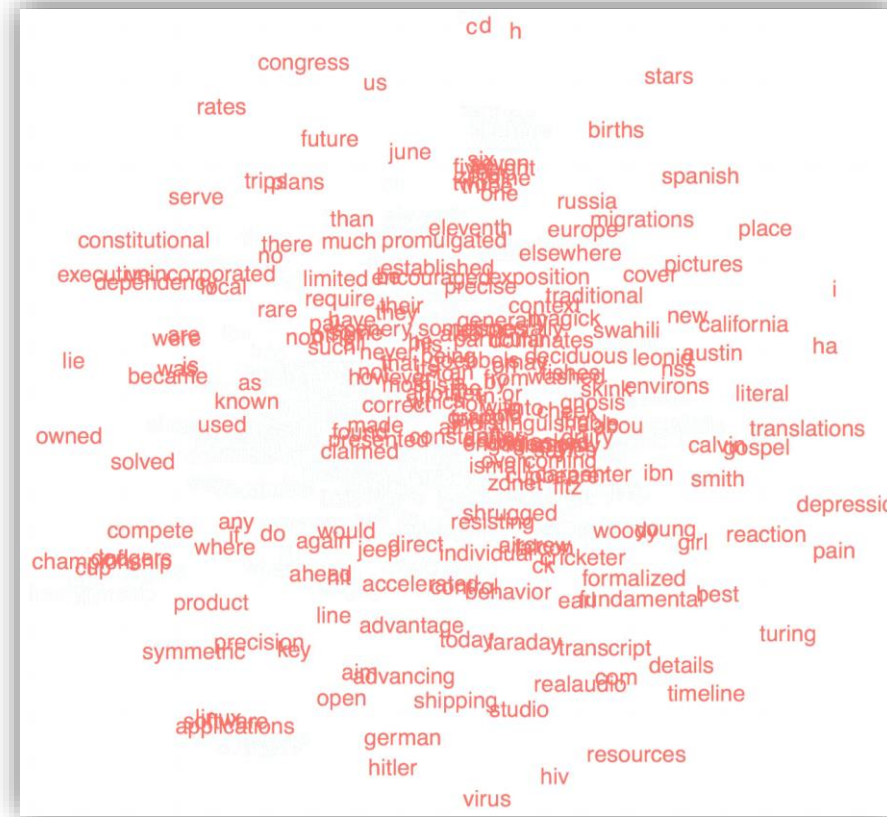
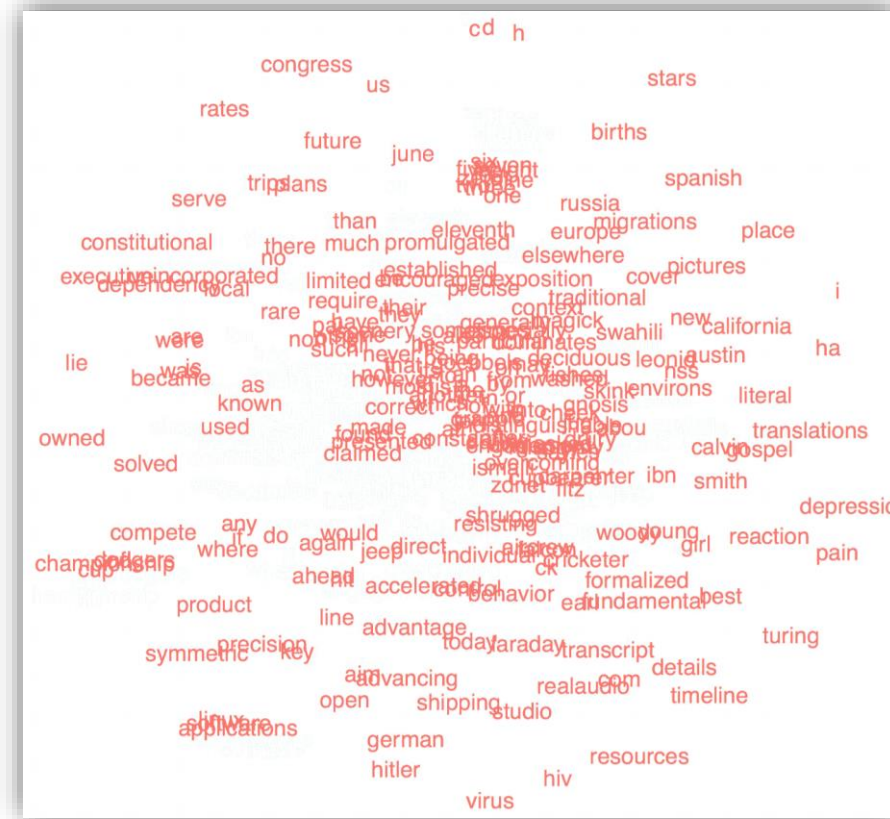


