## Problems with count vectorizers:

① each word is equidistant in its vector representation:
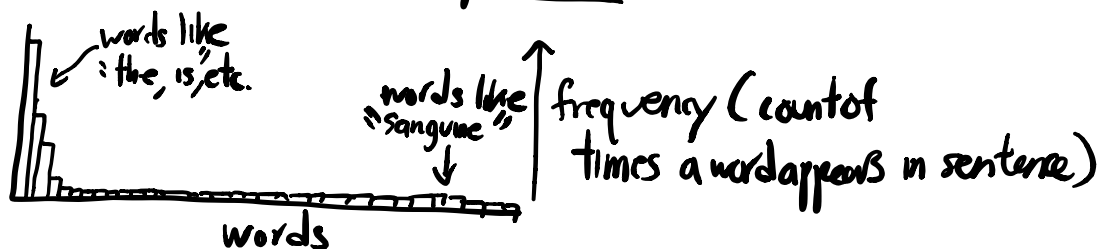
intuitively, we → dog: $[0\ 0\ 1]$
know dog should
be more similar      banana: $[1\ 0\ 0]$
to canine than → canine: $[0\ 1\ 0]$
banana.

all dot products are the same (0).

The angles made by a pair of vectors are all orthogonal:


Canine, banana, 90°, dog

② counts skew towards high frequency stopwords (the, is, on, I, like)
Remember the long tail distribution of words
we saw in Week 1 (Zipf's Law):


words like "the, is" etc.
words like "sanguine"
frequency (count of times a word appears in sentence)
words

## Problems with rule-based domain trees:


color
purple    red    green
violet   mauve   orchid

(ie. hyponyms and hypernyms)

- extremely labor intensive.
- slow to adapt to new words
- difficult to maintain
- "subjective"; two SMEs
(subject matter experts) can have two totally valid trees.

## So... what do we do?

We make a new assumption: the <u>distributional</u>
<u>hypothesis</u>: "You shall know a word by the
company it keeps" — JR Firth, (British linguist).

## Continuous Bag of Words

I used class notes to study for the test.

$$t\text{-}3 \quad t\text{-}2 \quad t\text{-}1 \quad t \quad t\text{+}1 \quad t\text{+}2 \quad t\text{+}3$$

context window $(m=3)$

$P(x_t = \text{study} \mid x_{t\text{-}3} = \text{class}) \times$

$P(x_t = \text{study} \mid x_{t\text{-}2} = \text{notes}) \times$

...

$\boxed{P(x_t = \text{study} \mid x_{t\text{+}3} = \text{test})}$

↖ probability that study is the
target word given that the context
word is test.

How would we find this probability?

|        | ant | art | bad | ... | study... | zebra |
|--------|-----|-----|-----|-----|----------|-------|

**these words represent the targets**

**these words represent the context**

| | ant art bad ... | | | | | |
|---|---|---|---|---|---|---|
| giraffe | 4 | 6 | 2 | | 5 | 7 |
| test | 3 | 2 | 5 | | (9) | 6 |
| zebra | 1 | 2 | 3 | | 4 | 0 |

$$\frac{count(test, study)}{count\ test} =$$

$$\frac{9}{3+2+5+9+6} = \frac{9}{25} = .36$$

However, in practice, we typically do not use CBOW, but rather skipgram, which is the opposite of CBOW:

CBOW

| | Inputs | Target |
|---|---|---|
| | context words | target word |
| Skipgram | target word | context words |

$$P(X_{t+3} = class \mid X_t = study)$$

$$P(X_{t-1} = t_6 \mid X_t = study) \times$$
$$\vdots$$
$$P(X_{t+3} = test \mid X_t = study)$$

$\Rightarrow$

"for each context word in the window"

$$\prod_{\substack{-M \le m \le M \\ m \ne 0}} P(X_{t+m} \mid X_t)$$

$M = 3$

context window size is 3.

We do this for **each** target word:

for each target word

for each context word

$$\prod_{t=1}^{T} \prod_{\substack{-M \le m \le M \\ m \ne 0}} P(X_{t+m} \mid X_t, \theta)$$

context word

target word

Word vectors (our parameters, will learn more)

This becomes our objective function that we optimize ( the higher this value the more "correct" our model is):

$$J(\theta) = \prod_{t=1}^{T} \prod_{\substack{-M \leq m \leq M \\ m \neq 0}} P(X_{t+m} | X_t, \theta)$$

So what exactly is $\theta$? It's actually two matrices:

context word vector          target word vector

ant
art
$$\begin{bmatrix} 4.3 & 0.1 & 0.4 \\ \\ \\ \boxed{-0.9 \;\; 1.2 \;\; 3.4 \ldots} \\ \\ \\ -0.2 \;\; 1.2 \;\; 2.1 \end{bmatrix}$$

study's context vector

study

zebra

$$\underset{D = 300}{\longleftrightarrow}$$

$$\begin{bmatrix} 1.2 & -1.1 & -3.2 \\ \\ \\ \boxed{1.2 \quad 0.4 \quad 1.5} \\ \\ \\ -1.1 \;\; -2.1 \;\; 3.2 \end{bmatrix}$$

study's target vector

$$\underset{D = 300}{\longleftrightarrow}$$

Each word has __both__ a __context__ vector and a __target__ vector.