# FEATURE SPACES III

Prof. Alexander Huth
10/19/2017

# HOMEWORKS

* Homework 1 will be graded/returned by next Thursday (10/26)

* Homework 2 out next Thursday (10/26)
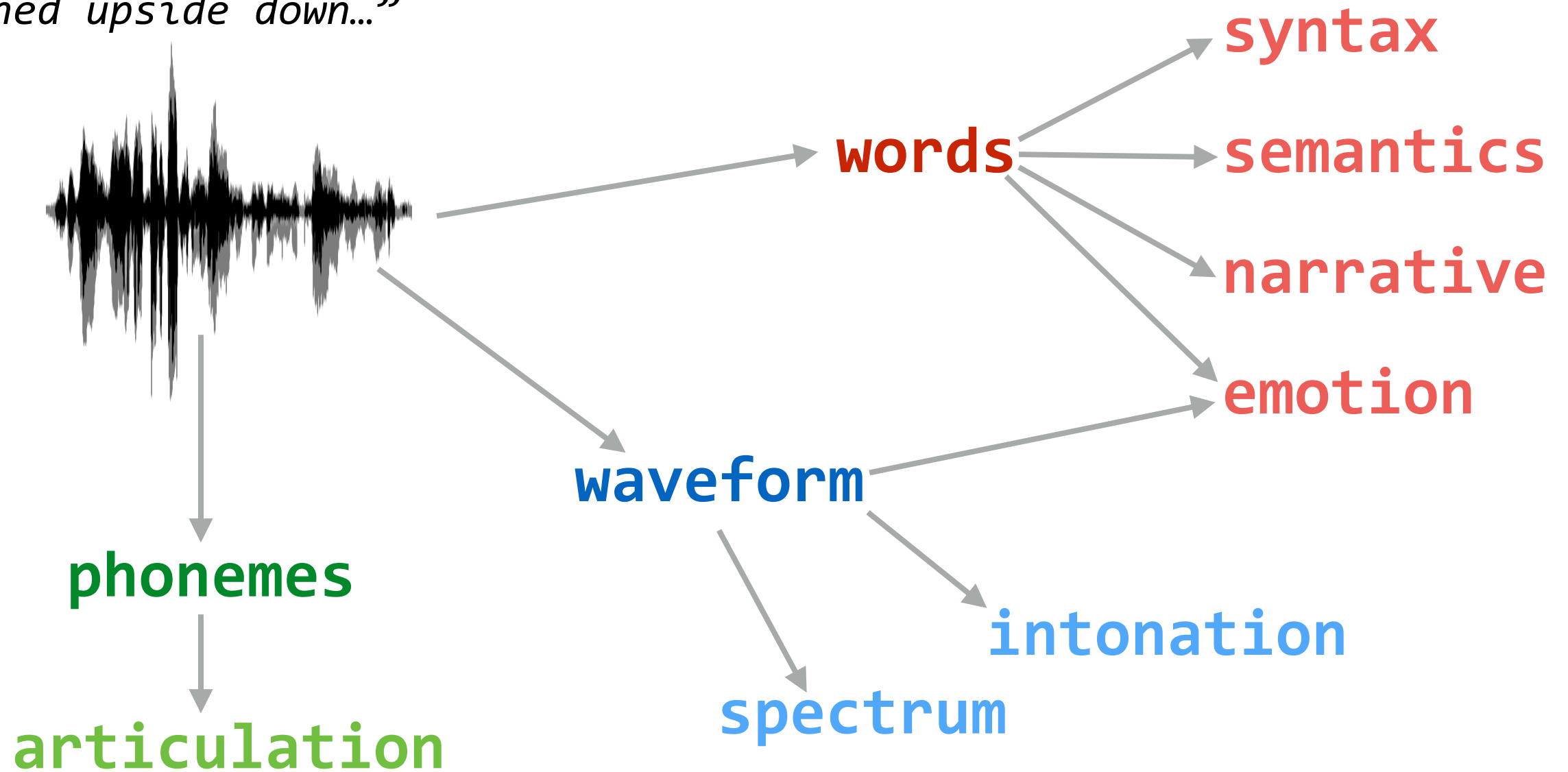
# SYSTEM IDENTIFICATION

