

Natural Language Processing

Lecture 05

Dirk Hovy

dirk.hovy@unibocconi.it

 @dirk_hovy

Today's Goals

- Introduce *n*-grams
- Learn how to compute *n*-gram probabilities with **Maximum Likelihood Estimation**
- Understand the **Markov assumption**
- Understand the effect of **smoothing**
- Learn how to use probabilistic language models for **inference** and **generation**
- Learn about **bag of words** (BOW) representations
- Learn about forms of **TF-IDF** and its possibilities

N-grams

"As Gregor Samsa awoke one morning from uneasy dreams, he found himself transformed in his bed into a gigantic insect-like creature."

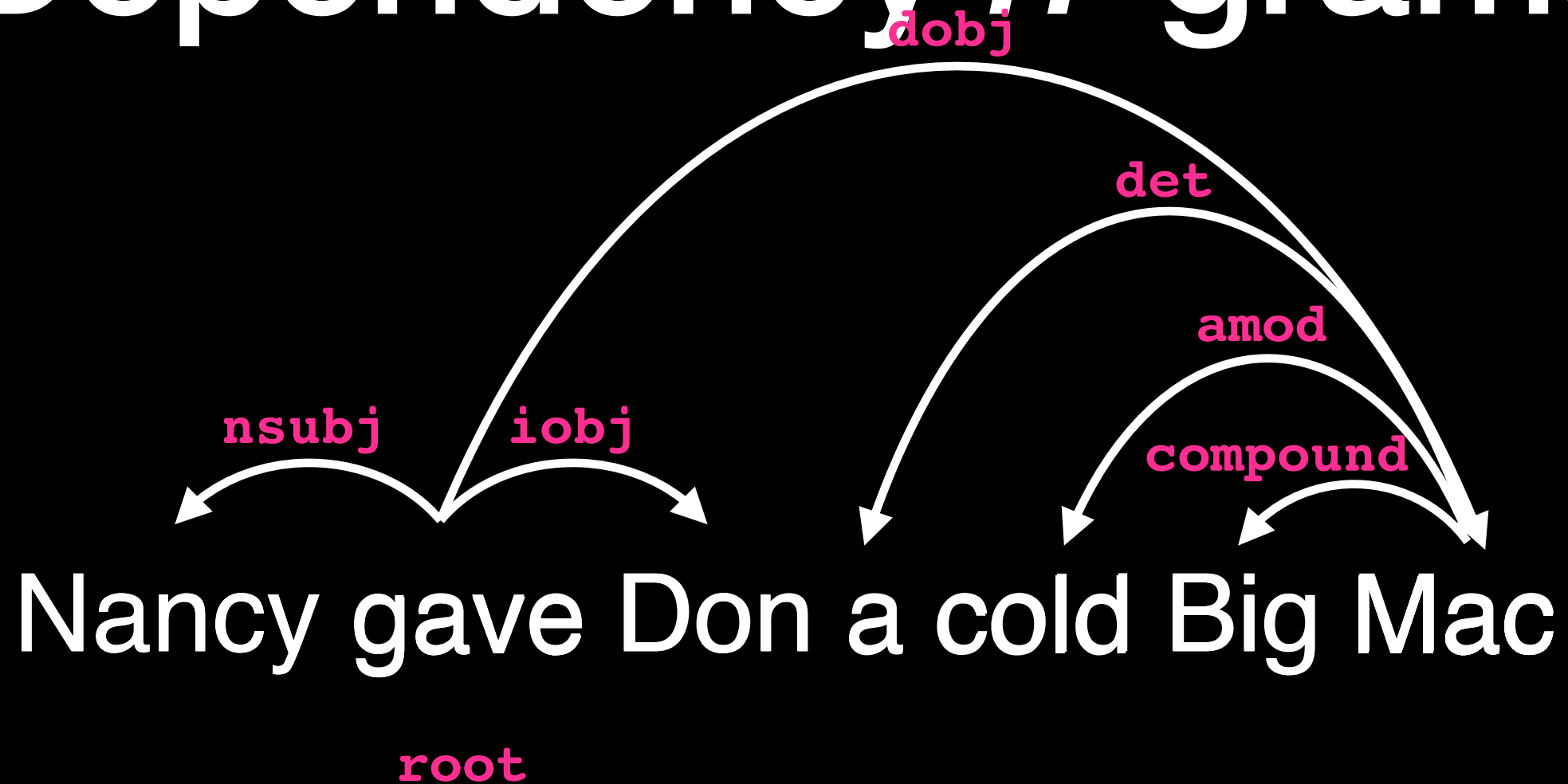
Unigrams As, Gregor, Samsa, awoke, one, morning, from, uneasy, dreams, ...

Bigrams As_Gregor, Gregor_Samsa, Samsa_awoke, awoke_one, one_morning, ...

Trigrams As_Gregor_Samsa, Gregor_Samsa_awoke, Samsa_awoke_one, awoke_one_morning, ...

4-grams As_Gregor_Samsa_awoke, Gregor_Samsa_awoke_one, Samsa_awoke_one_morning, ...

Dependency *n*-grams



Using n-grams in Language Models

What are LMs?

Ranking Sentences

LANGUAGE MODEL

$P(S)$

- I love to models language 0.034
- I love language models 0.62
- I love to language model 0.48

MOST LIKELY TO BE OBSERVED

Machine Translation



Text Generation



In a shocking finding, scientist discovered a herd of unicorns living in a remote, previously unexplored valley, in the Andes Mountains. Even more surprising to the researchers was the fact that the unicorns spoke perfect English.

GPT-3, WILL COME BACK LATER

The scientist named the population, after their distinctive horn, Ovid's Unicorn. These four-horned, silver-white unicorns were previously unknown to science.

Now, after almost two centuries, the mystery of what sparked this odd phenomenon is finally solved.

Dr. Jorge Pérez, an evolutionary biologist from the University of La Paz, and several companions, were exploring the Andes Mountains when they found a small valley, with no other animals or humans. Pérez noticed that the valley had what appeared to be a natural fountain, surrounded by two peaks of rock and silver snow.

Pérez and the others then ventured further into the valley. "By the time we reached the top of one peak, the water looked blue, with some crystals on top," said Pérez.

Language Models in Short

1. Break sentence into n -grams
2. Increase their counts
3. Compute probabilities
4. Multiply them together

Probability of a Sentence?

