

文脈を考慮した on-the-fly 型辞書埋め込みによる含意認識

西田 光甫

西田 京介

浅野久子

富田準二

日本電信電話株式会社 NTT メディアインテリジェンス研究所

nishida.kosuke@lab.ntt.co.jp

1 はじめに

文と文の含意関係の認識は、自然言語処理における最も重要な課題の1つである。質問応答、検索、要約といった様々な課題を本質的に解くためには、文の含意関係を理解することが求められる。

文と文の含意関係を認識する課題は含意認識 (Recognizing Textual Entailment, Natural Language Inference) と呼ばれ、2000年代から多くの研究者によって取り組まれている。特に近年は、深層学習技術の発展と大規模データセット SNLI [3] や MNLI [11] の公開によって、含意認識の研究が盛んに行われている。これらのデータセットは

- 文 P: A soccer game with multiple males playing.
- 文 H: Some men are playing a sport.

のような文ペアとラベルから構成され、文 P から文 H を推論できるかどうかを判定する。この例では、正解ラベルは「推論可能」、つまり「含意」である。

一方で、こうした大規模コーパスから学習した分類器は、人間の常識に相当するような外部知識を持たないという課題がある。例えば、自然言語は、数の少ない重要な単語は出現頻度が高くそれ以外の多くの単語は出現頻度が低い、long-tail な Zipfian 分布に従って現れることが知られている [12]。そのため、大規模コーパスによる素朴な学習では、低頻度語の学習を十分に行うことができない。

含意認識における外部知識の利用に関して、Bahdanau らが単語の辞書情報を Out of Vocabulary (OOV) の単語埋め込みに用いる手法を提案した [1]。Bahdanau らの手法は、評価実験によって、事前に獲得した単語埋め込みの語彙を限定した状況下での精度向上が確認された。また、Chen らは知識ベースの単語関係性情報を用いることで、SNLI データセットでの state-of-the-art の精度を達成した [4]。これらの手法は、文 P・文 H および外部知識を入力として on-the-fly に含意認識を行うネットワークを学習する。そのため、外部知識を導入した分類器は分類精度が高いだけでなく、テストデータにしか現れない単語があった場合も、新たな外部知識を導入することで頑健に分類することが期待される。

本稿では、単語の辞書情報を含意認識に利用する新しい手法を提案する。従来手法では、辞書情報を用いて入力文中の OOV に対応するベクトルを獲得する際、同じ定義文を持つ単語はどのような入力文脈中出现しても同じ単語ベクトルに変換されてしまう。そこで提案手法では、単語の辞書情報を入力文のエンコード後に用いる。さらに、単語の定義文の単語系列をエンコードしたベクトル系列を固定長ベクトルに変換する

際に、エンコード済みの入力文とのマッチングを取ることを考える。ベクトル化に入力文とのマッチングを反映することで、より含意認識に資する特徴量を残した形で、辞書情報をエンコードすることが期待される。

本稿では評価実験によって、提案手法が辞書情報を用いてベースラインモデルの分類精度を向上することを確かめた。

2 問題設定

本稿では含意認識、辞書を以下のように定義する。

定義 1 (含意認識). 含意認識は入力データから正しいラベルを出力する課題である。

- 入力: トークン系列で表される 2 つの文のペア (文ペア) および辞書。それぞれのトークン系列を前提文 (文 P) と仮定文 (文 H) と呼ぶ。
- 出力: 文 P と文 H の関係性ラベル。関係性ラベルは含意・矛盾・中立の 3 つからなる。
- データセットの語彙 V_I : 文 P と文 H に現れるトークンの、全ての文ペアに関する和集合。
- 単語埋め込みの語彙 V_E : 事前学習された単語埋め込みデータセットが持つ語彙。
- OOV: 含意認識データセットで現れるトークンのうち、事前学習された単語埋め込みを持たないトークンの集合。 $V_I - V_E$ に相当する。

定義 2 (辞書). 辞書 D は以下の構成要素を持つ。

- 見出し語 y : 任意のトークン。
- 定義文 D^y : 見出し語の意味を説明する文章。文章はトークン系列として表現されている。見出し語はただ一つの定義文を持つ。
- 語彙 V_D : 辞書 D に含まれる見出し語の集合。

