

# How Reliable is MT Evaluation?

He HUANG 22107447

October 29, 2023

## 1. Implement the Bleu metric.

```
1 from math import exp, log
2 import numpy as np
3 from collections import Counter
4
5 def ngrams(sequence, n):
6     """Generate n-grams from a sequence."""
7
8     sequence = sequence.split()
9     if len(sequence) < n:
10         return []
11
12     sequence = iter(sequence)
13     history = []
14
15     while n > 1:
16         history.append(next(sequence))
17         n -= 1
18     for item in sequence:
19         history.append(item)
20         yield tuple(history)
21         del history[0]
22
23 def brevity_penalty(ref, hyp):
24
25     """Brevity penalty"""
26     h = len(hyp.split())
27     r = len(ref.split())
28
29     if h > r:
30         return 1.
31     else:
32         return exp(1 - (h / r))
33
34 def bleu_score(ref, hyp, n):
35     """Implement BLEU score"""
36
37     m_p = []
38     weight = 1.0/n
39
40     for i in range(n):
41         # generate i-grams
42         ref_ngrams = Counter(ngrams(ref, i+1))
43         hyp_ngrams = Counter(ngrams(hyp, i+1))
```

```

45     # modified precision
46     common_ngrams = sum((min(hyp_ngrams[ngram], ref_ngrams.get(ngram, 0))
47                           for ngram in hyp_ngrams))
48     total = sum(hyp_ngrams.values())
49
50     if total == 0:
51         continue
52     else:
53         m_p.append(common_ngrams / total)
54
55     bp = brevity_penalty(ref, hyp)
56     # add epsilon to avoid log(0)
57     log_p = sum(weight * log((p_n) + 1e-10) for p_n in m_p)
58     bleu_compute = bp*np.exp(log_p)
59
60     return bleu_compute
61
62 def corpus_bleu(refs, hyps, n):
63     """Return the BLEU score in each translation and reference sample"""
64
65     bleu = 0
66     for _, (ref,hyp) in enumerate(zip(refs, hyps)):
67         bleu += bleu_score(ref, hyp, n)
68     return bleu / len(refs)

```

2. Using the WMT'15 test sets<sup>1</sup>, evaluate the performance of mBart and MarianMT. These two models can be easily used with the HuggingFace API<sup>2</sup>

```

1 import re
2 import html
3 import torch
4
5 en_fr_src_path = "test/newsdiscusstest2015-enfr-src.en.sgm"
6 en_fr_ref_path = "test/newsdiscusstest2015-enfr-ref.fr.sgm"
7
8 def sgm_to_text(filename):
9
10     """Read a SGM file and return the plain text."""
11
12     plain_text = []
13
14     with open(filename, "r") as sgm_file:
15         for line in sgm_file.readlines():
16             # convert html tags
17             line = html.unescape(line)
18             # for each line which starts with <seg>
19             if re.search(r"<seg id=.\+>.\+</seg>", line):
20                 pattern = re.compile(r"<seg id=\"\d+\">|</seg>")
21                 line = re.sub(pattern, "", line)
22                 line = line.replace("\n", "")
23                 plain_text.append(line)
24     return plain_text
25
26 src_en = sgm_to_text(en_fr_src_path)
27 ref_fr = sgm_to_text(en_fr_ref_path)

