

階層的トピックモデルに基づく 宿泊レビューの時間変化分析

第23回インタラクティブ情報アクセスと
可視化マイニング研究会 (SIG-AM)
2019/11/23(土)

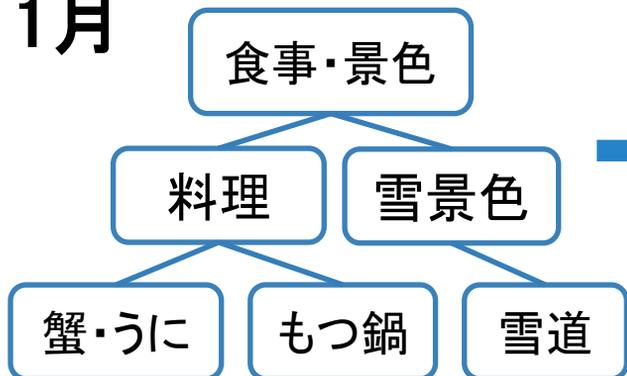
佐藤 裕次郎* (立命館大学)
山西 良典 (立命館大学)
西原 陽子 (立命館大学)

1 研究の概要

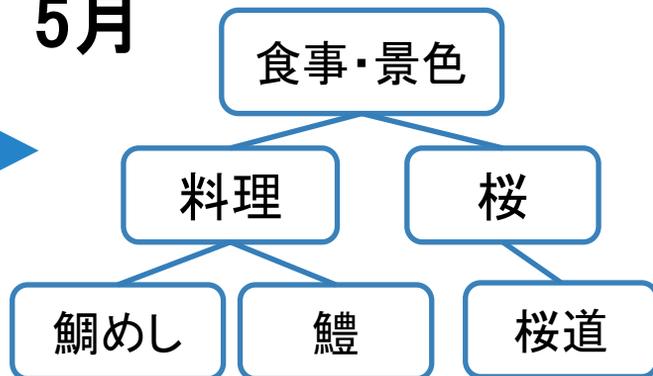
宿泊レビューから

季節の特徴を階層構造で抽出

1月



5月



1 研究背景

- 宿泊レビューはユーザの感想が得られる情報源
- 人気な宿泊施設はレビュー数が多く、欲しい情報が見つからない

クチコミ・お客さまの声(5861件)

