

人文社会ビジネス科学学術院 ビジネス科学研究群 2024年度 春C

テキストマイニングの実践

day 3

スケジュール

day 2

- 講義 – 自然言語処理の最新動向
- 講義 – テキストマイニングの手順
- 講義&演習 – データ理解

day 3

- 講義&演習 – 演習環境の準備
- 講義&演習 – テキスト解析 (1)
- 講義&演習 – テキスト解析 (2)

day 4

- 講義&演習 – テキスト分析 (1)
- 演習 – テキスト分析 (実践編)

day 5

- 演習 – テキスト分析 (実践編)

(前回) day 2 – レポート課題

- 以下を PDF ファイルで提出してください
 - データ集計により作成した「集計表」のキャプチャ (P.73~77) ※ページ番号は各スライド右下に記載
 - 作成した「集計結果の整理」の表 (P.78) ※ページ番号は各スライド右下に記載
- ※ 何らかの事情で上記2つを提出できない場合、本日の講義の感想を文章で記述してください

レポート形式	提出先	期限
PDF	manaba	次回～18:20

データ理解 — 集計例

①件数 (エリア別)

行ラベル	個数 / コメン
■ A_レジャー	5000
01_登別	1000
02_草津	1000
03_箱根	1000
04_道後	1000
05_湯布院	1000
■ B_ビジネス	5000
06_札幌	1000
07_名古屋	1000
08_東京	1000
09_大阪	1000
10_福岡	1000
総計	10000

②投稿者の傾向 (年代別x性別)

行ラベル	個数 / コメン	列ラベル	男性	女性	na	総計
10代			0.00%	0.05%	0.00%	0.05%
20代			0.70%	1.16%	0.00%	1.86%
30代			1.94%	2.60%	0.00%	4.54%
40代			4.54%	3.76%	0.00%	8.30%
50代			7.95%	4.17%	0.00%	12.12%
60代			6.34%	1.96%	0.00%	8.30%
70代			1.38%	0.36%	0.00%	1.74%
80代			0.07%	0.06%	0.00%	0.13%
na			0.00%	0.00%	62.95%	62.95%
120代			0.00%	0.01%	0.00%	0.01%
総計			22.92%	14.13%	62.95%	100.00%

③投稿者の傾向 (性別xカテゴリ別)

行ラベル	A_レジャー	B_ビジネス	総計
男性	22.52%	23.32%	22.92%
女性	15.98%	12.28%	14.13%
na	61.50%	64.40%	62.95%
総計	100.00%	100.00%	100.00%

- 男性の投稿者が多い(女性の倍程度) → 男性の観点によるコメントが多い

- 無回答(na)の中の分布が、表明した層と異なる(ある年代や性別に偏っている)可能性もある

データ理解 — 集計例

④投稿者の傾向 (性別xカテゴリーエリア別)

個数 / コメント	列ラベル	A_レジャー 集計					B_ビジネス 集計					総計	
行ラベル		01_登別	02_草津	03_箱根	04_道後	05_湯布院	06_札幌	07_名古屋	08_東京	09_大阪	10_福岡		
男性	A_レジャー	24.10%	22.80%	16.20%	27.10%	22.40%	22.52%	26.40%	25.30%	21.30%	20.50%	23.10%	22.92%
女性	A_レジャー	16.00%	16.00%	17.20%	11.40%	19.30%	15.98%	12.70%	12.90%	11.70%	11.80%	12.30%	14.13%
na		59.90%	61.20%	66.60%	61.50%	58.30%	61.50%	60.90%	61.80%	67.00%	67.70%	64.60%	62.95%
総計		100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%

- 男女差は、レジャーに比べてビジネスが大きい
- 男女差がレジャーで大きいのは道後(次いで登別や草津も大きい)

⑤投稿者の傾向 (年代別xカテゴリーエリア別)

個数 / コメント	列ラベル	A_レジャー 集計					B_ビジネス 集計					総計		
行ラベル		01_登別	02_草津	03_箱根	04_道後	05_湯布院	06_札幌	07_名古屋	08_東京	09_大阪	10_福岡			
10代	A_レジャー	0.10%	0.00%	0.10%	0.00%	0.00%	0.04%	0.00%	0.00%	0.10%	0.20%	0.00%	0.06%	0.05%
20代	A_レジャー	0.80%	3.40%	2.60%	1.60%	1.90%	2.06%	2.00%	2.50%	0.70%	1.60%	1.50%	1.66%	1.86%
30代	A_レジャー	4.20%	5.00%	4.40%	4.00%	6.30%	4.78%	3.40%	4.90%	4.80%	4.30%	4.10%	4.30%	4.54%
40代	A_レジャー	8.30%	9.00%	7.30%	6.70%	10.20%	8.30%	8.20%	8.80%	7.10%	9.20%	8.20%	8.30%	8.30%
50代	A_レジャー	13.40%	11.60%	8.70%	13.30%	11.60%	11.72%	15.10%	12.00%	11.80%	9.80%	13.90%	12.52%	12.12%
60代	A_レジャー	10.70%	8.10%	7.60%	10.20%	9.90%	9.30%	9.00%	8.50%	7.10%	6.10%	5.80%	7.30%	8.30%
70代	A_レジャー	2.30%	1.60%	2.30%	2.70%	1.50%	2.08%	1.20%	1.50%	1.40%	1.00%	1.90%	1.40%	1.74%
80代	A_レジャー	0.30%	0.10%	0.40%	0.00%	0.30%	0.22%	0.10%	0.00%	0.00%	0.10%	0.00%	0.04%	0.13%
120代	A_レジャー	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.10%	0.00%	0.00%	0.00%	0.00%	0.02%	0.01%
na		59.90%	61.20%	66.60%	61.50%	58.30%	61.50%	60.90%	61.80%	67.00%	67.70%	64.60%	64.40%	62.95%
総計		100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%

- あくまでも投稿者の傾向であって、旅行者の実態と一致するは限らない

データ理解 — 集計例

⑥投稿者の傾向 (同行者別)

- レジャーの中で一人が多いのは道後 →道後はもはや仕事で行く場所 (性別でも男性が多い)

- レジャーは家族が多く、ビジネスは一人が多い →出張は複数より単独が多い

個数 / コメント 行ラベル	A_レジャー 集計												B_ビジネス 集計			総計
	01_登別	02_草津	03_箱根	04_道後	05_湯布院	06_札幌	07_名古屋	08_東京	09_大阪	10_福岡						
一人	24.80%	13.70%	14.00%	44.40%	13.10%	22.00%	57.90%	64.30%	65.60%	57.70%	53.50%	59.80%				40.90%
家族	63.10%	64.30%	66.10%	42.70%	69.60%	61.16%	30.90%	24.10%	23.40%	29.40%	32.10%	27.98%				44.57%
恋人	4.80%	14.60%	11.20%	4.90%	8.30%	8.76%	4.40%	4.20%	4.10%	4.40%	3.90%	4.20%				6.48%
友達	5.30%	5.70%	7.10%	5.50%	7.80%	6.28%	4.40%	4.40%	4.80%	6.80%	7.60%	5.60%				5.94%
仕事仲間	1.50%	0.70%	0.50%	1.80%	0.60%	1.02%	1.90%	2.50%	1.30%	1.40%	2.40%	1.90%				1.46%
その他	0.50%	1.00%	1.10%	0.70%	0.60%	0.78%	0.50%	0.50%	0.80%	0.30%	0.50%	0.52%				0.65%
総計	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%				100.00%

⑦数値評価の構成 (総合別)

- 数値評価は、目的によらず高め →好評価しか投稿しない偏りがあるの可能性にも注意

- 高評価5は、レジャーがビジネスよりもやや多い

個数 / コメント 行ラベル	A_レジャー 集計												B_ビジネス 集計			総計
	01_登別	02_草津	03_箱根	04_道後	05_湯布院	06_札幌	07_名古屋	08_東京	09_大阪	10_福岡						
5	41.90%	48.40%	48.90%	49.60%	67.90%	51.34%	41.70%	36.70%	36.90%	41.90%	38.50%	39.14%				45.24%
4	41.30%	36.10%	36.90%	36.50%	22.50%	34.66%	41.90%	47.30%	41.70%	41.10%	40.90%	42.58%				38.62%
3	9.90%	9.90%	7.80%	9.20%	4.90%	8.34%	11.60%	12.00%	14.30%	11.70%	12.60%	12.44%				10.39%
2	4.30%	3.40%	4.30%	2.90%	3.20%	3.62%	3.20%	2.70%	4.60%	3.60%	5.00%	3.82%				3.72%
1	2.60%	2.20%	2.10%	1.80%	1.50%	2.04%	1.60%	1.30%	2.50%	1.70%	3.00%	2.02%				2.03%
総計	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%				100.00%

- レジャーの高評価5は、湯布院が多く、登別が少ない

- ビジネスの高評価5は、札幌・大阪が多く、名古屋・東京がやや少ない

データ理解 — 集計例

⑧-a 数値評価の平均 (エリア別×数値評価別)

行ラベル	平均 / サービス	平均 / 立地	平均 / 部屋	平均 / 設備・アメニ	平均 / 風呂			
■ A_レジャー	4.25	4.25	4.13	4.05	4.29	4.29	4.29	4.30
01_登別	4.07	4.21	3.95	3.90	4.34	4.08	4.16	
02_草津	4.23	4.22	4.07	3.97	4.32	4.20	4.25	
03_箱根	4.24	4.12	4.18	4.05	4.29	4.33	4.26	
04_道後	4.19	4.41	4.07	4.00	4.03	4.19	4.29	
05_湯布院	4.51	4.28	4.37	4.35	4.46	4.61	4.52	
■ B_ビジネス	3.98	4.30	4.01	3.88	3.74	4.05	4.13	
06_札幌	4.05	4.30	4.09	3.93	3.74	3.99	4.06	
07_名古屋	4.00	4.25	4.04	3.89	3.74	4.06	4.18	
08_東京	3.93	4.38	3.94	3.82	3.70	3.99	4.06	
09_大阪	4.01	4.35	4.05	3.93	3.82	4.06	4.18	
10_福岡	3.93	4.24	3.96	3.84	3.64	4.01	4.07	

- レジャーは、風呂や食事が、設備や部屋に比べて高評価

- 湯布院は、レジャーの中で、軒並み高評価が多い

⑧-b 数値評価の平均 (カテゴリ別×数値評価別)

行ラベル	平均 / サービス	平均 / 立地	平均 / 部屋	平均 / 設備・アメニ	平均 / 風呂	平均 / 食事	平均 / 総合
A_レジャー	4.25	4.25	4.13	4.05	4.29	4.29	4.30
B_ビジネス	3.98	4.30	4.01	3.88	3.74	4.05	4.13

- レジャーもビジネスも立地が評価される
- ビジネスは、立地がその他に比べて高評価

データ理解 — 集計例

⑨-a 数値評価の平均 (20~30代, 性別)

行ラベル	平均 / サービス	平均 / 立地	平均 / 部屋	平均 / 設備・アメニティ	平均 / 風呂	平均 / 食事	平均 / 総合
■ A_レジャー	4.48	4.32	4.35	4.33	4.22	4.46	4.46
男性	4.38	4.23	4.30	4.18	4.38	4.48	4.44
女性	4.55	4.37	4.38	4.42	4.45	4.44	4.47
■ B_ビジネス	4.18	4.36	4.19	4.10	3.93	4.41	4.28
男性	4.16	4.33	4.13	4.10	3.93	4.35	4.23
女性	4.19	4.39	4.24	4.10	3.93	4.45	4.32

- 20~50代はレジャーに対する風呂や食事、サービスの評価が概ね高い

⑨-b 数値評価の平均 (40~50代, 性別)

行ラベル	平均 / サービス	平均 / 立地	平均 / 部屋	平均 / 設備・アメニティ	平均 / 風呂	平均 / 食事	平均 / 総合
■ A_レジャー	4.29	4.33	4.15	4.10	4.32	4.35	4.34
男性	4.23	4.29	4.10	4.05	4.29	4.33	4.32
女性	4.38	4.38	4.21	4.17	4.36	4.37	4.36
■ B_ビジネス	4.00	4.30	4.07	3.91	4.00	4.33	4.28
男性	3.93	4.27	4.00	3.89	4.00	4.33	4.28
女性	4.13	4.36	4.19	4.02	4.02	4.37	4.32

- 年齢が高くなるに連れてレジャー・ビジネスとも評価が厳しくなる

⑨-c 数値評価の平均 (60~80代, 性別)

行ラベル	平均 / サービス	平均 / 立地	平均 / 部屋	平均 / 設備・アメニティ	平均 / 風呂	平均 / 食事	総合
■ A_レジャー	4.23	4.24	4.08	3.98	4.29	4.26	4.31
男性	4.21	4.20	4.08	3.98	4.27	4.27	4.27
女性	4.27	4.36	4.15	3.98	4.26	4.39	4.39
■ B_ビジネス	3.95	4.30	3.98	3.77	3.75	3.79	4.14
男性	3.91	4.24	3.98	3.77	3.75	3.79	4.13
女性	4.19	4.59	4.15	3.96	3.94	4.23	4.21

- 女性は立地に対する評価が高めで、60~80代で顕著（→期待が高い）

(参考) データ理解 — 集計結果の整理

観点	データの特徴	テキスト分析時に注意すべき点
年代別・性別	<ul style="list-style-type: none"> 約60%が年代や性別を表明していない 年代別では、目的によらず40~60代が多い 全体的に男性の投稿者が多い（女性の倍程度） レジャーに比べてビジネス方が男女差が大きい レジャーの中でも男女差が大きいのは道後 	<ul style="list-style-type: none"> レビュー観点がある年代や性別に偏っている可能性 → 偏りあるなら乱暴に一般化せず丁寧に見ていく方がよい 無回答("na")中には分布が異なる可能性 → "na"が多いなら属性による偏りがある可能性 無回答が多い場合は、属性に引っ張られない解釈を心がける
目的別	<ul style="list-style-type: none"> レジャーは家族が多い、ビジネスは一人が多い（出張は単独で宿泊するケースが多い） レジャーの中でも、道後は男性の一人客が多い（道後はもはや仕事で行く場所か） 	<ul style="list-style-type: none"> レビ… レビューの観点がカテゴリと一致していない可能性（道後→仕事）
数値評価 (総合)	<ul style="list-style-type: none"> 旅行目的によらず評価は高め レジャーがビジネスより評価が高め レジャーの中で高評価が多いのは湯布院 小さいのは ビジネスの中で高評価が多い屋がやや少ない 	<ul style="list-style-type: none"> 好評価しか投稿しない → コメントが好評価に偏っている可能性にも注意 旅行目的によって投稿の動機が異なっている可能性 属性に偏りがある場合は、乱暴に一般化しようとせず、属性ごとの分析を心がける
数値評価 (項目ごと)	<ul style="list-style-type: none"> レジャーの評価は、風呂や食 ビジネスの評価は、立地 > その他 レジャーの中で湯布院は軒並み高評価 レジャーもビジネスも立地は高評価（重視している） 	目的によって評価の観点や重みが異なっている可能性
その他	<ul style="list-style-type: none"> あくまでも楽天トラベルの特性であるので、旅行者の傾向として主張するためには別途裏付けが必要 	

演習環境の準備

無償で利用できる機械学習環境

- 近年、機械学習の教育・研究を目的とした研究用ツールが相次いで登場

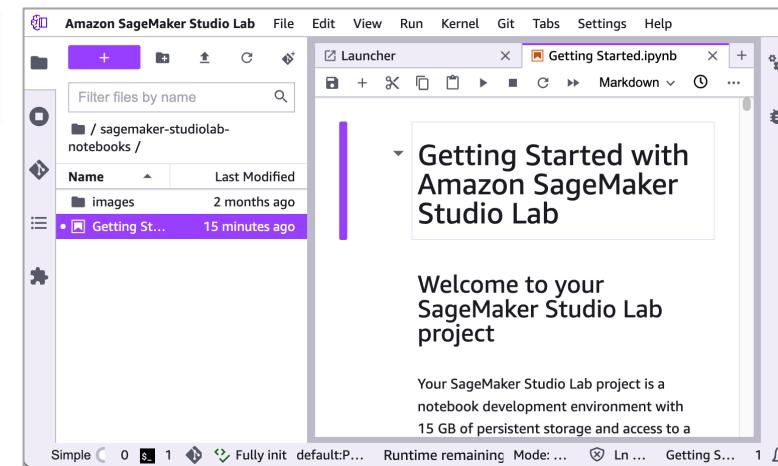
 Colaboratory

<https://colab.research.google.com>



 Amazon SageMaker Studio Lab

<https://studiolab.sagemaker.aws/>



演習で使用
↓

	Colab(無償版)	Studio Lab
GPU	T4(16GB)	T4(16GB)
最長実行時間	12時間	CPU:12時間 GPU:4時間
メモリ	12GB	15GB
ディスク	CPU:100GB GPU:78GB	15GB (永続化)
ターミナル	×	○
ランタイムの保存と再開	×	○
費用	無償	無償
その他	Googleアカウントが必要	AWSアカウントは不要 (クレカ不要)

SageMaker Studio Lab のアカウント作成

- <https://studiolab.sagemaker.aws/> にログインして、アカウントを作成してください

amazon
SageMaker Studio Lab

Sign in Request account

Request account

Request a free Amazon SageMaker Studio Lab account.

Enter your email*

Enter your first name

Enter your last name

Select your country

Enter your company or organization name

Select your occupation

Why are you interested in Amazon SageMaker Studio Lab?

Enter referral code **XXX-XXXX**

Submit request

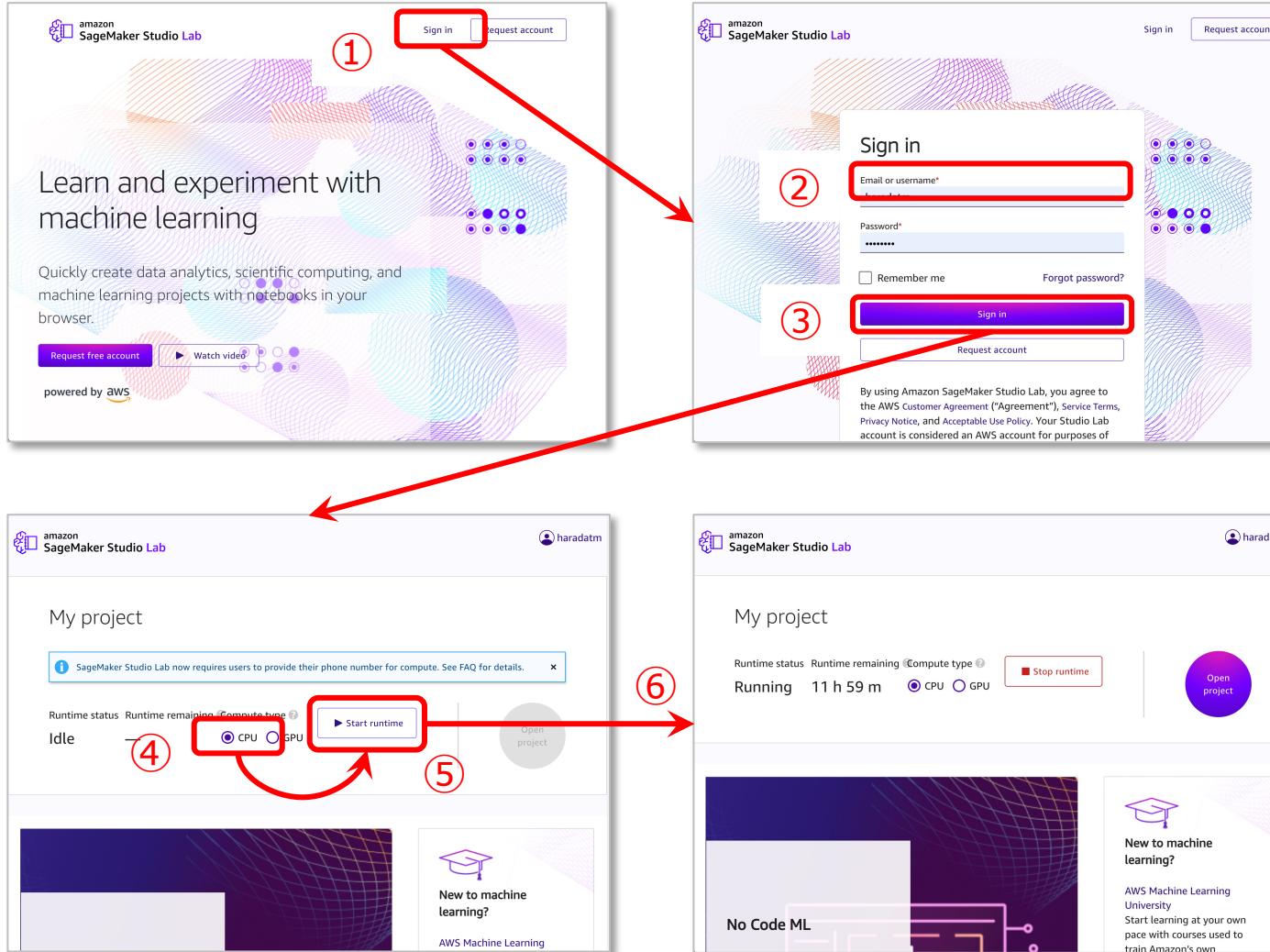
アカウント作成手順

1. [アカウント作成フォーム](#)からアカウントの申し込みを行う
注意: リファラルコードをアカウント作成フォームに忘れずに入力ください (受講者限定です)
2. 「Account request confirmed ...」のメールを受信し、メール内のリンクからアカウントを作成する
→ リクエストの受付はすぐにメールが届きます
3. 「Verify your email ...」のメールを受信し、メール内のリンクからメールアドレスを認証する
→ リファラルコードを利用している場合は2~3分以内に結果が届きます
4. 「Your account is ready ...」のメールを受信する
→ これで「Sign in」できます

※ リファラルコードの有効期間: 2024/6/27 ~ 2024/7/12

演習環境の準備

- <https://studiolab.sagemaker.aws/> にログインして、Runtimeを起動してください



1. ログインする

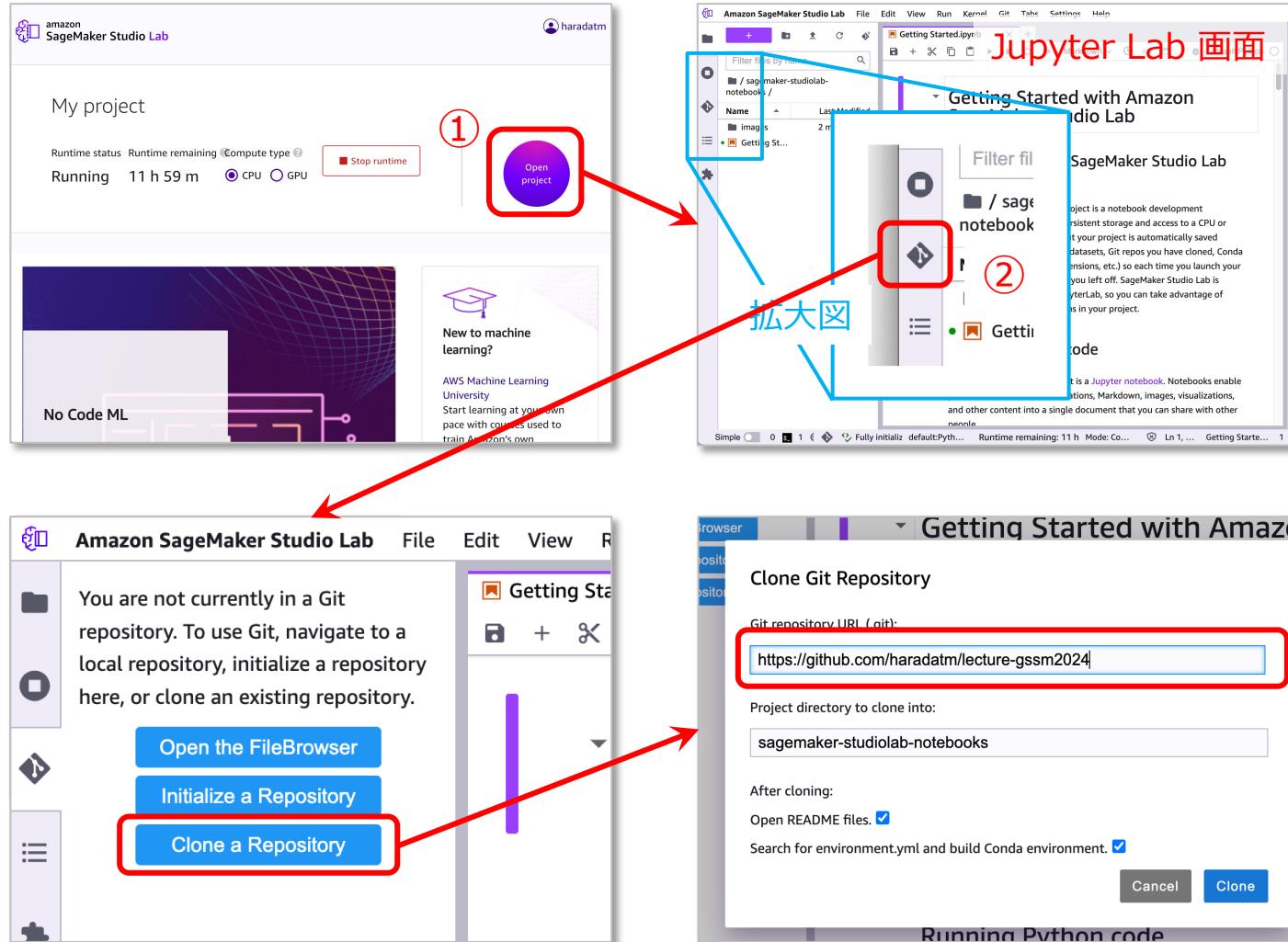
- ① 画面右上の「Sign in」ボタンを押す
- ② Eメールアドレス/ユーザー名、パスワードを入力する
- ③ 「Sign in」を押してプロジェクトのページを開く