* **Linearized model**
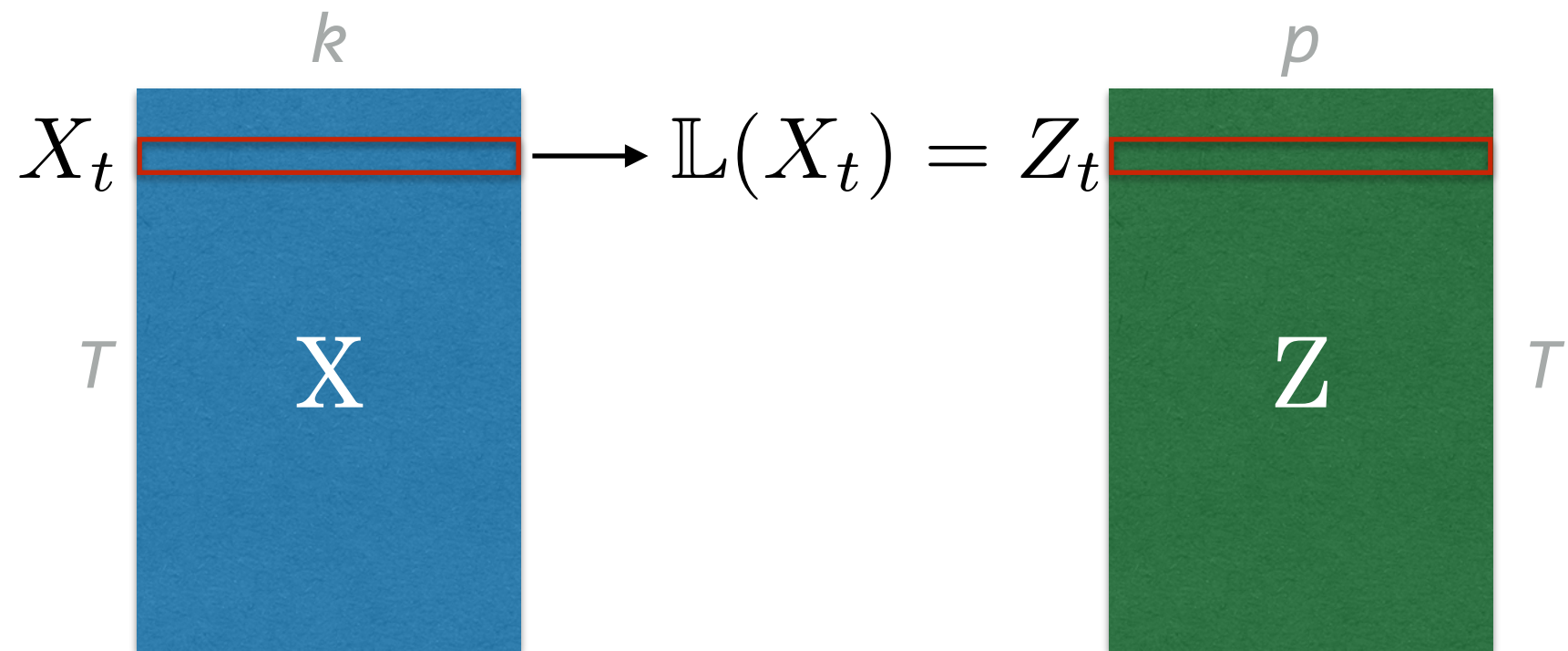
$$Y = \mathbb{L}(X)\beta$$