3 既存手法

3.1 深層学習による含意認識モデル

本節では一般的な含意認識モデル [3, 10, 5] の構成を示す。分類器には入力としてトークン系列のペア $\{X^s : s \in \{P, H\}\}$ が与えられる。トークン系列を $X^s := (x_1^s, \dots, x_{l_s}^s)$ と書く。 l_s は各トークンの系列長である。トークン系列やベクトル系列は同様の記法で表す。 $s = P, H$ はそれぞれ文 P, 文 H の場合を表す。

単語埋め込み層 単語埋め込み層は入力 X^s を受け取る。トークン y の単語埋め込みを $e(y) \in \mathbb{R}^{n_1}$ と書く。単語埋め込み層はベクトル系列

$$E^s := (e(x_1^s), \dots, e(x_{l_s}^s)) \in \mathbb{R}^{n_1 \times l_s}$$

を出力する。パラメータは文 P, 文 H で共通である。

文脈埋め込み層 ベクトル系列 E^s に対して RNN を作用させ、出力

$$C^s := \text{RNN}(E^s) \in \mathbb{R}^{n_2 \times l_s}$$

を得る。パラメータは文 P, 文 H に関わらず共通である。本稿では、単語埋め込み層と文脈埋め込み層を合わせてエンコーダと呼ぶ。

デコーダ デコーダはベクトル系列のペア $\{C^P, C^H\}$ を入力とし、ネットワークを作用させ分類ラベルのスコアを得る。デコーダの構成要素には、アテンション機構, RNN, 多層パーセプトロンをよく用いる。

3.2 定義文埋め込み機構

Bahdanau らの手法 [1] は単語の辞書情報を様々な自然言語処理タスクに用いる手法を提案した。さらに、含意認識・機械読解, 言語モデルのタスクで評価実験を行った。本稿では辞書情報に限定した定義文埋め込み機構として、含意認識のタスクに即して説明する。定義文埋め込み機構は文 P と文 H に対してそれぞれ作用する。入力として、トークン系列 X^s とエンコーダの文字埋め込み層の出力 E^s が与えられる。定義文埋め込み機構の出力は E'^s であり、エンコーダの文脈埋め込み層に E^s の代わりに入力する。 $E^s, E'^s \in \mathbb{R}^{n_1 \times l_s}$ である。定義文埋め込み機構は以下の層をもつ。

定義文抽出層 トークン系列 X^s のトークンのうち辞書 \mathcal{D} の語彙 V_D に含まれ、かつ OOV であるトークンの集合を V^s とする。トークン $y \in V^s$ に対して定義文 D^y が辞書 \mathcal{D} から得られる。定義文 D^y の系列長を m_y と書く。定義文抽出層はトークンの集合 V^s と定義文の集合 $\{D^y : y \in V^s\}$ を出力する。

定義文単語埋め込み層 この層は、通常の単語埋め込み層と同じパラメータを持つ。定義文抽出層の出力の要素 D^y に対してベクトル系列

$$E^y := (e(d_1^y), \dots, e(d_{m_y}^y)) \in \mathbb{R}^{n_1 \times m_y}$$

を出力する。

文脈埋め込み層 エンコーダ, デコーダとは異なる RNN を用いる。RNN の出力次元は n_1 である。単語埋め込み層の出力 E^y に対してベクトル系列

$$C^y := \text{RNN}(E^y) \in \mathbb{R}^{n_1 \times m_y}$$

を出力する。

出力層 $E_i'^s := \begin{cases} E_i^s + c_{m_y}^{x_i^s} & (x_i^s \in V^s) \\ E_i^s & \text{otherwise} \end{cases}$ を定義する。

ここで、 $c_{m_y}^{x_i^s}$ は文脈埋め込み層の RNN の出力の最終状態である。定義文埋め込み機構は $E'^s \in \mathbb{R}^{n_1 \times l_s}$ を出力する。含意認識モデルは、エンコーダの単語埋め込み層の出力 E^s の代わりに、定義文埋め込み機構の出力 E'^s をエンコーダの文脈埋め込み層に与える。

4 提案モデル

本稿では、文脈を考慮した定義文埋め込み機構を提案する。提案モデルには大きな新規性が 3 点ある。1

点目は、定義文埋め込み機構に定義文アテンション層を導入したことである。2 点目は、定義文埋め込み機構をエンコーダの後に用いることである。3 点目は、定義文埋め込みを OOV 以外のトークンにも用いることである。