2. Runtime を起動する

- ④ 「My Project」の「Select compute type」から「CPU」を選択する
- ⑤ 「Start runtime」を押す
- ⑥ 起動時に多要素認証を求められた場合、使用可能なデバイスで認証する

演習環境の準備

● Jupyter Lab を起動して、教材を開いてください



3. Jupyter Lab を起動する

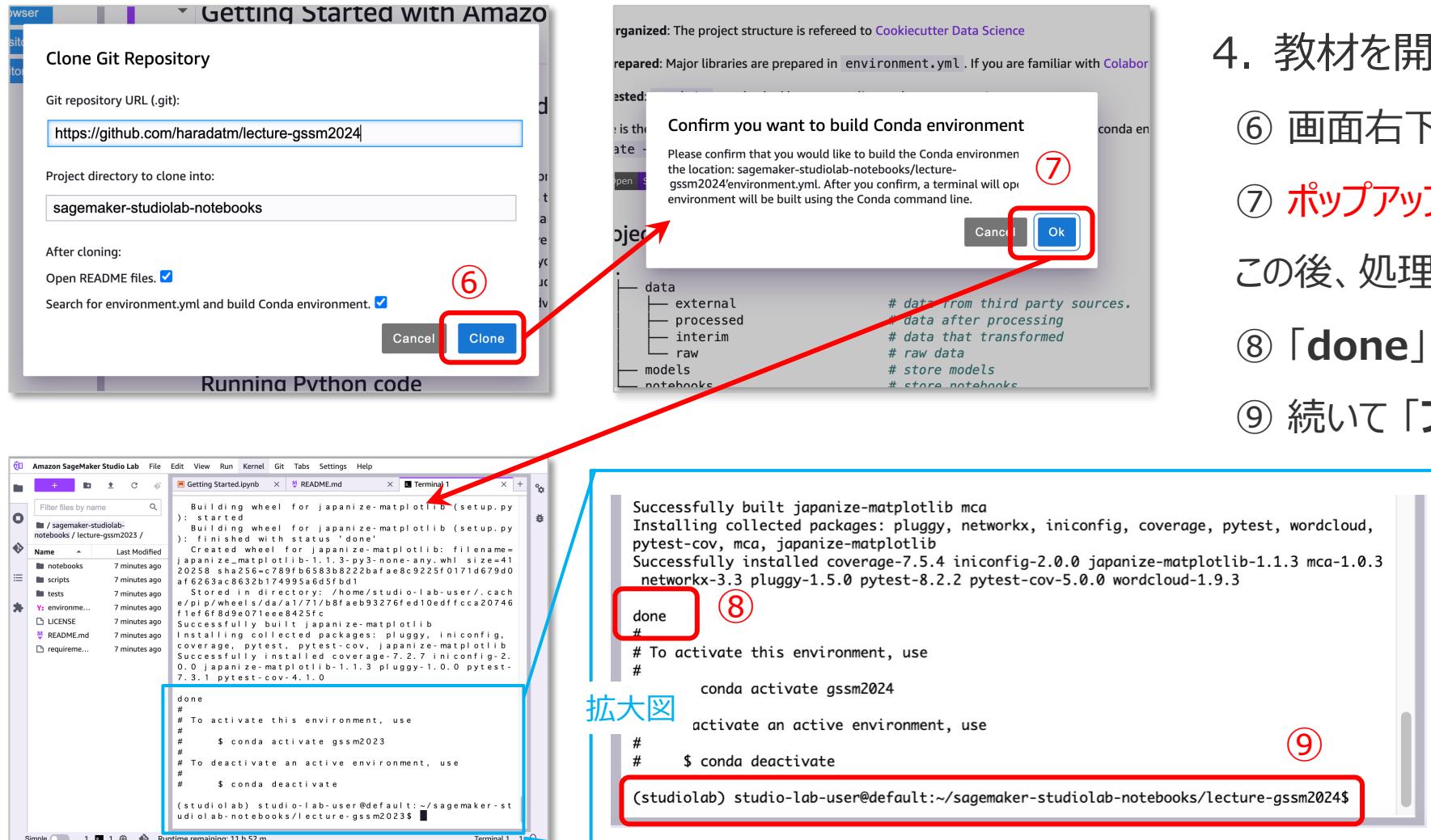
- ① ランタイムが開始したら「Open project」を押す

4. 教材を開く

- ② 「Git」  ボタンを押す
- ③ 「Clone a Repository」を押す
- ④ 「Git repository URL ...」に
<https://github.com/haradatm/lecture-gssm2024> を入力する
- ⑤ (任意) 「Project directory to ...」に保存先のパスを入力する
例) 「sagemaker-studiolab-notebooks」

演習環境の準備

● Jupyter Lab を起動して、教材を開いてください（続き）



4. 教材を開く（続き）

⑥ 画面右下の「Clone」を押す

⑦ ポップアップした画面で「OK」を押す

この後、処理完了まで**7分程度**待ちます

⑧ 「done」が表示される

⑨ 続いて「プロンプト」が表示されれば完了

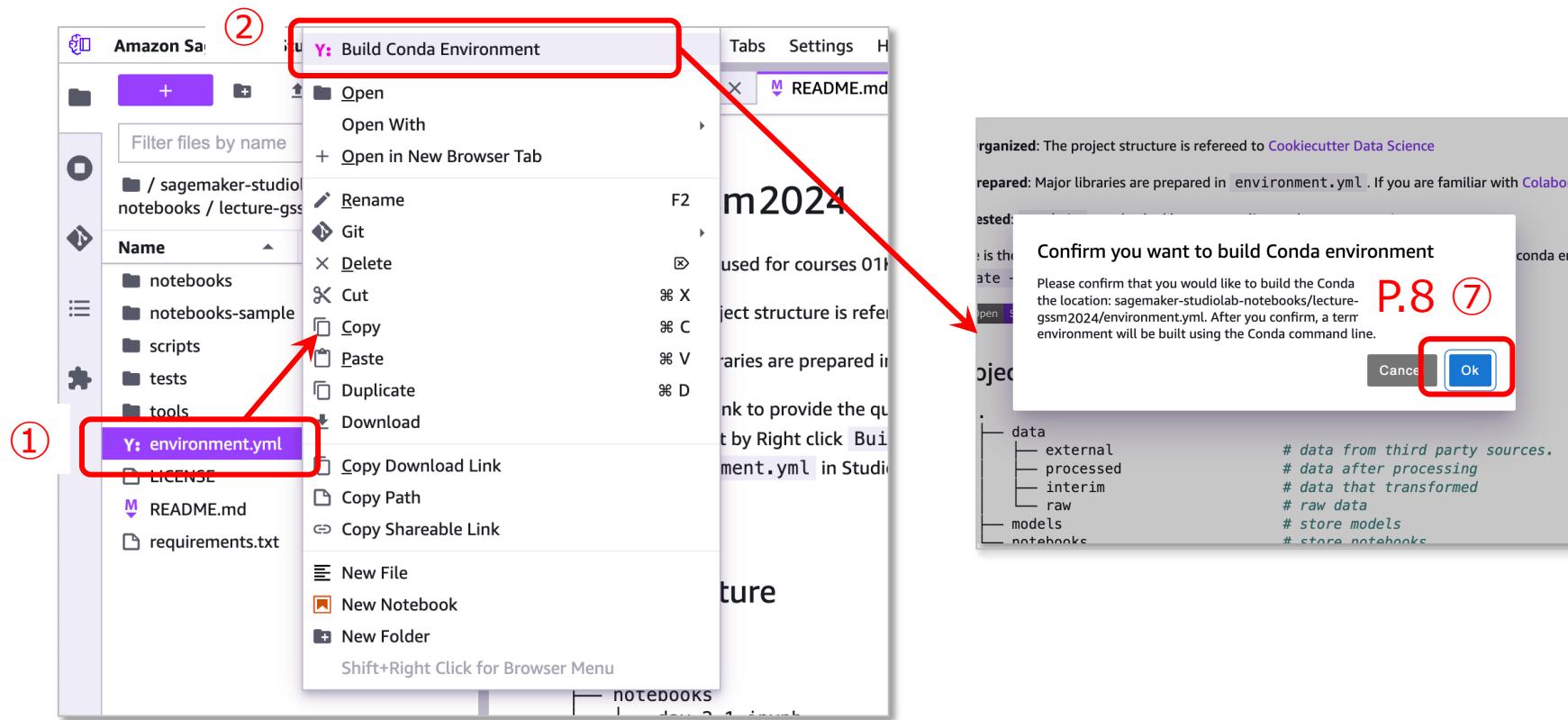
注意:

⑦で「OK」を押さなかった場合、
⑧や⑨の表示がされない場合は
後述の「トラブルシューティング」
から再開してください

(参考) ブラウザで「environment.yml」を選択する

- P.8 の⑦で「OK」を押さなかった場合や、⑧や⑨の表示がされない場合

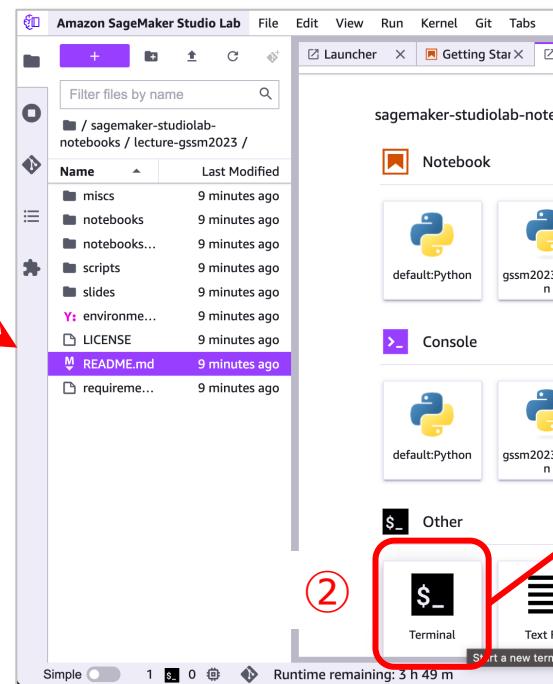
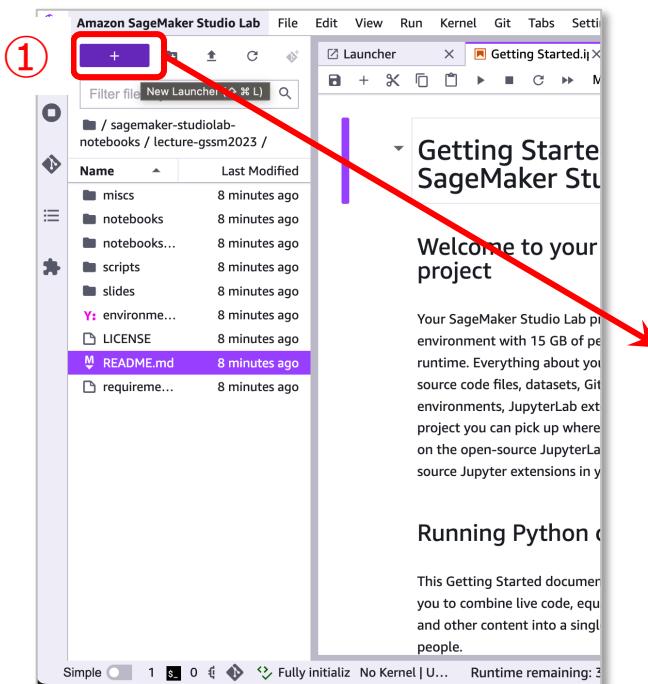
- ① 画面左の **File Browser** から「**environment.yml**」を選択する
- ② 右クリックメニューから「**Build Conda Environment**」を選択する → P.8 ⑦ へ



(参考) ブラウザによる開発環境構築

● (いちからやり直すために) プロジェクトを完全削除する **注意: この操作は慎重に行ってください**

- ① 画面上の **+** ボタンで **Launcher** を開く
- ② 画面上の **Launcher** から **Terminal** を開く
- ③ **Terminal** で右のコマンドを実行する
- ④ **Stop runtime** で Runtime を停止 → P.7 へ

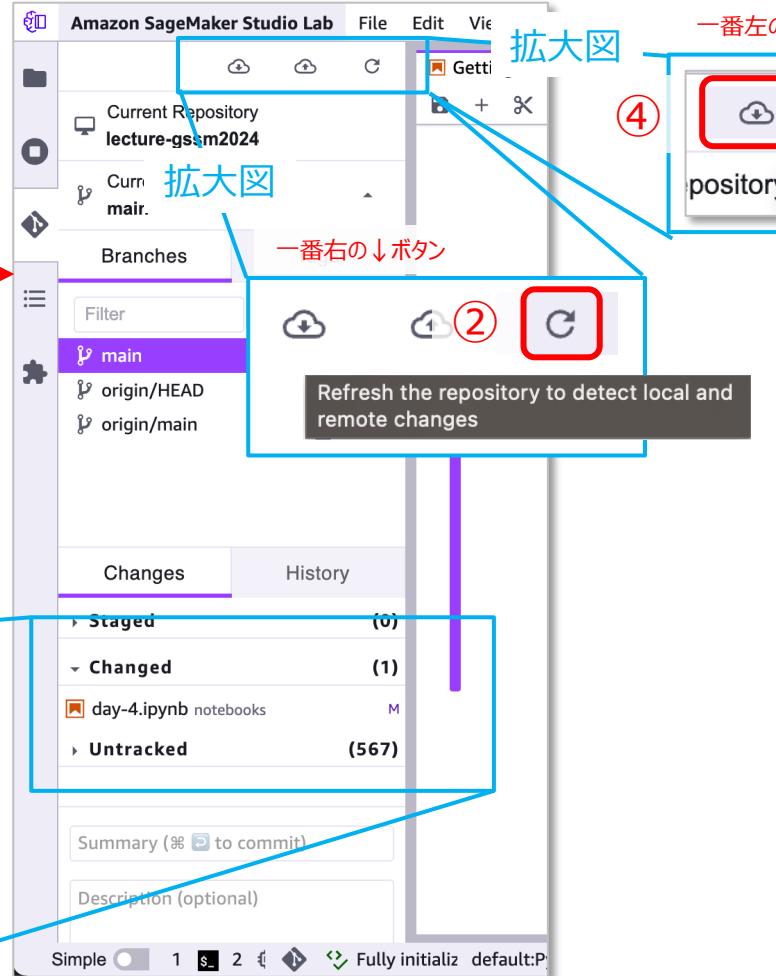
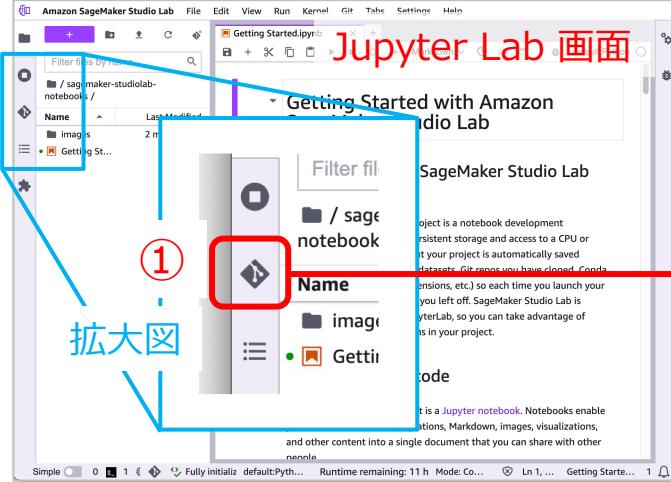


```
# ホームディレクトリへ移動  
cd  
# 仮想環境を削除  
conda remove -n gssm2024 --all # Proceed ([y]/n)? Y  
# プロジェクトを削除  
rm -fr sagemaker-studiolab-notebooks/lecture-gssm2024
```

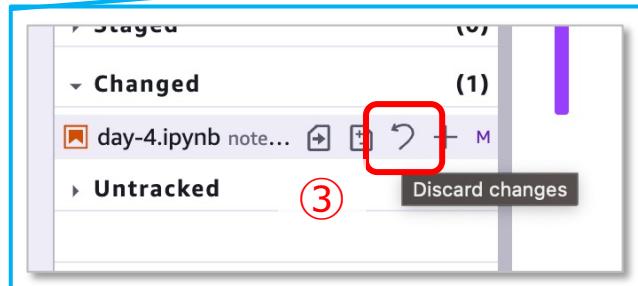
```
(studio1lab) studio-lab-user@default:~/sagemaker-studiolab-notebooks/lecture-gssm2024$ cd  
(studio1lab) studio-lab-user@default:~$ conda remove -n gssm2024 --all  
Remove all packages in environment /home/studio-lab-user/.conda/envs/gssm2024:  
## Package Plan ##  
environment location: /home/studio-lab-user/.conda/envs/gssm2024  
  
The following packages will be REMOVED:  
_libgcc_mutex-0.1-main  
_openmp_mutex-5.1-1_gnu  
altair-5.0.1-py311h06a4308_0  
appdirs-1.4.4-pyh3eb1b0_0  
asttokens-2.0.5-pyhd3eb1b0_0  
atk-1.0-2.36.0-ha0a79_0  
attr-22.1.0-py311h06a4308_0  
backcall-0.2.0-pyd3eb1b0_0  
beautifulsoup4-4.12.2-py311h06a4308_0  
blas-1.0-mkl  
:  
Proceed ([y]/n)? y  
Preparing transaction: done  
Verifying transaction: done  
Executing transaction: done  
(studio1lab) studio-lab-user@default:~$ rm -fr sagemaker-studiolab-notebooks/lecture-gssm2023  
(studio1lab) studio-lab-user@default:~$
```

(再掲) Jupyter Lab の教材を最新化するには

- 下記の手順で、教材を最新化(pull)することができます



(例) 競合がある場合のみ 拡大図



● 教材を最新化する

- ① 「Git」 ボタンを押す
- ② 「Refresh」 ボタンを押す
- ③ もし、競合がある場合(Changedが0でない場合)、
対象のファイルを手動で退避した後、
「Discard changes」 ボタンを押し
て変更を破棄する
- ④ 「Pull latest changes」 ボタンを押す
- ⑤ 画面の右下に「Successfully published」が表示されること確認する

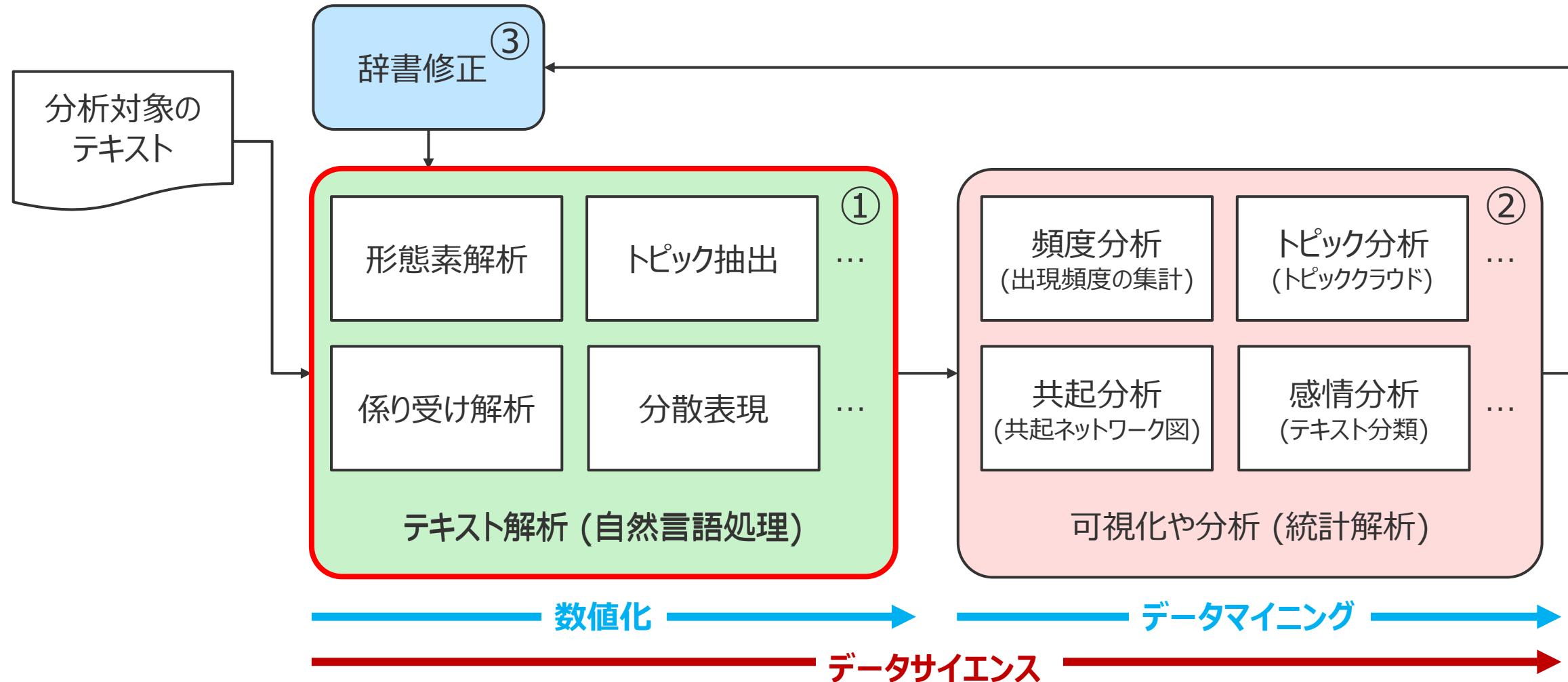
テキスト解析 (1)

(再掲) テキストマイニングの手順

- データをよく知る
 - データ件数や構成比を集計 → データを理解する
 - 旅行目的別の人気エリアは?
 - 同伴者別の人気エリアは?
 - 数値評価による人気エリアの差異は?
- テーマを設定する
 - 解決すべき課題を決める → 分析目的を明確にする
 - 数値評価が低い原因是?
 - 高評価の施設に学ぶ改善点は?
- テキスト分析に取り組む
 - これら課題を解決するために、テキスト分析を実施

テキスト分析の手順

①自然言語処理によりテキストを数値化する → ②統計解析や可視化を行う → ③結果を読み解きながら解析のための辞書を編纂する → 分析のサイクルを回していく(①へ)



代表的なテキスト解析器

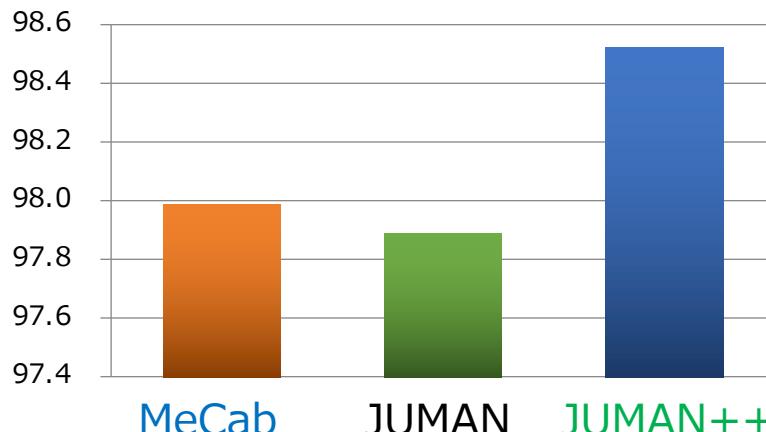
- 速度重視では **MeCab**、精度と出力情報の豊富さ重視では **JUMAN++** がお勧め

出典: <https://taku910.github.io/mecab/> をもとに加筆修正

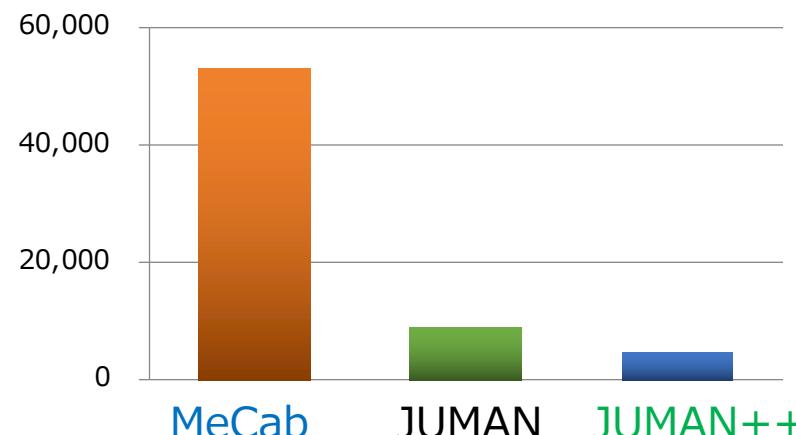
形態素解析器	ChaSen	MeCab	JUMAN	JUMAN++
コスト推定	HMM	CRF	人手	RNNLM
探索方法	接続コスト最小法 (ビタビアルゴリズム)			
係り受け解析	Cabocha	CaboCha		KNP

JUMAN++ 深層学習を使った手法で、自然な言葉の繋がりを考慮した

単語分割+品詞タグ付け精度 (F1)



