

日本語日常対話コーパスへの基礎解析アノテーション

赤間 怜奈^{1,2} 浅原 正幸³ 若狭 絢¹ 大村 舞³ 鈴木 潤^{1,2}

¹ 東北大学 ² 理化学研究所 ³ 国立国語研究所
{akama, aya.wakasa.c3, jun.suzuki}@tohoku.ac.jp,
{masayu-a, mai-om}@ninjal.ac.jp

概要

規範的な日常対話を収録した言語資源「日本語日常対話コーパス」に対し、基礎解析情報のアノテーションをおこなっている。具体的には、形態素解析および構文解析を施し、形態論情報として UniDic に基づく短単位・長単位形態論情報を、構文情報として文節境界と文節係り受け情報を付与した。さらに、これらの情報を用いて、多言語間で共通化された依存構造アノテーション仕様 Universal Dependencies 準拠の言語資源を構築した。本稿では、基礎解析アノテーションの手順ならびに進捗状況を報告するとともに、アノテーション済みデータの活用事例として、依存構造解析器を構築し既存解析器との比較によりその特性と有用性を紹介する。

1 はじめに

自然言語処理の研究・開発は、近年は分野として言語横断的な取り組みを推奨する動きはあるものの [1]、依然として英語を対象とした議論が主導的である。日本語を対象とした自然言語処理、たとえば、英語をはじめ他言語のうえで確立された最先端の技術や知見の日本語への適応可能性を検証すること、あるいは、日本語話者の文化や思想が反映された日本語特有の表現を満実に処理する方法論を確立することなどを議論していくためには、日本語言語資源の整備が必要不可欠である。

日本語日常対話コーパス (JDailyDialogue; JDD) は、規範的な日本語表現で構成される高品質なマルチターン日常対話コーパスである [2]。日本語の対話コーパスは音声対話の書き起こしを収録したものが多くなか、日本語日常対話コーパスは、Business Scene Dialogue [3], JPersonaChat, JEmpatheticDialogues [4] と同様、人が創作したテキスト対話を収録したものである。これらのコーパスは、人間同士の自然対話で観察されるリアルな言語現象 (フィラーや相

槌の多用、発話途中の話者交代など) を含まないが、近年の対話研究のデファクトスタンダードである英語対話データ DailyDialog [5], Persona-Chat [6], EmpatheticDialogues [7] の日本語版と見做すことができ、日英言語横断的研究の土台を成している。

日本語日常対話コーパスをより標準的で利便性の高い言語資源として整備することを目的として、コーパス内の対話に対して形態論情報や統語情報などの基礎解析情報のアノテーションを進めている。さらに、これらの情報に基づき、日本語日常対話コーパスを多言語間で共通化された依存構造アノテーション仕様である Universal Dependencies (UD) [8] に準拠した言語資源として再構築している。基礎解析情報が付与されている日本語対話コーパスには日本語話し言葉コーパス [9, 10] や日本語日常会話コーパス [11] があるが、これらはいずれも自然音声発話の書き起こしデータであり、日本語日常対話コーパスのような高品質テキスト対話に解析情報が付与されている言語資源は希少である。日本語の UD 言語資源には UD Japanese-GSD など [12, 13, 14] があるが、著作権の都合で表層形と合わせて利用可能な言語資源は限定的である。本稿で紹介する UD Japanese-JDD は、データの包括性に加えて規模と品質の観点からも、日本語の UD 言語資源としての高い有用性が期待できる。

本稿では、日本語日常対話コーパスに対して実施している基礎解析アノテーションの内容および手順を説明し、基礎解析結果の概要を報告する。また、基礎解析結果に基づいて構築した UD Japanese-JDD の活用の一例として依存構造解析器を作成し、既存解析器との比較によりその特性を明らかにする。

