

# Smoothing

What do we do with words that are in our vocabulary but appear in a test set in an unseen context (for example they appear after a word they never appeared after in training)?

— In this case zero probability to these unseen event

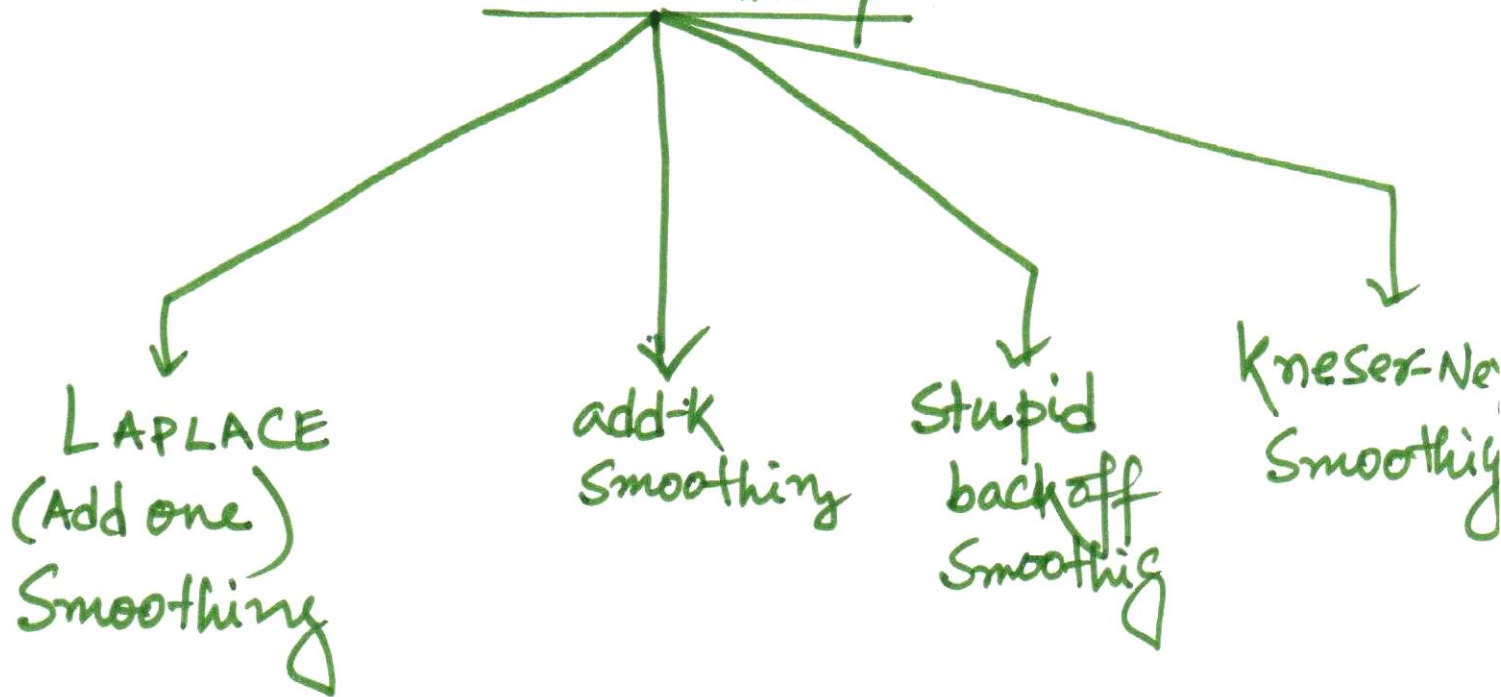
for example

Given Corpus:

<s> I am Henry </s>  
<s> I like college </s>  
<s> Do Henry like college </s>  
<s> Henry I am </s>  
<s> Do I like Henry </s>  
<s> Do I like college </s>  
<s> Do I like Henry </s>

Calculate (Big  
<s> I like a college  
</s>

# SMOOTHING



$\langle s \rangle$  I like a college  $\langle /s \rangle$

$$\equiv P(\langle s \rangle) \times P(I | \langle s \rangle) \times P(a | \text{like}) \times$$

$$P(\text{college} | a) \times P(\langle /s \rangle | \text{college})$$

$\equiv ?$

## LAPLACE SMOOTHING

The simplest way to do smoothing is to add one to all the n-gram counts, before we normalize them into probabilities.

- All the counts that used to be zero will now have a count of 1, the counts of 1 will be 2 and so on. This algorithm is called Laplace Smoothing.

$$P_{\text{Laplace}}(w_i) = \frac{C_i + 1}{N + V}$$

$V = \text{words in vocabulary}$