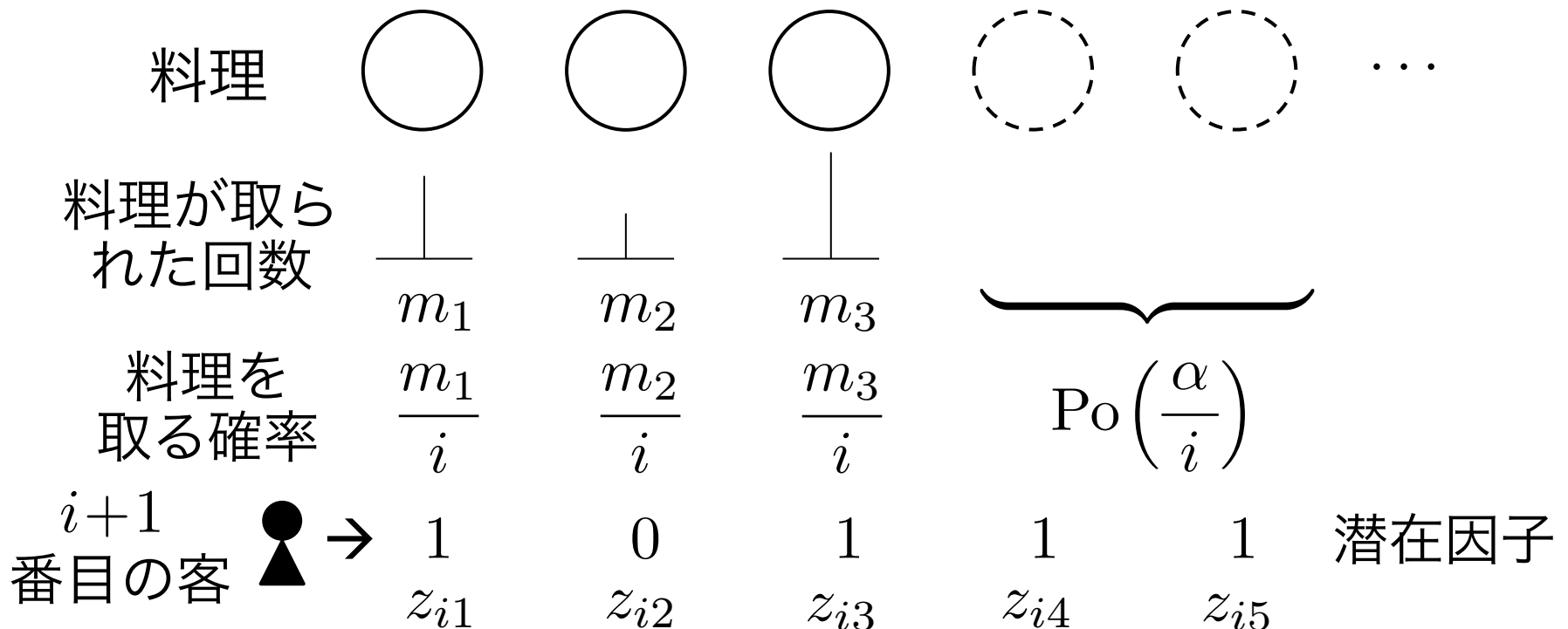


クラスタリング

■	□	□	□	□	□	□	□
□	□	□	□	□	□	□	□
□	□	□	■	□	□	□	□
□	■	□	□	□	□	□	□

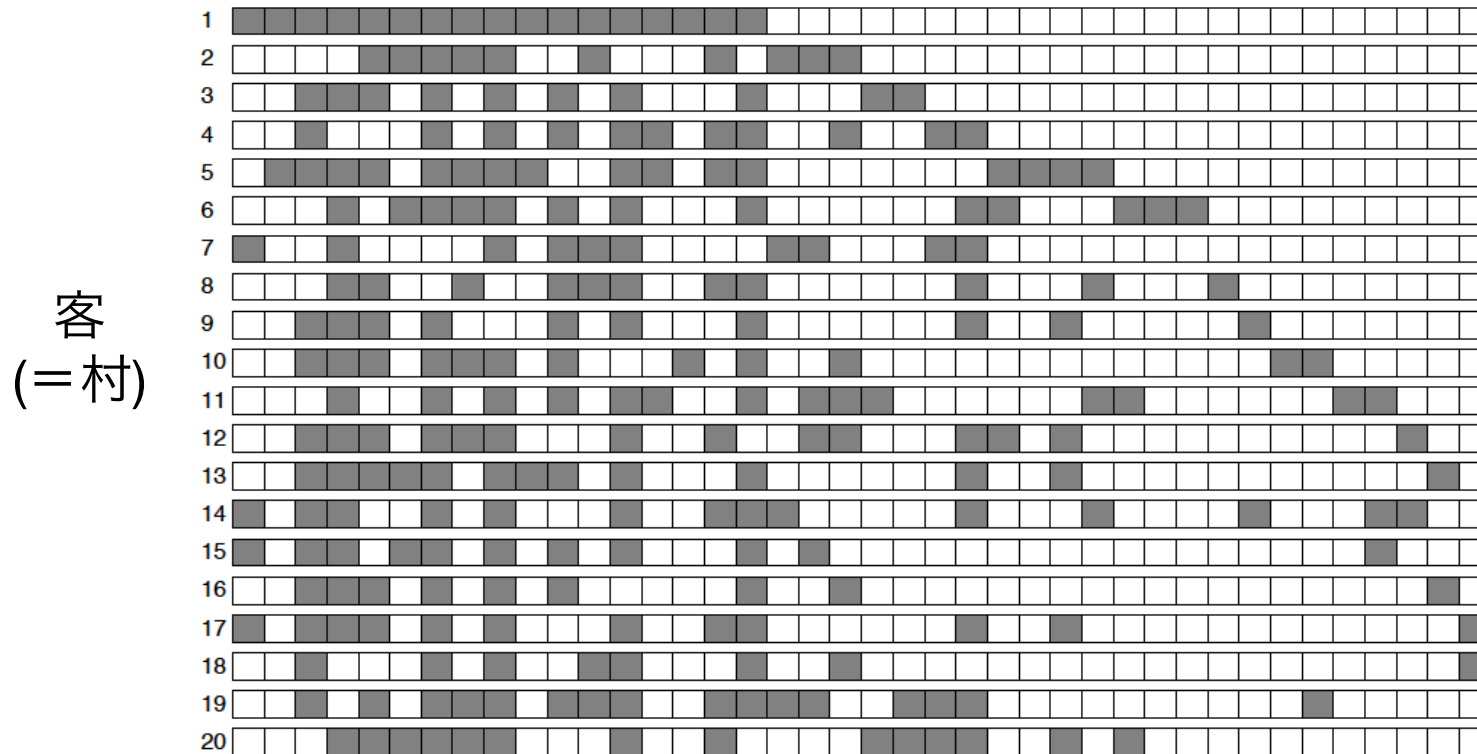
インド料理店過程(IBP) (Griffiths&Ghahramani 2005)

