

カウンセリングデータにおけるトピックモデルを用いた文書分類

A Text classification for Counseling Data Using Topic Model

単 壮 † 加藤 昇平 †
Zhuang Shan Shohei Kato

1 はじめに

近年、多くの人が仕事や勉強などからのストレスによる心理問題を抱えており、社会問題にもなっている [1]。一般的な治療方法は、カウンセラーのカウンセリングを受けることである [2]。しかし、カウンセリングを受けるには時間的、金銭的負担が少なくない。また、プラバイシーを他人に言い難いため、多くの人が治療を放棄する。従い、誰でも気軽に受けられるカウンセリングシステムが求められている。カウンセリングシステムは人間のカウンセラーに取って代われないけれども、ある程度患者の病気が重くなっているうちに患者の気分を改善し、重病にならないような効果を持っている。また、カウンセリングシステムはコンピュータやその他の電子通信手段を媒介とし、コストが少なく、日常の使用にも便利である。更に、システムの自身及び端末機の暗号化を活用されるため、ユーザのプラバイシーを保護することができる。しかし、既存のカウンセリングシステムは人間の言葉でよく用いられる同義語や類義語、婉曲表現といった言語の多様性により、そもそも利用者の発言意図を理解することが難しい。すなわち、精度の高い自然語処理システムの構築は困難である。ストレスによる心理問題を抱えている人に便利、有効なカウンセリングを与えるため、本研究はカウンセリングシステムの開発に注目する。ユーザの質問に対してより高い精度で答えを与えるカウンセリングシステムを構築するために、ユーザがどのような心理的問題を抱えているか分析することが重要だと考えられる。

2 関連研究

コンピュータやその他の電子通信手段を媒介とするカウンセリングやソーシャル・サポートの試みは今日多様な形態で試みられている。カウンセリングにおける媒体、このような変化は、従来のダイレクトな面談とは異なる様々な効果を生むことが指摘されている。しかし、コミュニケーションの媒体は変わっても、システムの両端に人がいる、すなわち、最終的には人と人との間で営まれる活動であることに変わりはない。それで、人でなく、コンピュータ自身にカウンセリングを代行させることはできるように、カウンセリングシステムが開発されている。本章では既存のカウンセリングシステムを紹介し、優劣を述べる。カウンセリングシステムは MIT の人工知能研究者であった Weizenbaum が開発した ELIZA [3] (イライザ) と呼ばれるコンピュータ・プログラムに端を発する。ELIZA は自然言語による対話システムであり、ユーザがキーボードから入力した文字列に対し、文字列で応答する。ELIZA 自身はユーザが求める質問に対して具体的に答えずに、ユーザに発言を続けさせるようなメッセージのみを返すだけで、会話が続けられるような仕組みである。それで、すごく単純なプログラムであるため、すぐ飽きるだろうと考えた開発者自身の予想に反し、ELIZA と会話した人達は皆、あたかも人間と会話をしているような錯覚に陥る。このシステムには、. 応答するコンピュータ・プログラムに対し親しみが生じる現象は "ELIZA 効果" と呼ばれている。この ELIZA プログラムはすぐに精神医療領域からの注目も受けることとなった。このコンピュー