```

<sup>1</sup><https://statmt.org/wmt15/translation-task.html>

<sup>2</sup>[https://huggingface.co/docs/transformers/model\\_doc/mbart](https://huggingface.co/docs/transformers/model_doc/mbart) and [https://huggingface.co/docs/transformers/model\\_doc/marian](https://huggingface.co/docs/transformers/model_doc/marian)

```

28
29 !pip install transformers
30 !pip install sentencepiece
31 !pip install sacremoses
32
33 device = torch.device("cuda") if torch.cuda.is_available() else torch.device("
    cpu")
34
35 # mBART
36 # define the mBart model and tokenizer
37 from transformers import MBartForConditionalGeneration, MBart50TokenizerFast
38
39 mbart_name = "facebook/mbart-large-50-many-to-many-mmt"
40 mbart_model = MBartForConditionalGeneration.from_pretrained(mbart_name).to(
    device)
41 mbart_tokenizer = MBart50TokenizerFast.from_pretrained(mbart_name, src_lang="
    en_XX")
42
43 # MarianMT
44 # define the MarianMT model and tokenizer
45 from transformers import MarianMTModel, MarianTokenizer
46 mt_name = "Helsinki-NLP/opus-mt-en-fr"
47 mt_model = MarianMTModel.from_pretrained(mt_name).to(device)
48 mt_tokenizer = MarianTokenizer.from_pretrained(mt_name)
49
50 def translation(model, model_name, data, tokenizer):
51
52     """Translate the English sentences into French"""
53
54     translated = []
55
56     for sent in data:
57         inputs = tokenizer(sent, return_tensors="pt", truncation=True, max_length
            =100, padding=True).to(device)
58
59         if model_name == "mbart":
60             generated_tokens = model.generate(**inputs, forced_bos_token_id=tokenizer.
                lang_code_to_id["fr_XX"])
61             decode_fr = tokenizer.batch_decode(generated_tokens, skip_special_tokens=
                True)
62         else:
63             generated_tokens = model.generate(**inputs)
64             decode_fr = [tokenizer.decode(g, skip_special_tokens=True) for g in
                generated_tokens]
65
66         translated.append(decode_fr)
67     return translated
68
69 # Results : a list of sentences
70 mbart_en_fr = translation(mbart_model, "mbart", src_en, mbart_tokenizer)
71 mbart_en_fr = [''.join(sent) for sent in mbart_en_fr]
72
73 mt_en_fr = translation(mt_model, "mt", src_en, mt_tokenizer)
74 mt_en_fr = [''.join(sent) for sent in mt_en_fr]
75
76 # Compute BLEU score
77 mbart_en_fr_bleu = corpus_bleu(ref_fr, mbart_en_fr, 4)
78 mt_en_fr_bleu = corpus_bleu(ref_fr, mt_en_fr, 4)
79
80 print(f"The BLEU score of mBART is {mbart_en_fr_bleu}")

```

```
81 print(f"The BLEU score of MT is {mt_en_fr_bleu}")
```

	mBART	MarianMT
en->fr	0.2269	0.2918

MarianMT obtains a BLEU score of 0.2918 and mBART of 0.2269. MarianMT model is better at translating English sentences into French than mBART.

### 3. Explain on an example why such permutations will never decrease the Bleu score.

BLEU simply counts matching n-grams, it doesn't consider their order in the sentence. Let's take a reference sentence "She loves reading books." And the translation sentence could be "Reading she books loves" or "She books reading loves." All these sentences might have similar uigram (i.e. "She", "loves", "reading", "books"), since BLEU doesn't penalize word order, these two translations might have similar BLEU scores.

### 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.

There are  $(n - b)!$  sentences generated by this principle.

```
1 # compute permutation
2 import math
3
4 def permutation(references, hypothesis):
5     total_nb = 0
6
7     for i in range(len(references)):
8         ref = references[i]
9         ref_grams = ngrams(ref, 2)
10
11         hyp = hypothesis[i]
12         hyp_len = len(hyp.split())
13         hyp_grams = ngrams(hyp, 2)
14
15         common = set(hyp_grams).intersection(set(ref_grams))
16         common_len = len(common)
17         total_nb += math.factorial(hyp_len - common_len)
18     return format(total_nb, 'e')
```

	mBART	MarianMT
permutation	1.447390e+75	2.350574e+78

On the sentences from WMT'15 test set, there are 1.447390e+75 permutations for mBART and 2.350574e+78 for MarianMT.

