Automatic Generation of Japanese Question-Answering Pairs

Hiroki Tanioka^{*1} Kaoru Kimura^{*1}

Kazuma Takaoka^{*2} Ryohei Nakatani^{*2} Yoshitaka Uchida^{*2}

*1 Tokushima University, Tokushima, Japan *² Works Applications, Tokushima, Japan {tanioka.hiroki, c501506035}@tokushima-u.ac.jp {takaoka_k, nakatani_r, uchida_yo}@worksap.co.jp

Introduction

What's the goal?

Question-answering pairs are needed for a lot of question-answering and chatbot systems. Because those systems are based on search engine. If many and high-quality question-answering pairs are indexed on the search engine, the question-answering system is expected to make a response to users.

Japanese Corpus Needed

The question-answering system needs lots of sophisticated question-answering corpora. There are already large- scale English question answer corpora available. However, there are no large-scale Japanese question-answer corpus. Therefore, we need more Japanese questionanswering corpora.

Machine Learning!!

A low-cost approach is needed for generating question-answering corpus. Furthermore, various types of question-answer corpora are needed for each domain.

Approach

Preparing Question-Answering Pairs

Step 1. Choose a category in the following 10 categories.

- Step 2. Choose 10 topics in the category.
- Step 3. Make 6 types of questions in each topic.
- Step 4. Describe 1 or more answers for each question.

| | | ~ | <u> </u> | |
|----------------|-----------|-----------|-----------|------------------|
| Category | Number of | Number of | Number of | |
| | Themes | Questions | Answers | |
| 学術(Academia) | 10 | 100 | 100 | |
| 技術(Technology) | 10 | 124 | 156 | |
| 自然(Nature) | 11 | 127 | 146 | |
| 社会(Society) | 11 | 108 | 132 | |
| 地理(Geography) | 10 | 102 | 100 | Total |
| 人間(Human) | 10 | 73 | 74 | |
| 文化(Culture) | 10 | 136 | 138 | Questions: 1,018 |
| 歴史(History) | 10 | 54 | 61 | Answers: 1 101 |
| マニュアル(Manual) | 10 | 74 | 74 | |
| シラバス(Syllabus) | 15 | 120 | 120 | |

Evaluation

BLEU and METEOR

BLEU is an algorithm for the evaluation machine translation output, based on matching reference cooccurrence in N-grams. Here, N = 4and BP_{RI}

Extracting Candidate Sentences

Extracting candidate sentences using Support Vector Machine (SVM) with Sudachi. The accuracy is 90.1%, and the precision is 41.4%.

| | Туре | SVM | Human | Target Sentence |
|-------------|------------------|----------|---|--|
| | True Positive | | ~ | 東京都の首長は、東京都知事である。 The head of Tokyo is the governor of Tokyo. |
| | | | ~ | 地震に対して、地殻が非常にゆっくりとずれ動く現象を地殻変動と呼ぶ。 Crustal deformation is a phenomenon in which the crust moves very slowly relative to an earthquake. |
| | False | | ~ | 魚料理や肉料理などの主菜など、イタリア料理のコースの流れがサブタイトルになっている。 A flow of course of Italian cuisine such as main dish such as fish dish and meat dish is subtitle. |
| Negative | | ~ | 1 メガパーセクは 326 万光年。 1 Mega persec is 3.26 million light-years. | |
| | √ False | ~ | | 2012 年現在、国際的な本初子午線として IERS 基準子午線が使用されている。 As of 2012, the IERS standard meridian is used as the international prime meridian. |
| Positive | ~ | | 東京都の議決機関は東京都議会である。 The decision-making body in Tokyo is the Tokyo Metropolitan Assembly. | |
| | | | | ツンドラ・タウン(Tundratown) 寒冷地域の動物たちが暮らすエリア。 Tundratown, an area is where cold region animals live. |
| Negative | | | | ディズニーのアニメーション映画で 10 億 ドルを突破したのは、『トイ・ストーリー3』、『アナと 雪の女王』に次いで史上 3 番目である。 It was the third one in history, after "Toy Story 3" and "Anna and the Snow Queen" that broke through \$ 1 billion in Disney's animated film. |
| * For Sudac | hi, analy | sing m | node A | , with the full version dictionary. For SVM, the kernel of SVM is RBF. The |
| parameters | are "c=1 | L000 -1 | w0 10 · | -w1 1". Train data contains 5.000 sentences including 188 sentences chosen |
| manually by | <i>i</i> human | Tost (| hata co | intains 2 000 sentences |
| manually Dy | numan | . iest (| | |

To make the index quickly, a machine learning technique can be used for generating question-answering pairs automatically.

It's a challenge to generate question sentences automatically.



Faculty of Engineering

Tokushima University

徳島人工知能NLP研究所

Learning Question Sentences

Attention Sequence to Sequence Model for learning to generate question sentences.



Here, CaboCha is employed for dependency parser. The following example includes chunks and tokens. Chunks are underlined.

Parsed sentence:

$$p_n = \frac{\sum_i S_i}{\sum_i T_i},$$

$$BLEU = BP_{BLEU} * \exp(\frac{1}{N} \sum_{n=1}^N \log p_n).$$

METEOR is a metric for the evaluation machine translation output, based on a F-mean score with precision P and recall R. α = 0.9, β = 3 and γ = 0.5. $F_{mean} = \frac{P \cdot R}{\alpha \cdot P + (1 - \alpha) \cdot R},$ $METEOR = F_{mean} \cdot (1 - \gamma (\frac{c}{u})^{\beta}).$

based on sentence similarity error.

$$PER = 1 - \frac{\sum_{i} d_{i}}{\sum_{i} m_{i}}$$

Dependency Error Rate (DER) is one of metrics using a dependency parser for evaluating Japanese fluency.

