

単語ベクトルの結合学習を用いた近現代語の意味変化の分析

相田 太一・小町 守（東京都立大学）

小木曾 智信（国立国語研究所）

高村 大也（産業技術総合研究所）

持橋 大地（統計数理研究所）

概要：コーパス間で意味の違いがある単語を検出するタスクは、主に word2vec や BERT などから得られる単語ベクトルによって行われる。ただし、言語学者や社会学者がこれらの手法を適用する実際の状況では、計算資源は限られており、BERT などのような計算コストの高い手法を導入することは難しい。そこで本研究では、限られた計算資源でこのタスクを実行するために計算コストの低い既存の手法を拡張する。実験より、拡張した手法が英語および日本語で既存の手法と同等またはそれ以上の結果を示すことを確認した。さらに、各手法の訓練時間を比較し、日本語のデータについて包括的な分析を行ったところ、拡張した手法が高速に学習し意味の変化した単語を適切に検出することを示した。

キーワード：単語ベクトル, 通時的変化, コーパス言語学

Analysis of Semantic Change in Modern Japanese Using Joint Learning of Word Vectors

Taichi Aida / Mamoru Komachi (Tokyo Metropolitan University)

Toshinobu Ogiso (National Institute for Japanese Language and Linguistics)

Hiroya Takamura (National Institute of Advanced Industrial Science and Technology)

Daichi Mochihashi (The Institute of Statistical Mathematics)

Abstract: The task of detecting words with semantic differences across corpora is mainly addressed by word representations such as word2vec or BERT. However, in the real world where linguists and sociologists apply these techniques, computational resources are typically limited. In this paper, we extend an existing simultaneously optimized model that can be trained on CPU to perform this task. Experimental results show that the extended models achieved comparable or superior results to strong baselines in both English corpora as well as Japanese corpora. Furthermore, we compared the training time of each model and conducted a comprehensive analysis of Japanese corpora.

Keywords: Word Embeddings, Historical Language Change, Corpus Linguistics

1. まえがき

異なる時代・分野間において、単語は異なる使われ方をすることがある。近年、この検出には単語ベクトルが広く用いられており、言語学・社会学だけでなく辞書学における分析の一助となることが期待されている[1]。本稿では、異なる時期・分野に対応した従来の単語ベクトル学習手法の問題点に対し、2つの拡張手法を提案する。実験では、提案した拡張手法の優れた性能および学習速度を示し、近現代日本語の文書間において意味の変化が予測された単語の網羅的な分析を行った。実験に用いたソースコードは <https://github.com/a1da4/pmi-semantic-difference> より公開予定である。

2. 関連研究

従来、異なる時代・分野を比較して異なる用例の単語を検出するために、頻度に基づく手法が使われていた[2]が、Mikolov らによって単語の情報をベクトルとしてより効果的に表現できる Word2Vec [3]が提案され、広く用いられるようになった。単語ベクトルを異なる時代・分野間の文書データに対応させるには alignment とよばれる手法が一般的である[4]。これは、各時期・分野のデータで単語ベクトルを独立に学習した後、線形変換により時期・分野間で対応づけを行う手法である。この手法は容易に導入できる反面、異なる時期・分野間のベクトルを線形で対応づけができるという強い仮定をおいている。そこで、この問題を回避するために、2つの手法が提案され

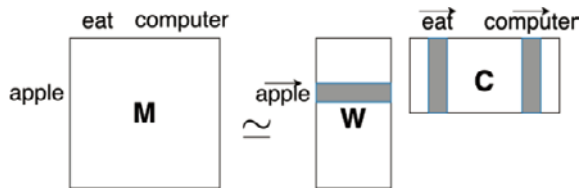


図 1 PMI-SVD の概要図.

た.