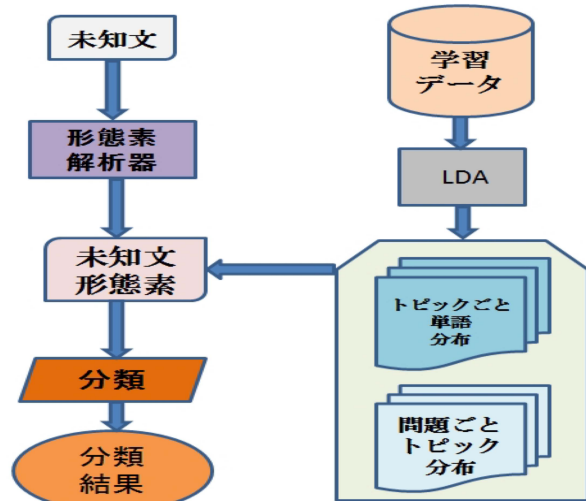


図1 Flow of the classification method

タ・プログラムを進展させれば、将来自動化された形の精神医療が可能になると真剣に考えたのである。研究者達は ELIZA 的なコンピュータ・カウンセリングのシステムを開発し、それがやがて臨床で使用できるようになるだろうと考えていた。もし臨床で使用可能になれば、コスト、マンパワー、時間的制約からの解放などの様々な点で精神医療に革新的な変化がもたらされるはずである。しかし、このアイデアは結局、実現することが出来なかった。

ELIZA 型のカウンセリング・システムが実現しなかった理由としては、精度の高い自然言語処理システムの構築の難しさに加え、利用者の発言意図を "理解" 出来ず、質問に対しての真ともうな答えを出せないと考えられるからだ。また、ELIZA 型のカウンセリングシステムはシステム主導で会話の自由度も乏しく、いかにも機械的で人間的温かみに欠ける。以上で、心理問題があっても、人間のような慰めて欲しい人に対して、単純に言葉を言い返すシステムは治療の効果が低いと考えられる。心理カウンセリングシステムは普通のコミュニケーションシステムと違い、人間とただ対話するだけではなく、ユーザの心理問題に対して、より人間的な知識を提供し、心理問題を解決することが必要となる [4]。ELIZA 型カウンセリングの不足している点について、改善するために Liu ら [5] は、より人間的な知識を求めることと、自然語処理システムの精度を高めることを注目し、新しいカウンセリングシステムを開発した。Liu らの調査によって、同じ心理問題がある患者の質問の表現が違っても、質問の本質がほぼ同じである。カウンセラーはそのような質問について、普段と同じく、あるいはよく似ている答えを与える。もし、大量な学習データがあれば、システムは圧倒的に多数のユーザの質問を答えられると考える。Liu らの研究では、彼は大量な心理問題と多くの人に認められる問題の答えをペアで収集し、学習データとする。Liu らが提案したカウンセリングシステムは、二つの部分に分けられている。

† 名古屋工業大学, Nagoya Institute of Technology

一つ目は学習部である．二つ目は判別部である．学習部では，収集したデータからキーワードを抽出し，キーワードによって答えの索引を作り，この索引とXMLの形式で保存しているデータと組み合わせてデータベースとする．判別部では，ユーザの質問の語義特徴量（関連性がある形態素）を抽出し，キーワードとする．そのキーワードを用い，索引によって学習されたデータベースからユーザの質問と最も似ている問題を選択し，ユーザから聞かれた問題の答えを出す．このようなカウンセリングシステムでは，キーワードの抽出はとても重要である．キーワードの抽出によって，システムがユーザに答えを正しく与えられるかどうかを決める．Liuらの研究では，「妻」「友達」という互いの関係を表す呼称あるいは，「後悔」「怒り」「悲しい」というユーザの感情を表す言葉が重要である．Liuらのキーワードに基づく検索技法を用いたカウンセリングシステムは確かにELIZA型のシステムより精度を高め，効果的な答えを出せるようになった．しかし，このようなシステムも限界がある．例えば，多くのキーワードは複数の類別の問題に出現する場合もある．それで，キーワードが一緒に全く違う類型の問題の答えを出す可能性が高くなるかもしれない．また，大量の心理問題からキーワードを取り出すにはすごく手数がかかる．以上から，心理問題を分類することと，問題からキーワードをよく選択できる文書分類技法が必要だと考えられる．

