

Text Analysis

TFIDF, REs, and Collocations

Dirk Hovy

dirk.hovy@unibocconi.it

@dirk_hovy

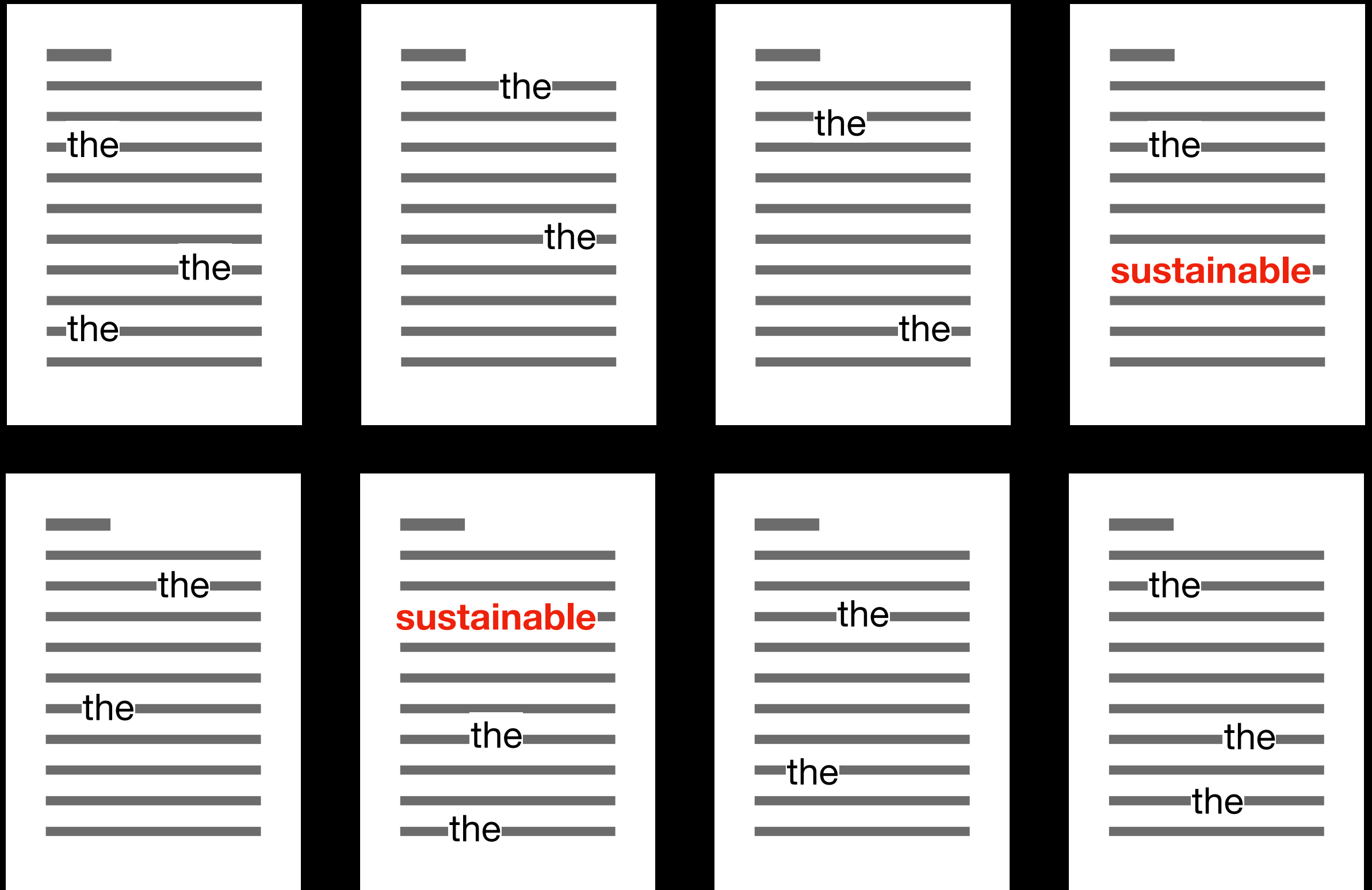


Today's Goals

- Learn about forms of **TF-IDF** and its possibilities
- Know when (and when not) to use **Regular expressions**
- Understand **PMI** and see why it can help find collocations

Finding Important Words: TF-IDF

Some Words are Just More Interesting...



