

# Multilingual NLP

## Lab 3 – How Reliable is MT Evaluation?

AN Ji

December 17, 2024

### 1 Re-evaluating the role of BLEU in machine translation research

#### 1.1 Implement the BLEU metric.

```
1 import numpy as np
2 from collections import Counter
3
4 # n-gram generator
5 def ngrams(seq, n):
6     """generate n-grams from a sequence using iterator & sliding window
7     Parameters -- seq: str, n: int
8     Returns -- generator of tuples or [] if seq size < n"""
9     seq = seq.split()
10    # extreme case: n > seq size
11    if n > len(seq):
12        return []
13    # general case: n <= seq size
14    seq_generator = iter(seq)
15    history = []
16    # initialize with first (n-1) items
17    while n > 1:
18        history.append(next(seq_generator))
19        n -= 1
20    # add next (1st item in remaining seq)
21    for item in seq_generator:
22        history.append(item)
23        yield tuple(history)
24        # delete oldest & move forward if any
25        del history[0]
26
27
28 def compute_corpus_bleu(refs, hyps, n=4, floor=1e-10):
29     """respectively aggregate sentence-level matching ngrams and individual sentence lengths
30     Params -- refs: list of str, hyps: list of str, n: int, floor: float
31     Returns -- float"""
32     hyps_ngrams, bingo_ngrams = [0]*n, [0]*n
33     refs_len, hyps_len = 0, 0 # rather than refs|hyps.split()
34
35     for ref, hyp in zip(refs, hyps):
36         # aggregate sentence length 1 by 1 to get corpus length
37         refs_len += len(ref.split())
38         hyps_len += len(hyp.split())
39
40     for i in range(1, n+1):
41         # aggregate 1 2 3 4-grams for each sentence pair
42         ref_igrams, hyp_igrams = list(ngrams(ref, i)), list(ngrams(hyp, i))
43         if hyp_igrams:
44             bag_ref, bag_hyp = Counter(ref_igrams), Counter(hyp_igrams)
45             bingo = sum((
46                 min(bag_hyp[ngram], bag_ref.get(ngram, 0))
47                 for ngram in bag_hyp
48             ))
49             bingo_ngrams[i-1] += bingo
50             hyps_ngrams[i-1] += sum(bag_hyp.values())
51
```

```

52 # set penalty
53 bp = 1 if hypos_len > refs_len else np.exp(1 - refs_len / hypos_len)
54
55 # calculate & smooth n-gram precisions
56 precisions = [
57     (bingo_ngrams[i] / hypos_ngrams[i]) if bingo_ngrams[i] > 0 else floor # no need to check
58     hypos_ngrams[i]
59     for i in range(n)
60 ]
61
62 return bp * np.exp(np.mean(np.log(precisions)))

```

## 1.2 Using the WMT'15 test sets, evaluate the performance of mBart and MarianMT. These two models can be easily used with the HuggingFace API. What can you conclude?

```

1 from bs4 import BeautifulSoup
2 from tqdm.notebook import tqdm
3 from transformers import MBartForConditionalGeneration, MBart50TokenizerFast, MarianMTModel,
4   MarianTokenizer
5
6 src_fr = '/content/drive/MyDrive/Colab Notebooks/MT/newsdiscusstest2015-fren-src.fr.sgm'
7 ref_en = '/content/drive/MyDrive/Colab Notebooks/MT/newsdiscusstest2015-fren-ref.en.sgm'
8
9 def sgm2dict(src_path, ref_path):
10     """parse .sgm files & extract sentences to dict"""
11     corpus = dict()
12     for filepath in [src_path, ref_path]:
13         with open(filepath, "r") as sgm:
14             soup = BeautifulSoup(sgm, "html.parser")
15             sentences = [
16                 seg.get_text().strip()
17                 for seg in soup.find_all("seg")
18             ]
19             key = filepath[-10:-4]
20             corpus[key] = sentences
21     return corpus
22
23 corpus = sgm2dict(src_fr, ref_en)
24 sample = corpus['src.fr'][:100]
25
26
27 def mbart_trans(src, batch_size=16, src_lang='fr_XX', tgt_lang='en_XX'):
28     """feed src sentences to mbart, return a list of translations"""
29     # Load mBART model & tokenizer
30     mbart = 'facebook/mbart-large-50-many-to-many-mmt'
31     mbart_model = MBartForConditionalGeneration.from_pretrained(mbart)
32     mbart_tokenizer = MBart50TokenizerFast.from_pretrained(mbart, src_lang=src_lang)
33
34     mb_hyps = []
35     with tqdm(total=len(src), desc="Translating batches", unit="sentence") as pbar:
36         for i in range(0, len(src), batch_size):
37             batch = src[i:i+16]
38             inputs = mbart_tokenizer(
39                 batch,
40                 return_tensors="pt",
41                 padding=True,
42                 truncation=False
43             )
44             pbar.set_postfix_str(f"Processing batch {i // batch_size + 1}")
45             outputs = mbart_model.generate(
46                 **inputs,
47                 forced_bos_token_id=mbart_tokenizer.lang_code_to_id[tgt_lang]
48             )
49             outputs_decoded = mbart_tokenizer.batch_decode(
50                 outputs,
51                 skip_special_tokens=True
52             )
53             pbar.update(len(batch))
54             mb_hyps.extend(outputs_decoded)
55
56     return mb_hyps
57

