

人間と BERT による文間接続関係の同定の比較と性能検証

相澤祐一, 鋤田大日, 笠間俊夫

みずほリサーチ&テクノロジーズ株式会社

{yuichi.aizawa, dainichi.sukita, toshio.kasama}@mizuho-rt.co.jp

概要

本研究は、計算機による文間接続関係の同定結果と人間による同定結果を比較することにより、計算機による誤判定の傾向および性能改善の課題を考察する。計算機実験では、大量のテキストデータを学習した BERT モデルにより、文間接続関係を同定させた。その結果、接続関係を順接に誤分類する傾向が見られた。一方、校正・校閲の経験者等の被験者に、計算機実験で用いたデータを使って接続関係を同定させた結果、接続関係を順接だけでなく、逆接にも誤分類する傾向が見られた。両実験の結果の比較から、被験者に比肩する精度の同定が可能な BERT モデルが構築できたことが分かった。また、接続関係間の違いをより明確に捕捉させる必要があることが示唆された。

1 はじめに

文間の関係推定は対話、要約、文章生成など様々なタスクで必要とされる処理であり、自然言語処理の重要な問題となっている。特に、接続関係は、論理的関係の把握及びその構成において、大きな役割を果たしていると考えられている[1]。

計算機に文間接続関係を同定させる試みとして、山本らは、単語と構文的要素に基づく大量のテキストデータを使用した統計的手法を提案した[2]。さらなるシステムの向上のため、斎藤らは、被験者に対し、テキスト全体を提示した条件と二文だけの情報を提示した条件の2通りにおける接続関係の推定結果から接続関係間の近さを定量化し、SVM (Support Vector Machine) を使った二値分類によって得られる近さと同じ傾向になることを示した。このように、計算機による分析と人間の認知を比較することは、計算機モデルをどのように向上させるかといった糸口を得るアプローチになる[3]。

また、同じ SVM の手法を用いて、続く二文に加えて前後各二文から品詞、文間類似度、係り受けの構文パターン等の素性を取り出し、多値分類を行ったものがある[4]。近年では、再帰的ニューラルネットワークを用いて文の概念ベクトルを計算、文間接続関係を RAE (Recursive AutoEncoder) を用いて同定する試みも行われている[5]。さらに、趙らは、BERT の Masked Language Model を用いて、少量データではあるが、マスクされた箇所に入る最も高い確率の接続詞を推定するという手法を提案している[6]。

このように、機械学習アプローチによる接続関係の同定に関する取り組みが進められてきたものの、深層学習等による SOTA が様々なタスクで達成されてからは、斎藤らが行ったような認知との比較が行われた研究は少ない。そこで本研究では、大量のテキストデータを学習した BERT モデルによる文間接続関係の同定結果と人間による同定結果を比較し、計算機による同定の傾向およびシステム改善の課題について考察する。

2 文間接続関係

本研究では、先行研究[1]と同様に、接続関係の同定にあたって、接続関係の種類を同じくする接続表現のいずれを選んでよいと仮定する。この仮定に基づき、接続関係の同定を、“接続関係種類の分類問題”と置き換えることとする。

接続関係の種類を定義するため、石黒[7]を参考に、形態素解析器 MeCab に辞書として登録されている接続詞を分類し、同文献で提案されている接続表現をそれぞれの種類に加えて表 1 を得た。なお、複数の接続関係の種類に見られる表現は、その表現が表れる文間接続関係の同定が困難になることから、取り扱わないこととした。

表 1 接続関係の種類

種類	接続表現の例
順接	そして、そこで、このため、そのため
逆接	しかし、ただ、だが、でも
対比	これに対して、一方、他方
並列	また、そのうえで、まして
列挙	まず、まずは、一つは
例示	たとえば、とりわけ、事実、実際
補足	なぜなら、ちなみに、ただし
換言	すなわち、ようするに、それよりも
結論	このように、結局、とにかく
転換	さて、ところで、そういえば

3 計算機実験

本節では、BERT モデルによる文間接続関係の同定実験について述べる。

3.1 データ

本研究では、文法誤りが混在したテキストデータを扱わないようにするため、校正・校閲が為された新聞記事データを用いる。具体的には、毎日新聞6年分の記事データ[8]から、表1で定義した接続表現が記載されている、51,790件の続く二文を取り出し、接続表現に該当する部分をマスクした。なお、モデル学習時のバイアスを避けるため、データ件数は、接続関係間で偏りがないようにしている。それぞれのペアには、正解となる接続関係をラベリングした。