- こうした潜在的因子の生成モデル
→ インド料理店過程 (Indian buffet process, IBP)
- 無限個の特徴を生成できる



IBPで生成された潜在因子の例

潜在因子

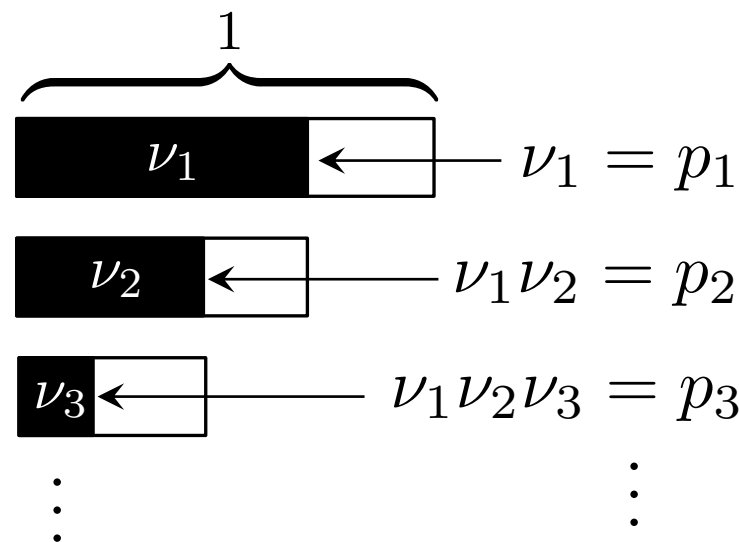
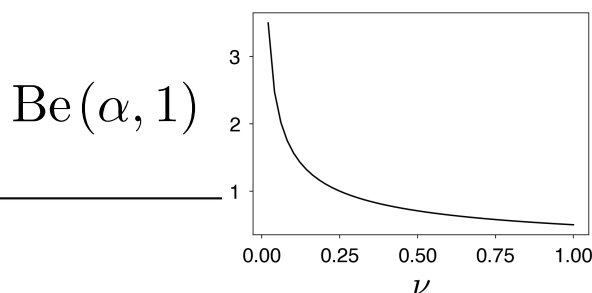


