

# 単語ベクトルに基づく記録文書の内容検索

三澤 虎遊汰<sup>†</sup> 斬 展<sup>‡</sup> 柴田 千尋<sup>†</sup> 田胡 和哉<sup>†</sup>

東京工科大学コンピュータサイエンス学部<sup>†</sup>

東京工科大学大学院バイオ・情報メディア研究科コンピュータサイエンス専攻<sup>‡</sup>

## 1 背景

近年、単語の同義語や上位下位語などの概念情報を登録したデータベースを用いた検索など様々なあいまい検索が存在する。そういった検索手法は概念情報によって類語を判断してしまう。そのため、ユーザの考える類語だけでなくそれ以上の類語で検索してしまう。これは複数人が編集する文書における検索手法には同義語や類語の判断が困難になるため適さない。そのため類語をクラスタリングしユーザの考える類語のみで検索することが必要となる。

本稿では検索時に複数の類語をクラスタリングすることで情報検索を行うツールを提案する。

## 2 提案手法

### 2.1 全体の流れ

図 1 に概念構成図を示す。ユーザの入力した単語に対し、複数の類語を抽出し、抽出された類語をグループ単位で分ける。例として「言語」という単語から「日本語、英語、Java、C」という 4 つの単語が類語として抽出された場合、国の言葉に関する「日本語、英語」とプログラミング言語に関連する「Java、C」のグループに分ける。そして複数のグループをユーザに提案し、選択したグループに含まれる全ての単語で検索を行いユーザに提供する。これによりユーザの求める情報のみの提供が可能となる。

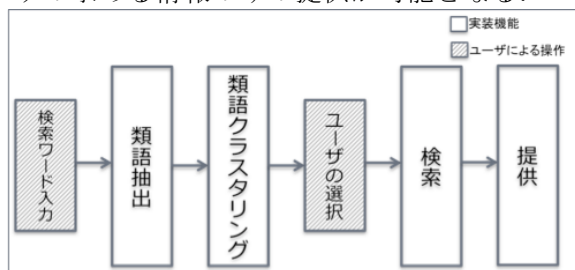


図 1 概念構成図

Concept search using vector representations of words from archive texts

<sup>†</sup> Koyuta MISAWA, Chihiro SHIBARA, Kazuya TAGO

School of Computer Science, Tokyo University of Technology

<sup>‡</sup> Kin TEN

Bionics, Computer and media science, Entrepreneurship program, Tokyo University of Technology Graduate school.

### 2.2 必要となる機構

図 1 の実現には以下の機構が必要となる。

#### ● 単語のベクトル化機構

文書データに出現する単語を文脈による単語の関係性からベクトル表現を獲得する。検索時に入力された単語はこの作成した単語ベクトルのリストから選択される。

#### ● 類語の抽出機構

機械が類語や同義語を評価し抽出する機構である。抽出された単語ベクトルと数値的に近いベクトルを抽出することで、類語や同義語として評価することが可能となる。(図 1「類語抽出」)

#### ● 類語クラスタリング機構

語群からそれぞれ関連する単語同士でクラスタを分ける機構である。抽出された全ての単語ベクトル同士の距離によりクラスタを設定する。(図 1「類語クラスタリング」)

#### ● 検索機構

選択したクラスタに含まれる全ての単語に対してキーワード検索を行い、マッチした箇所をユーザに提供する。(図 1「検索」)

## 3 実装

### 3.1 実装対象

東京工科大学のクラウドサービスセンター[1]では古谷によって開発された「板」[2]が運用されている。板はプロジェクト内での情報共有を行うテキストエディタであり複数人での同時編集が可能である。本提案手法は板上の検索システムとしてここに実装した。

### 3.2 各機構の実装方法

「単語のベクトル化機構」と「類語の抽出機構」は word2vecd で実装することで、単語を質の良い低次元ベクトルへと写像することができ、結果の高精度化および処理速度の低下の防止を実現できる。また、「クラスタリング機構」については K-means++法を用いて実装した。これにより類語として抽出されたベクトルを自動で複数クラスタに分ける。