Let's invent some **L**'s

# LINEARIZING TRANSFORMATIONS

\* Simplest version: time-invariant **L**

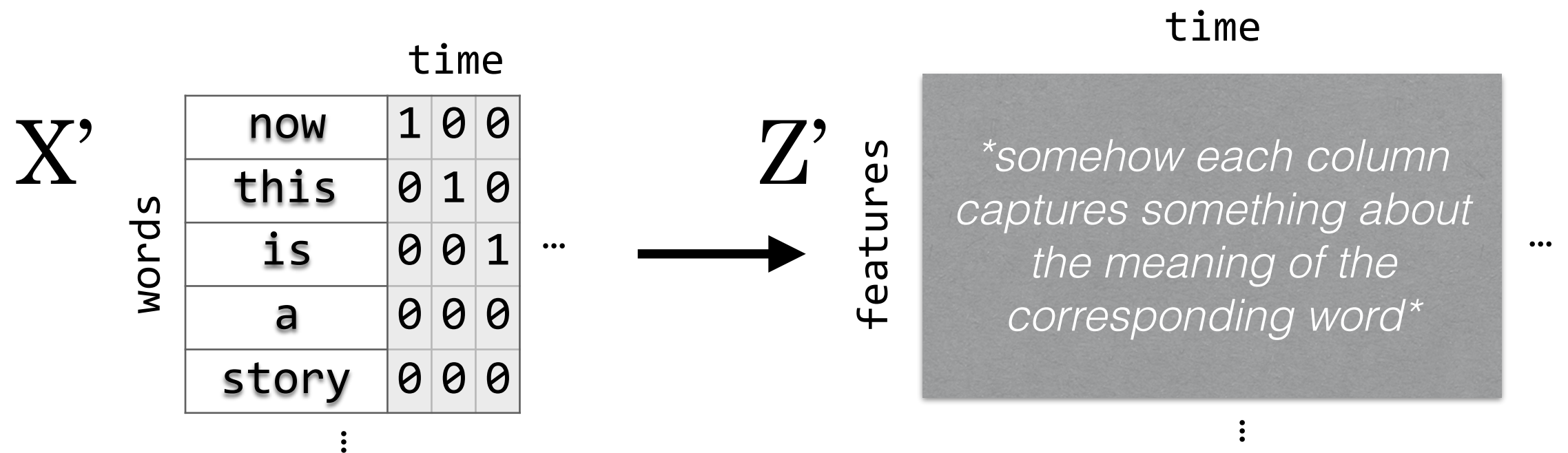

$$X_t \quad \longrightarrow \quad \mathbb{L}(X_t) = Z_t$$

# LEXICAL SEMANTICS

* Let's create an **L** that captures word-level semantic information

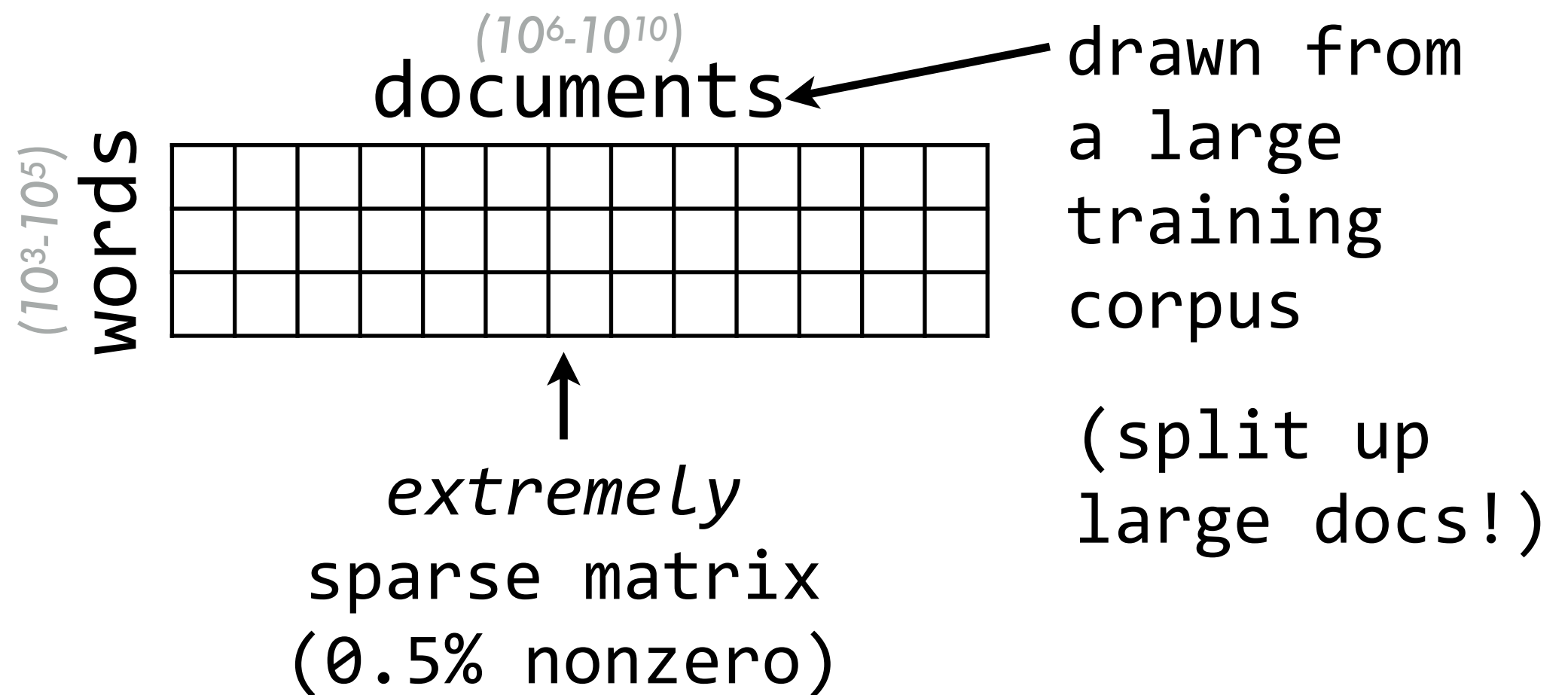* Unlike the ₐwful syntax models, this model will be *time-invariant*

# LEXICAL SEMANTICS

*1-hot vector*

*p-dim vector*

$$\mathbb{L}_{semantic}(X_t) = Z_t \in \mathbb{R}^p$$

X'

time

| words | | time | | |
|---|---|---|---|---|
| now | 1 | 0 | 0 | |
| this | 0 | 1 | 0 | |
| is | 0 | 0 | 1 | ... |
| a | 0 | 0 | 0 | |
| story | 0 | 0 | 0 | |

⋮

Z'

time

features

*somehow each column captures something about the meaning of the corresponding word*

...

⋮

# LEXICAL SEM. - LSA

* Latent Semantic Analysis (LSA)

$(10^6-10^{10})$

documents ← drawn from a large training corpus

$(10^3-10^5)$ words

↑
*extremely*
sparse matrix
(0.5% nonzero)

(split up large docs!)

Deerwester et al. (1988)

# LEXICAL SEM. - LSA

* Latent Semantic Analysis (LSA)