```

```

58
59 def marian_trans(src, batch_size=32):
60     """feed src sentences to marianMT, return a list of translations"""
61     # Load MarianMT model & tokenizer
62     marian = 'Helsinki-NLP/opus-mt-fr-en'
63     marian_model = MarianMTModel.from_pretrained(marian)
64     marian_tokenizer = MarianTokenizer.from_pretrained(marian)
65
66     mr_hyps = []
67     with tqdm(total=len(src), desc="Translating batches", unit="sentence") as pbar:
68         for i in range(0, len(src), batch_size):
69             batch = src[i:i+batch_size]
70             inputs = marian_tokenizer(
71                 batch,
72                 return_tensors="pt",
73                 padding=True,
74                 truncation=False
75             )
76             pbar.set_postfix_str(f"Processing batch {i // batch_size + 1}")
77             outputs = marian_model.generate(**inputs)
78             outputs_decoded = [
79                 marian_tokenizer.decode(output, skip_special_tokens=True)
80                 for output in outputs
81             ]
82             pbar.update(len(batch))
83             mr_hyps.extend(outputs_decoded)
84
85     return mr_hyps
86
87
88 mbart_hyps = mbart_trans(sample)
89 marian_hyps = marian_trans(sample)
90 sample_dict = {
91     'src.fr': sample,
92     'ref.en': corpus['ref.en'][:100],
93     'mbart_hyp.en': mbart_hyps,
94     'marian_hyp.en': marian_hyps
95 }
96
97 # ds = Dataset.from_dict(sample_dict)
98
99 mbart_bleu = compute_corpus_bleu(sample_dict['ref.en'], sample_dict['mbart_hyp.en'])
100 marian_bleu = compute_corpus_bleu(sample_dict['ref.en'], sample_dict['marian_hyp.en'])
101
102 print(f"MBart BLEU: {mbart_bleu}")
103 print(f"MarianMT BLEU: {marian_bleu}")
104
105 [Out]:
106 MBart BLEU: 0.32006793786190674
107 MarianMT BLEU: 0.34034764664241274

```

The above test is based on French-English translation using WMT'15 newsdiscusstest2015-fren test set. I only choose the first 100 of all 1500 sentences due to lack of enough RAM on Colab.

The results indicate strong overlaps between reference sentences and both mBart's and MarianMT's translations. With very close BLEU scores, both models well preserve n-gram precision and length similarity on the sample. In practice, MarianMT works much faster than mBart, achieving slightly higher score by 0.03.

**1.3 As noticed by [1], BLEU places no explicit constraints on the order that matching n-grams occur in. It is therefore possible, given a sentence, to generate many new sentences with at least as many n-gram matches by permuting words around bigram mismatches. Explain on an example why such permutations will never decrease the BLEU score.**

From the implementation of our BLEU function we can see that BLEU is a surface and static evaluation that ONLY takes into account the matching n-grams between two sentences, regardless of their word order. In other words, the difference between BLEU scores of two sentences can hardly reflect how semantically and syntactically correct two sentences are, since there's no penalty designed for

the purposes of these aspects.

For example, as an English speaker can easily tell, the following rearranged hypotheses are far from correct, but they all have excellent BLEU scores based on our implementation (with  $n = 2$ ). This case illustrates the uncorrelativeness between BLEU and human judgements, as well as its limit as an evaluation metric.

```
1 ref = "the quick brown fox jumps over the lazy dog and runs away"
2 hyp_1 = "the quick brown dog jumps over the lazy fox and runs away"
3 hyp_2 = "the lazy dog jumps over the quick brown fox and runs away"
4
5 hyp_1 BLEU: 0.7977240352174656
6 hyp_2 BLEU: 0.9045340337332909
```

#### 1.4 Given a sentence with $n$ words and $b$ bigram mismatches, how many sentences can you generate with this principle. Compute the number of sentences you will obtain on the WMT'15 test set.