PER and DER

Position independent word Error Rate (PER) is an evaluation metrics,

$$DER = 1 - \frac{\sum_{i} g_i}{\sum_{i} f_i}$$

Where g means the number of dependency pairs in a generated sentence, with matching reference dependency pairs in a correct sentence, and f means the number of dependency pairs in a correct sentence.

* c means the number of words in both a generated sentence and a correct sentence, and u means the number of words in the generated sentence, with matching reference words in the correct sentence.

* S_i is the number of N-grams in a generated sentence i with the matching reference cooccurrence in a correct sentence. T_i is the number of N-grams in a generated sentence i. BP_{BLEU} means a penalty when a sentence is too shorter than a correct sentence.

| Candidate sentence | Generated question sentence | Correct Sentence | BLEU | METEOR | PER * | DER * | Scoro |
|--|--|---|-------|--------|-------|-------|-------|
| 世界遺産:世界遺産は、 「顕著な普遍的価値」を有 する文化遺産や自然遺産な どであり、1972 年に成立し た世界遺産条約に基づき、 世界遺産リストに登録され た物件を指す。 World Heritage: The World Heritage are cultural heritage and natural heritage which have "remarkable universal value," and it refers to the property registered in the World Heritage list based on the World Heritage Convention established in 1972. | 世界遺産は何ですか? What is the World Heritage Site? | 世界遺産は、「顕著な普遍 的価値」を有する文化遺産 や自然遺産のことですか? Is the world heritage a cultural heritage or natural heritage having "Remarkable Universal Value"? | 0.275 | 0.353 | 0.273 | 0.182 | -1 |
| 世界遺産:世界遺産は、政 府間委員会である世界遺産 委員会の審議を経て決定さ れる。 World Heritage Site: The World Heritage Site is decided after deliberation by the World Heritage Committee which is an intergovernmental committee. | 世界遺産は政府の委員会で すか? Is the World Heritage a government committee? | 世界遺産は、どこの審議を 経て決定されますか? Which deliberation is the World Heritage decided after? | 0.156 | 0.309 | 0.375 | 0.200 | 0 |
| 東京都:東京都の首長は、 東京都知事である。 Tokyo: The head of Tokyo is the governor of Tokyo. | 東京都の首長は何ですか? What is the head of Tokyo? | 東京都の首長を何といいま すか? What do you call the head of Tokyo? | 0.181 | 0.366 | 0.600 | 0.300 | 2 |
| 地震:地震に対して、地殻 が非常にゆっくりとずれ動 く現象を地殻変動と呼ぶ。 Earthquake: A phenomenon in which the crust moves very slowly against earthquake is | 地震が地殻には何を出すに 何と言いますか? What does the earthquake say to what to put out on the crust? | 地震に対して、地殻が非常 にゆっくりとずれ動く現象 を何と呼ぶか? What do you call a phenomenon that the crust moves very slowly against earthquakes? | 0.041 | 0.325 | 0.526 | 0.063 |)-1 |

| The average scores of evaluation results to | 196 |
|---|-----|
| candidate sentences. | |

| Score |
|-------|
| 0.095 |
| 0.235 |
| 0.388 |
| 0.137 |
| |

* BLEU score is 2-gram based. METEOR score did not employ WordNet. PER

| 「東京都の首長は、東京都知事 (The head of Tokyo is the govern | である。」 or of Tokyo.) |
|--|--|
| Dependency pairs: | |
| 東京 (Tokyo) | →都 (city) |
| 都 (city) | →の (of) |
| <u>東京都の</u> (of Tokyo city) | → <u>首長は、</u> (The head,) |
| 首長 (head) | \rightarrow (t <subject></subject> |
| は <subject></subject> | \rightarrow , (,) |
| <u>首長は、</u> (The head,) | → <u>東京都知事である。</u> (is the governor of Tokyo.) |
| 東京 (Tokyo) | →都知事 (the governor of Tokyo) |
| 都知事 (the governor of Toyo) | →で(is) |
| で (is) | → ある (is) |
| ある (is) | → 。 (.) |

Conclusion

- Japanese question sentences are generated from labeled corpus.
- Extracting candidate sentences is high accuracy 90% using SVM. However, Generating question sentences is still low scores in each metric.
- DER is reasonable metric, because DER has larger reduction rate than PER to grammatical error. Also, DER shows a average score.
- Hence, further research to increase generating accuracy and detect answers is necessary with other examples and machine learning models.



Discussion

For qualitative evaluation, the generated sentences have grammatical errors and unnatural particles in Japanese.

A blue boxed sentence is scored **good** in each metric without DER.

A red boxed sentence is scored wrong in each metric without DER.

Score is calculated by counting good square (+1) and wrong square (-1) for each sentence without DER. Then, **DER is fit to Score**.

What is JQAC Library?

Japanese Question-Answering Corpus (JQAC) is a dataset, consisting of question-answering pairs, which is manually made by university students on a set of Japanese Wikipedia articles and some public documents.

(distributed under the CC BY-SA 4.0 license):

JQAC Library. (2018). Japanese Question-Answering Corpus. https://taniokah.github.io/jqac/

▣,∗▣ ⊡d-mi

* PER* = 1 – PER, and DER* is 1 – DER