処理速度 (文/秒)



学習・評価データ

京都大学テキストコーパス (NEWS),
京都大学ウェブ文書リードコーパス (WEB)

RNN言語モデルの学習

Webコーパス 1000万文

出所:

https://drive.google.com/file/d/1DVnrsWw4skRqC8jU6_RkeofOQEHFwctcview?usp=sharing

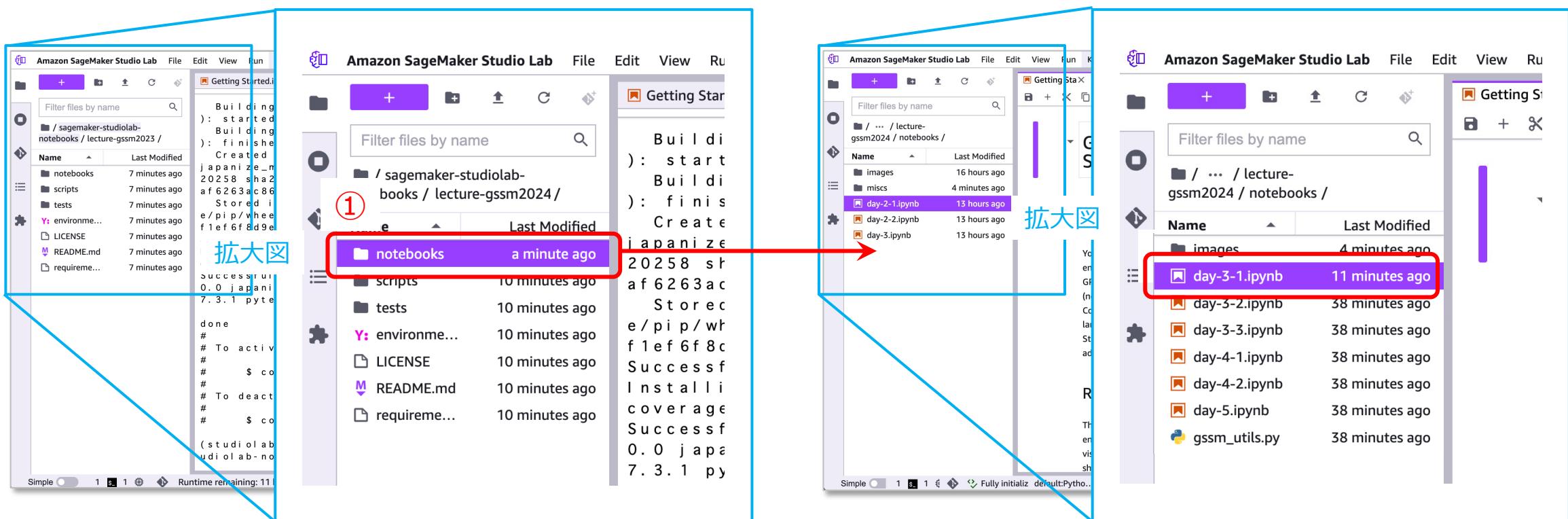
- Megagon Labs と国立国語研究所の共同研究結果として公開された OSS の日本語自然言語処理ライブラリ
 - ・「著作権表示」と「MIT ライセンスの全文」を記載する、という2条件のみで商用利用が可能
- **spaCy** (機械学習を組み込んだ自然言語処理ライブラリ) 上で動作するので、係り受け解析や固有表現抽出などの機能も利用可能
 - ・ 形態素解析には **Sudachi** (徳島人工知能NLP研究所)を利用、辞書は半年に約1回の頻度で更新されている (らしい)
 - ・ 20億文以上のWebテキストで事前学習した **Transformers** モデルも利用可能

ただし、高い利便性や機能の一方で処理が遅い (形態素解析でMeCabの10倍ぐらい)

演習 — Jupyter ノートブックの使い方

● day-3-1.ipynb を開いてください

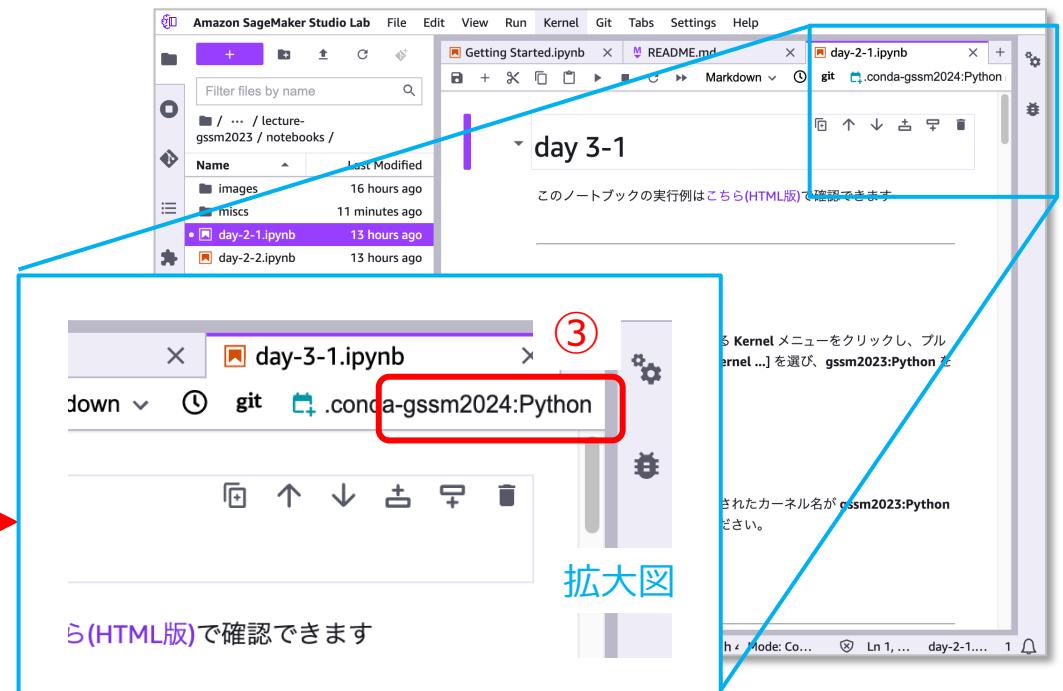
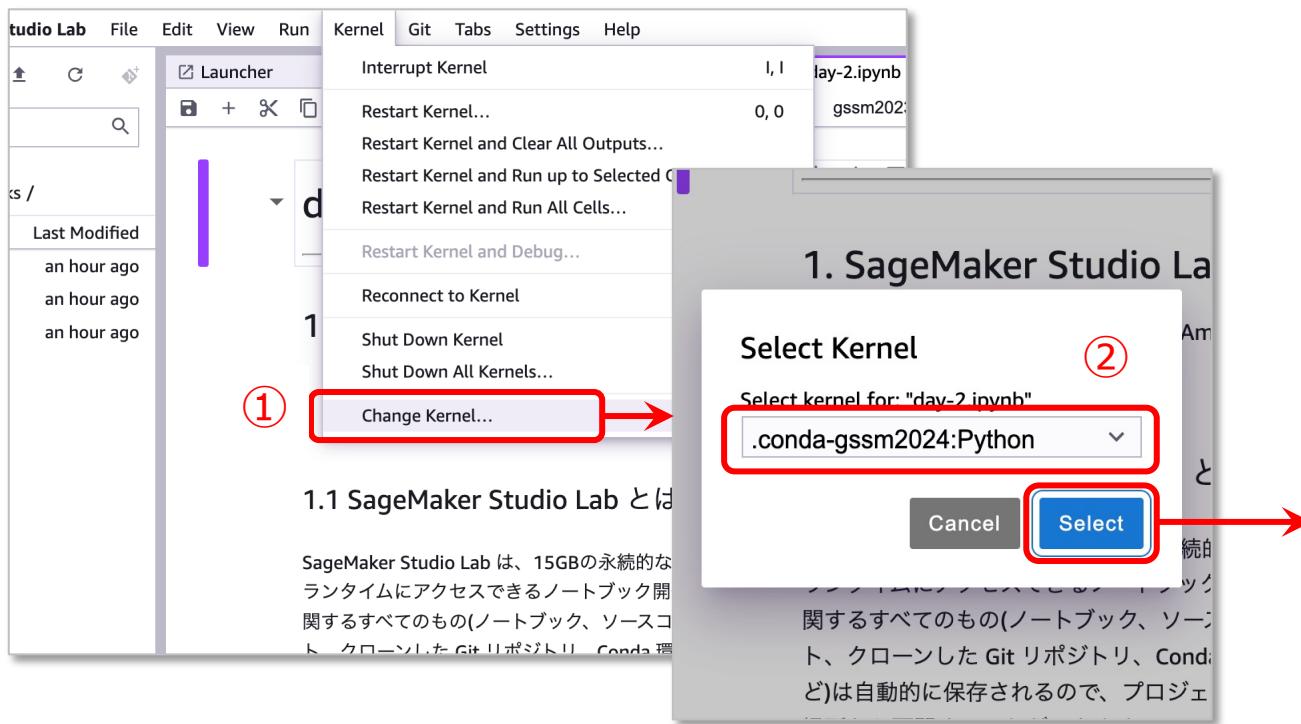
- ① 画面左の **File Browser** から ① **notebooks** をフォルダを開く
- ② さらに **day-3-1.ipynb** ノートブックを開く



演習 — Jupyter ノートブックの使い方

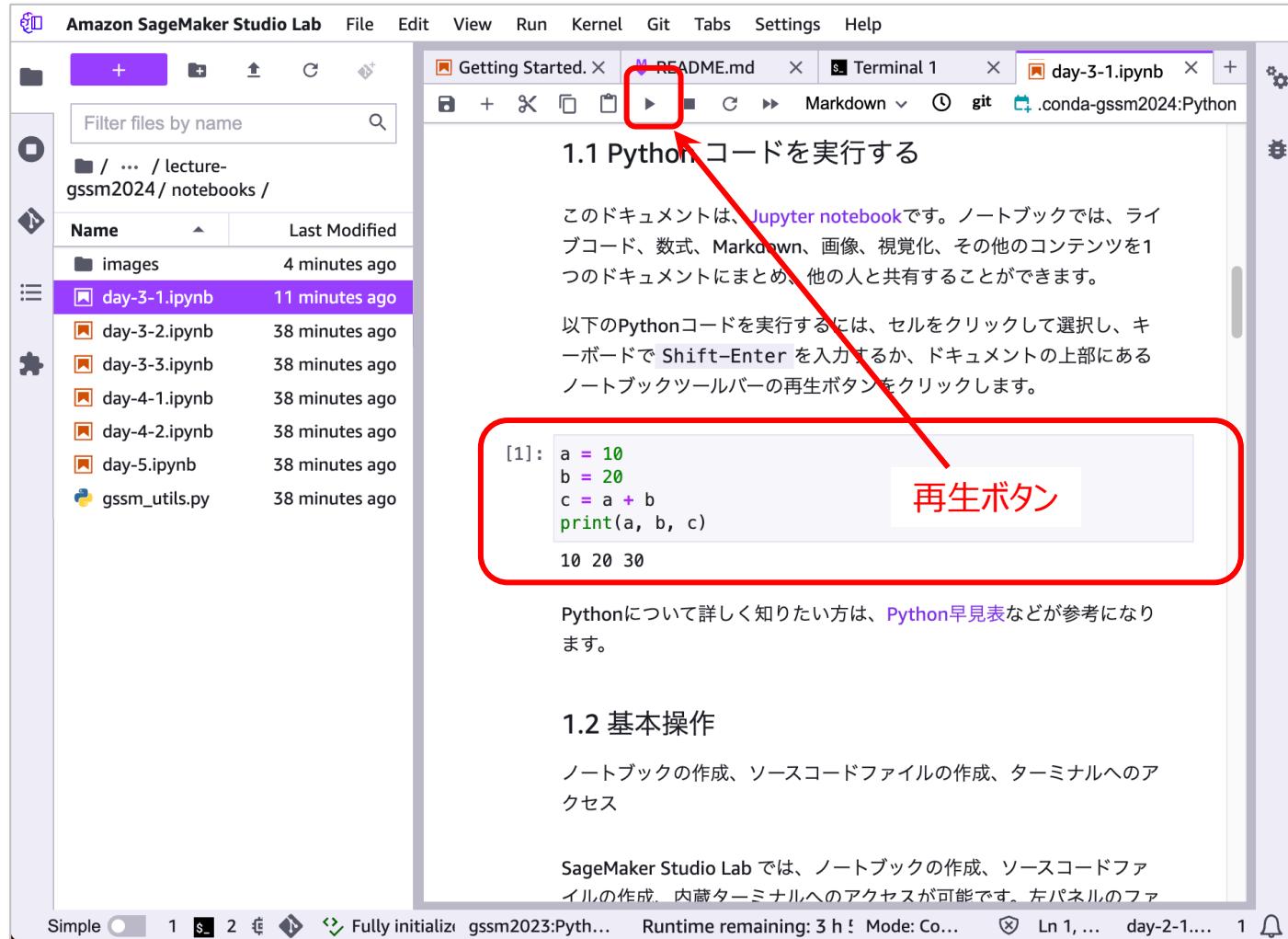
● カーネル **.conda-gssm2024:Python** を選択してください **!重要!**

- ① ページ上部の **Kernel** メニューから「**Change Kernel ...**」を選ぶ
- ② ポップアップ画面から「**.conda-gssm2024:Python**」を選択し、「**Select**」を押す
- ③ 右上隅にカーネル名「**.conda-gssm2024:Python**」が表示されていることを確認する



演習 — Jupyter ノートブックの使い方

● Python コードを実行する



- Jupyter ノートブックでは、Python のプログラムを実行することができます
(別のパッケージをインストールすることで、R 言語や C++ も実行できます)
- ノートブック上で、Pythonコードを実行するには、セルをクリックして選択し、キーボードで **Shift-Enter** を入力するか、ドキュメントの上部にあるノートブックツールバーの再生ボタンをクリックします

練習：左図の赤枠にあるセルをクリックし、セル内の計算を実行する

演習 — テキスト解析 (1)

● 形態素解析器 MeCab と、係り受け解析器 CaboCha をインストールする

The screenshot shows the Amazon SageMaker Studio Lab interface. On the left, a file browser displays a directory structure under 'lecture-gssm2024/notebooks/'. In the center, a Jupyter Notebook titled 'Getting Started.ipynb' is open. The notebook contains several cells, each with a red border and numbered steps:

- ① 2.1 MeCab のインストール (目安:約3分)**
[4]:
!bash .. /scripts/install_mecab.sh > install_mecab.log 2>&1
!tail -n 1 install_mecab.log
Successfully installed mecab-python-0.996
- ②** Successfully installed mecab-python-0.996 と表示されれば、インストール成功です。
- ③ 2.2 CaboCha のインストール (目安:約5分)**
[5]:
!bash .. /scripts/install_cabocha.sh > install_cabocha.log 2>
!tail -n 1 install_cabocha.log
Successfully installed cabocha-python-0.69
- ④** Successfully installed cabocha-python-0.69 と表示されれば、インストール成功です。
- ⑤ 2.3 Kernel のリスタート**
ページ上部のメニューにある **Kernel** メニューをクリックし、プルダウンメニューから [Restart Kernel ...] を選択してください。

A red arrow points to the play button icon in the toolbar above the notebook, labeled '再生ボタン' (Play button).

「2.1 MeCab のインストール」

- ① セルをクリックして選択し、再生ボタンを押す
この後、処理完了まで約3分程度待ちます
- ② 「Successfully ...」を確認する

「2.2 CaboCha のインストール」

- ③ セルをクリックして選択し、再生ボタンを押す
この後、処理完了まで約5分程度待ちます
- ④ 「Successfully ...」を確認する

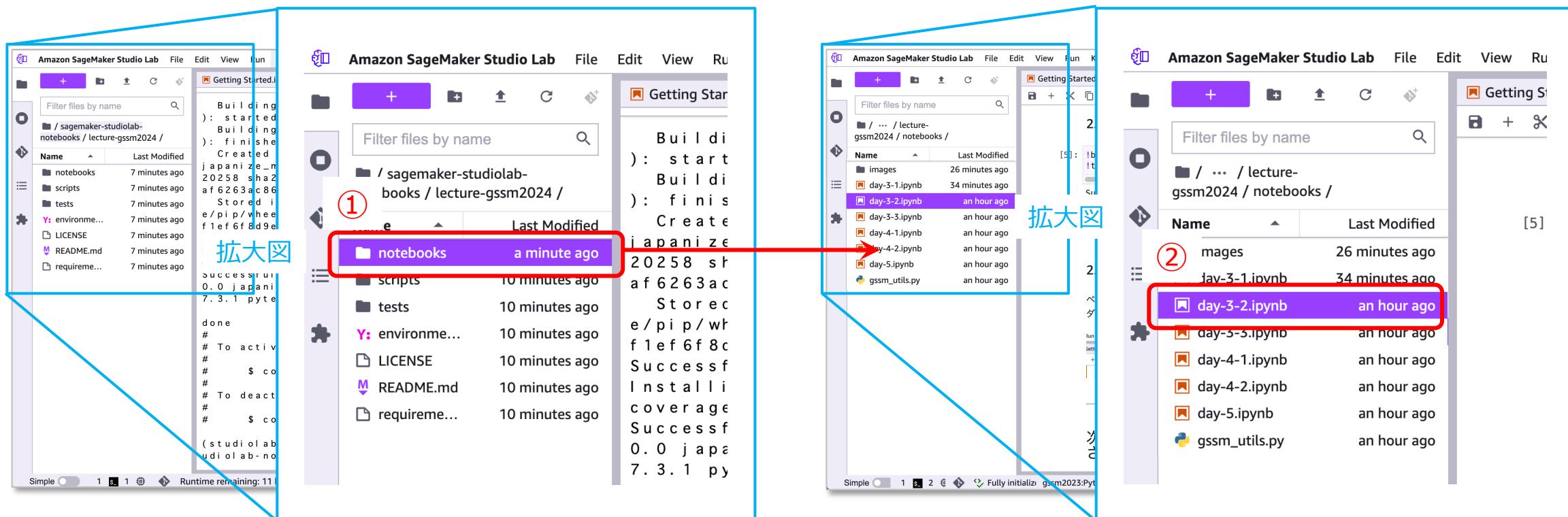
「2.3 Kernel のリスタート」

- ⑤ メニューバーにある **Kernel** メニューをクリックし、プルダウンメニューから [Restart Kernel ...] を選択する

演習 — テキスト解析 (1)

● day-3-2.ipynb を開いてください

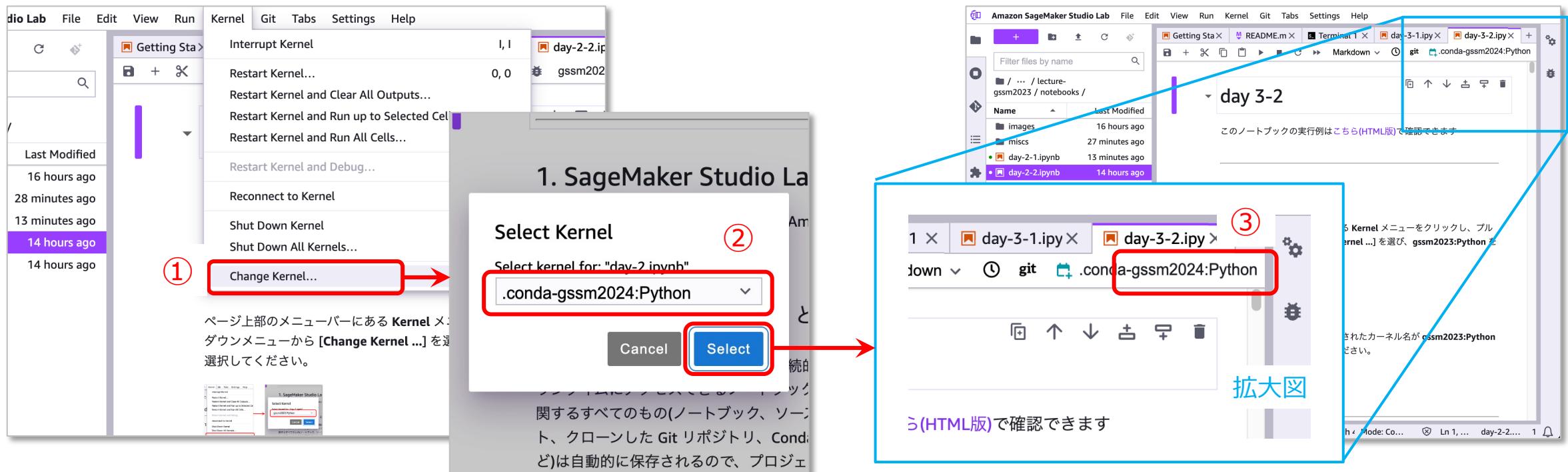
- ① 画面左の File Browser から ① notebooks をフォルダを開く (既に開いている場合はスキップ)
- ② 次に day-3-2.ipynb ノートブックを開く



演習 — テキスト解析 (1)

● カーネル **.conda-gssm2024:Python** を選択してください !重要!

- ① ページ上部の **Kernel** メニューから「**Change Kernel ...**」を選ぶ
- ② ポップアップ画面から「**.conda-gssm2024:Python**」を選択し、「**Select**」を押す
- ③ 右上隅にカーネル名「**.conda-gssm2024:Python**」が表示されていることを確認する



演習 — テキスト解析 (1)

● 形態素解析を行う (コマンドライン実行と同じ形式)

3.1 MeCab を使う

(1) そのまま出力してみる

①

```
import MeCab

tagger = MeCab.Tagger("-r ..//tools/usr/local/etc/mecabrc")
print(tagger.parse("今日はいい天気です"))
```

```
今日    名詞,副詞可能,*,*,*,*,今日,キヨウ,キヨー
は      助詞,係助詞,*,*,*,*,は,ハ,ワ
いい   形容詞,自立,*,*,形容詞・イイ,基本形,いい,イイ,イイ
天気   名詞,一般,*,*,*,*,天気,テンキ,テンキ
です   助動詞,*,*,*,特殊・デス,基本形,です,デス,デス
EOS
```

- ① セルをクリックして選択し、再生ボタンを押す
 - この方法では、コマンドライン実行した場合と同じ形式で出力されます
 - ただし、テキスト解析では、**テキストを数値化し、統計処理を行う必要**があります
 - そこで、**統計処理で扱いやすい DataFrame 型**(テーブル形式)に格納します → 次ページ

演習 — テキスト解析 (1)

● 形態素解析を行う (DataFrame 型に格納する)

②

```
import pandas as pd

node = tagger.parseToNode("今日はいい天気です")
features = []
while node:
    features.append(node.feature.split(','))
    node = node.next

columns = [
    "品詞", "品詞細分類1", "品詞細分類2", "品詞細分類3", "活用型", "活用形", "基本形",
    "読み", "発音",
]
pd.DataFrame(features, columns=columns)
```

[2]:

	品詞	品詞細分類1	品詞細分類2	品詞細分類3	活用型	活用形	基本形	読み	発音
0	BOS/EOS	*	*	*	*	*	*	*	*
1	名詞	副詞可能	*	*	*	*	今日	キョウ	キョー
2	助詞	係助詞	*	*	*	*	は	ハ	ワ
3	形容詞	自立	*	*	形容詞・イイ	基本形	いい	イイ	イイ
4	名詞	一般	*	*	*	*	天気	テンキ	テンキ
5	助動詞	*	*	*	特殊・デス	基本形	です	デス	デス
6	BOS/EOS	*	*	*	*	*	*	*	*

- ② セルをクリックして選択し、再生ボタンを押す
- この方法では、形態素解析器の出力を統計処理で扱いやすい DataFrame 型 (テーブル形式) に格納しています

練習: 入力文「**今日はいい天気です**」の内容を変更して、形態素解析(②)を行った結果を確認してください

演習 — テキスト解析 (1)

● 係り受け解析を行う（コマンドライン実行と同じ形式）

4.1 CaboCha を使う

(1) そのまま出力してみる

①

```
import CaboCha

cp = CaboCha.Parser("-r ../tools/usr/local/etc/cabocharc")
tree = cp.parse("今日はいい天気です")
print(tree.toString(CaboCha.FORMAT_LATTICE))
```

```
* 0 2D 0/1 -1.041733
今日 名詞,副詞可能,*,*,*,*,-,今日,キヨウ,キヨー
は 助詞,係助詞,*,*,*,-,は,ハ,ワ
* 1 2D 0/0 -1.041733
いい 形容詞,自立,*,*,-,形容詞・イイ,基本形,いい,イイ,イイ
* 2 -1D 0/1 0.000000
天気 名詞,一般,*,*,*,-,天気,テンキ,テンキ
です 助動詞,*,*,-,特殊・デス,基本形,です,デス,デス
EOS
```

- ① セルをクリックして選択し、再生ボタンを押す
- この方法では、コマンドライン実行した場合と同じ形式で出力されます
 - ただし、**係り元**や**係り先**の関係を把握するには、この出力形式でも、表形式でも直感的ではありません
 - そこで、**係り受け関係を確認し易いツリー形式**で出力します → 次ページ

演習 — テキスト解析 (1)

● 係り受け解析を行う（係り受けペアを抽出する）

②

```
# 構文木(tree)からチャンクを取り出す
def get_chunks(tree):
    chunks = {}
    key = 0
    for i in range(tree.size()):
        tok = tree.token(i)
        if tok.chunk:
            chunks[key] = tok.chunk
            key += 1
    return chunks

# チャンク(chunk)から表層形を取り出す
def get_surface(chunk):
    surface = ''
    begin = chunk.begin
    end = chunk.end
    for i in range(begin, end):
        surface += chunk[i]
    return surface
```