入力として、トークン系列 X^s と文脈埋め込み層の出力のペア $\{C^P, C^H\}$ が与えられる。出力は C'^s であり、デコーダに C^s の代わりに入力する。 $C^s, C'^s \in \mathbb{R}^{n_2 \times l_s}$ である。文脈を考慮した定義文埋め込み機構は以下の層をもつ。

定義文抽出層, 定義文単語埋め込み層 既存の定義文埋め込み機構と同様であるが、トークン集合 V^s として、トークン系列 X^s のトークンのうち辞書 \mathcal{D} の語彙 V_D に含まれる全てのトークンの集合を用いる。

文脈埋め込み層 文脈埋め込み層はエンコーダの文脈埋め込み層とモデル, パラメータともに共通である。単語埋め込み層の出力 E^y に対してベクトル系列

$$C^y := \text{RNN}(E^y) \in \mathbb{R}^{n_2 \times m_y}$$

を出力する。

定義文アテンション層 定義文アテンション層は、定義文と文 P, 文 H のアテンションを取り、定義文の固定長ベクトル表現を得ることを目的とする。定義文アテンション層の入力は辞書の見出し語 y の定義文, 文 s , 文 \bar{s} の文脈埋め込み $C^y, C^s, C^{\bar{s}}$ である。 $s \in \{P, H\}$ に対して \bar{s} は P, H のうち s と異なるものを表す。

$C^y \in \mathbb{R}^{n_2 \times m_y}, C^s \in \mathbb{R}^{n_2 \times l_s}$ に対し、アテンション行列を

$$A^{y,s} := \frac{1}{\sqrt{n_2}} C^{s\top} C^y \in \mathbb{R}^{l_s \times m_y}$$

とする。最大プーリング, 平均プーリングによって得られるアテンションベクトルをそれぞれ

$$a^{y,s,m} := \left(\max_i A_{ij}^{y,s} \right)_{j=1, \dots, m_y} \in \mathbb{R}^{m_y}$$

$$a^{y,s,a} := \left(\frac{1}{l_s} \sum_i A_{ij}^{y,s} \right)_{j=1, \dots, m_y} \in \mathbb{R}^{m_y}$$

と定義する。アテンションベクトルは定義文 D^y の各トークンが, 文 s とどの程度関係性をもつかを表すベクトルである。定義文 D^y とトークン系列 X^s に関してアテンションを取ったベクトル

$$h^{y,s,k} := \sum_i \frac{\exp[a_i^{y,s,k}]}{\sum_j \exp[a_j^{y,s,k}]} c_i^y \in \mathbb{R}^{n_2} \quad (1)$$

を求める。 c_i^y は文脈埋め込み層の RNN が出力する各状態である。 $k = m, a$ それぞれに対して計算する。

ベクトル $\alpha_1, \dots, \alpha_n$ のベクトル β に対する Enhancement 操作 [5] を

$$\text{Enh}(\alpha_1, \dots, \alpha_n; \beta) := [\beta, \alpha_1, \alpha_1 - \beta, \alpha_1 \odot \beta, \dots, \alpha_n, \alpha_n - \beta, \alpha_n \odot \beta]$$

で定義する。ここで \odot は要素積である。

文脈埋め込み層の RNN の最終状態は $c_{m_y}^y \in \mathbb{R}^{n_2}$ で

ある。定義文 D^y の文 P,H に関する定義文埋め込み

$$z^y := \text{Enh}(h^{s,y,m}, h^{s,y,a}, h^{\bar{s},y,m}, h^{\bar{s},y,a}; c_{m_y}^y)w \quad (2)$$

を求める。 w はパラメータである。定義文アテンション層は $z^y \in \mathbb{R}^{n_2}$ を出力する。

出力層 $c_i^{s'} := \begin{cases} c_i^s + z^{x_i^s} & (x_i^s \in V^s) \\ c_i^s & \text{otherwise} \end{cases}$ を定義する。定義文埋め込み機構は $C'^s \in \mathbb{R}^{n_2 \times l_s}$ を出力する。エンコーダの文脈埋め込み層の出力 C^s の代わりに、定義文埋め込み機構の出力 C'^s をデコーダに与える。