2 アノテーション内容と作業手順

日本語日常対話コーパスに形態論情報、文節係り受け情報を付与し、単語依存構造に変換した (図 1)。各工程の詳細と作業手順を以下で説明する。

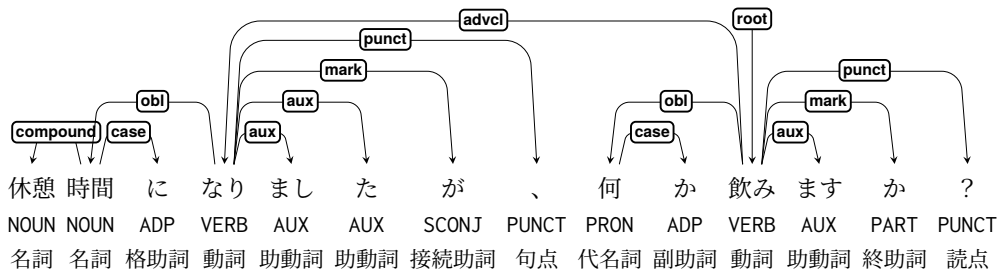


図1 UD Japanese-JDD の例. 短単位アノテーション. 発話は「日常生活」トピック内の対話より抜粋.

2.1 形態論情報の付与

形態論情報として、UniDic 品詞体系に基づく国語研短単位・国語研長単位形態論情報 [15] を付与した. まず、国語研短単位は短単位自動解析用辞書 UniDic と形態素解析器 MeCab を、国語研長単位は中長単位解析器 Comainu¹⁾ を用いて自動解析をおこなった. 専門の作業者が、形態論情報アノテーションシステム大納言 [16] を用いて自動解析結果を修正し、適切なラベルが付与されていることを確認した. 国語研長単位に基づいて、文節境界を付与した.

2.2 文節係り受け情報の付与

文節係り受け情報は、BCCWJ-DepPara [17] の基準に準じる. 先に係り受け解析器 CaboCha²⁾ を用いた自動解析をおこない、その後、専門の作業者がコーパス管理システム ChaKi.NET [18] を用いて解析結果を修正した. 一般に、文節係り受け解析は 1 文をひとつの単位としておこなわれる. しかし、日本語日常対話コーパスに含まれる対話は、1 発話がひとつの単位となっており、1 発話には 2 つ以上の文が含まれている場合も多い. このようなデータを適切に解析するために、文節境界に加えて、句点などを手がかりに改めて文境界を付与した.

2.3 誤記の修正・表現の正規化

前述の作業を実施する過程で、元の対話データに含まれる誤記等が発見された. たとえば、UniDic 内の適切な形態素を割り当てられない箇所は、誤字や一般的でない表現を含む可能性が高い. 発見された誤記等は、日本語日常対話コーパスの開発時と同様、既定の正書法 [19] に従って作業者が修正あるいは正規化した. また、適切な文節係り受けを付与できない箇所は、文法規則から逸脱した表現である可能性が高い. 発見された非文法的な表現は、可能な

1) <http://comainu.org/>

2) <https://taku910.github.io/cabocho/>

限り文意を保存しつつ作業者が柔軟に修正した.

2.4 単語依存構造への変換

形態論情報および文節係り受け情報に基づいて、Universal Dependencies (UD) に準拠した単語依存構造を構築した. UD は、他言語間で共通化されたアノテーション規則のもとでツリーバンクを作成する国際プロジェクトであり、これに則して記述された資源は言語横断的な解析も容易となる. 形態論情報・文節係り受け情報に UD の変換規則 [13] を適用し、国語研短単位 (Short Unit Word; SUW) および長単位 (Long Unit Word; LUW) に基づく依存構造に変換した. また、同規則により品詞体系を UniDic から UD 規定の Universal Part-of-Speech (UPOS) に変換し、さらに、UPOS に基づいて UD 規定の係り受けラベル DEPREL を割り当てた.

3 UD Japanese-JDD の特性

上述の手順で日本語日常対話コーパスを UD に変換した言語資源を UD Japanese-JDD と呼ぶ. 2024 年 1 月時点で UD Japanese-JDD に収録されている対話、すなわちアノテーション作業が完了している対話は 1,545 で、これは日本語日常対話コーパス全体 (5,261 対話) の約 29.4% に相当する. UD Japanese-JDD の特性を、以下に示す既存の日本語 UD 言語資源との比較を通して述べる.