HOW OFTEN WE

HAVE SEEN SENTENCE S

$$P(S) = \frac{c(S)}{\sum_{Z \in \mathcal{Z}} c(Z)}$$

...ALL POSSIBLE SENTENCES

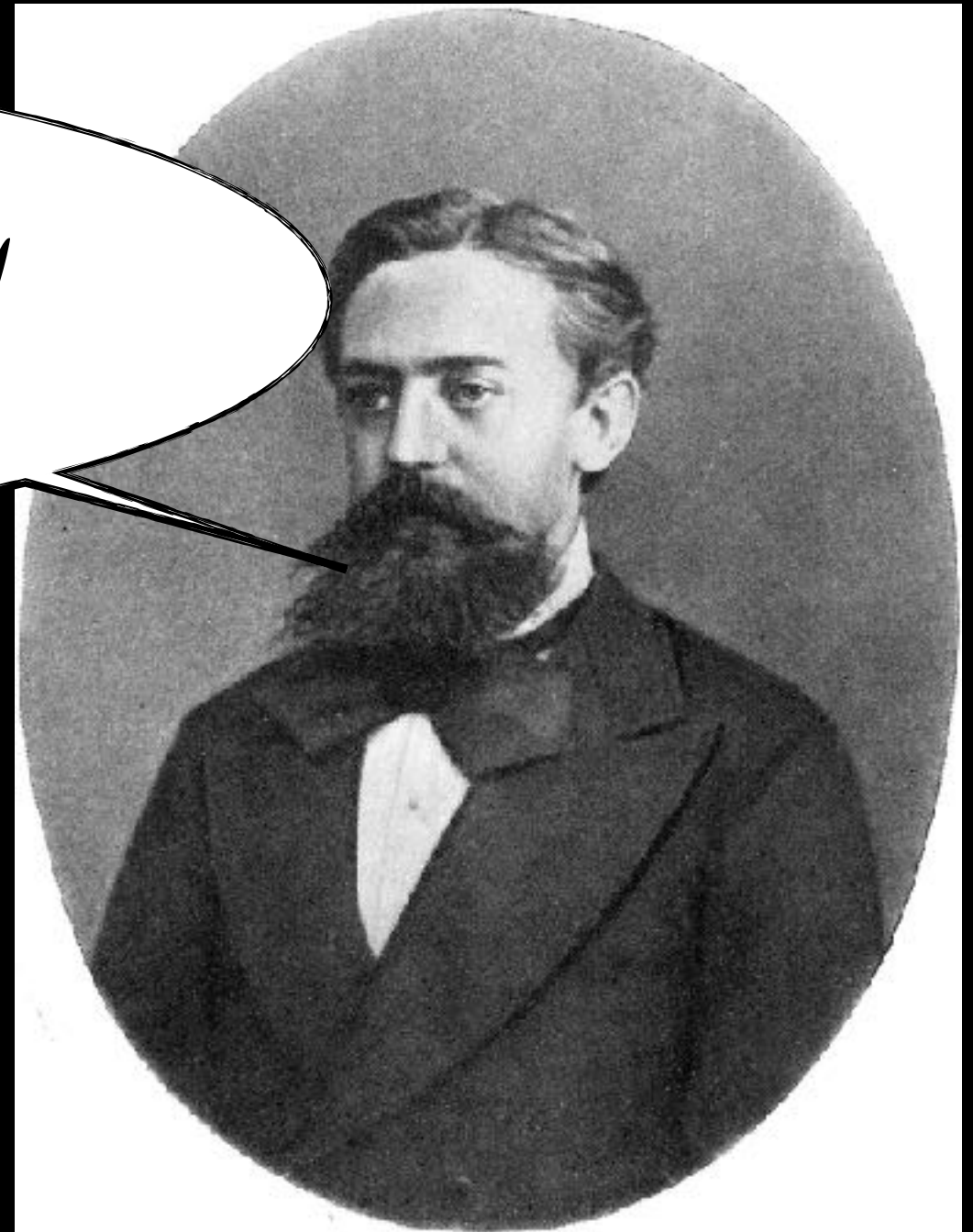
\times

$$= \frac{\quad}{\text{INFINITY}} = 0$$

Can We Make it Simpler?

BREAK IT DOWN!

$P(S) = P(w_1, w_2, \dots, w_n)$
*JOINT PROBABILITY
OF ALL THE WORDS*



Andrey Andreyevich Markov
(1856 – 1922)

Markov Assumption

BREAKING IT DOWN:

$$P(w_1, w_2, \dots, w_n) = \prod_{i=1}^N P(w_i | w_1, \dots, w_{i-1})$$

HISTORY

LIMITING THE HISTORY:

$$P(w_1, w_2, \dots, w_n) = \prod_{i=1}^N P(w_i | w_{i-k}, \dots, w_{i-1})$$

Markov Models:

UNIGRAM MODEL ($K=0$) N

$$P(w_1, w_2, \dots, w_n) \approx \prod_{i=1}^N P(w_i)$$

BIGRAM MODEL ($K=1$) N

$$P(w_1, w_2, \dots, w_n) \approx \prod_{i=1}^N P(w_i | w_{i-1})$$

TRIGRAM MODEL ($K=2$) N

$$P(w_1, w_2, \dots, w_n) \approx \prod_{i=1}^N P(w_i | w_{i-2}, w_{i-1})$$

A Trigram Model

* * The weather today is fine STOP

$$\begin{aligned} P(S) = P(w_1, \dots, w_n) = & P(\text{The} | * *) \\ & \times P(\text{weather} | * \text{The}) \\ & \times P(\text{today} | \text{The weather}) \\ \text{CHAIN RULE} & \times P(\text{is} | \text{weather today}) \\ & \times P(\text{fine} | \text{today is}) \\ & \times P(\text{STOP} | \text{is fine}) \end{aligned}$$

Count in 57m Tweets

MAXIMUM LIKELIHOOD ESTIMATION

12	The	weather	today	is	just
9	The	weather	today	is	so
9	the	weather	today	is	slightly
8	The	weather	today	is	perfect
5	The	weather	today	is	beautiful
4	The	weather	today	is	slightly
3	the	weather	today	is	so
3	the	weather	today	is	perfect
3	The	weather	today	is	nearly
3	the	weather	today	is	bitter
3	The	weather	today	is	absolutely
2	The	weather	today	is	wonderful
2	The	weather	today	is	beyond
2	The	weather	today	is	amazing
2	The	weather	today	is	a
1	the	weather	today	is	worth
1	the	weather	today	is	weird
1	The	weather	today	is	too
1	the	weather	today	is	the
1	The	weather	today	is	that
1	the	weather	today	is	that
1	The	weather	today	is	splendid
1	THE	WEATHER	TODAY	IS	SO
1	the	weather	today	is	simply
1	The	weather	today	is	sickening
1	The	weather	today	is	seriously
1	The	weather	today	is	pretty
1	the	weather	today	is	pretty
1	The	weather	today	is	Perrfff
1	the	weather	today	is	PERFECT

(MLE)

Dealing with the Unknown: Smoothing

Many Counts are Still 0

* * The weather today is fine **STOP**

$$\begin{aligned} P(S) = P(w_1, \dots, w_n) &= P(\textit{The} | * *) \\ &\times P(\textit{weather} | * \textit{The}) \\ &\times P(\textit{today} | \textit{The weather}) \\ &\times P(\textit{is} | \textit{weather today}) \\ c(\textit{fine}) &= 0 \times P(\textit{fine} | \textit{today is}) \\ &\times P(\textit{STOP} | \textit{is fine}) \end{aligned}$$

