

統計的翻訳評価尺度に基づく日英翻訳文の訳質分析

鏑木 元*¹ 安田 圭志*² 山本 博史*² 匂坂 芳典*^{1,2,3}

*¹ 早稲田大学大学院 国際情報通信研究科

*² NiCT/ATR

*³ 早稲田大学 ことばの科学研究所

E-mail : *¹ hjm-tsubaki@asagi.waseda.jp

*^{2,3} {keiji.yasuda, hirofumi.yamamoto, yoshinori.sagisaka}@atr.jp

あらまし 統計的自動翻訳に用いる評価尺度を用い、日本人学習者日英翻訳文の客観的学習評価の可能性を調べた。評価尺度には2つの統計量を考えた。1つは日英翻訳文の英語らしさを測る指標として、単語Nグラムを用いた正解文に対する生起確率比を、他の1つは翻訳らしさを測る指標として、単語対応の翻訳確率生起確率比を利用した。日英翻訳文に対して人が付与した主観評価値とこれらの統計量との対応関係を分析した。この結果、翻訳らしさを測る指標に主観評価値と明確な対応関係が観察でき、この分析結果より、翻訳らしさを測る指標の客観的翻訳評価への利用可能性が示された。

キーワード 統計的自動翻訳、Nグラム、翻訳確率

Translation quality analysis of Japanese-English translation sentence based on statistical translation evaluation criterion

Hajime Tsubaki*¹ Keiji Yasuda*² Hirofumi Yamamoto*² Yoshinori Sagisaka*^{1,2,3}

*¹ GITI, Waseda University

*² NiCT/ATR

*³ Language and Speech Science Research Laboratories, Waseda University

E-mail: *¹ hjm-tsubaki@asagi.waseda.jp

*^{2,3} {keiji.yasuda, hirofumi.yamamoto, yoshinori.sagisaka}@atr.jp

Abstract A possibility of objective study rating of Japanese learners' Japanese-English translation sentences was examined using evaluation criteria for statistical machine translation. Two statistic values were designed for the evaluation. One index was a ratio of N-gram between a learner's translation sentence and a correct sentence to measure English characteristics in each sentence. The other index was a ratio of word translation probability between the two sentences to measure characteristics as the translation. Correspondence relation between subjective evaluation score for Japanese-English translated sentences and the statistic values of the indexes were analyzed. As for the index of translation characteristics, a clear correlating relationship was observed. The result showed an application possibility of machine translation evaluation criteria to measure translation characteristics of translation in objective translation evaluation.

Key words statistical machine translation, N-gram, translation probability

1. はじめに

我々は、第二言語の生成、知覚能力の測定、自動評定についての研究を進めている。器械による音声言語能力の自動評定は、評定の省力化や評定者によらない評価の可能性等、有用性が期待できる。しかしながら、音声言語データの準備等には労力と時間がかかり、少ないデータを基に研究が緒に着いた段階である。人間の翻訳、作文に関する言語能力の自動測定についてはこれまで、日本人英語の通じやすさの測定[1]、翻訳システムの自動評価手法の人間による翻訳の評価への適用[2]、英語コミュニケーション能力の自動測定[3]、自由英作文からの言語取得度の推定[4]等の研究がなされてきている。これらの研究結果からは、統計的な特徴量に基づく翻訳・作文能力の自動測定の可能性が伺える。

本研究では、日英翻訳を対象として、統計的自動翻訳で用いられる2つの統計的尺度についての利用可能性を検証した。日英翻訳文の英語らしさを示す単語Nグラム確率値、翻訳の妥当性を示す単語対応の翻訳確率値について、日英翻訳文に対して人が付与した主観評価値との相関関係を分析した。これらの尺度の翻訳評価への利用可能性、並びに、評定基準となる正解翻訳文が不要な翻訳評価の実現可能性を検討した。