*($10^6$-$10^{10}$)*

documents

($10^3$-$10^5$) words

$a_{ij}$ term i, doc j

Often the entries are normalized
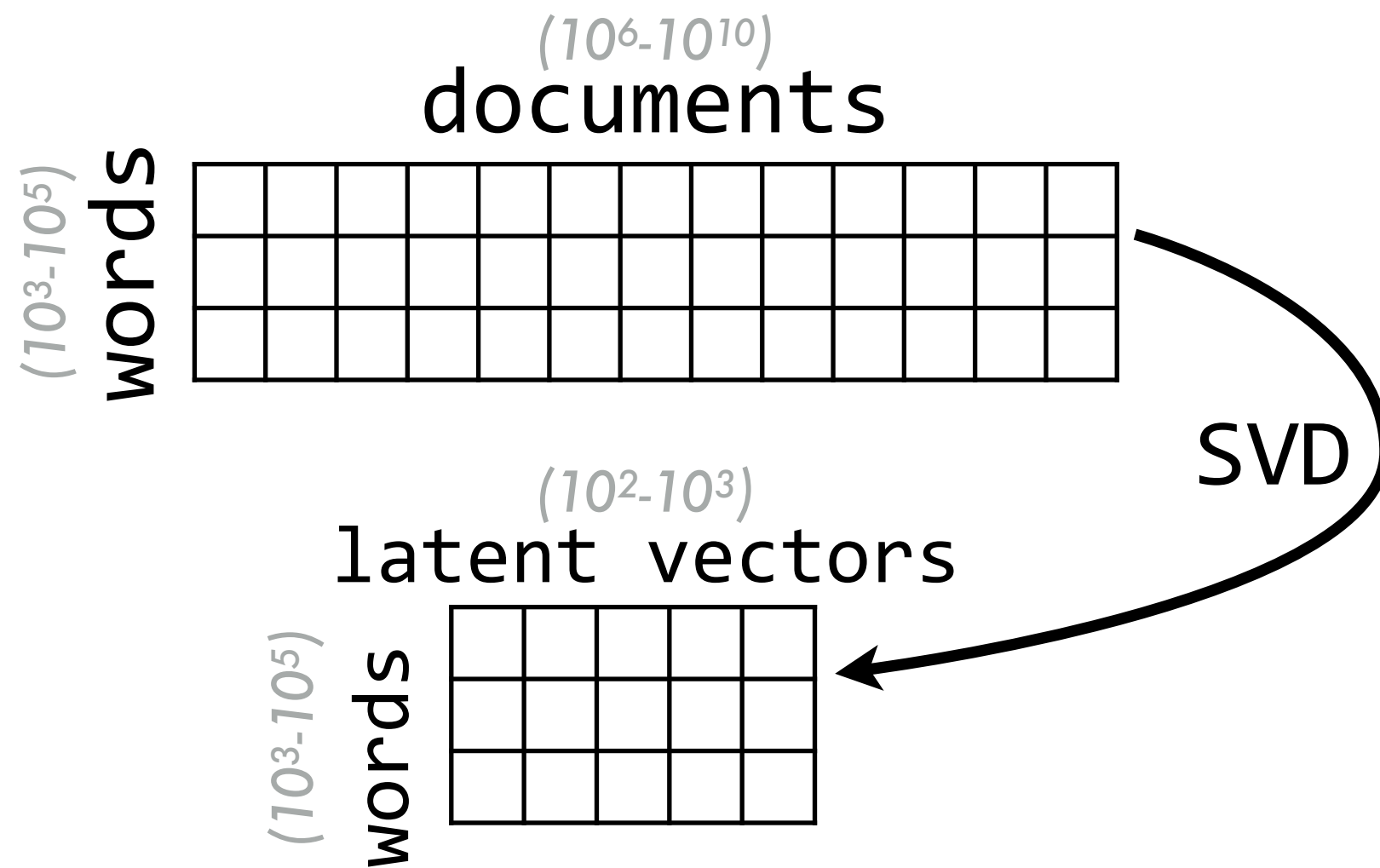
$$a_{ij} = \frac{\text{tf}_{ij}}{\log_2 \frac{n}{1+\text{df}_i}}$$
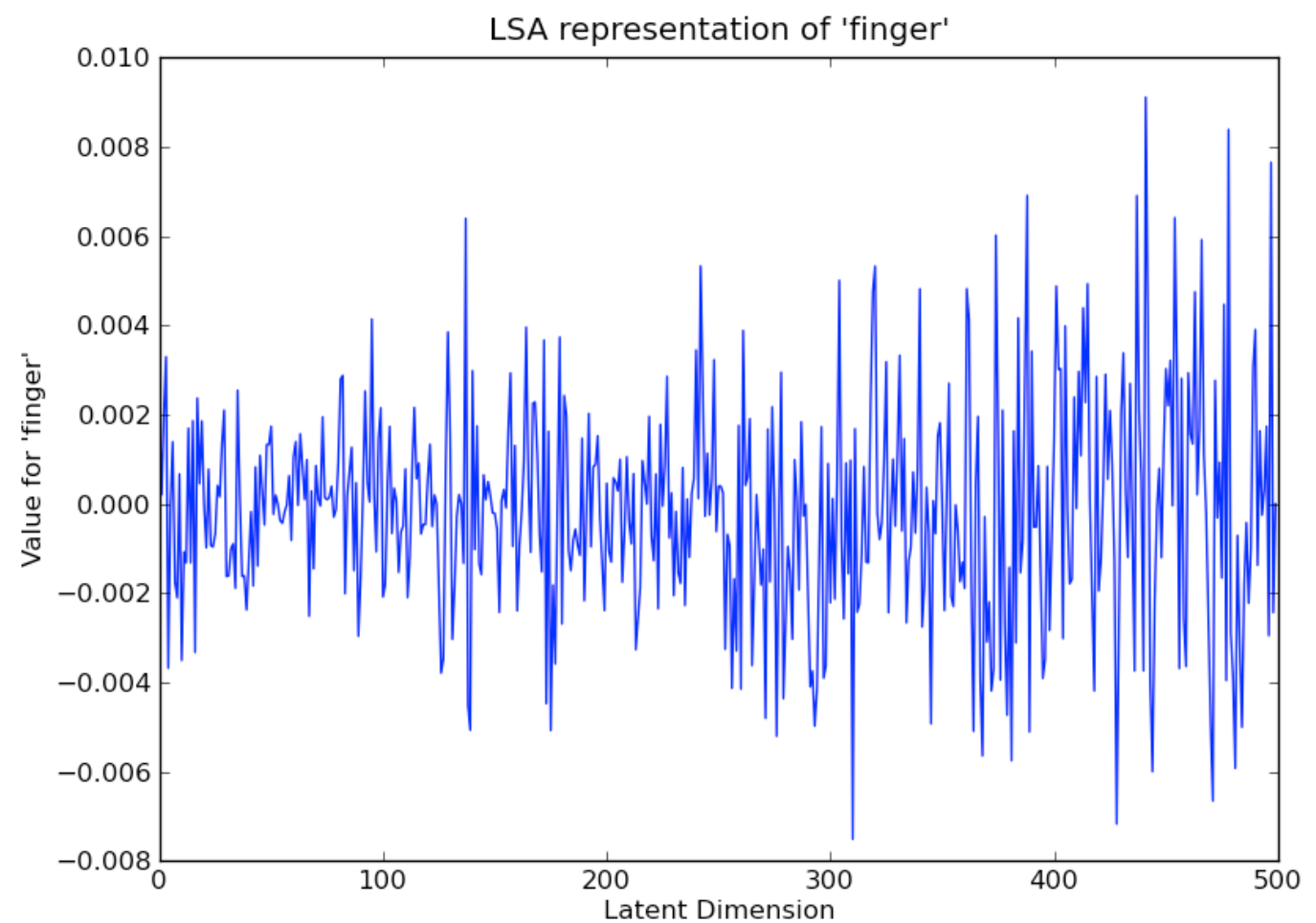
— # word $i$ in doc $j$
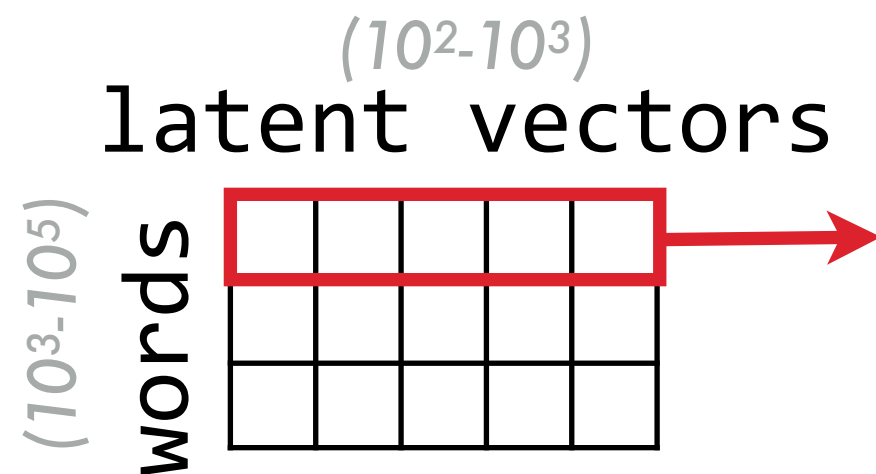
— # docs

