



**Lesson5:**  
**Generative Models for Text on the Web**  
**Unit3:**  
**Evaluating the Quality of a generative Model**

Rene Pickhardt

Introduction to Web Science Part 2  
Emerging Web Properties

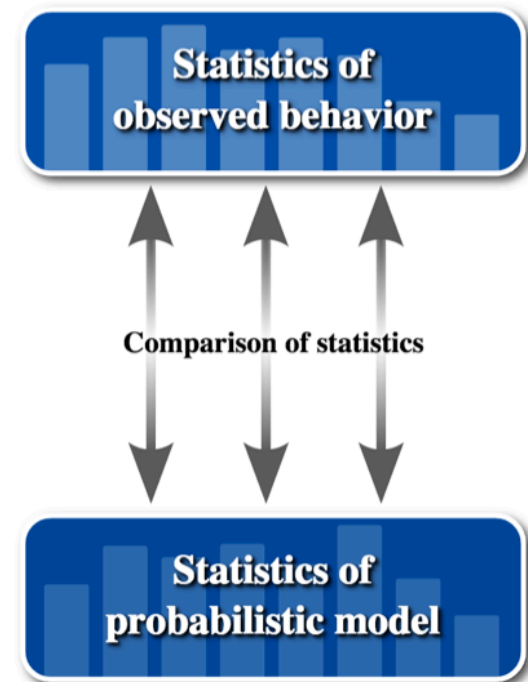
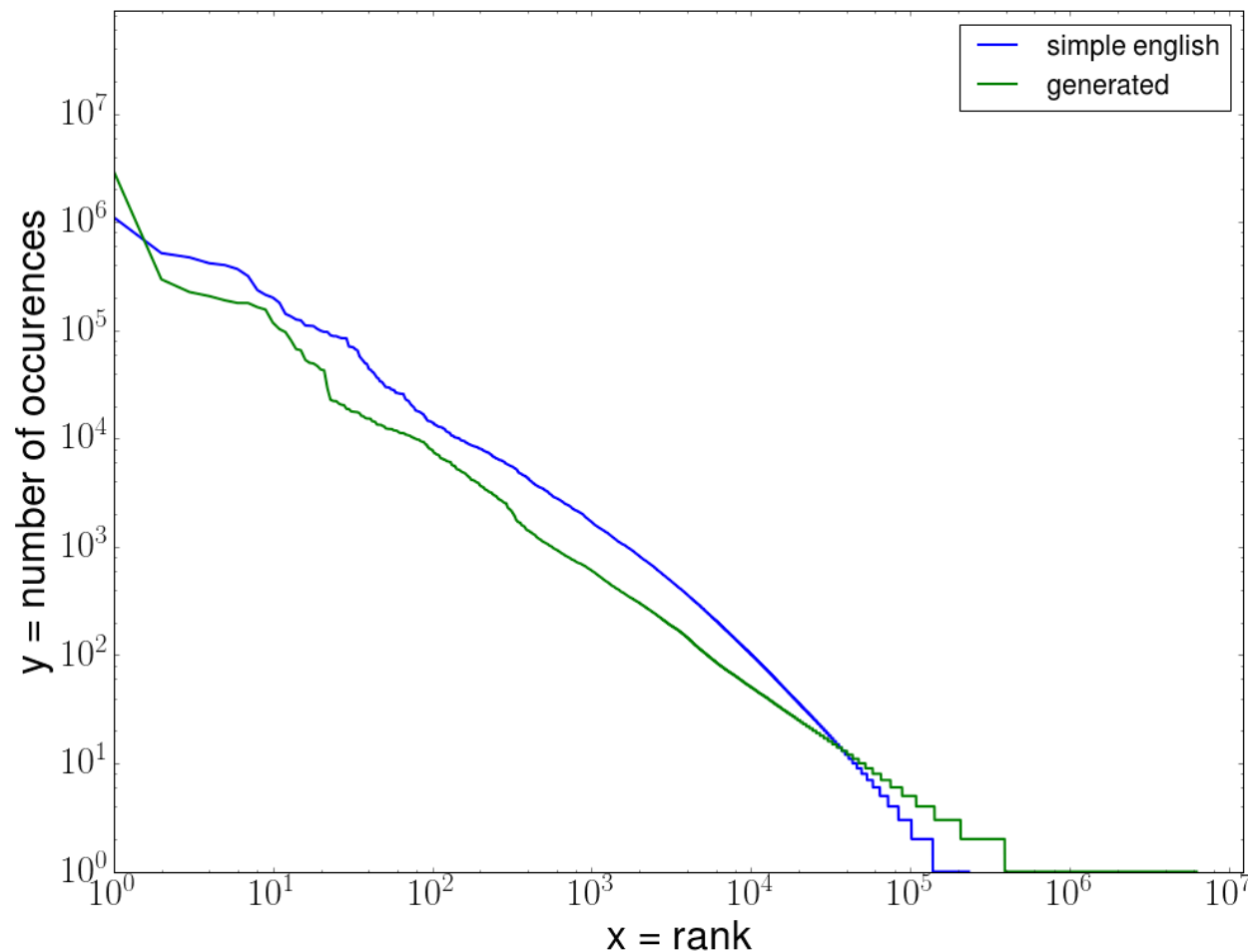
**WeST**   
People and Knowledge Networks

## Completing this unit you should

- See that it makes sense to compare statistics
- Understand that comparing statistics is not a well defined task
- Be aware of the fact that very different models could lead to the same statistics