← 繰り返し呼ばれる処理などをまとめて関数として定義したもの

③

```
tree = cp.parse("今日はいい天気です")
chunks = get_chunks(tree)

for from_chunk in chunks.values():
    if from_chunk.link < 0:
        continue
    to_chunk = chunks[from_chunk.link]

    from_surface = get_surface(from_chunk)
    to_surface = get_surface(to_chunk)

    print(from_surface, '→', to_surface)
```

今日は → 天気です
いい → 天気です

② セルをクリックして選択し、再生ボタンを押す（③より前に一度実行しておく）

③ セルをクリックして選択し、再生ボタンを押す

- この方法では、係り受け解析器の出力を**係り元**と**係り先**の**関係**を持つ単語のペアを抽出しています

練習：入力文「**今日はいい天気です**」の内容を変更して、係り受け解析(③のみ)を行った結果を確認してください

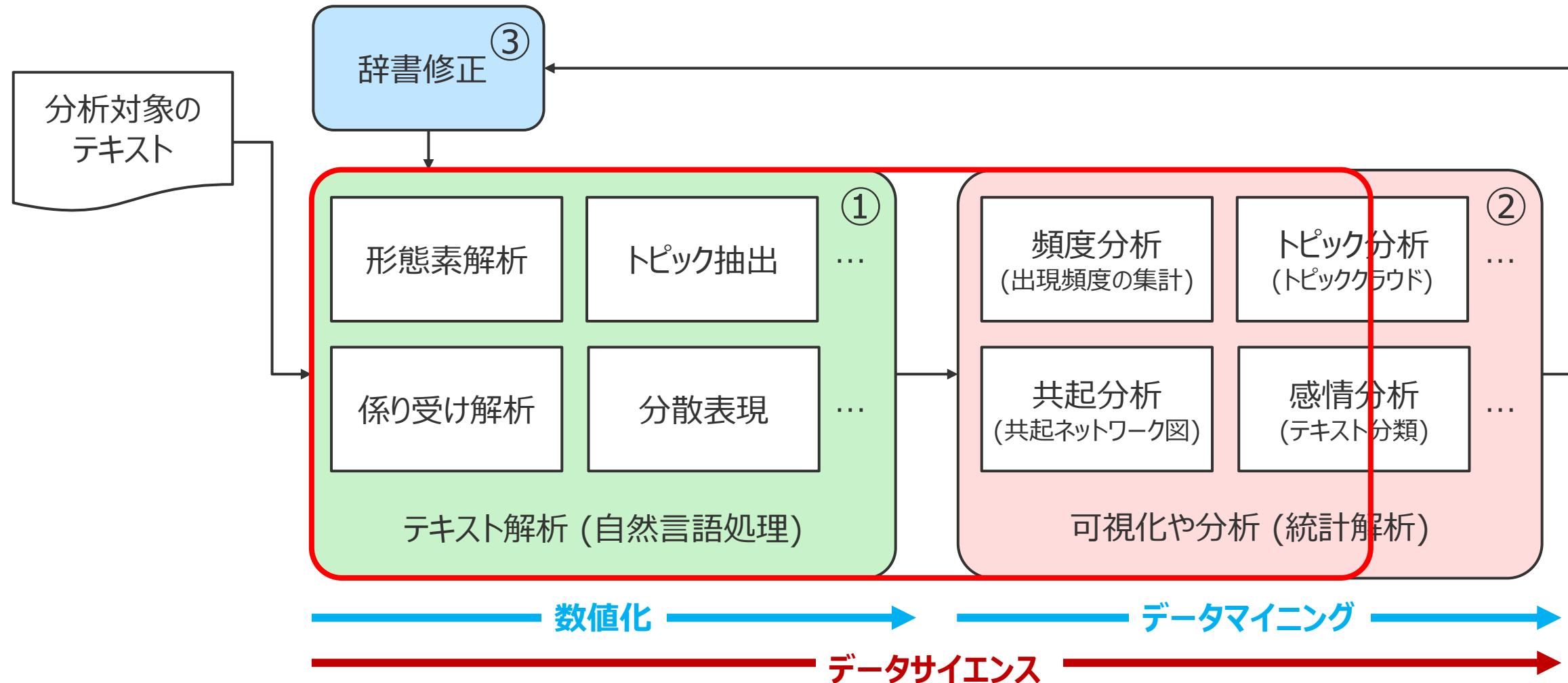
テキスト解析 (2)

(再掲) テキストマイニングの手順

- データをよく知る
 - データ件数や構成比を集計 → データを理解する
 - 旅行目的別の人気エリアは?
 - 同伴者別の人気エリアは?
 - 数値評価による人気エリアの差異は?
- テーマを設定する
 - 解決すべき課題を決める → 分析目的を明確にする
 - 数値評価が低い原因是?
 - 高評価の施設に学ぶ改善点は?
- テキスト分析に取り組む
 - これら課題を解決するために、テキスト分析を実施

テキスト分析の手順

①自然言語処理によりテキストを数値化する → ②統計解析や可視化を行う → ③結果を読み解きながら解析のための辞書を編纂する → 分析のサイクルを回していく(①へ)



- 社会調査データを分析する目的で開発されたフリー(~~商用可能~~)のツール

- 高機能かつ~~商用可能~~でフリー
- Rを用いた多変量解析と可視化
- 実装されている分析手法
 - 階層的クラスター分析
 - 多次元尺度構成法(MDS)
 - 対応分析
 - 共起ネットワーク
 - 自己組織化マップ
 - 文書のクラスター分析
 - トピックモデル (LDA)

論文検索サービスも提供 → <http://khcoder.net/bib.html>

研究事例リスト

KH Coderを用いたご研究の成果を発表された際には、書誌情報をフォームにご記入いただけますと幸いです。

出版年 :

著者名 :

キーワード :

ヒット件数 : 0200 / 6135

KH Coderを用いた研究事例のリスト 6135件

※2023/6/16 現在

→1646→2042→2695→3741件→4554件→昨年5355件→6135件)

(参考) 2023年12以降のライセンスや料金

2023年12月 無償公開が終了しました

	Starting	Base	Master	Pro
通常ライセンス 販売価格(税込) ※下段は、バージョンアップに追従できるアップデートのサブスクリプション費用(1年間)	無料	59,950 (43,780)	196,900 (43,780)	396,000 (43,780)
アカデミック・ライセンス 販売価格(税込) ※下段は、バージョンアップに追従できるアップデートのサブスクリプション費用(1年間)	無料	24,750 (18,480)	69,850 (18,480)	—
インストール可能台数 ※ライセンス保有者が管理するPCに限る	1台	2台	2台	2台

	Starting	Base	Master	Pro
無料で分析を始められる 一部機能制限ありの公開テスト(beta)版 ・分析対象はデータファイルの最初の100件目まで ・強制抽出語と無視する語の指定はそれぞれ1語のみ ・分析対象の品詞タイプを選択できない	○	—	—	—
機能制限なしの正式版 外部プラグイン「文錦@シリーズ」購入による機能追加も可能	—	○	○	○
アップデートのサブスクリプション(1年間有効) KH Coderのバージョンアップに追従できる便利なバッチを提供	—	○	○	○
インストール／セットアップのサポート KH Coderの導入段階のトラブルシューティング	—	—	○	○
セミナー受講券(1年半以内有効) 「計量テキスト分析公式セミナー」初級編・ステップアップ編に各1回参加可	—	—	1名	1名
使い方のQ&A対応 使い方に関する困り事やご質問にメールでサポート	—	—	—	6回/年

出典: <https://www.screen.co.jp/as/solution/khcoder>

共起ネットワーク

抽出語またはコードを用いて、出現パターンの似通ったものを線で結んだ図、すなわち共起関係を線（edge）で表したネットワークを描く機能です。



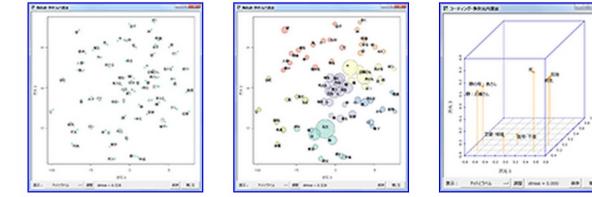
共起の程度が非常に強いものだけを線で結んだ図

やや弱い共起関係も描画に含め、自動的にグループ分け（色分け）

出現数が多い語ほど大きく、また共起の程度が強いほど太い線で描画

多次元尺度構成法 (MDS)

同じく抽出語またはコードを用いての、多次元尺度構成法です。



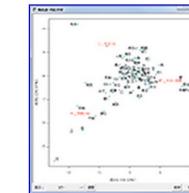
2次元の解

New! クラスタリングと色分け

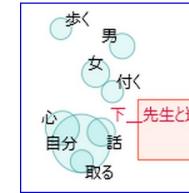
3次元の解

対応分析

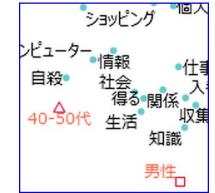
同じく抽出語またはコードを用いての、対応分析です。



同時布置図



New! バブルプロット



複数の外部変数を用いた多重対応分析

分析手法

説明

共起ネットワーク

- 同時に出現した単語同士をネットワークで結んで図示したもの
- 同時に出現したかといった共起の有無を集計し、ネットワークを作成
- 関係の強さ Jaccard 係数で評価し、媒介性やグラフクラスタリングを使ってサブグラフも検出できる

多次元尺度構成法 (MDS)

- 出現パターンの似た単語同士を近くに置くよう図示したもの
- 出現パターンとは、ある単語がどの文書に出現したかといった関係を単語ベクトルとして表現したもの
- 似ている(=距離が近い)の計算は Jaccard、ユークリッド、コサイン距離のいずれかで求める

対応分析 (コレ спинデンス分析)

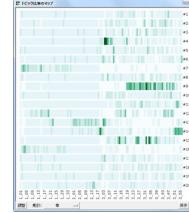
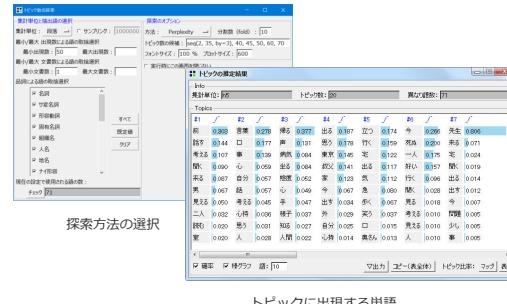
- 出現パターンの似た単語や外部変数を近くに置くよう図示したもの
- 単語と単語または外部変数が同時に出現した頻度をクロス集計し、相関が最大になるような2軸でプロット
- PCA が元の情報をそのまま可視化するのに対して、対応分析は似ているものを近くに表示する
- 外部変数も同時にプロットできる

(参考) KH Coder — 分析手法 (2)

2023年12月 無償公開が終了しました

トピックモデル (LDA)

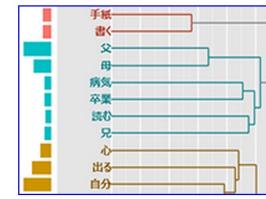
文書ごとにトピックの出現割合を表示したり、各トピックに高い確率で出現する語を表示できます。



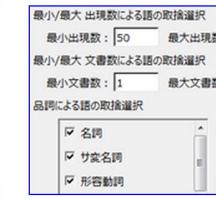
文書ごとのトピック比率

階層的クラスター分析

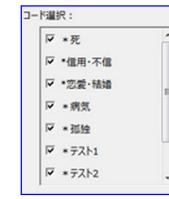
抽出語の階層的クラスター分析を行い、デンドログラムを表示します。抽出語だけでなくコーディング結果（コード）についても、同じように分析を行えます。



New! デンドログラム



抽出語は出現数や品詞で選択



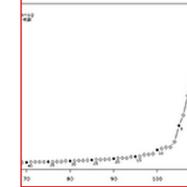
コードはチェックボックスで直接選択

文書のクラスター分析

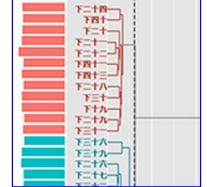
文書の分類を行うクラスター分析です。



クラスター分析の結果画面



併合水準のプロット。クラスター数5付近から併合水準が急上昇。10でも少し上がっているので、この場合クラスター数は11が良いか。



文書のデンドログラム。左の棒グラフは各文書の長さをあらわす。なお、文書数が500を超える場合、デンドログラムは表示不可。

分析手法

トピックモデル (LDA)

- 文書が複数のトピックを持つと仮定、文書ごとにトピックの出現割合、各トピックに高確率で出現する語を表示
- R の topicmodels パッケージに含まれる LDA 関数(ギブスサンプリング)を利用 (乱数のシードは固定)
- トピックモデルは教師なし学習**のため、コーディングルールで単語を集約するよりも客観性が高い

階層的クラスター分析

- 出現パターンの似た**単語同士をグルーピング(クラスタリング)**して、樹形図にしたもの
- 出現パターンは、ある単語がどの文書に出現したかといった関係を単語ベクトルとして表現したもの
- 似ている(=距離が近い)の計算は Jaccard、ユークリッド、コサイン距離のいずれかで求める

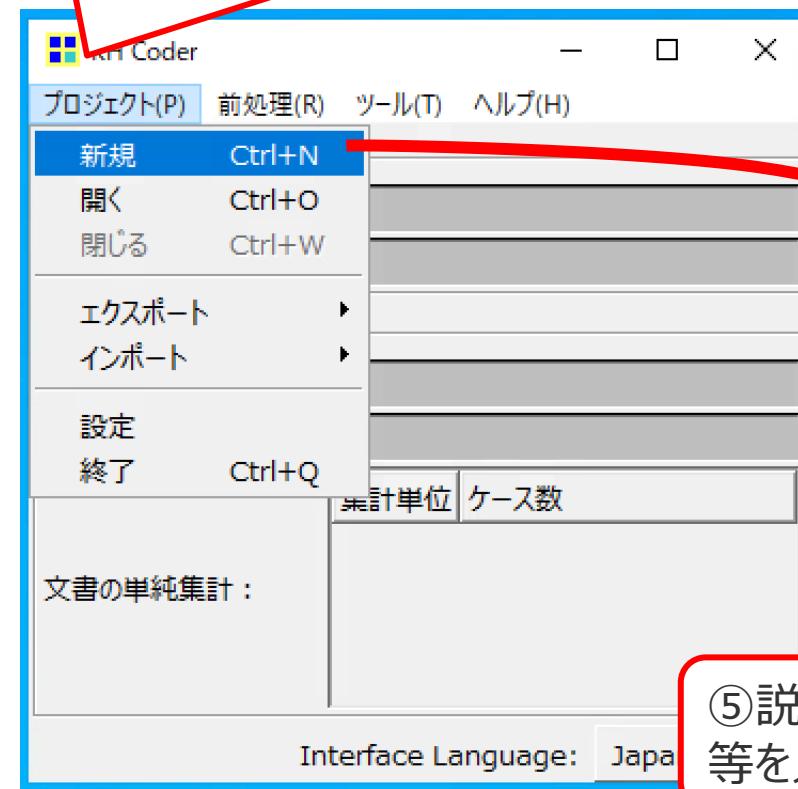
文書のクラスター分析

- 似た**文書同士をグルーピング(クラスタリング)**して、樹形図にしたもの
- 各文書は、文書中に出現する単語の有無でベクトル化した文書ベクトルで表現
- 似ている(=距離が近い)の計算は Jaccard、ユークリッド、コサイン距離のいずれかで求める
- いわゆる Ward法、群平均法、最遠隣法で階層クラスタを作成する

説明

● プロジェクトの作成

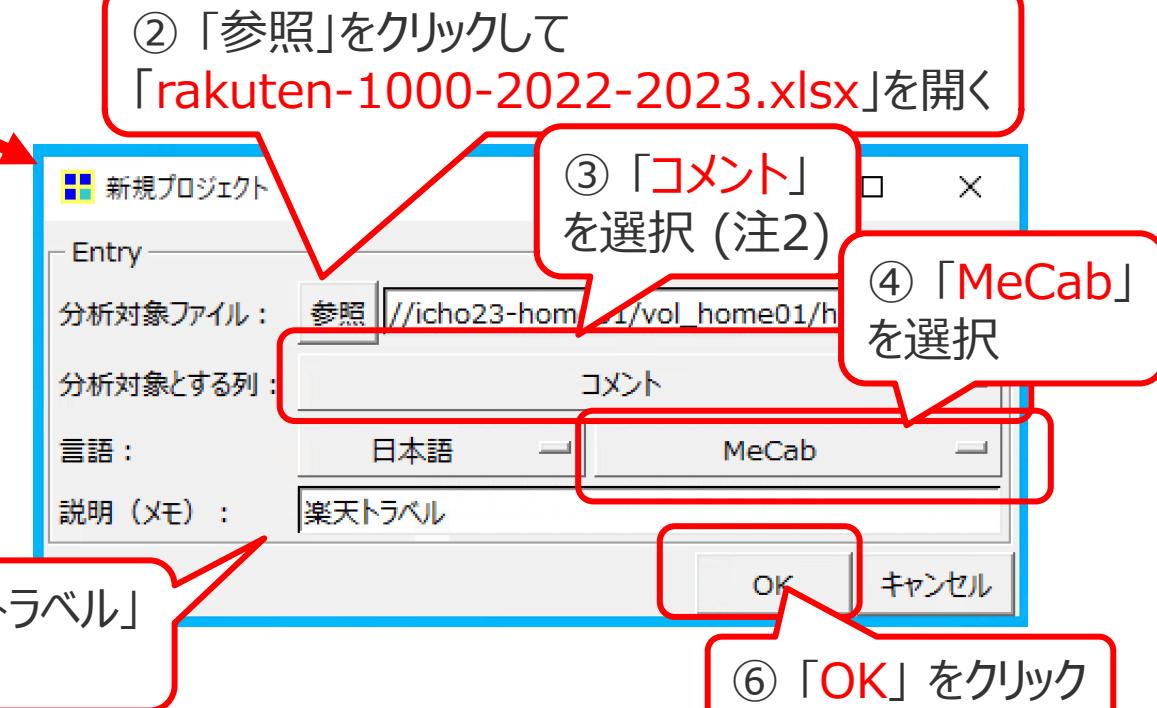
①メニューから「プロジェクト」「新規」を選択 (注1)



注1: 次回 KH Coderを起動した時は「新規」ではなく
「開く」を選択します

注2: ②のファイル選択後,ここに「テキスト」等の
選択項目が表示されるまで数分がかかります

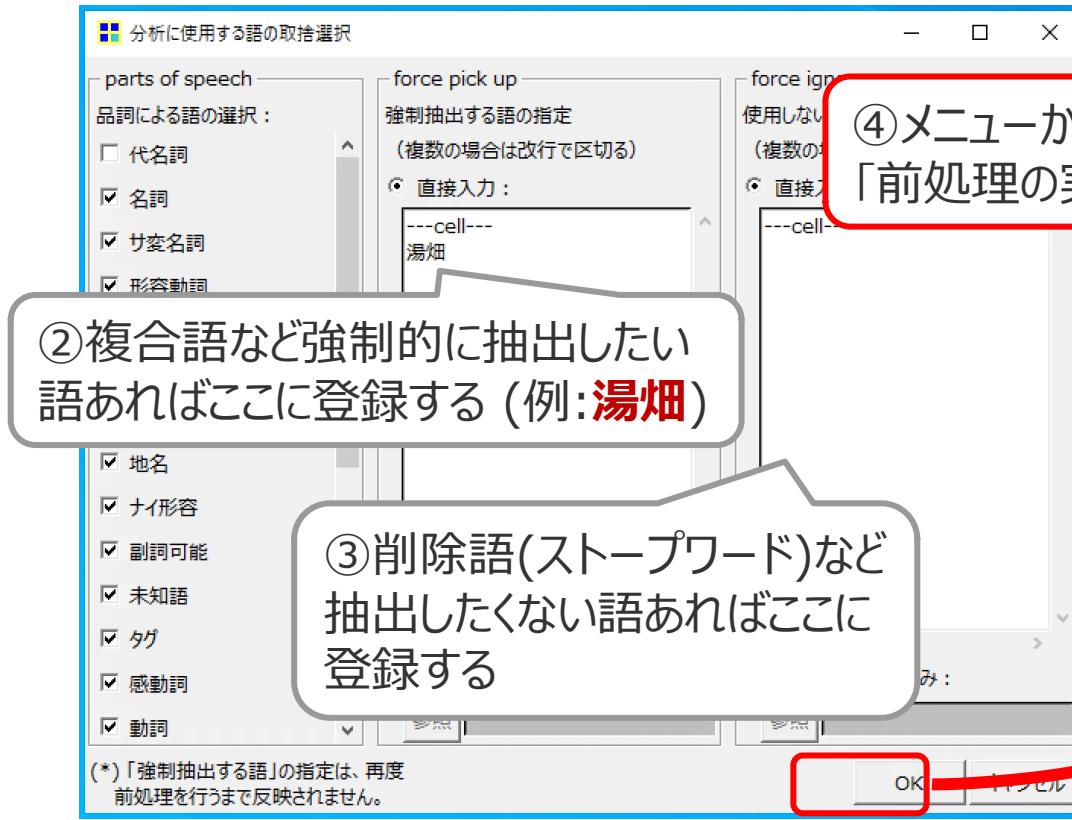
②「参照」をクリックして
「rakuten-1000-2022-2023.xlsx」を開く



⑤説明「楽天トラベル」
等を入力

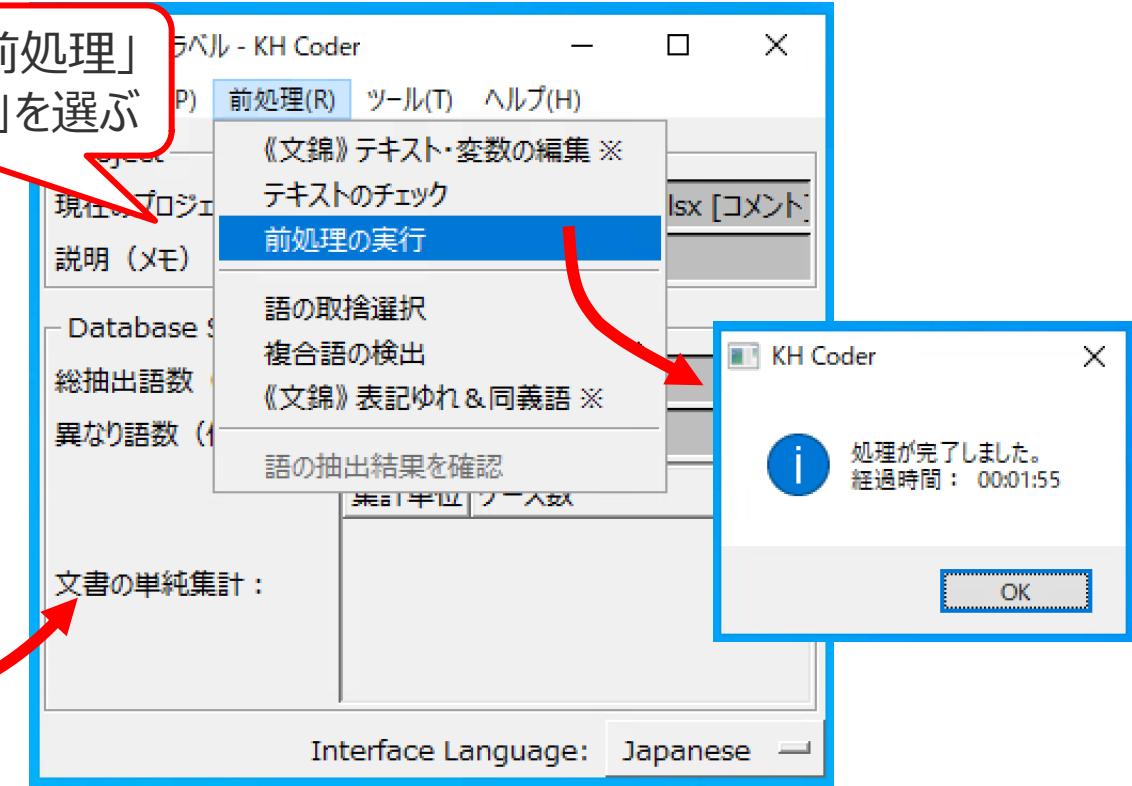
● 前処理(形態素解析)の実行

①メニューから「前処理」「語の取捨選択」を選ぶ



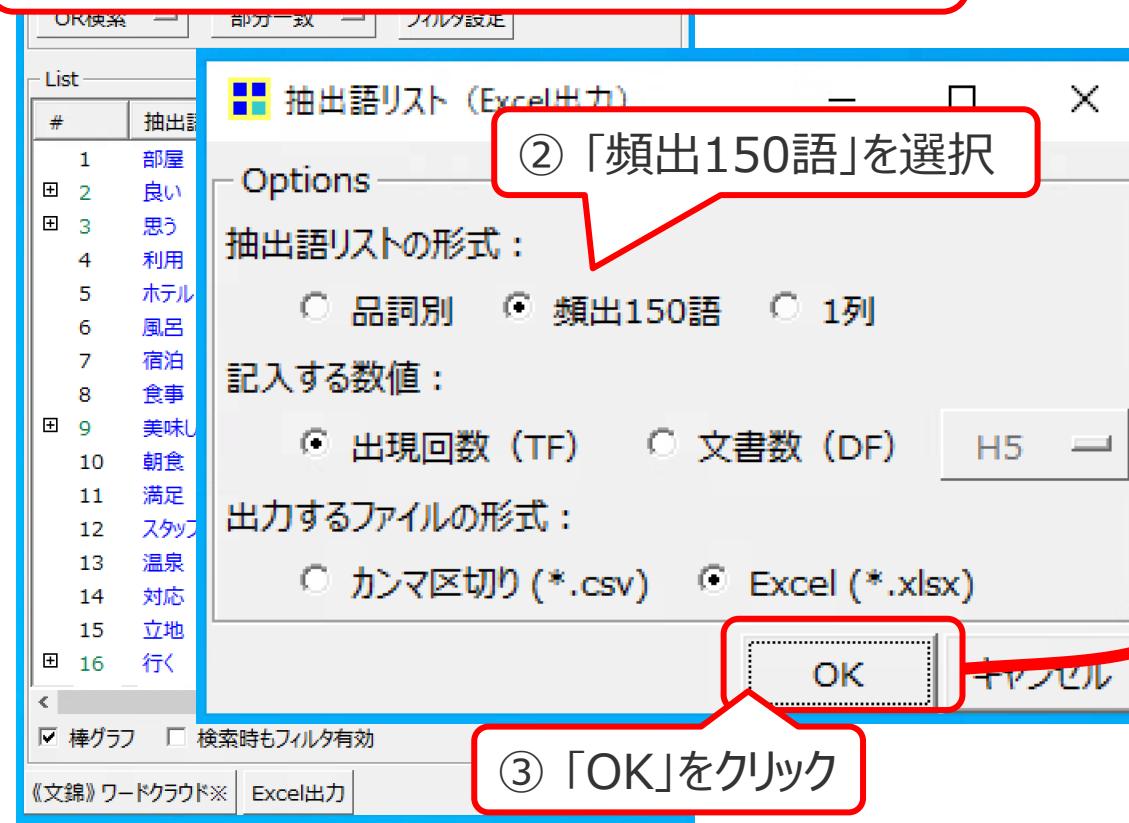
④メニューから「前処理」「前処理の実行」を選ぶ

- 注1: EXCELファイルを読み込んで分析する場合,あらかじめ「---cell---」が入力されています
- 注2: メニューから「前処理」「複合語の検出」を選ぶと,複合語候補の一覧を出力できます



