

# 説明の難易度と講義内容との一致度を考慮した講義用語検索

竹村 直規<sup>1</sup>      小林 伸行<sup>2</sup>      椎名 広光<sup>3</sup>

<sup>1</sup> 岡山理科大学大学院 総合情報研究科 情報科学専攻

<sup>2</sup> 山陽学園大学 総合人間学部 生活心理学科

<sup>3</sup> 岡山理科大学 総合情報学部 情報科学科

i15im03tn@ous.jp<sup>1</sup>, koba\_nob@sguc.ac.jp<sup>2</sup>, shiina@mis.ous.ac.jp<sup>3</sup>

## 1 はじめに

近年、日本への外国人留学生は増加傾向にある [1]. 海外から日本へ留学するにはある程度の日本語能力が必要となる. この日本語能力の証明として日本語留学試験 (EJU) [2], 日本語能力試験 (JLPT) [3] といった試験が存在する. このような日本語試験の受講を入学条件としている大学も少なくない. 大学教育においても, 留学生向けの講義の工夫がなされおり, 報告がなされている [4]. しかし, 大学における専門的な知識を教授する講義では, 日本語試験で出題されるような日常用語よりも, より特殊な用語を用いることが多い. そういった専門的な用語が頻出した場合, 留学生によっては講義の理解が難しいと考えられる.

そこで, 講義に連動して専門用語を検索する際に, 日本語の理解も併せて考えるために母国語ではなく理解がしやすい日本語での検索結果の出力が必要であると考えられる.

一方, 日本語の難易度に関連する研究としては, 文章の難易度の推定 [5] や語彙の平易化システム [6] については, 提案されている.

本研究では, 日本語を母語としない学生に向けて講義に則した検索を行う際, Wikipedia[7] や WordNet[8], Web 検索の結果を, 説明文の難易度や講義との適合度を示す指標による順位付けを行うシステムを提案する.

## 2 使用したデータ

講義の単語を検索する検索データベースと, 対象となる講義について述べる.

### (1) 検索データベース

検索データベースとして, Wikipedia, Wikipedia Simple English 版 [9], WordNet の 3 つを使用した.

Wikipedia では曖昧性の低減を, WordNet では表記ゆれの補完や網羅性を高める.

### (2) 講義データ

講義データは, 岡山理科大学総合情報学部情報科学科で行われている講義のうち 1 年次開講科目の “プログラミング基礎” と 2 年次開講科目の “データベース” の VOD 講義の字幕データ及び講義資料を使用した.

## 3 講義用語検索の概要

検索システムとしては, 検索表示機能, リーダビリティの判定, 講義との適合度の項目からなっている.

### (1) 検索表示機能 (図 1)

(1-1) Wikipedia の検索結果の提示: 同義語の場合は, #転送及び#REDIRECT 先の記事ページを検索結果とする.

(1-2) WordNet の検索結果の提示: WordNet では, 単語にリンクされている複数の synset 全てを提示する.

(1-3) Wikipedia の simple english の提示: Wikipedia の simple english ページの記事を取得し, それを Microsoft Translator Text API[10] によって機械翻訳したものを表示する.

(1-4) 検索エンジンの結果: Microsoft が提供する Bing Search API[11] を使用した. 本研究では検索結果の上位 5 件を表示する.

### (2) 検索表示結果の難易度判定

検索結果については, 意味情報の記述がよりやさしいか難しいかといった読みやすさの難易度の判定を決定木及び  $k$  近傍法による機械学習によって行っている.

### (3) 講義との適合度

複数の語義を持つ多義性を解消するために単語の意味情報を, 講義で使われる単語から TF-IDF を用いた特徴抽出を行うことで, 講義の内容に適合した表示





図 3: 学習用データの各特徴量

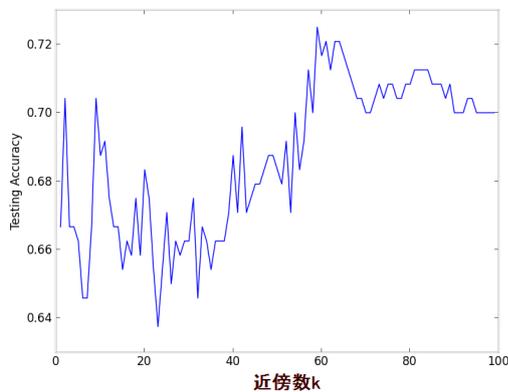


図 4:  $k$  近傍法の近傍数

## 5 多義性の解消のための適合順の提示

語には多義性があり、Wikipedia, WordNet, 検索エンジンを使用した意味情報も多義性に合わせた記述になっている。そこで、講義に連動させた意味の提示を行う。