# The Strange Geometry of Skip-Gram with Negative Sampling: A story of geometric observations



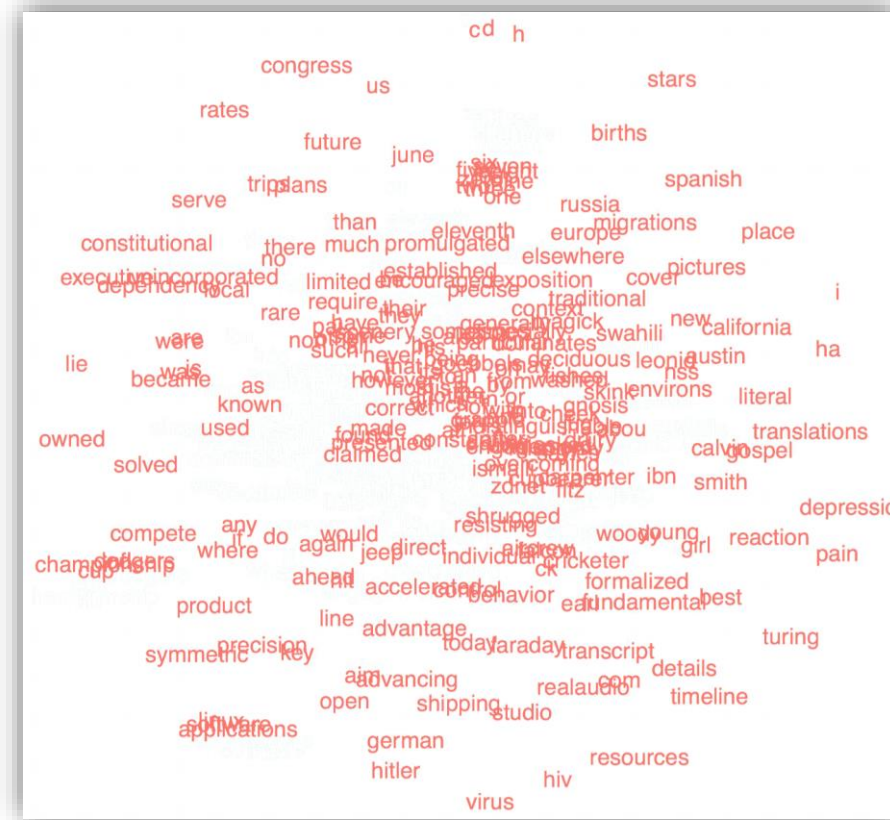
David Mimno and Laure Thompson  
Cornell University