2. 翻訳評価のための統計的情報量分析

統計的言語翻訳では、次式に示すように2つの特徴量を用い、翻訳が行われる。

$$\begin{aligned}\hat{e} &= \arg \max_{\text{all candidates}} p(e | j) \\ &= \arg \max_{\text{all candidates}} p(j | e)p(e)\end{aligned}$$

$p(e)$:対象言語文(英語) e の単語Nグラム確率
 $p(j|e)$:元文(英語) e が対象言語文(日本語) j に変換される翻訳確率

翻訳結果を得るための特徴量において、 $p(e)$ は、英語らしさを評価するための特徴量として理解できる。ま

た、 $p(j|e)$ は、翻訳らしさを評価する特徴量と考えることができる。ここでは、これらの特徴量を用いた客観評価を考え、主観評価値との対応関係を分析した。

$p(e)$ には、単語3グラム確率を考え、学習者翻訳文の $p(e)$ と正解翻訳のそれとの比を評価指標に用いた。 $p(j|e)$ には、単語対応の翻訳確率を用い、同様に学習者翻訳文の $p(j|e)$ と、正解翻訳のそれとの比を指標として用いた。

3. 分析実験

2つの指標と主観評価値との対応関係を分析した。分析データには、ATRで収録された旅行会話基本表現データ(以下、BTEC)を用いた。内訳は次の通りである。

- (1) BTEC 日本語会話文に対する、プロの翻訳者による日英翻訳文 162318 文。この対訳コーパスで言語モデル(英語)、翻訳モデルを構築。
- (2) BTEC 日本語会話文の日本語文 510 文、それらに対するプロの翻訳者、日本人英語学習者の日英翻訳及び発話したものの書き下し文 11220 文を作成。決まり文句等を除いた日本語文 473 文、それらに対する日英翻訳文 10406 文を評価対象データに使用。翻訳文の構成は以下の通り。

- ・日本語文 1 文あたり、プロの翻訳者が1つの日英翻訳文を生成(計 473 文)
- ・日本語文 1 文あたり、21 名の日本人学習者が日英翻訳文を生成(計 9933 文)
- ・日本人学習者の日英翻訳文には、1 人のネイティブ評定者が表1の基準に基づき、5段階の主観評価値を付与
(BTEC 日本語会話文 1 文に対し、21 名の日本人学習者が訳出した日英翻訳文の組を以後、テストセットとする。実験では、473 セットを使用)

日本人学習者が日英翻訳文 9933 文について、2つの指標を算出し、表1の訳質レベル毎にそれぞれ集計、平均値を算出する。これら2つの指標の平均値と訳質レベルとの対応関係を図1、2に示す。

表1 学習者英文の評定基準

| 訳質レベル | 評定基準 |
|-------|--|
| S ランク | 原文の情報が漏れ無く翻訳されており、訳出に文法的な間違いがない。使われている語彙もネイティブから見て自然である。 |
| A ランク | 使われている語彙はネイティブから見て不自然であるが、原文の情報が漏れ無く翻訳されており、訳出に文法的な間違いがない。 |
| B ランク | 原文のあまり重要でない情報が一部漏れていたり、訳出に文法的な間違いが若干あるが、容易に理解できる。 |
| C ランク | 原文の重要な情報が漏れていたり、訳出に文法的な間違いが大分あって、かなり崩れた訳出であるが、良く考えれば理解出来る。 |
| D ランク | 重要な情報が誤訳されており、理解不能である。 |