定義文埋め込み機構の概要は Algorithm1 である。以上の議論では含意認識の課題に即して説明をした。しかし、提案手法は式 (2) の引数を任意の数に設定できるので、文 \bar{s} の数に制約がない。そのため、質問応答、機械翻訳といったテキストを入力として持つ幅広いタスクに適用することが可能である。

Algorithm 1 定義文埋め込み機構

Input: $X^s \in V_1^{l_s}, C^s \in \mathbb{R}^{n_2 \times l_s}, C^{\bar{s}} \in \mathbb{R}^{n_2 \times l_{\bar{s}}}$

Output: $C'^s \in \mathbb{R}^{n_2 \times l_s}$

- 1: $V^s, \{D^y : y \in V^s\} \leftarrow$ 定義文抽出層 (x^s)
 - 2: **for all** y in V^s **do**
 - 3: $E^y \leftarrow$ 単語埋め込み層 (d^y)
 - 4: $C^y \leftarrow$ 文脈埋め込み層 (E^y)
 - 5: $z^y \leftarrow$ 定義文アテンション層 ($C^y, C^s, C^{\bar{s}}$)
 - 6: **end for**
 - 7: $C'^s \leftarrow$ 出力層 ($C^s, \{z^y : y \in V^s\}$)
-

5 実験

本章では、提案する定義文埋め込み機構について含意認識タスクにて評価を行った結果について述べる。

5.1 比較手法

ベースラインモデルとして Enhanced Sequential Inference Model (ESIM) [5], 3.2 節にある Bahdanau らの手法の 2 つと比較した。Bahdanau らの手法と提案手法はそれぞれ ESIM にモジュールを追加した。なお、Bahdanau らの手法は OOV の単語埋め込みを補完することが目的であるため、定義文埋め込みを用いる単語の集合は $(V_1 \cap V_D) - V_E$ である。提案手法は辞書情報による分類精度向上が目的であるため、定義文埋め込みを用いる単語の集合は $V_1 \cap V_D$ である。

5.2 実験設定

データセットは Multi-Genre Natural Language Inference (MNLI) [11] を用いた。トークン化は Python の `str.split()` 関数を行った後、小文字への統一と句読点など一部記号を削除する前処理を行った。単語埋め込みに事前学習された 300 次元 GloVe 840B ベクトル [9] を用いた。OOV の単語は正規分布からランダムにサンプリングした。単語埋め込みは学習中固定した。

エンコーダとデコーダに用いる RNN には、2 層双方向 SRU [7] を用いた。双方向 SRU の出力の次元数 $n_2 = 200$ とし、活性化関数に tanh 関数を用いた。デ

コーダのアテンションを $\frac{1}{\sqrt{n_2}}$ でスケールした。ドロップアウト率は 0.45 とし、既存研究 [5] と同じ層で用いた。定義文埋め込み機構ではドロップアウトを行わなかった。

訓練は 1 つの GPU で行った。ミニバッチサイズは 32 とした。最適化は Adam [6] を用い、第 1 モメンタムを 0.9、第 2 モメンタムを 0.999 とした。初期学習率は 0.0004 とし、減衰率は 0.25 とした。訓練データから学習を行い、スケジューリングで学習率を減衰させ、開発データで評価をした。

5.3 辞書

辞書として WordNet [8] の語彙と定義文を用いた。WordNet を Python の `str.split()` 関数で行った後、小文字への統一と句読点など一部記号を削除する前処理を行った。1 つの見出し語に定義文が複数ある多義語では、WordNet に提供されている語義の出現頻度で降順に、5 つの定義文をつなげた文章を定義文とした。定義文の系列長が 2 以上の単語を語彙として採用した。また、Natural Language Toolkit [2] のストップワードを語彙から取り除いた。

5.4 評価

実験 1: OOV が多い状況下での含意認識精度比較 既存研究 [1] に即して、単語ベクトルの語彙を制限した状況での精度 (正解率) の比較をした。単語ベクトルの語彙を意図的に制限することで、OOV の単語が多い状況下での辞書情報の精度への影響を調べることができる。単語埋め込みの語彙 V_E に訓練データで出現頻度が高い 3000 語のみを使用した。残りの単語は OOV とした。辞書の語彙は WordNet の全語彙とした。