- UD Japanese-CEJC [14]: 日本語日常会話コーパス [11] のコアデータに基づく. 自然音声対話の転記で、完全でない文を多く含む. 句読点が含まれない. 表層形は契約者のみ利用可.
- UD Japanese-GSD [12]: Wikipedia テキストに基づく. 文単位のみ、文書レベルの単位は定義されない. 表層形はオープンライセンスで公開.
- UD Japanese-BCCWJ [13]: 現代日本語書き言葉均衡コーパス [20] に基づく. 書籍・雑誌・新聞・白書・ブログなど多様な媒体の書き言葉を収録. 表層形は契約者のみ利用可.

表 1 日本語 UD 言語資源の統計情報。「平均長」は 1 文 (発話) あたりの単語数。*JDD は 2024 年 1 月時点のもの。

コーパス	文書 (対話) 数	文 (発話) 数	文節数	短単位 (SUW)		長単位 (LUW)	
				単語数	平均長	単語数	平均長
UD Japanese-JDD*	1,545	16,691	61,468	182,463	10.9	159,680	9.6
UD Japanese-CEJC	52	59,319	136,071	256,885	4.3	231,774	3.9
UD Japanese-GSD	N/A	8,100	65,966	193,654	23.9	150,243	18.5
UD Japanese-BCCWJ	1,980	57,109	425,751	1,253,903	21.9	995,632	17.4

表 2 品詞ラベル UPOS の分布。

	JDD	CEJC	GSD	BCCWJ
ADJ	3.40%	3.69%	1.98%	2.14%
ADP	16.87%	13.58%	21.62%	20.03%
ADV	3.03%	6.74%	1.22%	1.51%
AUX	17.30%	13.21%	10.93%	9.74%
CCONJ	0.35%	1.71%	0.42%	0.41%
DET	0.50%	0.56%	0.51%	0.48%
INTJ	0.65%	10.74%	0.01%	0.07%
NOUN	18.08%	14.86%	30.05%	29.24%
NUM	0.69%	1.67%	2.67%	3.11%
PART	4.35%	8.49%	0.65%	1.18%
PRON	1.96%	3.76%	0.57%	0.90%
PROPN	0.93%	1.39%	3.69%	2.87%
PUNCT	13.48%	0.00%	9.93%	11.69%
SCONJ	6.26%	6.68%	4.13%	4.49%
SYM	0.03%	0.00%	0.67%	1.53%
VERB	12.14%	9.86%	10.96%	10.57%
X	0.00%	3.05%	0.00%	0.03%

表 3 係り受けラベル DEPREL の分布。

	JDD	CEJC	GSD	BCCWJ
<i>acl</i>	2.05%	2.19%	3.61%	3.62%
<i>advcl</i>	4.06%	4.03%	3.72%	3.85%
<i>advmod</i>	2.29%	4.84%	1.18%	1.43%
<i>amod</i>	0.10%	0.10%	0.23%	0.25%
<i>appos</i>	0.00%	0.00%	0.00%	0.00%
<i>aux</i>	13.18%	9.12%	8.90%	7.56%
<i>case</i>	16.15%	12.72%	21.33%	19.65%
<i>cc</i>	0.35%	1.66%	0.42%	0.41%
<i>ccomp</i>	0.38%	0.34%	0.20%	0.22%
<i>compound</i>	4.87%	3.97%	14.19%	14.67%
<i>cop</i>	1.82%	1.98%	1.26%	1.20%
<i>csbj</i>	0.09%	0.09%	0.08%	0.11%
<i>csbj:outer</i>	0.00%	0.01%	0.00%	0.00%
<i>dep</i>	0.00%	0.05%	0.04%	0.99%
<i>det</i>	0.50%	0.55%	0.51%	0.48%
<i>discourse</i>	0.27%	2.85%	0.01%	0.03%
<i>dislocated</i>	0.00%	0.00%	0.00%	0.00%
<i>fixed</i>	5.65%	4.15%	4.45%	4.26%
<i>mark</i>	10.50%	14.20%	4.06%	5.04%
<i>nmod</i>	2.70%	3.03%	6.70%	6.92%
<i>nsubj</i>	4.38%	2.42%	4.02%	3.69%
<i>nsubj:outer</i>	0.42%	0.17%	0.23%	0.18%
<i>nummod</i>	0.48%	0.98%	1.45%	1.16%
<i>obj</i>	2.15%	0.49%	2.74%	2.62%
<i>obl</i>	4.98%	5.76%	6.55%	5.41%
<i>punct</i>	13.48%	0.00%	9.93%	11.69%
<i>reparandum</i>	0.00%	1.21%	0.00%	0.00%
<i>root</i>	9.15%	23.09%	4.18%	4.55%