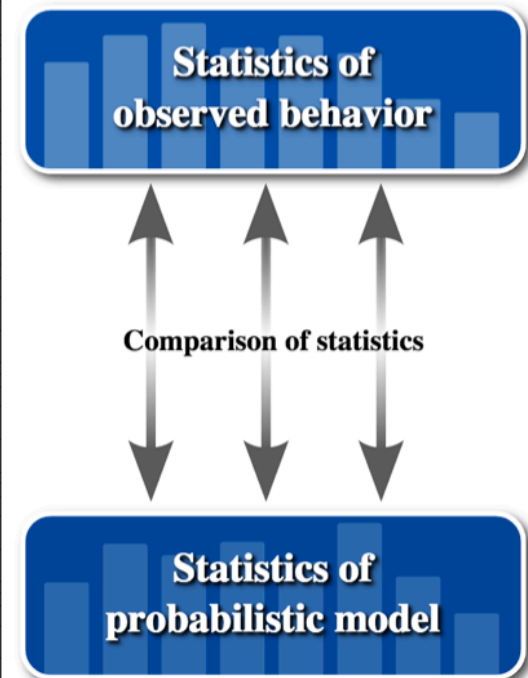
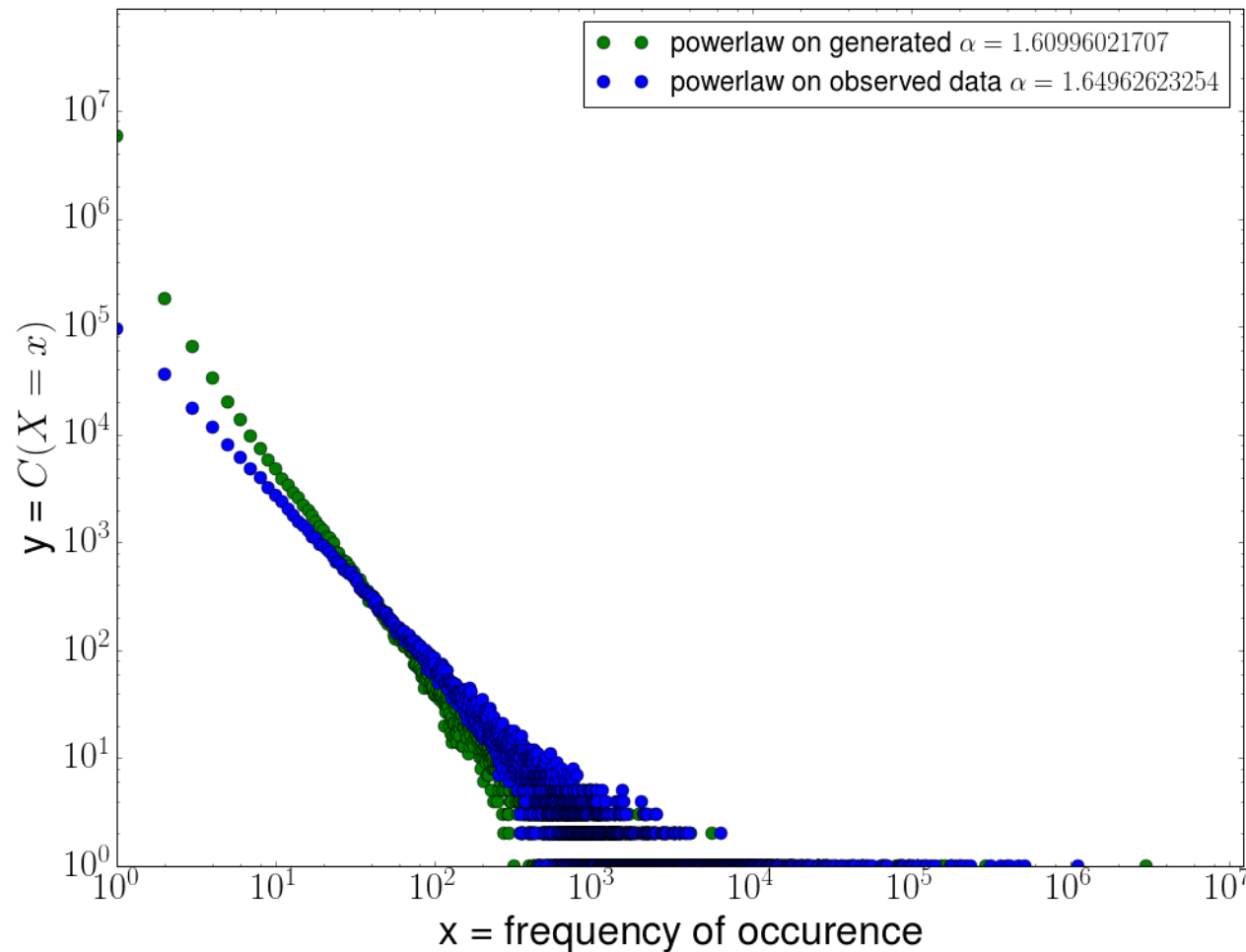
# Plotting the Zipf distribution

Word frequencies depending on word rank on (Simple) English Wikipedia



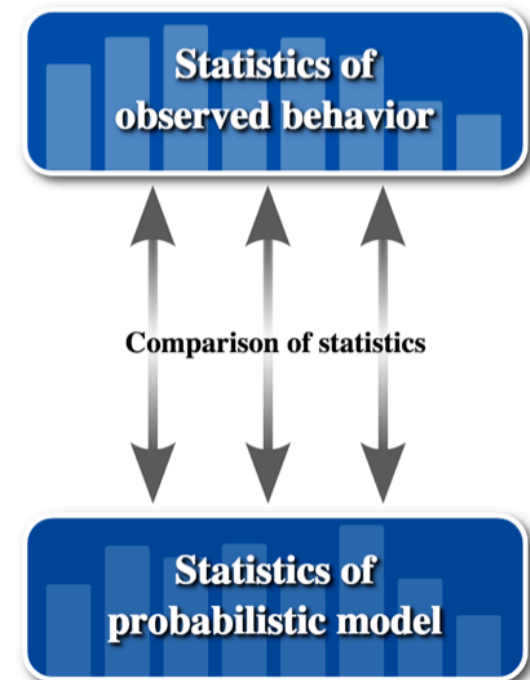
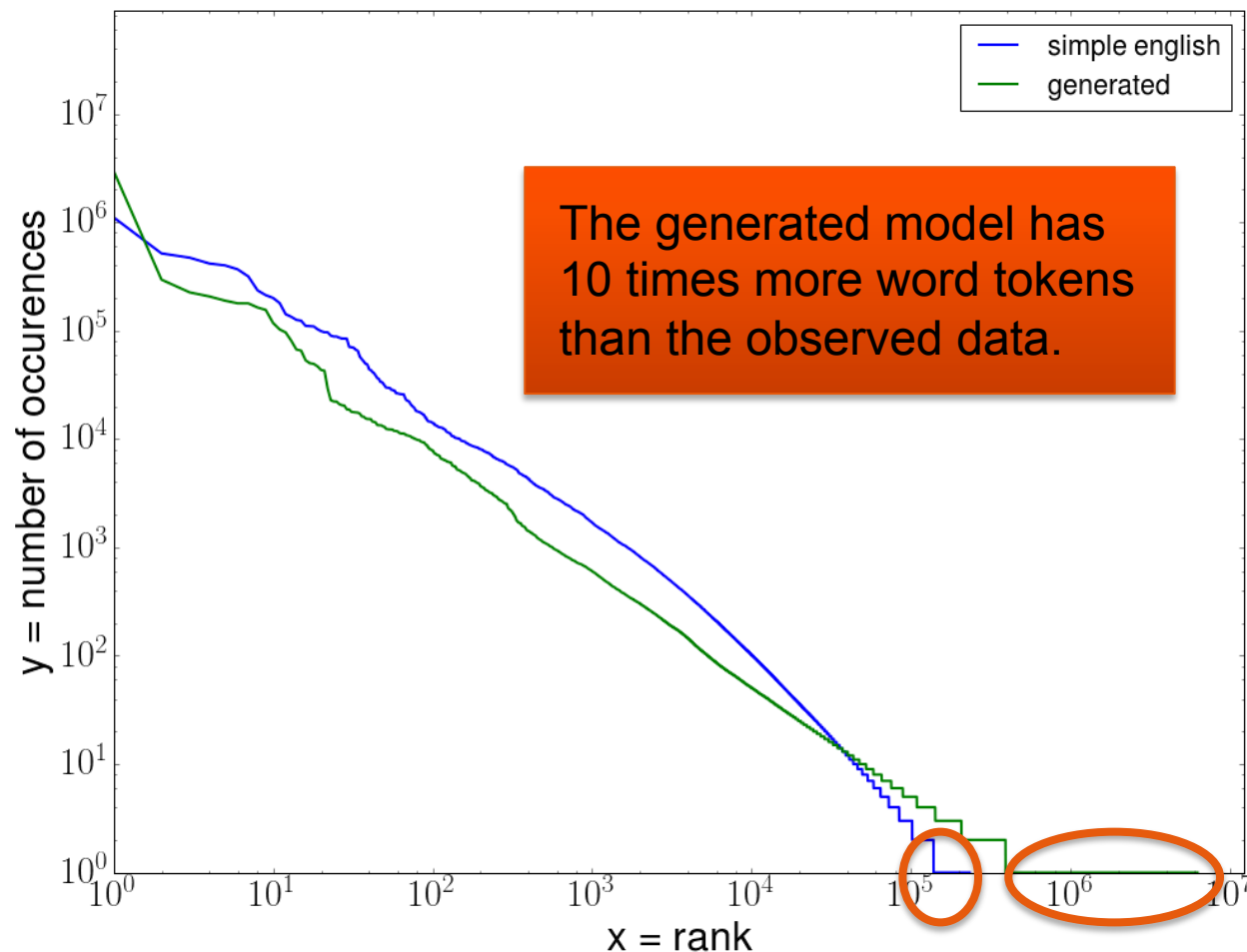
# Powerlaw plot and exponent look good

