

## 第9分科会

# 授業評価アンケートの 自由記述の自動分類とその応用

### 報告者

松河 秀哉 氏      東北大学 高度教養教育・学生支援機構 講師

### コーディネーター

根岸 千悠 氏      京都外国語大学 共通教育機構 講師



## 授業評価アンケートの自由記述の自動分類とその応用

コーディネーター

京都外国語大学 共通教育機構 講師 根岸 千悠

---

---

### ○本分科会のねらい

文部科学省の調査によると、授業評価アンケートはほとんどの大学で実施されている。しかし、自由記述のデータはどのように分析すれば良いのだろうか。本分科会では、「授業評価アンケートの自由記述を自動分類するためのWebシステム」を開発された先生からご講演をいただいた後、参加者に実際にシステムでの分析を体験していただき、授業評価アンケートの自由記述の活用方法について模索することを目的とした。

### ○報告の概要

前半は、松河秀哉先生（東北大学）から「授業評価アンケートの自由記述を自動分類するためのWebシステム」の開発経緯と特性に関して講演いただいた。授業評価の自由記述のようなテキストデータは、手動やテキストマイニングツールなどによる分析が行われているが、それぞれにメリットや限界がある。そこで松河先生は、それぞれの限界を踏まえて、トピックモデルを用いて授業評価の自由記述を分析するWebシステムを開発している。

後半は、本Webシステムの具体的な使用方法を説明していただいた。コーディネーターの根岸が簡単にデモンストレーションをしながら、参加者には事前に準備いただいたデータ、もしくは、体験用の架空のデータを使ってシステムを体験していただいた。その後、分析結果をもとに、本システムのさらなる活用方法などについて2～3名程度でディスカッションをしていただき、全体で共有した。

講演に加えて、実践的な体験や参加者同士のディスカッション、質疑応答が行われ、活発な会となった。

### ○報告に対する質疑ならびに全体討議の内容

質疑は、システムの体験中に個別で行われたほか、全体でも実施された。システムの使い方に関する具体的な質問や、授業ではなくワークショップの自由記述でも活用できるか、といった質問があった。松河先生からは、本システムは、授業評価アンケートの自由記述のデータをもとに開発されているため、ワークショップの自由記述では適当なカテゴリが付与されない場合があると考えられるが、授業に近いワークショップであれば参考になる結果が出る可能性はあるとの回答があった。

また参加者同士のディスカッション後の全体共有では、今後本システムを使用して、年度間での比較をしたいという声や、他の教育データとの組み合わせによる分析を行いたいという声が寄せられた。

## スライド1

授業評価アンケートの  
自由記述の自動分類と  
その応用

東北大学 高度教養教育・学生支援機構  
松河秀哉

## スライド2

## 自己紹介

- 名前：松河 秀哉(まつかわ ひでや)
- 経歴：大阪大学大学院人間科学研究科修了。大阪大学全学推進機構助教をへて現職。仙台在住8年目。
- 所属：東北大学 高度教養教育学生支援機構・講師
- 専門：教育工学

## スライド3

## 自己紹介

- これまでの研究
  - 電子掲示板の教育利用(園と家庭の連携システム、初期のLMS開発)
  - 大規模教育データの分析(教育データのデータマイニング、授業評価アンケートの自由記述の分析)
  - 遠隔教育(コロナ時にはオンライン授業支援も行いました)など
- 電子掲示板に蓄積されるやりとりのデータを何とかしたいというのが、テキストデータの分析に手をだした動機。そのときに使った形態素解析と、その後関わったデータマイニングが組み合わさって、テキストマイニング的な研究が増え、授業評価アンケートの自由記述の自動分類を夢見ながら現在に至ります。

## スライド4

## 授業評価アンケートの自由記述の分類

- 自由記述を分類・分析しろといわれたらどうしますか？
- 手動？
- テキストマイニングツール？

## スライド5

## 手動でのテキストデータの分類

- 人間が文章をひたすら読み、似た内容の文章をグループ化する、もっとプリミティブであるが、柔軟な方法。
- メリット
  - 分類するグループを非常に柔軟に設定・作成できる
  - 個々の文章がどのグループに属するかについても(分類者のスキル・熟練度にもよるが)文脈を読みながら非常に柔軟かつ高精度に分類できる
- デメリット
  - 分類から分類者の知識・スキル・信念の影響を排除することは不可能なので、妥当性・信頼性の保証が難しい
  - 分類対象の文章の件数が増えてくると、負荷が急激に増大  
数百件が限界？