● 頻出語を確認する

- ①メニューから「ツール」「抽出語」「抽出語リスト」
→右下「EXCEL出力」ボタンを選択

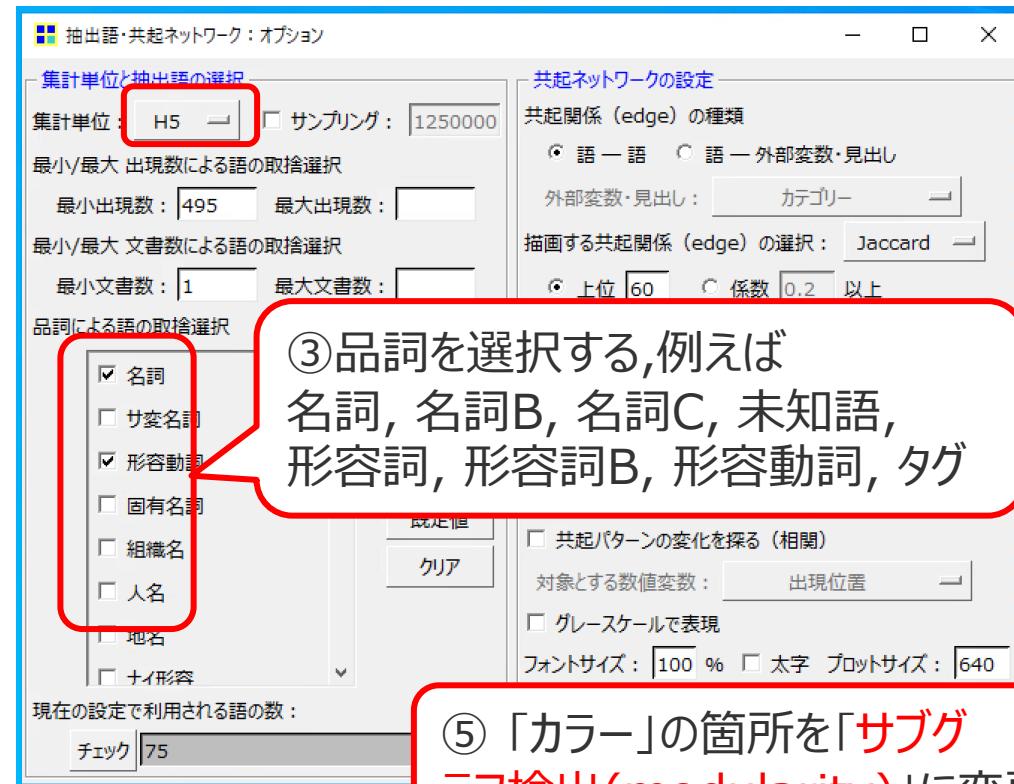


A	B	C	D	E	F	G	H
1	抽出語	出現回数	抽出語	出現回数	抽出語	出現回数	
2	部屋	6689	子供	661	プラン	389	
3	良い	4302	過ごす	657	見える	388	
4	思う	3976	家族	648	機会	387	
5	利用	3481	予約	636	設備	387	
6	ホテル	2831	過ごせる	626	旅館	386	
7	風呂	2702	駐車	613	置く	384	
8	宿泊	2649	素晴らしい	612	きれい	377	
9	食事	2447	月	611	歩く	368	
10	美味しい	2249	バス	610	湯	359	
11	朝食	2172	丁寧	610	施設	345	
12	満足	1785	アメニティ	609	無料	345	
13	スタッフ	1712	清潔	556	新しい	340	
14	温泉	1705	入れる	544	楽しい	335	
15	対応	1603	使う	536	掃除	335	
16	立地	1374	初めて	523	気持ち	328	
17	行く	1334	無い	521	雰囲気	328	
18	広い	1314	人	520	女性	323	
19	綺麗	1193	バイキング	515	シャワー	321	
20	宿	1171	嬉しい	515	建物	316	
21	大変	1157	ベッド	514	高い	316	
22	少し	1156	他	504	問題	316	
23	残念	1155	親切	503	全体	314	
24	最高	1118	種類	502	大きい	313	

● 共起ネットワークの作成(1)

①メニューから「ツール」「抽出語」「共起ネットワーク」を選ぶ

②「集計単位」として「H5」を選んで「OK」をクリック



⑤「カラー」の箇所を「サブグラフ検出(modularity)」に変更

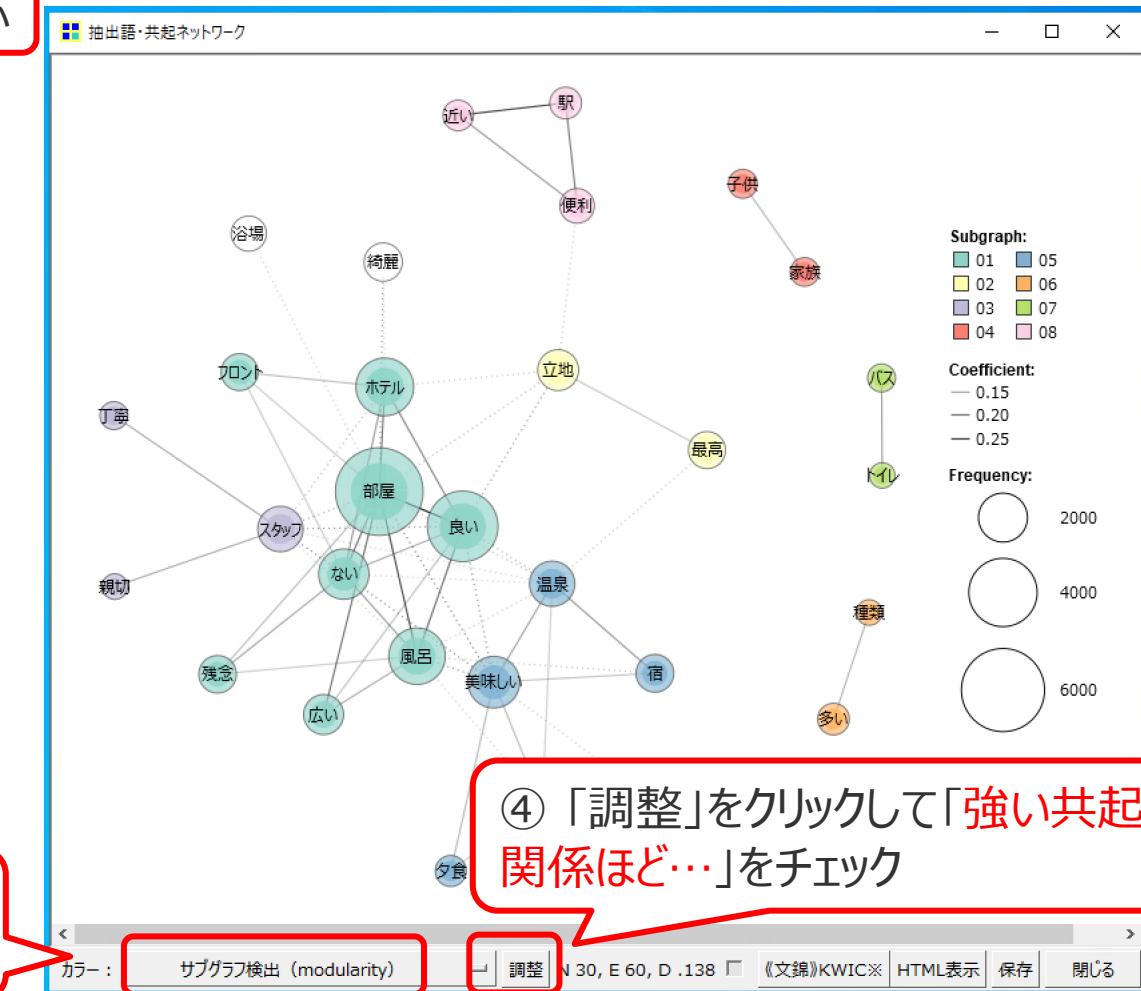


表 A.1 KH Coder の品詞体系

KH Coder 内の品詞名	茶筌の出力における品詞名
名詞	名詞-一般（漢字を含む 2 文字以上の語）
名詞 B	名詞-一般（平仮名のみの語）
名詞 C	名詞-一般（漢字 1 文字の語）
サ変名詞	名詞-サ変接続
形容動詞	名詞-形容動詞語幹
固有名詞	名詞-固有名詞-一般
組織名	名詞-固有名詞-組織
人名	名詞-固有名詞-人名
地名	名詞-固有名詞-地域
ナイ形容	名詞-ナイ形容詞語幹
副詞可能	名詞-副詞可能
未知語	未知語
感動詞	感動詞またはフィラー
タグ	タグ
動詞	動詞-自立（漢字を含む語）
動詞 B	動詞-自立（平仮名のみの語）
形容詞	形容詞（漢字を含む語）
形容詞 B	形容詞（平仮名のみの語）
副詞	副詞（漢字を含む語）
副詞 B	副詞（平仮名のみの語）
否定助動詞	助動詞「ない」「まい」「ぬ」「ん」
形容詞（非自立）	形容詞-非自立（「がたい」「つらい」「にくい」等）
その他	上記以外のもの

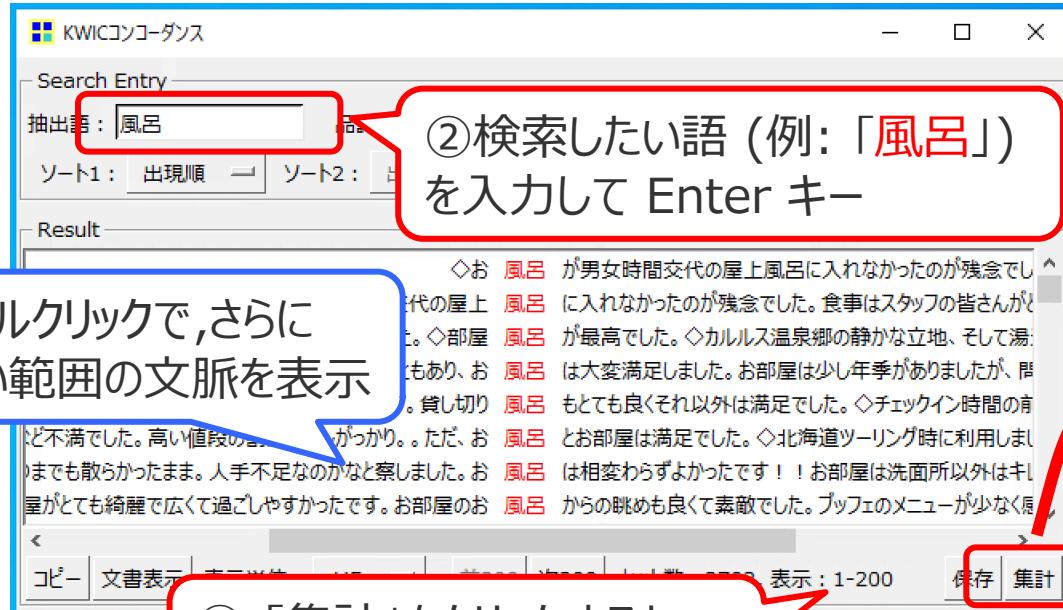
出典: KH Coder 3 リファレンス・マニュアル

注: どの品詞を選択すべきかは、分析対象のデータや分析目的により異なります。

分析結果を確認しながら、適宜、適切な品詞選択を検討することが重要です。

● 前後文脈を確認する

- ①メニューから「ツール」「抽出語」「KWICコンコーダンス」を選ぶ



- ③「集計」をクリックするとコレーション統計(右)を開く

注: 共起ネットワーク上で「風呂」をクリックすると①②と同じ操作となります(V3以降)

「右1」は右側の1つ目(=直後)に出現していた回数

The screenshot shows the Collocation Statistics window with '風呂' entered in the search field. A callout bubble ④ points to the text '表示する語の品詞を選択(例: 形容詞, 形容詞B, 形容動詞)' (Select the part of speech of the word to be displayed (example: Adjective, AdjectiveB, Adjective Verb)). Another callout bubble ⑤ points to the '右合計' button at the bottom right of the window. The table lists various words and their collocation statistics, including the count of occurrences to the right of the target word.

N	抽出語	品詞	合計	左合計	右合計	左5	左4	左3	左2	左1	右1	右2	右3	右4	右5	スコア
1	広い	形容詞	190	41	149	8	5	13	13	2	1	104	20	18	6	81.0
2	良い	形容詞	222	77	145	40	9	14	13	1	6	48	43	21	27	77.4
3	最高	名詞	105	13	92	5	3	3	2			44	12	22	8	42.8
4	部屋	名詞	439									22	20	18	172.7	
5	トイレ	名詞	117									1	4	4	60.2	
6	露天風呂	名詞	87									9	13	13	30.6	
7	風呂	名詞	130									31	14	20	35.6	

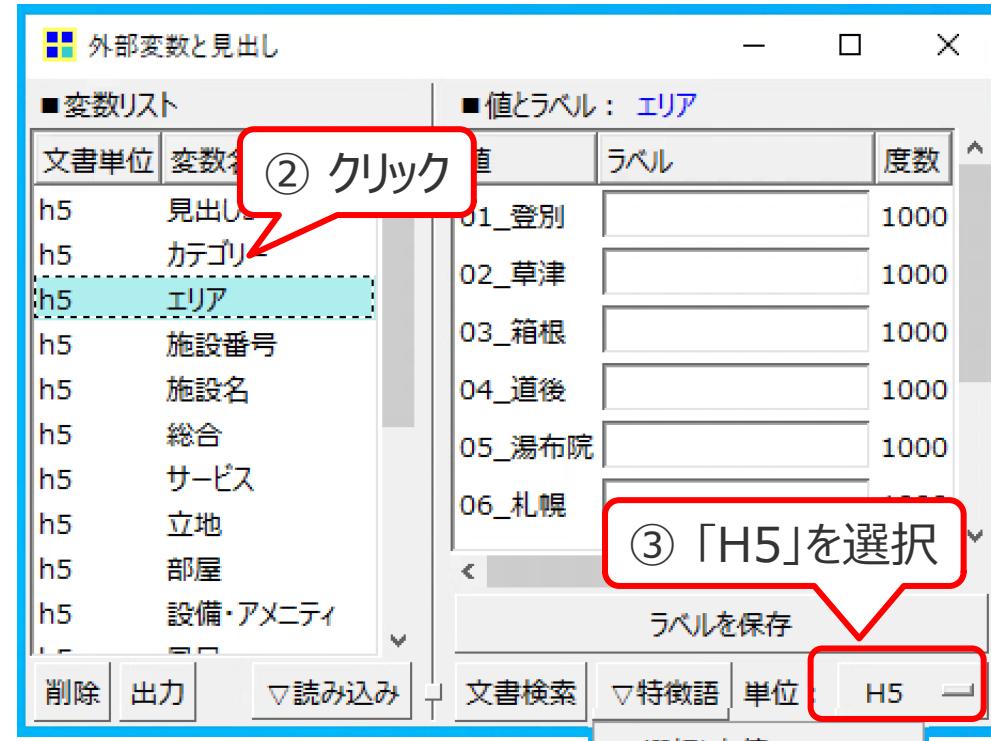
「広い」は「風呂」の2語後に 104 回出現

- ④表示する語の品詞を選択(例: 形容詞, 形容詞B, 形容動詞)

- ⑤「右合計」でソート

● 外部変数を利用する

- ① メニューから「ツール」「外部変数と見出し」を開く



- ④ 「特徴語」「一覧(Excel形式)」を選択

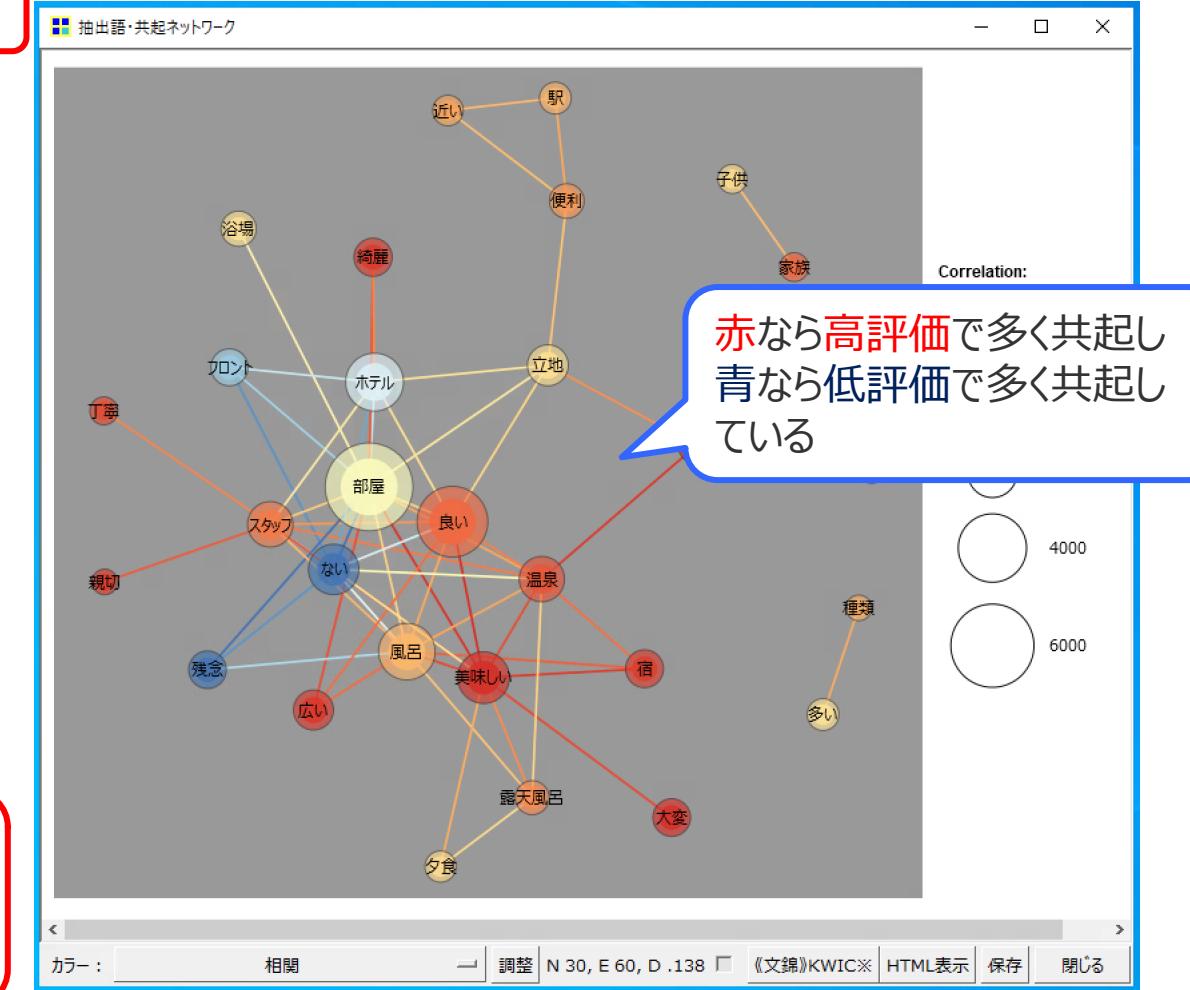
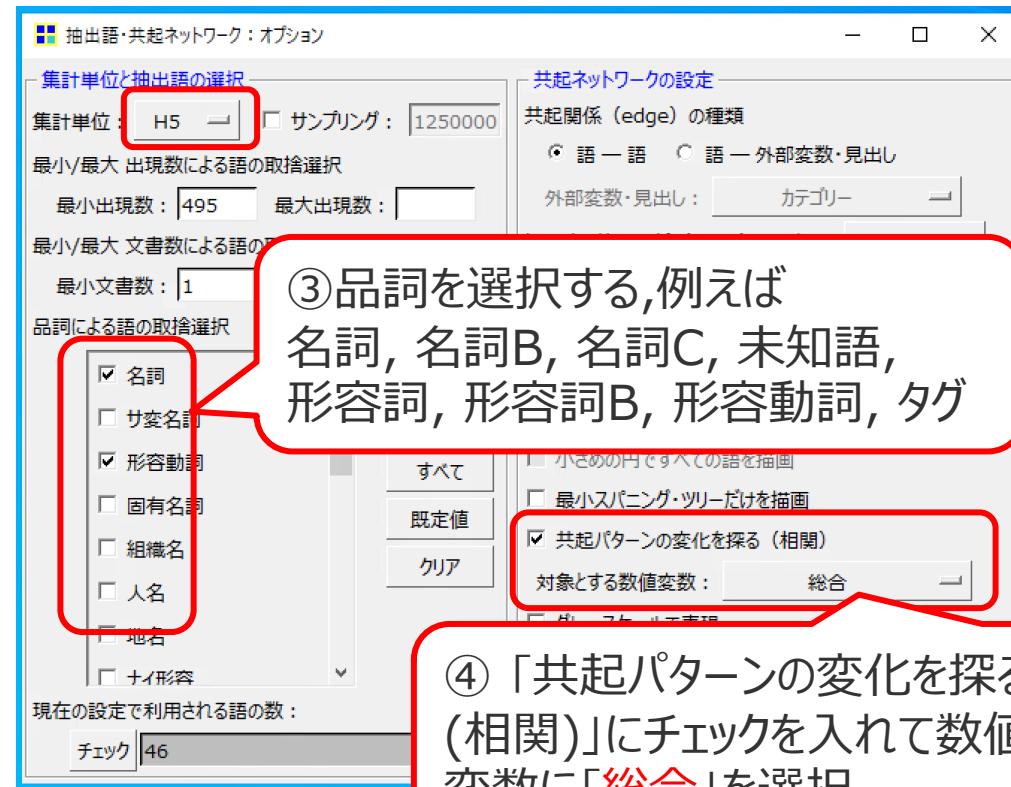
A	B	C	D	E	F	G	H	I	J	K
1										
2	01_登別		02_草津		03_箱根		04_道後			
3	風呂	.115	湯畑	.327	美味しい	.136	温泉	.109		
4	温泉	.107	温泉	.136	露天風呂	.134	立地	.082		
5	美味しい	.094	風呂	.126	風呂	.116	最高	.066		
6	良い	.093	宿	.120	部屋	.109	広い	.063		
7	バギング	.090	美味しい	.102	良い	.106	浴場	.059		
8	残念	.078	良い	.100	温泉	.102	よい	.058		
9	ない	.077	部屋	.096	宿	.097	フロント	.057		
10	夕食	.076	最高	.090	スタッフ	.096	大変	.057		
11	種類	.075	夕食	.085	夕食	.095	夕食	.055		
12	露天風呂	.074	ない	.074	ない	.083	便利	.055		
13	05_湯布院		06_札幌		07_名古屋		08_東京			
14	宿	.180	ホテル	.092	ホテル	.086	駅	.102		
15	美味しい	.144	立地	.077	便利	.072	ホテル	.086		
16	露天風呂	.135	便利	.077	駅	.070	便利	.078		
17	風呂	.127	綺麗	.071	綺麗	.069	立地	.077		
18	温泉	.124	浴場	.070	フロント	.066	近い	.071		
19	最高	.114	フロント	.065	立地	.065	綺麗	.064		
20	スタッフ	.110	広い	.063	近い	.059	快適	.063		
21	家族	.104	快適	.056	アメニティ	.056	コンビニ	.059		
22	部屋	.099	駅	.056	快適	.055	フロント	.055		
23	良い	.097	ベッド	.055	コンビニ	.051	アメニティ	.052		
24	09_大阪		10_福岡							
25	ホテル	.108	ホテル	.090						
26	駅	.096	便利	.087						
27	便利	.080	立地	.082						
28	立地	.074	駅	.074						
29	綺麗	.072	フロント	.072						
30	フロント	.067	綺麗	.067						
31	快適	.064	トイレ	.064						
32	広い	.064	コンビ	.064						
33	近い	.064	よい	.064						
34	アメニティ	.054	快適	.054						

