

CS 383 – Computational Text Analysis

Lecture 19 Generating text

Adam Poliak

03/29/2023

Slides adapted from Philipp Koehn

Machine Learning in a nutshell

In a ML model, what are we training?

- **Parameters!**

How do we train parameters in supervised learning?

train parameters == figure out values for the parameters

- Update weights by using them to make predictions and seeing **how far off our predictions** are
 - **Loss function!**

Algorithm to learn weights?

- **SGD**
- Others exist but not covering them

Outline

Generating text

Evaluating Generated Text

When do we even want to generate text?

Machine Translation

Summarization

ChatBots

...

Generating text so far

Greedy approach:

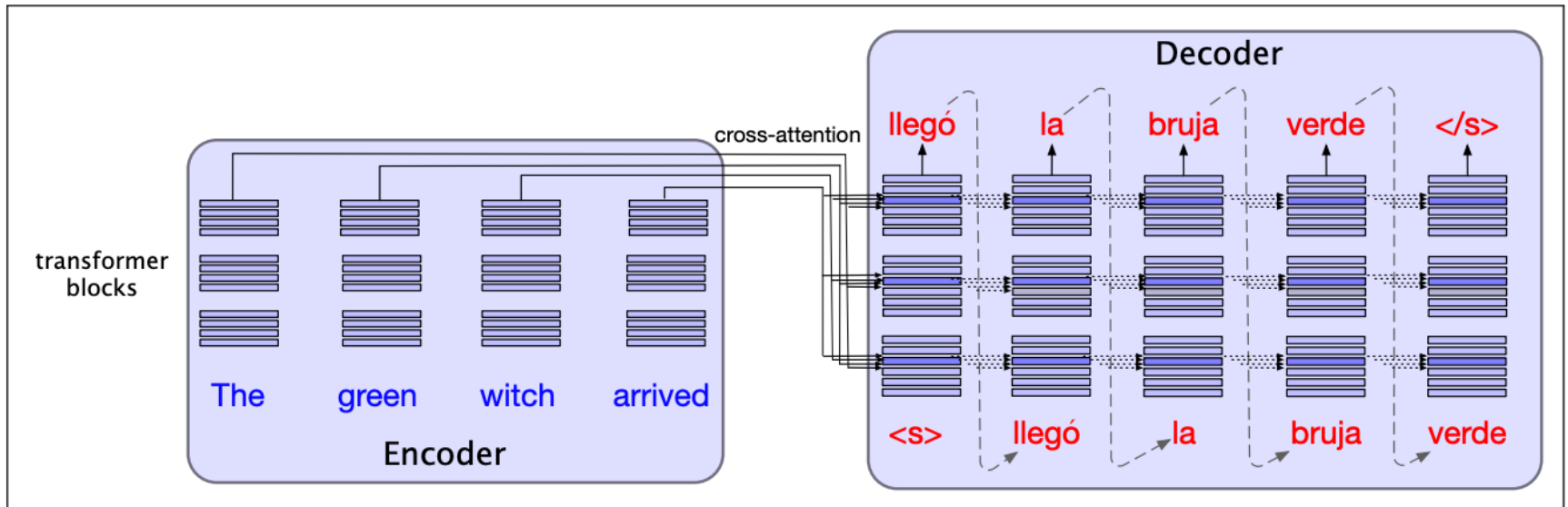
- At each time step, choose the word with the highest probability
- Optionally can condition on the previous generated words
 - Like in Language modeling!

Side: how do we prove correctness of greedy approaches?

Hint: use induction

<https://jeffe.cs.illinois.edu/teaching/algorithms/book/04-greedy.pdf>

Encoder-decoder for MT



Why might Greedy not be optimal?

$$\hat{y}_t = \operatorname{argmax}_{w \in V} P(w|y_1 \dots y_{t-1})$$

Choosing the best local word to generate, not necessarily the best sentence/paragraph globally to generate

Possible solutions?

Beam Search Decoding!