— # docs with word $i$

"tf-idf"

{ij}}{\log_2
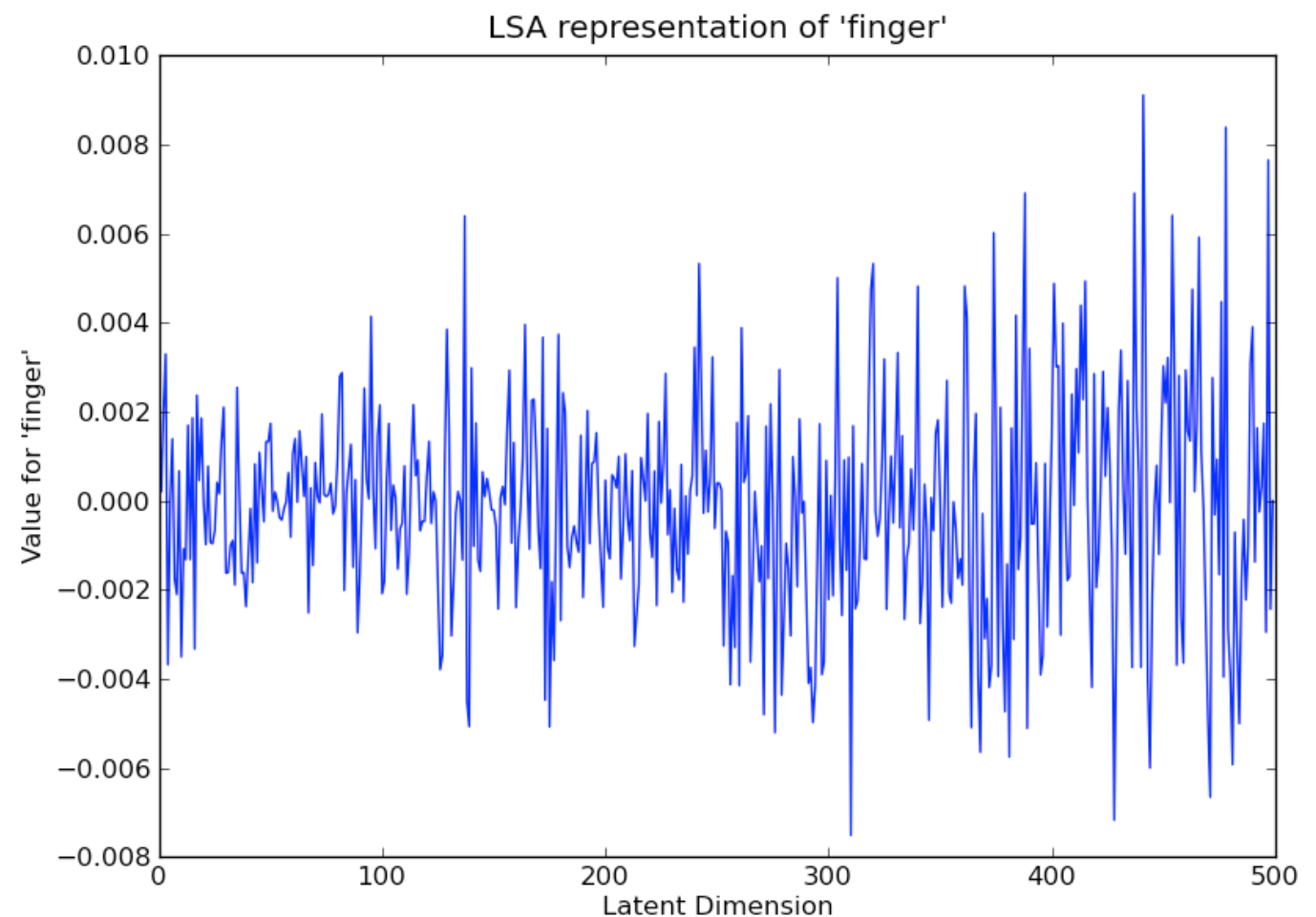}

# LEXICAL SEM. - LSA
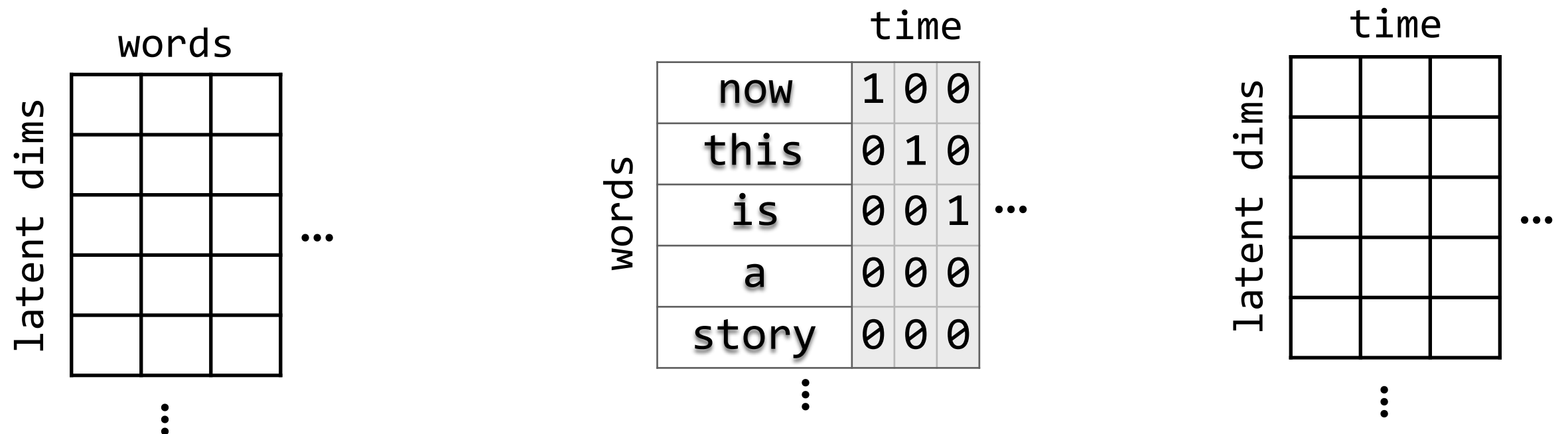
* Latent Semantic Analysis (LSA)

# LEXICAL SEM. - LSA

# LEXICAL SEM. - LSA

(10²-10³)

latent vectors

(10³-10⁵) words



0.61- *thumb*
0.52- *fingers*
0.51- *forefinger*
0.49- *tip*
0.47- *wrist*
0.41- *elbow*
0.41- *throat*
0.41- *lips*
0.40- *string*

LSA representation of 'finger'

# LEXICAL SEM. - LSA

$$E \cdot X' = Z'$$

embedding matrix' * word matrix' = semantic stimulus matrix'



words

latent dims