3 提案手法の流れ

図1に分類の流れを示す．提案手法では，心理問題に関する質問を収集し，学習データとする．心理問題に関する質問の類型により，問題のタグを作り，類別を分けて保存する．次にLDAを用いて学習データを分析する．学習データが類別のタグがついているため，ある類別のトピック分布と各トピックの中でどのような単語がよく登場するのを知ることが出来る．このようにLDAのパラメータを試し，最も良い問題ごとのトピック分布とトピックごとの単語分布を得るまで，学習データを訓練する．未知文を分類する時，まず形態素解析[6]を用いて，未知文を形態素単位に分割する．また，分類の結果に影響を与えないように，ストップワードを取り除く．最後，ストップワードをなくした未知文の形態素とLDAにより，学習データを訓練された結果を用いて未知文を分類する．本稿では形態素単位に分けられた語を単語として扱う．単語ごとにその単語がどのトピックに属するかを判別し，そのトピックがどの問題タグに属するかを判別することで，単語ごとに問題タグを判別する．なお，判別には生起確率を用いた．対象の単語の生起確率が最も高いトピックをその単語が属するトピックと判別する．トピックも同様に，対象のトピックの生起確率が最も高いタグをそのトピックが属するタグと判別する．これを未知文中の全ての単語に対して行い，未知文中で最も多く出現した問題タグをその未知文の問題タグとする．

4 Latent Dirichlet Allocation

LDA[7]は，文書の生成過程を確率的にモデル化したトピックモデルの一つであり，一つの文書中に複数のトピックが存在することを表現できる潜在的意味解析手法である．LDAによって，文書を構成するトピックの多項分布及び各トピックを構成する単語の多項分布を表現することができる．ここで，LDAによる文書を生成する確率的なモデルと本研究が使っている文書のLDAに生成された単語分布とトピック分布の構成を紹介する．

4.1 Unitgram Model

ある文章 W を構成する確率は W の全ての単語の出現確率の積だと考えられる．式(1)に文書 W の構成確率グラフィカルモデルを示す．

$$p(w) = \sum_{n=1}^N p(w_n) \quad (1)$$

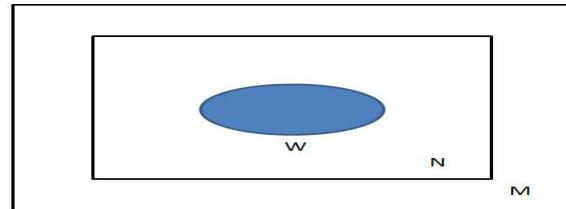


図2 Unitgram Model

4.2 Mixture of Unigrams

文章はただ一つのトピックを持っていることを仮定すれば，ある文章 W を構成する確率はトピック Z によって単語が選択された確率の積だと考えられる．式(2)に文書 W の構成確率グラフィカルモデルを示す．

$$p(w) = \sum_z p(z) \prod_{n=1}^N p(w_n | z) \quad (2)$$

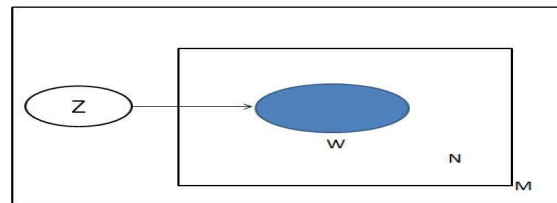


図3 Mixture of Unigrams

4.3 LDA 生成モデル

LDAは文章は複数のトピックを持っていることを仮定する．LDAにおいて文書は，まずその文書のトピック分布に従いトピックが選択され，そのトピックの単語分布に従って単語が選択される過程で生成されると仮定される．図4にLDA生成モデルのグラフ表現を示す．

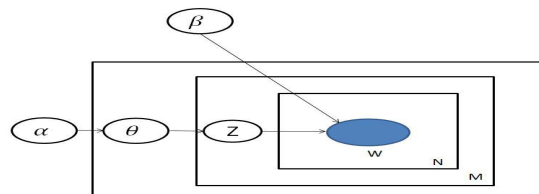


図4 Graphical model representation of LDA

α はトピックの事前分布の k 次元のハイパラメーターであり， β はトピックの単語分布のパラメーターである．また， N

表 1 タグ恋愛のトピック分布

| TAG. 恋愛 | |
|---------|----------|
| TOPIC | 生起確率 |
| Top.14 | 0.742518 |
| Top.32 | 0.243423 |
| Top.46 | 0.003237 |
| Top.2 | 0.002921 |
| Top.18 | 0.002763 |
| Top.3 | 0.002684 |
| Top.26 | 0.002658 |
| Top.38 | 0.002342 |

... ..