PowerLaw on Simple Wikipedia and Generated Words



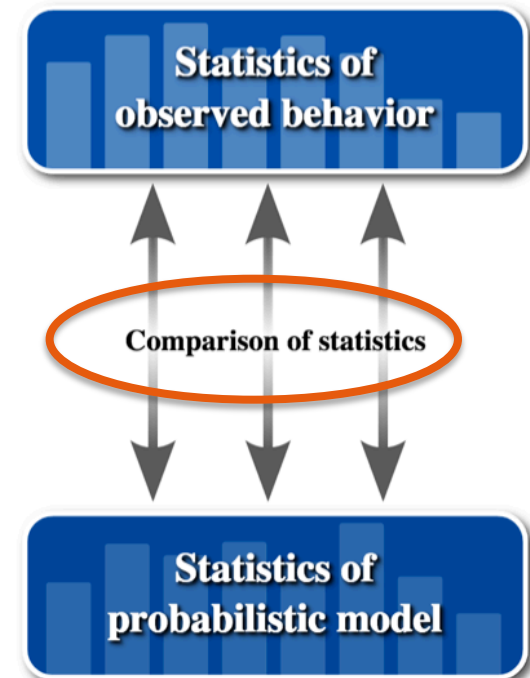
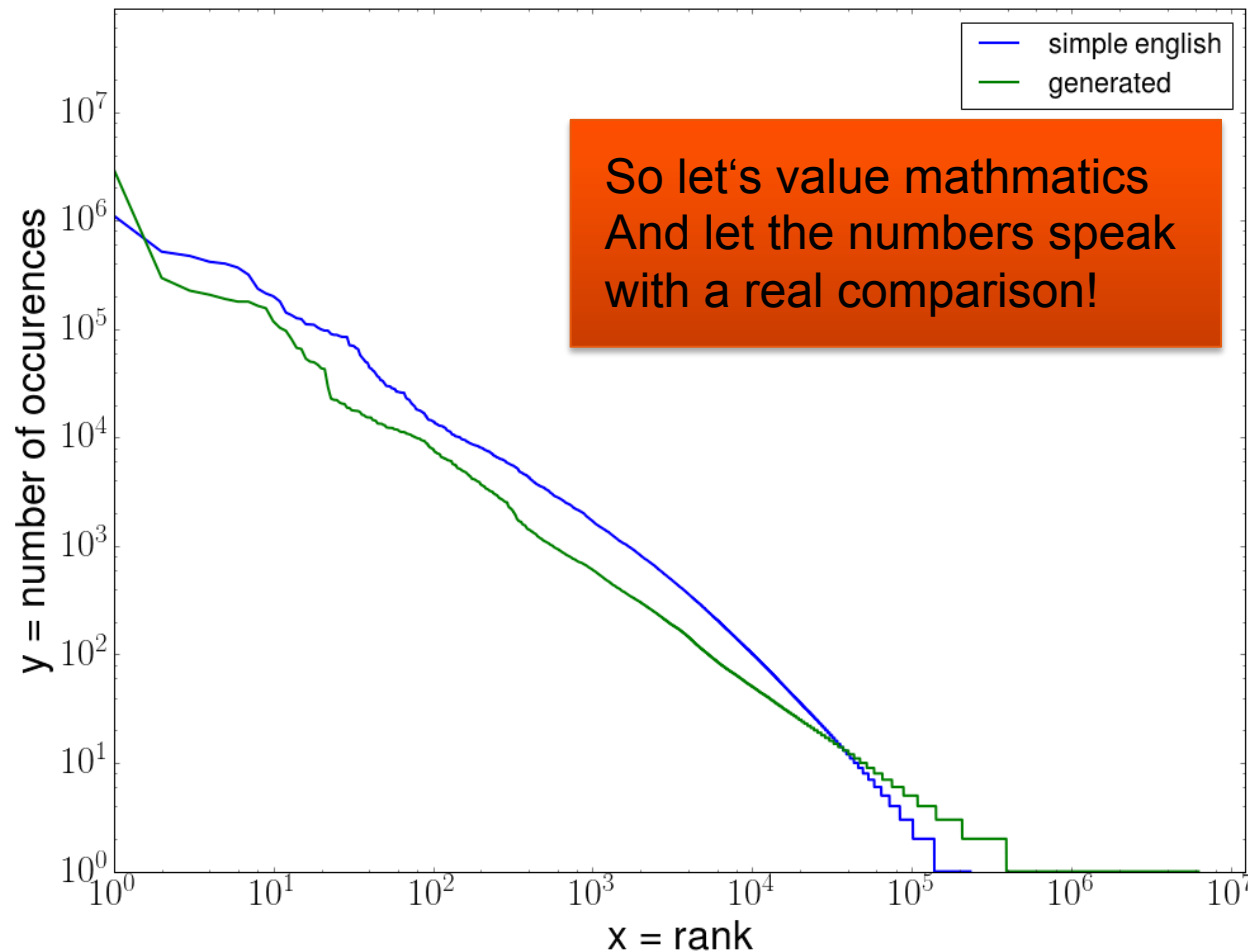
# Remember not to get fooled by the log scale