各エリアの特徴語を10件ずつ
一覧(数値はJaccard係数)

● 共起ネットワークの作成(2)

①メニューから「ツール」「抽出語」「共起ネットワーク」を選ぶ

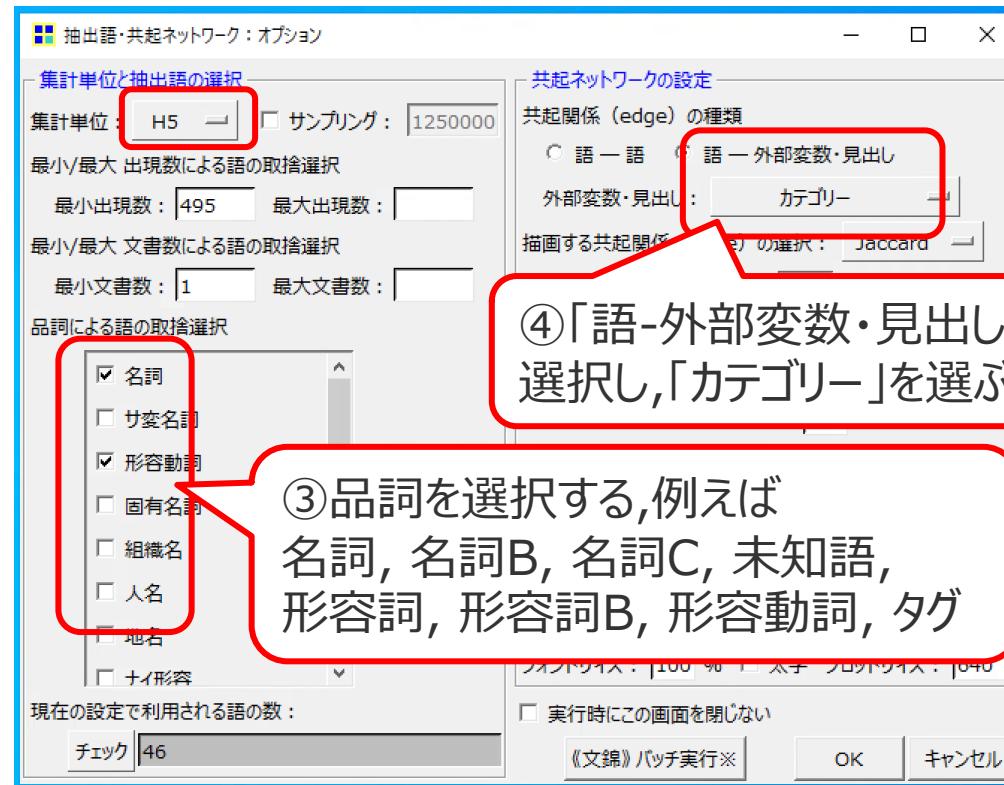
②「集計単位」として「H5」を選んで「OK」をクリック



● 共起ネットワークの作成(3)

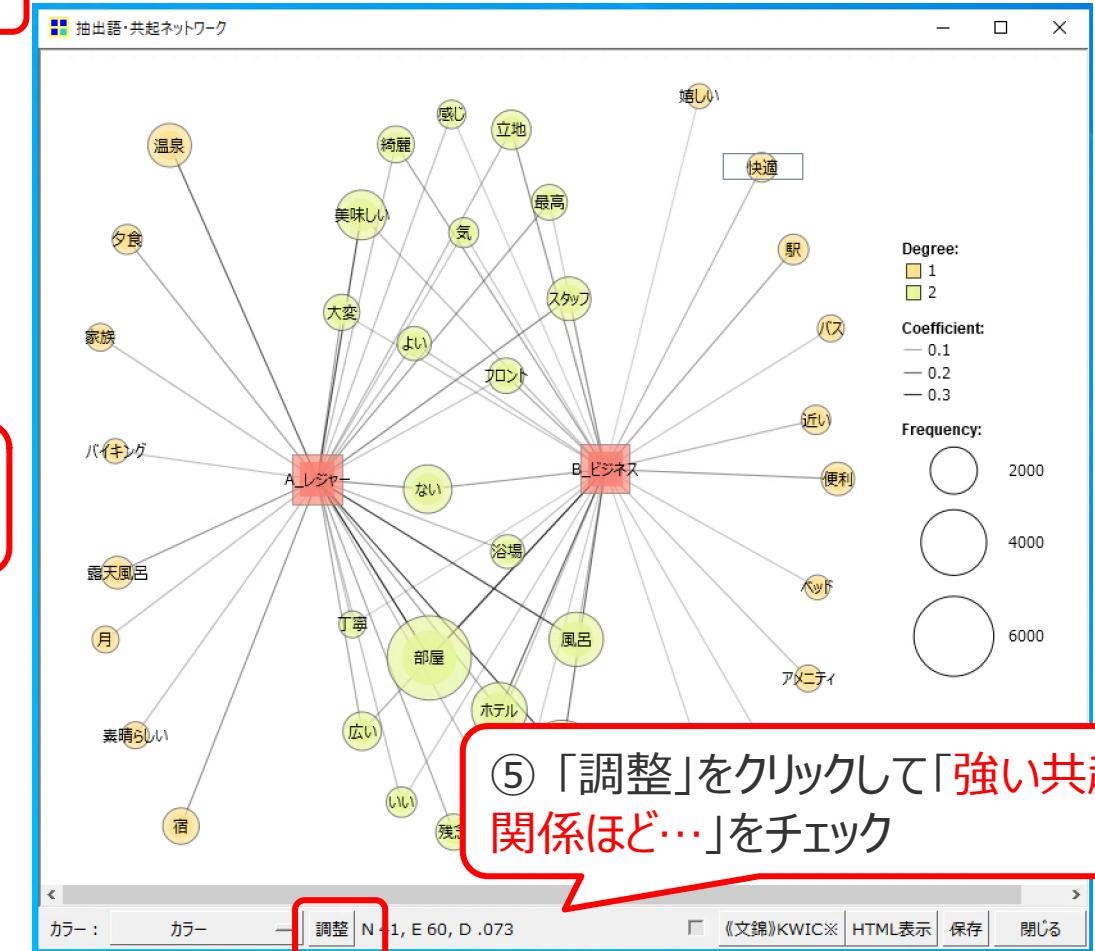
①メニューから「ツール」「抽出語」「共起ネットワーク」を選ぶ

②「集計単位」として「H5」を選んで「OK」をクリック



③品詞を選択する、例えば
名詞, 名詞B, 名詞C, 未知語,
形容詞, 形容詞B, 形容動詞, タグ

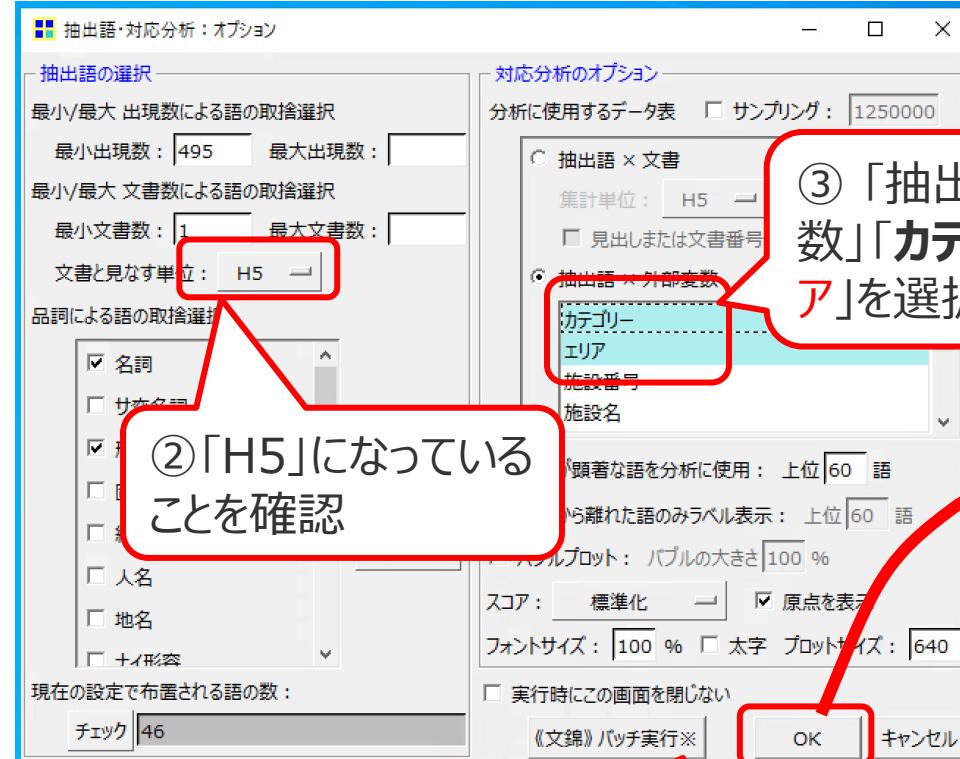
④「語-外部変数・見出し」を選択し、「カテゴリー」を選ぶ



⑤「調整」をクリックして「強い共起
関係ほど…」をチェック

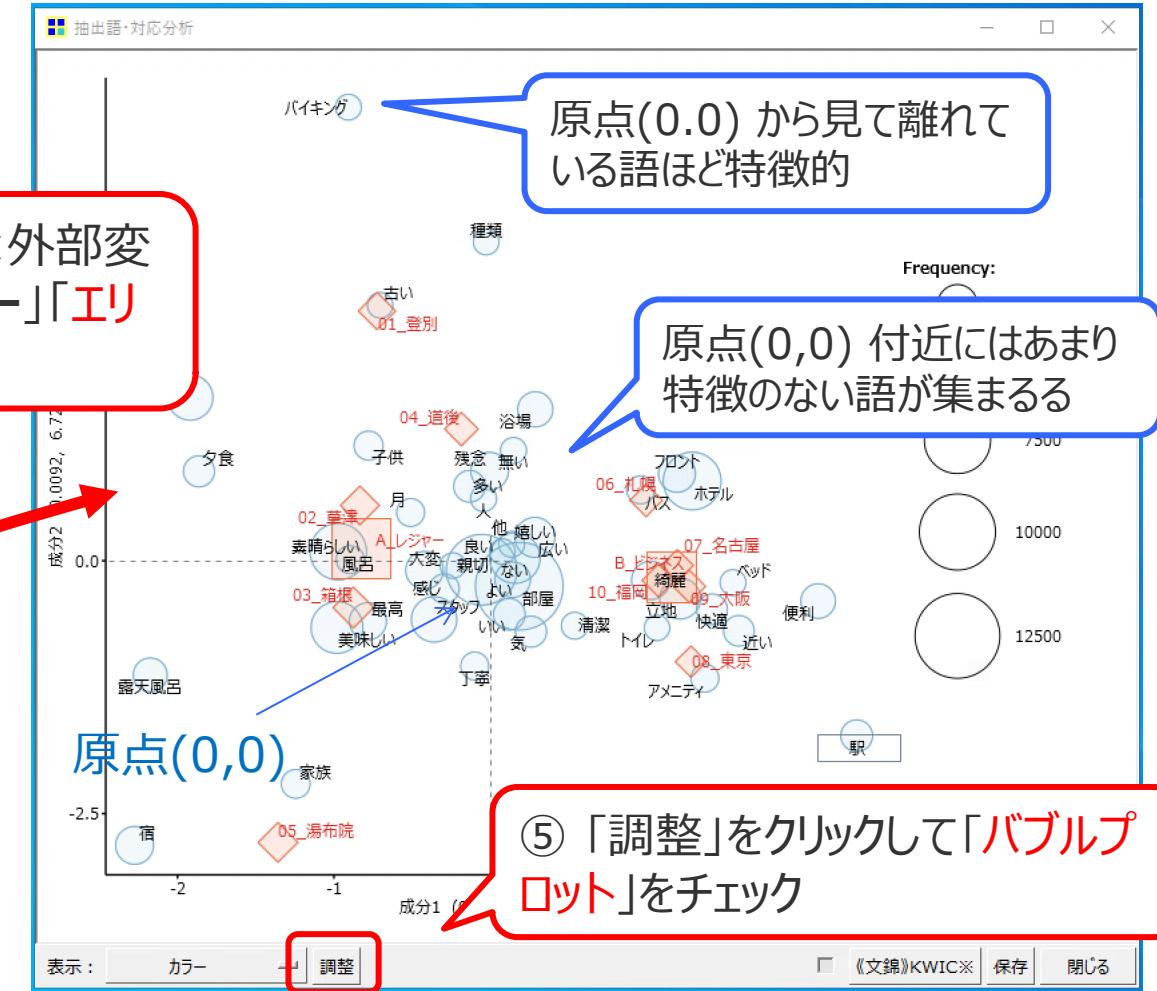
● 対応分析による探索(1)

- ① メニューから「ツール」「抽出語」「対応分析」を選ぶ



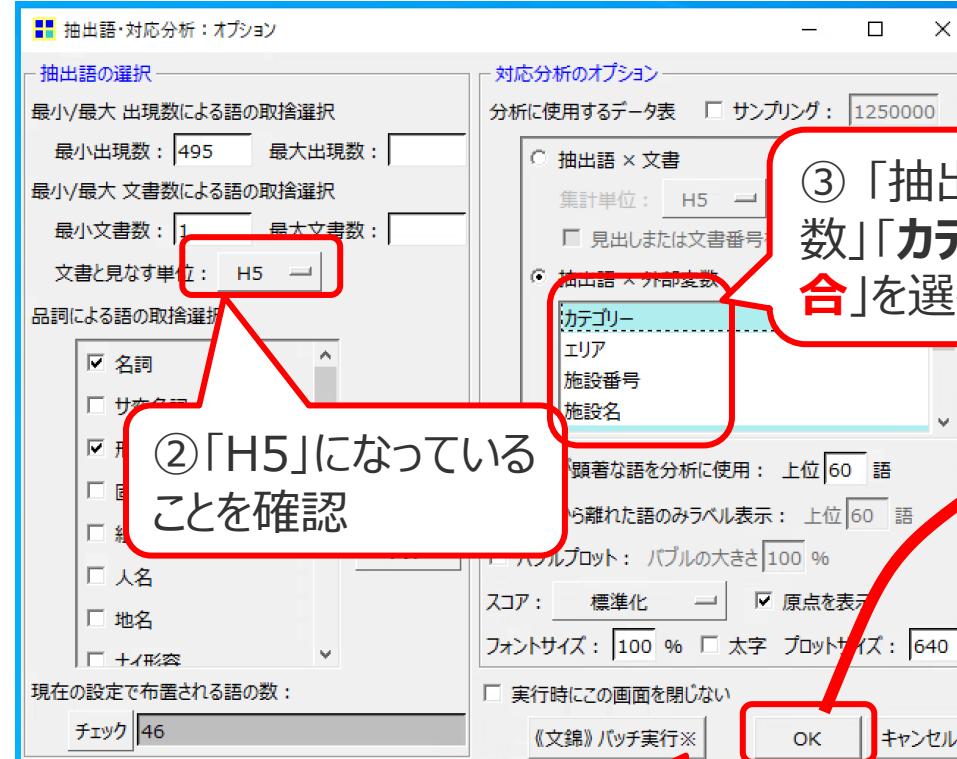
④ 「OK」をクリック

③ 「抽出語×外部変数」「カテゴリー」「エリア」を選択



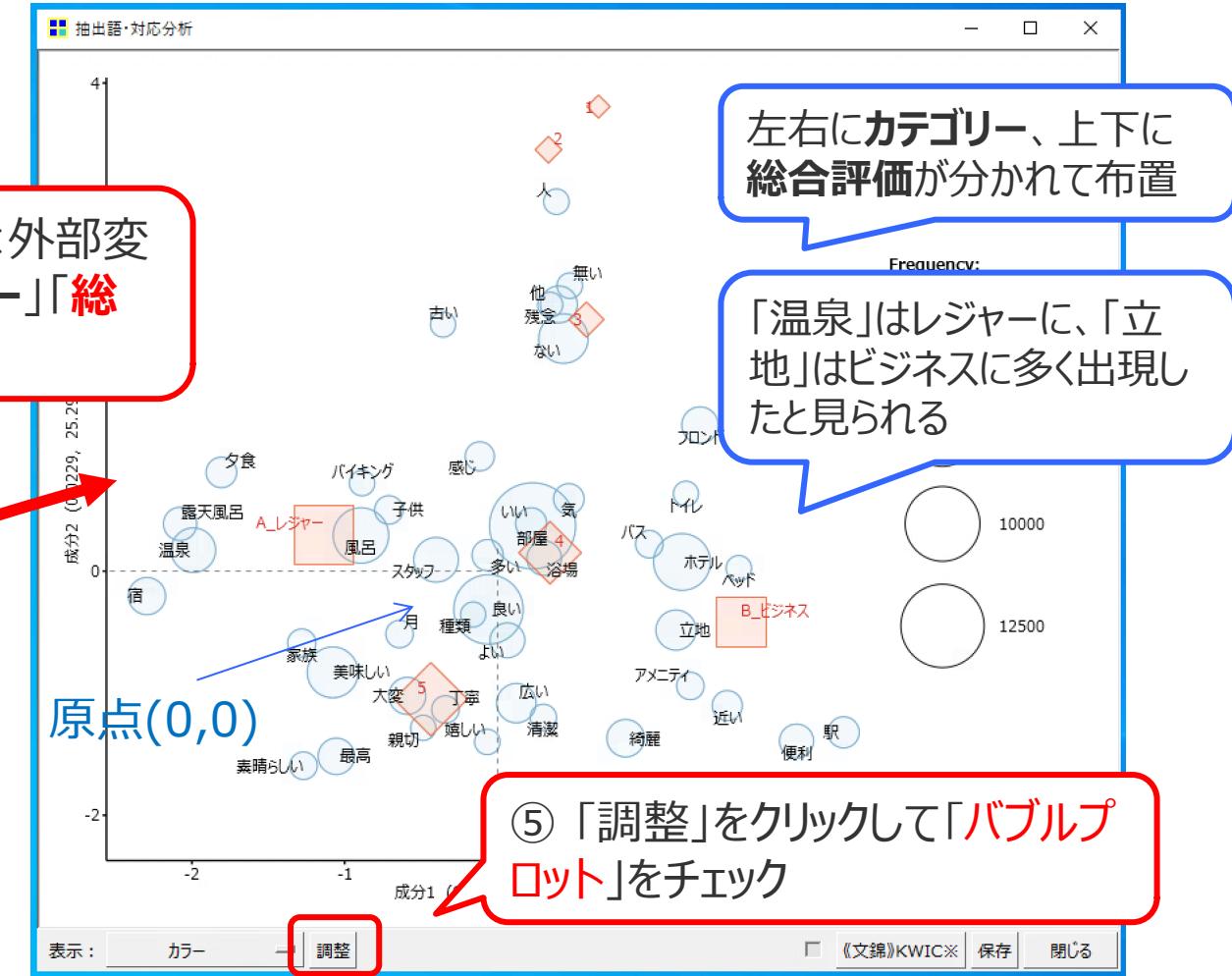
● 対応分析による探索(2)

- ① メニューから「ツール」「抽出語」「対応分析」を選ぶ



③ 「抽出語×外部変数」「カテゴリー」「総合」を選択

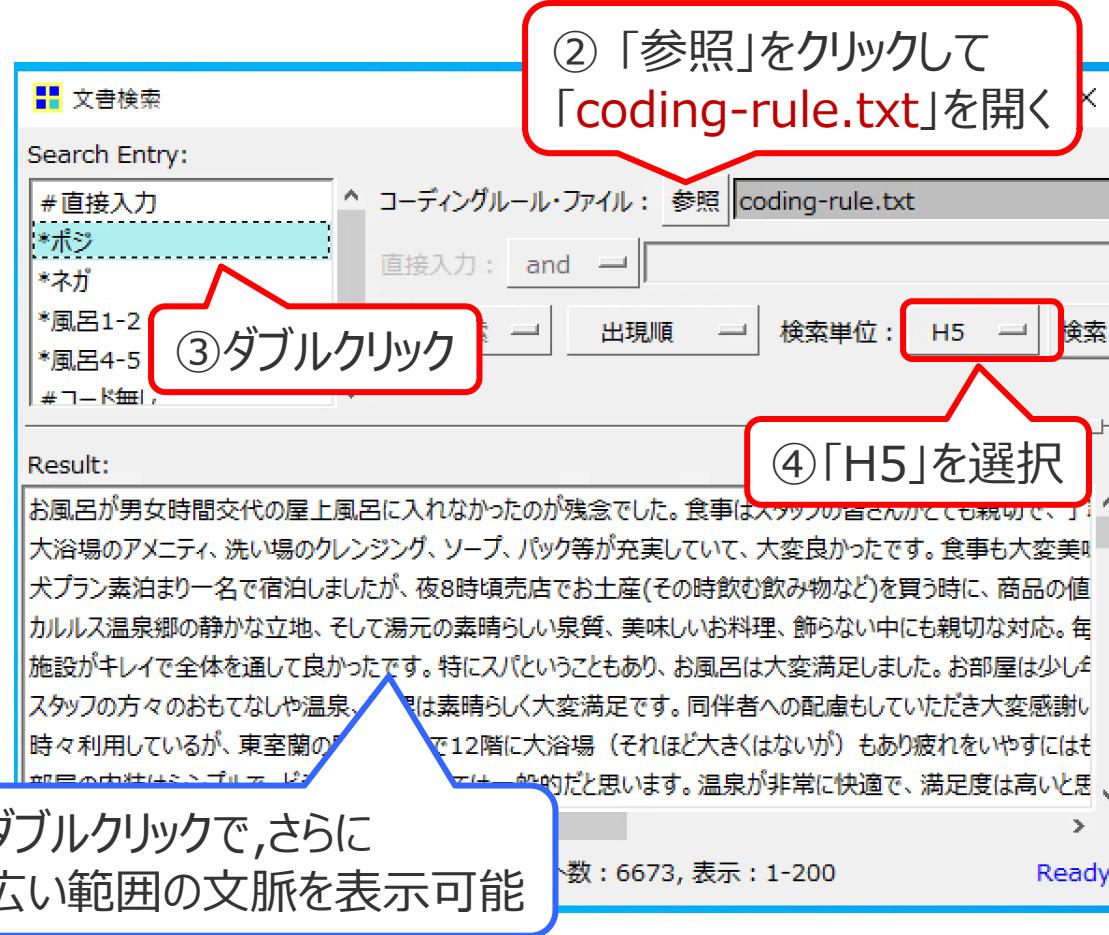
④ 「OK」をクリック



⑤ 「調整」をクリックして「バブルプロット」をチェック

●コーディングルール（語ではなくコンセプトを数える方法）

- ①メニューから「ツール」「文書」「文書検索」を選ぶ



coding-rule.txt の中身

*ポジ

良い or 美味しい or 広い or 多い or 素晴らしい or 嬉しい or 気持ちはよい or 楽しい or 近い or 大きい or 気持ち良い or 温かい or 早い or 優しい or 新しい or 暖かい or 快い or 明るい or 美しい or 可愛い

*ネガ

古い or 無い or 高い or 悪い or 小さい or 狹い or 少ない or 寒い or 遅い or 熱い or 欲しい or 暑い or 冷たい or 遠い or 臭い or 暗い