...

...

time

|  | 1 | 0 | 0 |
|---|---|---|---|
| now | 1 | 0 | 0 |
| this | 0 | 1 | 0 |
| is | 0 | 0 | 1 |
| a | 0 | 0 | 0 |
| story | 0 | 0 | 0 |

words

...

time

latent dims

...

...

REMINDER FROM A
FEW WEEKS AGO...

# TIKHONOV REGRESSION

* this is equivalent to **TIKHONOV REGRESSION** on the **WORDS** with a prior determined by the **WORD EMBEDDING**

$$\frac{1}{\sigma^2}\Sigma_\beta = (C^T C)^{-1} = E^T E$$

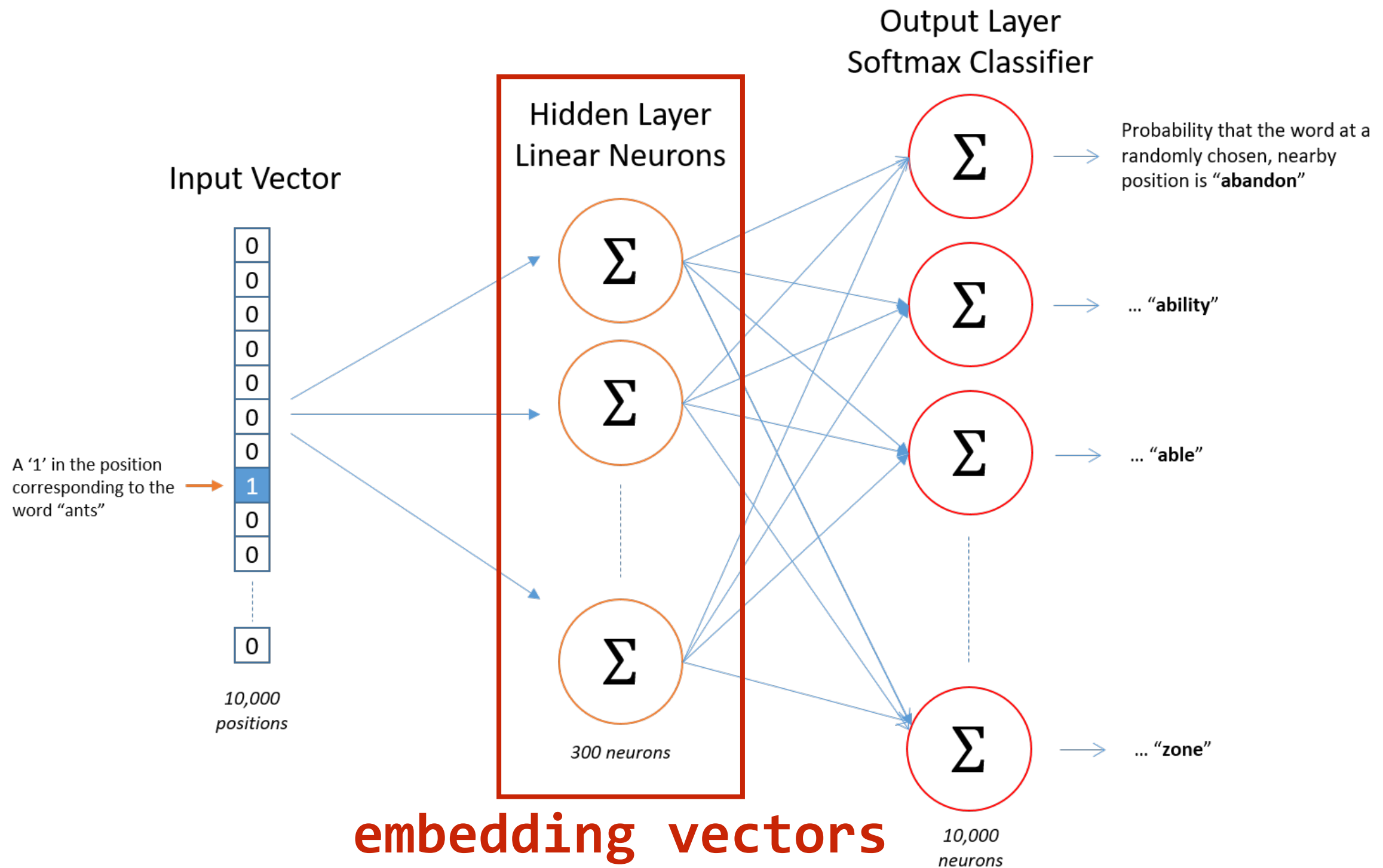**PRIOR COVARIANCE**     **INVERSE OF PENALTY INNER PRODUCT**     **EMBEDDING INNER PRODUCT**