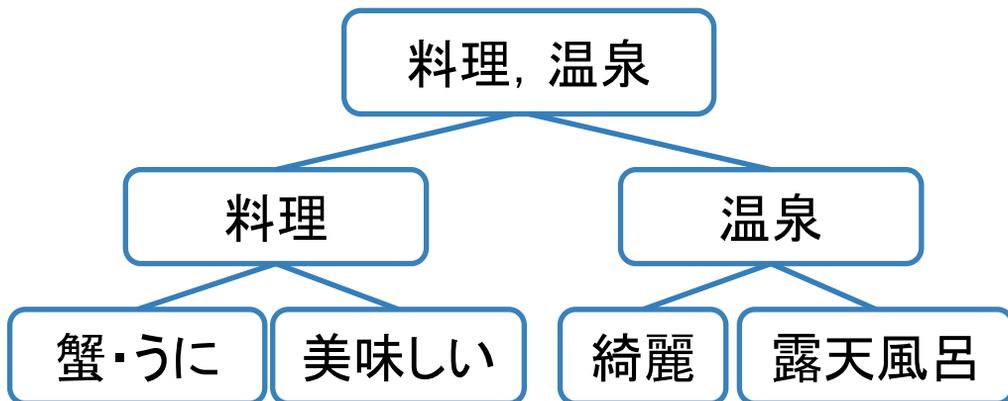
このうち最近の**7,8件**
しか表示されない

- 季節の変化
- 泊まる予定の季節

1 研究目標

- 宿泊施設の特徴
 - カテゴリーの階層構造

この構造で**季節の特徴**を抽出



1. 研究背景・目標

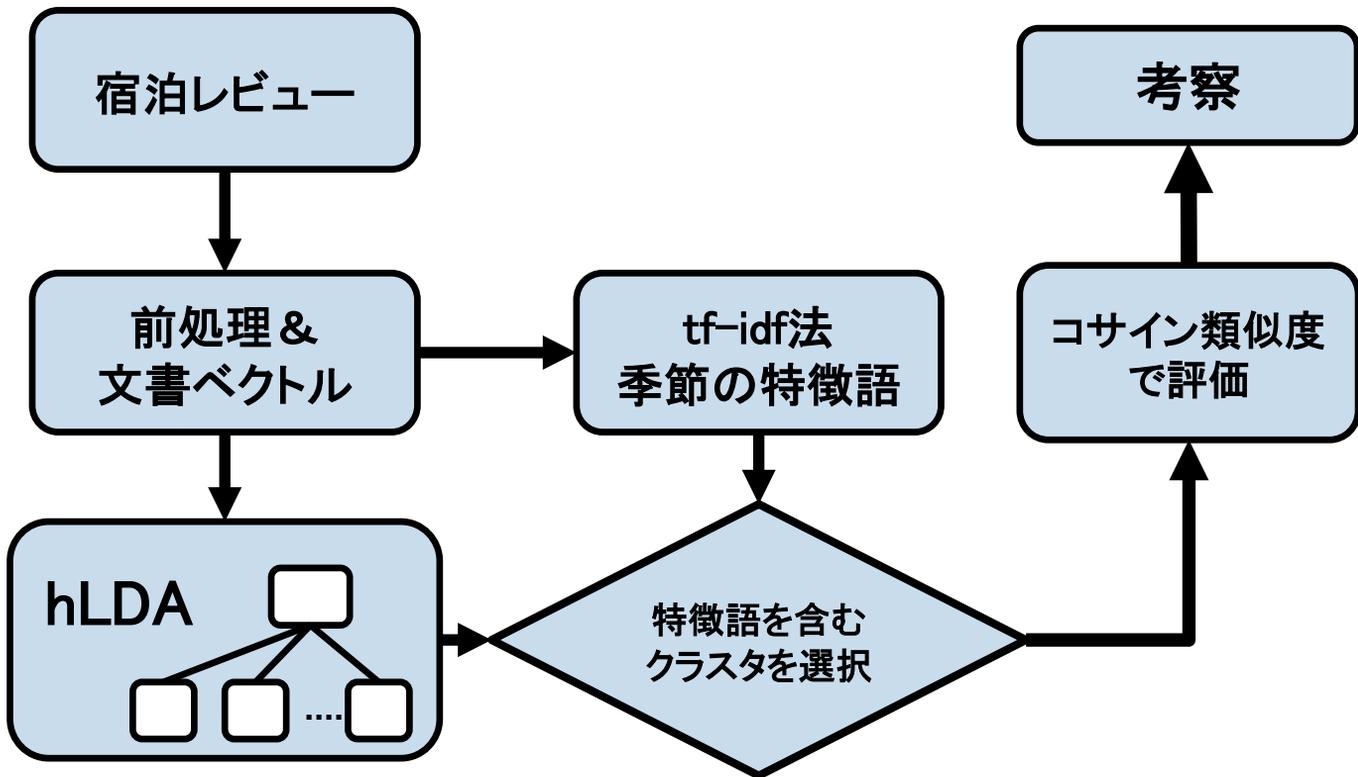
2. 分析手法・評価方法

3. 抽出結果

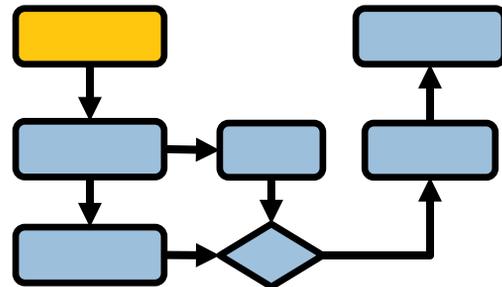
4. 考察

5. 結論・展望

2 分析フロー



2 入力データ



■ 楽天トラベルのレビューデータ

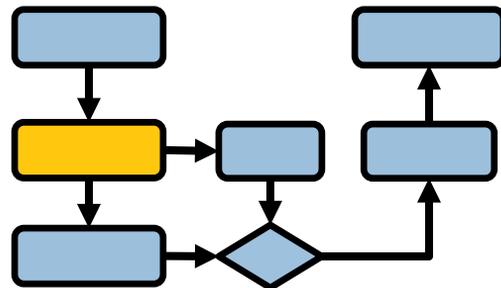
□ 1996年～2016年に収録

□ レビュー数上位10%から選出

◆ 粟津温泉 旅亭懐石 のとや (3,118件)