Search Tree

Nodes: states

The potential word to generate

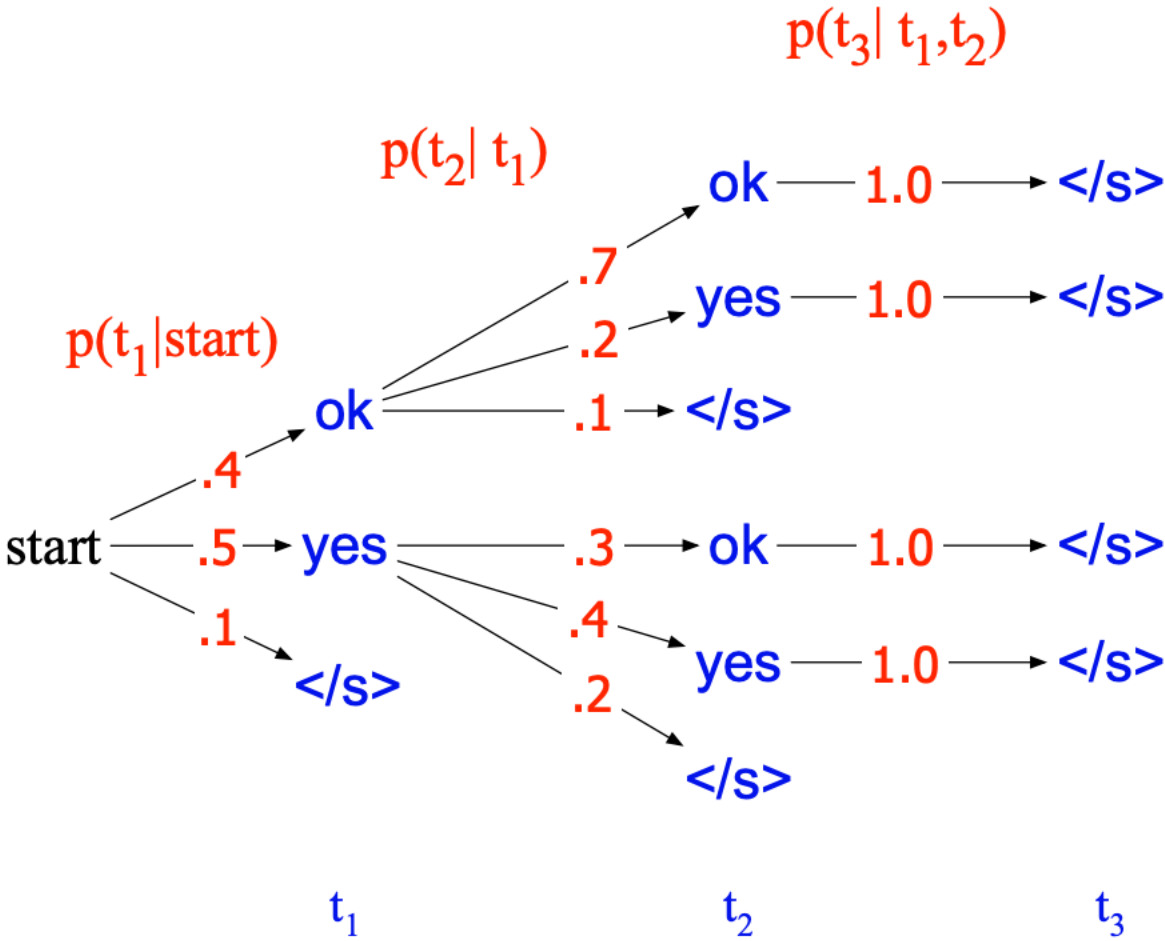
Branches/edges: actions

Edges have weights

Score for each word

Probability of generating each word

Search Tree Example



How long to search through the tree?

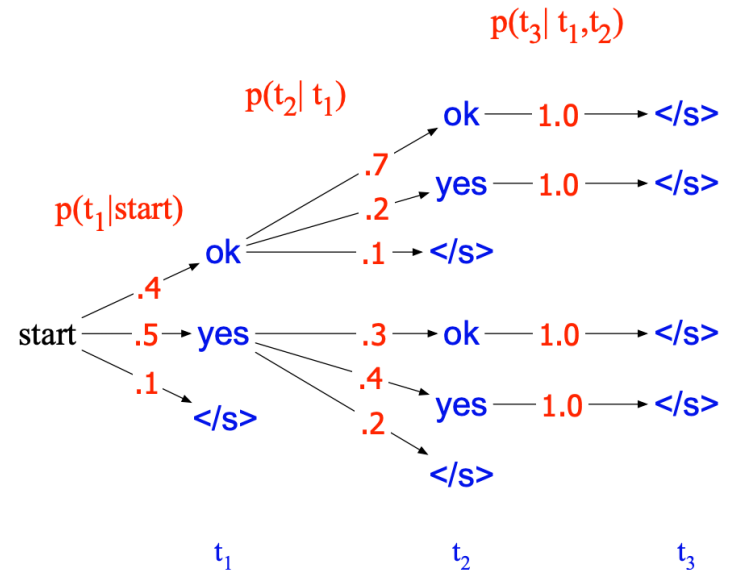
We need to consider all paths

How many paths are there?

