

An Experimental Study of Automatic Metrics for Machine Translation Evaluation

K.Sourabh¹, S.M Aaqib², Minu Bala³

¹Assistant Professor, Dept. of Computer Science, GGM Science College, Jammu

²Assistant Professor, Dept. of Computer Science, Amar Singh Science College, Srinagar

³Sr. Assistant Professor Dept. of Computer Science GGM Science College

ABSTRACT

Machine Translation has gained popularity over the years and has become one of the promising areas of research in computer science. Due to a consistent growth of internet users across the world information is now more versatile and dynamic available in almost all popular spoken languages throughout the world. From Indian perspective importance of machine translation become very obvious because Hindi is a language that is widely used across India and whole world. Many initiatives have been taken to facilitate Indian users so that information may be accessed in Hindi by converting it from one language to other. In this paper we have studied various available automatic metrics that evaluate the quality of translation correlation with human judgments. . Also we have tested various automatic metrics available for their performance against human evaluation.

Keywords: Machine Translation, Evaluation, Bleu, Meteor, Ter, GTM.

I. INTRODUCTION

Machine translation is a sub-field of computational linguistics that investigates the use of software to translate text or speech from one language to another. As internet is now flooded with multilingual information for global community, research and development giving space to Machine Evaluation plays a major role in the field of Natural Language Processing. Many MT tools like Google translate, Bing, Systran SDL etc are providing online services to translate text from one language to another. Evaluation done manually (by humans) is the most reliable way for evaluating MT systems but it is subjective, expensive, time-consuming and involves human labor that cannot be reused. In general, evaluation can be understood as judgment on the value of a public intervention with reference to defined criteria of this judgment. Automatic MT evaluation metrics play a prominent role in the evaluation of MT systems. Many automatic measures have been proposed to facilitate fast and cheap evaluation of MT systems, the most widely used of which is BLEU, NIST METEOR etc. For Hindi language evaluation METEOR-Hindi is one of the promising metrics which has gained

popularity. The measure of evaluation for metrics is correlation with human judgment. This is generally done at two levels, at

the sentence level, where scores are calculated by the metric for a set of translated sentences, and then correlated against human judgment for the same sentences. In this paper we have studied various automatic metrics available which correlate with human judgment.

II. MANUAL TRANSLATION EVALUATION

Evaluation plays a very important role in examining the quality of MT output. Manual evaluation is very time consuming and prejudiced, hence use of automatic metrics is made most of the times. Some parameters taken into consideration for manual evaluation are *Rating Adequacy fluency ranking and post editing*.

Denkowski & Lavie, 2010 [2] Snover et al., 2006[3]

III AUTOMATIC TRANSLATION EVALUATION

Due to the high costs, lack of repeatability, subjectivity, and slowness of evaluating machine translation output using human judgments; automatic machine translation evaluation metrics were developed. Automated MT evaluation metrics are fast, scalable, and consistent which makes them very efficient to use but most of the times not reliable. Automated MT metric needs to correlate quality with respect to human translator, and produce reliable results for similar translations of the same content. Automated measures judge the output (candidate text) of a MT system against reference text. There are a number of automatic MT evaluation metrics: WER, TER, BLEU, NIST, METEOR, GTM and the list go on. Mostly all automatic metrics are based one of the following methods to calculate scores.

- **Edit Distance Based:** Number of insertions, deletions and substitutions that are being made to change candidate into reference are counted
- **Precision Based:** Total numbers of matched unigrams are divided by the total length of candidate
- **Recall Based:** Total number of matched unigrams is divided by the total length of reference
- **F-measure Based:** Both precision and recall scores are used collectively

Problem with Precision and Recall method

Source: Teachers are responsible for institutional development

Reference: Teachers must account for overall institutional growth

- Precision correct / output-length = $3/6 = 50\%$
- Recall correct / reference-length = $3/7 = 43\%$
- F-measure $\text{precision} \times \text{recall} / (\text{precision} + \text{recall})/2 = .5 \times .43 / (.5 + .43)/2 = 46\%$

Example 2:

SYSTEM A: “Teachers are responsible for institutional development”

Reference: “Teachers must account for **overall** institutional growth”

SYSTEM B: “for **overall** institutional growth Teachers must account”

Method	System A	System B
Precision	50%	100%
Recall	43%	100%
f-measure	46%	100%

Table 1

It can be seen that for two sentences with almost same meaning system B achieves 100% score whereas system A shows variations.