◆ シーサイトホテル舞子ヒラ神戸 (2,770件)

2 文書ベクトル化

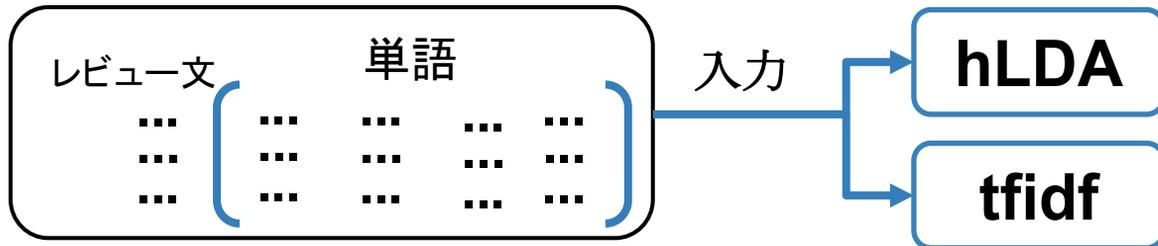


- レビューを月毎に分割

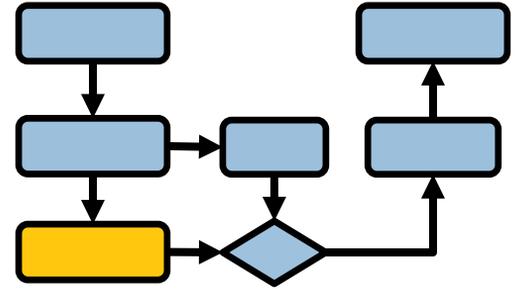
12ヶ月の文書集合

1月のレビュー

- 文書ベクトルに変換

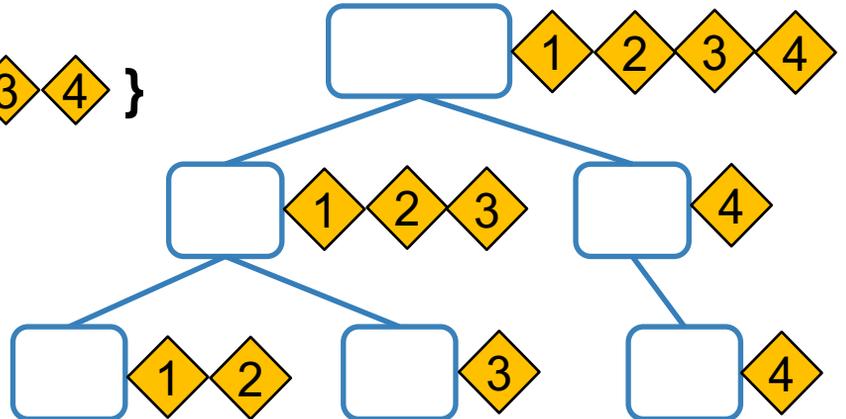


2 hLDA(1/2)

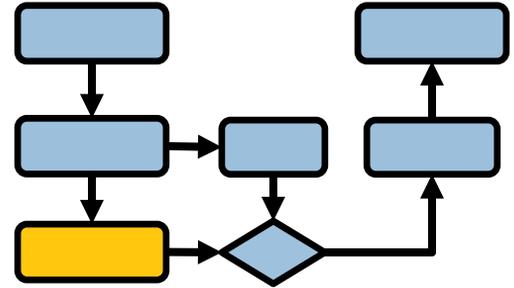


- hLDAによる階層型トピック分析
✓ トピックモデル(LDA)の拡張

Topics = {     }



2 hLDA (2/2)