1 つ目は、Yao らによる動的な単語ベクトル (Dynamic Word Embeddings: DWE) である[5]. これは、制約条件を設定して各時期の単語ベクトルを同時に学習することで、線形の対応づけを回避する方法である。しかし、この学習手法には 3 つのハイパーパラメータが存在し、その設定に敏感であるため、膨大なハイパーパラメータ探索が必要である。

2 つ目は、Dubossarsky らによる、調査対象の単語だけ時期・分野で区別して単語ベクトルを学習する方法 (Temporal Referencing: TR) である[6]. これは、分析対象の時期・分野の異なる文書を 1 つの大きな文書として扱い、事前に用意した調査対象の単語リストにある単語のみ時期・分野を区別して単語ベクトルを学習する方法である。この手法はこれまでに指摘されていた線形の対応づけや膨大なハイパーパラメータ探索が不要となる一方、事前に調査対象となる単語のリストを用意する必要がある。

上記の手法以外にも、近年の深層学習の発展により、与えられた文脈から単語ベクトルを直接獲得できる BERT [7]などの事前学習済み言語モデルも提案されている。しかし、BERT による文脈を考慮した手法は Word2Vec などによる文脈を考慮しない手法と比較して膨大な計算資源を必要とするため、本研究では Word2Vec に基づいた文脈を考慮しない手法に着目する。

また、従来の研究では、英語の “gay” のように意味変化が自明な単語に対する分析が多い。そこで、本研究では各単語ベクトルから得られた意味変化の自明でない単語にも着目し、網羅的な分析を行った。

3. 手法

本研究では、TR に対して 2 つの拡張を行った。単語ベクトルは TR と同様に、Word2Vec の skip-gram with negative sampling との等価性が示されている PMI-SVD を用いた[8]. 概要図を図 1 に示す。これは、調査対象の単語とその周辺語で頻度行列を作成し、自己相互情報量を計算した行列 M を行列分解することで Word2Vec と同様の単語ベクトルの集合 W 、その周辺語ベクトルの集合 C を獲得できる手法である。

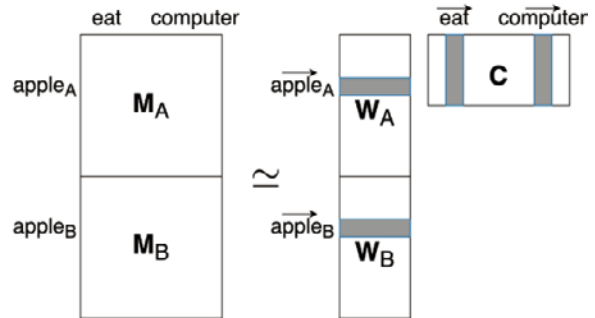


図 2 PMI-SVD_{joint} の概要図.

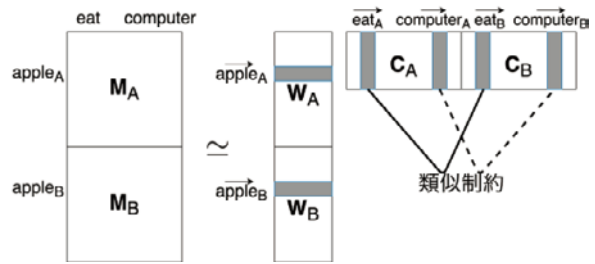


図 3 PMI-SVD_c の概要図.

1 つ目は、語彙中の全ての単語を対象単語とすることである (PMI-SVD_{joint}). 概要図を図 2 に示す。任意のコーパス A, B に対して、共起行列から自己相互情報量の行列 M_A, M_B を計算し、それらを結合してまとめて行列分解することで単語ベクトル W_A, W_B を獲得する。TR では事前に調査対象となる単語のリストが必要であるが、実際に任意の文書間に適用する際にはそのようなリストが存在しない場合が多い。そこで、通常の単語ベクトルを学習する時に選定する語彙に含まれる全ての単語を調査対象の単語とみなすことで、この問題を解消する。

2 つ目は、周辺単語のベクトルも時期・分野で区別して学習することである (PMI-SVD_c). 概要図を図 3 に示す。TR が採用している単語ベクトルは、図 2 のように、学習の際に周辺単語の時期・分野間における変化を考慮していない。そこで、制約条件を設けて周辺単語も時期・分野間の変化を考慮 (C_A, C_B) し、ベクトルの学習を行う。この手法は DWE の簡略版に近く、3 つのハイパーパラメータを持つ DWE に対して PMI-SVD_c はハイパーパラメータを 1 つだけ持っている。4 節では両者の性能および学習時間の比較を行った。

