



HOKKAIDO UNIVERSITY

| | |
|------------------|---|
| Title | 楽天データの解析について |
| Author(s) | 白井, 康之; 小山, 聡; 櫻井, 祐子 他 |
| Description | ERATO湊離散構造処理系プロジェクト : 2010年度初冬のワークショップ (ERATO合宿) . 2010年11月29日 (月) ~12月1日 (水) . 札幌北広島クラッセホテル. |
| Relation | 2010年度科学技術振興機構ERATO湊離散構造処理系プロジェクト講究録. p.434-436. |
| Issue Date | 2011-06 |
| Doc URL | https://hdl.handle.net/2115/48344 |
| Type | conference presentation |
| File Information | 12.sakurai.pdf |



楽天データの解析について

白井康之⁽¹⁾, 小山聡⁽²⁾, 櫻井祐子⁽¹⁾, 川原純⁽¹⁾, 鶴間浩二⁽¹⁾

⁽¹⁾ JST ERATO 湊離散構造処理系プロジェクト
⁽²⁾ 北海道大学 大学院情報科学研究科

November 2010

1

楽天データ公開

【2010年 楽天データ公開】

- 第3回楽天研究開発シンポジウムの一環として、今年度より試行。
- 学術研究機関を対象にデータを公開
- 25件程度の申し込み (うち、実際に分析を行っているのは10件程度か?)

【公開データ】

- 楽天市場の全商品データ (5000万商品)
- 楽天トラベルの施設データ (11,468施設)、レビューデータ (35万レビュー)
- 楽天GORA (ゴルフ) の施設データ (1,669施設)、レビューデータ (32万レビュー)



<http://rit.rakuten.co.jp/dr/index.html>

November 2010

2

楽天データ公開 (商品データ)

| Column | Sample |
|--------|---|
| 店舗コード | rakutenstore |
| 商品ID | 12345678 |
| 商品名 | 【Rakuten T-Shirt】楽天ロゴ入り 着心地満点のTシャツ 体を締め付けない伸びる生地でデザイン |
| 商品説明文 | 上質の素材を使用し、シルエットにも末を使ったデザインです。ストリートだけでなく、有名百貨店でも取り扱われるようになりました。人気急上昇のTシャツです！ベーシックなシルエットでありながら、年齢、性別を問わず楽しんでいただけるデザインになります。オンにもオフにも、デノのみサイズにもコーディネートに最適です。動きやすいだけでなく、着心地にもこだわりました。 |
| 商品URL | http://www.rakuten.co.jp/rakutenshop/12345/98765/ |
| 商品画像 | http://thumbnail.image.rakuten.co.jp/9_0mail/rakutenshop/cabinet/tshirt_1234_56_78/img1234567890.jpg?_src=64x64 |
| 商品価格 | 4,500 |
| ジャンルID | 403882 |

<http://rit.rakuten.co.jp/dr/index.html>

商品数 5000万, 店舗数 7万

November 2010

3

楽天データ公開 (トラベルデータ)

| Column | Sample | Column | Sample |
|----------|---|------------|----------|
| 施設ID | 121212 | 施設番号 | 98887766 |
| ユーザー投稿本文 | 値段のわりには、きれいなホテルでした。以前泊まったときより、きれいに改装されたようで、入ったときほびっくりしました。しかし、以前と同じフロントでは丁寧な対応もして頂き、気分よく泊まることができました。接客には十分満足しましたが、駐車場が少し混みかたり、コンビニが遠くになかったりしたところ少し残念です。朝食はバイキングでしたが、品数も多く、おいしかったです。レストランスタッフの方をみんなでもなかなか出てこられなかったのは残念でした。 | ニックネームID | 12345 |
| 施設番号 | 98887766 | 目的 | レジャー |
| 分類 | 宿泊・情報 | 同業者 | 家族 |
| プランID | 242424 | 評価1 (立場) | 4 |
| プランタイトル | 新館・禁煙・朝食バイキング！平日お得プラン！ | 評価2 (部屋) | 4 |
| 部屋種類 | 2 | 評価3 (食事) | 4 |
| 部屋名前 | 新館・禁煙・セミダブルルーム | 評価4 (風呂) | 0 |
| 施設投稿本文 | この度は、ご利用いただき誠にありがとうございます。また、貴重なご意見をお寄せ頂きありがとうございます。今年初めに改装しましたが、お蔭の消さまてありがとうございます。レストランスタッフの対応に満足しては、大変申し訳ございません。従業員の方を徹底して教育します。またのご利用をお待ちしております。 | 評価5 (サービス) | 3 |
| | | 評価6 (設備) | 3 |
| | | 評価7 (総合) | 4 |