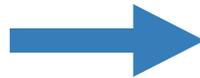


- hLDAで宿泊施設のカテゴリを階層構造で抽出できる
 - ✓ 本研究では名詞が分類される

仮説

宿泊レビュー

hLDA



食事・風呂

料理

風呂

蟹・うに

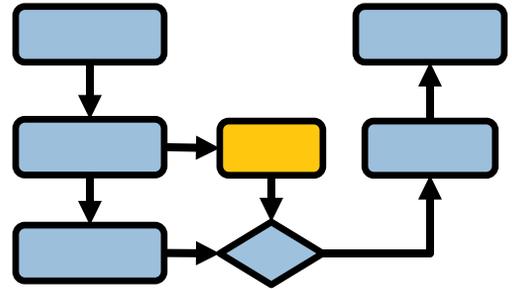
美味しい

綺麗

露天

2

tf-idf法



✓ tfidf値: 文書に含まれる単語の重要度

→つまり本研究では**季節(月)の特徴量**と言える

12個のベクトル



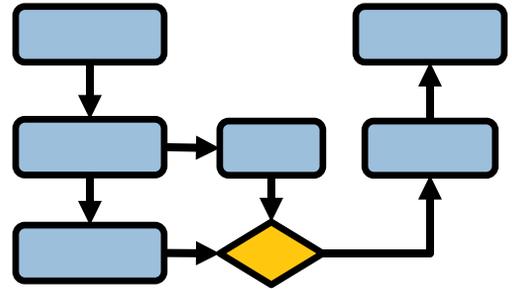
tf-idf法



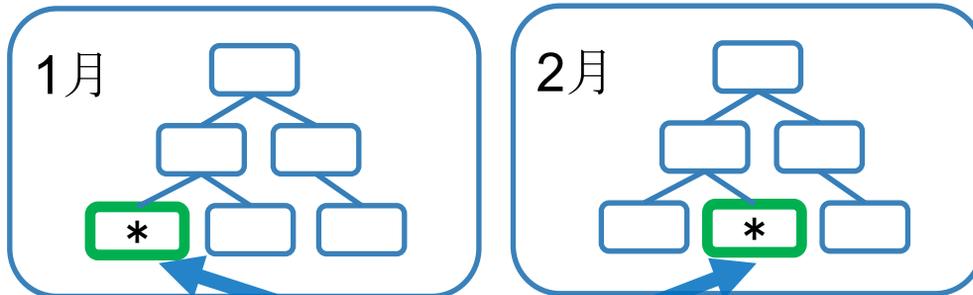
季節の特徴語

各月のtfidf値の大きい
上位5件の名詞集合

2 クラスタの選択

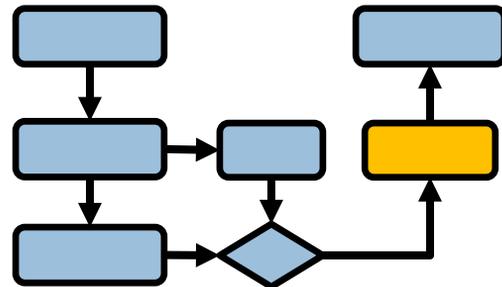


- 季節の特徴語を含むクラスタの抽出



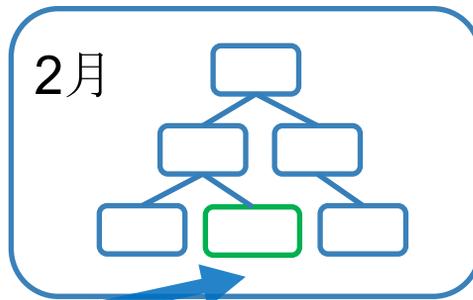
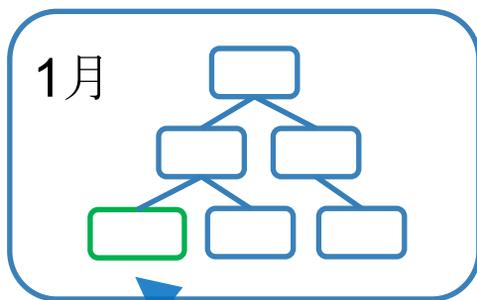
tf-idf法により抽出された
季節の特徴語 *

2 コサイン類似度(1/2)



■ コサイン類似度による評価

✓ コサイン類似度: ベクトル間の類似度



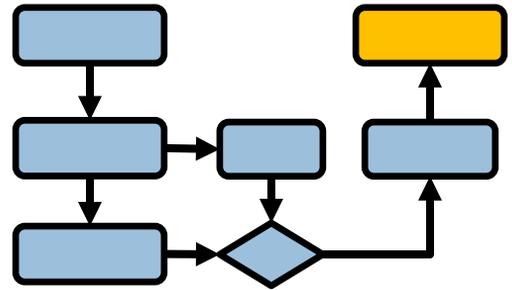
...