Since bigrams are not necessarily distinct within a hypothesis sentence, we CANNOT directly calculate the number of matching bigrams given  $b$  mismatches. So here we switch to considering there are  $b$  matching bigrams.

In this updated case, if a given hypothesis has  $k$  words, from Callison-Burch et al. [1] we know that the number of identically scored sentences generated from the given hypothesis can be computed using the following formula:

$$\text{Number of identically scored sentences} = (k - b)!$$

```
1 # Number of permuted sentences generated from both model hypotheses
2 # based on our sample of size 100
3
4 def fake_hyps(refs, hyps):
5     """given an mbart/marian hypothesis, calculate the number of permuted sentences that score
6         identically"""
7     total = 0
8     for ref, hyp in zip(refs, hyps):
9         ref_bigrams = set(ngrams(ref, 2))
10        hyp_bigrams = set(ngrams(hyp, 2))
11        b = len(hyp_bigrams & ref_bigrams) # number of matched bigrams
12        k = len(hyp.split()) # number of words in hyp
13        total += math.factorial(k - b) # (k-b)! of current hyp
14    return f'{total:.2e}'
15
16 fake_hyps(refs, mbart_hyps)
17
18 [Out]:
19 '1.03e+40'
20
21 fake_hyps(refs, marian_hyps)
22
23 [Out]:
24 '8.23e+33'
```

We will obtain 1.03e+40 permuted sentences for the 100 hypotheses of mBart and 8.23e+33 for those of MarianMT.

#### 1.5 Why does this result question the use of BLEU as an evaluation metric.

As mentioned in 1.3, there are lots of syntactic and semantic aspects that BLEU is incapable to account for. On the one hand, it is possible to have sentences which receive identical Bleu scores but are judged by humans to be worse; on the other hand, it is also possible to have a higher Bleu score without any genuine improvement in translation quality [1].

## 1.6 SACREBLEU is an implementation of BLEU that aims to provide “hassle-free computation of shareable, comparable, and reproducible BLEU scores”. Evaluate the two previous systems using SACREBLEU. What can you conclude?

```
1 !pip install sacrebleu
2
3 sacre_mbart = sacrebleu.corpus_bleu(mb_hyps, [refs], lowercase=True)
4 sacre_marian = sacrebleu.corpus_bleu(mr_hyps, [refs], lowercase=True)
5
6 print(f"MBart sacreBLEU: {sacre_mbart}")
7 print(f"MarianMT sacreBLEU: {sacre_marian}")
8
9 # Recall
10 # MBart BLEU: 0.32006793786190674 (multiplied by 100 => 32.01)
11 # MarianMT BLEU: 0.34034764664241274 (multiplied by 100 => 34.03)
12
13 [Out]:
14 MBart sacreBLEU: BLEU = 36.17 65.2/41.7/29.3/21.5 (BP = 1.000 ratio = 1.027 hyp_len = 1732 ref_len
    = 1686)
15
16 MarianMT sacreBLEU: BLEU = 38.13 67.2/44.7/31.3/22.5 (BP = 1.000 ratio = 1.023 hyp_len = 1725
    ref_len = 1686)
```

We notice that both SACREBLEU scores are **4 points higher** than our BLEU scores. This may be explained by various strategies used in both methods.

First, SACREBLEU apply normalization techniques to both references and hypotheses, such as removing special characters and punctuation, and we also pass `lowercasing=True` in the scoring function, which we did not apply in our BLEU implementation, and this does allow SACREBLEU to recognize more matching bigrams (which we have verified by switching the `lowercasing` setting) and return higher scores.

A second reason may be that SACREBLEU has its own tokenization methods, while we only used `.split()` to produce tokens based on spaces, which is not as robust as the former. For example if token HIV occurs in a reference, then token HIV. in a hypothesis wouldn't be taken into account when counting unigram matches using our BLEU function, while the normalization step allows SACREBLEU to tokenize HIV. as HIV, thus counting it as a matching unigram and yielding higher precision.

## 1.7 Using SACREBLEU and your own implementation of BLEU compute the score achieved:

- when considering the “raw” translation hypotheses and references;
- when the translation hypotheses and references have been tokenized in subword units;
- when the translation hypotheses and references have been tokenized in characters (this amounts to adding a space between each character of the references and of the translation hypotheses).