## スライド6

## 手動でのテキスト分類

- 信頼性・妥当性を高める方向はいくつか考えられる
    - 複数人で分類する
    - 分類用カテゴリの生成にも一定の手順を用いる
      - 例えばk法とか
        - 大変さは増える一方ですね。
  - 件数増加に伴う負荷の増大への対応
    - ランダムサンプリング
      - 大体の傾向は変わらないけれど・・・
      - サンプリングでも大変な割に信頼性についてとやかく言われがち
- はっきりいって最終手段。できればやりたくない。

## スライド7

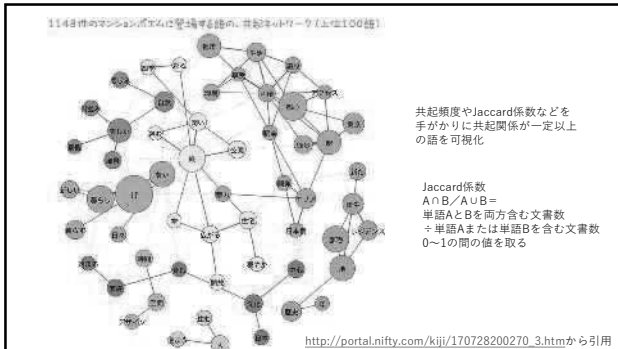
### テキストマイニングツールを使った分析

- KHCoderは有料化されてしまいましたが、単語と単語の共起関係を考慮に入れた分析を見たことがある人は多いでしょう
  - 共起ネットワーク分析
  - クラスタ分析
  - コレスポネンス分析 など
- メリット
  - 全体としてどんな内容が書かれているかなんとわかる
- デメリット
  - 単語数が増えてくると分析が困難に
  - 個々の文章にどんな内容が書かれているかは、直接は分からない
  - 本当にその内容が書かれていることを納得させるのは結構面倒くさい

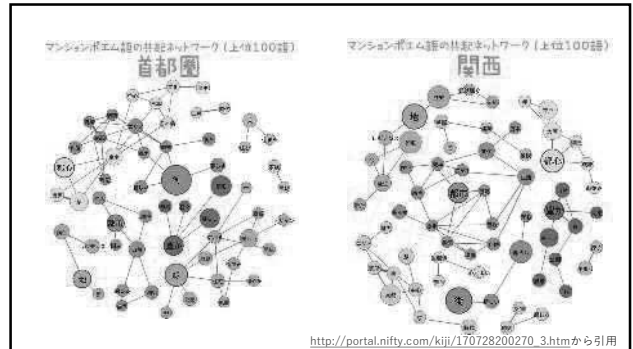
## スライド8



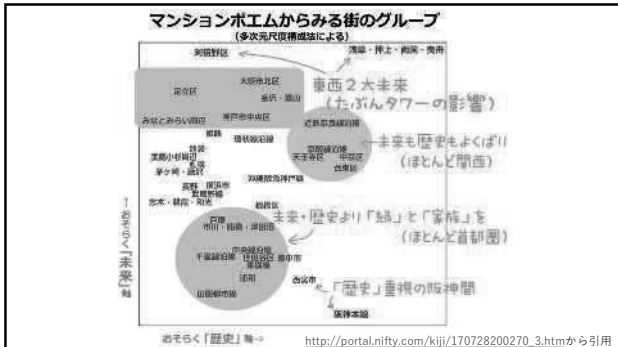
## スライド9



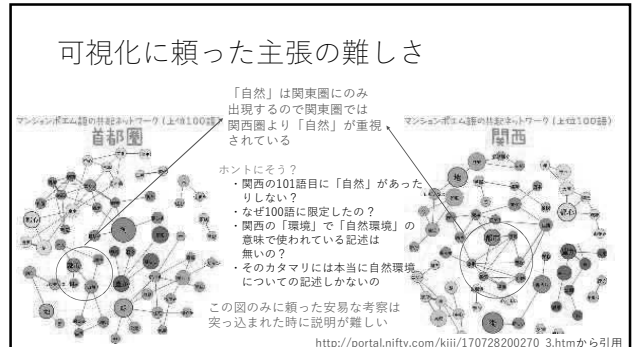
## スライド10



## スライド11



## スライド12



## スライド13

### 可視化に頼った主張の難しさ

- どうしてその100語を選んだの？
  - これについては、頻度順にとか、tf-idfの値が大きい順に…と言えなくはない。
- 101語目にキーとなる語がないとなぜ言えるの？
  - これはなかなか反論しづらい。
  - 究極的には全単語を使えば回避出来るけど、ネットワークが密すぎて解釈困難になる。