Automatic machine translation evaluation started with the introduction of BLEU then followed by NIST, GTM, ROUGE, CIDEr METOER and many others like [4], Blanc Ter Rose Amber Lepor Port and Meteor Hindi. Few of the metrics are studied below.

A. *Word error rate (WER)*

derived from the Levenshtein distance, word error rate can be computed as: $WER = \frac{S+D+I}{N}$ where $N = S+D+C$ where S is the number of substitutions, D is the number of deletions, I is the number of insertions, C is the number of the corrects, N is the number of words in the reference ($N = S+D+C$) [5]

B. *Translation Error Rate (TER)*

An error metric for machine translation that measures the number of edits required to change a system output into one of the references. [Snover, M. (2006)].[3]

C. *GTM (General text Matcher)*

Turian et al the evaluation score is obtained by sharing corresponding words between MT output and mentioned output, Not only on precision and recall but it is also based upon harmonic mean of both, known as F-measure calculated as

$$F - Measure = \frac{2PR}{P + R}$$

D. *BLEU (Bilingual Evaluation Understudy)*

Proposed by, Papineni in 2000 [6] the metric is one of the most popular in the field. The fundamental idea behind the metric is that "the closer a machine translation is to a professional human translation, the better it is" Papineni et al.

(2002). N-grams in the candidate translation are matched with n-grams in the reference text, where 1-gram (unigram) is a token and a bigram assessment would be each word pair. The comparison is made despite of word order. BLEU is not perfect, but offers five convincing benefits: [7]

- Calculation is quick and inexpensive
- Easy to understand
- Language independent
- Correlates highly with human evaluation
- Widely adopted

N-gram precision, Clipping and Brevity Penalty are main components of BLEU [8]

BLEU uses tailored n-gram precision a brevity penalty is introduced to compensate difference in the length of candidate and reference translations. Since the precision of 4-gram is many times 0, the BLEU score is generally computed over the corpus than on the sentence level. [4] Many up gradations have been made on the basic BLEU like Smoothed BLEU, BLEU deconstructed etc. to offer enhanced results.

Score calculation method for Blue can be:

$$BLEU = BP \cdot \exp \left(\sum_{n=1}^N W_n \log P_n \right)$$

$$BP = \begin{cases} 1 & (if c > r) \\ e^{(1-\frac{r}{c})} & (if c \leq r) \end{cases}$$

BP (Brevity Penalty), N is length of N grams used to compute P_n and P_n Modified

n gram precision

E. *NIST (National Institute of Standards and Technology)*

[Doddington 2002] **NIST** A modification of BLEU has been adopted by NIST for MT. Based on the score of Bleu, attempt is made to compute particular n-gram's usefulness i.e. how informative an n-gram is candidate text by giving it more weight depending upon its rareness. Instead of n-gram precision the information gain from each n-gram is taken into account Additionally, BP(Brevity Penalty) calculation varies somewhat as little disparity in translation text length don't affect the general score as much as in BLEU. It uses Arithmetic mean rather than geometric mean [9]. Bad correlation on sentence level with respect to human judgment still remains a problem.

F. *The METEOR (Metric for Evaluation of Translation with Explicit Ordering)*

[Satanjeev Banerjee, Lavie (2005)][10] for evaluation of machine translation output is based on the harmonic mean of unigram precision and recall, with recall weighted higher than precision. It can also make use of features such as stemming and synonymy matching which are not present in other metrics. Contrary to BLEU which aims to achieve good

correlation with human judgment at the corpus level, METEOR was designed to produce a good correlation at the sentence or segment level.

METEOR addresses several limitations in IBM's BLEU metric. METEOR supports not only matching between words that are identical, but can also match words that are simple morphological variants and synonyms of each other. The results reported by [banerjee et al] demonstrate that all of the individual components included within METEOR contribute to improved correlation with human judgments. In particular, METEOR is shown to have statistically significant better correlation compared to unigram-precision, unigram recall and the harmonic F1 combination of the two.

Score calculation in METEOR

$$F_{mean} = \frac{10PR}{R + 9P}$$

Where First unigram precision (P) is computed as the ratio of the number of unigrams in the system translation that are mapped (to unigrams in the reference translation) to the total number of unigrams in the system translation. Similarly, unigram recall (R) is computed as the ratio of the number of unigrams in the system translation that are mapped (to unigrams in the reference translation) to the total number of unigrams in the reference translation.