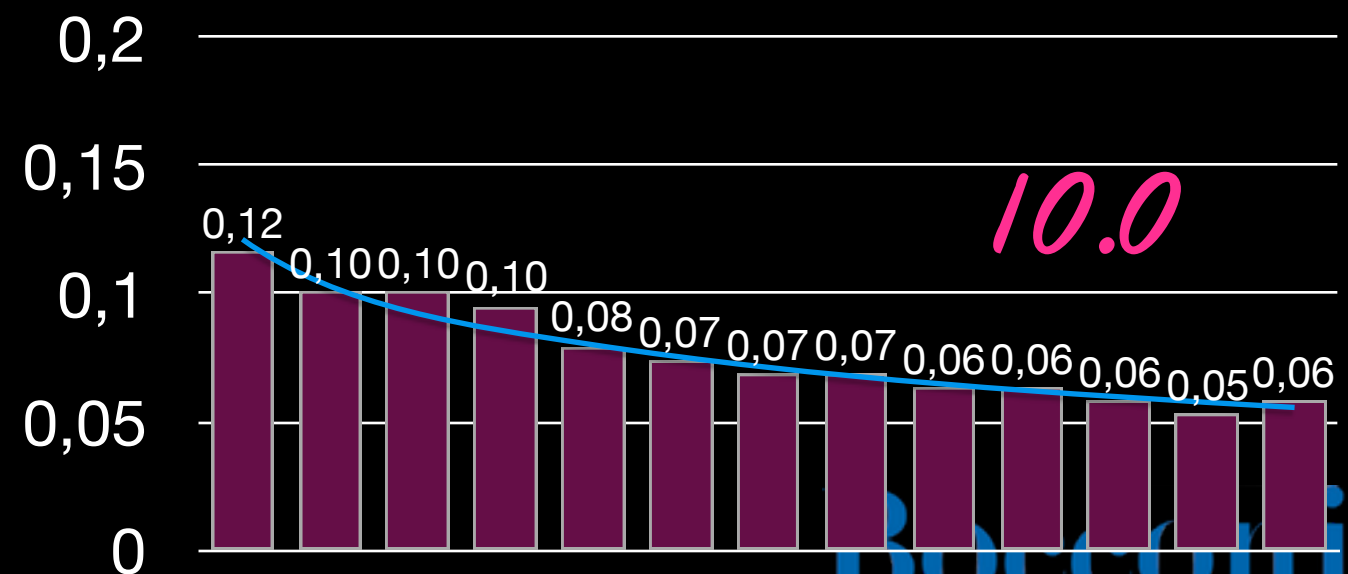
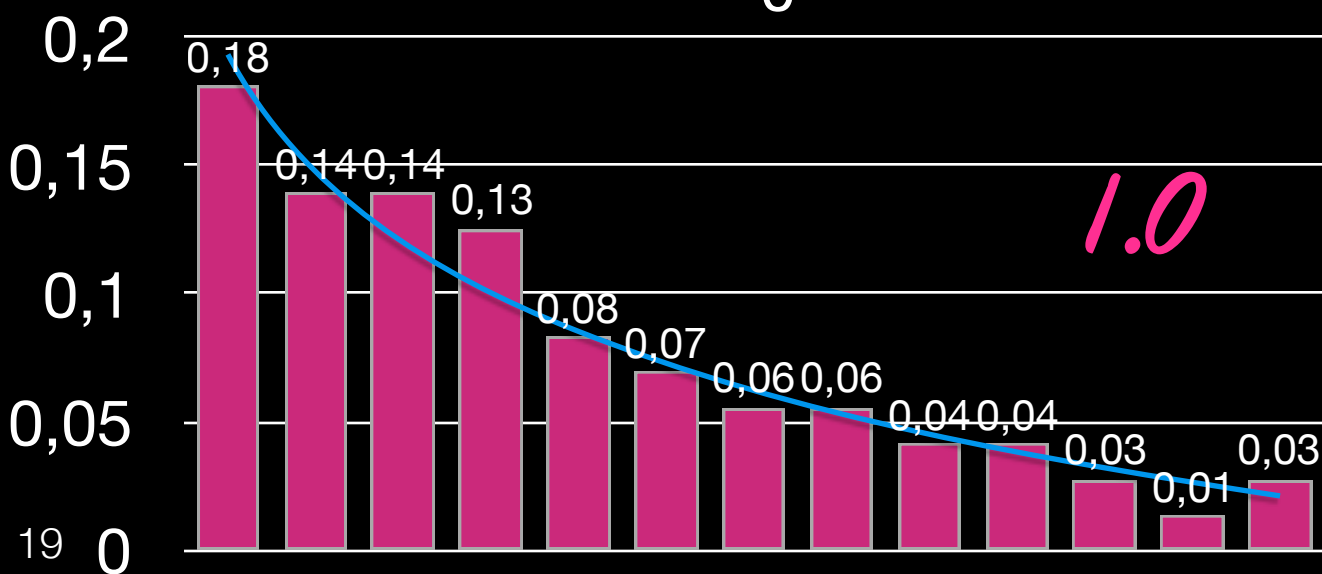
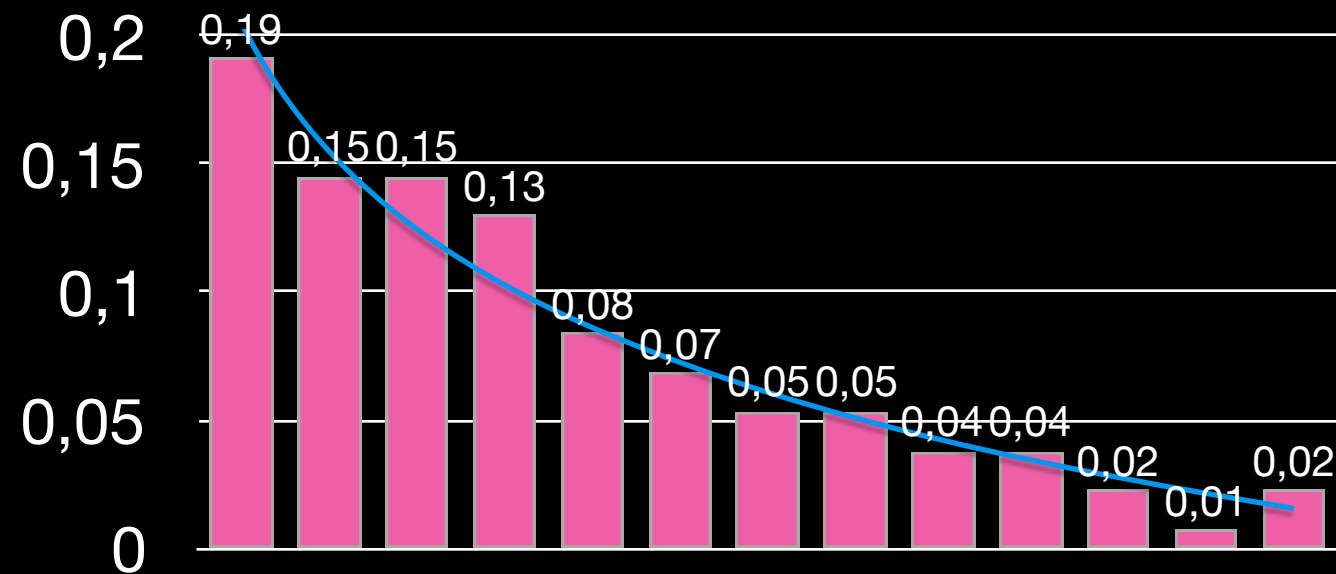
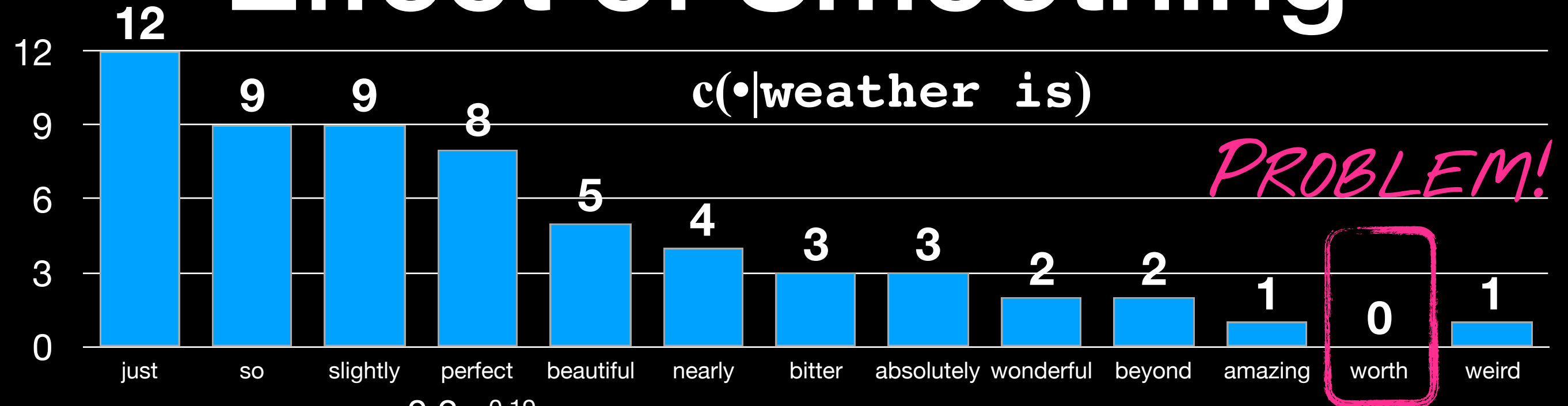
Add-one (Laplace) smoothing