密な部分は同じ単語がいろんな単語と繋がっている=同じ語が違う文脈で使われている可能性も高く、話題の全体像を正しく把握するには重要なのに有効活用できない

## スライド14

### 可視化に頼った主張の難しさ

- 関西の「環境」で「自然環境」の意味で使われている記述は無いの？
    - こうした指摘には、関西の中で「環境」を含む記述を全部抽出して何件あるか確認し、さらに内容を目視して「自然環境」について書かれた記述が何件あるか調べて、十分低い割合であることを示したり、さらに、関東のなかで、「環境」を含む文を全部抽出して、同様のことを行い、「自然環境」に対する記述の割合が十分に高いことを示す必要がある。
  - そのカタマリには本当に自然環境についての記述しかないの？
    - こちらもカタマリ内の複数の単語を全部含む文を抽出して、同様な確認が必要。
- その話題が「存在する」ことを納得してもらうには、結局結構な手間がかかる

## スライド15

### テキストマイニングツールを使った分析の限界

- 分析対象とする単語数をあまり増やせない
- 全体に存在する話題はある程度把握できても、一件ごとの記述がどの話題に対応するかは直接は分からない
- ある話題に対応した記述が「どの程度存在するか」を確かめるには結局人手に頼ることになり、かなり手間がかかる

丁寧にやれば問題無いが、安易に使うと後が大変  
 全体の話題と、各記述がどの話題かが同時に分かる手法はないものか  
 →トピックモデル

## スライド16

### トピックモデルとは

- 文章に含まれる話題推定のモデル
  - 情報検索技術のLatent Semantic Analysis(潜在意味解析)が源流←行列の特異値分解
  - それを確率化したpLSI(Probabilistic Latent Semantic Indexing)をへて
  - さらにそれをベイズ化したLDA(Latent Dirichlet Allocation)潜在ディレクレ分配モデルへと進化してきた。

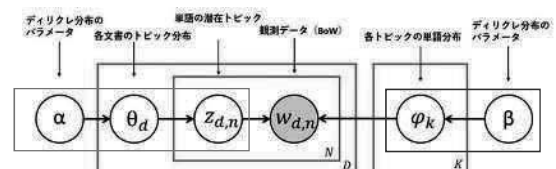
## スライド17

### トピックモデル

- LDA(潜在ディレクレ分配モデル)
  - 話題=トピックは複数の単語の分布として規定される
  - 1つの文章が複数の話題を持つと仮定する。
    - 文章全体の中にどのようなトピックがあるか
    - 各文章にどのようなトピックが含まれるかの両方が一度にわかる
  - 複数のアルゴリズムがあるが、MCMC(マルコフ連鎖モンテカルロ)法の一つであるギブスサンプリングでベイズ推定を行う手法が主流
  - 正解データが不要な教師なし学習の一種
  - 分類精度も十分にある(個人的感想)

## スライド18

### LDAのグラフィカルモデル



K個の各トピックは一定の割合で各単語を含んでいる  
 各文章は一定の割合で各トピックを含んでおり、その割合に応じた数の各トピックの単語から構成されている。  
 観測データに合致するようパラメータ推定すれば、それらの一定の割合がわかり、各トピックや各文章がどのような性質が判明する。

[https://qiita.com/K\\_Noguchi/items/2f0579ea51f5329a4008](https://qiita.com/K_Noguchi/items/2f0579ea51f5329a4008) から引用

## スライド19

### トピックモデル(LDA)の分析手順

- 分析対象となるテキストデータを整える
- Rのなど分析に対応したソフトで、整えたデータ読み込む
- トピック数等のパラメータを設定し与えて分析を実行する
- 分析結果を解釈し、トピックに名前を付けたり、トピック名の妥当性の検証などを行う