Word frequencies depending on word rank on (Simple) English Wikipedia



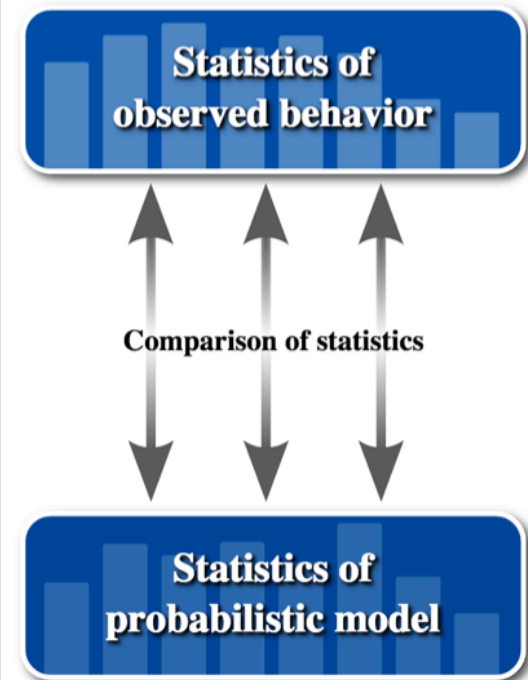
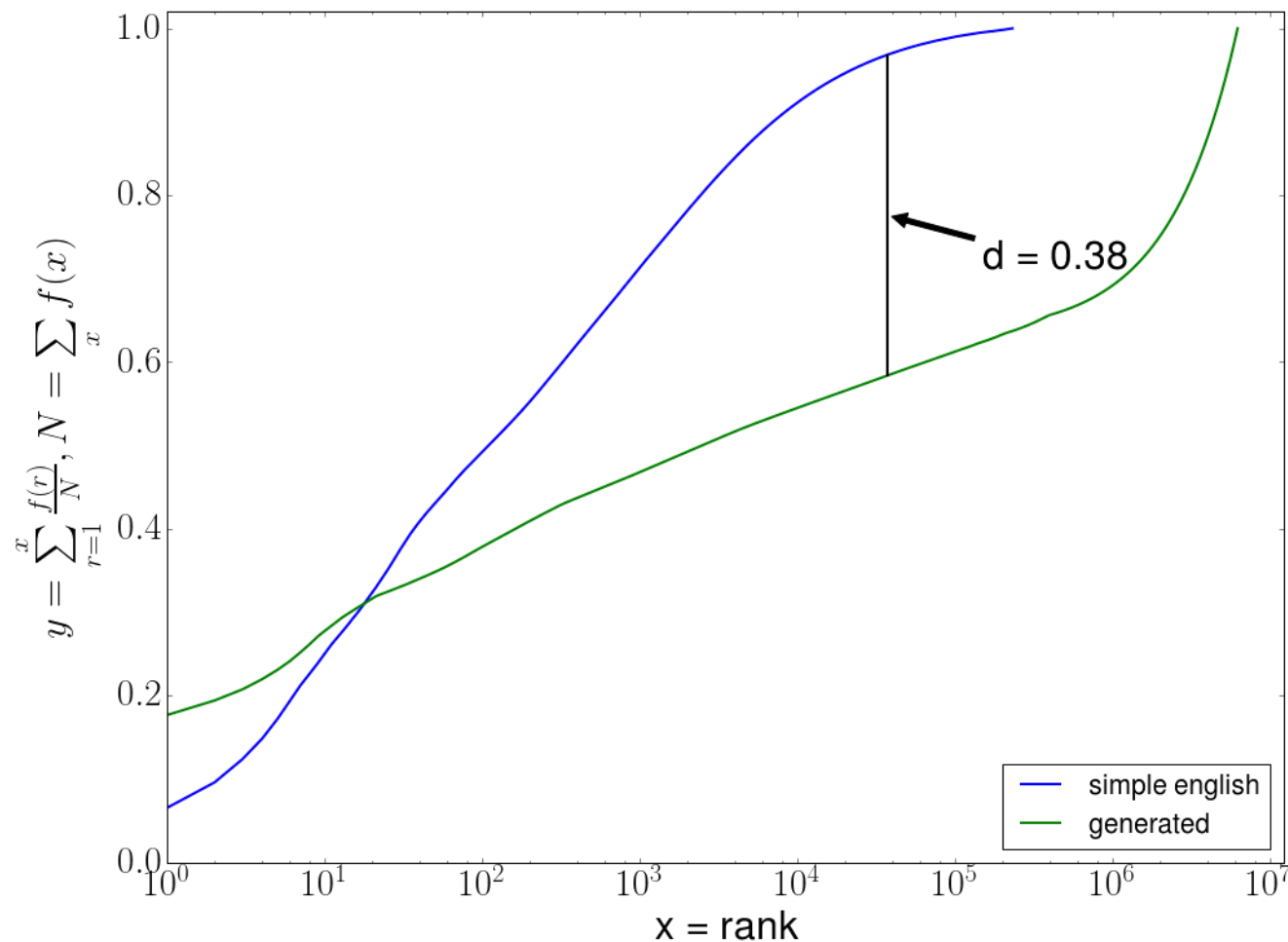
# Who remembers a good distance measure?

Word frequencies depending on word rank on (Simple) English Wikipedia



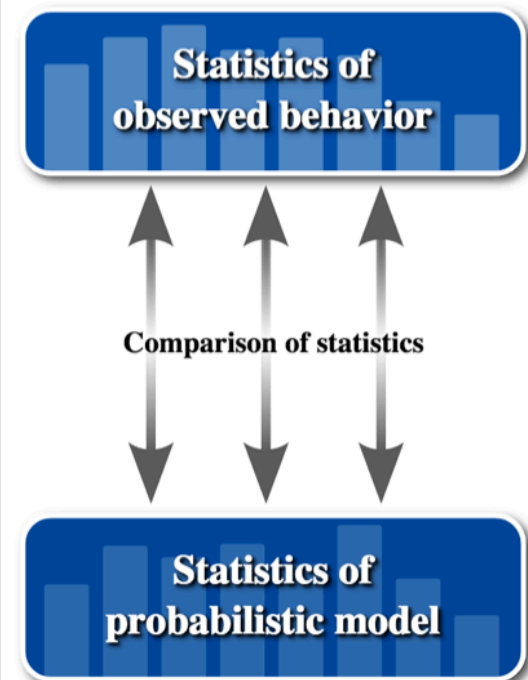
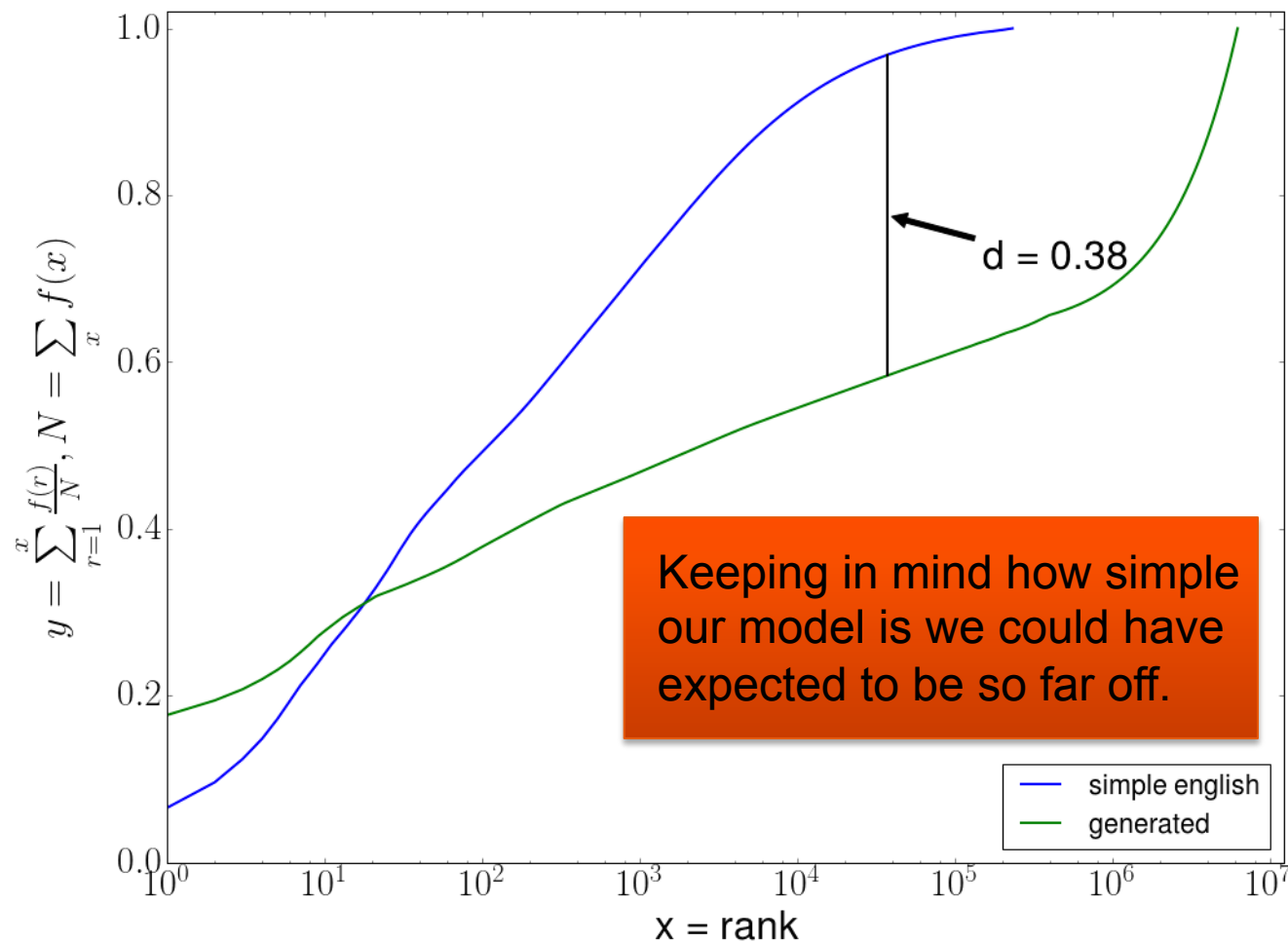
# OUCH! That's a surprise – Or isn't it?