*JUST PRETEND
YOU'VE SEEN IT!*



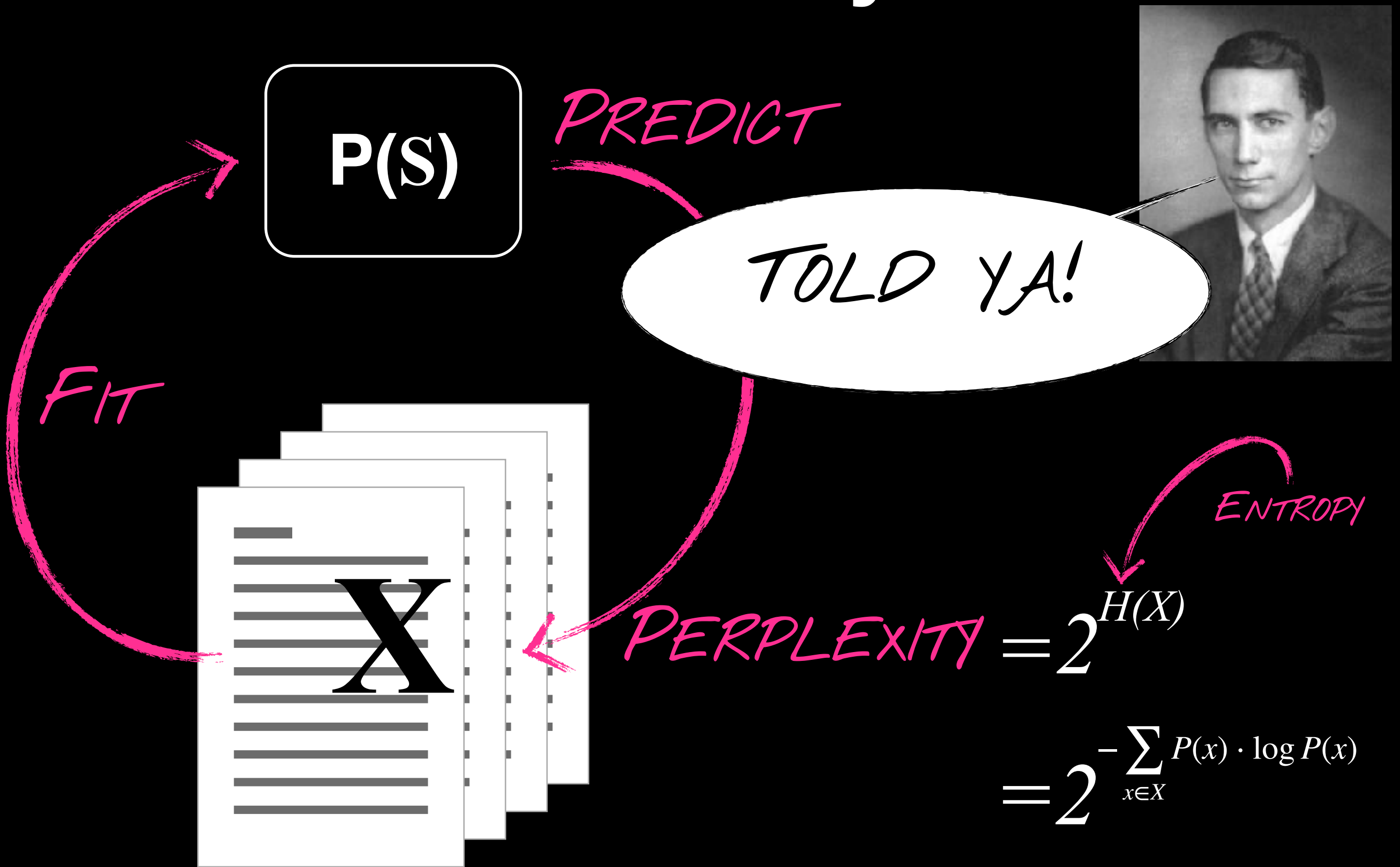
Pierre-Simon, marquis de Laplace
(1749 – 1827)

Effect of Smoothing



Evaluating LMs

How Good is My Model?