<http://rit.rakuten.co.jp/dr/index.html>

施設数 1万,
レビュー数 35万

November 2010

4

楽天データ公開 (ゴルフ場データ)

| Column | Sample | Column | Sample |
|-------------|--|------------|---|
| コースID | 9876 | ゴルフ場ID | 1234 |
| コースナンバー | 1 | ゴルフ場名 | サンパシフィックコース |
| ゴルフ場ID | 1234 | 施設番号1 | 200 |
| コース名 | 5477 | 施設番号2 | 30005 |
| ホール | 4 | 位置 | 埼玉県東武東上線 |
| ホール名 | 3 | 施設年開業 | 200904 |
| ホール名のパー数の合計 | 25 | 施設コメント | 穴場100のフラットなコースで、インターバルも短く、練習しやすいコースになっています。100mの最大な長さのなかで、プレーを楽しめます。初心者にもオススメのコースです。朝一混雑の日の受付も行ってきます。 |
| ハンディキャップ1 | 7 | ゴルフ場コメント | 毎月1日は女性限定のレディースデーを開催しております。詳しくはHPをご覧ください。毎週月曜1日曜日は休館です。 |
| ハンディキャップ2 | 15 | ゴルフ場ID | 200779 |
| ホール説明1 | 左にゆるやかにドッグレッグ、右に急なアップダウン。右に急なアップダウン。右に急なアップダウン。右に急なアップダウン。 | ゴルフ場ID | 1212 |
| ホール説明2 | フェアウェイが狭いアップダウン。右に急なアップダウン。右に急なアップダウン。右に急なアップダウン。 | 施設番号ニックネーム | user001 |
| | | 総合評価 | 3 |
| | | コストパフォーマンス | 3 |
| | | コストパフォーマンス | 4 |
| | | コース/環境性 | 4 |
| | | 食事/設備性 | 5 |
| | | 設備が充実 | 2 |
| | | タイトル | 美しくプレーできるコース |
| | | コメント | きれいに手入れされたグリーンで、気持ちよくプレーできます。初心者の私には少し難しかったですが、練習してまたチャレンジしたいと思います。尚、朝食、朝食も多く |
| | | プレー日 | 2010/01 |

楽天GORA: クチコミデータ

November 2010

5

楽天市場の現状のサービス

- 購入履歴の管理
- 閲覧履歴の管理
- 上記2つをもとにした推薦商品の提示
- 全文検索機能 (and/or/not検索)
- ジャンル・価格絞り込み機能



November 2010

6

2. プラントタイトル作成支援システムの構築 分析の対象



November 2010

13

2. プラントタイトル作成支援システムの構築 (1/2)

- ・ **プラントタイトル作成支援の必要性**
 - ホテルは50文字以内でプラントタイトルを作成
例) [東京タワーの夜景に乾杯~]~タワービュークリスマス~
IN17:00/OUT10:00(朝食付)
 - アピールしたい単語を入れることが一般的
 - ユーザはプラントタイトルからホテルを選定する場合も多い
- ・ **楽天公開データ(分析対象)**
 - プラントタイトル, ホテル名, ユーザ評価(7項目:総合評価, サービス, 立地, 部屋, 設備・アメニティ, 風呂, 食事)など
 - タイトル総数:84,000 件余
- ・ **プラントタイトル作成支援システム**
 - プラントタイトルに含まれる単語とユーザ評価(項目別)の相関関係を考慮し, ホテルが重視する評価項目に応じたプラントタイトルの作成支援を行う。