- $\alpha = 10$ のIBPから生成された因子行列
 - この場合は、村は独立 → 地理的相関を入れない

空間的インド料理店過程

- どうやって、地理的に近い客(この場合は村)の間の相関を入れるか?
→ IBPのstick-breaking表現 (Teh+ 2007)
- 長さ1の棒を確率的に折ることで、離散的な確率を生成する (元々はベータ分布を使用)

$$p(z_{ik} = 1) = \prod_{j=1}^k \nu_j,$$
$$\nu_j \sim \text{Be}(\alpha, 1)$$

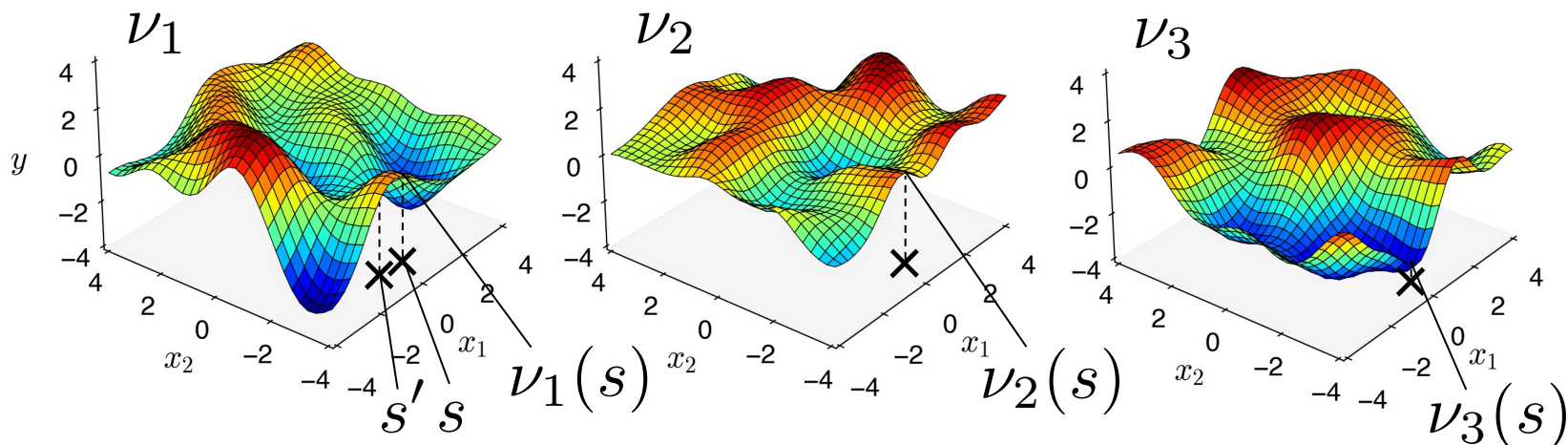


空間的インド料理店過程 (2)

- ロジスティック棒折り過程 (Ren+, JMLR2011)の方法を使うことで、**ガウス過程で棒を折る場所を決定 (=空間相関)**

$$p(z_{ik}(s_i) = 1) = \prod_{j=1}^k \sigma(\nu_j(s_i)), \quad \sigma(x) = \frac{1}{1 + e^{-x}},$$