3.1 基礎統計

表 1 に各日本語 UD 言語資源の統計情報を示す。すべての言語資源で単位が共通している文節数と単語数に着目すると、JDD は、現時点で既に GSD と同程度の規模である。コーパス内のすべての対話にアノテーションが完了すると、現在の約 3 倍強の規模になることが予想される。JDD の 1 文 (発話) あたりの平均単語長は、話し言葉 (CEJC) と書き言葉 (GSD, BCCWJ) の中間の大きさであった。

3.2 ラベル分布

品詞ラベル 表 2 に各 UD 言語資源における短単位 SUW に付与された品詞ラベル UPOS の分布を示す。対話データである JDD と CEJC は、GSD や BCCWJ と比べて複合名詞 NOUN が少ない。同じ対話であっても、JDD は CEJC と比べて、副詞 ADV, 接続詞 CCONJ, 感動詞 INTJ の値が小さい。これは JDD がフィラーや言い淀みをあまり含まないことに起因する。JDD では規範的な表現が用いられているため、対話であっても格助詞の省略が少なく、格助詞・係助詞 ADP が比較的大きい値である。また、JDD は、

助動詞 AUX の値がとくに大きい。書き言葉である GSD や BCCWJ と比較しても大きな値である。このことから、JDD では丁寧語の「です」「ます」や複合動詞が多用される傾向にあることが見て取れる。

係り受けラベル 表 3 に各 UD 言語資源における短単位 SUW に付与された係り受けラベル DEPREL の分布を示す。JDD は、品詞ラベルと同様、助動詞 *aux* の割合がとくに大きい。形容詞節 *acl*, 副詞節 *advcl* の分布は、書き言葉である GSD や BCCWJ よりも話し言葉である CEJC に近い。つまり、JDD には節が少ない文 (単文) が比較的多く含まれていることが示唆される。一方で、名詞主語 *nsubj* と目的語 *obj* の割合は、CEJC よりも明らかに多く、GSD や BCCWJ に近い。このことから、JDD は主語や目

的語を省略しない表現が多く含むことがわかる。

4 活用事例：依存構造解析器

UD 言語資源の活用事例のひとつとして、依存構造解析器の構築が挙げられる。本章では、UD Japanese-JDD を用いて依存構造解析器を構築し、他の解析器との比較を通してその特性を分析する。

4.1 実験設定

解析器の作成 既存研究 [14] に倣い、自然言語処理ライブラリ spaCy³⁾ を用いて解析器を作成した。解析モデルとして spacy-transformers を採用した。これは、Transformers ベースの事前学習モデル⁴⁾と解析コンポーネント⁵⁾との間で損失勾配を共有しながら同時学習をおこなうモデルである。UD Japanese-JDD を訓練データ⁶⁾としてモデルを学習し、解析器を作成した。比較のために、UD Japanese-GSD および UD Japanese-CEJC を訓練データ⁷⁾とした解析器も作成した。

評価方法 作成した依存構造解析器を用いて対象のテキストデータを解析し、その精度を F 値として算出した。評価観点には次の 6 つとした：単語分割精度 (Tokens)、UD 品詞ラベルの予測精度 (UPOS)、UniDic 形態論情報の予測精度 (XPOS)、語彙素の予測精度 (Lemmas)、依存関係の予測精度 (Unlabeled Attachment Score; UAS)、依存関係に加えて UD 係り受けラベルの予測精度 (Labeled Attachment Score; LAS)。いずれも短単位に基づく解析。解析対象には、JDD、CEJC、GSD のテストデータを用いた。