Cumulative word probabilities depending on word rank



# OUCH! That's a surprise – Or isn't it?

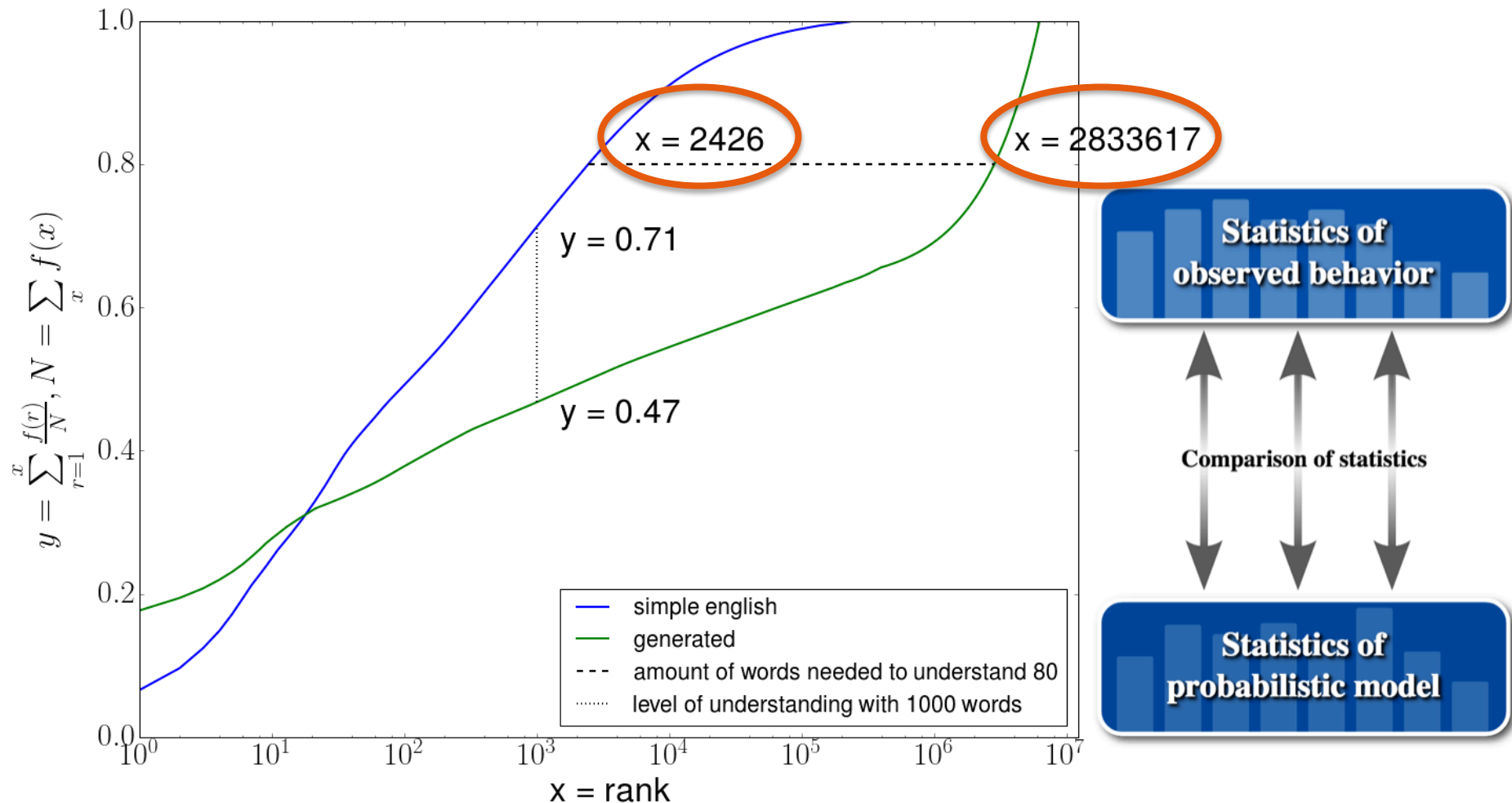
Cumulative word probabilities depending on word rank





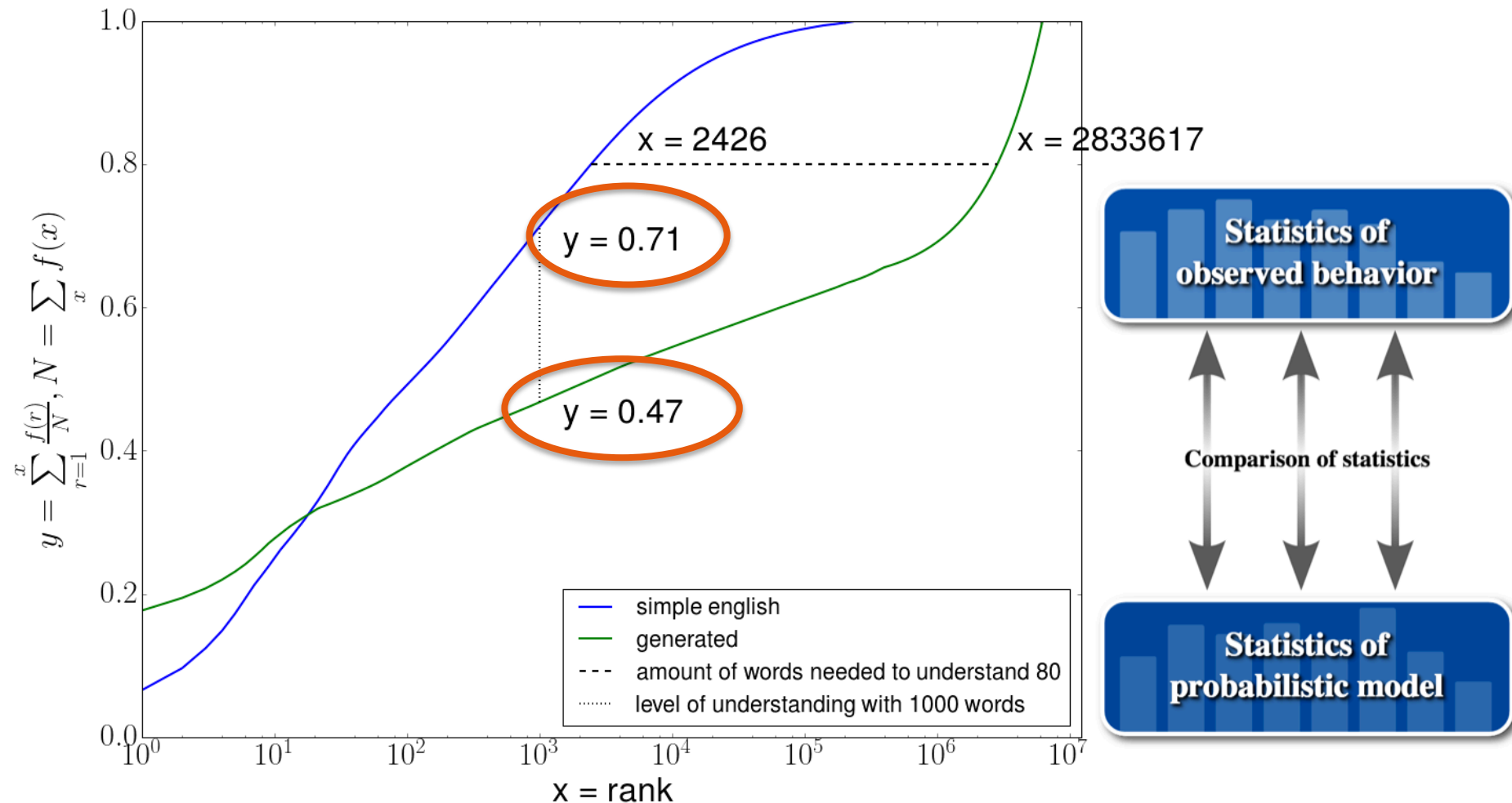
# To understand 80% you need to know 1150 times more words in the generated language

