

# 読みコーパスの整備と語義曖昧性解消による 読みと意味の頻度表の作成

2025/2/27

東京農工大学 古宮研究室 大井恵奈

# 目次

- 本研究の背景と目的
- コーパス(文例集)の内容
- それぞれのコーパスへの読みと語義番号の付与方法
- 対応表の作成方法
- 結果と考察

# 背景

- 近年、音声認識システムや読み上げシステムの需要は高まっている
  - 視覚障がい者の方のための機能
  - 機械音声で作成された動画
- 特に、日本語のシステムでは漢字の読みの正確性は大きな課題

# 背景

例

- ・ 「今日」のおやつは「最中」です.



# 背景

例

- ・ 「今日」のおやつは「最中」です.



きょう

こんにち



~~さいちゅう~~

もなか

# 背景

例

- ・ 「今日」のおやつは「最中」です.



- ・ 「きょう」のおやつは「もなか」です.

→話者の当日のおやつが「もなか」である.

- ・ 「こんにち」のおやつは「もなか」です.

→最近よく食べられているおやつは「もなか」である.

どちらの読みも正しいが意味が異なる

# 関連研究

語義曖昧性解消(Word-sense disambiguation : WSD)

…文中のある単語が、どの語義を表すのかを判定する  
自然言語処理のタスク

I grow plants. ←植物？工場？

現在は様々な言語で、教師あり学習と大規模な事前学習済み  
モデルを使用することで一定の性能を達成している.  
(Shumidmanら<sup>[1]</sup>, 浅田ら<sup>[2]</sup>)

# 目的

- ・ 『日本経済新聞記事オープンコーパス』, 『日本語話し言葉コーパス』 を対象に, 読みと意味を付与したデータを整備
- ・ 日本語の微妙な読みと意味の関係を形態素レベルで明らかにするため, 読みと意味の対応表を作成



# 日本経済新聞記事オープンコーパス<sup>[3]</sup>

- 2013年1-2月の日本経済新聞朝夕刊 96記事から作成
- 国立国語研究所によりUniDic 形態論情報(短単位・長単位)・文節係り受けを付与
- 加藤らによって語義情報として『分類語彙表』の分類番号が人手で付与されている

短単位形態素数	33,346
長単位形態素数	24,379
文節数	10,627
文数	1,333
記事数	96

日経が提供するデータ



日本経済新聞記事  
オープンコーパス

日本経済新聞記事オープンコーパスは、日本経済新聞の朝夕刊(2013年1~2月)から選出した約100本の記事を元に、国立国語研究所が作成した日本語の書き言葉コーパスです。

記事データに対して国立国語研究所が形態論情報(短単位・長単位)と文節係り受け情報のアノテーションを人手でおこなっています。

コーパスデータとその元となった約100本の記事データを無償で公開するものです(無償での利用は研究利用に限ります。商用利用については、日本経済新聞社と有償契約が必要です)。

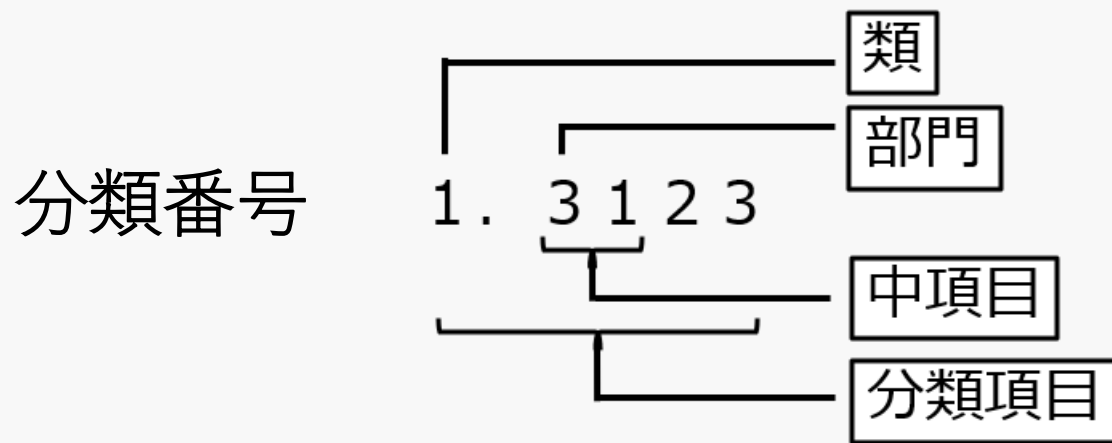
# 日本語話し言葉コーパス (Corpus of Spontaneous Japanese : CSJ)<sup>[4]</sup>