Using LMs for Generation

Take a Random Walk

PROPORTIONATELY



Pick a random word w from $P(\cdot \mid * *)$

$H = [* , * , w]$

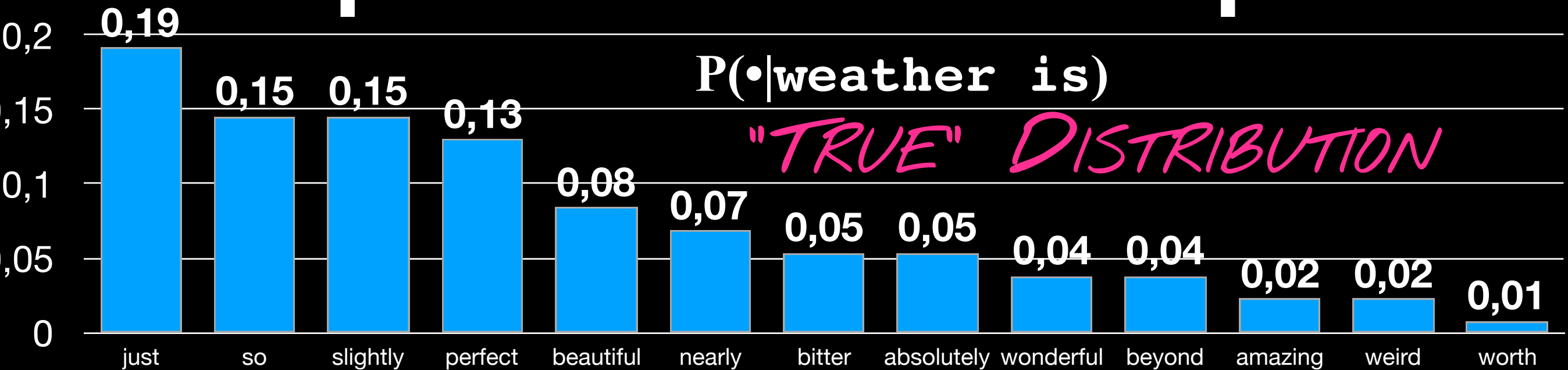
While $H[-1]$ is not STOP:

 Pick a random word w from $P(\cdot \mid H[-2:])$

$H += [w]$

return H

Proportionate Samples



20

15

10

5

0

1 SAMPLE

20

15

10

5

0

10 SAMPLES

20

15

10

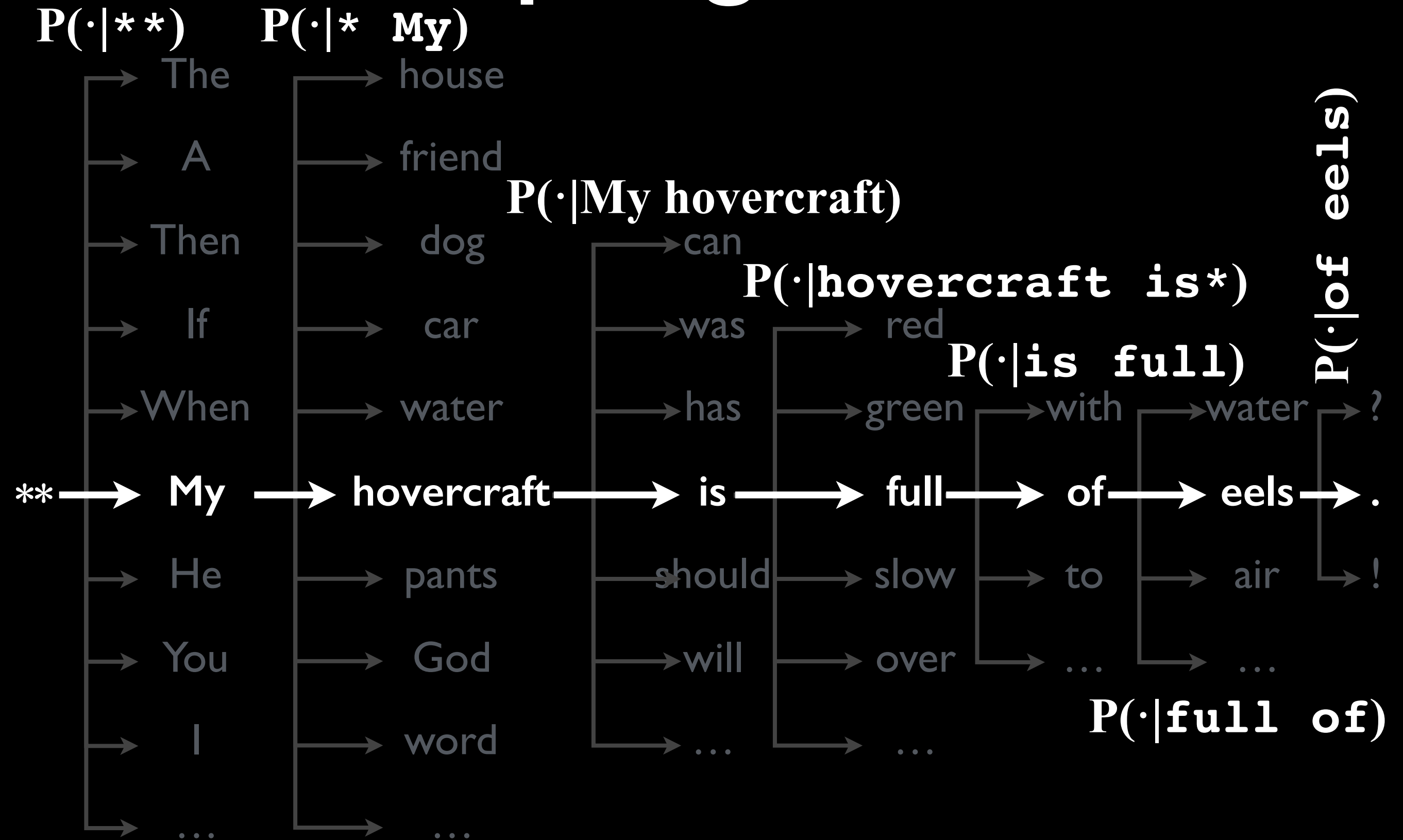
5

0

100 SAMPLES

24

Sampling Words



Language Model in Short

1. Break sentence into n -grams *MARKOV HORIZON*
2. Increase their counts *SMOOTHING*
3. Compute probabilities *MLE*
4. Multiply them together *MARKOV ASSUMPTION*

Using n-grams in Text Analysis

Ham or Spam?

From: offr4u@rsph.com
Subject: Unique wealth offerings
To: dirk.hovy@unibocconi.it

Greetings dear friend

We have an amazing offer 4U: Click here to get access to a free consultation for serious wealth benefits! Urgent: offer expires soon.