Cumulative word probabilities depending on word rank



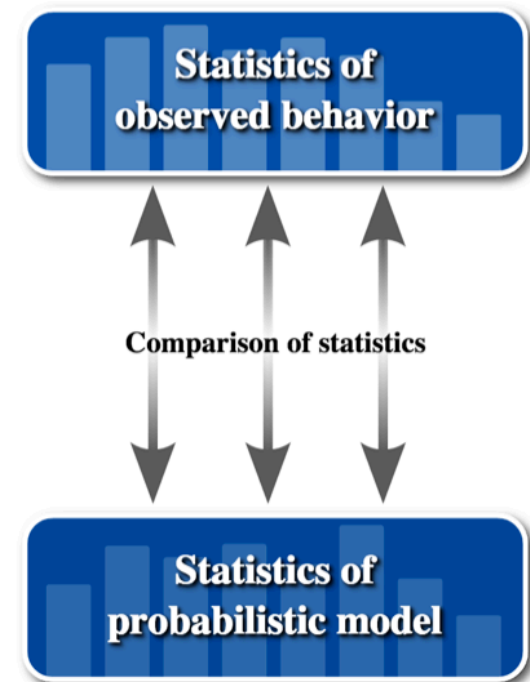
# Level of expressiveness with 1000 known words also differs drastically

Cumulative word probabilities depending on word rank



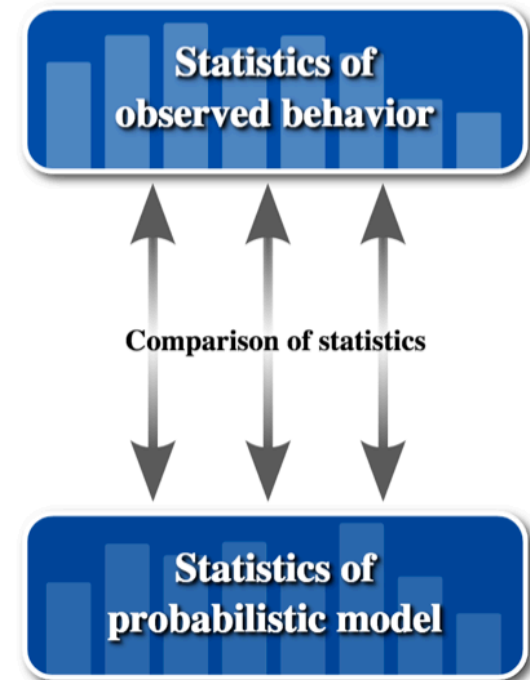
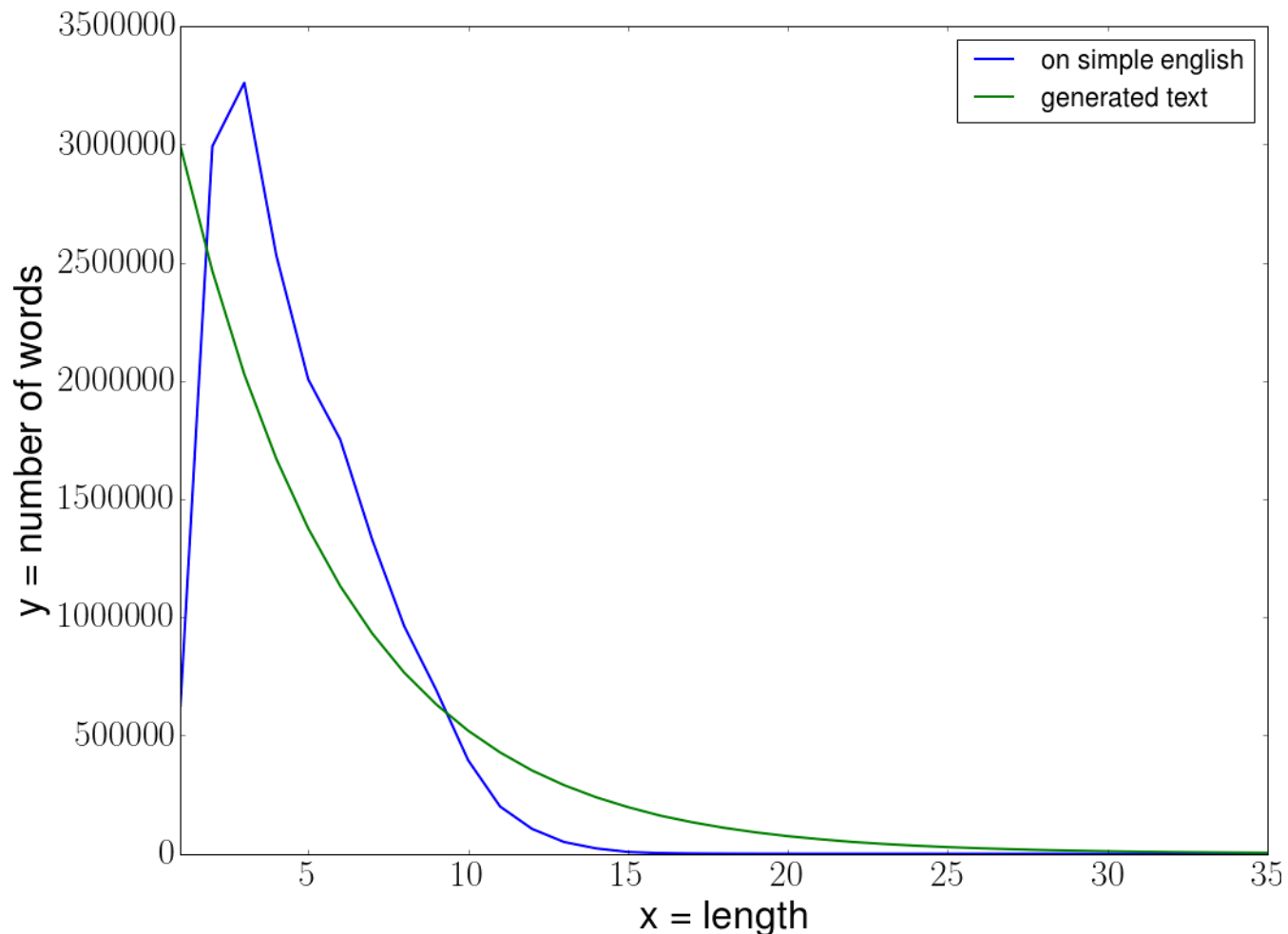
## More descriptive statistics for text modelling

- Word length distribution
- Document frequency
- Sentence distribution
- Number of documents
  - Can't be used in since we didn't generate documents