- 日本語の自発音声を大量にあつめて多くの研究用情報を付加した話し言葉研究用のデータベース
- 音声の書きおこしのデータのみで，語義情報が付与されたデータは非公開

転記テキスト	752万語
短単位形態素数	752万
長単位形態素数	631万
文節数	50万

# 分類語彙表<sup>[5]</sup>

- 語を意味によって分類・整理したシソーラス（類義語集）
- レコード ID 番号／見出し番号／レコード種別／類／部門／中項目／分類項目／**分類番号**など15項目から成る



本研究では、これを意味として扱う

# コーパスと整備方法

- 日本経済新聞記事オープンコーパス  
語義番号は付与されているため、読みを付与する  
+ クラウドソーシング  
+ 専門家によるアノテーション
- 日本語話し言葉コーパス  
音声の書き下しテキストに語義情報を付与する

# 『日本経済新聞記事オープンコーパス』 漢字の読みのクラウドソーシング調査

各例文の【 】の部分の可能な読みをすべて選択してください

全国に約2【0】00店ある三菱電機の系列販売店で12月末から発売した。

☐ ゼロ ☐ レイ  
☐ 他の読みがある

金融機関が中小企業からの条件変更要請を受けると、借り手企業との【間】で合意した経営改善計画に基づき返済期限が延長され、月々の返済額が減額される。

☐ アイダ ☐ カン  
☐ 他の読みがある

98年発売のパソコン「iMac」、【2】001年には携帯プレーヤー「iPod」。

☐ ニ ☐ ニイ  
☐ フタ ☐ フツ  
☐ プタ ☐ プタ  
☐ 他の読みがある

自公両党は低【所得】層ほど負担感が重いとされる「逆進性」の対策として、軽減税率が有力だとの見方では一致している。

☐ ショトク ☐ トコロエ  
☐ 他の読みがある

サッカー【元】日本代表選手で、現在は横浜F・マリノスのアンバサダーである波戸康広さんとの出会いは2年前のこと。

☐ ガン ☐ ゲン  
☐ モト ☐ 他の読みがある

短単位 33,346語のうち、UniDic で  
読みの曖昧性がある 4,260語を対象

すべての可能な読みと「他の読み  
がある」を選択肢として展開

1表現あたり 2人ずつ10回調査

2024/01/01 11:00-13:10 に実施

異なり 1,171 人が参加

2回答あたり2円相当の謝礼

# 調査概要

- 例文内の単語に対して調査(複数回答可)
  - 例
    - ・ ー→イチ/イッ/ヒト/ビト/他の読みがある
    - ・ 厳しい→イカメシイ/キビシイ/他の読みがある
- 一つの例文に対して複数の単語を対象に読みを質問
  - 直後、再びネイマールが打ったミドルを【止める】と、目を大きく見開いた川島がほえた。
  - 直後、再びネイマールが【打っ】たミドルを止めると、目を大きく見開いた川島がほえた。