表 2 トピック 14 とトピック 38 の単語分布

| TOPIC.14 | | TOPIC.38 | |
|----------|----------|----------|----------|
| 単語 | 生起確率 | 単語 | 生起確率 |
| 彼女 | 0.116900 | 美しい | 0.004272 |
| 別れ | 0.043869 | 間接的 | 0.002899 |
| 彼 | 0.039412 | 我がまま | 0.002898 |
| 愛 | 0.021050 | 見証 | 0.002894 |
| 恋愛 | 0.020210 | 方法 | 0.002892 |
| 好き | 0.020188 | 条件 | 0.002887 |
| 私たち | 0.014798 | 上手 | 0.002884 |
| 私 | 0.014798 | 仕方ない | 0.002883 |
| 相談 | 0.013178 | 見る | 0.002882 |

... ..

は文中の単語数, M (定数) は文書の個数である. α と β が与えられたとき, トピックの混合分布 θ , N 個のトピック集合 z , N 個の単語集合 w は式 (3) によって与えられる.

$$P(\theta, z, w|\alpha, \beta) = p(\theta|\alpha) \prod_{n=1}^N p(z_n|\theta) p(w_n|z_n, \beta) \quad (3)$$

4.4 単語分布とトピック分布

LDA モデルを用いて文書の単語分布とトピック分布を得ることが出来る. 本研究では, Liu らの研究を参考に, BAIDUZHIDAO という中文問答知識プラットフォームでよく質問された 3 種類の心理問題「恋愛」「人間関係」「自己認識」を収集した. ここで, 恋愛に関する問題文は 387 問, 同僚やクラスメートや友達などの人間関係に関する問題文は 395 問, 自己認識に関する問題は 215 問, 計 995 問の問題文を収集した. データによって, LDA のパラメータをより良く設定すれば, 分類の精度を高められる. 本研究ではデータベースの問題文数と単語数によって, $\alpha = \frac{50}{\text{トピック数}}, \beta = 0.1, k = 50$ と設定し, 収集したデータの分析を行った. 単語分布では, あるトピックに選択された単語は生起確率の高い順で並べる. 単語分布と同じ, トピック分布では, ある問題の類別に選択されたトピックも生起確率の高い順で並べる. 表 1 に恋愛のタグのトピック分布の一部を示す. 表 2 に恋愛のタグに生起率最も高いトピック分布の 14 番トピックと生起率最も低いトピック分布の 38 番トピックの単語分布の一部を示す. 14 番のトピックでは, 「彼女」「別れ」「好き」という恋愛の問題をよく代表できる単語が出現した. 38 番のトピックでは, 「美しい」「我がまま」という単語は直接に恋愛のことを表しず, しかし, 恋愛の問題に登場する可能性が高い単語が出現した.

類別の判定例を示す. ここで, 未知文を形態素解析によって, 「自信」「友達」「成功」「信じない」「失敗」等の形態素単位に分割した. トピックごとの単語分布により, 「自信」という単語最も高い生起率でトピック 1 に属する. トピック 1 は最も高い生起率で自己認識の類別に属する. これで, 「自信」という単語により, 未知文が自己認識の類別に投票する.