# Ideally...



# Ideally...

Words span  
the entire  
K-dimensional  
space

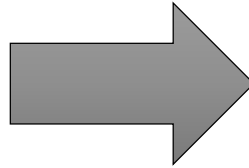




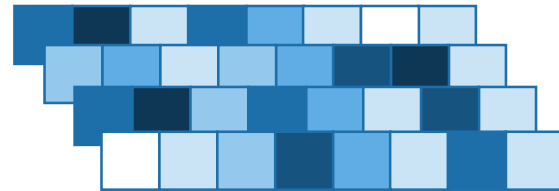
# Word Embeddings



**WIKIPEDIA**  
*The Free Encyclopedia*

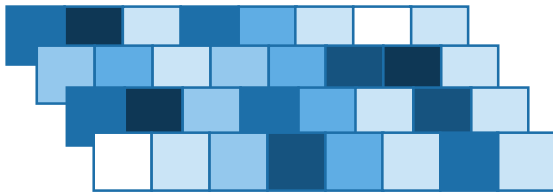


**Dense Vectors**

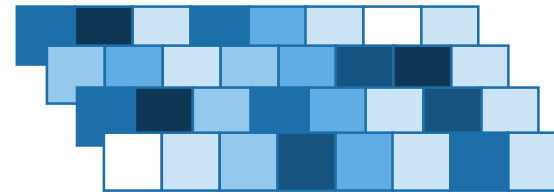


# Skip-Gram with Negative Sampling (SGNS)

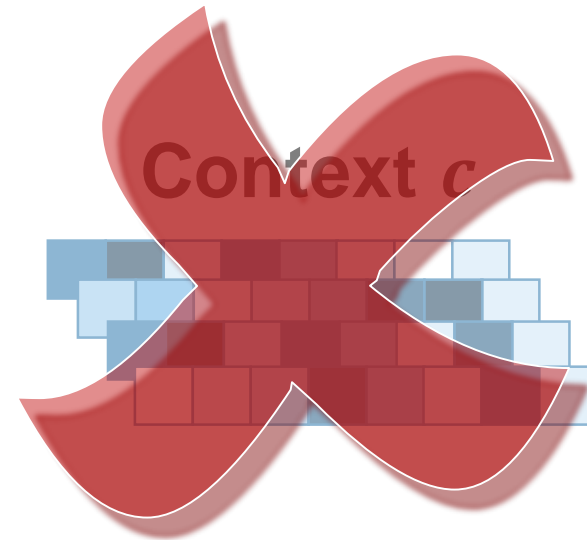
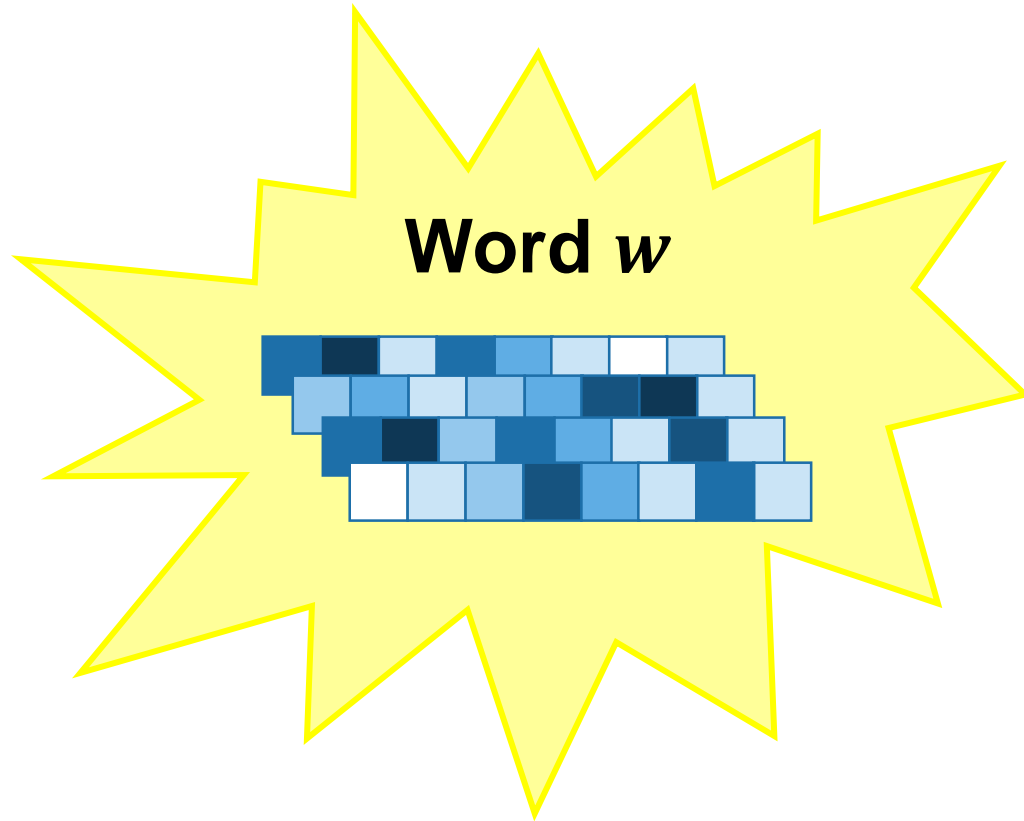
**Word  $w$**



**Context  $c$**



# Skip-Gram with Negative Sampling (SGNS)



# SGNS: Skip-Gram Model

The brown fox jumps over the lazy dog.





# SGNS: Skip-Gram Model

The brown fox **jumps** over the lazy dog.



# SGNS: Skip-Gram Model

The brown fox jumps over the lazy dog.

Context Window Size = 2

# SGNS: Skip-Gram Model

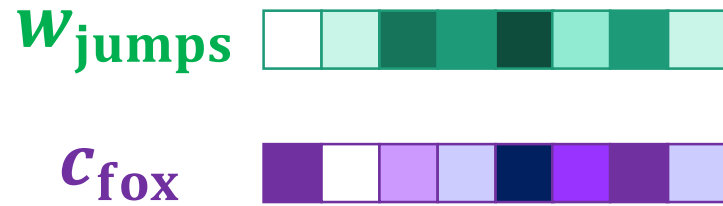
The brown fox jumps over the lazy dog.

Context Window Size = 2

jumps → { brown, fox, over, the }

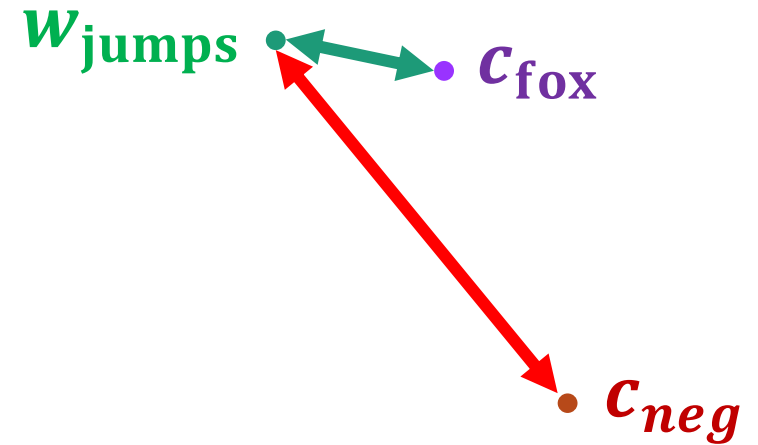
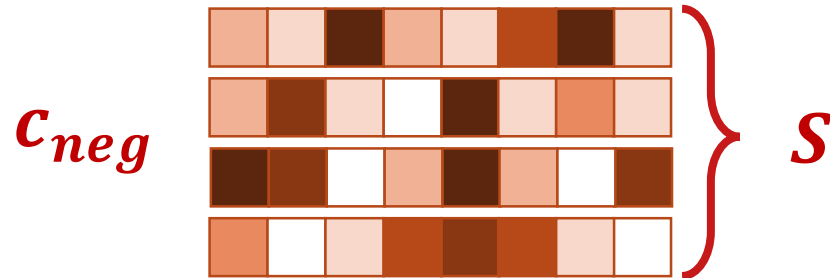
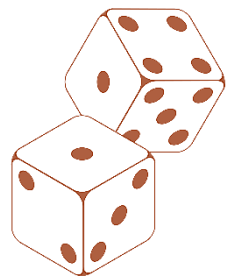
# SGNS: Negative Sampling