$$(\nu_j(s_1), \dots, \nu_j(s_N)) \sim \text{GP}(s_1, \dots, s_N)$$



プロジェクトで作成したデータ

- 100 Word Lists : 100個の意味 x 村の方言の行列
 - 空間的な離散データ、列ごとに語彙が異なる

	A	B	C	D	E	F	AI	AJ	AK	AL	AM	AN	AO	AP	AQ	AR											
1	OBJECTID	VillagesComCodes	VillagesComCor	VillagesCor	ComCodeGrp	ComGrou	e	at	in	na	the	nō	gen	N	oqō	tl	oqori	th	oyā	that	iko	you	koya	hin	keda	us	keimami
2	1	Tuvana village			Ono	Lau	i	a	we	kaiīiei	xxx	xxx	iko	koikoyaiko	kīketaketa	keimami											
3	2	Matavualevu Settlement			Ono	Lau	i	a	we	kaiīiei	xxx	xxx	iko	koikoyaiko	kīketaketa	keimami											
4	3	Matokana village			Ono	Lau	i	a	we	kaiīiei	xxx	xxx	iko	koikoyaiko	kīketaketa	keimami											
5	4	Nukuni village			Ono	Lau	i	a	we	kaiīiei	xxx	xxx	iko	koikoyaiko	kīketaketa	keimami											
6	5	Lovoni village			Ono	Lau	i	a	we	kaiīiei	xxx	xxx	iko	koikoyaiko	kīketaketa	keimami											
7	6	Doi village			Ono	Lau	i	a	we	kaiīiei	xxx	xxx	iko	koikoyaiko	kīketaketa	keimami											
8	7	Vatoo village		LAU	Laucake	Lau	i	a	o-	iei, ī	iqore, iqoya	maiei, maī	iko	koikoya	keiketa, ke	keimami											
9	8	Ogea Village		LAU	Laucake	Lau	i	a	o-	iei, ī	iqore, iqoya	maiei, maī	iko	koikoya	keiketa, ke	keimami											
10	9	Muanaicake Village			Laucake	Lau	i	a	o-	iei, ī	iqore, iqoya	maiei, maī	iko	koikoya	keiketa, ke	keimami											
11	10	Muanaicake Village			Laucake	Lau	i	a	o-	iei, ī	iqore, iqoya	maiei, maī	iko	koikoya	keiketa, ke	keimami											
12	11	Matanuku Village	RAVITAKI	KADAVU	Ravitaki	Kadavu	i	na	nō-	kā	kari	kacei	iko	kia	keda	kēmī											
13	12	Burelevu Village	RAVITAKI	KADAVU	Ravitaki	Kadavu	i	na	nō-	kā	kari	kacei	iko	kia	keda	kēmī											
14	13	Muani Village		KADAVU	Ravitaki	Kadavu	i	na	nō-	kā	kari	kacei	iko	kia	keda	kēmī											
15	14	Levuka Village	NAKUKELEVU	KADAVU	Nabukelevu	Kadavu				kea	meri	mē	iko	kaia	kēdā												
16	15	Kabariki Village	NAKUKELEVU	KADAVU	Nabukelevu	Kadavu				kea	meri	mē	iko	kaia	kēdā												
17	16	Nasau Village	NAKUKELEVU	KADAVU	Nabukelevu	Kadavu				kea	meri	mē	iko	kaia	kēdā												
18	17	Qaliira Village	NAKUKELEVU	KADAVU	Nabukelevu	Kadavu				kea	meri	mē	iko	kaia	kēdā												
19	18	Nadaviqelevu Village	NAKUKELEVU	KADAVU	Nabukelevu	Kadavu				kea	meri	mē	iko	kaia	kēdā												
20	19	Naividamu Village			Laucake	Lau	i	a	o-	iei, ī	iqore, iqoya	maiei, maī	iko	koikoya	keiketa, ke	keimami											
21	20	Cevai Village	TAVUKI	KADAVU	Tavuki	Kadavu	i	na	nō-	kā	kari	kacei	iko	kia	keda	kēmī											
22	21	Namanusa Village	RAVITAKI	KADAVU	Ravitaki	Kadavu	i	na	nō-	kā	kari	kacei	iko	kia	keda	kēmī											
23	22	Nabukelevuirua Village	NAKUKELEVU	KADAVU	Nabukelevu	Kadavu				kea	meri	mē	iko	kaia	kēdā												
24	23	Lomati Village	NAKUKELEVU	KADAVU	Nabukelevu	Kadavu				kea	meri	mē	iko	kaia	kēdā												
25	24	Talaulia Village		KADAVU	Nabukelevu	Kadavu				kea	meri	mē	iko	kaia	kēdā												
26	25	Mokoisa Village		KADAVU	Ravitaki	Kadavu	i	na	nō-	kā	kari	kacei	iko	kia	keda	kēmī											
27	26	Dagai Village		KADAVU	Nabukelevu	Kadavu				kea	meri	mē	iko	kaia	kēdā												
28	27	Wailevu Village	RAVITAKI	KADAVU	Ravitaki	Kadavu	i	na	nō-	kā	kari	kacei	iko	kia	keda	kēmī											
29	28	Baidamudamu Village	TAVUKI	KADAVU	Tavuki	Kadavu	i	na	nō-	kā	kari	kacei	iko	kia	keda	kēmī											
30	29	Waisomo Village		KADAVU	Tavuki	Kadavu	i	na	nō-	kā	kari	kacei	iko	kia	keda	kēmī											
31	30	Natumua Village	TAVUKI	KADAVU	Tavuki	Kadavu	i	na	nō-	kā	kari	kacei	iko	kia	keda	kēmī											

方言解析の目的関数

- $\pi_{iml}(s)$: 意味 m を表すのに、 i 番目の村 (座標 s) で方言 l が使われる確率

$$\pi_{iml}(s) = \frac{\exp\left(\eta_{ml} + \sum_{k=1}^K \theta_{kml} z_{ik}(s)\right)}{\sum_{\ell=1}^{C_m} \exp(\dots)}$$