例えば、データベースやプログラミングの講義でよく使われる単語“配列”に対して WordNet では、10 個の意味を用意している (表 1)。このうちデータベースやプログラミングの講義では、音楽に関連する意味が含まれていることは不自然であるため、これらをそのまま提示するのは好ましくないと考えられる。そこで、講義に則して、意味の提示順を変えることとした。意味の提示順については、TF-IDF を利用して、適合度の指標を決めている。

### 5.1 適合度の処理手順

一科目の講義全体や各回の講義での単語の適合度を、TF と IDF から作られる TF-IDF を利用して求める。次にその手順を示す。

表 1: WordNet での単語“配列”の表示結果

単語	記述	一致度
配列	整理や分類のための整った構造。	0.1798
	物事が論理的順番または連続するパターンで続く連続した配置関係。	0.0854
	ある順番に並べる。	0.0496
	順番に、論理的にあるいは、線上のように配列する。	0.0454
	グループの個々の要素を論理的にわかりやすい配列したもの。	0.0389
	順序正しく置くこと。	0.0248
	人や物の生来の、あるいは後天的な気質、または特徴的な傾向。	0.0093
	物の置き方の空間的性質。	0.0049
	規則正しい取り決め。	0
音楽作品を編曲する行為。	0	

- (1) 講義資料をテキスト形式に変換し、形態素に分割する。
- (2) 対象文書での単語の比率を出現比率を TF として、ストップワードとしている単語の一部を除外した単語集合に対して計算する。

$$TF(t, d) = \frac{\text{対象文書での検索語の出現回数}}{\text{対象文書の単語総数}}$$

ストップワードは、“もの”、“する”、“ある”、“こと”、“から”、“この”、“いる”、“ない”、“よう” のようなひらがなで、かつ名詞や動詞を除いている。

- (3) 逆文書頻度 (IDF) を下記の式で計算する。

$$IDF(t) = \log\left(\frac{\text{文書の総数}}{\text{単語 } t \text{ を含む文書数} + 1}\right)$$

- (4) TF-IDF を、TF と IDF から求める。

$$TF\text{-}IDF(t, d) = TF(t, d) \cdot (IDF(t) + 1)$$

- (5) 適合度の計算

(5-1) WordNet の場合は、検索語の WordNet の同義語リンクに含まれている単語の TF-IDF の合計を、適合度とする。検索結果としては、適合度の高い順に概念記述の高い順に記述する。

(5-2) Wikipedia, 検索エンジンの場合は、記事概要部に現れる単語の TF-IDF の合計を適合度とする。

## 6 難易度と適合度を合わせた総合指標

検索結果の読み取りやすさと講義との適合度を合わせた総合的な順位付けを行う。総合的な指標には、第 4 章の  $k$  近傍法による簡易クラス予測値と第 5 章の TF-IDF 値の重み  $\alpha$  と  $\beta$  とする線形和を用い、講義で重要となる単語に対して人手で評価順と最も近くなるように重みを設定している。