November 2010

14

2. プラントタイトル作成支援システムの構築 (2/2)

| 単語 | プラン数 | 投稿数 | 施設数 | 出現回 |
|------|-------|--------|------|-------|
| プラン | 56299 | 187853 | 8913 | 58906 |
| の | 47729 | 177338 | 6502 | 52782 |
| に | 47729 | 177513 | 6504 | 52788 |
| は | 31828 | 155892 | 6360 | 45473 |
| と | 22759 | 143848 | 5350 | 44013 |
| を | 26917 | 108117 | 6141 | 34007 |
| で | 23782 | 84588 | 6164 | 29025 |
| が | 21763 | 83908 | 4780 | 26884 |
| を | 15640 | 68825 | 3876 | 25500 |
| は | 15689 | 76867 | 4410 | 23633 |
| ★ | 12671 | 75336 | 3641 | 22941 |
| 付 | 20596 | 65085 | 5350 | 22087 |
| 朝食 | 21106 | 82400 | 4900 | 21782 |
| 指定 | 18455 | 72580 | 5332 | 20317 |
| と | 15095 | 55804 | 4098 | 18898 |
| は | 16412 | 56332 | 5282 | 17907 |
| と | 11812 | 51516 | 4208 | 16299 |
| で | 13894 | 43287 | 4663 | 15114 |
| ポイント | 9928 | 40833 | 3459 | 12257 |
| と | 8208 | 38418 | 2700 | 11985 |
| は | 10150 | 33870 | 4083 | 11148 |
| ◆ | 5668 | 37804 | 1463 | 10391 |
| 付 | 9776 | 36640 | 3552 | 10157 |
| 付 | 8403 | 28769 | 3678 | 9731 |
| 付 | 9037 | 34111 | 3343 | 9720 |

- ・ **現状**
 - 形態素解析を行い, 単語毎に出現回数や項目別評価点(平均, 分散)を分析
 - 単語数:17,000 語余
- ・ **今後の方針**
 - 得点に影響する単語の組合せ(positive/negative synergy)の抽出手法の検討
 - ・ とりあえず二単語間
 - 50文字制限の下で, 大きい値以上の評点を持つ単語の組合せの検出方法の検討

November 2010

15

3. 多面的数値・テキスト評価データに基づく推薦システム 分析の対象



November 2010

3. 多面的数値・テキスト評価データに基づく推薦システム (1/2)

■ 従来の協調フィルタリングによる推薦システム

ユーザのアイテムに対する評価値

Item 1 Item 2 Item 3 Item 4

| | | | | |
|--------|---|---|---|---|
| User 1 | 5 | 1 | 3 | 1 |
| User 2 | ? | 3 | 4 | 3 |
| User 3 | 4 | 2 | 4 | ? |
| User 4 | 3 | 4 | ? | 3 |

- ◆ 既知の評価値を基に未知の評価値を予測
- ◆ 協調フィルタリング:「似ているユーザの評価値は似ている」と仮定して未知の評価値を予測

■ 従来の方法の問題点

- ◆ 既知の評価値の相関をユーザ間の類似度として用いる
- ◆ 評価値が似ているからといって, ユーザの嗜好が似ているとは限らない
- ◆ 同じホテルに付けた評価値「5」は, あるユーザは「立地」を, 別のユーザは「設備」を重視して評価したかもしれない

November 2010

17

3. 多面的数値・テキスト評価データに基づく推薦システム (2/2)

■ 楽天トラベルデータ(ホテル評価)の特徴

- ◆ 総合評価だけでなく「立地」「設備」「食事」等の多面的な評価値を付与
- ◆ 数値評価に加えて感想や苦情がレビューテキストとして与えられている

■ 目標:複数側面での(多面的な)評価値とテキストの内容を考慮して推薦精度を向上

■ 方針:

- ◆ 各側面での評価値に加え, レビューでの語の使用回数も「側面」として行列に保持
- ◆ これらをまとめた「テンソル補完」の問題として定式化
- ◆ 語の組合せやフレーズを考慮する場合, ZDD等の技術が必要となる

| | Item 1 | Item 2 | Item 3 | Item 4 |
|--------|-------------|-------------|--------|--------|
| User 1 | 「立地」に関する評価値 | | | |
| User 2 | 「食事」に関する評価値 | | | |
| User 3 | ? | 「設備」に関する評価値 | | |
| User 4 | 4 | ? | 3 | 4 |
| | 3 | 2 | 1 | 5 |
| | 1 | 3 | ? | 3 |

| | Item 1 | Item 2 | Item 3 | Item 4 |
|--------|-------------|----------------|--------|--------|
| User 1 | 「地下鉄」の使用回数 | | | |
| User 2 | 「おいしい」の使用回数 | | | |
| User 3 | 0 | 「ブロードバンド」の使用回数 | | |
| User 4 | 2 | 0 | 1 | 0 |
| | 1 | 0 | 2 | 1 |
| | 1 | 0 | 0 | 2 |

November 2010

18