

Challenges in N-Gram Language Modeling

N-Gram Language Modeling - Challenges

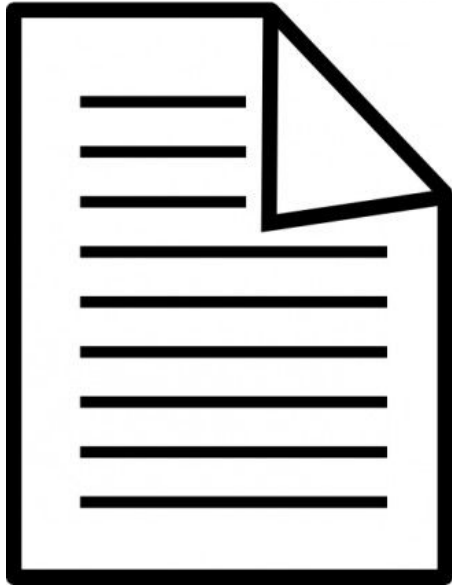
- $N = 3$ (Tri-gram model)
- $P(\text{"Jack built that house"})$



N-Gram Language Modeling - Challenges

- $N = 3$ (Tri-gram model)
- $P(\text{"Jack built that house"}) =$
 $P(\text{"Jack"}) \cdot P(\text{"built"}|\text{"Jack"}) \cdot P(\text{"that"}|\text{"Jack built"}) \cdot P(\text{"house"}|\text{"built that"})$

N-Gram Language Modeling - Challenges



“Jack” appears 50 times

Total words = 500

N-Gram Language Modeling - Challenges



“Jack” appears 50 times

Total words = 500

$$P(\text{“Jack”}) = 50/500 = 0.1$$

N-Gram Language Modeling - Challenges



10 sequences start with "Jack"

Total words = 500

N-Gram Language Modeling - Challenges



10 sequences start with "Jack"

Total words = 500

$$P(\text{"Jack"}) = 10/500 = 0.02$$

N-Gram Language Modeling - Challenges

- $N = 3$ (Tri-gram model)
- $P(\text{"Jack built that house"}) =$
 $P(\text{"Jack"}) \cdot P(\text{"built"}|\text{"Jack"}) \cdot P(\text{"that"}|\text{"Jack built"}) \cdot P(\text{"house"}|\text{"built that"})$

N-Gram Language Modeling - Challenges

- $N = 3$ (Tri-gram model)
- $P(\text{"Jack built that house"}) =$
 $P(\text{"Jack"}) \cdot P(\text{"built"}|\text{"Jack"}) \cdot P(\text{"that"}|\text{"Jack built"}) \cdot P(\text{"house"}|\text{"built that"})$
- Add a start token
 - $P(\text{"<s> Jack built that house"})$

N-Gram Language Modeling - Challenges

- Input = “Jack built that house”



N-Gram Language Modeling - Challenges

- Input = “Jack built that house”
- Model output = “Jack built that house in two years...”



N-Gram Language Modeling - Challenges

- Input = “Jack built that house”
- Model output = “Jack built that house in two years in...”



N-Gram Language Modeling - Challenges

- Input = “Jack built that house”
- Model output = “Jack built that house in two years in India...”



N-Gram Language Modeling - Challenges

- Input = “Jack built that house”
- Model output = “Jack built that house in two years in India is...”



N-Gram Language Modeling - Challenges

- Input = “Jack built that house”
- Model output = “Jack built that house in two years in India is...”
- Model does not know when to stop text generation

N-Gram Language Modeling - Challenges

- $N = 3$ (Tri-gram model)
- $P(\text{"Jack built that house"}) =$
 $P(\text{"Jack"}) \cdot P(\text{"built"}|\text{"Jack"}) \cdot P(\text{"that"}|\text{"Jack built"}) \cdot P(\text{"house"}|\text{"built that"})$
- Add Start Token
 - $P(\text{"<s> Jack built that house"})$
- Add End Token
 - $P(\text{"<s> Jack built that house </s>"})$

N-Gram Language Modeling - Challenges

- N = 3 (Tri-gram model)

- $P(\text{"Jack built that house"}) =$

$P(\text{"Jack"}) \cdot P(\text{"built"}|\text{"Jack"}) \cdot P(\text{"that"}|\text{"Jack built"}) \cdot P(\text{"house"}|\text{"built that"})$

- Add Start Token

- $P(\text{"<s> Jack built that house"})$

- Add End Token

- $P(\text{"<s> Jack built that house </s>"})$

$P(\text{"Jack"}) \cdot P(\text{"built"}|\text{"Jack"}) \cdot P(\text{"that"}|\text{"Jack built"}) \cdot P(\text{"house"}|\text{"built that"}) \cdot P(\text{"</s>"}|\text{"that house"})$

N-Gram Language Modeling - Challenges

Limitations of N-Gram LM

- The higher the N, the better is the model usually. But this leads to lots of computation overhead that requires large computation power.