$$V^T$$

Solution:

Beam search!

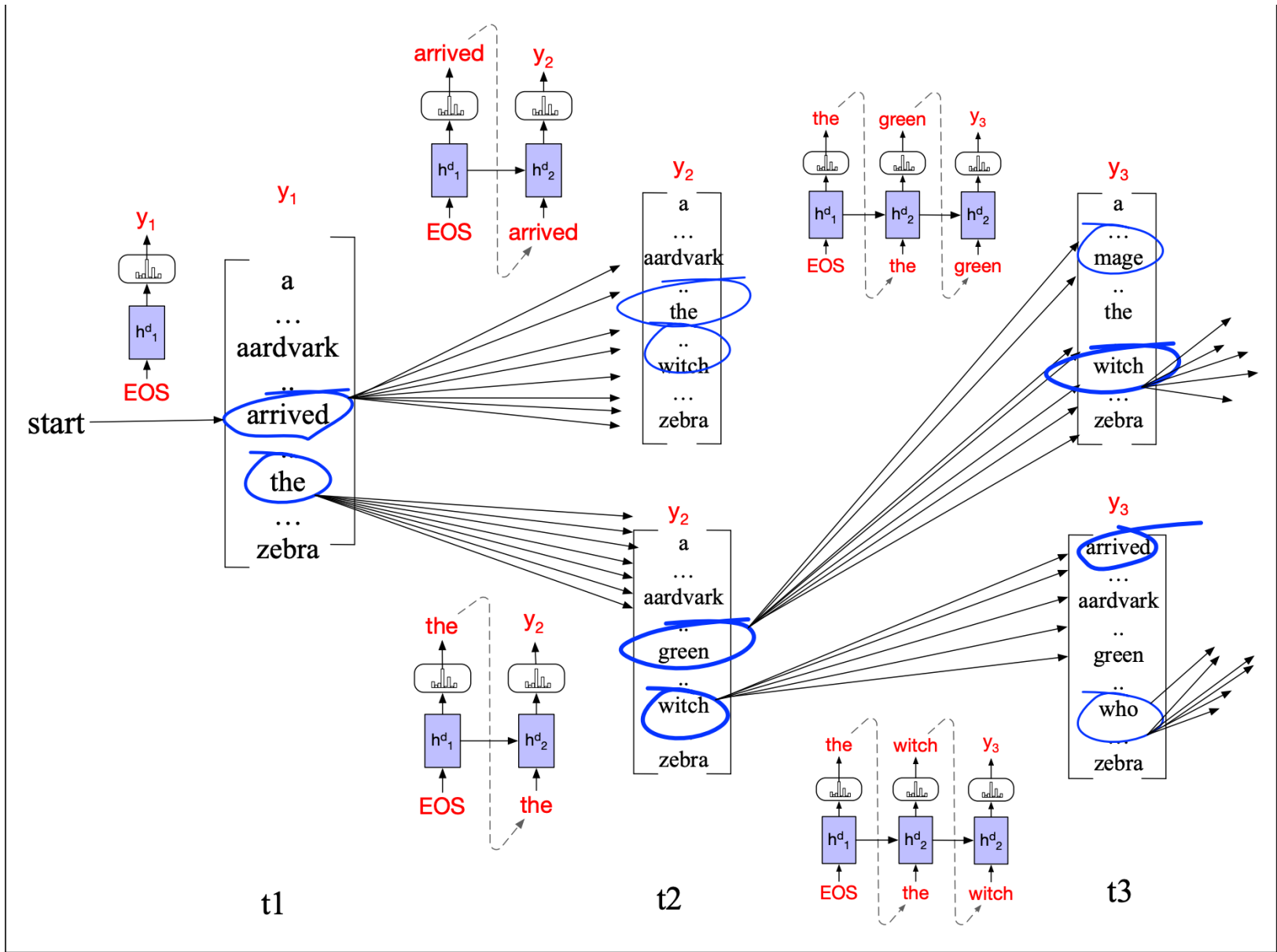


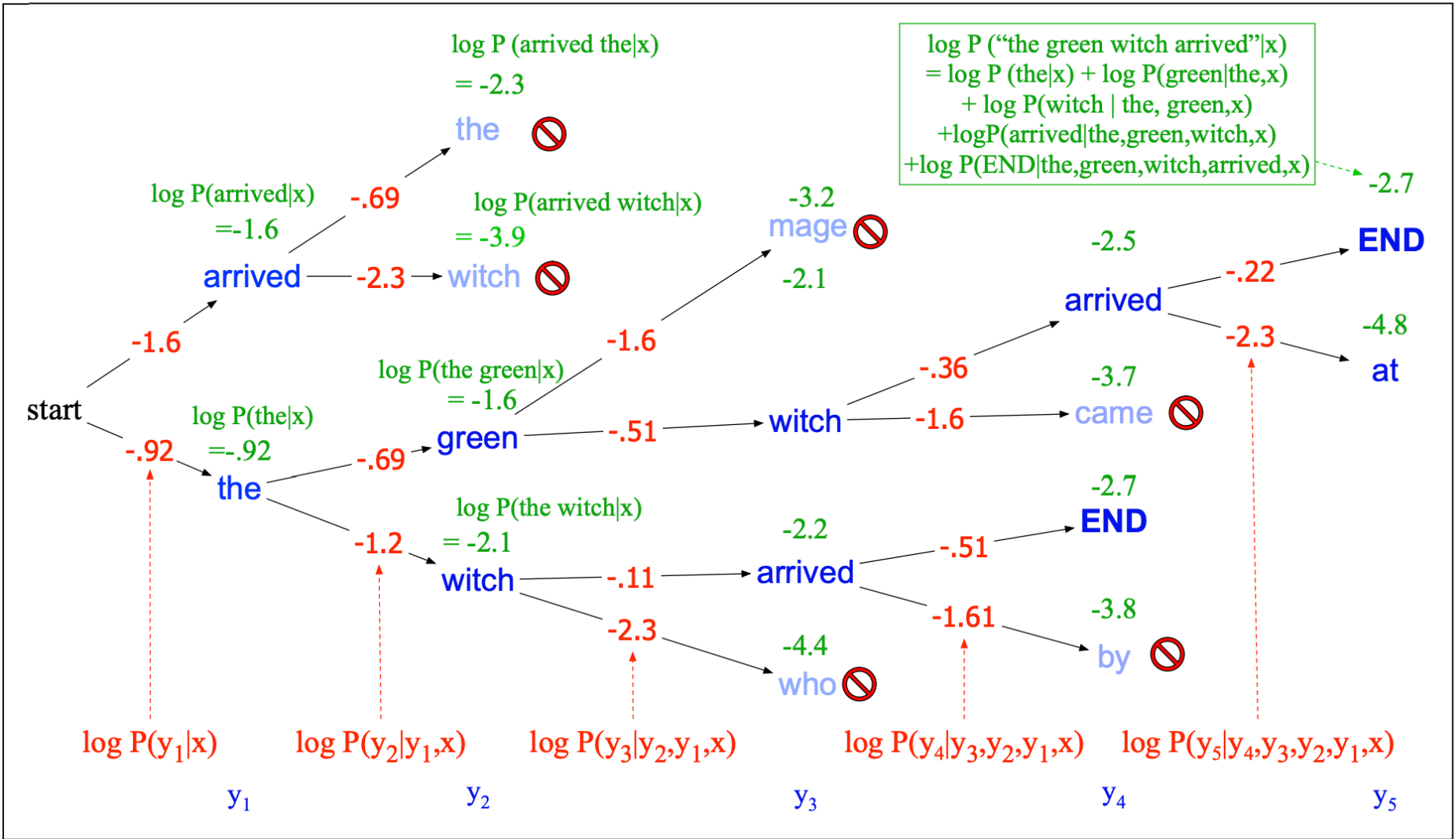
Beam search

Only keep k vocabulary terms at each time step
beam width

Apply softmax over vocabulary, keep just top k terms
search frontier
k-hypotheses

Incrementally repeat this until we generate `</s>`





Scoring hypotheses

$$\begin{aligned} \text{score}(y) &= \log P(y|x) \\ &= \log (P(y_1|x)P(y_2|y_1,x)P(y_3|y_1,y_2,x)\dots P(y_t|y_1,\dots,y_{t-1},x)) \\ &= \sum_{i=1}^t \log P(y_i|y_1,\dots,y_{i-1},x) \end{aligned}$$

What issues might there be?