4. 実験

2 つの時期の文書において、語彙全体から意味の変化した単語を抽出する実際の状況に近いタスクを設定し、提案した 2 つの拡張手法 PMI-SVD_{joint}, PMI-SVD_c と以下 4 つの既存手法を

例：置き換え率 0.8 の時
猫→寿司 へと変化する単語を生成

戦前	戦後
吾輩は猫である。	私は猫を飼っています。
鮨の寿司を食べる。	日本の 寿司 は海外でも人気だ。 猫 とれたての魚を使った 寿司 だ。

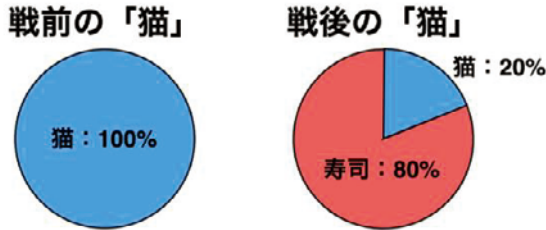


図4 意味が「猫」から「寿司」に変化する単語「猫」の擬似的な生成。

比較した。

- Word2Vec_{align} : 各時期で単語ベクトル Word2Vec を独立に学習した後, Hamilton らによる線形対応づけ[4]を行った。
- PMI-SVD_{align} : 各時期で単語ベクトル PMI-SVD を独立に学習した後, Hamilton らによる線形対応づけを行った。
- DWE : 各時期で同時に学習する手法. 3つのハイパーパラメータはそれぞれ 10^{-3} から 10^3 までの 7 通りの中から探索した (PMI-SVD_c も同様に探索を行った)。
- BERT : huggingface で公開されている事前訓練済みの BERT-base (日本語版は東北大が提供するモデル) を用い, 各時期の単語を代表するベクトルは平均によって獲得した。

今回は日本語と英語の 2 つの言語で評価を行った。日本語では, 近代雑誌コーパスに「中央公論」「文藝文集」を追加したデータを戦前と戦後で 2 つに分けて調査した。英語では, Corpus of Historical American English (COHA) の 1900 年代と 1990 年代のデータを用いた。各時期で 100 回以上出現する名詞・動詞・形容詞・副詞を調査対象の単語とし, 日本語では UniDic の語彙素で単語のまとめ上げを行った。評価では, 同じ単語における異なる時期のベクトルの余弦類似度に基づいて語彙中の単語を並び替え, 評価セットに含まれる単語の逆順位の平均を算出した。また, 意味の変化した単語を検出する様子を可視化するために, 再現率でも評価を行った。