クラスタ間の類似度

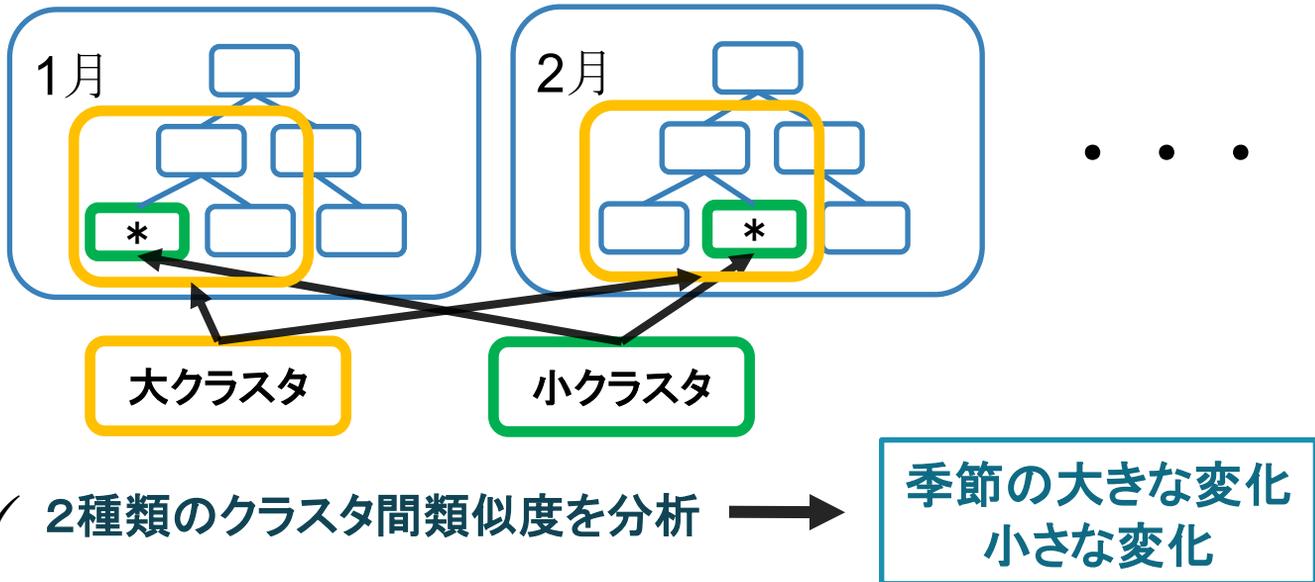


高い: 似た特徴を持つ
低い: 季節的な特徴が変化した

2 コサイン類似度(1/2)



■ 類似度の変化量を評価



1. 研究背景・目標

2. 分析手法・評価方法

3. 抽出結果

4. 考察

5. 結論・展望

3

hLDAの結果(1/2)

hLDA, 1月, 栗津温泉 旅亭懐石 のとや

料理, お部屋, 蟹, 満足, 部屋,
宿泊, 仲居, 風呂, 食事, 温泉

清潔, 割烹, きれい, 鍋, 程度, 旅行,
温泉, 季節, 得

雪道, 規制, 客室, 北陸道, 影響,
雪, 体調, 調節, 温度, きれいな

蟹, 食事,
合掌造り, 見学,
甘味, 見事,
ガニ, 立派,
越前, 建物

利用, 炬燵,
栗津温泉, 一泊,
都会, 喧騒

感心, 雪見,
従業員, 清潔感,
心配り, 最初,
以前, 自然

暖房,
醍醐味,
意見,
夏休み,
一泊

3 hLDAの結果(2/2)