Co-occurrence **jumps**, **fox**:



# SGNS: Negative Sampling

Co-occurrence **jumps**, **fox**:



# Experimental Setup

## Corpus



**WIKIPEDIA**  
*The Free Encyclopedia*

## Embeddings

word2vec  
GloVe

## Parameters

Vector Size: 50  
Window Size: 5

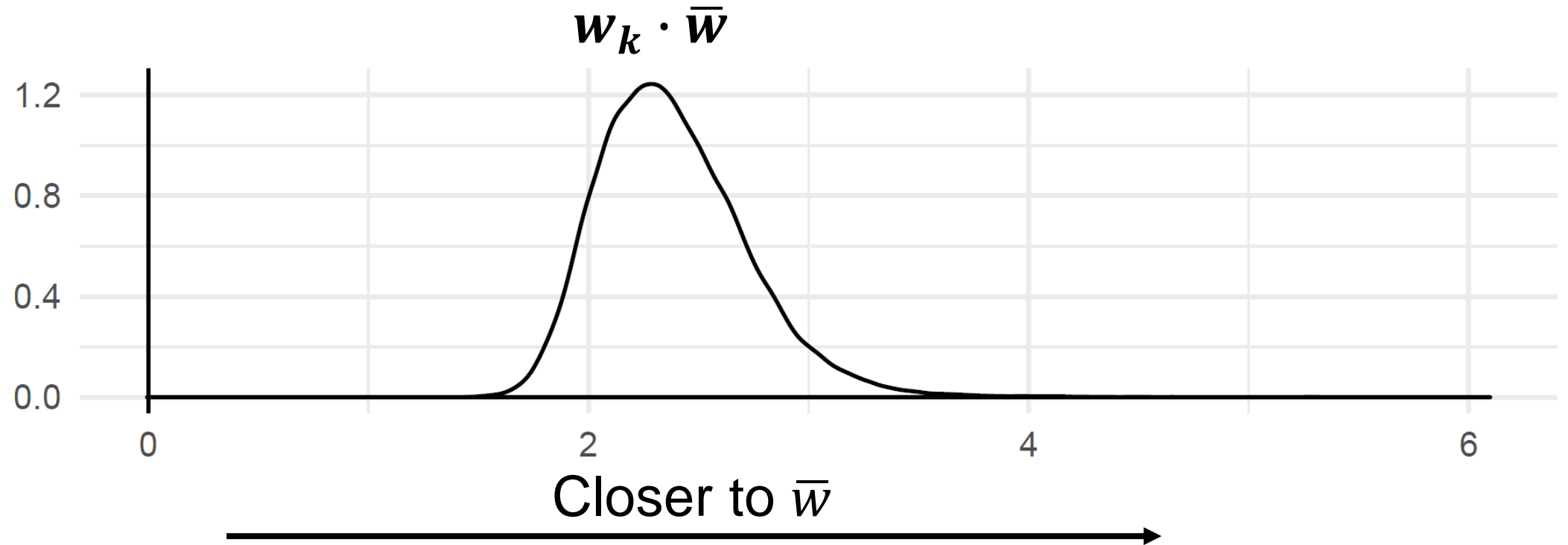
## Mean word vector

$$\bar{w} = \text{mean} \left( \begin{array}{cccccccc} \blacksquare & \square & \square & \square & \square & \square & \square & \square \\ \square & \square & \square & \blacksquare & \square & \square & \square & \square \\ \square & \square & \square & \square & \square & \square & \blacksquare & \square \end{array} \right)$$

Observation #1:

SGNS vectors arrange along a primary axis

# SGNS vectors point toward mean word vector



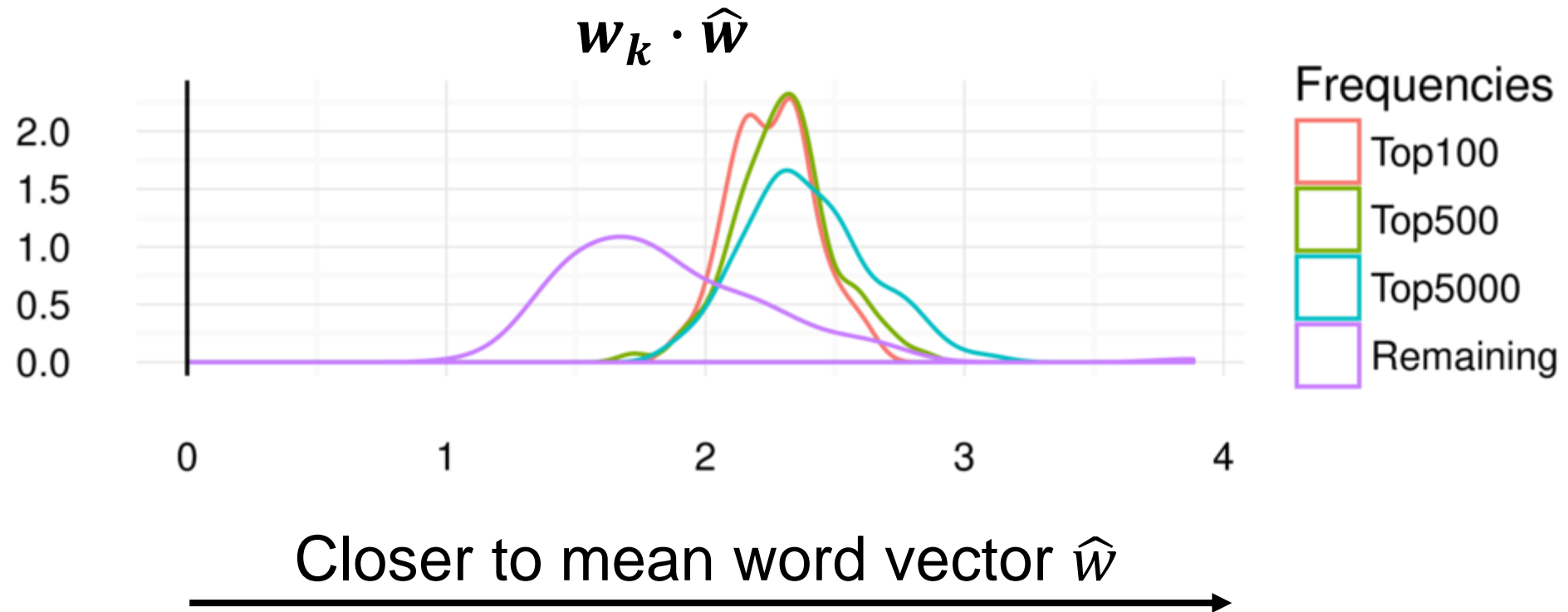


# Artifact of word frequency?

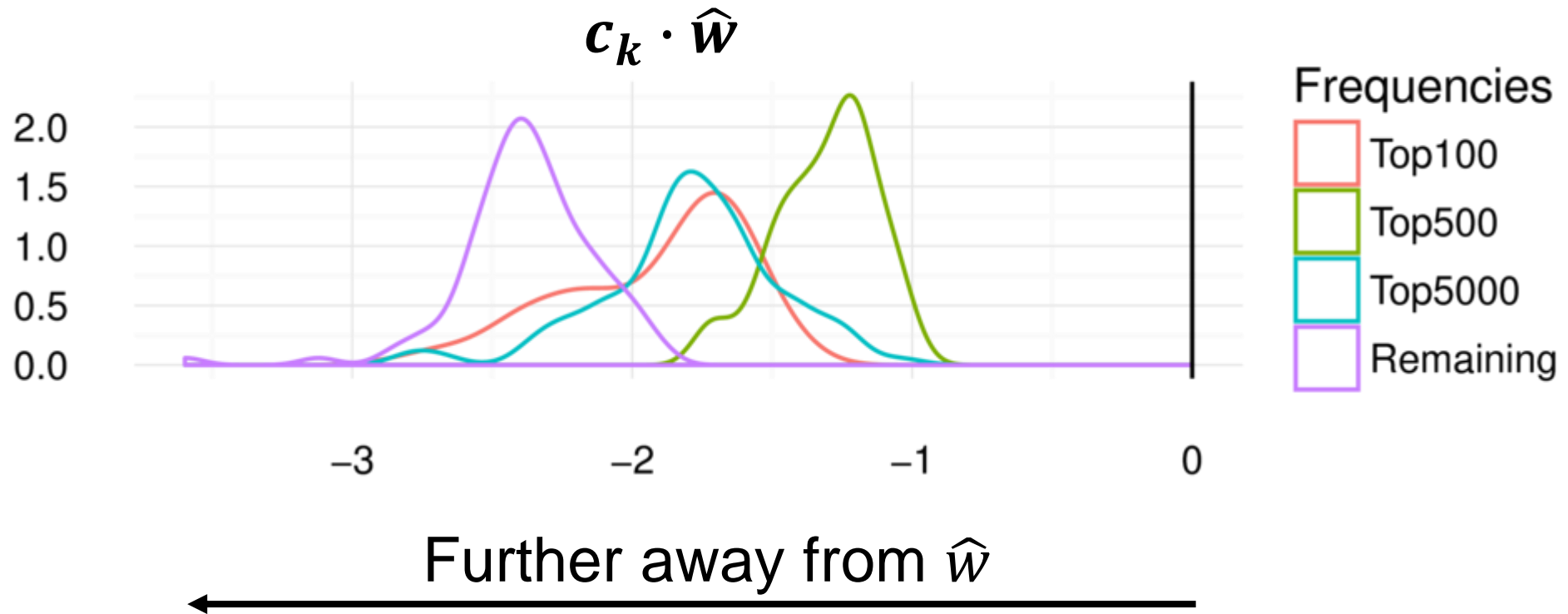
## 4 Frequency Levels:

- Ultra-high (1–100)
  - High (101–500)
  - Moderate (501–5000)
  - Low (5001+)
- 
- Sample 100 from each
  - Use sample mean vector  $\hat{w}$  instead of global mean  $\bar{w}$

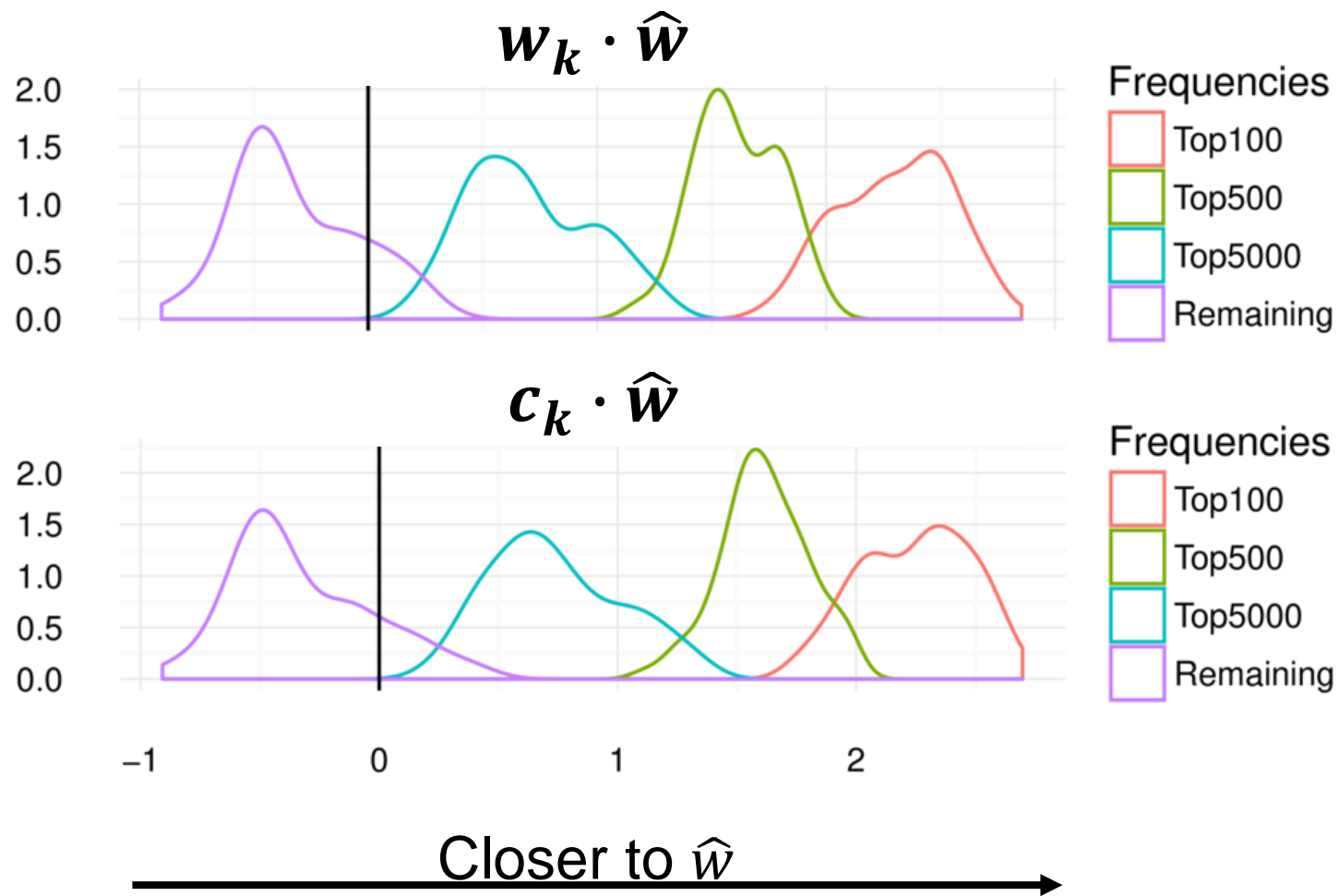
...true for all frequency classes



...and away from the context vectors



# Not true for GloVe

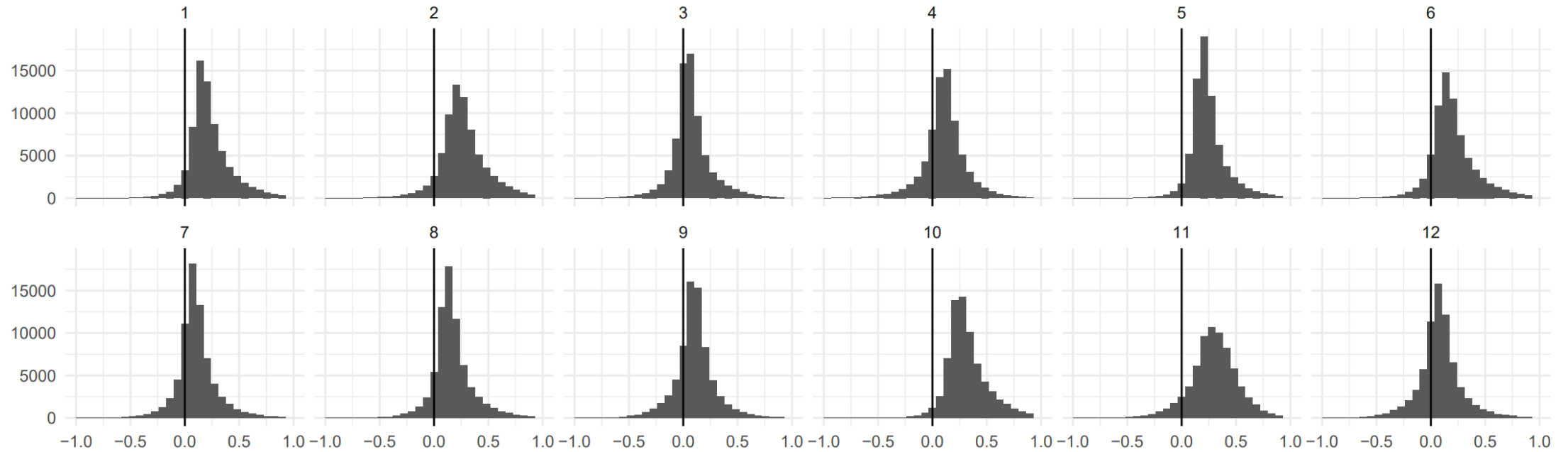