まず, 意味変化する単語を擬似的に生成して評価を行った。擬似的な生成の様子を図 4 に示す。2つの時期間で意味が猫から寿司へと変化する

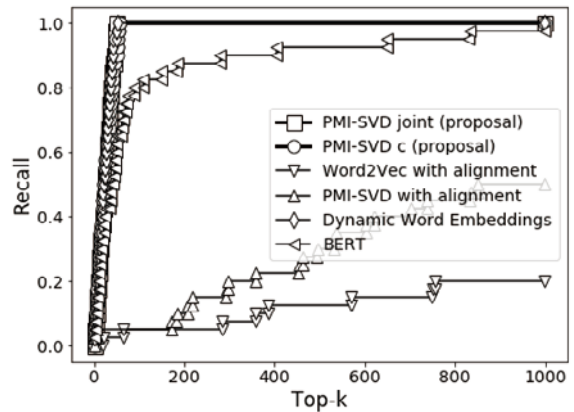


図5 日本語における擬似的に生成した意味の変化する単語を検出の様子 (再現率)。

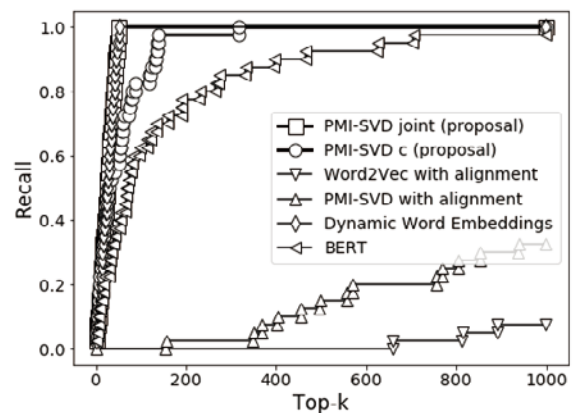


図6 英語における擬似的に生成した意味の変化する単語を検出の様子 (再現率)。

単語「猫」を擬似的に生成するために, 新しい方の時期の文書中に出現する単語「寿司」を「猫」に置き換え, 最終的に「寿司」の文脈で単語「猫」が出現する頻度の割合が任意の置き換え率となるように調整した。今回の実験では, 置き換え率を 1.0 とし, ある単語が 2つの文書において完全に異なる文脈で使われるように設定した。また, 置き換える単語のペアは, 単語ベクトルの余弦類似度が両方の文書において 0.01 以下であるものからランダムに 50 ペア抽出した。本実験では, 擬似的に生成した 50 単語のうち 10 単語を DWE と PMI-SVD_c のハイパーパラメータ探索に, 40 単語を評価に用いた。

表1 擬似的に生成した意味の変化する単語を用いた平均逆順位。

手法	日本語	英語
Word2Vec _{align}	0.0022	0.0004
PMI-SVD _{align}	0.0171	0.0010
DWE	0.0913	0.0835
BERT	0.0776	0.0590
PMI-SVD _{joint}	0.0737	0.0933
PMI-SVD _c	0.0781	0.0870

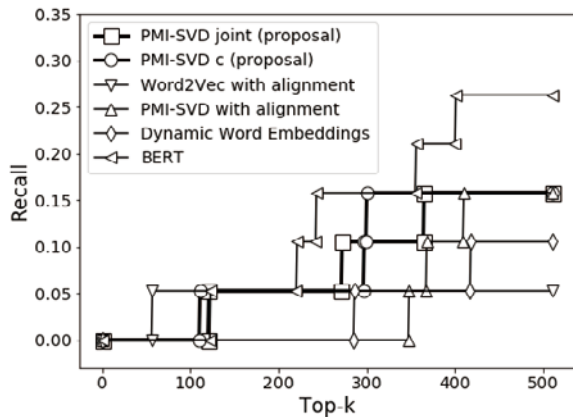


図7 日本語における実際に意味変化が報告されている単語を検出する様子 (再現率).

擬似的に生成した単語での実験結果を表1に、意味の変化した単語を検出する様子 (再現率) の推移を図5,6に示す. 表1より, 提案手法が既存手法と比べて同等またはそれ以上の性能であることが確認できる. また, 図5,6からも提案手法が意味変化した単語を上位で捉えていることがわかる.

ここで, Hamilton らの線形で対応づけを行う手法 (Word2Vec_{align}, PMI-SVD_{align}) について着目すると, 表1より他の手法と比べて平均逆順位が低いことがわかる. 図5,6の再現率の推移でも他の手法に比べて意味変化した単語の検出に苦戦していることから, 線形で対応づけを行う手法は, 同じ単語でも意味や用例が大きく離れる文書間における分析には適していないと考える.

次に, 実際に意味の変化が報告されている単語を用いて評価を行った. 日本語では間瀬らが作成したリスト[9]をハイパーパラメータ探索と評価に, 英語では Tahmasebi らのリスト[10]をハイパーパラメータ探索に, Kulkarni らが作成したリスト[11]を評価に用いた.

表2 実際に意味が変化した単語を用いた平均逆順位.

手法	日本語	英語	学習時間
Word2Vec _{align}	0.00137	0.00040	6m22s
PMI-SVD _{align}	0.00091	0.00100	3m26s
DWE	0.00058	0.00047	30h20m
BERT	0.00163	0.00250	2h23m
BERT-tiny	0.00078	0.00100	12days
BERT-mini	0.00119	0.00135	2weeks
PMI-SVD _{joint}	0.00131	0.00186	2m58s
PMI-SVD _c	0.00120	0.01045	26m01s

結果を表2に, 意味の変化した単語を検出する様子を図7,8に示す. 表2および図7,8より, 提案した2つの拡張手法と事前訓練済み BERT が両言語において他の手法を上回っている事が確認できる. また, 英語のデータにおける各単語べ