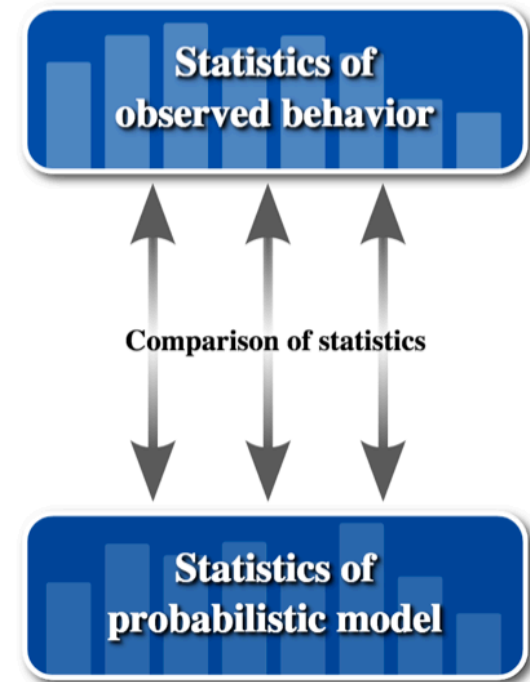
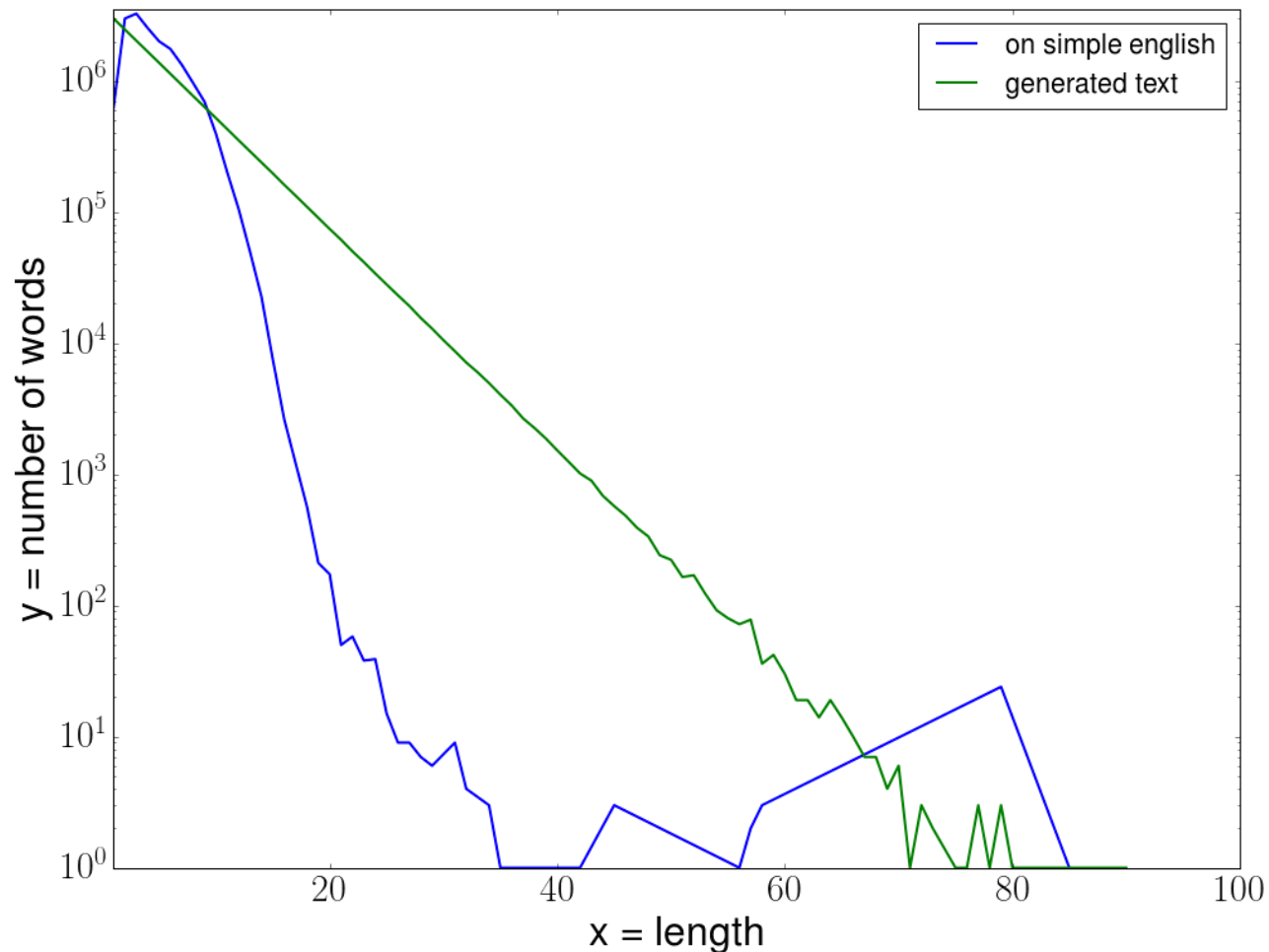
# Generated text has too many short and long words. Too few between 3 and 9 characters

Length distribution of Words



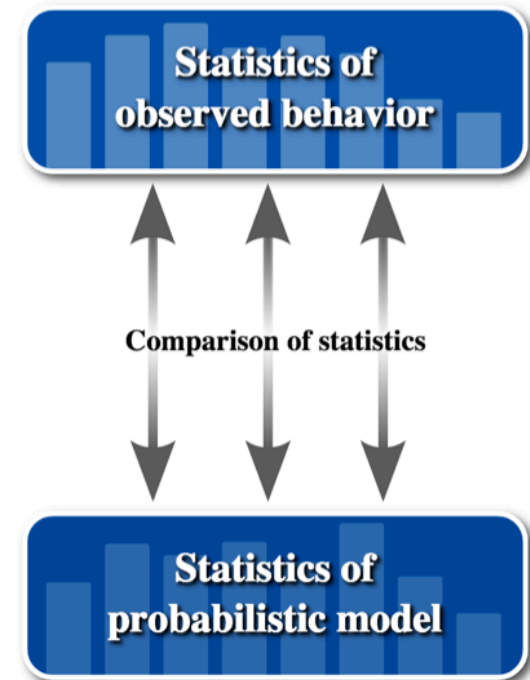
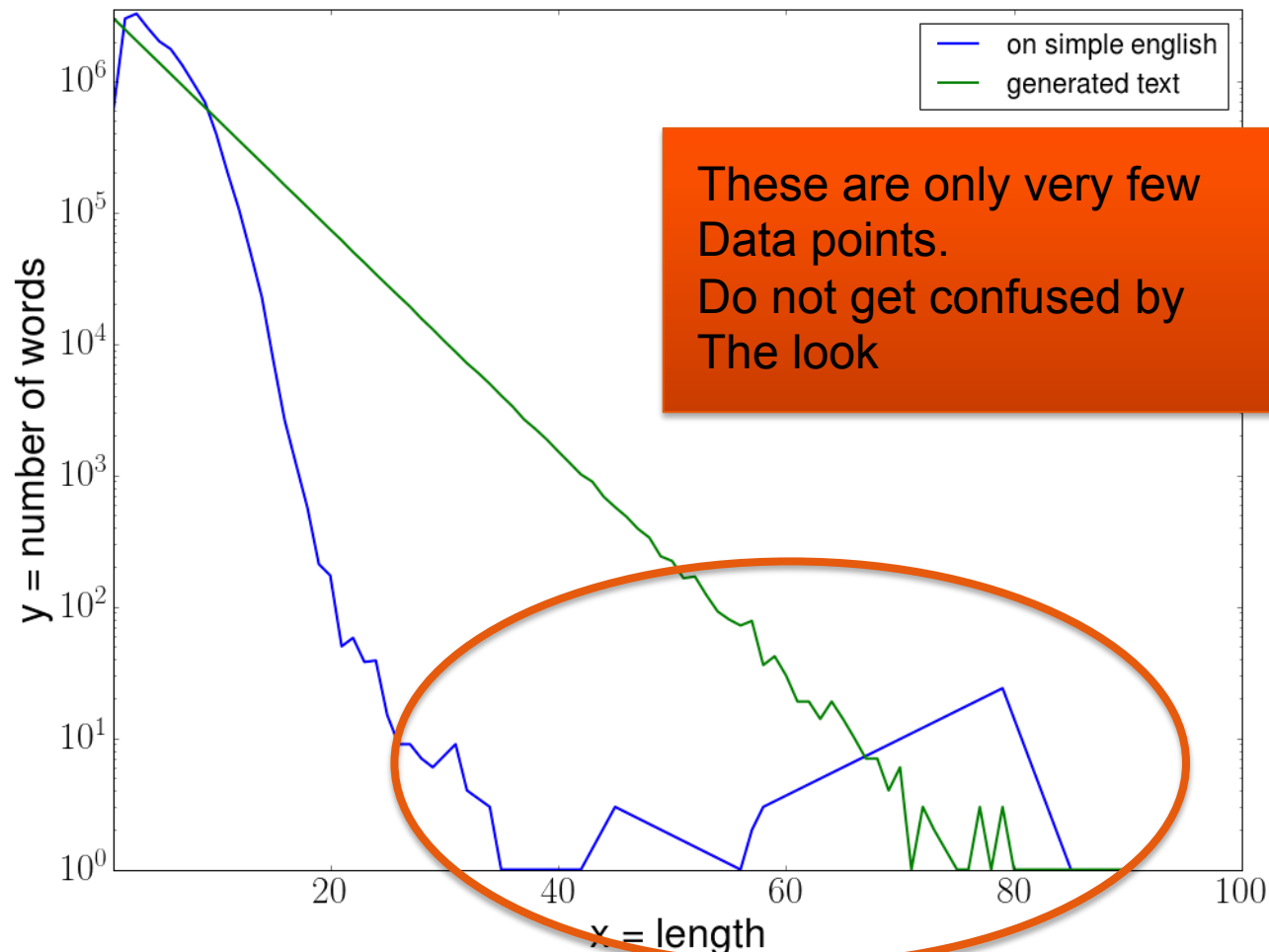
# Generated curve is a perfect line on a plot with log y-axis. Exponential decrease!