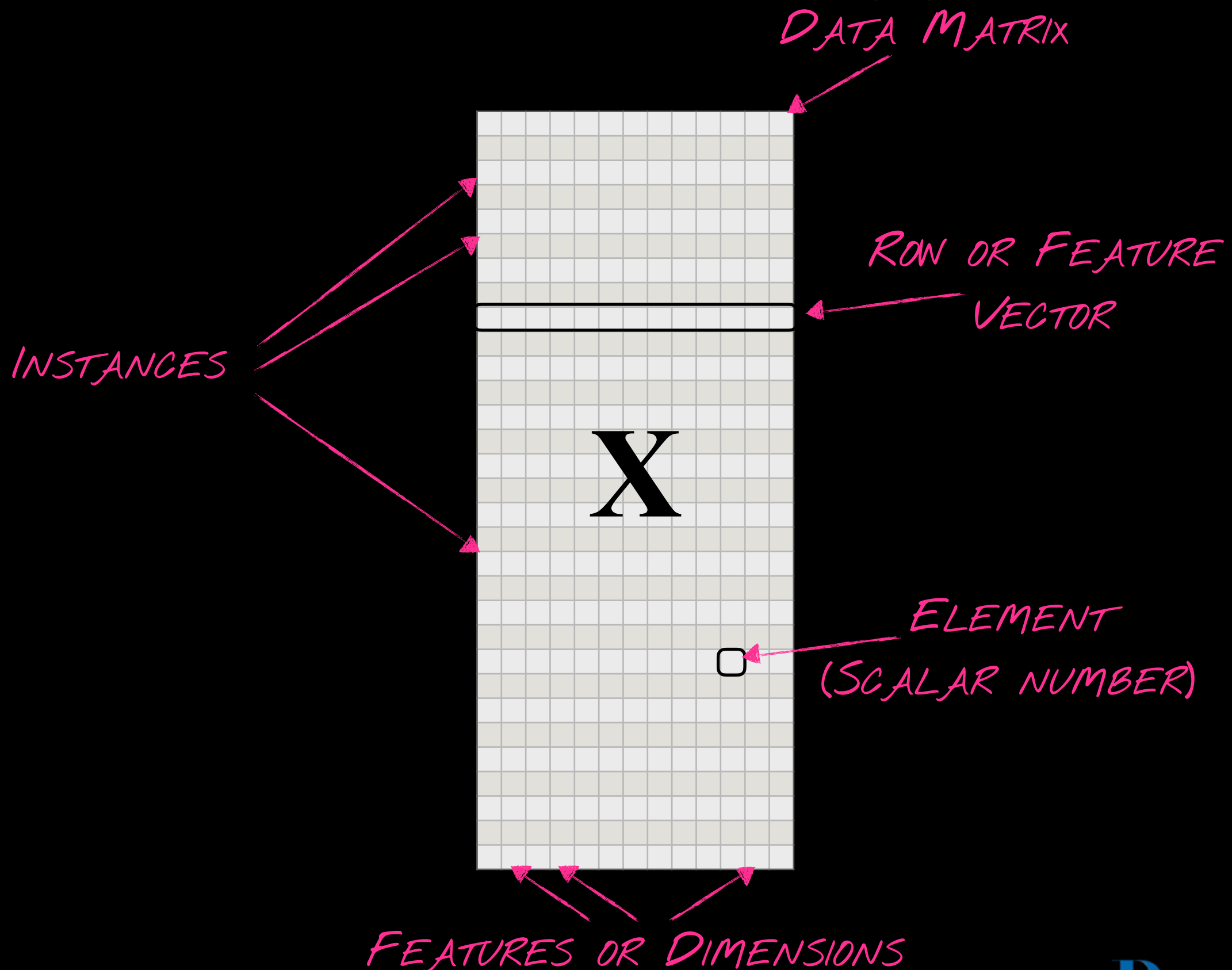
Works guaranteed! Triple your income.

Spam terms:

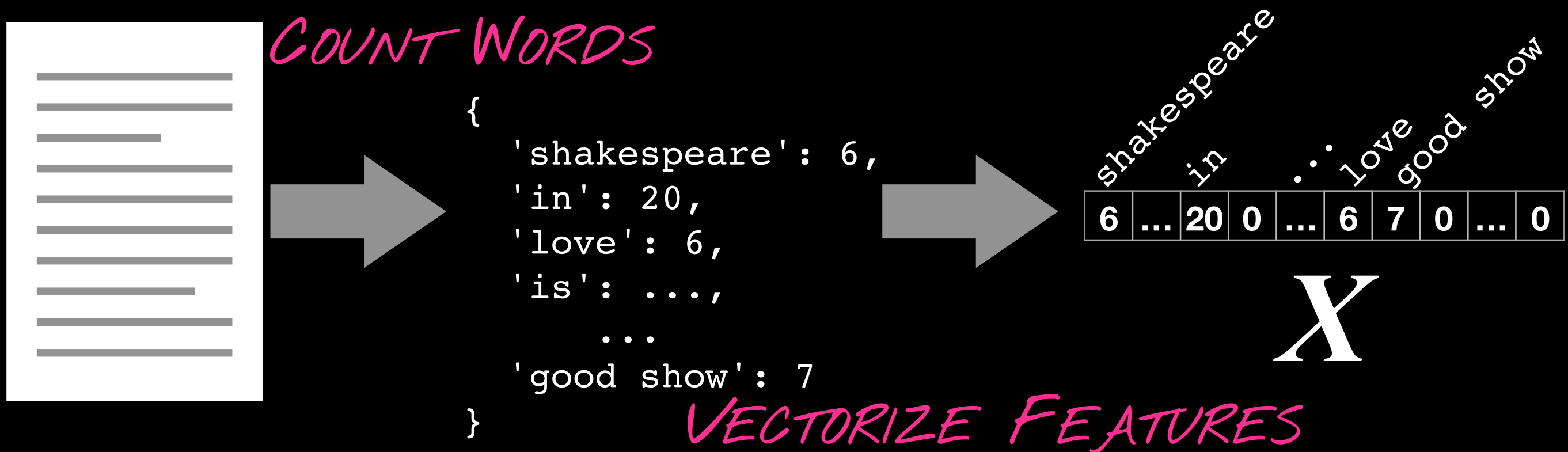
- 4U
- click
- amazing
- free
- guarantee
- offer
- urgent
- dear friend
- income
- serious

Discrete Representations

Terminology



Bags of words (BOW)



Quiz!

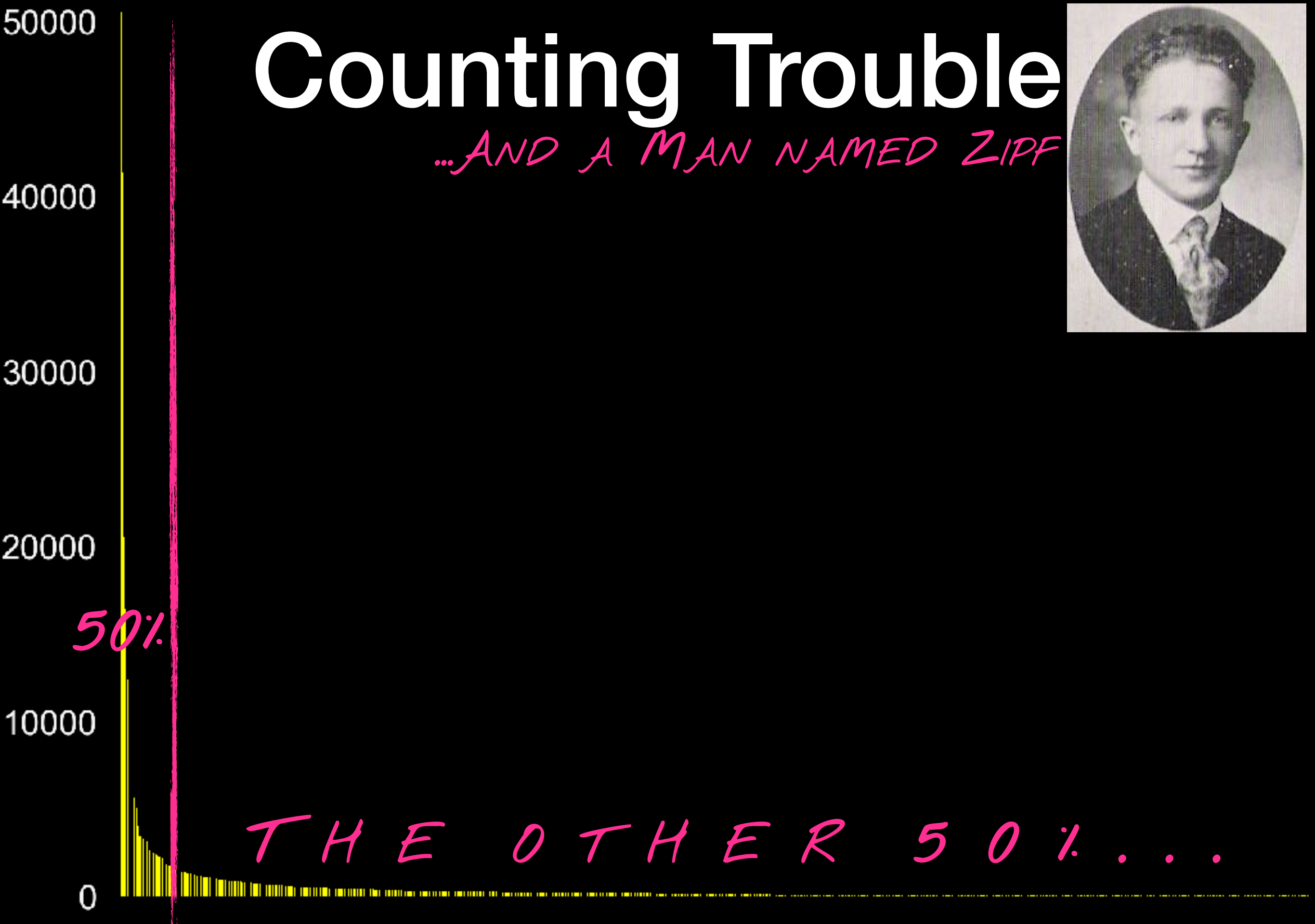
What happens if we allow every possible word to constitute a feature?

