

テキストマイニングによる不満から読み取る鉄道に関する問題

ベイズ学習による投稿者の不明な属性の推定

王 帥・三道弘明 関西学院大学

背景・目的

- データマイニング技術がビジネスに活用されつつある
- ベイズ学習の適用可能性
- テキストマイニングの研究結果に基づいて社会に助言
- ベイズ学習による不明な情報を推測できるかどうかの検証

データの分類

トイレ トイレ + 汚い
女子トイレ + 混雑

施設 ゴミ箱 + 臭い
階段 + 長い

人・マナー 荷物 + 邪魔
席 + 譲らない

電車施設 エアコン + 寒い
コンセント + ない

ダイヤ 電車本数 + 少ない
終電 + 早い

料金 料金 + 高い
運賃 + 下げる

コーディングルール

- 共通する概念を一つのキーワードに集約
- 施設：駅構内のごみ箱が臭いやホームまでの階段が長いなど駅に対する不満
- ダイヤ：終電が早いや電車の本数が少ないなどダイヤに対する不満

*トイレ (トイレ | 手洗い | 便所 | おしっこ | 小便) & (汚い | 古い | ゴミ | 臭い | におい | 暗い | 増やす) | 洗面所

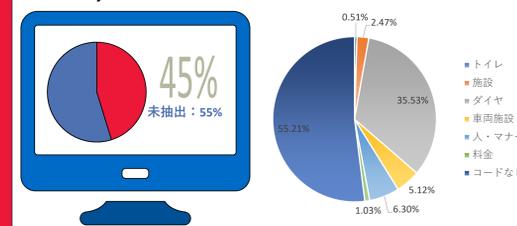
*施設 (階段 & (長い | 狭い | 混雑 | 多い | 急 | つらい | 辛い | 危ない)) | (ベビーカー | 車椅子 & 不便 | 専用) | (エレベーター & (狭い | 混雑 | 混む | 小さい | (広い & 欲しい | ほしい) | 臭い)) | (エスカレーター & (遅い | 狭い)) | (場所 & 座る & ない) | 困る & (駅 & 寒い | 暑い | 冷房 | 暖房) | 電波 & (駅 | ホーム))

*ダイヤ ((朝 | 帰宅 | 電車 | 通勤) & ラッシュ) | 乗り換え & 時間 | ((本数 | 電車 | 車両) & (少ない | 増やす | 欲しい | ほしい | 混む)) | ((特急 | 急行 | 快速 | 列車) & 増やす | 少ない) | (終電 & 早い | 遅い | ない) | 子連れ & 車両 | 電光 | 待ち時間 | 振替 | 輸送 | スピード & (遅い | 電車 | 上げる) | 職員

*車両施設 (電車 & 寒い | 暑い | 冷房 | 暖房 | 空調 | クーラー | 臭い | 匂い | におい | ニオイ | 空気) | クーラー & (効く | 寒い | 暑い) | ((電車 | 車内) & (狭い | 寒い | 暑い | 広い & (欲しい | ほしい))) | 手すり | 吊革 | 足元 & (暖房 | ヒーター | 熱い | 暖かい) | 隙間 & (ホーム | 広い | 危ない | 怖い) | 電波 & (駅 | ホーム) | 揺れ & (激しい | 電車 | 大きい | 酔う | 酷い) | 網棚 & (高い | 邪魔 | 狭い | 位置) | コンセント & (窓側 | 座席 | 席 | 新幹線 | 車両 | 窓際) | ガラス & (割れる | 曇る | 汚い | 汚れる)

集計結果

- 駅・電車サブカテゴリ 186,595レコードから
- 83,576レコードが抽出され
- 83,576レコードが6つに帰属



ベイズによる投稿者不明な属性の推定

不明データサンプル

G	H	I	J	K	L	M	N	O
product	text	crated_at	gender	age	birth	job	prefecture	arital_stat
トイレ	ない駅のトイレ	42291.69	女	不明	不明	学生	大阪府	不明
カビ	ないのか、しばらく体	42291.92	女	不明	不明	社員 (事務)	神奈川県	不明
己慮	に加え開く行為は公共	42292.14	不明	不明	不明	不明	不明	不明
微妙	な遅延困る乗換うまく	42292.07	不明	不明	不明	社員 (事務)	大阪府	不明
自動	改札機を持って左手	42292.03	不明	不明	不明	社員 (その他)	不明	不明
えどりの	感、どういこと	42292.59	男	不明	不明	社員 (技術)	兵庫県	不明
電車	とかつらすぎ	42292.04	不明	不明	不明	不明	不明	不明

データの整理

- 不明な属性 (外部変数) 性別・結婚歴・年収・年齢・職業・地域...etc..