4.2 結果と考察

表 4 に、UD 言語資源を活用して作成した依存構造解析器の各設定における解析精度（一部抜粋）を示す。⁸⁾ここでは、UD を構成するラベル（品詞 UPOS、係り受け DEPREL）に関連する 3 つの評価観点（UPOS、UAS、LAS）から依存構造解析タスクにおける UD Japanese-JDD の特性を議論する。

訓練データとして JDD を用いて作成した解析器は、いずれのテストデータでも他の解析器を下回るあるいは同程度の解析精度であった。特に、CEJC

表 4 UD 言語資源に基づく依存構造解析器の解析精度。

訓練データ	解析対象	UPOS	UAS	LAS
CEJC	CEJC	93.35	88.13	86.17
GSD	CEJC	79.77	80.69	74.68
JDD	CEJC	70.34	71.34	62.82
CEJC	GSD	84.27	79.54	70.41
GSD	GSD	96.99	91.25	90.29
JDD	GSD	88.59	78.05	71.75
CEJC	JDD	83.65	91.61	78.07
GSD	JDD	97.27	93.81	92.95
JDD	JDD	96.03	92.34	89.80
CEJC + GSD	JDD	97.74	94.18	93.24
CEJC + JDD	JDD	95.13	90.57	87.76
GSD + JDD	JDD	95.39	91.25	88.86
CEJC + GSD + JDD	JDD	94.83	90.50	87.80

を解析したときの LAS は 62.82 と低い値であった。これは、JDD は規範的な対話であることに対し、自然発話である CEJC は不規則な係り受け関係を多く含むことが要因のひとつと考えられる。

JDD を解析対象とした場合、CEJC や GSD を対象とした場合よりも解析器の精度が高くなる傾向にあった。JDD は、その性質上、形態や統語構造の面でも規範的なものが多く例外的な現象をあまり含まないため、多くの解析器にとって解析し易いデータであると考えられる。JDD の解析において最も高い精度を示したのは、CEJC と GSD の 2 種類を組み合わせる訓練データとして用いた解析器であった。このことから、JDD は、形態や統語構造の面で話し言葉と書き言葉の両方を併せ持つ、あるいはそれらの中間的な性質を有していることが示唆された。

5 おわりに

規範的な日常対話を収録した言語資源「日本語日常対話コーパス」に対して基礎解析情報のアノテーションをおこなっている。本稿では、アノテーションの内容と作業手順を説明し、単語依存構造への変換までのすべての基礎解析が完了した対話データで構築した日本語 UD 言語資源「UD Japanese-JDD」を紹介した。既存資源との比較分析により、UD Japanese-JDD は、表層形が利用可能な日本語 UD 言語資源としては現時点で既に最大級の規模であること、形態および統語構造的な特徴として書き言葉と話し言葉双方の性質を併せ持っていること、例外的な言語現象が少ない規範的なデータで解析が容易であることを確認した。完成した UD Japanese-JDD は主に研究用途として公開することを予定している。

3) <https://github.com/explosion/spaCy/tree/v3.7.2>

4) <https://huggingface.co/cl-tohoku/bert-base-japanese-v2>

5) Transformers, Morphologizer, Parser, Ner の順で配置

6) 対話・トピックごとに訓練。開発: テストの比率が 8:1:1 になるよう均衡性を保持してデータを分割した。UD は v2.12.

7) 各データの既定の分割に従った。UD はいずれも v2.12.

8) 完全版は付録 A に掲載。

謝辞

本研究は JSPS 科研費 JP22K17943, JST ムーンショット型研究開発事業 JPMJMS2011-35 (fundamental research), 国立国語研究所共同研究プロジェクトの支援を受けたものです。

