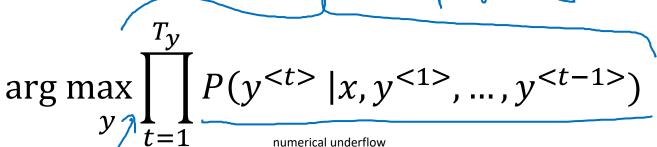


## Sequence to sequence models

## Refinements to beam search

Length normalization

P(y(1) - y(x) / x) P(y(1) / x)



too small for the floating part representation in your computer to store accurately

it unnaturally tends to prefer very short translation, it tends to prefer very short outputs

a strictly monotonically increasing function  $\log P(y|x) \leftarrow P(y|x) \leftarrow$ 

$$\arg\max_{y} \sum_{t=1}^{T_y} \frac{\log P(y^{< t>} | x, y^{< 1>}, ..., y^{< t-1>})}{\log P(y^{< t>} | x, y^{< 1>}, ..., y^{< t-1>})}$$

So by taking logs, you end up with a more numerically stable algorithm that is less prone to rounding errors, numerically rounding errors, or to really numerical underflow.

normalize take the average of the log of the probability of each word.

COA

And this significantly reduces the penalty for outputting longer translations

 $\sum_{t=1}^{T_y} \log P(y^{< t>} | x, y^{< 1>}, ..., y^{< t-1>})$ normalized log probability objective
normalized log likelihood objective

this is a heuristic or this is a hack, there isn't a great justification for it but people have found this work well.

p(ycty)(x,ycs,,,,cty-1))

this is somewhere in between full normalization and no normalization alpha is another hyperparameter of algorithm that you can tune to get the best result

d=0

**Andrew Ng** 

## Beam search discussion

large B: better result, slower small B: worse result, faster

Beam width B?

Unlike exact search algorithms like BFS (Breadth First Search) or DFS (Depth First Search), Beam Search runs faster but is not guaranteed to find exact maximum for arg max P(y|x).

y

And there are some simple things you can compute to give you guidance on what you need to work on improving your algorithm.