How can you explain these results?

```
1 def tokenize_corpus(sentences, method):
2     if method == 'subword':
3         tokenizer = BertTokenizer.from_pretrained("bert-base-multilingual-cased")
4         updated_sentences = [' '.join(tokenizer.tokenize(sent)) for sent in sentences]
5     elif method == 'character':
6         updated_sentences = [' '.join(list(sent)) for sent in sentences]
7     else: updated_sentences = sentences
8     return updated_sentences
9
10 original_sample = {
11     'refs': refs,
12     'mbart': mb_hyps,
13     'marian': mr_hyps
14 }
15
16 methods = ['raw', 'subword', 'character']
17
```

```

18 # Dict for tokenized sample
19 tokenized_sample = {}
20 for method in methods:
21     for k, v in original_sample.items():
22         tokenized_sample[f'{method}_{k}'] = tokenize_corpus(v, method)
23
24 scores = {
25     ('mbart', 'bleu'): [],
26     ('mbart', 'sacrebleu'): [],
27     ('marian', 'bleu'): [],
28     ('marian', 'sacrebleu'): [],
29 }
30
31 # bleu
32 for model in ['mbart', 'marian']:
33     for method in methods:
34         scores[(model, 'bleu')].append(
35             float("{:.2f}".format(
36                 compute_corpus_bleu(
37                     tokenized_sample[f'{method}_refs'],
38                     tokenized_sample[f'{method}_{model}']
39                 ) * 100 # for convenient comparison
40             ))
41         )
42
43 # sacrebleu
44 for model in ['mbart', 'marian']:
45     for method in methods:
46         scores[(model, 'sacrebleu')].append(
47             float("{:.2f}".format(
48                 sacrebleu.corpus_bleu(
49                     tokenized_sample[f'{method}_{model}'],
50                     [tokenized_sample[f'{method}_refs']],
51                     # tokenize='none'
52                 ).score
53             ))
54         )
55
56 scores
57
58 [Out]:
59 {('mbart', 'bleu'): [32.01, 38.15, 63.37],
60 ('mbart', 'sacrebleu'): [35.05, 43.34, 63.37],
61 ('marian', 'bleu'): [34.03, 40.33, 65.95],
62 ('marian', 'sacrebleu'): [37.02, 46.21, 65.95]}

```

|        |           | raw          | subword      | character |
|--------|-----------|--------------|--------------|-----------|
| model  | metric    |              |              |           |
| mbart  | bleu      | 32.01        | 38.15        | 63.37     |
|        | sacrebleu | <b>35.05</b> | <b>43.34</b> | 63.37     |
| marian | bleu      | 34.03        | 40.33        | 65.95     |
|        | sacrebleu | <b>37.02</b> | <b>46.21</b> | 65.95     |

Table 1: Scores under various tokenization strategies using default settings of SACREBLEU

|        |           | raw          | subword      | character |
|--------|-----------|--------------|--------------|-----------|
| model  | metric    |              |              |           |
| mbart  | bleu      | 32.01        | 38.15        | 63.37     |
|        | sacrebleu | <b>32.01</b> | <b>38.15</b> | 63.37     |
| marian | bleu      | 34.03        | 40.33        | 65.95     |
|        | sacrebleu | <b>34.03</b> | <b>40.33</b> | 65.95     |

Table 2: Scores under various tokenization strategies using `tokenize='none'` in SACREBLEU

From Table 1 we notice that when tokens become more and more atomized, both BLEU and SACRE-BLEU scores will be raised considerably, and they become identical at character-level tokenization.

In SACREBLEU we also tried disabling entirely the default tokenization method<sup>1</sup> (`tokenize='none'`). Table 2 shows how SACREBLEU scores change accordingly.

It is interesting that SACREBLEU returns the same scores as those of BLEU on all situations. Based on this observation we can conclude that without more tokenization, both metrics rely on spaces as token boundaries and thus leading to identical scores. The application of default 13a tokenization (Table 1) implies more matching n-grams, obviously raising scores by about 3 points on raw cases and about 5 points on subword cases for both models' hypotheses.

## References

- [1] Callison-Burch, C., Osborne, M., and Koehn, P. (Apr. 2006). "Re-evaluating the Role of Bleu in Machine Translation Research". In: *11th Conference of the European Chapter of the Association for Computational Linguistics*. Ed. by McCarthy, D. and Wintner, S. Trento, Italy: Association for Computational Linguistics, pp. 249–256. URL: <https://aclanthology.org/E06-1032>.
- [2] Post, M. (Oct. 2018). "A Call for Clarity in Reporting BLEU Scores". In: *Proceedings of the Third Conference on Machine Translation: Research Papers*. Belgium, Brussels: Association for Computational Linguistics, pp. 186–191. URL: <https://www.aclweb.org/anthology/W18-6319>.

---

<sup>1</sup>SACREBLEU's tokenization method defaults to '13a'. See [their GitHub documentation](#).