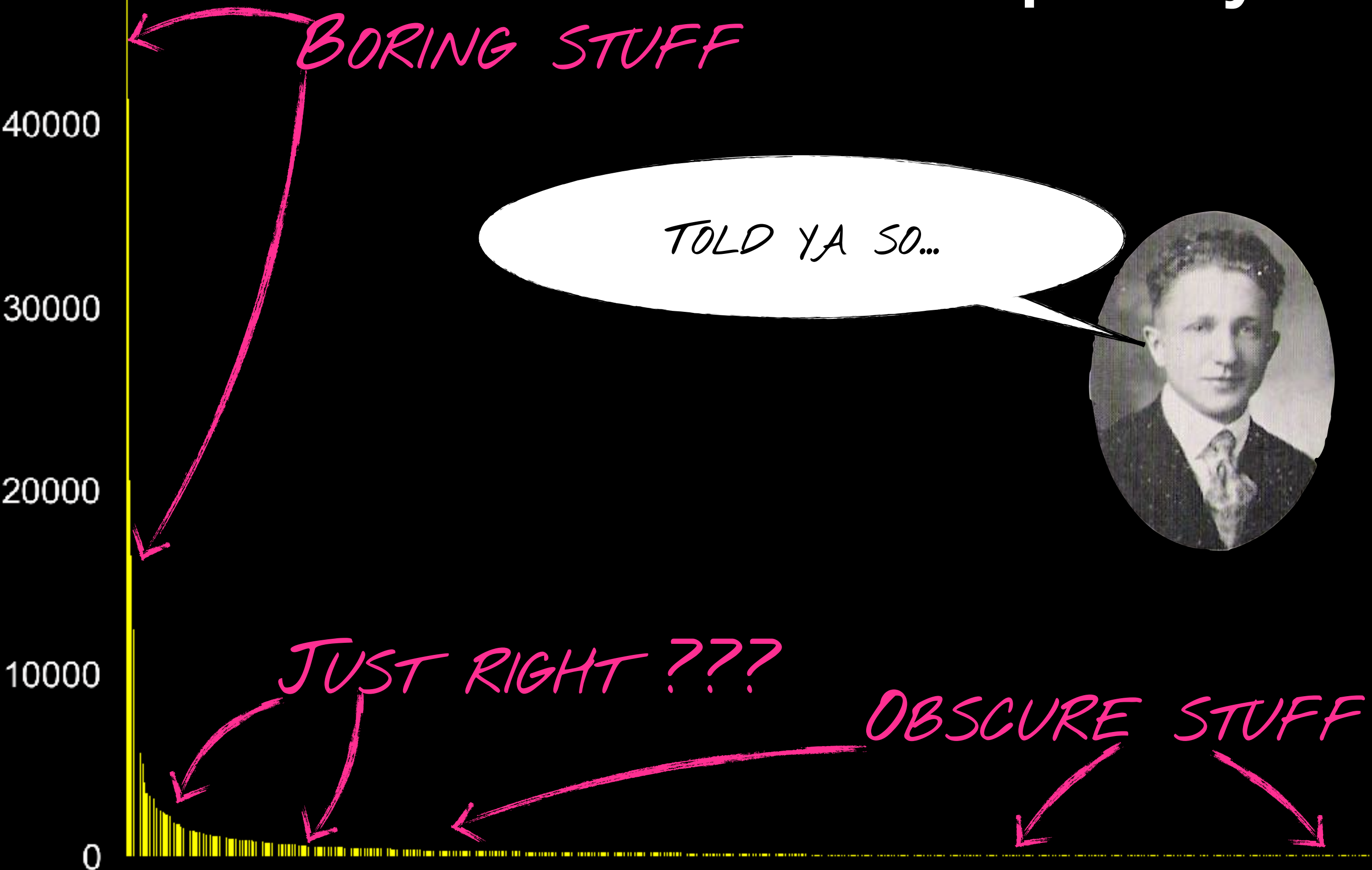
Karen Spärck Jones

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



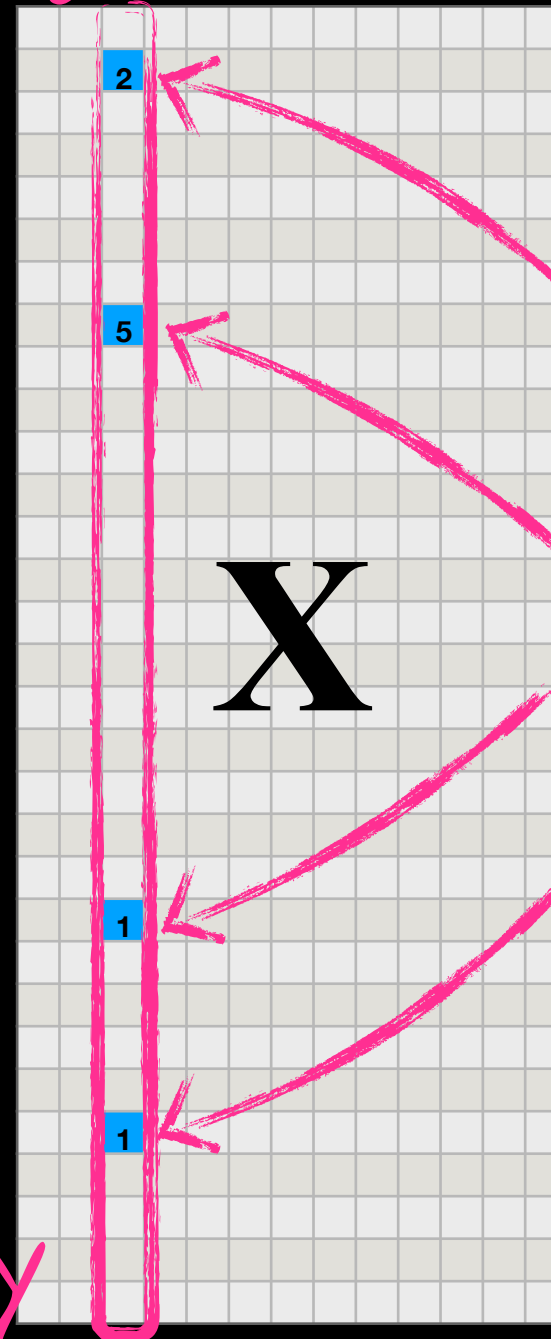
Problems with Term Frequency



Document and Term Frequency

FEATURE

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



X

TERM FREQUENCY
(SUM): 9 TF

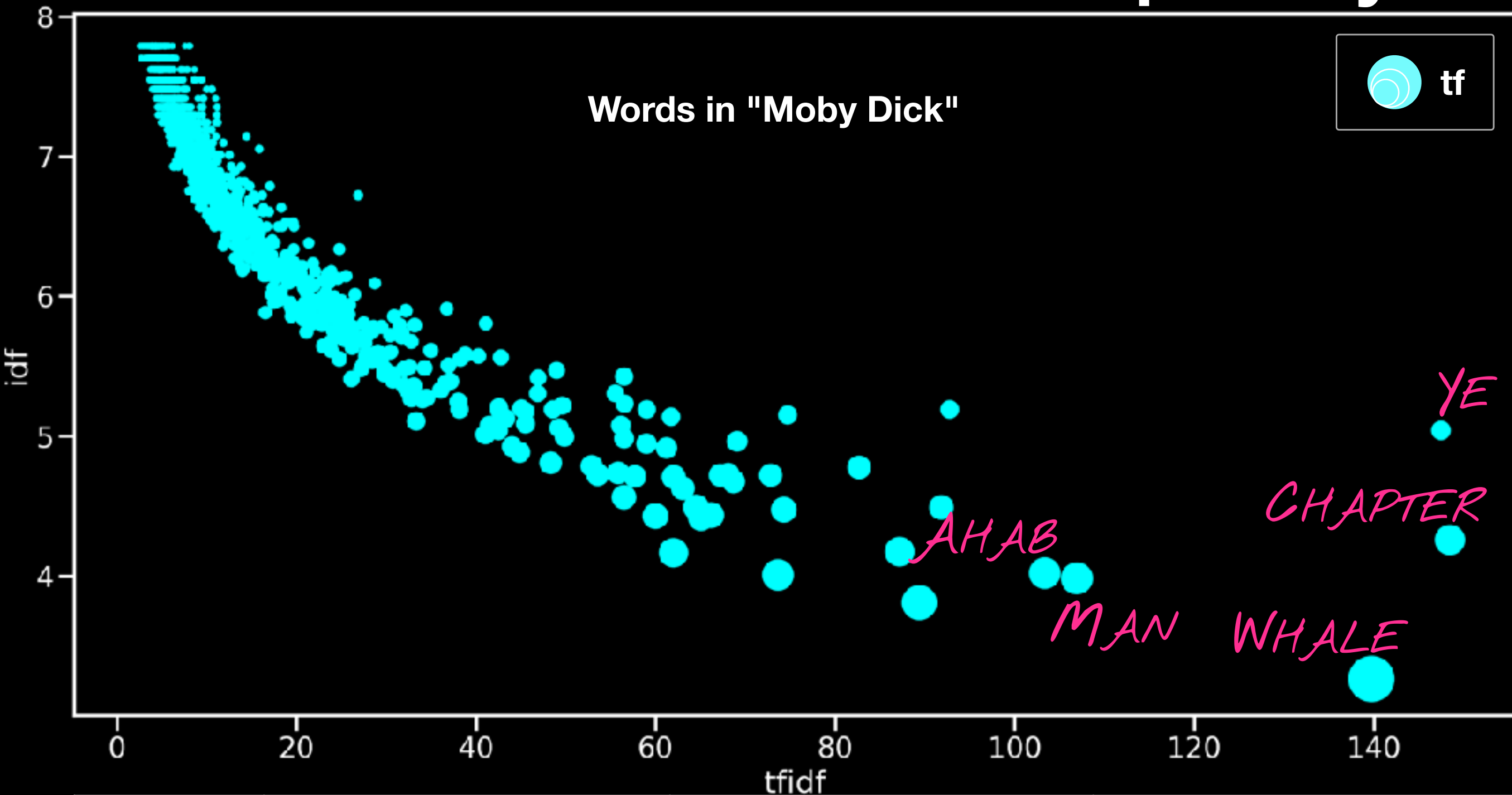
Putting it Together

HOW OFTEN WE
SAW THE WORD

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

ADJUSTED BY
HOW MANY
DOCUMENTS

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

Flexible Matches: Regular Expressions

The promise...

WHENEVER I LEARN A NEW SKILL I CONCOCT ELABORATE FANTASY SCENARIOS WHERE IT LETS ME SAVE THE DAY.

OH NO! THE KILLER MUST HAVE FOLLOWED HER ON VACATION!



BUT TO FIND THEM WE'D HAVE TO SEARCH THROUGH 200 MB OF EMAILS LOOKING FOR SOMETHING FORMATTED LIKE AN ADDRESS!