3.2 モデル

モデルは、東北大学から提供されている BERT Pre-trained Model (<https://github.com/cl-tohoku/bert-japanese/releases/tag/v2.0>) を用いて、図1に示す構成とした。

続く二文を[SEP]トークンで結合したうえで、Pre-trained Model にインプットし、全結合層で処理して得られる出力値と、One-hot encoding した正解データを Binary Cross Entropy Logitloss で比較することで、文間接続関係を学習させた。なお、Optimizer を AdamW, 学習率を 2×10^{-5} とした。

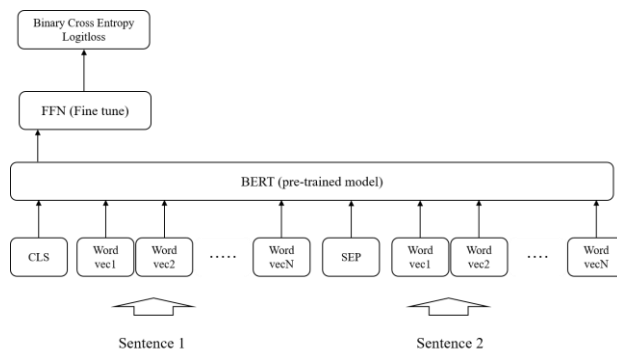


図 1 BERT モデル構成

3.3 実験結果

3.1 で用意したデータの約 75% を用いて、5 分割交差検証を実施した。テストデータとして用意した残りの約 25% の 12,948 件のデータに対して、BERT モデルが接続関係を同定した結果を混同行列として表 A-1 に示す。全体の正解率は、0.52 となった。

表 A-1 において、並列が正解となる接続関係を対比と誤分類する傾向、補足や対比が正解となる接続関係を逆接と誤分類する傾向が見られた。

続いて、接続関係ごとの Recall, Precision, および F 値を算定し表 2 を得た。表 2 から、列挙について、Recall, Precision, および F 値が全接続関係の中で高いことが分かった。また、接続関係を順接に誤分類する傾向があることから、Precision と F 値は順接が最も低くなっていることが確認された。

4 被験者実験

本節では、被験者による接続関係の同定実験について述べる。

4.1 実験設定

実験には、校正・校閲の経験者や国語教員免許の保有者など 22 名が参加した。

被験者には、BERT モデルのパフォーマンスの評価に使用したテストデータのうち、接続関係 1 種類あたり 40 件、合計 400 件を取り出し、接続表現を () でマスクした二文として提示、() に入るのに相応しい接続関係の候補を回答させた。なお、回答にあたっては、上位 3 つまでの接続関係の候補を選択できるようにした。

表 2 BERT モデルの分類結果 (Recall, Precision, F 値)

	順接	逆接	対比	並列	列挙	例示	補足	換言	結論	転換
Recall	0.51	0.47	0.53	0.42	0.61	0.53	0.44	0.51	0.56	0.59
Precision	0.39	0.44	0.52	0.64	0.69	0.47	0.59	0.54	0.51	0.52
F 値	0.44	0.45	0.53	0.51	0.65	0.50	0.53	0.50	0.53	0.56

表 3 被験者の回答結果 (Recall, Precision, F 値)

	順接	逆接	対比	並列	列挙	例示	補足	換言	結論	転換
Recall	0.56	0.71	0.52	0.48	0.45	0.43	0.32	0.42	0.41	0.33
Precision	0.31	0.45	0.53	0.47	0.65	0.51	0.46	0.66	0.43	0.61
F 値	0.39	0.55	0.52	0.46	0.51	0.45	0.37	0.50	0.41	0.41

4.2 実験結果

表 A-2 は、被験者が第 1 候補として回答した接続関係を混同行列として表したものである。表 A-2 から、計算機実験と同様、接続関係を順接であると誤回答する傾向が見られた。また、補足や対比が正解となる接続関係を逆接と誤回答する傾向も確認された。一方で、BERT モデルで見られる並列を対比として誤分類する傾向は被験者実験では見られなかった。