- 不明な属性のない データ数

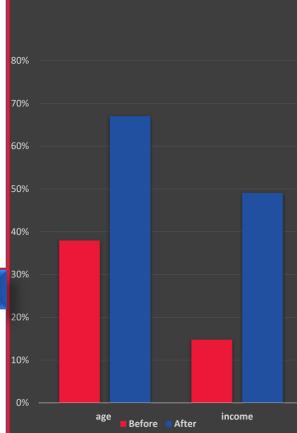
- 91399→76201
- 約2割弱が「欠損データ」



セグメントの見直し

年齢	年収
10	100
20	200
30	300
40	400
50	500
60	700
70	800
900	
1000	
1200	
1500	
2000	

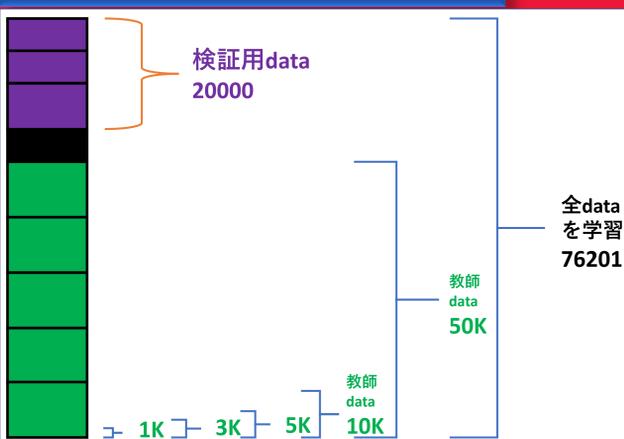
セグメントを見直した結果



全データの推測

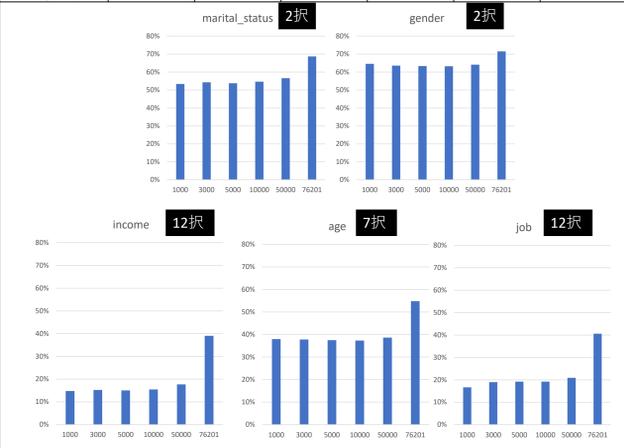
- 1,000件だけ学習
- 91,399の不明を含む data不明の部分推測
- クロス集計, 対応分析