→滅茶苦茶面倒くさい

## スライド20

### トピックモデルの分析用データの準備

面倒くさい

↓ 表現が同じで品詞が違う言葉を区別するために、形態素と品詞をつなげている  
※語根、スペース区切りの分かち書きだけでもOK

## スライド21

### Rでごちゃごちゃ

Rの環境で以下のコマンドを打つ

```
>install.packages("lda") # ldaパッケージをインストール。初回のみ必要
>install.packages("topicmodels") # topicmodelsパッケージをインストール
>library(lda)
>doc <- scan("documents.txt")
>result <- lexicalize(doc) # 分析用の Document-Term Matrix を作成
>library(topicmodels) # topicmodels パッケージを読み出し
>dtm <- ldaformat2dtm(result$documents, result$vocab) # Document-Term Matrix を topicmodels パッケージ用に交換
```

面倒くさい

## スライド22

### トピックに名前を付ける

- 得られた情報を整理して、トピックごとにこのような表を作る

トピック	各トピックに属する確率が高い自由記述の例(単語)	各トピックに属する確率が高い自由記述の例(ランダム抽出)	件数
1	学習すべき文法について、文法事項中心のドイツ語知識より先行	した。	
2	ドイツ文法 中級者 文化の紹介 英語知識 外国人労働者 二語的	う一回者にとっては単語も文法もわかっていなかった。難易度について気づくことができました。	1059
3	授業内容 日本語 単語 単語 スペル 練習 単語 単語 単語	ドイツ語の文法だけでなくドイツの文化などについてもできておきました。	
4	授業内容 日本語 単語 単語 スペル 練習 単語 単語 単語	授業内容は興味深いものであったが、ドイツ語知識は本文法に理解のよりである。あまり文法事項を習得せず、また、もう一方のドイツ語の単語も覚えきれず、授業で習得している単語で授業が行われるので、勉強が苦手になってきた。単語の覚えが早いです。日本語とドイツ語の単語の覚えが早いです。日本語とドイツ語の単語の覚えが早いです。	
5	授業内容 日本語 単語 単語 スペル 練習 単語 単語 単語	小テストなど定期的にやってくれるので復習しやすいと感じます。ドイツ語で先生が話すのを聞いて理解できないので授業後のコミュニケーションを大切にしたいです。	
6	授業内容 日本語 単語 単語 スペル 練習 単語 単語 単語	授業内容が面白く、先生が丁寧で分かりやすいです。授業内容が面白く、先生が丁寧で分かりやすいです。	

## スライド23

### トピックモデルの利点

- 一度ちゃんと分類できるモデル作っておけば、そのモデルで新しいデータを分類できる
  - ただし、例えばコロナ前のデータで作ったモデルだと、オンライン授業に対する感想が皆無だったので、コロナ後の感想には対応できないという問題もある。
  - 授業の性質が変わる事件はもう、そんなに起こらないと思うので、最近のデータもいれてモデルを作ってしまうと、使い回せる
- ということで、面倒くさい部分は全部私の方でやりました。

## スライド24

### 授業評価アンケートの自由記述用モデル

- 2022年から2020年までの、4大学38万件の授業評価アンケートのデータを収集
- トピックモデル(LDA)で分析。160トピックを抽出・命名
- オンライン授業の分析も対応

トピック	トピック名	トピックを構成する単語
1	授業内容	授業内容、授業内容、授業内容
2	授業内容	授業内容、授業内容、授業内容
3	授業内容	授業内容、授業内容、授業内容
4	授業内容	授業内容、授業内容、授業内容
5	授業内容	授業内容、授業内容、授業内容
6	授業内容	授業内容、授業内容、授業内容

