

# NGRAM Data Acquisition

1. Download small data set
  - Lacking performance
2. Download big data set
  - a. ca. 9 TB to download = 270 GB after aggregation
  - b. ca. 36 days processing time
3. Request Backend of Ngram Viewer
  - a. Slow -> Caching
  - b. Access to big data without download
  - c. Simple switch between languages
  - d. Vulnerable to changes by Google

Google Books Ngram Viewer

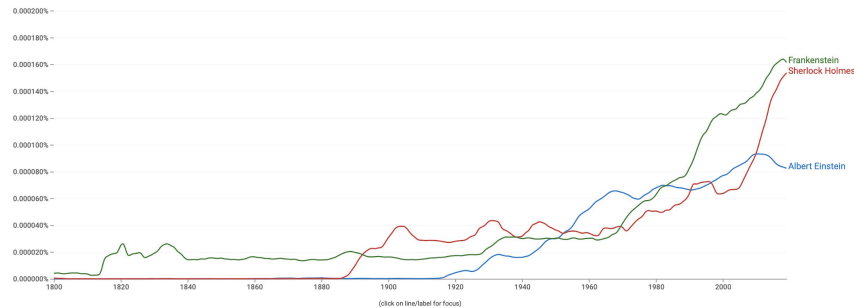
Albert Einstein, Sherlock Holmes, Frankenstein

1800 - 2019

English (2019)

Case-Insensitive

Smoothing



# Error Detection using N-GRAM Language Models

- n-gram models with up to 4-grams
- assumption: an error is detectable by the previous 3 words
- calculating an error value for every word if it is below a threshold mark it as error
- smoothing for unseen n-grams
- using “stupid backoff” a smoothing technique for huge language models
- according to the authors of the paper  $\lambda=0.4$  works well

$$S(w_i | w_{i-k+1}^{i-1}) = \begin{cases} \frac{\text{count}(w_{i-k+1}^i)}{\text{count}(w_{i-k+1}^{i-1})} & \text{if } \text{count}(w_{i-k+1}^i) > 0 \\ \lambda S(w_i | w_{i-k+2}^{i-1}) & \text{otherwise} \end{cases}$$

# N-GRAM Language Model Example

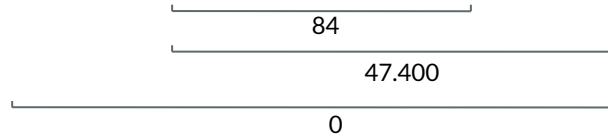
\_START\_ The boy live on top of the hill.

	147.759.931	2.061.498	3.451.293
107.219.041			149
	143.660		
120.428			
		0	
	0		

$$S = \frac{\text{count}(\text{\_START\_ The boy live})}{\text{count}(\text{\_START\_ The boy})} \Rightarrow \frac{\lambda \text{count}(\text{The boy live})}{\text{count}(\text{The boy})} \Rightarrow \frac{\lambda \lambda \text{count}(\text{boy live})}{\text{count}(\text{boy})} \Rightarrow \frac{0.4 * 0.4 * 149}{2.061.498} \approx 0.0000115644$$

# Smoothing Example

`_START_` The boy lives on top of the hill.



$$S = \frac{\text{count}(\text{boy lives on top})}{\text{count}(\text{boy lives on})} \Rightarrow \frac{\lambda \text{count}(\text{lives on top})}{\text{count}(\text{lives on})} \Rightarrow \frac{0.4 * 84}{47.400} \approx 0.0007148936$$

# Error correction: The BERT Model



- Machine Learning model published 2018 by Google
- Transfer learning algorithm
- Mask Language Modelling:

Fill-Mask

Mask token: [MASK]

I like to play [MASK] when the weather is nice.

Compute

Computation time on cpu: 0.053 s

golf	0.098
football	0.077
hockey	0.040
here	0.039
tennis	0.037

# Error correction: Implementation



- BERT Model from the NLP company Hugging Face
- Training data: 11.000 books and the complete English Wikipedia

## Correction process:

1. Replace the wrong word by a mask token
2. Receive the twenty most likely suggestions from the BERT model
3. Evaluate the twenty suggestions using N-Grams and word distance
4. Show the four remaining suggestions to the user