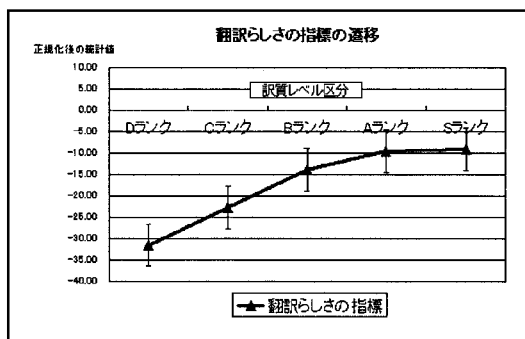


図1 翻訳らしさの指標と訳質レベルとの対応

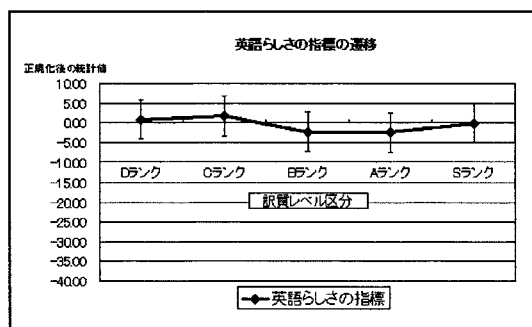


図2 英語らしさの指標と訳質レベルとの対応

4. 分析結果と考察

4.1 訳質レベルと統計的情報量の関係

翻訳らしさの指標は、図1のように訳質レベルが上がるとに従い、平均値は増加傾向を示した。この対応関係から、この指標の翻訳評価への利用可能性が示された。

分析実験での翻訳評価は、評価対象翻訳文を、表1の評価基準に基づいて正解翻訳と比較する絶対評価である。そのため、その基準に合致した正解翻訳が必要になる。一方、学習者の翻訳文同士の評価、すなわち翻訳文の相対評価の場合、この指標で評価対象翻訳文を順位付けすることで、それが可能になると考えられる。図1の対応関係からは、この可能性も示めされた。

英語らしさの指標は、図2のように訳質レベルにかかわらず、横ばいの傾向を示した。この傾向は、人間が日英翻訳文を作る場合、英語の語順を意識して作成しているため、単語Nグラム確率において差異が生じないことを反映していると考えられる。

4.2 正解翻訳を用いる評価手法との比較

評価基準としての正解翻訳が翻訳評価に及ぼす影響を調べるために、正解翻訳文を用いる翻訳評価手法との比較を行った。評価手法には、BLEUを用いた。BLEUは機械翻訳による翻訳文の自動評価手法で、人間による正解翻訳を評価基準とし、単語Nグラムの重なりに基づき、翻訳システムの生成した翻訳文を評価する。スコアは次式で求められる。

$$S_{BLEU} = BP \times \exp\left(\sum_{n=1}^N w_n \log P_n\right)$$

$$BP = \begin{cases} 1 & \text{if } c > r \\ \exp(1-r/c) & \text{if } c < r \end{cases}$$

cは、翻訳システム結果の文長
rは、人間による正解翻訳の文長
 P_n は、修正Nグラム精度

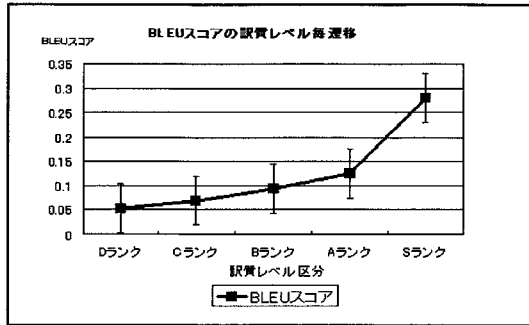


図3 BLEU スコアと訳質レベルとの対応

翻訳システムの結果と人間による正解翻訳を、学習者の翻訳文とプロの翻訳者による正解翻訳文に置き換え、スコアを算出し、訳質レベルとの対応関係を見た。

図3のように、訳質レベルが上がるに従い、BLEU スコアの平均値は増加傾向にあった。注目すべきはAランクとSランクとの境界で、翻訳らしさの指標では誤差の重なりが大きい部分(図1)が、BLEU スコアでは明確に区別された。Sランクの翻訳文のほうがAランクのそれらよりも、正解翻訳に用いられている語彙を多く含んでいることが、この区別を生じさせていると考えられる。