Analytics
Vidhya

N-Gram Language Modeling - Challenges

Limitations of N-Gram LM

- The higher the N, the better is the model usually. But this leads to lots of computation overhead that requires large computation power.
- **Sparsity Problem:** N-grams are a sparse representation of language. It will give zero probability to all the words that are not present in the training corpus.

N-Gram Language Modeling - Challenges

- Training Set
 - ...find the differences
 - ...find the station
 - ...find the jacket



N-Gram Language Modeling - Challenges

- Training Set
 - ...find the differences
 - ...find the station
 - ...find the jacket
- Unseen Set
 - ...find the cat



N-Gram Language Modeling - Challenges

- Training Set
 - ...find the differences
 - ...find the station
 - ...find the jacket

- Unseen Set

- ...find the cat

- $P(\text{"cat"} \mid \text{"find the"}) = 0$



N-Gram Language Modeling - Challenges

- **Laplace Smoothing:** Add one to the occurrence of each word



N-Gram Language Modeling - Challenges

- **Laplace Smoothing:** Add one to the occurrence of each word
- Just add one to all the counts



N-Gram Language Modeling - Challenges

- **Laplace Smoothing:** Add one to the occurrence of each word
- Just add one to all the counts

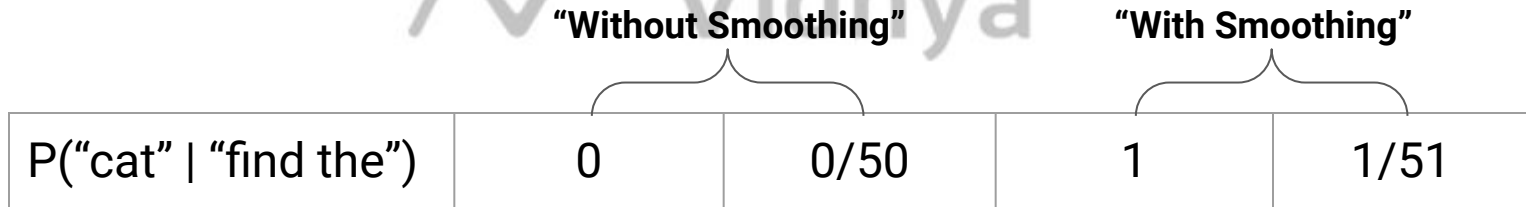
 Analytics Vidhya

"Without Smoothing"

$P(\text{"cat"} \mid \text{"find the"})$	0	0/50
--	---	------

N-Gram Language Modeling - Challenges

- **Laplace Smoothing:** Add one to the occurrence of each word
- Just add one to all the counts



	"Without Smoothing"		"With Smoothing"	
P("cat" "find the")	0	0/50	1	1/51

N-Gram Language Modeling - Challenges

- **Backoff:** Switch to shorter context/history as an approximation
 - $P(\text{"offer"} | \text{"denied the"}) \rightarrow P(\text{"offer"} | \text{"the"}) \rightarrow P(\text{"offer"})$



N-Gram Language Modeling - Challenges

- **Backoff:** Switch to shorter context/history as an approximation
 - $P(\text{"offer"} | \text{"denied the"}) \rightarrow P(\text{"offer"} | \text{"the"}) \rightarrow P(\text{"offer"})$
- **Interpolation:** Use weighted sum of unigrams, bigrams, and trigrams
 - $$P(\text{word}_n | \text{word}_{n-1} \text{ word}_{n-2}) = a_1 P(\text{word}_n | \text{word}_{n-1} \text{ word}_{n-2}) + a_2 P(\text{word}_n | \text{word}_{n-1}) + a_3 P(\text{word}_n)$$
 - $\sum a_i = 1$

N-Gram Language Modeling - Challenges

- “There she built a ? ”



N-Gram Language Modeling - Challenges

- “There she built a ? ”
- “Alice went to the beach. There she built a ? ”



N-Gram Language Modeling - Challenges

- “There she built a ? ”
- “Alice went to the beach. There she built a sandcastle ”





Thank You