村の潜在的特徴

- 例: Naividamu村 (南緯18.1267°, 東経177.5200°) で、“him/her”を表すのに ‘koikoya’ を使う確率 (他の候補: ‘kia’, ‘kaia’, ‘koya’, ‘kua’)
- K は最大の潜在特徴数 (ここでは $K=10$) で、自動学習
- 観測データ全体の対数尤度を最大化:

$$L = \sum_{i=1}^N \sum_{m=1}^M \sum_{\ell=1}^{C_m} \log \pi_{iml}(s_i)$$

推論

- 目的関数にガウス過程とそこからロジスティック回帰の両方が含まれている
→ 共役ではないので、Pólya-Gamma補助変数法 (Scott 2013) を使用
- 効率的なGibbs samplerが得られる
- 観測点の数が多い場合があるので、最近傍ガウス過程 (Datta+ 2016, JASA) で効率的に近似
- 詳細は、我々のarXiv論文をご参照ください
<https://arxiv.org/abs/2409.01943>

学習されたフィジー語の潜在因子



Factor 1



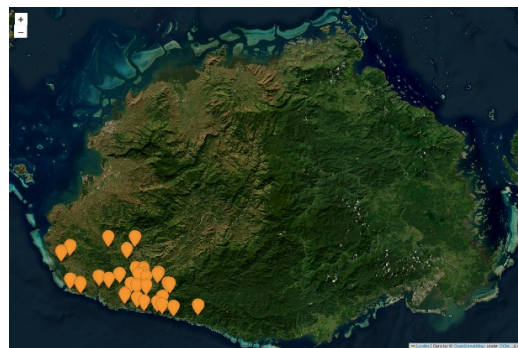
Factor 2



Factor 3



Factor 4



Factor 5



Factor 6



Factors 7,8

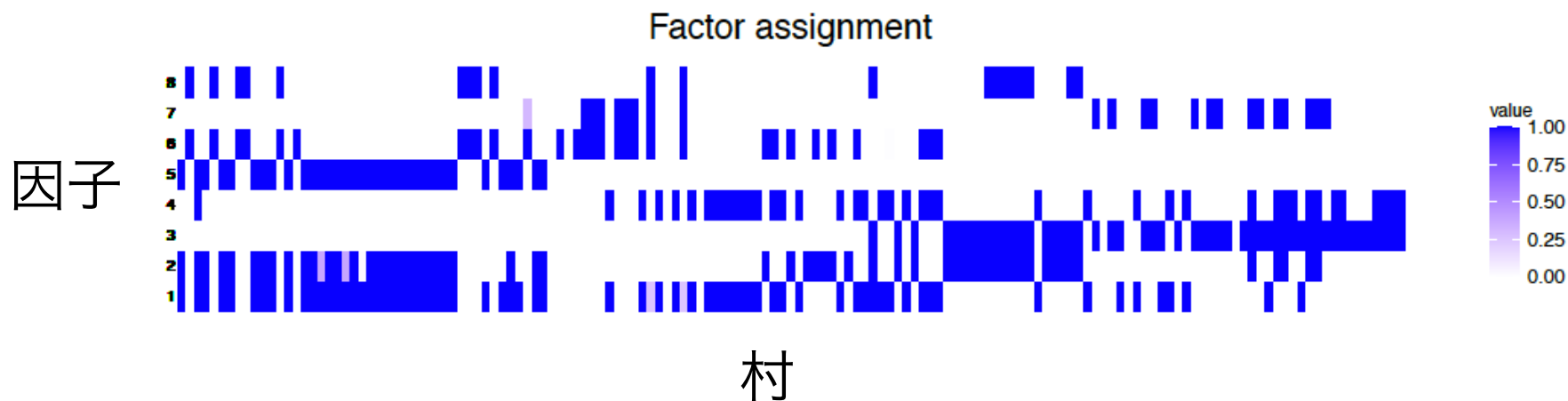
PythonのFolium
を使用してプロット

フィジー語の潜在因子 (全体表示)