# 問題点と修正

一定数誤った読みを付与する作業者がいる

各解答で外れ値（1件のみ/20件中）を選択した作業者の回答に対して、外れ値を回答することにより有効票を点数を初期値 1.0 として 0.95 を乗ずる

修正した有効票を割合に変換（足して 1.0 になるように修正）

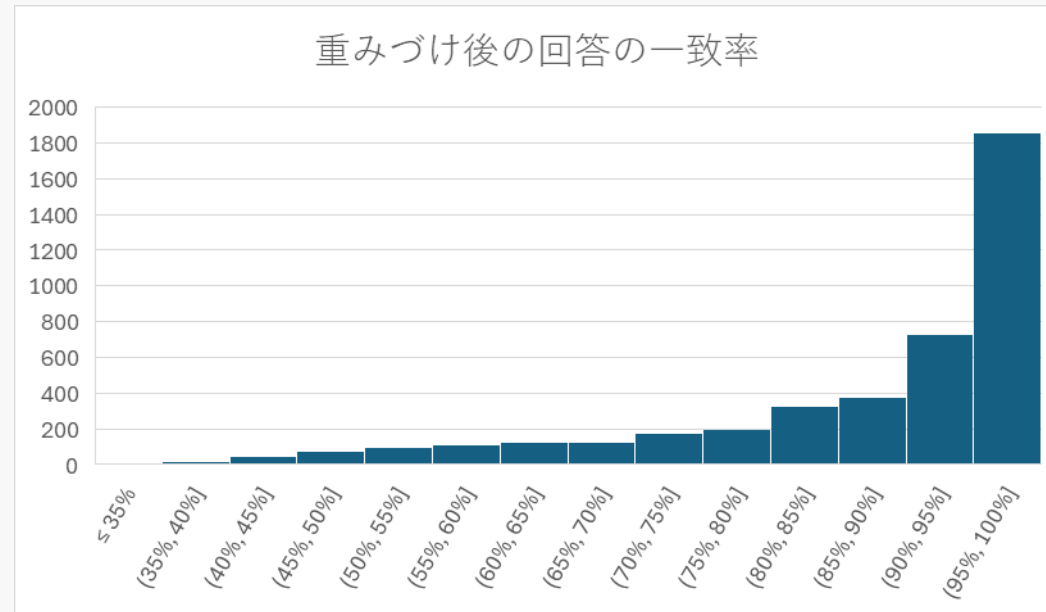
→外れ値の回答が多い回答者の重みが軽くなるような調整

回答が一致したものの件数は448件

# 回答結果と割合

- 一番多かった回答の割合(重み修正済み)が**0.95**以上→1857件
  - 横河電、【経常】(ケイジョウ)益 4 2 % 増、4 ～ 1 2 月 8 9 億円、資源国向け堅調。…95%
- 一番多かった回答の割合(重み修正済み)が**0.90**以上→2585件
  - ■石油資源開発 1 6 日、【米】(ベイ)テキサス州で新型原油「シェールオイル」の鉱区権益を追加取得したと発表した。…90%

今回は、割合が一番高かった回答を読みとして採用





# 『日本経済新聞記事オープンコーパス』 漢字の読みの人手によるアノテーション

コーパス内の全ての形態素に対して，Unidicを用いてすべての可能な読みと「他の読みがある」を選択肢として展開し，人手による付与を行う．

※

読みは複数選択可

表層系	読み仮名
その	ソノ
後	アト，ゴ，ノチ
に	ニ
税率	ゼイリツ
が	ガ

表層系	読み仮名
許	キヨ，シュー
其亮	キリョウ，チーリャン

短単位では付与できない  
場合，長単位に付与

表層系	読み仮名
1	センキュウヒャクゴジュウ
9	↓
5	↓
0	↓
年	ネン
大会	タイカイ

略語は一般的な読み仮名  
を付与

表層系	読み仮名
日本	ニッポン，ニホン
の	ノ
DF	ディフェンダー
を	ヲ
抜き去り	ヌキサリ

# 『日本語話し言葉コーパス』への語義番号付与

浅田ら<sup>[2]</sup>の提案手法により，日本語BERTのFine-tuningを利用して，コーパス内の単語全てを対象とするWSD(all-words WSD)をCSJに行った．

CSJ

表音	表層	…	分類番号
マズ	先ず		3.1650
ハッピョウ	発表		1.3140
ハ	は		—
サイショ	最初		1.1651

新たに付与

# コーパスと整備方法(再掲)

- 日本経済新聞記事オープンコーパス  
語義番号は付与されているため、読みを付与する  
+ クラウドソーシング  
+ 専門家によるアノテーション
- 日本話し言葉コーパス  
音声の書き下しテキストに語義情報を付与する

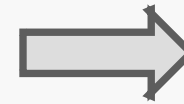
# 対応表の作成

- 以上の整備したコーパスから，読みと語義の組を抽出
- 複数読みがある場合は，全ての組を作成

その【後】に税率が10%を超える  
【後】……1.1650



ゴ/アト/ノチ



作成する読みと語義の組  
ゴー1.1650  
アトー1.1650  
ノチー1.1650

- 全ての読みと語義の組を数え，表にする

後	3.1670	1.1650	1.1670
ゴ		+1	
アト		+1	
ノチ		+1	

# 対応表(例:日本語話し言葉コーパス)

掘る	2.3822
ホル	1

両端	1.1960
リョウタン	4
リョウハシ	1

	読みが一つ	読みが複数
語義が一つ	6102	68
語義が複数	493	70

用法	1.1960	1.3081
ヨウホウ	1	5

略	3.1920	1.1251
リャク	0	1
ホボ	130	0

# 対応表の件数

日本経済新聞記事オープンコーパス

- クラウドソーシング

	読みが一つ	読みが複数
語義が一つ	375	32
語義が複数	50	21

- 人手によるアノテーション

	読みが一つ	読みが複数
語義が一つ	371	33
語義が複数	50	24

日本語話し言葉コーパスの例

	読みが一つ	読みが複数
語義が一つ	6102	68
語義が複数	493	70

# 考察－「後」

日本語話し言葉コーパス

後	1.1643	1.1650	1.1670	1.1740	3.1670
コウ	0	0	1	0	0
ゴ	0	9	63	0	0
アト	7	39	297	10	14
ノチ	0	2	28	0	1