*風呂1-2

<>風呂-->1 | <>風呂-->2

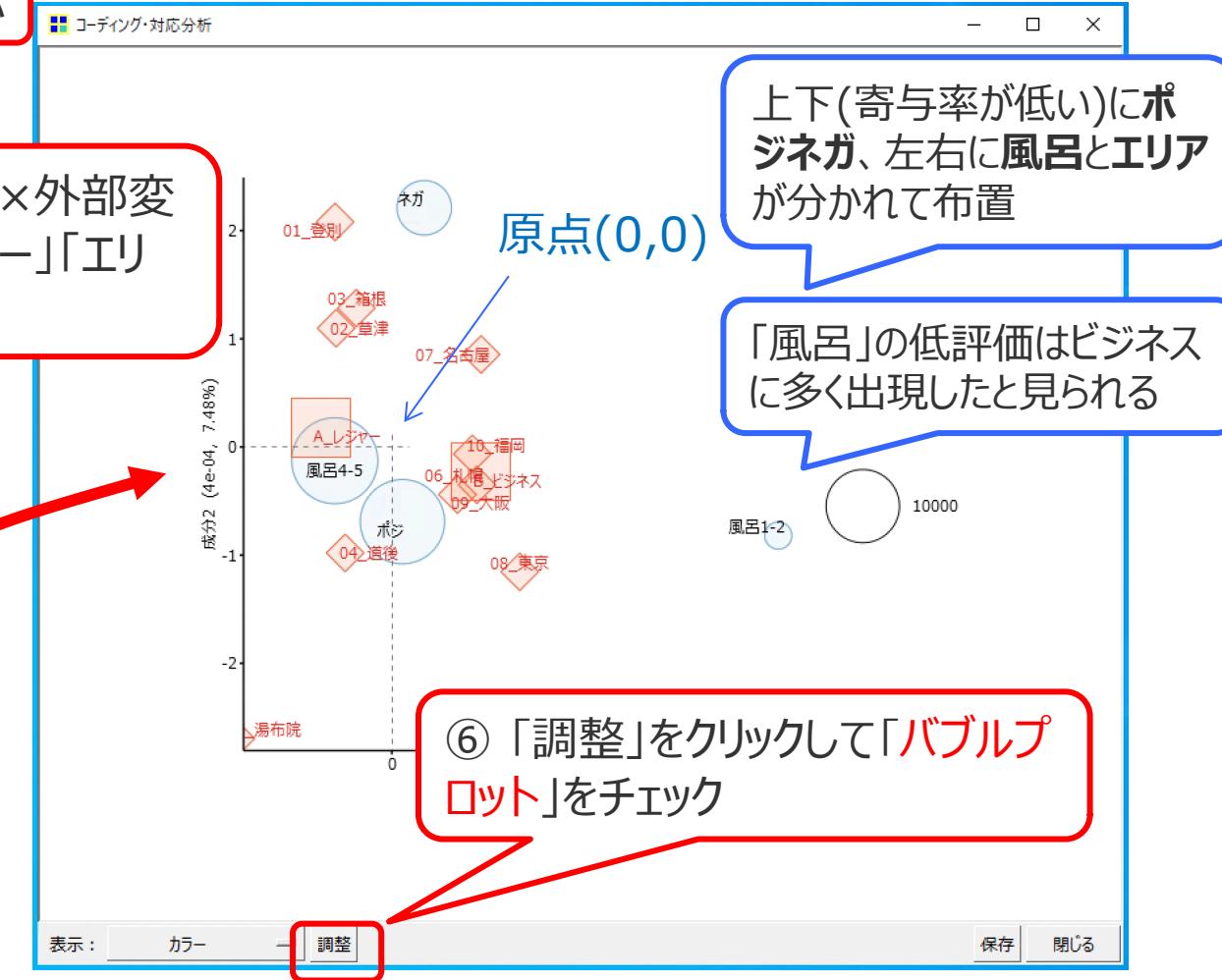
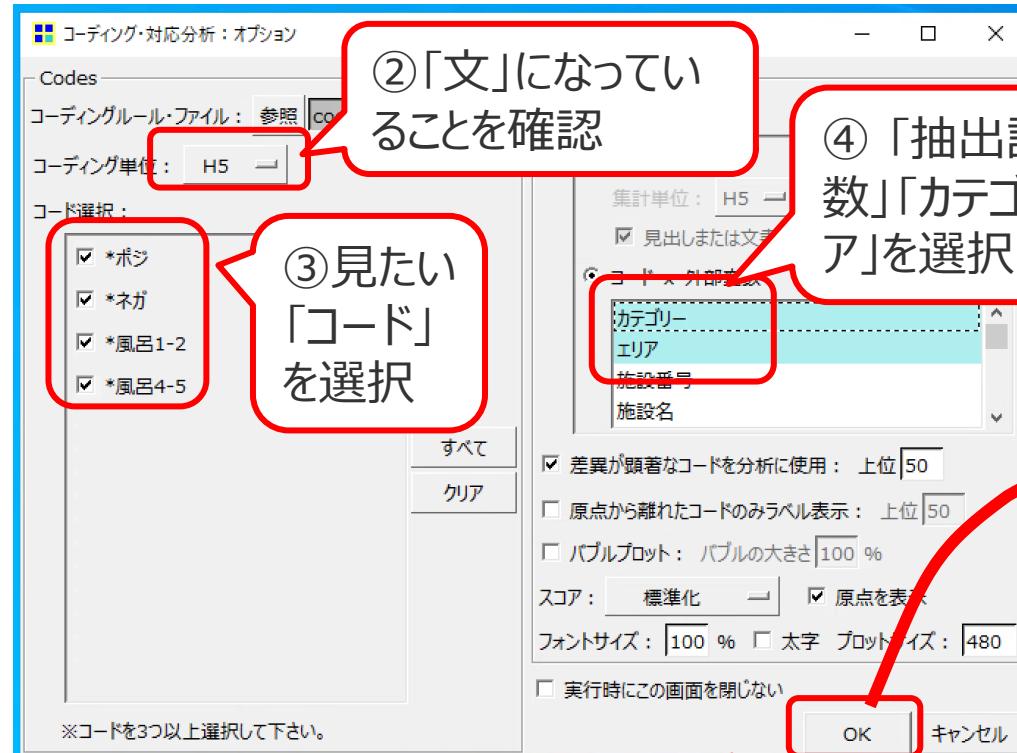
外部
変
数

*風呂4-5

<>風呂-->4 | <>風呂-->5

● 対応分析による探索(3)

- ① メニューから「ツール」「コーディング」「対応分析」を選ぶ



● クロス集計

① メニューから「ツール」「コーディング」「クロス集計」を選ぶ

② 「参照」をクリックして
「coding-rule.txt」を開く

⑤ 「集計」を
クリック

Entry

コーディングルール・ファイル：参照 coding-rule.txt

セル内容：度数とパーセント

コーディング単位：H5 クロス集計：エリア 集計

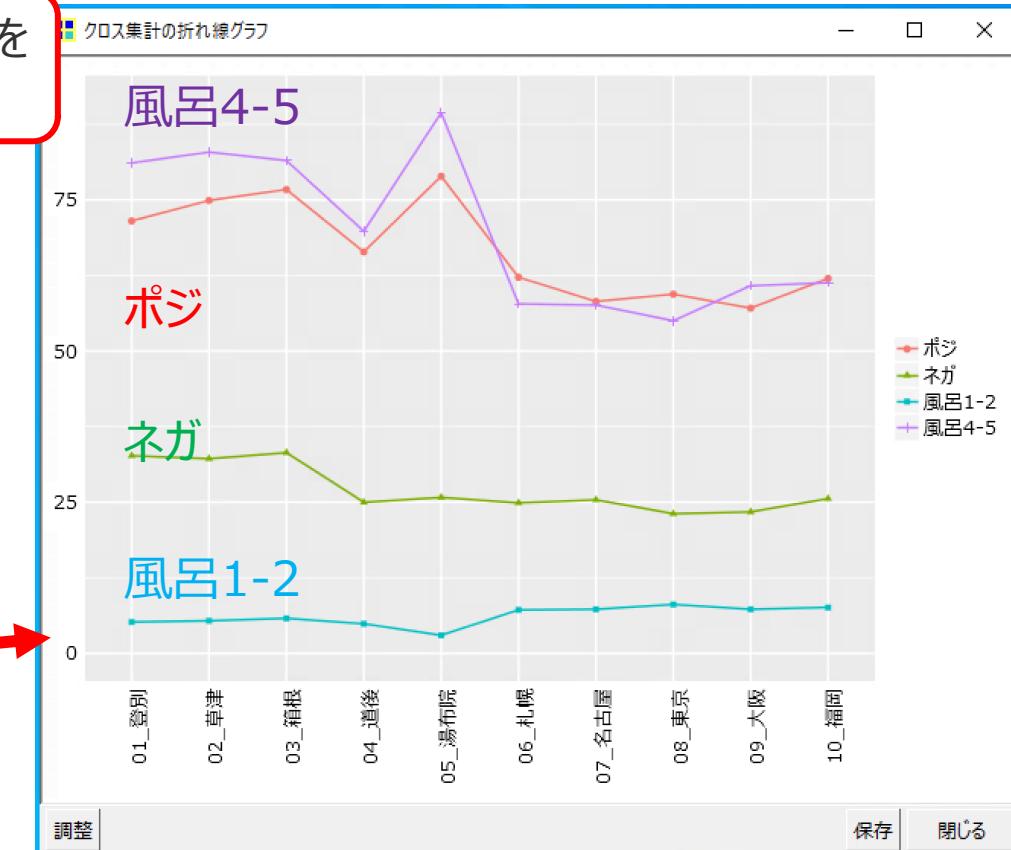
Result

	*ポジ	**	ネガ	ケース数
01_登別	72 (7.20%)	578 (57.80%)	1000	
02_草津	73			
03_箱根			27	
04_道後			25	
05_湯布院			25	
06_札幌			25	
07_名古屋			25	
08_東京			25	
09_大阪	81			
10_福岡	73			
合計	6673 (66.73%)	2713 (27.13%)	618 (6.18%)	10000
カイ2乗値	269.915**	70.988**	39.179**	703.615**

Ready. マップ：ヒート バブル 折れ線：すべて 選択 コピー（表全体）

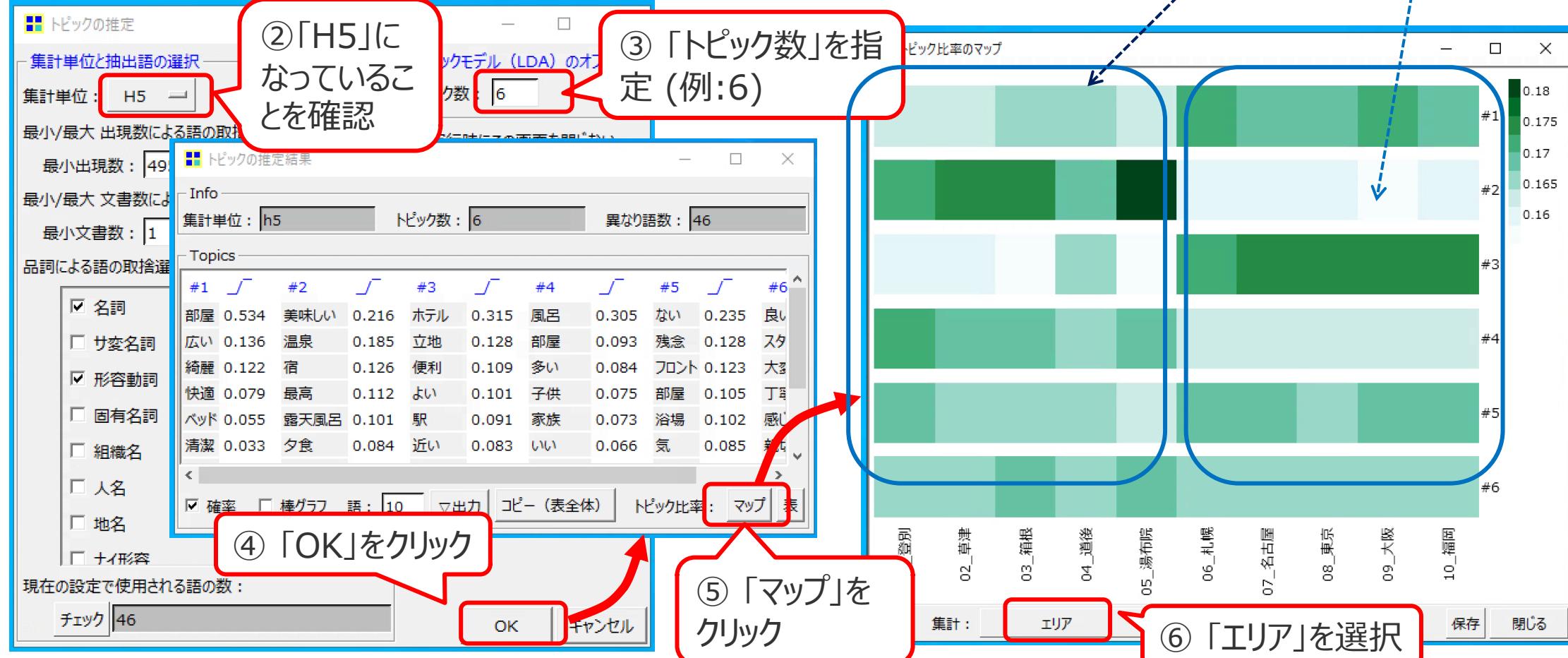
⑥ 「すべて」をクリック

注：プロット左側のラベルは表示されません



● トピックモデルによる分析

- ① メニューから「ツール」「文書」「トピックモデル」「トピックの推定」を選ぶ



KH Coder の解析・分析手法

- 「単語と単語」、「外部変数と単語」の関係に注目した分析が得意

- 特徴的な単語を見つける

- 特定の文書に特徴的な単語を見つける → TF・IDF
→ その文書に特に頻出するが、他の文書ではそれほどではない

- 特徴的な関係を見つける

- 関係性のある単語と単語と見つける → **共起ネットワーク(Jaccard係数)**
例) 単語:「風呂」と 単語:「広い」に関係がありそう
 - 関係性の強い単語と外部変数を見つける → **対応分析(カイ2乗値)**
例) 外部変数:「レジヤー」と 単語:「風呂」に関係がありそう

KH Coder で使われるデータ表

- 「文書-抽出語」表 [【行】ある文中に出現する単語の数を要素とする文ベクトル
【列】全文中に出現する単語の数を要素とする単語ベクトル]

「文書-抽出語」頻度表

「抽出語-抽出語」共起頻度表（共起ネットワーク）

	部屋	良い	ホテル	風呂	美味しい	ない	スタッフ	温泉	よい	立地	広い	綺麗だ	宿	大変だ	残念だ	最高
部屋	4568	1961	1091	1273	1068	1012	783	706	676	691	955	817	461	512	635	526
良い	1961	3701	892	1026	896	692	712	638	339	701	558	494	406	439	468	356
ホテル	1091	892	2041	391	354	506	411	235	267	356	291	339	72	207	283	198
風呂	1273	1026	391	2143	598	503	376	359	333	317	414	280	334	277	325	294
美味しい	1068	896	354	598	1962	379	432	429	229	215	299	255	300	308	249	270
ない	1012	692	506	503	379	1629	337	296	269	289	260	185	202	193	341	159
スタッフ	783	712	411	376	432	337	1411	275	216	201	182	184	206	214	203	175
温泉	706	638	235	359	429	296	275	1298	210	176	221	120	293	197	182	246
よい	676	339	267	333	229	269	216	210	1205	238	167	135	124	124	158	113
立地	691	701	356	317	215	289	201	176	238	1328	167	186	120	139	157	243
広い	955	558	291	414	299	260	182	221	167	167	1186	209	122	153	152	128
綺麗だ	817	494	339	280	255	185	184	120	135	186	209	1132	76	120	119	117

「外部変数-抽出語」クロス集計表（対応分析）

	部屋	良い	ホテル	風呂	美味しい	スタッフ	温泉	よい	立地	広い	綺麗な宿	大変だ	残念だ	最高
A_レジヤー	2398	2046	757	1535	1430	880	888	1188	631	518	631	459	769	652
B_ビジネス	2170	1655	1284	608	532	749	523	110	574	810	555	673	57	355
01_登別	447	409	194	323	255	187	148	222	114	38	125	88	76	120
02_草津	488	434	181	352	274	180	154	275	117	155	114	109	195	136
03_箱根	548	436	134	326	355	202	212	212	133	57	140	93	161	137
04_道後	416	349	191	181	174	130	135	225	137	176	129	87	59	108
05_湯布院	499	418	57	353	372	181	239	254	130	92	123	82	278	151
06_札幌	452	346	255	121	129	151	114	38	103	166	129	142	10	87
07_名古屋	434	310	241	116	97	144	102	18	133	141	84	138	11	58
08_東京	441	338	240	106	99	131	99	12	104	166	98	128	14	69
09_大阪	431	317	297	135	88	162	93	20	121	161	132	144	9	62
10_福岡	412	344	251	130	119	161	115	22	113	176	112	121	13	79

カイ2乗値

- カイ2乗値は「無関係でない」度合いを測る尺度 → カテゴリと変数間の関連性を測定

$$\text{カイ2乗値} = \frac{(\text{観測度数}_{ij} - \text{期待度数}_{ij})^2}{\text{期待度数}_{ij}}$$

「観測度数」: カテゴリと変数に従ってクロス集計された度数

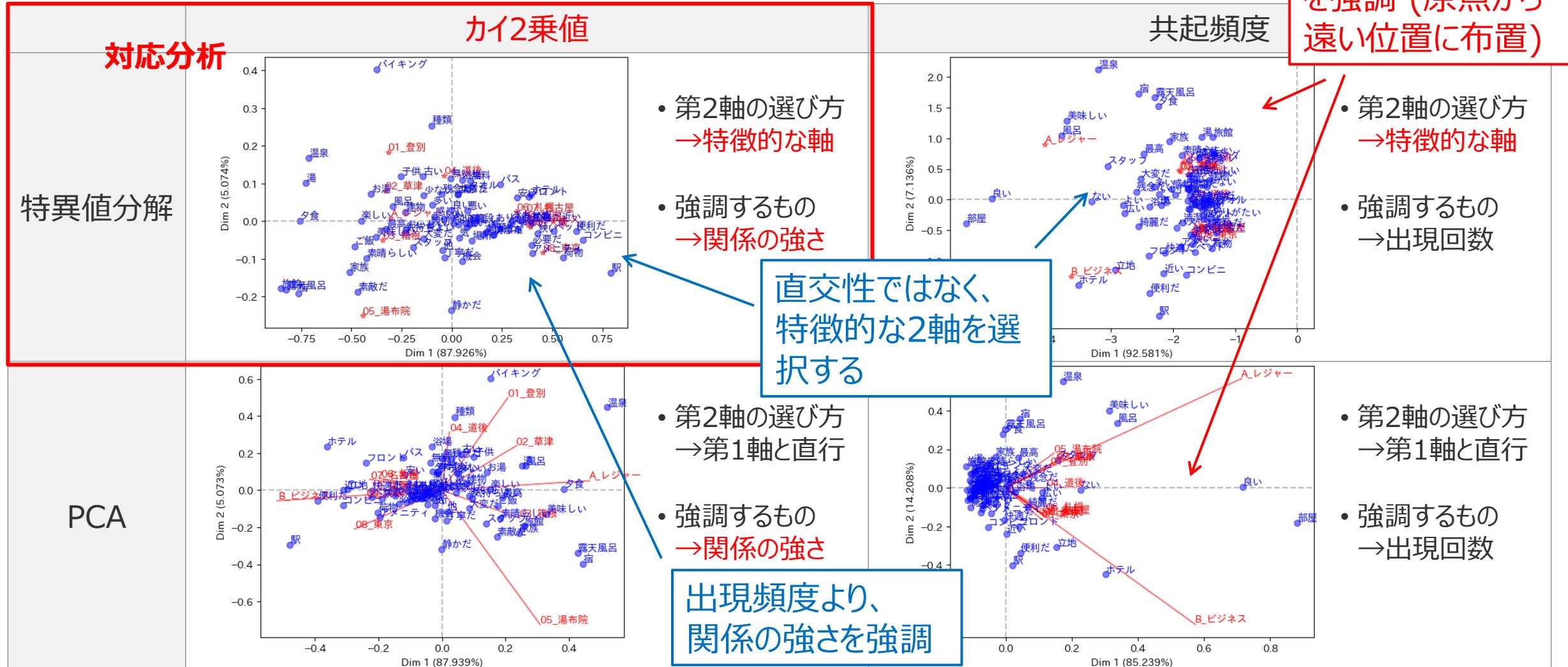
「期待度数」: 変数が互いに独立している場合に期待される度数

「観測度数 - 期待度数」: 実際の度数と独立と期待される度数の差

- カイ2乗値も大きい → カテゴリと変数間の関係が**期待より強い**を示す

クロス集計表 (観測度数)						期待度数					観測度数-期待度数					カイ2乗値				
	A	B	C	D	E	合計		A	B	C	D	E	合計		A	B	C	D	E	合計
地質学	3	19	39	14	10	85	地質学	3.310	13.668	33.103	13.775	21.143	85.000	地質学	-0.310	5.332	5.897	0.225	-11.143	0.000
生物化学	1	2	13	1	12	29	生物化学	1.129	4.663	11.294	4.700	7.214	29.000	生物化学	-0.129	-2.663	1.706	-3.700	4.786	0.000
科学	6	25	49	21	29	130	科学	5.063	20.905	50.628	21.068	32.337	130.000	科学	0.937	4.095	-1.628	-0.068	-3.337	0.000
動物学	3	15	41	35	26	120	動物学	4.673	19.296	46.734	19.447	29.849	120.000	動物学	-1.673	-4.296	-5.734	15.553	-3.849	0.000
物理学	10	22	47	9	26	114	物理学	4.440	18.332	44.397	18.475	28.357	114.000	物理学	5.560	3.668	2.603	-9.475	-2.357	0.000
工学	3	11	25	15	34	88	工学	3.427	14.151	34.271	14.261	21.889	88.000	工学	-0.427	-3.151	-9.271	0.739	12.111	0.000
微生物学	1	6	14	5	11	37	微生物学	1.441	5.950	14.410	5.996	9.204	37.000	微生物学	-0.441	0.050	-0.410	-0.996	1.796	0.000
植物学	0	12	34	17	23	86	植物学	3.349	13.829	33.492	13.937	21.392	86.000	植物学	-3.349	-1.829	0.508	3.063	1.608	0.000
統計学	2	5	11	4	7	29	統計学	1.129	4.663	11.294	4.700	7.214	29.000	統計学	0.871	0.337	-0.294	-0.700	-0.214	0.000
数学	2	11	37	8	20	78	数学	3.038	12.543	30.377	12.641	19.402	78.000	数学	-1.038	-1.543	6.623	-4.641	0.598	0.000
合計	31	128	310	129	198	796	合計	31.000	128.000	310.000	129.000	198.000	796.000	合計	0.000	0.000	0.000	0.000	0.000	0.000

- 対応分析は、カイ2乗値を利用して、関係の強さを強調して表現できる

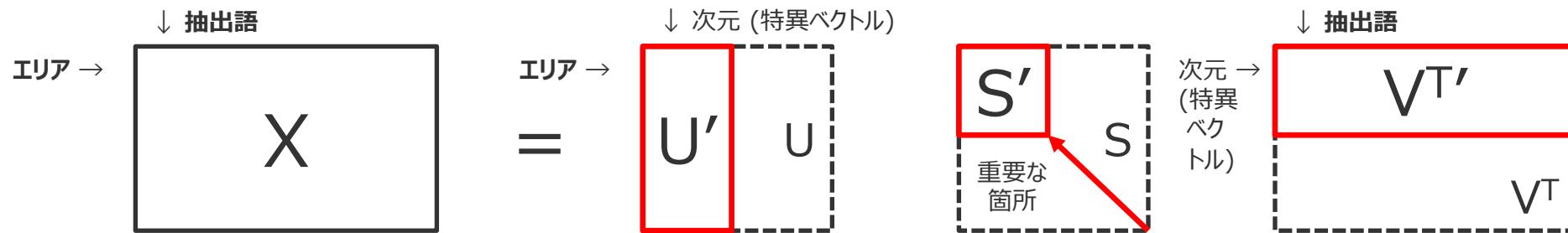


特異値分解

● 特異値分解 $X = USV^T$

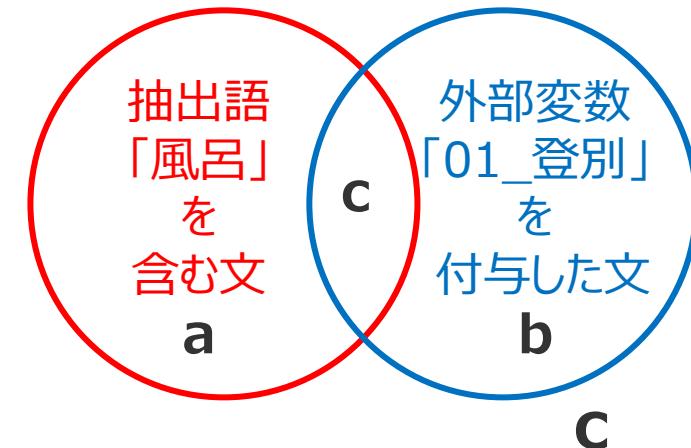
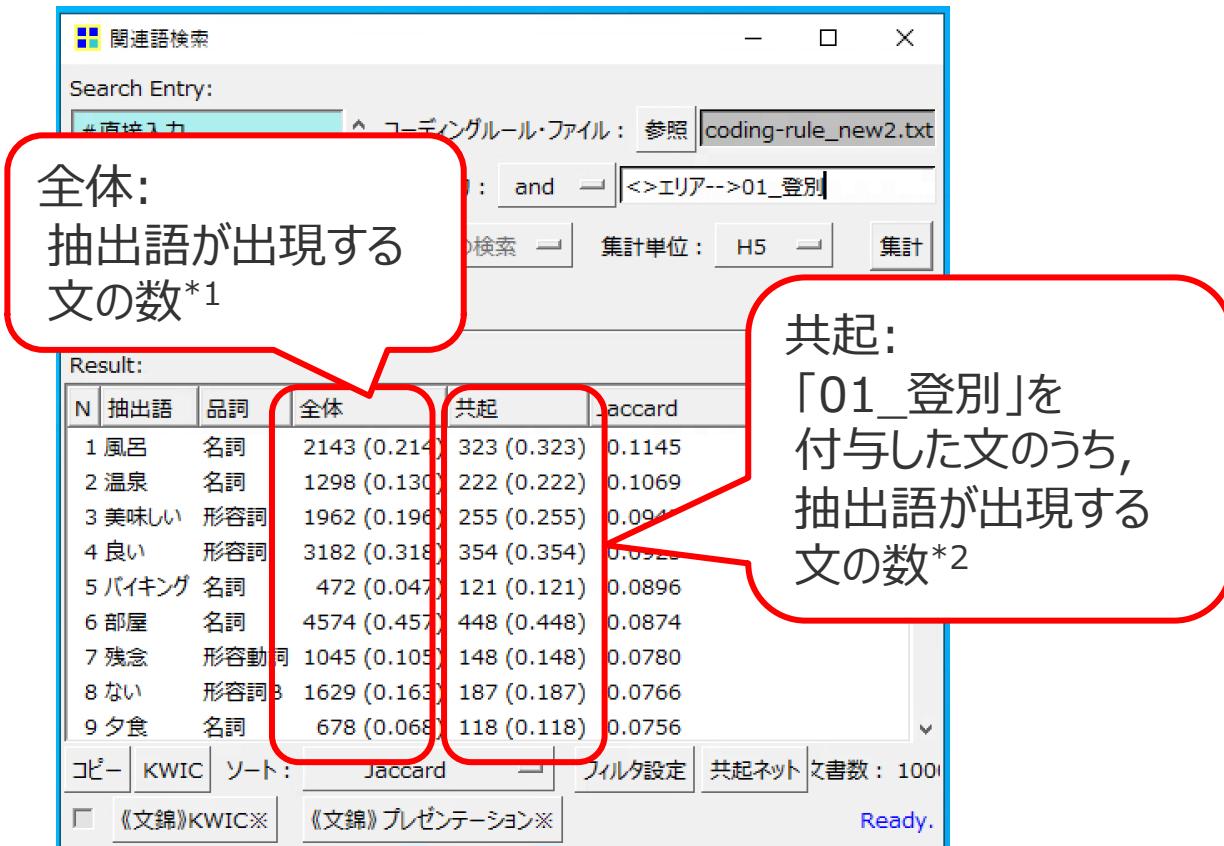