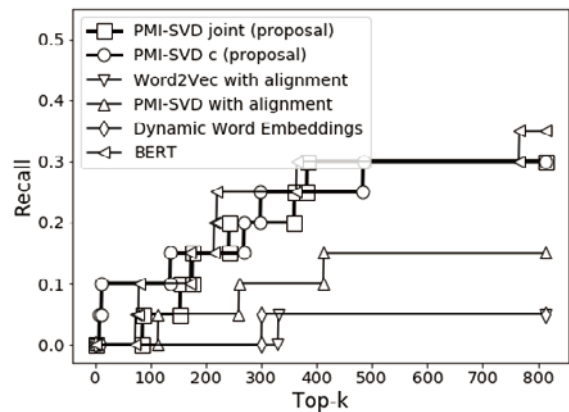


図8 英語における実際に意味変化が報告されている単語を検出する様子 (再現率).

クトル獲得手法の学習時間を比較すると, PMI-SVD_{joint} がどの既存手法よりも効果的に学習している事がわかる.

ここで, DWE とその簡略版に近い PMI-SVD_c についての比較を行う. 両者は共にハイパーパラメータ探索が必要であるが, PMI-SVD_c は DWE に比べて探索時間を大幅に削減し, 同等以上の性能を示している.

これまでの実験では, 大規模な文書で事前に訓練された BERT を用いて評価を行ってきた. そこで, 本実験で扱う文書のみを用いて BERT モデルの再訓練を行い, 平均逆順位および学習時間の比較を行った. 今回はデータ量の問題から, 公開されている BERT-base よりもパラメータ数の少ない BERT-mini, さらに少ない BERT-tiny モデルを訓練した. 表2より, 提案した2つの手法が BERT-tiny, BERT-mini モデルの両方を上回る性能を発揮することが確認できる. また, BERT モデルは大規模な計算資源で2週間程度の学習時間を必要とするのに対し, 提案手法は CPU のみで高速に学習できていることがわかる.

5. 分析

ここでは, 4節の定量的な評価で高い性能を示した PMI-SVD_c と BERT について, 日本語のデータにおける定性的な分析を行った. まず, 各手法が意味変化の可能性が高いと予測した上位10単語を分析した.

結果を表3に示す. どちらの手法も「行い」や「欠け」のように異なる使われ方をしている単語を検出できている事がわかる. 特に, BERT は「若く」「ふれ」「在り」「幼稚」のように比較的意味変化の度合いが大きい単語を, PMI-SVD_c は「おまけ」「反面」のように意味変化の度合いが小さい単語を検出している. これは, BERT が与

えられた文全体から単語ベクトルを学習し、PMI-SVD_c は対象単語から前後数単語の頻度情報から単語ベクトルを学習しているからであると考えられる。

表3 日本語データにおいて意味変化の可能性が高いと予測された上位10単語（1文字の単語は除く）。

BERT		PMI-SVD _c	
単語	説明	単語	説明
若く	匹敵→年齢	行い	振舞い→実行
触れ	言及→触る	かねて	以前→同時
行い	振舞い→実行	おまけ	追加→減額
公明	公正→組織名	無論	(副詞)
思い	感情→思考	年中	官職→1年
削除	文字の削除	キー	音楽→鍵
在り	物理→概念	欠け	物理→概念
参議	官職→議員	皆無	全然ないこと
欠け	物理→概念	馬場	人名、芝生
幼稚	幼い→幼稚園	反面	反対→一方

次に、表3より検出した単語「欠け」、および間淵らのリストに含まれる単語「了解」について、各時期のベクトルに対して余弦類似度の高い周辺単語を分析した。それぞれの結果を表4,5に示す。

表4 実際に意味が変化した「欠け」の周辺5単語（1文字の単語は除く）。

BERT		PMI-SVD _c	
戦前	戦後	戦前	戦後
マイナス	欠如	切り	有し
決まり	乏しい	切ら	欠如
構え	不足	諦め	富ん
重み	崩れ	箸	づけ
当て	破れ	つける	把握

表5 実際に意味が変化した「了解」の周辺5単語（1文字の単語は除く）。

BERT		PMI-SVD _c	
戦前	戦後	戦前	戦後
承諾	承諾	理解	承諾
承知	承知	納得	承知
納得	承認	推測	納得
理解	同意	判断	同意
断定	納得	断定	理解