分類語彙表(一部抜粋)

分類番号	類・部門・中項目・分類項目	用例
1.1643	体・関係・時間・未来	先, 先週, 今後, 明日
1.1650	体・関係・時間・時間的順序	次, 初, 順, 優先
1.1670	体・関係・時間・時間的	食前, 食後, 前日
1.1740	体・関係・空間・左右・前後・たてよこ	前方, 門前
3.1670	相(形容詞的)・関係・時間・時間的前後	まだ, 以後, もう

日本経済新聞記事オープンコーパス

- クラウドソーシング

後	1.1650	1.1670	3.1670
ゴ	1	1	1
アト	0	2	1
ノチ	1	0	1

- 人手によるアノテーション

後	1.1650	1.1670	3.1670
ゴ	1	1	1
アト	2	2	1
ノチ	2	0	1

# 考察－「後」

日本語話し言葉コーパス

左右・前後  
・たてよこ

未来 時間的順序 時間的

時間的前後

後	1.1643	1.1650	1.1670	1.1740	3.1670
コウ	0	0	1	0	0
ゴ	0	9	63	0	0
アト	7	39	297	10	14
ノチ	0	2	28	0	1

分類語彙表(一部抜粋)

分類番号	類・部門・中項目・分類項目	用例
1.1643	体・関係・時間・未来	先, 先週, 今後, 明日
1.1650	体・関係・時間・時間的順序	次, 初, 順, 優先
1.1670	体・関係・時間・時間的	食前, 食後, 前日
1.1740	体・関係・空間・左右・前後・たてよこ	前方, 門前
3.1670	相 (形容詞的)・関係・時間・時間的前後	まだ, 以後, もう

日本経済新聞記事オープンコーパス

クラウドソーシング

時間的順序 時間的 時間的前後

後	1.1650	1.1670	3.1670
ゴ	1	1	1
アト	0	2	1
ノチ	1	0	1

人手によるアノテーション

時間的順序 時間的 時間的前後

後	1.1650	1.1670	3.1670
ゴ	1	1	1
アト	2	2	1
ノチ	2	0	1



# 考察－「後」 日本経済新聞記事オープンコーパス

- クラウドソーシング

	時間的順序	時間的	時間的前後
後	1.1650	1.1670	3.1670
ゴ	1	1	1
アト	0	2	1
ノチ	1	0	1

- 人手によるアノテーション

	時間的順序	時間的	時間的前後
後	1.1650	1.1670	3.1670
ゴ	1	1	1
アト	2	2	1
ノチ	2	0	1

人手によるアノテーションの方が、複数の読みを回答することが多い

日本経済新聞記事オープンコーパスでは、時間的關係を表す「後」は「ノチ」と読まない

# 考察－「後」コーパス間の違い

日本経済新聞記事オープンコーパス  
クラウドソーシング

時間的順序 時間的 時間的前後

後	1.1650	1.1670	3.1670
ゴ	1	1	1
アト	0	2	1
ノチ	1	0	1

日本語話し言葉コーパス

左右・前後  
・たてよこ

未来 時間的順序 時間的

／ 時間的前後

後	1.1643	1.1650	1.1670	1.1740	3.1670
コウ	0	0	1	0	0
ゴ	0	9	63	0	0
アト	7	39	297	10	14
ノチ	0	2	28	0	1

話し言葉コーパス(口語)では時間的關係の「後」を「ノチ」と読むこともある

- ・ 書き言葉と話し言葉の違い
- ・ コーパスの大きさの違い

# まとめ

- 日本経済新聞記事オープンコーパスには読みをクラウドソーシングと人手によるアノテーションで付与した
- 日本語話し言葉コーパスには分類番号を自動で付与した  
→二つのコーパスに語義と読みを整備
- 整備したコーパスから，日本経済新聞記事オープンコーパスからは478件，CSJからは6,733件の語について，読みと意味のペアの出現数の表を作成
- コーパスによって，異なる読みの分布が得られた  
→読みの辞書作成や，読み推定システムへの利用