### 5. Why this result question the use of Bleu as an evaluation metric.

In many languages, the word order and semantic meaning of a sentence is crucial for the evaluation of the quality of translation. A higher BLEU score does not necessarily mean a better translation. Because there are also scores for permuted gibberish.

6. `sacreBleu`<sup>3</sup> 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 import sacrebleu
4
5 # en->fr
6 # mBART
7 sacrebleu_mbart_en_fr = sacrebleu.corpus_bleu(mbart_en_fr,[ref_fr])
8 # MarianMT
9 sacrebleu_mt_en_fr = sacrebleu.corpus_bleu(mt_en_fr,[ref_fr])
10 print(f"mBART sacrebleu : {sacrebleu_mbart_en_fr}")
11 print("\n")
12 print(f"MarianMT sacrebleu : {sacrebleu_mt_en_fr}")
```

```
1 Output:
2
3 mBART sacrebleu : BLEU = 32.59 60.8/38.5/26.3/18.3 (BP = 1.000 ratio = 1.016
4     hyp_len = 28435 ref_len = 27975)
5
6 MarianMT sacrebleu : BLEU = 38.38 65.6/44.3/31.9/23.6 (BP = 0.998 ratio = 0.998
7     hyp_len = 27929 ref_len = 27975)
```

This score has been multiplied by 100, so if I follow my own implementation, the score is actually 0.33 for mBART and 0.38 for MarianMT. Compared to my score (0.2269 for mBART and 0.2918 for MT), the `sacreBleu` score is a bit higher with the same translations. In my implementation BLEU, I separated the sentences on space, `sacreBleu` takes another way to tokenization. This may affect the calculation by these two metrics.

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<sup>4</sup> ;
- 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 # Raw text
2 print(f"mBART BLEU score : {mbart_en_fr_bleu}")
3 print(f"MarianMT BLEU score : {mt_en_fr_bleu}")
4 print("\n")
5 print(f"mBART sacrebleu : {sacrebleu_mbart_en_fr}")
6 print(f"MarianMT sacrebleu : {sacrebleu_mt_en_fr}")
7
8 Output:
```

---

<sup>3</sup><https://github.com/mjpost/sacrebleu>

<sup>4</sup>You can use either one of the tokenizer provided by HuggingFace or train your own model, for instance on Europarl data.

```

9 mBART BLEU score : 0.22685186354146247
10 MarianMT BLEU score : 0.2918108984560062
11
12 mBART sacrebleu : BLEU = 32.59 60.8/38.5/26.3/18.3 (BP = 1.000 ratio = 1.016
    hyp_len = 28435 ref_len = 27975)
13 MarianMT sacrebleu : BLEU = 38.38 65.6/44.3/31.9/23.6 (BP = 0.998 ratio = 0.998
    hyp_len = 27929 ref_len = 27975)

```

```

1 # subword units
2 from transformers import AutoTokenizer
3
4 def subword_tokenization(sentences):
5     tokenizer = AutoTokenizer.from_pretrained("bert-base-multilingual-cased")
6     tokens = []
7
8     for sent in sentences:
9         tokenization = tokenizer.tokenize(sent)
10        subword_sent = ' '.join(tokenization)
11        tokens.append(subword_sent)
12    return tokens

```

```

1 # character level
2 def char_tokenization(sentences):
3     return [' '.join(list(sent)) for sent in sentences]

```

For mBART:

model	score	raw	subword	character
mBART	BLEU	22.69	36.82	64.46
	sacreBleu	32.59	45.23	64.06

For MarianMT:

model	score	raw	subword	character
MarianMT	BLEU	29.18	44.38	70.48
	sacreBleu	38.38	50.82	68.85

It is obvious that when the tokenisation unit becomes smaller, the scores of both BLEU and sacreBleu will increase. When we split sentences into smaller units (words to subwords and characters), they will have more opportunities to match more n-grams. In this way, the score will become higher.