表 3 は、被験者それぞれの回答結果の Recall, Precision, および F 値の平均値を算出したものである。計算機実験と同様、順接の Precision が低いという共通点が見られるが、表 2 と比較して逆接の Recall が大きく異なることが分かった。図 2 に、同じ 400 問に対する BERT の分類結果と、被験者による回答結果の分布を示す。図 2 から、被験者のほうが、逆接をより多く選択していることが明らかとなった。

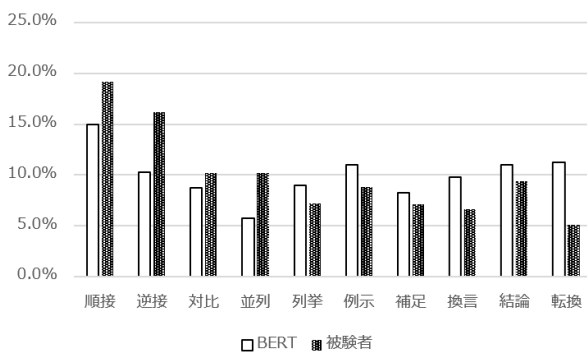


図 2 モデルと被験者による予測・回答の分布

F 値については、逆接を除いて、BERT モデルが被験者よりも高く、また、被験者回答の平均正解率は 0.46 となった。このことから、BERT モデルは、被験者に比肩する精度の同定が可能であることが示唆された。

以下に、結果をまとめる。

- BERT モデルによる接続関係の同定性能は、被験者に比肩することが分かった。
- BERT モデル、被験者ともに、接続関係を順接として誤分類してしまう傾向、および、補足や対比を逆接に誤分類する傾向が見られた。
- 被験者の回答では、接続関係を逆接として誤回答する傾向がより強く見られた。BERT モデルで見られる並列を対比と誤分類する傾向は被験者実験では見られなかった。

5 考察

本節では、さらなる性能改善に向けた考察を行う。

まず、計算機実験及び被験者実験において、接続関係を順接と誤分類する傾向や、補足や対比の接続関係を逆接と誤分類する傾向がある点については、順接の接続関係が構造上、他の接続関係を包含あるいは代替する可能性、補足・対比の接続関係が逆接と近い可能性が考えられる。これを確認するためには、それぞれの接続関係の組み合わせごとに、Sentence-BERT を使った深層距離学習による二値分類を行い、先行研究[3]のように、接続関係間の距離の近さを可視化し、被験者実験の結果とも比較することが有効であろう。

続いて、接続関係を逆接として誤分類する傾向が被験者実験で見られた点について述べる。ここで、

被験者回答結果に対して、選択候補1位が正解した場合は1点、2位は1/2点、3位は1/3点と点数を付け、400問を点数の高い順に並び替えた。上位50問を取り出すと、50問中18問(36%)を逆接が占めており、かつこの18問については、問題にはよるが、19名~22名が第1候補に逆接と回答できていることが分かった。全体的に被験者が接続関係を逆接として誤回答する傾向を踏まえると、そもそも人間は逆接の関係を期待して文章を認知している可能性がある。なぜならば、接続詞の頻度の調査結果において、多くのジャンルの文章において、逆接が頻出することが分かっており[9]、一定の認知的バイアスが入っていることが考えられるからである。

しかしながら、BERTモデルは、この18問のうち11問しか正しく同定できていない。たとえば、22名中20名の被験者が逆接であると正答している以下の文章例において、BERTモデルは転換であると誤分類しており、前後の文脈や意味を捕捉しきれていない可能性があることが分かった。これを確認するためには、逆接の接続関係を持つ二文と、転換の二文とで、BERTモデルのAttentionの違いを観察する必要があると考えられる。

<文章例> (正解：逆接)

岸田文雄外相は「予算編成作業があり、十分な受け入れ態勢がとれなかった」と述べた。()、訪日は10月時点で決定しており、延期要請の理由としては説得力を欠き、その対応は不可解だ。

さらに、BERTモデルが、並列の接続関係を対比と誤分類する傾向について考察する。文献[7]では、並列と対比は、共に整理の機能を有するとされるが、BERTモデルは対比を並列と誤分類する傾向よりも、並列を対比と誤分類する傾向が強いことが表A-1から読み取れる。ここで、被験者実験の点数上位5名の並列の接続関係に対する回答と、同じく下位5名の回答を比較したところ、下位5名のほうが対比を選択する傾向があることが示唆された。このことから、BERTモデルにとっても、対比構造を把握する方が、並列関係を見出すよりも平易である可能性がある。これを確認するためには、対比の接続関係を持つ続く二文と並列の二文とで、BERTモデルのAttentionの違いを観察すること、および、システムが対比と並列をより区別して捕捉できるようにする必要があると考えられる。具体的には、先述のよう