まず、意味変化の自明でない単語「欠け」は表3より「物理的な欠損」から「概念的な欠損」という変化が見られた単語である。表4より、どちらの手法も戦前に「物理的欠損」、戦後に「概念的欠損」に関する周辺単語を検出できていることがわかる。

次に、間淵らによって取り上げられた意味変化の自明な単語「了解」は「理解」という意味から「承知」へと変化する単語であり、表5よりど

らの手法も各時期で対応する意味の単語を周辺語として獲得できている。しかし、BERTでは変化前の戦前に変化後の「承知」の意味を持つ単語「承諾」、「承知」が出現してしまっている。これは、BERTが大規模なデータで事前訓練される過程で現代語に強い偏りを持ってしまったのだと考えられる。

以上の結果より、本研究で提案した拡張手法は、比較的高速に学習して既存の手法と同等以上の性能を示すだけでなく、意味変化が自明な単語だけでなく自明でない単語についても分析可能であることを示した。

6. あとがき

本研究では、時期や分野の異なる文書間で意味や用例の異なる単語を検出するタスクにおいて、既存の単語ベクトル学習手法に対して2つの拡張手法を提案した。この手法は従来の手法の問題点を解消し、大規模な計算資源を必要とするBERTと比べCPUのみで軽量かつ高速に学習可能である。日本語と英語での実験より、提案した拡張手法は従来の手法よりも効果的に学習し、優れた性能を発揮する事を示した。また、網羅的な分析の結果、提案した拡張手法は意味変化が既知である単語だけでなく、自明でない単語についても検出可能である事を示した。

本研究で提案した拡張手法は言語学や社会学への導入を期待している。今後は、2つの文書間ではなくさらに複数の文書間で意味や用例の変化する単語を適切に捉える手法を模索するとともに、言語学者や社会学者が手軽に扱える効果的な手法も検討したい。

謝辞

本研究は国立国語研究所の共同研究プロジェクト「現代語の意味の変化に対する計算的・統計力学的アプローチ」、同「通時コーパスの設計と日本語史研究の新展開」およびJSPS科研費19H00531, 18K11456の研究成果の一部を報告したものである。

また、本論文の作成にあたり、丁寧に指導して下さった東京都立大学の岡照晃特任助教に感謝する。

参考文献

- [1] Kutuzov, A., Øvrelid, L., Szymanski, T., et al.: Diachronic word embeddings and semantic shifts: a survey. Proc. COLING, pp. 1384-1397, (2018).
- [2] Sagi, E., Kaufmann, S., and Clark, B.: Semantic density analysis: Comparing word meaning across time and phonetic space. Proc. GEMS, pp.104-111, (2009).
- [3] Mikolov, T., Chen, K., Corrado, G., et al.: Efficient estimation of word representations in vector

- space, Proc. *ICLR*, (2013).
- [4] Hamilton, W. L., Leskovec, J. and Jurafsky, D.: Diachronic word embeddings reveal statistical laws of semantic change, Proc. *ACL*, pp.1489-1501, (2016).
- [5] Yao, Z., Sun, Y., Ding, W., et al.: Dynamic word embeddings for evolving semantic discovery, Proc. *WSDM*, pp.673-681, (2018).
- [6] Dubossarsky, H., Hengchen, S., Tahmasebi, N., et al.: Time-out: Temporal referencing for robust modeling of lexical semantic change, Proc. *ACL*, pp.457-470, (2019).
- [7] Devlin, J., Chang, M. W., Lee, K., et al.: BERT: Pre-training of deep bidirectional transformers for language understanding, Proc. *NAACL*, pp.4171-4186, (2019).
- [8] Levy, O. and Goldberg, Y.: Neural word embedding as implicit matrix factorization, Proc. *NIPS*, pp.2177-2185, (2014).
- [9] 間淵洋子, 小木曾智信: 近現代日本語の意味変化分析のための単語データセット構築の試み, 言語処理学会第27回年次大会発表論文集, pp.1166-1170, (2021).
- [10] Tahmasebi, N., and Thomas, R.: Word Sense Change Testset, Zenodo (online), available from <<https://zenodo.org/record/495572>> (accessed 2021-10-22).
- [11] Kulkarni, V., Al-Rfou, R., Perozzi, B., et al.: Statistically significant detection of linguistic change, Proc. *WWW*, pp.625-635, (2015).