Longer sentences will have higher scores

Solution:

Normalize by length of generated sentence

Algorithm

function BEAMDECODE(c , $beam_width$) **returns** best paths

$y_0, h_0 \leftarrow 0$

$path \leftarrow ()$

$complete_paths \leftarrow ()$

$state \leftarrow (c, y_0, h_0, path)$;initial state

$frontier \leftarrow \langle state \rangle$;initial frontier

while $frontier$ contains incomplete paths **and** $beamwidth > 0$

$extended_frontier \leftarrow \langle \rangle$

for each $state \in frontier$ **do**

$y \leftarrow \text{DECODE}(state)$

for each word $i \in Vocabulary$ **do**

$successor \leftarrow \text{NEWSTATE}(state, i, y_i)$

$extended_frontier \leftarrow \text{ADDTOBEAM}(successor, extended_frontier,$
 $beam_width)$

for each state in $extended_frontier$ **do**

if state is complete **do**

$complete_paths \leftarrow \text{APPEND}(complete_paths, state)$

$extended_frontier \leftarrow \text{REMOVE}(extended_frontier, state)$

$beam_width \leftarrow beam_width - 1$

$frontier \leftarrow extended_frontier$

return $completed_paths$

function NEWSTATE($state$, $word$, $word_prob$) **returns** new state

function ADDTOBEAM($state$, $frontier$, $width$) **returns** updated frontier

if $\text{LENGTH}(frontier) < width$ **then**

$frontier \leftarrow \text{INSERT}(state, frontier)$

else if $\text{SCORE}(state) > \text{SCORE}(\text{WORSTOF}(frontier))$

$frontier \leftarrow \text{REMOVE}(\text{WORSTOF}(frontier))$

$frontier \leftarrow \text{INSERT}(state, frontier)$

return $frontier$

Beam search decoding

What data structure might you use to store hypotheses?

- Priority Queue
 - Maintain the order of the best hypotheses

Beam search in action:

<http://mt-class.org/jhu/stack-decoder/>

Other decoding methods

A*

Hill-Climbing

Greedy approach

Finite State Transducers

Outline

Generating text

Evaluating Generated Text

How would you determine if a generated text is *good*

Why cannot we not use accuracy (or other stats from the confusion matrix)?

Compare the generated text to some human generated text

Comparing to human generated text

Why is this a hard problem?

many different generated text as acceptable →
semantic equivalence / similarity

How might you do this?

Metrics:

BLEU

Embedding approaches (e.g. BertScore)

BLEU

N-gram overlap between machine translation output and reference translation

Compute precision for n-grams of size 1 to 4

Add brevity penalty (for too short translations)

$$\text{BLEU} = \min \left(1, \frac{\text{output-length}}{\text{reference-length}} \right) \left(\prod_{i=1}^4 \text{precision}_i \right)^{\frac{1}{4}}$$

Typically computed over the entire corpus, not single sentence

BLEU Example

SYSTEM A: Israeli officials responsibility of airport safety
2-GRAM MATCH 1-GRAM MATCH

REFERENCE: Israeli officials are responsible for airport security

SYSTEM B: airport security Israeli officials are responsible
2-GRAM MATCH 4-GRAM MATCH

Metric	System A	System B
precision (1gram)	3/6	6/6
precision (2gram)	1/5	4/5
precision (3gram)	0/4	2/4
precision (4gram)	0/3	1/3
brevity penalty	6/7	6/7
BLEU	0%	52%

Plan for next week

Statistical Inference

- Quantifying uncertainty
 - Hypothesis testing
 - Null hypothesis
 - p-value
 - Confidence interval
 - Bootstrapping
- Forecast (time-series prediction)