IT'S HOPELESS!

EVERYBODY STAND BACK.



I KNOW REGULAR EXPRESSIONS.



Is it an (Email) Address?

- notMyFault@webmail.com ✓
- smithie123@gmx ✗
- Free stuff@unibocconi.it ✗
- mark_my_words@hotmail;com ✗
- truthOrDare@webmail.in ✓
- look@me@twitter.com ✗
- how2GetAnts@aol.dfdsfgfdsgfd ✗



Simple Matching

sequence	Matches
e	any single occurrence of e
at	<u>a</u> t, r <u>a</u> t, m <u>a</u> t, s <u>a</u> t, c <u>a</u> t, <u>a</u> ttack, <u>a</u> ttention, l <u>a</u> ter

Quantifiers

	Means	Example	Matches
*	0 or more	cooo*l	cool, coool
+	1 or more	hello+	hello, helloo, hellooooooooo
?	0 or 1	fr?og	fog, frog

Special Characters

	Means	Example	Matches
.	any single character	.e1	ee1, Ne1, ge1
\n	newline character (line break)	\n+	One or more line breaks
\t	a tab stop	\t+	One or more tabs
\d	a single digit [0-9]	B\d	B0, B1, ..., B9
\D	a non-digit	\D.t	' t, But, eat
\w	any alphanumeric character	\w\w\w	Top, WOO, ash, bee, ...
\W	non-alphanumeric character		
\s	a whitespace character		
\S	a non-whitespace character		
\	"Escapes" special characters to match them	.+ \.com	abc.com, united.com
^	the beginning of the input string	^...	First word in line
\$	the end of the input string	^\n\$	Empty line

Classes

	Means	Example	Matches
[abc]	Match any of a, b, c	<code>[bcrms]at</code>	<code>bat, cat, rat, mat, sat</code>
[^abc]	Match anything BUT a, b, c	<code>te[^]+s</code>	<code>tens, tests, teens, texts, terrors...</code>
[a-z]	Match any lowercase character	<code>[a-z][a-z]t</code>	<code>act, ant, not, ... wit</code>
[A-Z]	Match any uppercase character	<code>[A-Z]...</code>	<code>Ahab, Brit, In a, ..., York</code>
[0-9]	Match any digit	<code>DIN A[0-9]</code>	<code>DIN A0, DIN A1, ..., DIN A9</code>