● S の特異値が小さいものを削る



Jaccard 系数

- Jaccard 系数は、共起の強さを測る尺度 (KH Coderで標準的に使用)
 - どちらの語も含まない文書を無視 → 言語のようなスパースデータ分析に向いている



$$\text{Jaccard 系数} = \frac{c}{a+b+c}$$

抽出語「風呂」の場合:

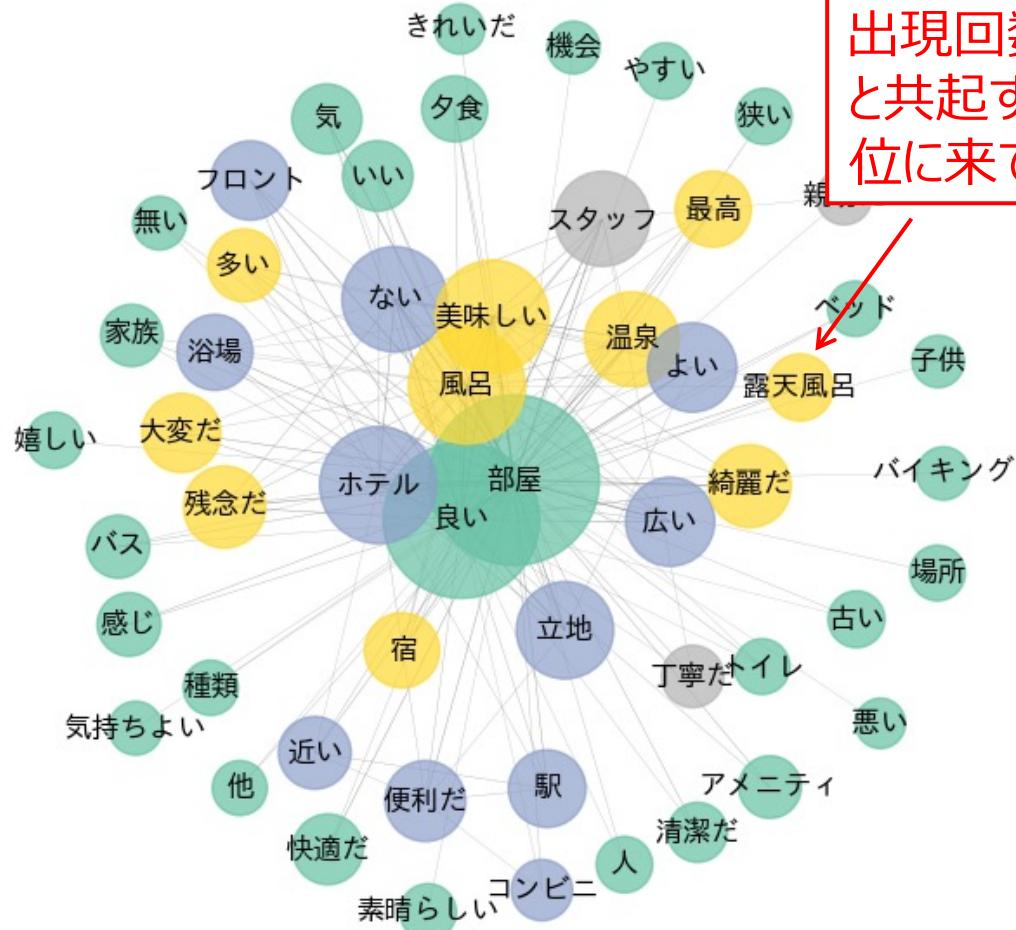
$$c = 323 \text{ ("共起"列の値)}$$

$$a = 2143 \text{ ("全体"列の値)} - 323 = 1820$$

$$b = (323 / 0.323) - 323 = 677$$

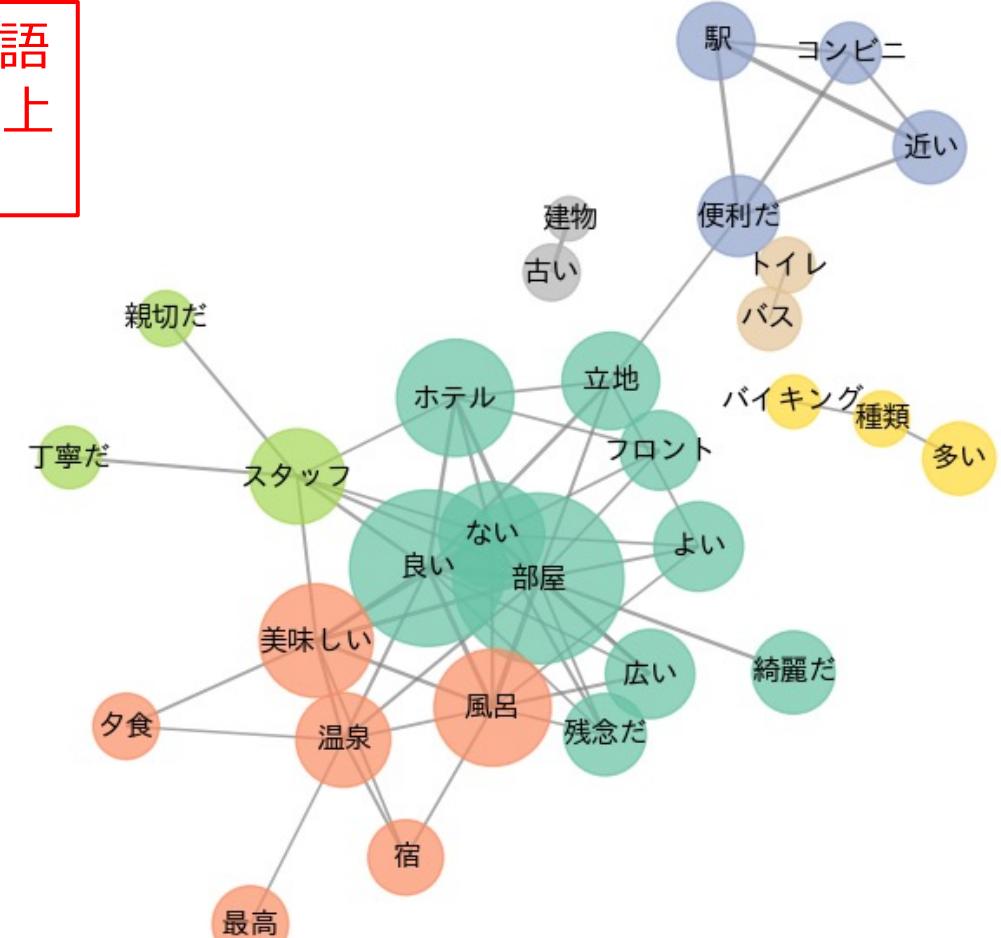
^{*1} 括弧内はデータ全体に対する割合(前提確率) ^{*2} 括弧内は「01_登別」を付与したデータに対する割合(条件付き確率)

- 共起頻度の高いエッジを残す



出現回数の高い語
と共に起するだけで上
位に来てしまう

- Jaccard 係数が上位のエッジを残す

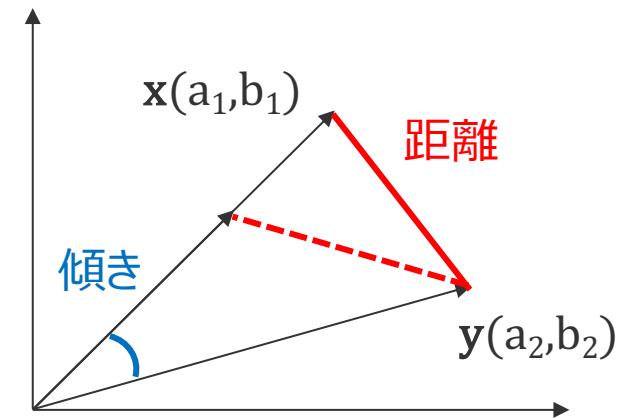


その他の尺度

- 出現パターンが似てる を測る = ユークリッド距離、コサイン距離
- 1つひとつの文が長く、各文中での語の出現回数の大小が重要なケースに向く
(語が1回出現したか、10回出現したかを区別したい)

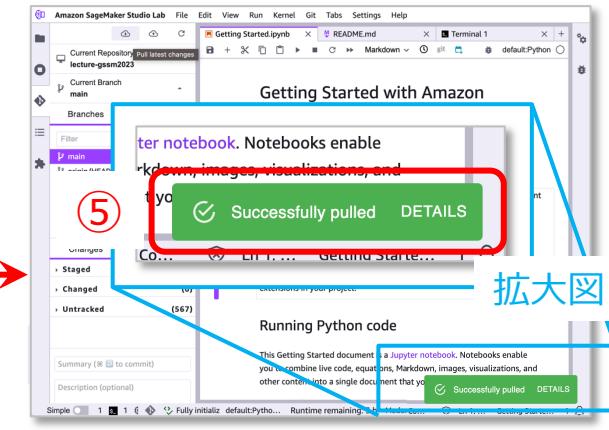
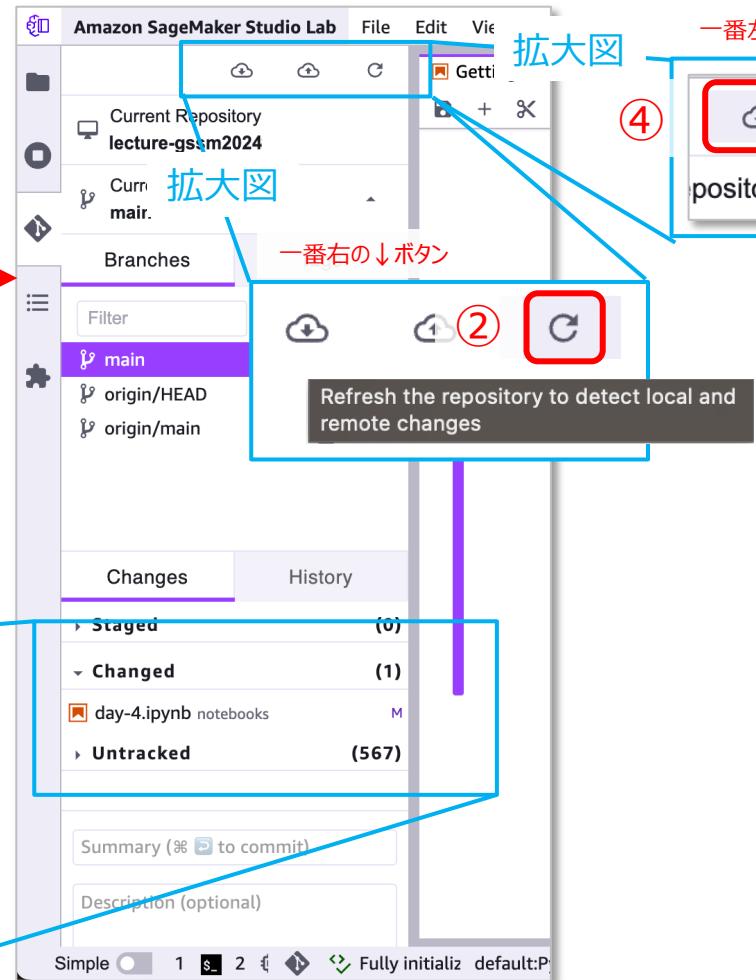
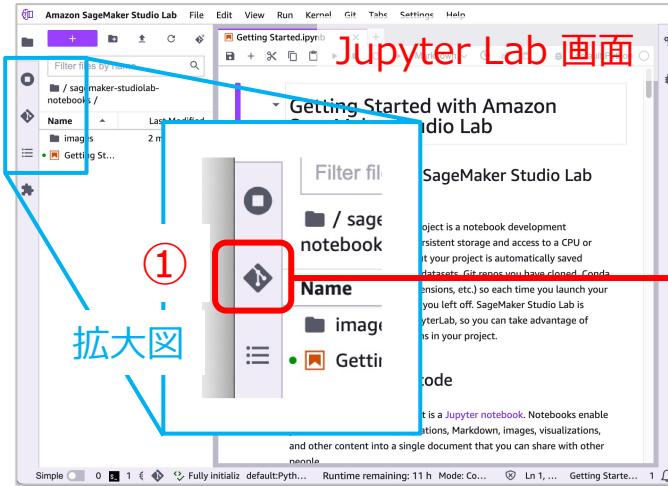
ユークリッド距離	コサイン距離
サイズ(出現回数の大小)の差まで見る場合向き	傾きが似ているかどうかだけを見る場合向き
$d(x, y) = \sqrt{\sum (x_i - y_i)^2}$	$d(x, y) = \frac{\sum x_i y_i}{\sqrt{\sum x_i^2 \sum y_i^2}}$

※ x, y はそれぞれの単語ベクトル (単語の出現パターン)

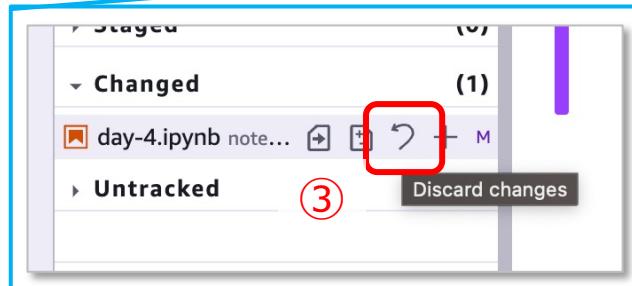


演習 — テキスト解析 (2)

● Jupyter Lab を起動して、教材を最新化(pull)してください



(例) 競合がある場合のみ 拡大図



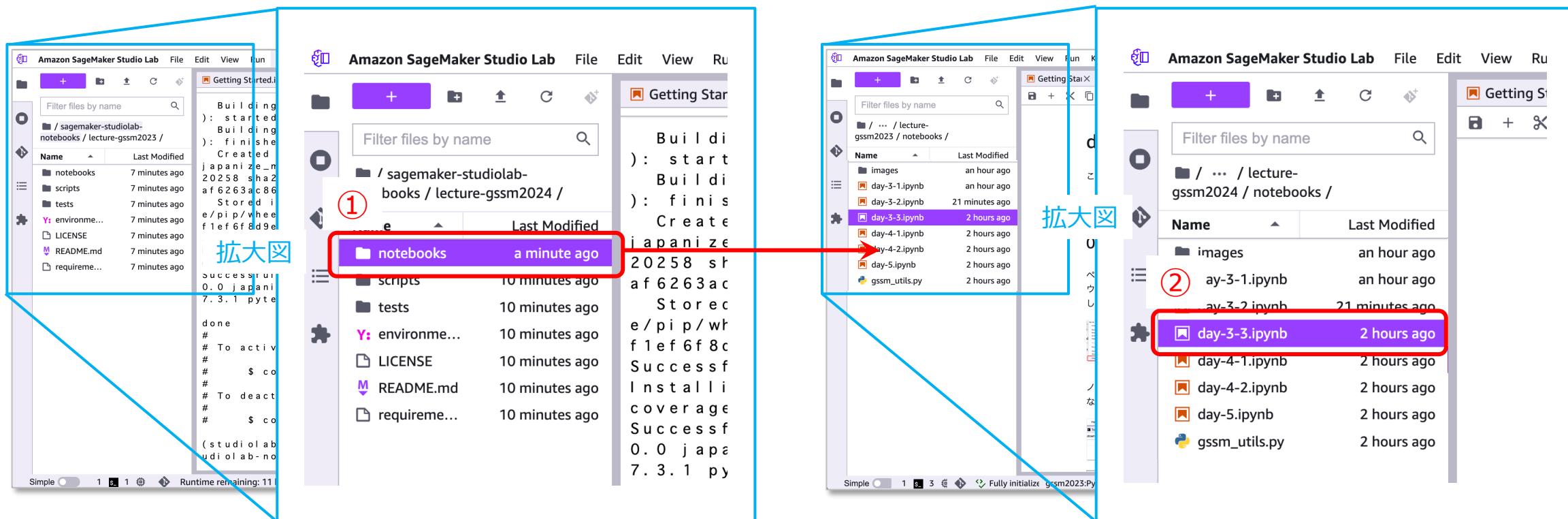
● 教材を最新化する

- ① 「Git」 ボタンを押す
- ② 「Refresh」 ボタンを押す
- ③ もし、競合がある場合(Changedが0でない場合)、
対象のファイルを手動で退避した後、
「Discard changes」 ボタンを押し
て変更を破棄する
- ④ 「Pull latest changes」 を押す
- ⑤ 画面の右下に「Successfully published」が表示されること確認する

演習 — テキスト解析 (2)

● day-3-3.ipynb を開いてください

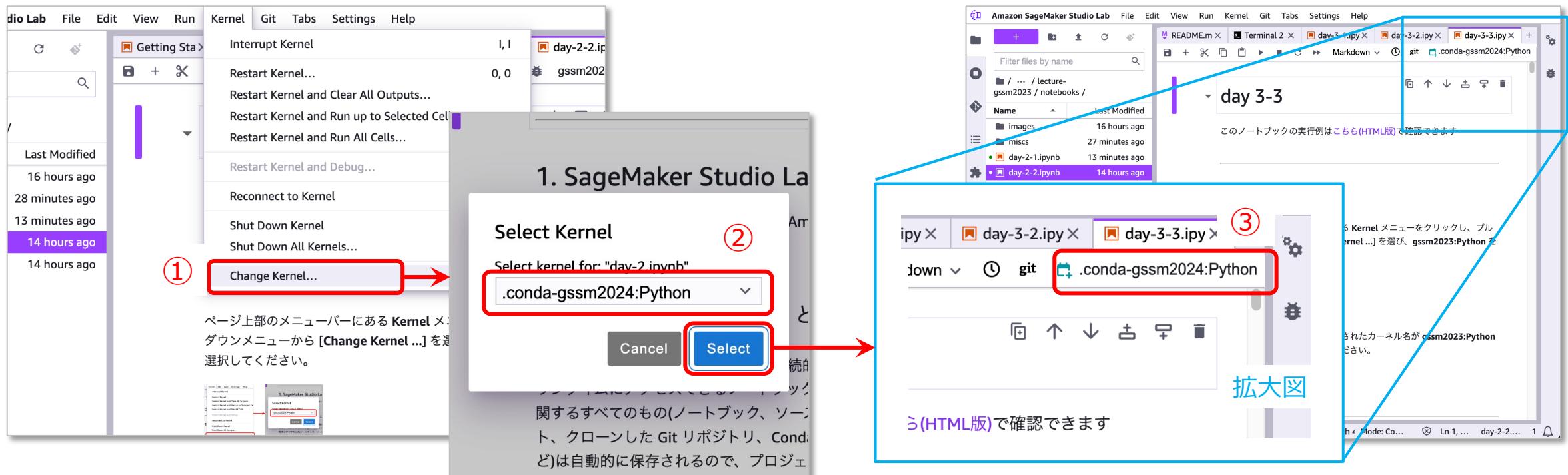
- ① 画面左の **File Browser** から ① **notebooks** をフォルダを開く (既に開いている場合はスキップ)
- ② 次に **day-3-3.ipynb** ノートブックを開く



演習 — テキスト解析 (2)

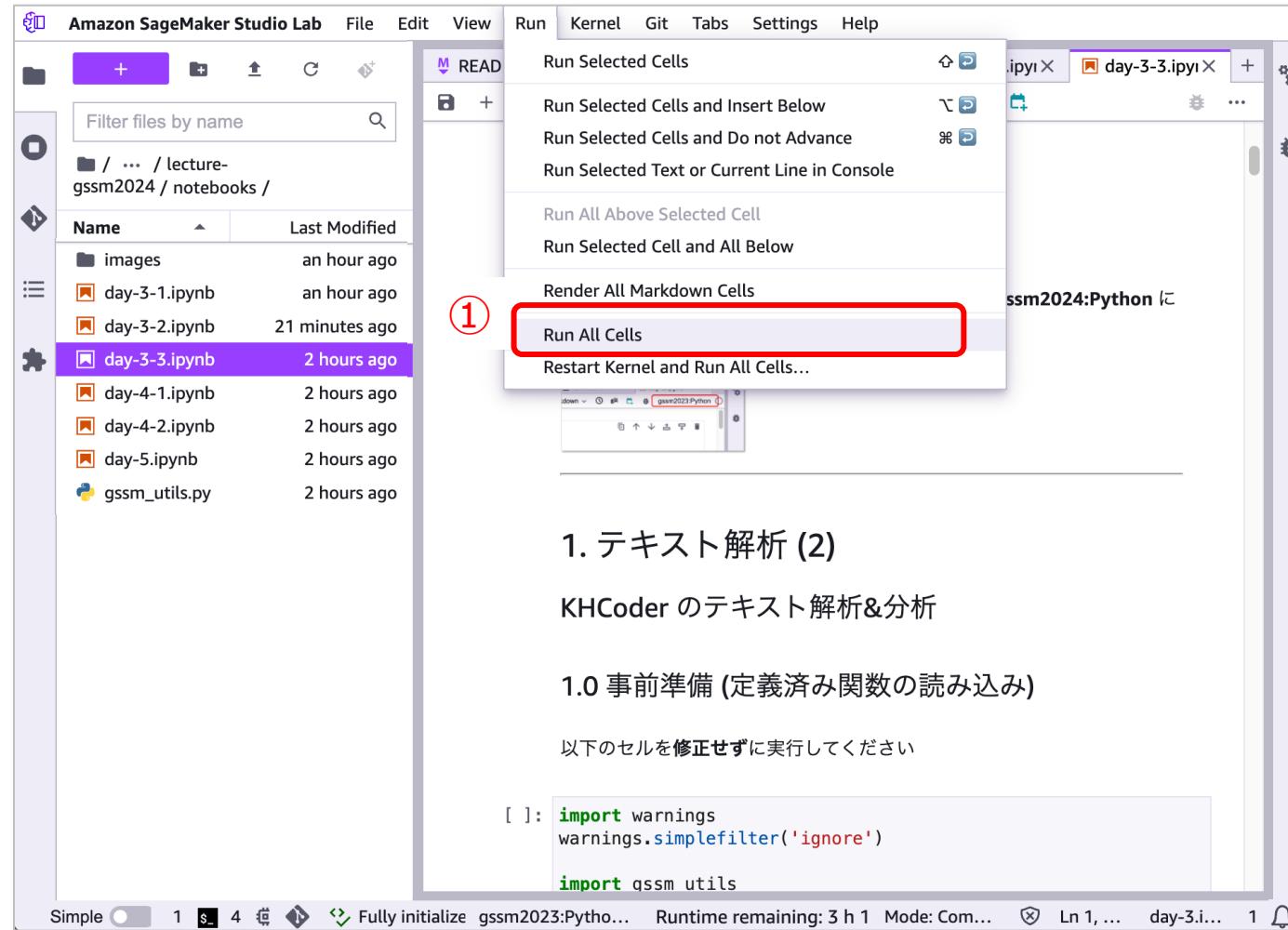
● カーネル **.conda-gssm2024:Python** を選択してください !重要!

- ① ページ上部の **Kernel** メニューから「**Change Kernel ...**」を選ぶ
- ② ポップアップ画面から「**.conda-gssm2024:Python**」を選択し、「**Select**」を押す
- ③ 右上隅にカーネル名「**.conda-gssm2024:Python**」が表示されていることを確認する



演習 — テキスト解析 (2)

● KH Coder のテキスト解析&分析 (day-3-3.ipynb)



演習:

- ① ページ上部の Run メニューから「Run All Cells」を選ぶ

この後、Step-by-step で解説します

day 3 – レポート課題

- 以下を PDF ファイルで提出してください
 - ノートブック **day-3-3.ipynb** の末尾にある「【演習】2021~2022 データセット」に従って、別のデータセット (rakuten-1000-**2021-2022.xlsx.zip**) で作図した「共起ネットワーク図」と「対応分析プロット」のキャプチャ

※ 何らかの事情で上記のキャプチャを提出できない場合、本日の講義の感想を文章で記述してください

レポート形式	提出先	期限
PDF	manaba	次回～18:20