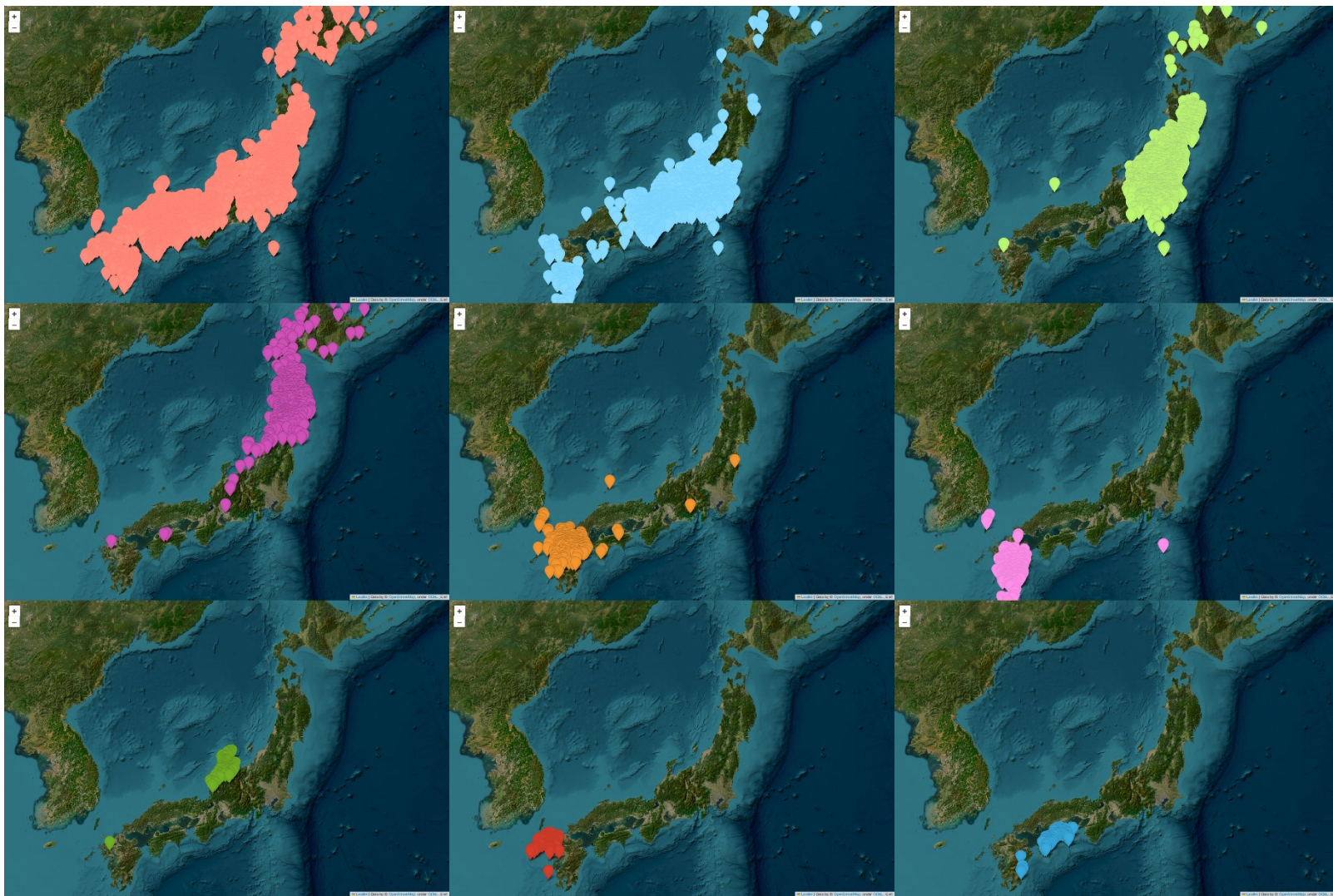
Using ArcGIS tiles in Python!

村と潜在因子の関係



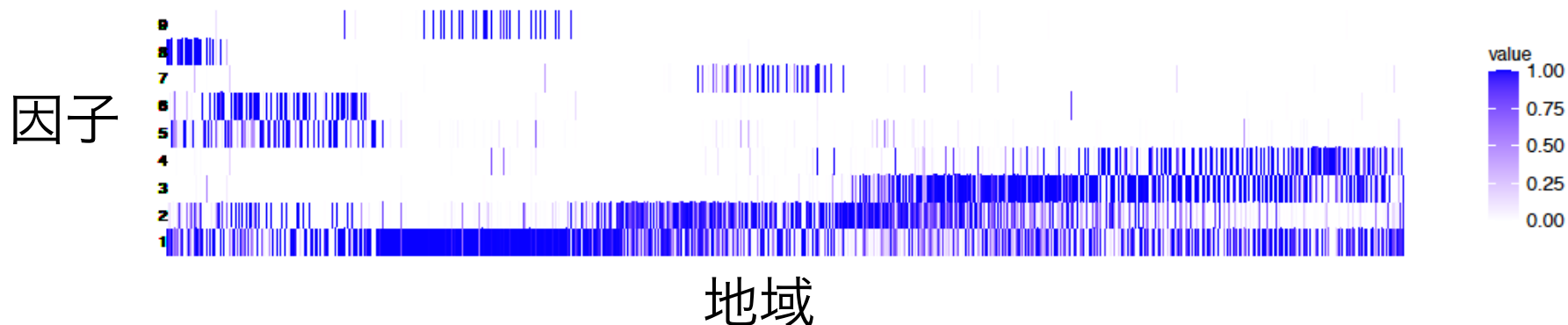
- それぞれの村が、複数の潜在因子を持っている
 - 横軸の村はだいたい地理的に近い順になっているので、空間相関があることがわかる
 - “例外”も少しあるが、それらが複数の村で共有されている