* **i.e. the prior covariance between two words' weights is equal to the dot product of their embedding vectors**
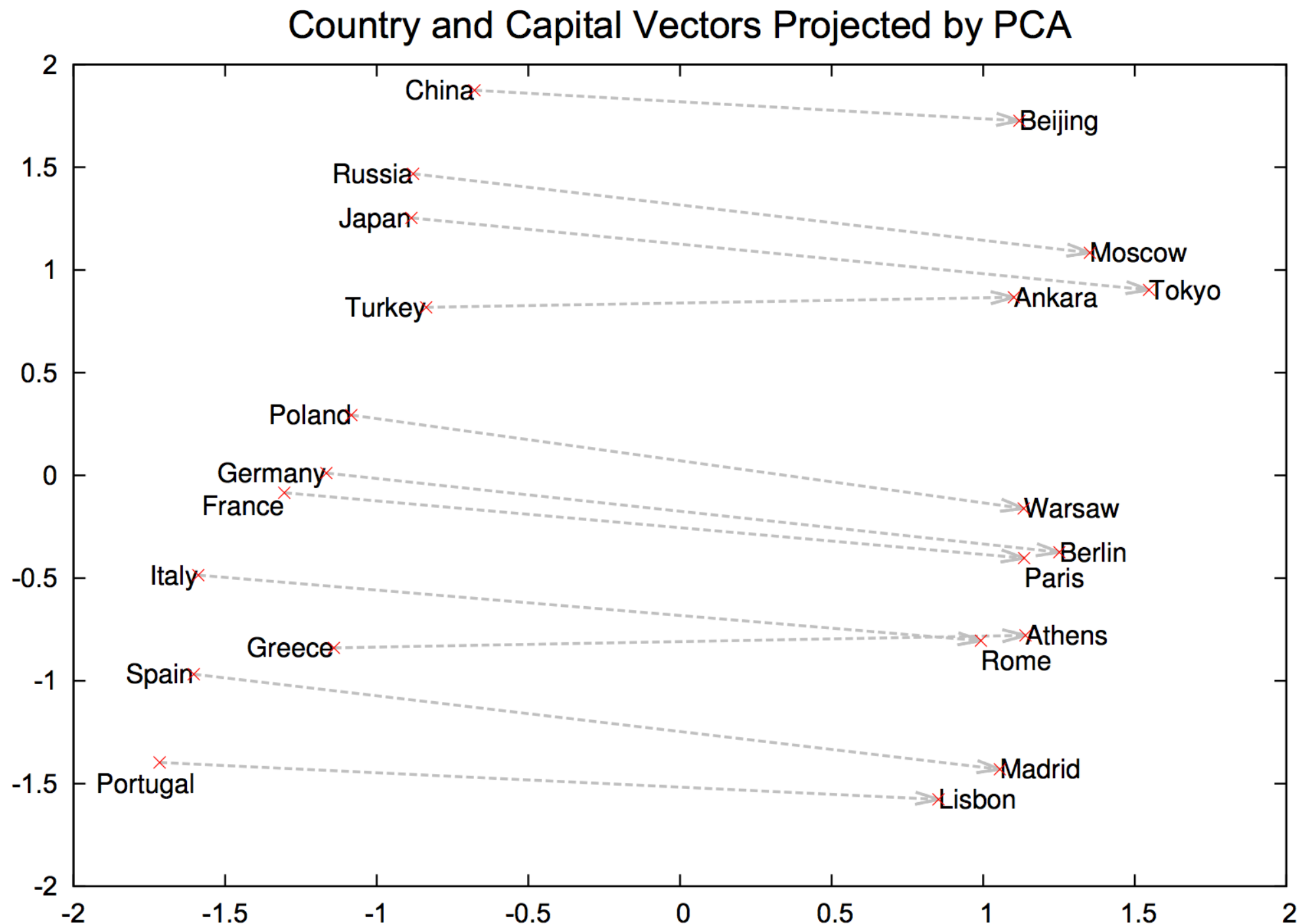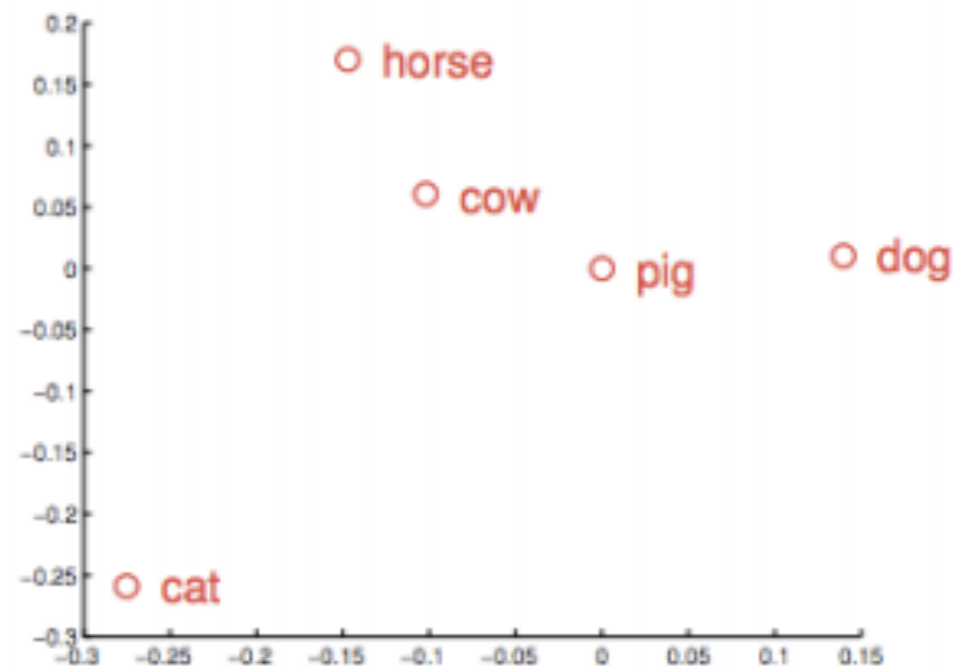
# SEMANTICS - WORD2VEC