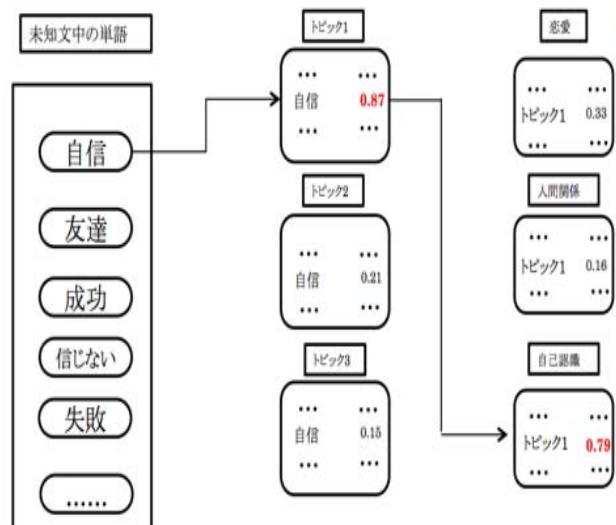


図 5 The example of classification

これを未知文中の全ての単語に対して行い, 最後の投票によって, 未知文の類別を判断する.

5 分類手法

本稿では形態素単位に分けられた語を単語として扱う. 単語ごとにその単語がどのトピックに属するかを判別し, そのトピックがどの問題タグに属するかを判別することで, 単語ごとに問題タグを判別する. なお, 判別には生起確率を用いた. 対象の単語の生起確率が最も高いトピックをその単語が属するトピックと判別する. トピックも同様に, 対象のトピックの生起確率が最も高いタグをそのトピックが属するタグと判別する. これを未知文中の全ての単語に対して行い, 未知文中で最も多く出現した問題タグをその未知文の問題タグとする. 図 5 に

6 分類実験

6.1 提案手法による分類実験

提案手法の有効性を確認するため, LDA による文書分類実験を行った. サンプル調査法により, 995 問から問題タグ毎にランダムにそれぞれ 100 問ずつ選出し, Leave-one-out 交差検定方法により, 分類実験を 300 回行った.

6.2 TF・IDF による分類実験

提案手法の有効性を確認するため, 改良した TF・IDF[?] による比較実験を行う. 提案手法による分類実験と同じ, 比較実

表3 提案手法による分類実験の結果

| 問題タグ分類結果 | 恋愛 | 人間関係 | 自己認識 | 正答率 |
|----------|----|------|------|-------|
| 恋愛 | 83 | 14 | 4 | 83.0% |
| 人間関係 | 12 | 80 | 8 | 80.0% |
| 自己認識 | 7 | 15 | 78 | 78.0% |

総正答率:80.3%

表4 TF・IDFによる分類実験の結果

| 問題タグ分類結果 | 恋愛 | 人間関係 | 自己認識 | 正答率 |
|----------|----|------|------|-------|
| 恋愛 | 67 | 25 | 8 | 67.0% |
| 人間関係 | 19 | 64 | 17 | 64.0% |
| 自己認識 | 16 | 27 | 57 | 57.0% |

総正答率:62.6%

験もサンプリング調査法により、995問から問題タグ毎にランダムにそれぞれ100問ずつ選出し、Leave-one-out交差検定方法により、分類実験を300回行った。

6.2.1 改良したTF・IDFの式

Kuangら[8]はTF・IDFの問題に対して、文書の類別の違いを考慮したTF・IDFの式を提案した。式(4)にKuangらが提案した式を示す。

$$TF \cdot IDF \cdot C_i = tf_{ij} \times idf_i \times C_i \\ = \frac{D_{ij}}{\sum_k D_{kj}} \times \log \frac{N}{n_i} \times \frac{1}{n_i - m_i + 1} \quad (4)$$