Observation #2:

SGNS vectors are mostly non-negative

# Latent dimensions skew “positive”



# ...inefficient use of K-dimensional space?

Preserve Semantic Properties:

1. Dropping “negative” entries

$$w'_k = \max(0, w_k * \text{sign}(\bar{w}_k))$$

2. Subtracting mean vector

$$w' = w - \bar{w}$$

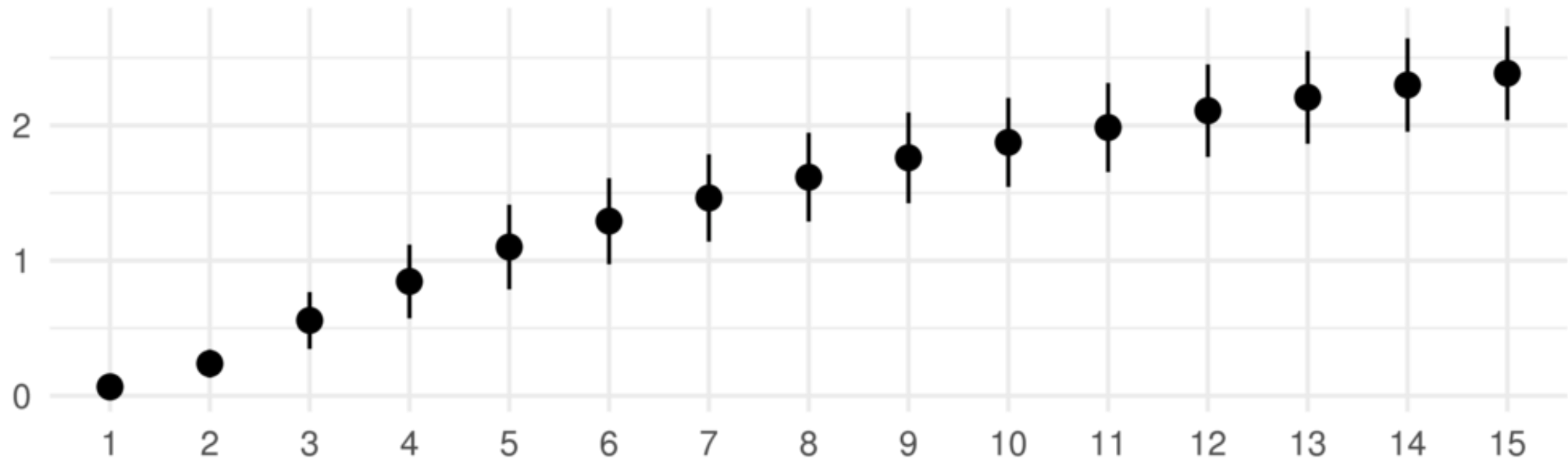


Observation #3:

Negative sampling affects SGNS geometry

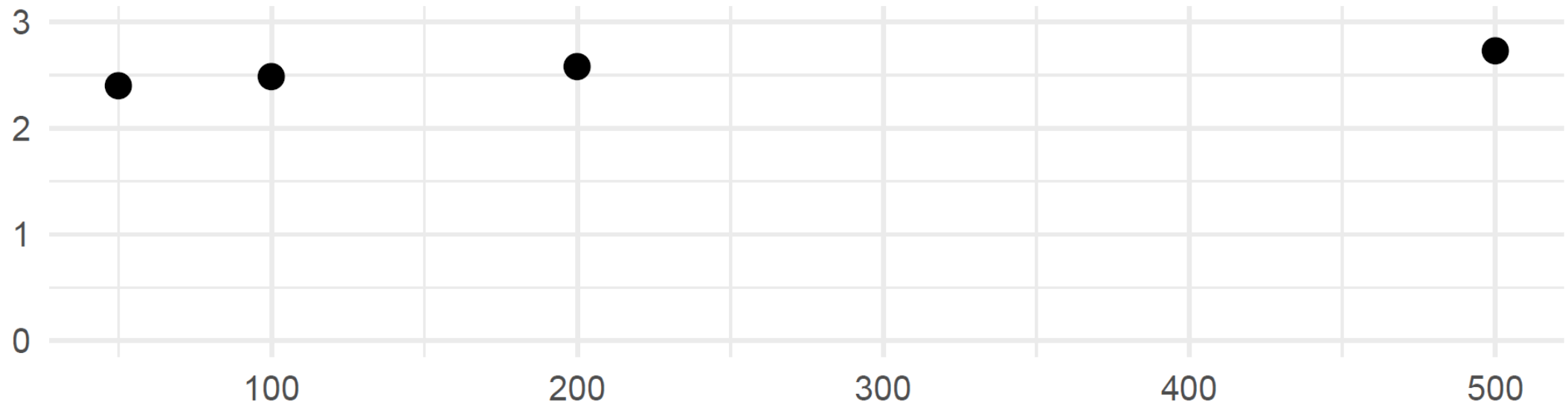
# More negative samples, better alignment

$Avg(w_k \cdot \bar{w})$  vs. # of Negative Samples



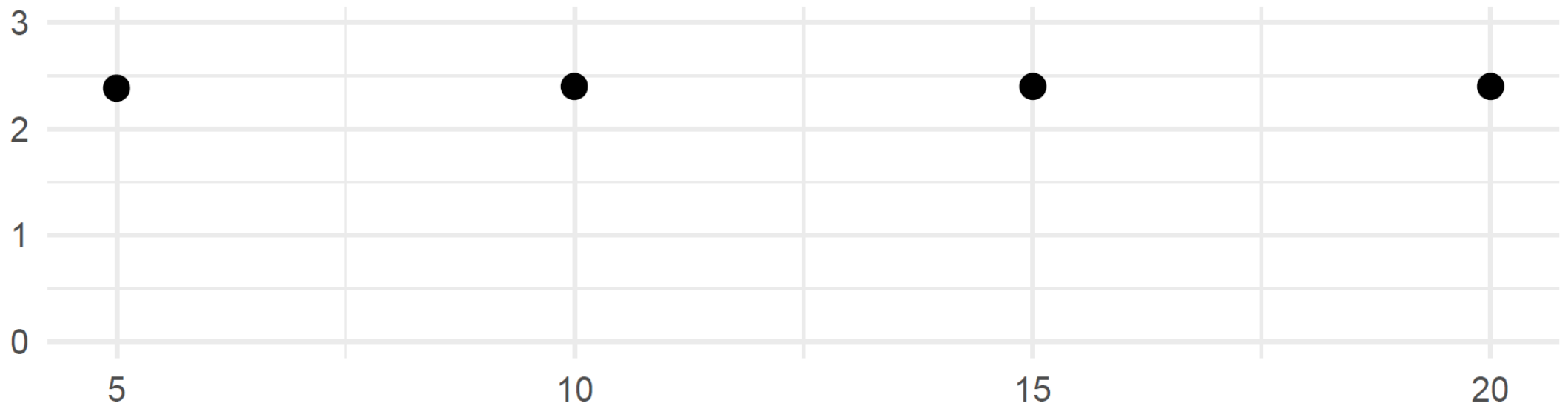
...this is not seen in other parameters

$Avg(w_k \cdot \bar{w})$  vs. Vector Size



...this is not seen in other parameters

$Avg(w_k \cdot \bar{w})$  vs. Window Size



# Thank You!