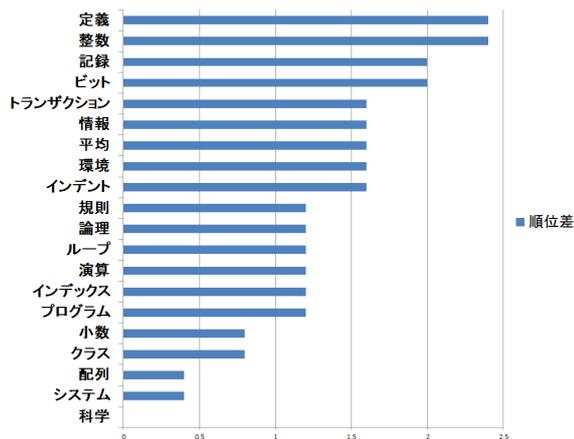


図 5: 検索結果の順位付け評価

$$\text{総合指標} = \alpha \cdot \text{簡易クラス予測値} + \beta \cdot \text{TF-IDF}$$

講義“プログラミング基礎”で使用された単語のうち 20 個を Bing Search による検索結果上位 5 件と人手による評価順位との差を図 5 に示す。講義で使用される単語全体の順位差平均差は 1.264 であった。“配列”, “プログラム”, “ループ”などのプログラミングの講義でよく使用される単語は順位差が小さくなっている。一方で, “定義”, “整数”, “記録”といったプログラミングの講義には出ないが, 特徴を上手く捉えられない単語は順位差が大きくなっている。今回の実験では, 講義との適合度である TF-IDF 値の幅が大きいため, その部分の影響が大きくなっている。全体の順位差は小さくなったが, 講義との適合度に大きく影響される結果になったのではないかと考えられる。

## 7 おわりに

本研究では, 説明の難易度及び講義の適合度に即した専門用語の検索システムを開発した。留学生にもわかり易い意味情報を表示する方法として, 決定木と  $k$  近傍法を用いた難易度判定を行った。また, 意味情報の剪定として, TF-IDF を用いた適合度により, 講義の内容に合わせた検索結果の表示を行っている。最後に総合的な順位付けを行った結果, 全体の順位差は抑えられたが一部の単語では上手く順位付けすることができなかった。しかし, そういった一部の手く順位付けできなかった単語群は講義との一致度が低いため, 最終的な単語の表示結果の順番を変える指標になると考えられる。

## 参考文献

- [1] 文部科学省, “留学生 30 万人計画の実現に向けた留学生の住環境支援の在り方に関する検討会報告書”, [http://www.mext.go.jp/b\\_menu/houdou/26/08/\\_ics\\_files/afiedfile/2014/08/29/1350840\\_01\\_1.pdf](http://www.mext.go.jp/b_menu/houdou/26/08/_ics_files/afiedfile/2014/08/29/1350840_01_1.pdf)
- [2] 日本留学試験, [http://www.jasso.go.jp/ryugaku/study\\_j/eju/index.html](http://www.jasso.go.jp/ryugaku/study_j/eju/index.html)
- [3] 日本語能力試験: <http://www.jlpt.jp/index.html>
- [4] 古本裕子, 苗田 敏美, 松下美知子, 専門教育における留学生の日本語-日本人学生との比較を通じた分析-, 金沢大学留学生センター紀要 9, pp.21-33, 2006.
- [5] 近藤陽介, 松吉俊, 佐藤理史. 教科書コーパスを用いた日本語テキストの難易度推定. 言語処理学会第 14 回年次大会発表論文集, pp.1113-1116, 2008.
- [6] 梶原 智之, 山本 和英, 語釈文を用いた小学生のための語彙平易化, 情報処理学会論文誌, Vol.56, No.3, pp.983-992, 2015.
- [7] Wikipedia, [https://en.wikipedia.org/wiki/Main\\_Page](https://en.wikipedia.org/wiki/Main_Page)
- [8] WordNet, <https://wordnet.princeton.edu/>
- [9] Simple English Wikipadia, [https://simple.wikipedia.org/wiki/Main\\_Page](https://simple.wikipedia.org/wiki/Main_Page)
- [10] Microsoft Translator Text API, <https://www.microsoft.com/en-us/translator/translatorapi.aspx>
- [11] Bing Search API, <http://azure.microsoft.com/ja-jp/services/cognitive-services/search>
- [12] Altman, N. S. (1992). “An introduction to kernel and nearest-neighbor nonparametric regression”. The American Statistician. 46 (3): pp. 175-185.
- [13] やさしい日本語ニュース, <http://www3.nhk.or.jp/news/easy/index.html>
- [14] 読売オンライン, <http://www.yomiuri.co.jp>
- [15] 徳弘 康代, 日本語学習のためのよく使う漢字 2100, 三省堂, 2010.