ここで m_i はある類別に単語 i を含む文書の総数である。改良したTF・IDFの式は分類の性能を高めた。

TF・IDFはよくコサイン類似度と一緒に使って文書間の関連度を計算する。コサイン類似度とは、ベクトル空間モデルにおいて、文書同士を比較する際に用いられる類似度計算手法である。コサイン類似度はベクトル同士の成す角度の近さを表現するため、三角関数の普通のコサインの通り、1に近ければ類似しており、式5によって与えられる。

$$\cos \theta = \frac{\sum_{i=1}^n (A_i * B_i)}{\sqrt{\sum_{i=1}^n A_i^2} * \sqrt{\sum_{i=1}^n B_i^2}} \quad (5)$$

6.2.2 分類の流れ

本稿では、文書内の全単語に対するTF・IDF・ C_i のベクトルを作成し、 $TFIDF'$ と表す。分類手法の流れとしては、まず未知文の $TFIDF'$ を算出し、データベース内の各文書 $TFIDF'$ とのコサイン類似度を算出する。コサイン類似度が高い文書を上位10個選出し、10個中最も多いタグを未知文のタグとする。

6.3 実験の結果

実験結果を表3と表6.3に示す。提案手法と比較対象の問題タグ毎に100個の問題が分類された分布を表す。提案手法は問題タグごとの正答率は全て75%を超えており、総正答率も80.3%で判別できていることが分かる。

7 まとめと今後の課題

本稿では、精度の高いカウンセリングシステムを構築するため、既存カウンセリングシステムの不足点を述べ、文書分類の必要性を確認する上にトピックモデルを用いた文書分類法を提案した。また、分類実験において、三種類の心理問題の総正答率は80.3%を超えることで、提案手法が「恋愛」「人間関係」「自己認識」などの複雑な心理問題を高い能力で分類できるこ

とが示唆された。今後の課題としてはより多種類の問題に対応する実験を行うとともに、分類された問題文に対応する回答文を出力できる新たな手法を開発し、カウンセリングシステムを実装したい。また、カウンセリングシステムをifbot[9]などの音声、表情を出せるロボットに移植し、擬人化なロボット・カウンセリングシステムを実装したい。

謝辞

本研究は、一部、文部科学省科学研究費補助金(課題番号25280100、および、25540146)の助成により行われた。

参考文献

- [1] 石村光資郎 and 田中暢子 and 加藤千恵子 and 土田賢省 and 後藤武秀: オンラインカウンセリングによる海外赴任者のメンタルヘルスの改善に関する考察, 工業技術: 東洋大学工業技術研究所報告, Vol. 32, No. 13499955, pp. 7275 (2010).
- [2] 大石由起子, 木戸久美子, 林典子: ピアサポート・ピアカウンセリングにおける文献展望, 口県立大学社会福祉学部紀要, Vol. 13, pp. 107121 (2007).
- [3] Weizenbaum, J.: ELIZAA Computer Program For the Study of Natural Language Communication Between Man and Machine, Commungicats of the ACM, Vol. 9, No. Nonr-4102(01), pp. 3635 (1966).
- [4] 藤野博: 擬人化エージェントによるカウンセリング・システム構築の試み, <http://www.u-gakugei.ac.jp/hfujino/botmama/botmama.html>, pp. 1316 (2009).
- [5] Liu, Y., Liu, M., Lu, Z. and Song, M.: Extracting Knowledge from On-Line Forums for Non-Obstructive Psychological Counseling Q・A System, International Journal of Intelligence Science, Vol. 2, No. 220066, pp. 4048 (2012).
- [6] W. Che, Z. H. L. and Liu, T.: LTP: A Chinese Language Technology Platform, COLING '10 Proceedings of the 23rd International Conference on Computational Linguistics, pp. 1316 (2010).
- [7] D. M. Blei, A. Y. N. and Jordan, M. I.: Latent dirichlet allocation, J. Mach. Learn, Vol. 3, pp. 9931022 (2003).
- [8] Kuang, Q. and Xu, X.: Improvement and Application of TFIDF Method Based on Text Classification, Internet Technology and Applications, 2010 International Conference on, pp. 14 (2010).
- [9] <http://www.ifoo.co.jp/sub7.html>.