表 1 が実験 1 の結果である。MNLI は 10 のドメインから成る。5 つのドメインは訓練データ・開発データともに含まれ、matched ドメインと呼ばれる。残り 5 つのドメインは開発データのみに含まれ、mismatched ドメインと呼ばれる。実験 1 では、双方のドメインで提案手法が最も高い分類精度を示した。よって、OOV の単語が多い状況下で提案手法がベースラインモデルの精度を向上することを確認した。本実験のように、3000 語に V_E を限定した場合、基本的な語以外は OOV となるため、matched と mismatched のドメインの違いによる大きな差が見られなかった。

表 1: 単語埋め込みの語彙を限定した場合の分類精度

	Total	matched	mismatched
ESIM	67.4	66.9	68.0
Bahdanau	68.2	68.1	68.2
提案手法	69.5	69.5	69.5

実験 2: 辞書の語彙数を増やした時に含意認識精度は向上するか MNLI データセットにおいて辞書の語彙数を変化させたときの分類精度の変化を調べた。辞書の語彙は、訓練データ・開発データでの出現頻度が高い語から順に使用した。辞書の語彙数が 0 のとき、Bahdanau らの手法と提案手法は ESIM と一致する。

図 1 が実験 2 の結果である。図 1 の右端は全ての語彙を用いたときの結果である。図 1 では、辞書の語彙を増やすほど精度が向上する傾向が見られる。Bahdanau らの手法は、語彙数が小さいときは精度の向上が起きない。これは、語彙数が小さい辞書には OOV が含まれ

表 2: 各ドメインでの分類精度。順に MNLI, MNLI の matched ドメイン, MNLI の mismatched ドメイン, 全ての文ペアのうち OOV が現れる文ペア, matched ドメインに含まれる 5 つのドメイン, mismatched ドメインに含まれる 5 つのドメインにおける精度。OOV 数は開発データの各ドメインで現れる OOV の単語の異なり数。

	Total	Matched	Mismatched	OOV	Telephone	Travel	Fiction	Government	Slate	Letters	Face to Face	9/11	OUP	Verbatim
OOV 数	1999	1050	984	1999	34	367	161	225	287	159	47	158	205	430
ESIM	71.0	70.5	71.4	68.9	71.0	71.9	69.1	74.2	66.5	74.5	69.9	71.0	71.2	70.5
Bahdanau	71.5	71.3	71.6	69.1	73.0	70.9	70.4	75.1	67.3	74.7	71.5	70.5	71.0	70.1
提案手法	72.0	71.3	72.7	69.7	73.3	72.0	69.7	75.5	66.2	75.2	72.4	70.4	73.2	72.0

ないため, Bahdanau らの手法が ESIM と一致することに因る。提案手法は全ての語彙数で最も高い精度を達成している。特に, 語彙数が小さい時の Bahdanau らの手法と提案手法の差は, 頻出語における定義文埋め込みの有用性を示唆している。

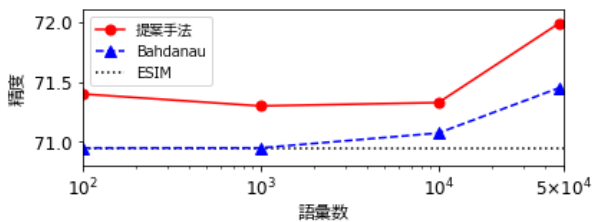


図 1: 辞書の語彙数での精度変化

実験 3: OOV 数と分類精度に相関関係があるか 実験 1, 実験 2 の結果では, 定義文埋め込みの有用性が単語が OOV であることに依存しているかどうか明らかになっていない。MNLI データセットではドメイン毎に OOV の数が違うため, OOV 数と分類精度の関係性を調べることができる。そこで実験 3 では, ドメイン毎の OOV 数と分類精度を調べた。