参考文献

- [1] Pratik Joshi, Sebastin Santy, Amar Budhiraja, Kalika Bali, and Monojit Choudhury. The state and fate of linguistic diversity and inclusion in the NLP world. In **Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (ACL)**, 2020.
- [2] 赤間怜奈, 磯部順子, 鈴木潤, 乾健太郎. 日本語日常対話コーパスの構築. 言語処理学会第 29 回年次大会発表論文集, pp. 108–113, 3 2023.
- [3] Matīss Rikters, Ryokan Ri, Tong Li, and Toshiaki Nakazawa. Designing the business conversation corpus. In **Proceedings of the 6th Workshop on Asian Translation (WAT)**, pp. 54–61, 2019.
- [4] Hiroaki Sugiyama, Masahiro Mizukami, Tsunehiro Arimoto, Hiromi Narimatsu, Yuya Chiba, Hideharu Nakajima, and Toyomi Meguro. Empirical analysis of training strategies of transformer-based japanese chat systems. In **2022 IEEE Spoken Language Technology Workshop (SLT)**, pp. 685–691, 2023.
- [5] Yanran Li, Hui Su, Xiaoyu Shen, Wenjie Li, Ziqiang Cao, Shuzi Niu, and Hong Kong. DailyDialog: A Manually Labelled Multi-turn Dialogue Dataset. In **Proceedings of the The 8th International Joint Conference on Natural Language Processing (IJCNLP)**, pp. 986–995, 2017.
- [6] Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur Szlam, Douwe Kiela, and Jason Weston. Personalizing Dialogue Agents: I have a dog, do you have pets too? In **Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (ACL)**, Vol. 1, pp. 2204–2213, 2018.
- [7] Hannah Rashkin, Eric Michael Smith, Margaret Li, and Y. Lan Boureau. Towards Empathetic Open-domain Conversation Models: A New Benchmark and Dataset. In **Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics (ACL)**, pp. 5370–5381, 2019.
- [8] Marie-Catherine de Marneffe, Christopher D. Manning, Joakim Nivre, and Daniel Zeman. Universal Dependencies. **Computational Linguistics**, Vol. 47, No. 2, pp. 255–308, 2021.
- [9] Kikuo Maekawa. Design, compilation, and some preliminary analyses of the corpus of spontaneous Japanese. **Spontaneous speech: Data and analysis**, Vol. 3, pp. 87–108, 2004.
- [10] Kiyotaka Uchimoto, Ryoji Hamabe, Takehiko Maruyama, Katsuya Takanashi, Tatsuya Kawahara, and Hitoshi Isahara. Dependency-structure annotation to corpus of spontaneous Japanese. In **Proceedings of the Fifth International Conference on Language Resources and Evaluation (LREC)**, pp. 635–638, 2006.
- [11] Hanae Koiso, Haruka Amatani, Yasuharu Den, Yuriko Iseki, Yuichi Ishimoto, Wakako Kashino, Yoshiko Kawabata, Ken’ya Nishikawa, Yayoi Tanaka, Yasuyuki Usuda, and Yuka Watanabe. Design and evaluation of the Corpus of Everyday Japanese Conversation. In **Proceedings of the Thirteenth Language Resources and Evaluation Conference**, pp. 5587–5594, 2022.
- [12] Masayuki Asahara, Hiroshi Kanayama, Takaaki Tanaka, Yusuke Miyao, Sumire Uematsu, Shinsuke Mori, Yuji Matsumoto, Mai Omura, and Yugo Murawaki. Universal Dependencies version 2 for Japanese. In **Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC)**, pp. 1824–1831, 2018.
- [13] Mai Omura and Masayuki Asahara. UD-Japanese BC-CWJ: Universal Dependencies annotation for the Balanced Corpus of Contemporary Written Japanese. In **Proceedings of the Second Workshop on Universal Dependencies**, pp. 117–125, 2018.
- [14] Mai Omura, Hiroshi Matsuda, Masayuki Asahara, and Aya Wakasa. UD-Japanese-CEJC: Dependency relation annotation on corpus of everyday Japanese conversation. In **Proceedings of the 24th Annual Meeting of the Special Interest Group on Discourse and Dialogue**, pp. 324–335, 2023.
- [15] 伝康晴, 小木曾智信, 小椋秀樹, 山田篤, 峯松信明, 内元清貴, 小磯花絵. コーパス日本語学のための言語資源: 形態素解析用電子化辞書の開発とその応用. 日本語科学, Vol. 22, pp. 101–123, 2007.
- [16] 小木曾智信, 中村壮範. 『現代日本語書き言葉均衡コーパス』形態論情報アノテーション支援システム的设计・実装・運用. 自然言語処理, Vol. 21, No. 2, pp. 301–332, 2014.
- [17] Masayuki Asahara and Yuji Matsumoto. BCCWJ-DepPara: A syntactic annotation treebank on the ‘Balanced Corpus of Contemporary Written Japanese’. In **Proceedings of the 12th Workshop on Asian Language Resources**, pp. 49–58, 2016.
- [18] Masayuki Asahara, Yuji Matsumoto, and Toshio Morita. Demonstration of ChaKi.NET – beyond the corpus search system. In **Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: System Demonstrations**, pp. 49–53, 2016.
- [19] 一般社団法人共同通信社. 記者ハンドブック 第 14 版 新聞用字用語集. 2022.
- [20] Kikuo Maekawa, Makoto Yamazaki, Toshinobu Ogiso, Takehiko Maruyama, Hideki Ogura, Wakako Kashino, Hanae Koiso, Masaya Yamaguti, Makiro Tanaka, and Yasuharu Den. Balanced corpus of contemporary written Japanese. **Language resources and evaluation**, Vol. 48, No. 2, pp. 345–371, 2014.