教師データと検証データの用意

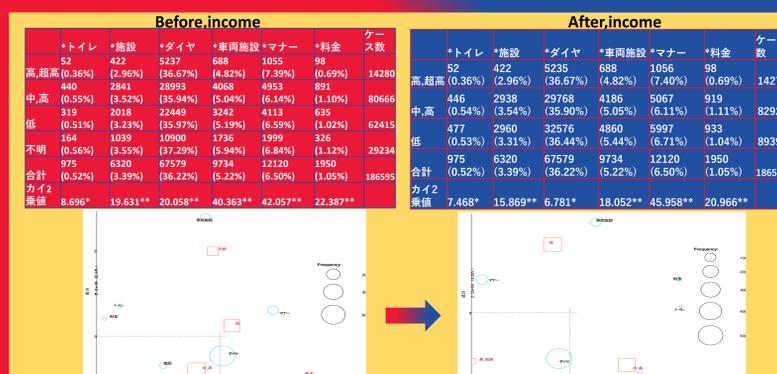
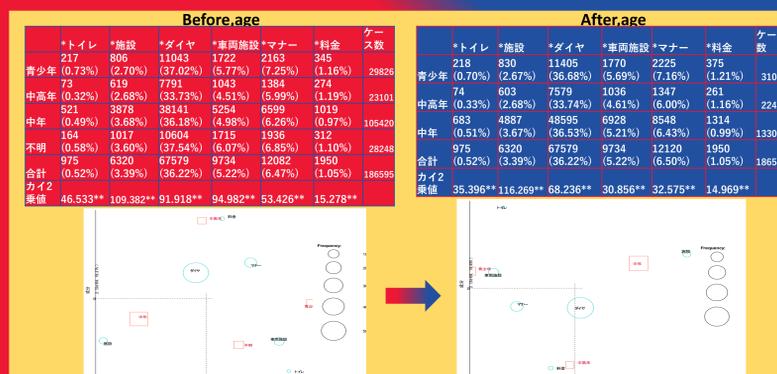
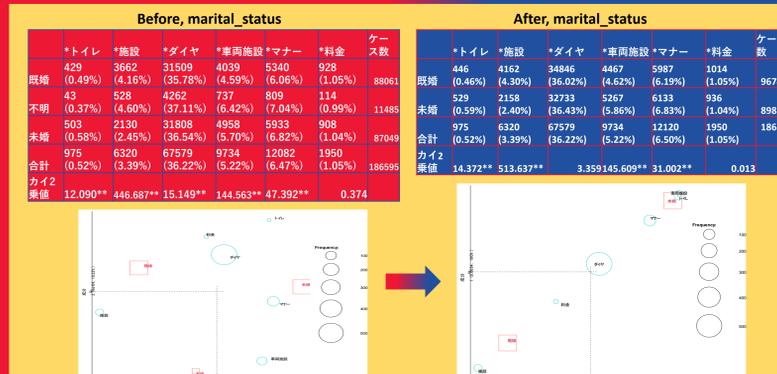
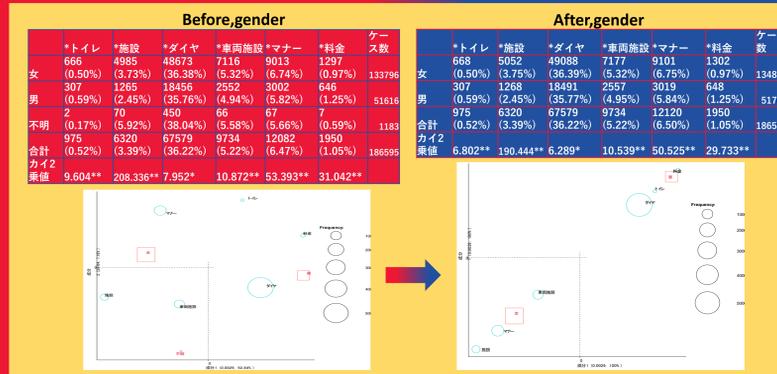


データの集計

回数×正解	1000	3000	5000	10000	50000	76201
男×男	454	924	1200	1551	1970	3160
女×女	12464	11794	11467	11091	10849	11149
女×男	6391	5921	5645	5294	4875	3685
男×女	691	1361	1688	2064	2306	2006
正解の確率	65%	64%	63%	63%	64%	72%
不正解の確率	35%	36%	37%	37%	36%	28%



外部変数による不明を推測する前後の変化



結果・考察

- ダイヤに対する不満が圧倒的に多い 女 > 男, 青少年 > 中年 > 中高年,
- 本研究においては、1000件のテキストデータと50000件のデータの予測結果は大した差がない 1000件のテキストデータで外部変数の傾向がわかる
- 適切なセグメンテーションにより学習効果が向上
- ベイズ学習により推測後、有意でなくなるデータがある 結婚歴×ダイヤ
- 欠損データを推定することによって、テキストマイニングの結果はより簡潔, 明確になった

参考文献

- 小木しのぶ, テキストマイニングの技術と動向, 計算機統計学, 28(1), 31-40, 2015
- 樋口 耕一, テキスト型データの計量的分析: 二つのアプローチの峻別と統合-理論と方法, 19(1), 101-115, 2004.
- 奥村 学, ソーシャルメディアを対象としたテキストマイニング, 電子情報通信学会 基礎・境界ソサイエティ, Fundamentals Review 6(4), 285-293, 2013.
- 打田 裕樹, Webユーザーレビューにおける評価情報の時系列変化の可視化, 知能と情報: 日本知能情報ファジィ学会誌: journal of Japan Society for Fuzzy Theory and Intelligent Informatics 22(3), 377-389, 2010-09-30.
- 高野 敦子, 因果関係に着目した口コミWebサイトからの評価表現抽出, 人工知能学会論文誌 24(3), 322-332, 2009.
- 樋口 耕一, 社会調査のための計量テキスト分析-内容分析の継承と発展を目指して, ナカニシヤ出版, 2014.