Expensive computation, and vectors have too many zeros.

Limit to most frequent/informative words!

Counting Trouble

...AND A MAN NAMED ZIPF



Finding Important Words: TF-IDF

Some Words are Just More Interesting...

the _____

the _____

the _____

_____ the _____

_____ the _____

the _____

the _____

the _____

sustainable _____

sustainable _____

the _____

the _____

sustainable _____

the _____

the _____

the _____

the _____

the _____

the _____

the _____

Karen Spärck Jones

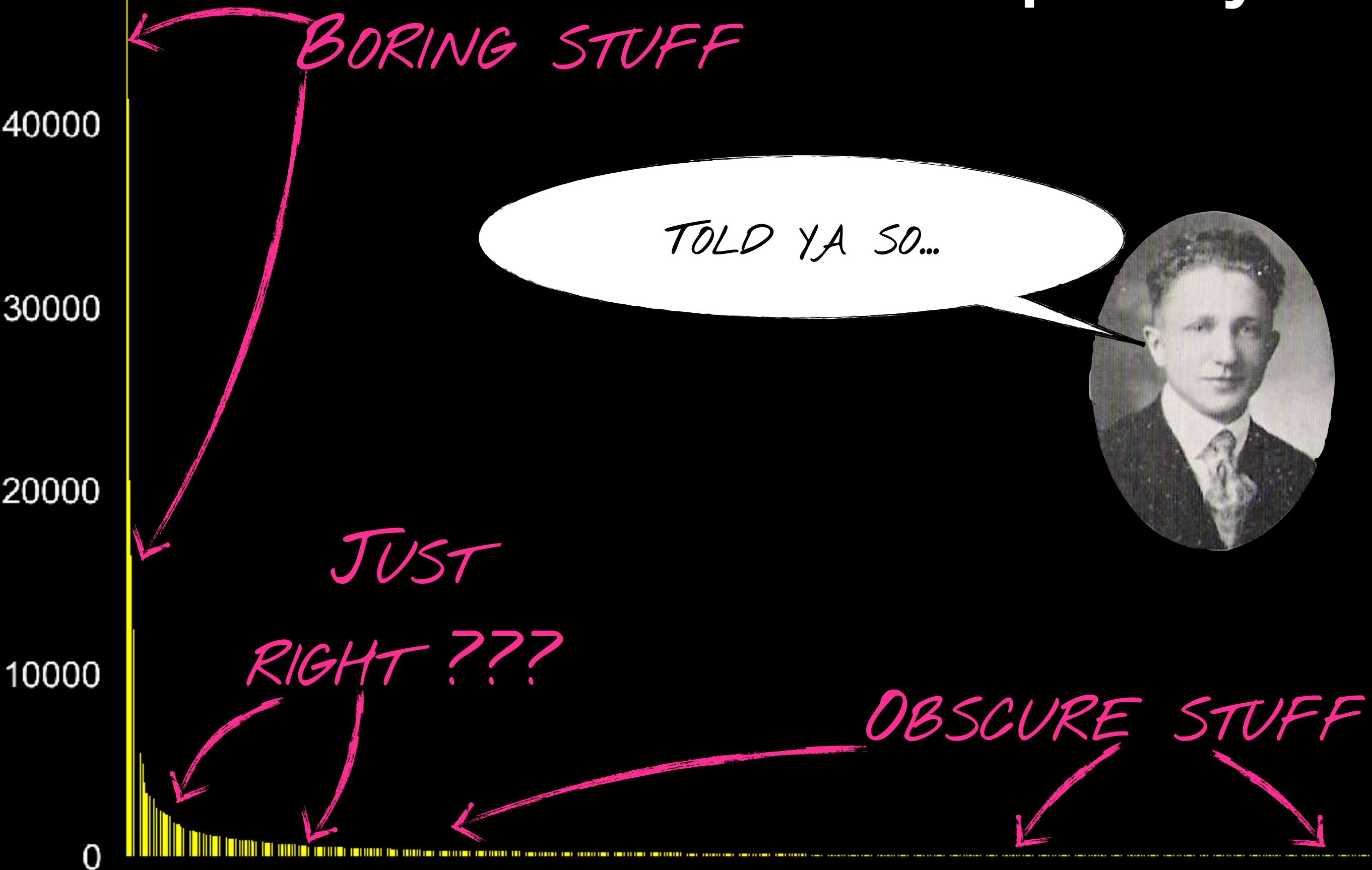
THANKS, KAREN!