Groups

	Means	Example	Matches
(abc)	Match abc	<code>.(ar).</code>	<code>hard, cart, fare, ..</code>
(ab c)	Match ab OR c	<code>(ab c)ate</code>	<code>abate, cate</code>

Matching Addresses

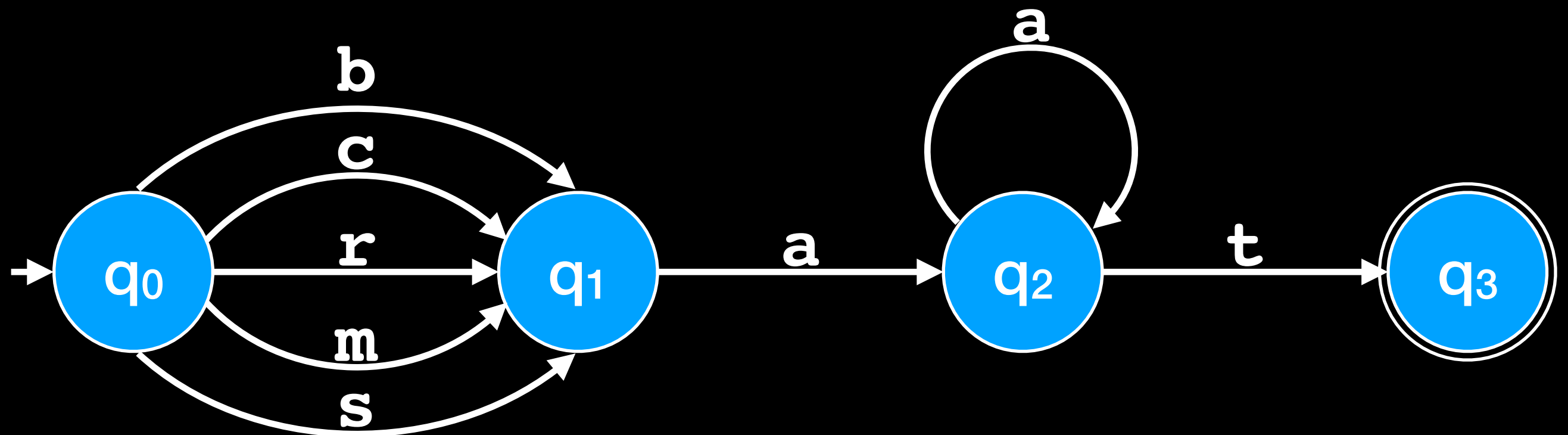


$^{\wedge}[\text{A-Za-z0-9_}\backslash.\text{-}]^{\text{+}}@\text{[A-Za-z0-9_}\backslash.\text{-}]^{\text{+}}\backslash.\text{[A-Za-z0-9_]}[\text{A-Za-z0-9_}]^{\text{+}}\text{\$}$

A (W|w)ord of [Ww]arning



RegEx as Automata



[bcrms] a+t

Telling Neighbors: Pointwise Mutual Information

Some are not like the Others



Mutual Informativity

HOW WELL CAN WE GUESS THE BLANK?

social _____

and _____

_____ media

_____ the

Pointwise Mutual Information

CHANCE OF SEEING THEM TOGETHER

$$PMI(x, y) = \log \frac{P(x, y)}{P(x)P(y)}$$

...SEEING EITHER

x	y	c(x)	c(y)	c(xy)	P(x)	P(y)	P(x, y)	PMI(x; y)
moby	dick	83	83	82	0.0003	0.0003	0.0003	3.48
captain	ahab	327	511	61	0.0013	0.0020	0.0002	1.97
white	whale	280	1150	106	0.0011	0.0045	0.0004	1.93
under	the	119	14175	45	0.0005	0.0553	0.0002	0.83
is	a	1690	4636	110	0.0066	0.0181	0.0004	0.56

$$c(X) = 256,149$$

$$c(XY) = 256,148$$

Wrapping up

Take home points

- **TF-IDF** finds "bursty" words: medium frequency overall, but concentrated in few documents
- **Regular expressions** allow us to search for flexible patterns
- **PMI** tells us how likely one word is to occur with/without another to find **collocations**