In order to compute this penalty, unigrams are grouped into the fewest possible chunks, where a chunk is defined as a set of unigrams that are adjacent in the hypothesis and in the reference. The longer the adjacent mappings between the candidate and the reference, the fewer chunks there are. A translation that is identical to the reference will give just one

chunk. The penalty p is computed as follows $P = 0.5 * \left(\frac{\#Chunks}{\#UnigramsMatched} \right)^3$

The final score for a segment is calculated as M below. $M = F_{mean}(1 - P)$

G. *METEOR-Hindi*

Ankush Gupta et.al [11] developed METEOR-Hindi, an automatic evaluation metric for a machine translation system where the target language is Hindi. METEOR-Hindi is a modified version of the metric METEOR, containing features specific to Hindi. Appropriate changes are made to METEOR's alignment algorithm and the scoring technique. METEOR, does not support Hindi by default, as it requires Hindi specific tools for computing synonyms, stem words, etc. additional modules listed below are added to METEOR to make well-organized for Hindi.

- Local Word Group (LWG) match
- Part-of-Speech (POS) and Clause match

METEOR-Hindi achieved high correlation with human judgments significantly outperforming BLEU.

III EXPERIMENTAL RESULTS

In order to evaluate various metrics we have performed tests on various sentences taken from subtitle files of movies. The Hindi sentences and their equivalent English translated sentences are taken from Hindi and English subtitles. The English

sentence is taken as candidate. Example given below with an assumption that source sentence from subtitle file is human evaluated and is correct

Source Hindi from movie “Taare Zameen Par” Hindi Subtitle

हर बच्चे की अपनी खूबी होती है, अपनी काबिलियत होती है, अपनी चाहत होती है

Source English from “Taare Zameen Par” English Subtitle

Every child has his own capabilities, his own desires, his own dreams.

Ref 1 (Google MTS 1): Every child has his own talent, his ability, he wants to be

Ref 2(Bing MTS 2): Every child has its own grasses, its prove, the desire

Ref 3(SDL MTS 3): Every child is their abilities, its spoke, are its aid

METRIC	Document	MTS 1	MTS 2	MTS 3
		Ref Mean	Ref Mean	Ref Mean
BLEU	1	0.26	0.35	0.24
	2	0.30	0.41	0.32
METEOR	1	0.71	0.61	0.58
	2	0.69	0.62	0.49
GTM	1	0.64	0.59	0.51
	2	0.58	0.53	0.42
TER	1	0.34	0.42	0.41
	2	0.45	0.47	0.38

Table 2

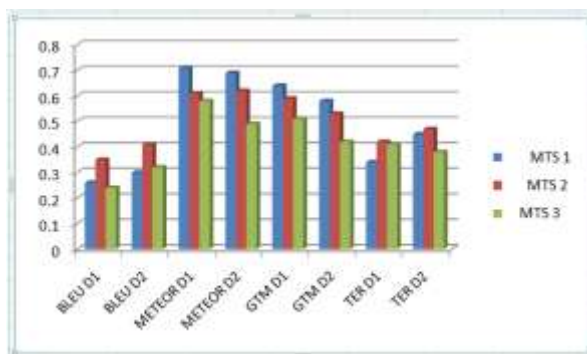


Figure 1

Figure: 1 shows comparison of four metrics under evaluation for single reference. It can be observed that METEOR performs best followed by GTM as compared to BLEU and TER. Also MT system 1 performs best in case of METEOR and GTM where as MT system 2 performs best in case of BLEU and TER.

IV Problems with BLEU/NIST metric

Reported by [Xingyi Song et.al] [7] a short document or sentence, there is a high probability of obtaining zero tri-gram or 4-gram precision, which makes the overall BLEU score equal zero due to the use of geometric mean..BLEU shows good performance for corpus level comparisons over which a high number of ngram matches exist. However, at a sentence-level the n-gram matches for higher n rarely occur [12]. As a result, BLEU performs poorly when comparing individual sentences. [Xinlei Chen.et.al]

BLEU supports multiple references, which makes it hard to obtain an estimate of recall. Therefore, recall is replaced by the BP, but BP is a poor substitute for recall.

BLEU with only uni-gram precision has the highest adequacy correlation (0.87), while adding higher order n-gram precision factors decreases the adequacy correlation and increases fluency

Reported by [ankush etal] BLEU is not an appropriate metric for English-Hindi evaluation because of Meaningless Sentence-level Score, Only Exact Matches (morphological variants not considered, Lack of recall and Geometric Averaging of n-grams.