A UD 言語資源に基づく依存構造解析器の性能評価

表 5 に、UD 言語資源を活用して作成した依存構造解析器の各設定における解析精度を示す。評価観点、単語分割精度 (Tokens)、UD 品詞ラベルの予測精度 (UPOS)、UniDic 形態論情報の予測精度 (XPOS)、語彙素の予測精度 (Lemmas)、依存関係の予測精度 (UAS)、依存関係に加えて UD 係り受けラベルの予測精度 (LAS) の全 6 つで、いずれも F 値で精度を算出した。spaCy を用いると、解析対象が同一のデータの場合は訓練データによらず Tokens, XPOS, Lemmas がほぼ同じ結果になった。今後、spaCy 内部で何が起きているかについて検討する。表中の太字は各観点における最良の値を示す。

表 5 UD 言語資源に基づく依存構造解析器の解析精度 (完全版).

訓練データ	解析対象	Tokens	UPOS	XPOS	Lemmas	UAS	LAS
CEJC	CEJC	95.41	93.35	89.28	86.13	88.13	86.17
GSD	CEJC	95.41	79.77	89.28	86.13	80.69	74.68
JDD	CEJC	95.41	70.34	89.28	86.13	71.34	62.82
CEJC+GSD	CEJC	95.41	93.47	89.28	86.13	88.42	86.64
CEJC+JDD	CEJC	95.41	84.05	89.28	86.13	75.63	68.18
GSD+JDD	CEJC	95.41	68.57	89.28	86.13	71.82	60.86
CEJC+GSD+JDD	CEJC	95.41	82.65	89.28	86.13	75.93	67.64
CEJC	GSD	98.09	84.27	96.91	95.16	79.54	70.41
GSD	GSD	98.09	96.99	96.91	95.16	91.25	90.29
JDD	GSD	98.09	88.59	96.91	95.16	78.05	71.75
CEJC+GSD	GSD	98.09	97.11	96.91	95.16	91.15	90.14
JDD+CEJC	GSD	98.09	86.80	96.91	95.16	75.44	68.84
JDD+GSD	GSD	98.09	92.61	96.91	95.16	84.25	80.46
JDD+CEJC+GSD	GSD	98.09	92.17	96.91	95.16	82.40	78.31
CEJC	JDD	98.57	83.65	97.81	93.04	91.61	78.07
GSD	JDD	98.57	97.27	97.81	93.04	93.81	92.95
JDD	JDD	98.57	96.03	97.81	93.04	92.34	89.80
CEJC+GSD	JDD	98.57	97.74	97.81	93.04	94.18	93.24
CEJC+JDD	JDD	98.57	95.13	97.81	93.04	90.57	87.76
GSD+JDD	JDD	98.57	95.39	97.81	93.04	91.25	88.86
CEJC+GSD+JDD	JDD	98.57	94.83	97.81	93.04	90.50	87.80