■ hLDA, 8月, シーサイドホテル舞子ビラ神戸

部屋, ホテル, 利用, 宿泊, 満足,
景色, 大変, 朝食, 最高, 残念

近所, 評価, バイキング, 価格,
ロケーション, 舞子駅, チェックアウト

タオル, 表示,
雰囲気,
リゾート, 子供,
エアコン, 音,
家族連れ

コンビニ, 仕事,
部屋, 混雑,
チェックイン,
別館, 評価

受付, 眺め, 有料, 感動, 海水浴,
海水浴場, 船, メニュー

今回, 子達,
海, 目的,
海水浴場,
気持ち

限定, 全区,
料金, 客,
大浴場, 種類,
接客, 眺め

3

tfidf法の結果

- ・ 宿泊レビューから抽出したtfidf値が高かった上位5件の名詞

✓ 宿泊施設の繁忙期の月を表示

	粟津温泉 旅亭懐石 のとや	シーサイトホテル舞子ヒヅ神戸
1月	お正月, 年越し, 年始年末, 餅搗き, 最上級	お正月, お年玉, お菓子, 周年, スヘンシャル
5月	コールテンウィーク, 鯛, 春 の, 愛, 連休	室温, コールテンウィーク, 週末, クレート, エステ
8月	お盆, 合掌, 道場六三郎, 鮎, 納涼	プール, 夏休み, 海水浴場, 海水浴, バーベキュー
12月	クリスマス, めし, 突破, 名物, 最上級	ルミナリエ, クリスマス, 冬, おもてなし, スヘンシャル

1. 研究背景・目標

2. 分析手法・評価方法

3. 抽出結果

4. 考察

5. 結論・展望

4 考察(1/3)

- 「栗津温泉 旅亭懐石 のとや」
- “鯛”(5月)を含むクラスタの類似度

	1月～3月	4月	5月	6月	7月	8月	9月	10月	11月, 12月
大クラスタ	0.0	0.029	1.0	0.046	0.018	0.018	0.021	0.028	0.0
小クラスタ	0.0	0.099	1.0	0.199	0.099	0.099	0.099	0.105	0.0

✓ 4月, 6月, 10月の類似度が高い → 季節による変化が緩やか

✓ 小クラスタより大クラスタの方が, 変化の詳細がわかる

4月, 6月に共起した名詞
鯛, 鯛茶漬け, 香り, 懐石



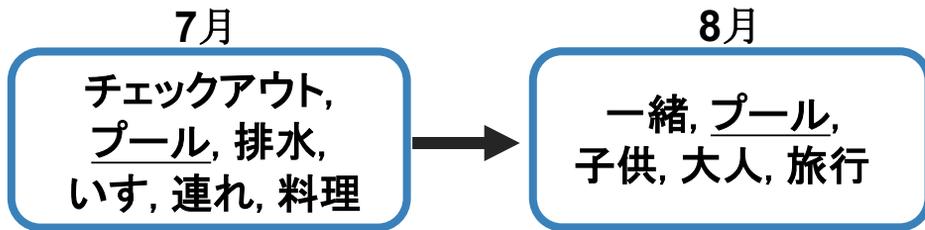
4月～6月にかけて
“鯛”が注目

4 考察(2/3)

- 「シーサイドホテル舞子ビラ神戸」
- “プール”(8月)を含むクラスタの類似度

	1月～5月	6月	7月	8月	9月	10月～12月
大クラスタ	0.0	0.110	0.042	1.0	0.018	0.0
小クラスタ	0.0	0.099	0.099	1.0	0.099	0.0

- ✓ 7月, 9月の類似度が低い → 季節の変化が激しい
- ✓ 7月, 8月小クラスタの比較 → 8月は家族旅行が人気



4 考察(3/3)

- ・「粟津温泉 旅亭懐石 のとや」
 - ✓ “鯛”が4月から6月にかけて人気
 - ✓ 食事に関する注目度が高い
- ・「シーサイドホテル舞子ビラ神戸」
 - ✓ プールは6月～9月に楽しめる
 - ✓ 8月は家族旅行が人気