V CONCLUSION

This paper we have presented different approaches of evaluation. Here we provided the performance results of automatic evaluation metrics and compared their performances for various translation systems. It can be concluded that METOER performs best as compared to BLEU, GTM and TER. While well known, shortcomings have been noted in BLEU as of late, most remarkably the absence of solid sentence-level scores. Further, it isn't appropriate for assessment of English-Hindi MT frameworks in view of the properties of Hindi, for example, rich morphology and relative free word orderings. With a specific end goal to beat the shortcomings of BLEU, a few measurements were proposed, for example, METEOR, GTM, TER. METEOR is the most reasonable for assessment of English-Hindi MT, as it offers immense flexibility in encoding parameters that show nature of understanding the translated text.

Since automatic evaluation metrics do not always correspond to human judgment. This is necessary to resolve whether future algorithms are actually improving, or whether they are merely over fitting to a specific metric.

REFERENCES

- [1] Philipp Koehn, Christof Monz ,”*Manual and Automatic Evaluation of Machine Translation between European Languages*” School of Informatics University of Edinburgh ,Department of Computer Science Queen Mary, University of London. Proceeding StatMT '06 Proceedings of the Workshop on Statistical Machine Translation Pages 102-121
- [2] Michael Denkowski and Alon Lavie Language ,”*Choosing the Right Evaluation for Machine Translation: an Examination of Annotator and Automatic Metric Performance on Human Judgment Tasks*” Proceedings of the Ninth Biennial Conference of the Association for Machine Translation in the Americas <https://www.cs.cmu.edu/~mdenkows/pdf/mteval-amta-2010.pdf>
- [3] Matthew Snover Bonnie Dorr Richard Schwartz, Linnea Micciulla, and John Makhoul, “*A Study of Translation Edit Rate with Targeted Human Annotation*” Proceedings of association for machine translation in the Americas, pp 223-231.
- [4] Aditi Kalyani, Hemant Kumud Shashi Pal Singh Ajai Kumar,” *Assessing the Quality of MT Systems for Hindi to English Translation*” International Journal of Computer Applications (0975 – 8887) Volume 89 – No 15, March 2014
- [5] Klakow, Dietrich; Jochen Peters (September 2002). "Testing the correlation of word error rate and perplexity". *Speech Communication*. 38 (1-2): 19–28. doi:10.1016/S0167-6393(01)00041-3. ISSN 0167-6393
- [6] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu,” *BLEU: a Method for Automatic Evaluation of Machine Translation* “. Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics (ACL), Philadelphia, July 2002, pp. 311-318
- [7] Jason Brownlee “*A Gentle Introduction to Calculating the BLEU Score for Text in Python* “.November 20, 2017 in Natural Language Processing” Online <https://machinelearningmastery.com/calculate-bleu-score-for-text-python/>
- [8] Xingyi Song and Trevor Cohn and Lucia Specia,” *BLEU deconstructed: Designing a Better MT Evaluation Metric*” University of Sheffield Department of Computer Science Proceedings of the 14th International Conference on Intelligent Text Processing and Computational Linguistics (CICLING)
- [9] Doddington, George. (2002),”*Automatic evaluation of machine translation quality using n-gram co-occurrence statistics*”.138-145. 10.3115/1289189.1289273
- [10] Satanjeev Banerjee Alon Lavie ,”*METEOR: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments*” Institute Language Technologies Institute Carnegie Mellon University. Proceedings of Workshop on Intrinsic and Extrinsic Evaluation Measures for MT and/or Summarization at the 43rd Annual Meeting of the Association of Computational Linguistics (ACL-2005), Ann Arbor, Michigan, June 2005.
- [11] Ankush Gupta and Sriram Venkatapathy and Rajeev Sangal ,” *METEOR-Hindi : Automatic MT Evaluation Metric for Hindi as a Target Language*”. Language Technologies Research Centre, IIIT-Hyderabad, Hyderabad, India.

Proceedings of ICON-2010:8th International conference on Natural language processing, Macmillan Publishers, India.

- [12] Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam Saurabh Gupta, Piotr Dollar, C. Lawrence Zitnick, "Microsoft COCO Captions: Data Collection and Evaluation Server". CoRR 2015 Vol: abs/1504.00325