



deeplearning.ai

# Sequence to sequence models

---

## Bleu score (optional)

If there are multiple great answers, how do you measure accuracy?  
The way this is done conventionally is through something called the Bleu score.

# Evaluating machine translation

French: Le chat est sur le tapis.

Reference 1: The cat is on the mat. ←

Reference 2: There is a cat on the mat. ←

MT output: the the the the the the the.

Precision:

Modified precision:

Bleu  
bilingual evaluation understudy

What Bleu does, is given a machine generated translation, it allows you automatically compute the score that measures how good is that machine translation.

# Evaluating machine translation

French: Le chat est sur le tapis. the word "the", gets credit up to twice.

Bleu  
bilingual evaluation understudy

→ Reference 1: The cat is on the mat.   
 2 appears

→ Reference 2: There is a cat on the mat.

→ MT output: the the the the the the the

Precision:  $\frac{7}{7}$

Modified precision:  $\frac{2}{7}$    
 ← Count<sub>clip</sub>("the")   
 ← Count("the")

# Bleu score on bigrams pairs of words appearing next to each other.

Example: Reference 1: The cat is on the mat. ←

Reference 2: There is a cat on the mat. ←

MT output: The cat the cat on the mat. ←

	Count	Count <sub>clip</sub>	
the cat	2 ←	1 ←	
cat the	1 ←	0	4
cat on	1 ←	1 ←	<hr/>
on the	1 ←	1 ←	6
the mat	1 ←	1 ←	
	↑		

# Bleu score on unigrams

Example: Reference 1: The cat is on the mat.

Reference 2: There is a cat on the mat.

→ MT output: The cat the cat on the mat. ( $\hat{y}$ )

$$p_1, p_2 = \underline{1.0}$$

$$p_1 = \frac{\sum_{unigram \in \hat{y}} \text{count}_{clip}(unigram)}{\sum_{unigram \in \hat{y}} \text{count}(unigram)}$$

*Handwritten notes:* "unigram" with an arrow pointing to the numerator's variable, and "count(unigram)" written below the denominator's term.

$$p_n = \frac{\sum_{ngram \in \hat{y}} \text{count}_{clip}(ngram)}{\sum_{ngram \in \hat{y}} \text{count}(ngram)}$$

*Handwritten notes:* "n-gram" with an arrow pointing to the numerator's variable, and "count(n-gram)" written below the denominator's term.

# Bleu details

$p_n$  = Bleu score on n-grams only

$p_1, p_2, p_3, p_4$

Combined Bleu score:

$$BP \exp\left(\frac{1}{4} \sum_{n=1}^4 p_n\right)$$

We actually adjust this with one more factor called the BP penalty.

BP = brevity penalty

is an adjustment factor that penalizes translations systems that output translations that are too short

$$BP = \begin{cases} 1 & \text{if } \underline{MT\_output\_length} > \underline{reference\_output\_length} \\ \exp(1 - MT\_output\_length/reference\_output\_length) & \text{otherwise} \end{cases}$$

So the blue score is useful single real number evaluation metric to use whenever you want your algorithm to generate a piece of text. And you want to see whether it has similar meaning as a reference piece of text generate by humans.

This is not used in speech recognition because there is usually one ground truth. And you just use another measures to see if you got the speech transcription on pretty much exactly word for word correct.

But for image caption, for multiple captions for a picture could be about equally good or for machine translations, there are multiple translations but equally good.

The Bleu score give you a way to evaluate that automatically and therefore speed up your development.

[Papineni et. al., 2002. Bleu: A method for automatic evaluation of machine translation]

Andrew Ng