4.3 テストセット数と統計的情報量の関係

翻訳らしさの指標について、訳質レベルとの明確な対応関係が確認でき、翻訳評価に利用できる可能性が見出せた。しかし、翻訳文1文、すなわち1つのテストセットから得られる指標のみで翻訳評価を行うことは、誤差が大きいため困難である。そこで、どのぐらいの量のテストセットを用意すれば、評価が可能になるのかを見るために、テストセット数を徐々に増加させ、翻訳らしさの指標と訳質レベルとの対応関係、分散の変化を調べた。

図4のように、テストセット数が少ない場合、翻訳らしさの指標と訳質レベルとの対応関係が安定しない。だが、テストセット数の増加に伴い、対応関係は安定していくことが分かる。さらに表2からは、ランク内分散が小さくなっていることがわかる。この結果は、一定量のテストセットから求めた翻訳らしさの指標を用いれば、学習者の翻訳能力を評価できる可能性を示している。

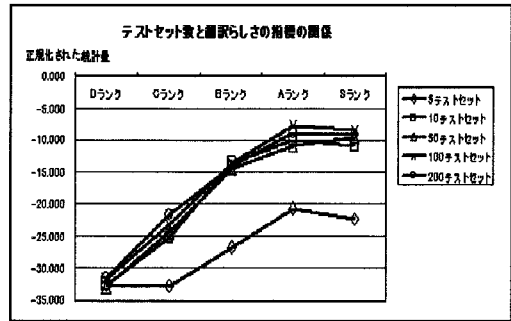


図4 テストセット数と訳質レベルの関係

表2 テストセット数とランク内分散

| テストセット数 | ランク内分散 |
|---------|--------|
| 10 | 14.34 |
| 50 | 10.74 |
| 100 | 9.66 |
| 200 | 9.61 |

5. まとめ

統計的自動翻訳で用いられる統計量を、学習者の翻訳文評価に用いる可能性について検討を行った。英語らしさの指標、翻訳らしさの指標を考案し、それぞれについて翻訳文の客観評価への適用可能性について調べた。翻訳らしさの指標について、人の主観評価における訳質レベルに対し規則性を示した。この結果、この指標は翻訳評価に有用性を示した。

分析では、正解文との絶対評価を前提に、その基準となる正解翻訳を用いて尺度化を計った。一方で、学習者の翻訳文同士の相对比较を考えると、翻訳らしさの指標で翻訳文を順位付けすることで可能になる。さらに翻訳モデルから正解文の翻訳確率に代わる特徴量を算出し、それを基に得られた翻訳らしさの指標が利用できれば、正解文を事前に準備することなく、翻訳文を評価できることも考えられる。従来の客観尺度化の議論では、テスト文に対する正解は必須であるとされていた。本研究で考案した評価尺度の利用は、評価基準となる正解文を必ずしも必要としない評価法としての可能性に道を開いた。

謝辞

日頃より統計的翻訳に関するご指導をいただき隅田室長をはじめとする ATR 音声言語コミュニケーション研究所の皆様にご感謝致します。

参考文献

- [1] 和泉絵美, 内元清貴, 井佐原均 “日本人英語の通じやすさに関する研究”, 言語処理学会 第12回年次大会, S1-4 pp16-19, 2006
- [2] 山本誠一, 菅谷史昭, 安田圭志, 隅田英一郎, “音声翻訳技術開発の経験に基づく外国語能力評価法の提案”, 電子情報通信学会技術報告書, pp30-31, 2003
- [3] 安田圭志, 隅田英一郎, 山本誠一, 柳田益造, 前川喜久雄, 菅谷史昭 “英語コミュニケーション能力の自動測定技術の提案”, 情報処理学会研究報告 pp65-70, 2003
- [4] 坂田浩亮, 新保仁, 松本裕治, “コーパスを用いた言語習得度の推定”, 情報処理学会研究報告 pp113-119, 2007