(O

2017　2019

The Corpus of Kansai Vernacular Japanese

The Corpus of Kansai Vernacular Japanese

HEFFERNAN, KEVIN

1,500,000

138
Mecab
98
Google Web

The corpus of Kansai Vernacular Japanese includes speakers ranging in age from 15 years old to 80 years old. This large age range allows researchers to explore how language is used differently by each generation. Such an approach is useful for doing research language change and standardization.

The two primary objectives of this research project were to create a corpus of Kansai vernacular Japanese, and to make that corpus available on the internet. In total, 138 sociolinguistic interviews were conducted. Each interview was transcribed, and checked for errors. The transcriptions were parsed and tagged with part of speech data using Mecab. The tagged data was checked by students hired for this job, and mistakes were corrected. I estimate the final accuracy rate is about 98%
The final data, along with supporting documents such as a description of the transcription methods, are available on a google website. The data is shared under a creative commons license. Users may be downloaded and used free of charge. However, users are prohibited from using the data for profit.

linguistics

Japanese  dialect  corpus

One of the goals of Variationist linguistics is to explore the way language changes over time. A popular methodology for doing so is the apparent-time method. This method requires that the researcher gather data from speakers who span a wide age range, from as young as teenagers to elderly people. The researcher then observes how people of different age groups speak differently. As a simple example, consider the use of the word            . This word has increased in usage among younger people. Furthermore, carefully observing the differences in usage between younger and older speakers shows that the grammatical changes of

parallel the English word *like*. This parallelism suggests that there is a universal mechanism common to all languages of the world that controls the way words change into discourse markers such as            and *like*.

When this project was started, there was very little data natural conversation data available for research, and the available data did not cover a wide age range. Furthermore, this project collected conversations between speakers of the Kansai dialect. Previously, there was no large scale publicly available dataset of Kansai speech.

The two primary objectives of this research project were to create a corpus of Kansai vernacular Japanese, and to make that corpus available on the internet for others to use. Furthermore, a wide age range of speaker data was collected in order to allow for apparent-time studies of the way that spoken Japanese, and in particular Kansai dialect, was changing from one generation to the next.

In order to make the data easy to use by others, it is shared under a creative commons license. This license allows anyone to download the data free of charge and use it for non-commercial purposes. Also, if the user changes the data (for example, improves it), then he or she must share the data free of charge with others.

In total, 138 sociolinguistic interviews were conducted. The interviews were conducted by students attending a university in Kansai for course credit. The students were trained by researcher to conduct an interview following established sociolinguistic methods. The primary goal of such an interview is to have the interviewee relax and speak in a very casual style. All of the interviewers and interviewees are native Kansai Japanese speakers. Each interview was transcribed, and checked for errors. The transcriptions were parsed and tagged with part of speech data using Mecab. The tagged data was checked by students hired for as part-time workers, and mistakes were corrected. The final accuracy rate of the part of speech information is about 98%. Furthermore, each line of data is marked for the speaker (either interviewer or interviewee) to allow researchers to easily sort out the speech of the student interviewers from the interviewees.

The primary result of this project is the interview data. The final data, along with supporting documents such as a description of the transcription methods, are available on a google website. Following is a sample of ten lines of data. The entire data collection is over two million lines long.

```
s,      ,     ,*,*,*,*,        ,           ,
s,      ,        ,     ,*,*,*,  ,  ,
s,      ,    ,*,*,*,*,   ,      ,
s,      ,        ,     ,*,*,*,   ,  ,
s,      ,    ,*,*,            ,        ,     ,           ,
s,      ,      ,*,*,*,*,         ,        ,
s,      ,    ,*,*,     ,        ,         ,        ,
s,         ,*,*,*,            ,          ,     ,          ,
s,         ,*,*,*,           ,        ,  ,  ,
s,      ,     ,*,*,*,*,   ,  ,
```

Other than that data, several research articles and a book were also produced. The research articles investigated two areas: language change and language processing. With regards to language change, the following topics were reported on: changes in the grammatical nature of the word        ; changes in *u-onbin* in Kansai dialect; and the spread of lexical phrases of such as                from a large urban area to a small rural area.

With regards to language processing, the following topics were reported on. The omission of the case particles      and      in spoken Japanese; and the processing of the negative forms of verbs in spoken Japanese. Examining these language phenomena in details shows us how language is processed by our memory.

Finally, *The grammar of Kansai vernacular Japanese* was published by the K.G. University Press. This book is an introduction to Kansai Japanese written in English with Japanese examples, an aim at intermediate level learners of Japanese. The book covers a wide range of grammar of spoken Japanese, including Kansai dialect, but also grammar used throughout Japan. Some examples of the Kansai dialect grammar covered by the book are the copula *ya*, the verbal negative suffixes -    , -    , -    , and the word        . Some examples of general spoken Japanese covered by the book are      (as in            ) and the emphasis of adjectives by double consonants (as in                ). Each grammatical point is illustrated with several examples of natural spoken Japanese taken from the corpus data. Furthermore, each grammatical point is rated for frequency of usage based on the rate of usage in the corpus data. This rating system allows beginner students to focus on the most commonly-occurring phrases only.

| 2 | 2 | 2 | 1 | |
|---|---|---|---|---|
| Heffernan Kevin  Sato Yo | | | | 3 |
| Relative frequency and the holistic processing of morphology | | | | 2017 |
| Asia-Pacific Language Variation | | | | 67 94 |
| DOI<br>10.1075/aplv.3.1.04hef | | | | |
| | | | | |

| Heffernan Kevin  Hiratuka Yusuke | 3 |
|---|---|
| Morphological relative frequency impedes the use of stylistic variants | 2018 |
| Asia-Pacific Language Variation | 200 231 |
| DOI<br>10.1075/aplv.16009.hef | |
| | |

| 7 | 0 | 7 |
|---|---|---|

| Heffernan, Kevin and Yo Sato |
|---|
| The varying complexity of the syntactic role of nouns: Evidence from Japanese corpora |
| 4th Asia Pacific Corpus Linguistics Conference |
| 2019 |

| Sato, Yo, and Kevin Heffernan |
|---|
| How distinctive is your corpus? A metric for the degree of deviation from the norm. |
| 4th Asia Pacific Corpus Linguistics Conference |
| 2019 |

Yo Sato and Kevin Heffernan

Creating dialect sub-corpora by clustering: a case in Japanese for an adaptive method

Eleventh International Conference on Language Resources and Evaluation

2019

Heffernan, Kevin, Imanishi, Yusuke, Honda, Masaru, and Sato, Yo

Case particle omission correlates with object syntactic complexity: Evidence from the Corpus of Kansai Vernacular Japanese.

Exploiting Parsed Corpora: Applications in Research, Pedagogy, and Processing

2017

Heffernan, Kevin

The diffusion of lexical bundles from an urban center to a rural community in Japan

16th International Conference on Methods in Dialectology

2017

Heffernan, Kevin, and Sato, Yo

Age-related patterns in lexical bundle usage: Evidence from a corpus of vernacular Japanese

Corpus Linguistics International Conference

2017

Sato, Yo, Heffernan, Kevin, Kishie, Shunsuke, and Hattori, Kota

Quantifying 'standardness' of the language use in a locality: a study with Twitter data

Corpus Linguistics International Conference

2017

| 1 | |
|---|---|
| Kevin Heffernan | 2019 |
| | 174 |
| The Grammar of Kansai Vernacular Japanese | |

The Corpus of Kansai Vernacular Japanese
https://sites.google.com/view/kvjcorpus

| | | | |
|---|---|---|---|
| | | | |