スライド25

モデルを組み込んだwebシステムの開発



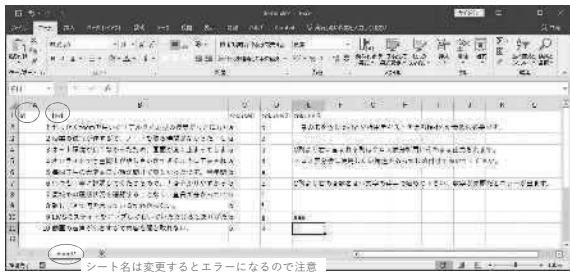
スライド26

Webシステムのその他の機能

- ネガポジ分析
  - 入力したテキストが、ネガティブな記述なのか、ポジティブな記述なのかを分析。
  - トピックモデルの分析単体では分からない部分に分かる。
- fasttextを用いた高速分析
  - トピックモデルによる分類結果を教師データとして学習したモデル
  - 分析がすぐ終わる。
  - 160種類のトピックは多すぎて煩雑なので、このモデルでは、似たトピックをまとめて67種類に絞っている
  - 本日はこれを使ってください。

スライド27

Webシステムで使うデータの準備



スライド28

Webシステムの使い方

IDとパスワードを使ってログインします



スライド29

Webシステムの使い方

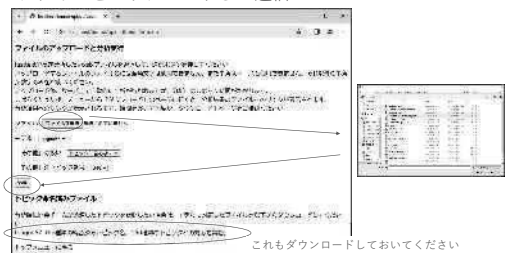
ファイルをアップロードしてfasttextで高速分析実行をクリック



スライド30

Webシステムの使い方

ファイルをアップロードして送信





スライド31

### Webシステムの使い方

結果をダウンロードをクリック

このスクリーンショットは、ウェブブラウザで表示されたメニュー画面です。メニュー項目には「プロフィールをアップロードしてトピックを再分析する」、「プロフィールをアップロードして再分析を行う」、「プロフィールとアップロードしたChatGPTの履歴を表示する」、「結果をダウンロード」、「プロフィール」があります。この中で、「結果をダウンロード」のボタンが赤い円で囲まれています。

スライド32

### トピックIDとトピック名の関係

分析後ダウンロードする分析結果のファイルでは、トピックは全てトピックIDで示されているため、トピック名が必要な場合は、このファイルを活用し、Vlookupなどを使ってトピック名を追加してください。

このスクリーンショットは、分析結果のファイルに含まれるトピックIDとトピック名に関するデータを示しています。右側の注釈には、トピックがすべてIDで示されているため、必要に応じてこのファイルを用いてVlookupなどでトピック名を追加する必要があると説明されています。

スライド33

### Webシステムの使い方

結果が出てきたらリンクをクリックしてファイルをダウンロード

ダウンロードされるのはエクセルファイルです。ブラウザのバージョンによっては、上記のような警告があるので、「保存」をクリックしてください。

このスクリーンショットは、ウェブブラウザでダウンロードされたファイルのダウンロードダイアログボックスを示しています。ダイアログには「download.xlsx」というファイル名が表示されています。右側の注釈には、ブラウザのバージョンによっては警告が出る可能性があるため、「保存」をクリックする必要があると説明されています。

スライド34

### 結果ファイルの見方

一番確率が高いトピック(アンダーバー区切り)基本的にlabel1が入るprob1の値とprob2やprob3の値がほとんど同じ(9割以内)の場合はそれもここに入るS40(その他)を無視する設定の場合、label1がS40なら無視されて、label2の値がここに入る

そのテキストに、そのトピックがあれば1、無ければ0  
縦に足して、テキスト数で割れば、そのトピックの言及率がわかる

このスクリーンショットは、分析結果のExcelファイルの表示画面を示しています。表にはトピックID、prob1、prob2、prob3などの列があります。右側の注釈には、トピックの確率や言及率の計算方法について説明されています。

スライド35

### 結果ファイルの見方

「挿入」でカウントするため、複数のトピックを持つテキストはトピック以外の列は同じ情報のまま複数行繰り返される

クロス表を作るには、属性列と、topics列を選択した状態で「挿入」→「ピボットテーブル」をクリック

このスクリーンショットは、Excelでクロス表を作成するための手順を示しています。注釈には、「挿入」メニューから「ピボットテーブル」を選択してクロス表を作成する方法が説明されています。

スライド36

### 結果ファイルの見方

このスクリーンショットは、Excelでクロス表を作成した後の表示画面を示しています。注釈には、クロス表の作成方法が説明されています。