1. 研究背景・目的

2. 分析手法・評価方法

3. 抽出結果

4. 考察

5. 結論・展望

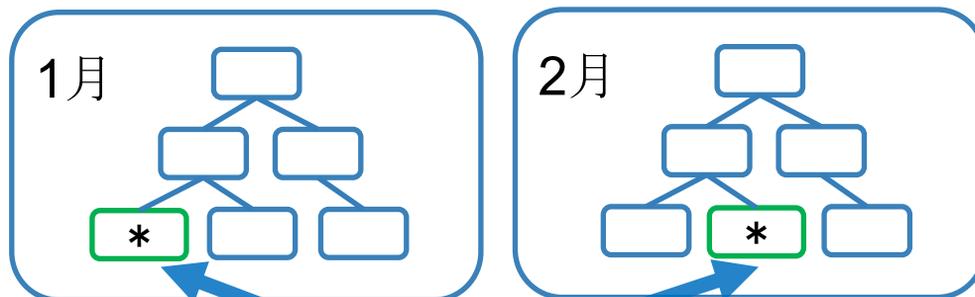
5 課題・展望

- hLDAの精度
 - ハイパーパラメータの調整
- 地域における特徴量の差
- 宿泊施設決定支援へのモデル化

5 まとめ

1. 宿泊レビューの季節の特徴を階層構造で抽出することを目的とした
2. hLDAとtfidf法の組み合わせによる分析
3. その結果を類似度により評価
4. それぞれのクラスタには季節差が存在
5. 今後は地域における特徴量の差を抽出

- 季節の特徴語を含むクラスタを抽出

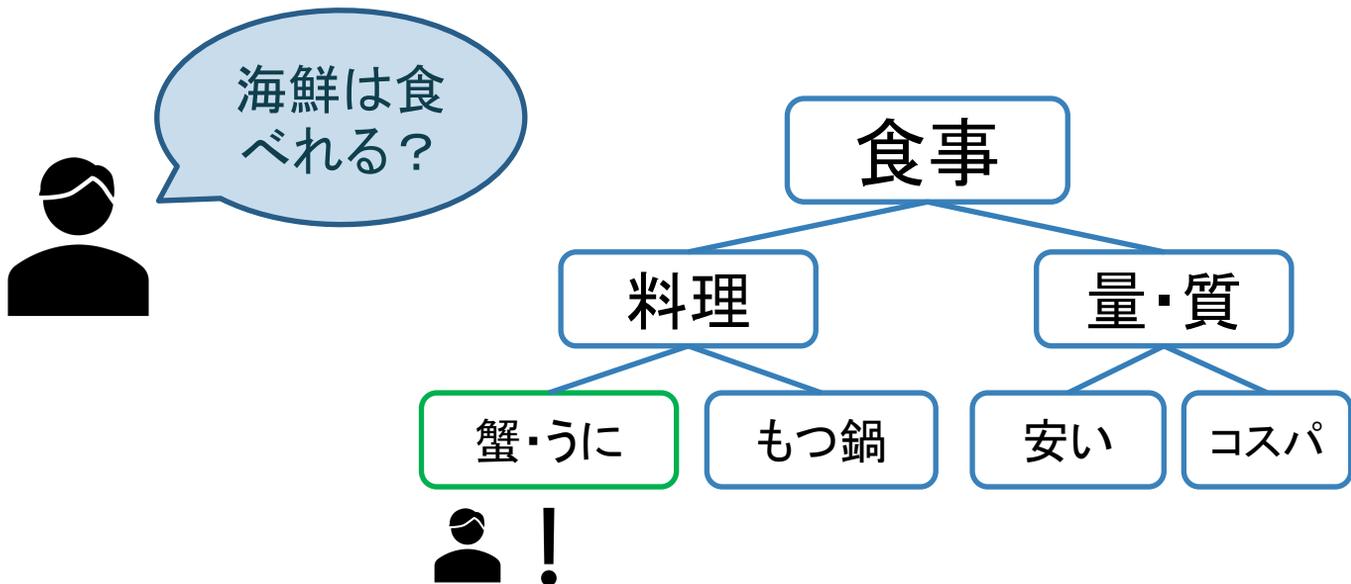


• • •

tfidf法により抽出された季節の特徴語*

1 研究目標

- カテゴリの細分化により，欲しい情報が手に入る
 - 宿泊施設決定要因になる



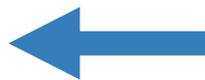
1 研究目標

1. 宿泊レビューにおけるカテゴリを階層構造で提示

2. 宿泊レビュー内の季節的な特徴を抽出



Aホテル
7月, 海



提示



- tfidf法により抽出された季節の特徴語*

ラスタを抽

...

5. 結論・展望



1. 研究背景・目的

2. 分析手法・評価方法

3. 抽出結果

4. 考察

5. 結論・展望