1935–2007

- Became a teacher before starting CS career at Cambridge
- Laid the foundation for modern NLP, Google Search, text classification
- Campaigned for more women in CS
- Namesake of prestigious CS prize



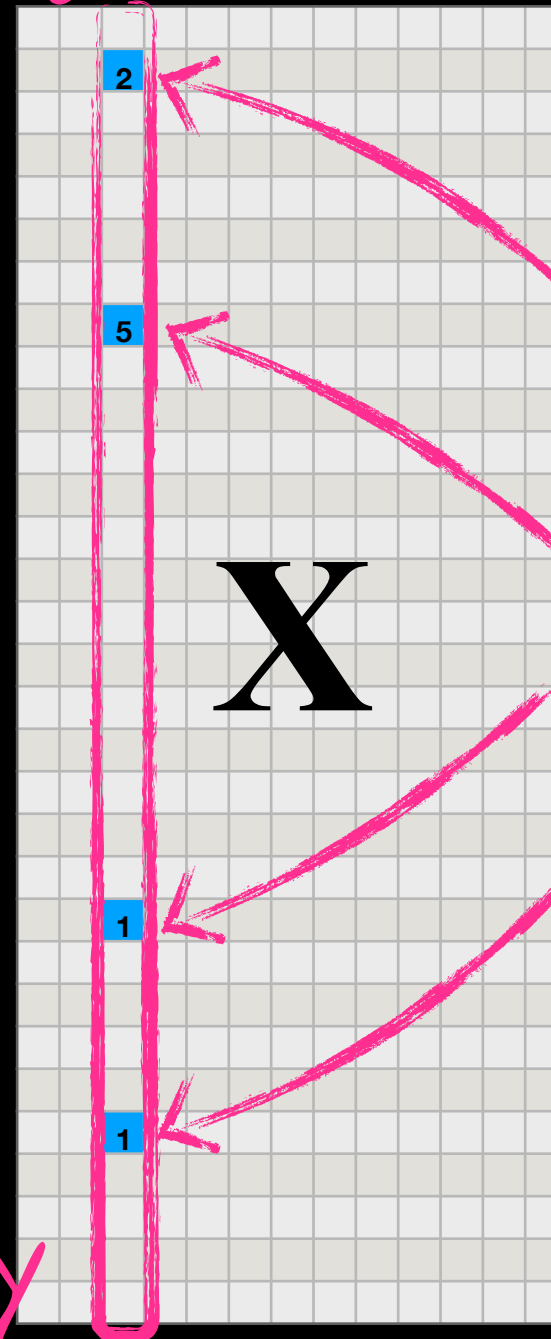
Problems with Term Frequency



Document and Term Frequency

FEATURE

$$IDF = \log \frac{N}{df(w)}$$



DOCUMENT
FREQUENCY
(COUNT): 4

TERM FREQUENCY
(SUM): 9 TF

Putting it Together

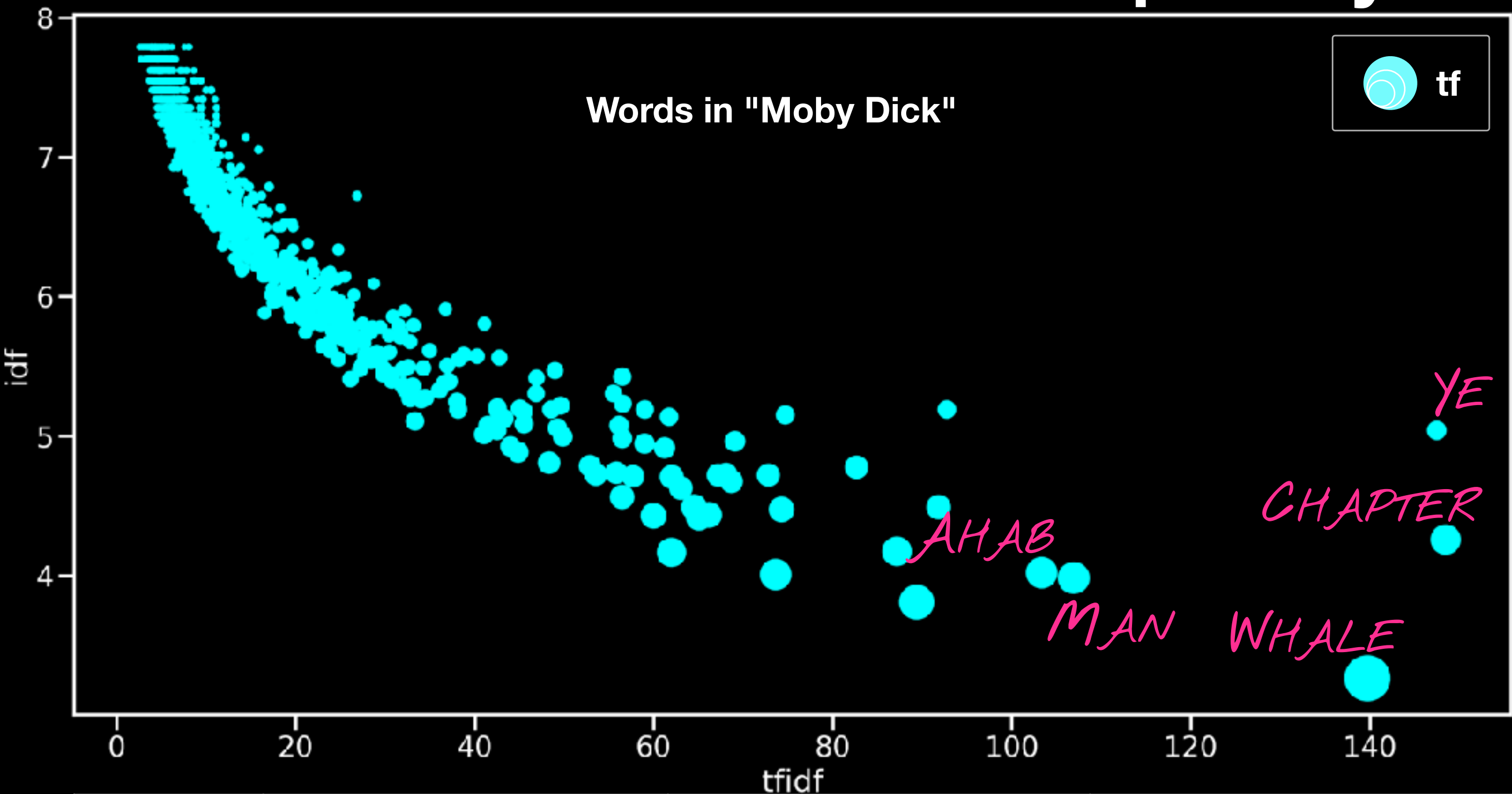
HOW OFTEN WE
SAW THE WORD

$$\textbf{TFIDF}(w) = \textbf{TF}(w) \cdot \log \frac{N}{df(w)}$$

ADJUSTED BY
HOW MANY
DOCUMENTS

0

Document and Term Frequency



word	tf	idf	tfidf
ye	467	4.257380	148.497079
chapter	171	5.039475	147.504638
whale	1150	3.262357	139.755743
man	525	3.982412	106.932953
ahab	511	4.019453	103.357774

Variants

	TF
binary	<i>1 if word in D, else 0</i>
raw	$c(\text{word}, D)$
relative	$c(\text{word}, D) / \text{len}(D)$
smooth	$\log(c(\text{word}, D) + 1)$

	IDF
regular	$\log \frac{N}{df(\text{word})}$
smooth	$\log \frac{N}{df(\text{word}) + 1} + 1$

Wrapping up

Take home points

- **Language Models** assign a probability to any sentence, can be used for **text generation**
- The **Markov assumption** breaks sentence probability into a **chain** of word conditional probabilities
- **Markov order** determines the size of the conditional n -grams
- **Smoothing** helps address the problem of unseen words
- Words and texts can be represented as **sparse, discrete** feature vectors over counts
- **TF-IDF** finds "bursty" words: medium frequency overall, but concentrated in few documents