Length distribution of Words



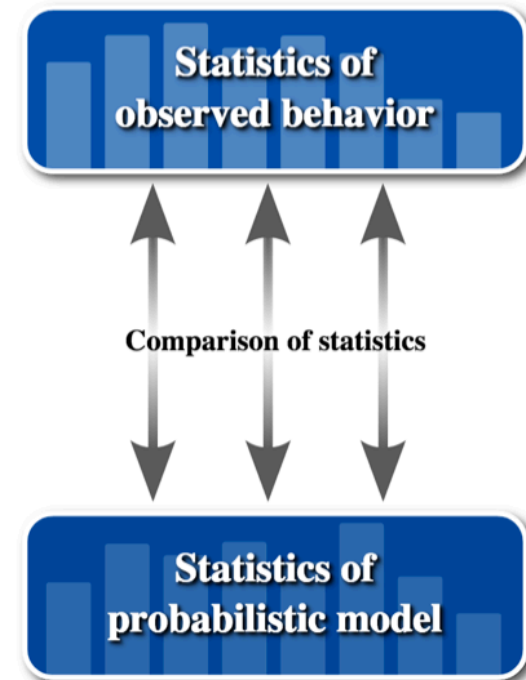
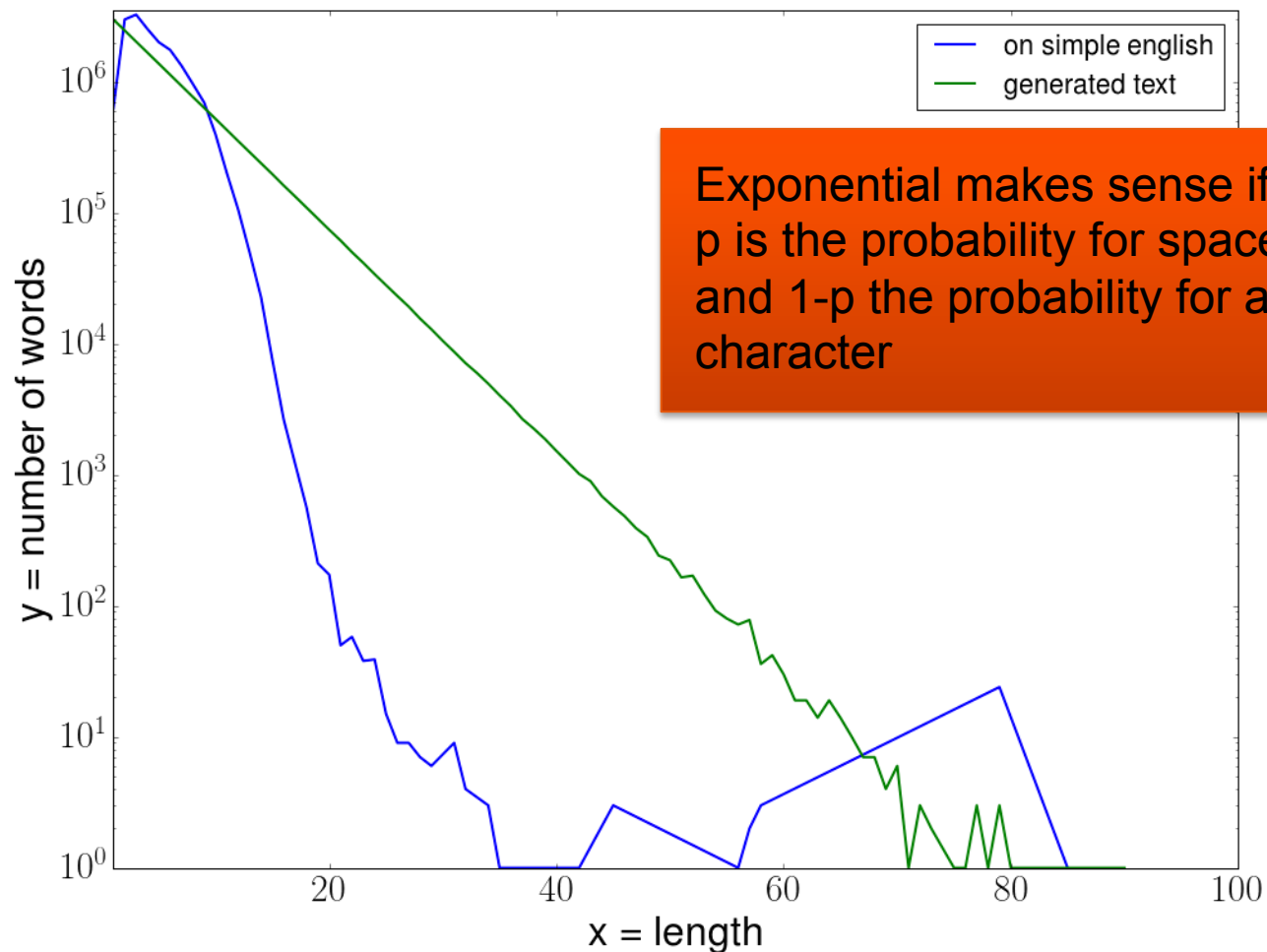
# Generated curve is a perfect line on a plot with log y-axis. Exponential decrease!

Length distribution of Words



$$P(\text{len}(\text{word}) = n) = (1 - p)^n p \sim c^n, 0 < c < 1$$

Length distribution of Words





# Thank you for your attention!



Contact:

Rene Pickhardt  
Institute for Web Science and Technologies  
Universität Koblenz-Landau  
[rpickhardt@uni-koblenz.de](mailto:rpickhardt@uni-koblenz.de)

**WeST**   
People and Knowledge Networks



# Copyright:

- **This Slide deck is licensed under creative commons 3.0. share alike attribution license. It was created by Rene Pickhardt. You can use share and modify this slide deck as long as you attribute the author and keep the same license.**
- By ArchonMagnus (Own work) [CC BY-SA 4.0 (<http://creativecommons.org/licenses/by-sa/4.0>)], via Wikimedia Commons