表 2 が実験 3 の結果である。まず, データセット全体を含む多くのドメインで提案手法が最も高い分類精度を達成した。さらに OOV が現れる文ペアにおける精度を見ると, 他の文ペアと同程度で精度が向上している。ドメイン毎の OOV 数に注目した場合も, 例えば Telephone, Face to Face のドメインは OOV 数が小さいにも関わらず大きな精度の向上が見られる。10 のドメインをサンプルとして OOV 数と分類精度の相関係数を計算すると, ESIM, Bahdanau らの手法, 提案手法ではそれぞれ -0.07, -0.38, -0.21 であり, p 値は 0.86, 0.28, 0.56 であった。OOV の数と性能に有意な関係が見られない理由について以下考察する。本実験の設定では GloVe で用意される語彙が MNLI の語彙を広くカバーするため, 非 OOV の方が OOV に比べて多く出現する。このため, 定義文埋め込みが及ぼす影響は OOV よりも非 OOV に対して大きくなる。よって, 辞書に含まれる全ての語に対して辞書埋め込みを行う提案手法は, Bahdanau らの手法に比べて精度を改善することができたと考える。

その他の考察 上記の実験では, 多くの設定で提案手法がベースラインモデルの精度を上回った。提案手法は定義文埋め込み機構内に定義文アテンション層が存在するため, 定義文内の特定の箇所注目することができる。そのため, 定義文が長文である場合や多義を表す複数文である場合に優位性を持つことが期待される。今回は WordNet の頻出 5 定義を定義文として用いたため, 語義曖昧性解消の効果があったと考えられる。

6 おわりに

本稿では以下の特徴を持つ辞書埋め込み機構を含意認識モデルに導入する手法を提案した。

- 含意認識のモデルに辞書情報を利用することができる。特に, 定義文埋め込み機構をエンコーダの後に用いたこと, 定義文埋め込み機構に定義文アテンション層を追加したことによって, 分類する文ペアの情報を利用して辞書情報をベクトル化することができる。
- 含意認識に限らず, 質問応答, 機械翻訳といったテキストが入力として与えられる自然言語処理タスクに幅広く適用できる。

また, 評価実験によって以下の知見が得られた。

- 含意認識のタスクにおいて, 辞書情報の利用に関して state-of-the-art のアプローチである Bahdanau らの手法を, 提案手法が分類精度で上回った。
- 含意認識において辞書情報が有用である。さらに, 辞書情報は高頻度語, 低頻度語どちらにおいても有用である。

今後の課題としては, 評価実験の拡充が挙げられる。含意認識で state-of-the-art を獲得しているモデルの他, 含意認識以外の自然言語処理のタスクのモデルに辞書埋め込み機構を適用した場合の評価が必要である。

参考文献

- [1] D. Bahdanau, T. Bosc, S. Jastrzebski, E. Grefenstette, P. Vincent, and Y. Bengio. Learning to compute word embeddings on the fly. *CoRR*, abs/1706.00286, 2017.
- [2] S. Bird, E. Klein, and E. Loper. *Natural Language Processing with Python*. O'Reilly Media, Inc., 2009.
- [3] S. R. Bowman, G. Angeli, C. Potts, and C. D. Manning. A large annotated corpus for learning natural language inference. In *EMNLP*, 2015.
- [4] Q. Chen, X. Zhu, Z. Ling, D. Inkpen, and S. Wei. Natural language inference with external knowledge. *CoRR*, abs/1711.04289, 2017.
- [5] Q. Chen, X. Zhu, Z. Ling, S. Wei, H. Jiang, and D. Inkpen. Enhanced LSTM for natural language inference. In *ACL*, 2017.
- [6] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2014.
- [7] T. Lei and Y. Zhang. Training RNNs as fast as CNNs. *CoRR*, abs/1709.02755, 2017.
- [8] G. A. Miller. WordNet: A lexical database for english. *Commun. ACM*, 38(11):39–41, 1995.
- [9] J. Pennington, R. Socher, and C. D. Manning. GloVe: Global vectors for word representation. In *EMNLP*, pp. 1532–1543, 2014.
- [10] T. Rocktäschel, E. Grefenstette, K. M. Hermann, T. Kočiský, and P. Blunsom. Reasoning about entailment with neural attention. *CoRR*, abs/1509.06664, 2015.
- [11] A. Williams, N. Nangia, and S. R. Bowman. A broad-coverage challenge corpus for sentence understanding through inference. *CoRR*, abs/1704.05426, 2017.
- [12] G. K. Zipf. Human behavior and the principle of least effort: An introduction to human ecology, 1949.