に、それぞれの接続関係の組み合わせごとに、Sentence-BERTを使った深層距離学習による二値分類を行い、違いを捕捉させることが有効であろう。

6 まとめ

本研究では、計算機による文間接続関係の同定結果と校正・校閲の経験者等被験者による同定結果を比較し、被験者に比肩する同定性能のモデルを構築できたことが分かった。また、計算機にも被験者にも共通して見られる誤分類の傾向があることが分かった。このような知見は、人間が作成した文章の自動評価にも一定の貢献があると考えられる。今後は、システムに接続関係間の違いを捕捉させる取り組みや文章構造に対するモデルのAttentionの観察を通じて、より性能のあるモデルを構築していく予定である。

参考文献

1. 山本和英, 齋藤真実. 用例利用型による文間接続関係の同定. 自然言語処理, Vol.15, No.3, 2008.
2. 齋藤真実, 山本和英, 関根聡. 大規模テキストを用いた2文接続関係の同定. 言語処理学会第12回年次大会, pp.969-972, 2006.
3. 齋藤真実, 山本和英, 関根聡. 文間接続関係の自動同定のための人間による同定分析, NL174-12, pp.65-70, 2006.
4. 若山裕介, 内海彰. SVMを用いた接続関係の同定, 人工知能学会全国大会論文集, 2012.
5. 大塚淳史, 平野徹, 宮崎千明, 東中竜一郎, 牧野俊朗, 松尾義博. Recursive AutoEncoderを用いた文間の接続関係推定. 人工知能学会全国大会論文集, 2015.
6. 趙一, 曹銳, 白静, 馬ブン, 新納浩幸. BERTのMasked Language Modelを用いた二文間の接続関係の推定. 言語資源活用ワークショップ発表論文集, Vol.5, pp.181-188, 2020.
7. 石黒圭. 「接続詞」の技術. 実務教育出版, 2016.
8. 毎日新聞記事データ集 2017年版, 2018年版, 2019年版, 2020年版, 2021年版, 2022年版.
9. 石黒圭, 阿保きみ枝, 佐川祥予, 中村紗弥子, 劉洋. 接続表現のジャンル別出現頻度について, 一橋大学留学生センター紀要第12号, 2009.

A 付録

表 A-1 BERT モデルの分類結果 (混同行列)

		予測										
		順接	逆接	対比	並列	列挙	例示	補足	換言	結論	転換	総計
正解	順接	661	115	37	29	47	76	24	88	139	79	1,295
	逆接	152	610	120	14	7	86	80	60	93	73	1,295
	対比	111	156	690	88	14	77	75	22	29	33	1,295
	並列	108	33	254	544	34	123	55	28	68	47	1,294
	列挙	111	24	8	39	789	88	13	66	96	61	1,295
	例示	122	60	68	45	75	685	41	54	59	86	1,295
	補足	88	184	97	44	29	79	569	56	35	114	1,295
	換言	117	86	25	17	56	73	37	662	95	127	1,295
	結論	132	76	10	18	73	71	18	95	721	81	1,295
	転換	85	52	15	17	26	98	55	101	77	768	1,294
	総計	1,687	1,396	1,324	855	1,150	1,456	967	1,232	1,412	1,469	12,948

表 A-2 被験者の回答結果 (混同行列)

		回答										
		順接	逆接	対比	並列	列挙	例示	補足	換言	結論	転換	総計
正解	順接	497	65	58	39	18	26	35	30	97	15	880
	逆接	42	626	62	17	16	12	48	6	41	10	880
	対比	95	164	458	95	4	19	18	6	10	11	880
	並列	191	35	62	418	32	51	43	11	21	16	880
	列挙	174	31	17	50	397	75	27	25	70	14	880
	例示	163	19	33	93	51	376	57	17	52	19	880
	補足	96	255	86	47	7	39	285	17	22	26	880
	換言	86	43	22	54	37	70	55	373	102	38	880
	結論	228	67	24	41	44	28	6	71	357	14	880
	転換	118	118	72	42	26	81	50	30	56	287	880
	総計	1,690	1,423	894	896	632	777	624	586	828	450	8,800