「日本言語地図」で学習された潜在因子



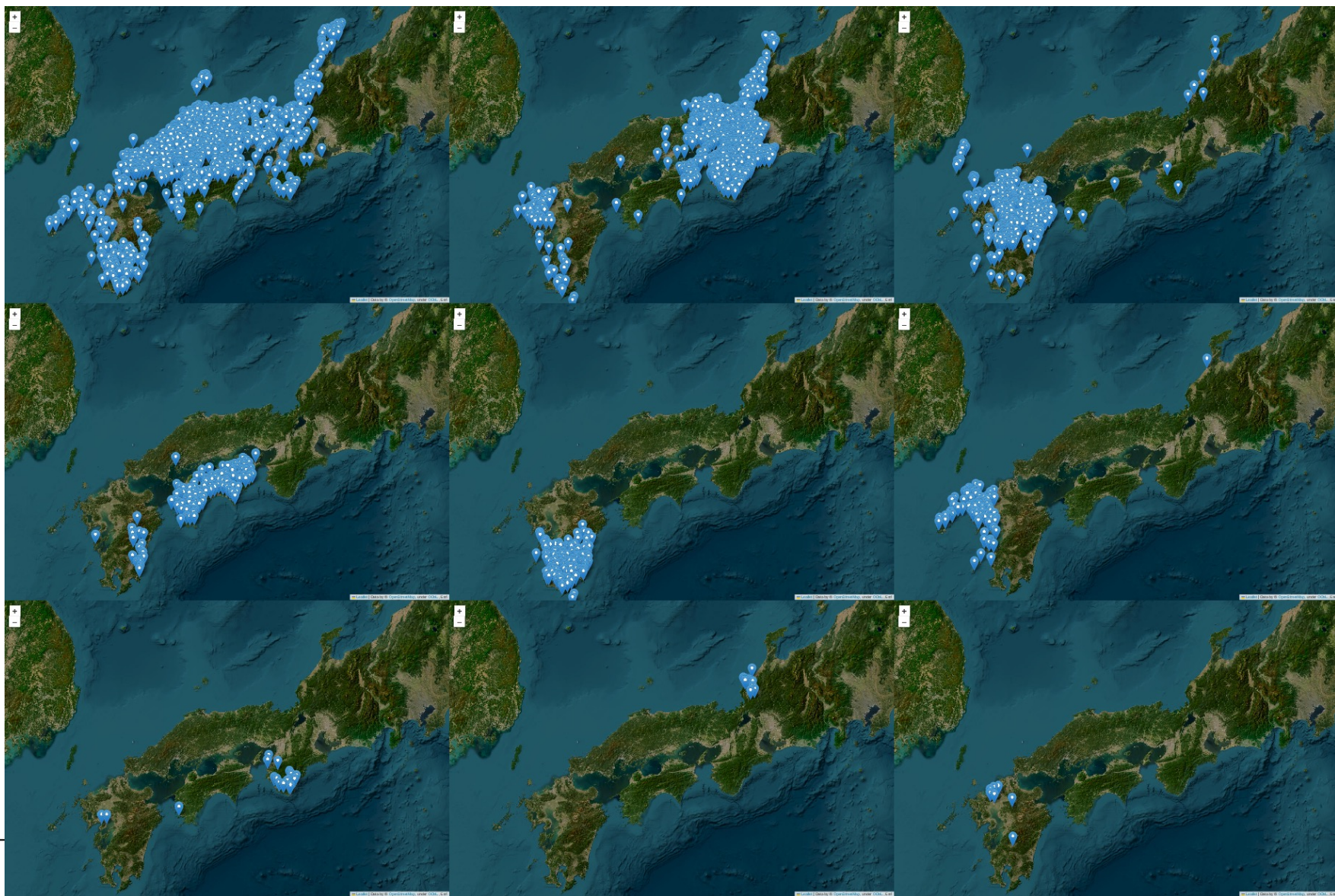
日本の地域と潜在因子の関係

Factor assignment

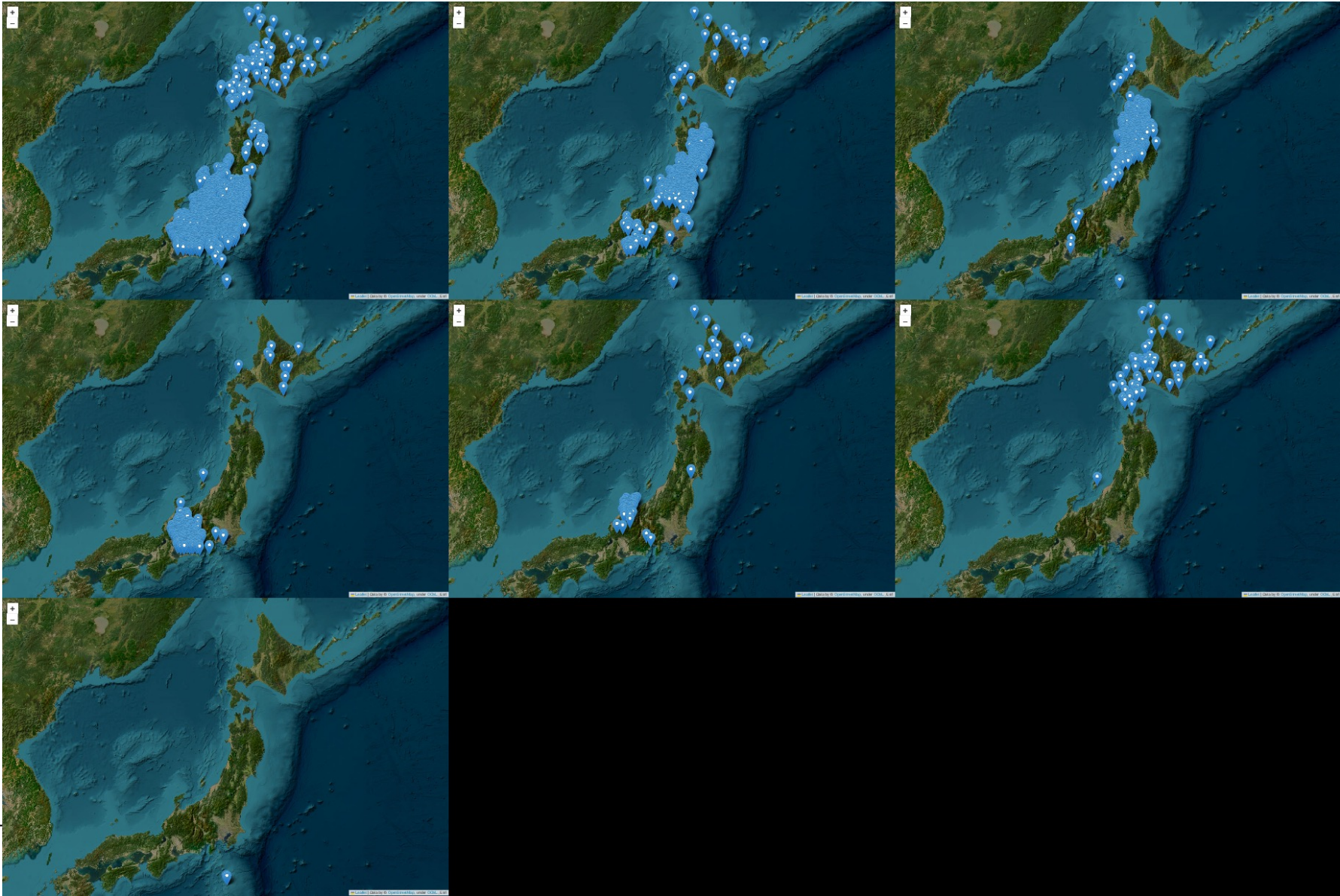


- 日本の各地域が、異なる潜在因子に確率的に所属している
- 植民された北海道、比較的単純な東日本、より歴史のある西日本の様子が読み取れる

「日本言語地図」の分析 (西日本)



「日本言語地図」の分析 (東日本)



まとめ

- 国立国語研究所などでも蓄積されてきた方言のデータを統計的にモデル化するために、空間的インド料理店過程を提案し、潜在的な因子を自動的に推定した
 - IBP+ガウス過程+ロジスティック回帰
- 統計的には、空間相関のある離散データを扱うための空間統計学の拡張
- 得られた因子の解釈には、これまでの言語学での知見と対照する必要がある (今後の課題)

おまけ

安全ではありません - chasen.org/~daiti-m/textmodel

岩波書店 確率と情報の科学シリーズ
『統計的テキストモデル』
持橋大地 著
(2025年発売予定)

News:

- 初校で発見されたバグを修正した, 最終版y版を更新しました. (2025/2/20)

- 岩波書店より、春頃に発売予定
<http://chasen.org/~daiti-m/textmodel/>