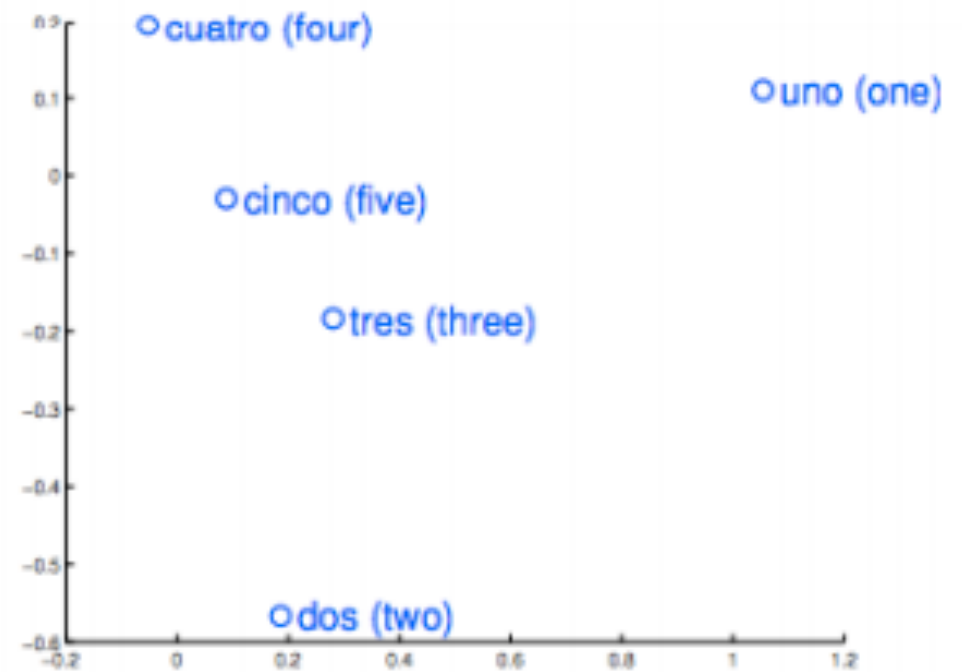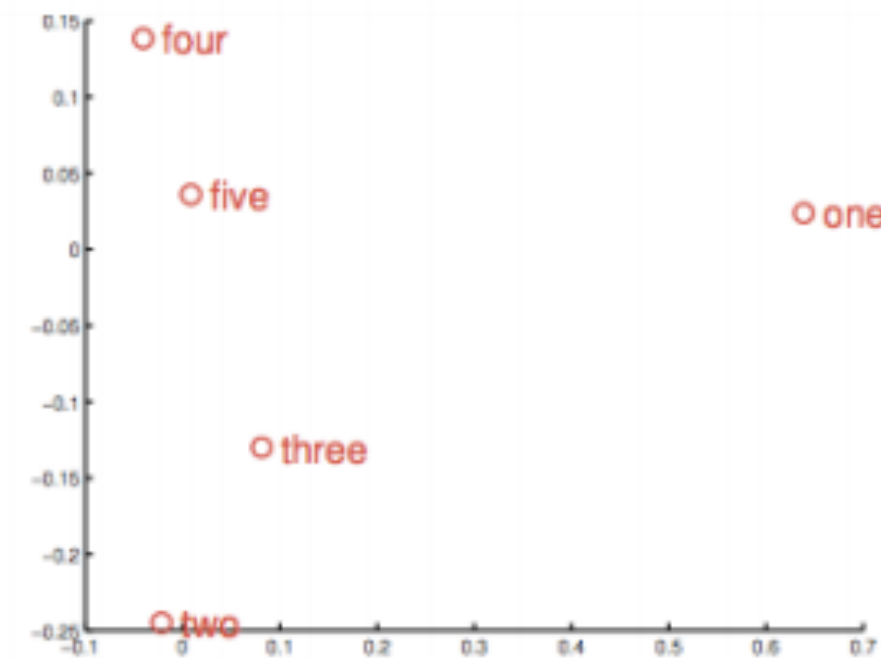Input    projection    output

w(t) → w(t-2), w(t-1), w(t+1), w(t+2)

Mikolov et al. (2013)

# SEMANTICS - WORD2VEC

# SEMANTICS - WORD2VEC



Country and Capital Vectors Projected by PCA

Mikolov et al. (2013)

# SEMANTICS - WORD2VEC

# SEMANTICS - ENGLISH-1000

…
difficult
…
husband
…
potato
…
remember
…

CO-OCCURRENCE
MATRIX

TARGET WORDS

985 w.

FULL LEXICON

10,470 words
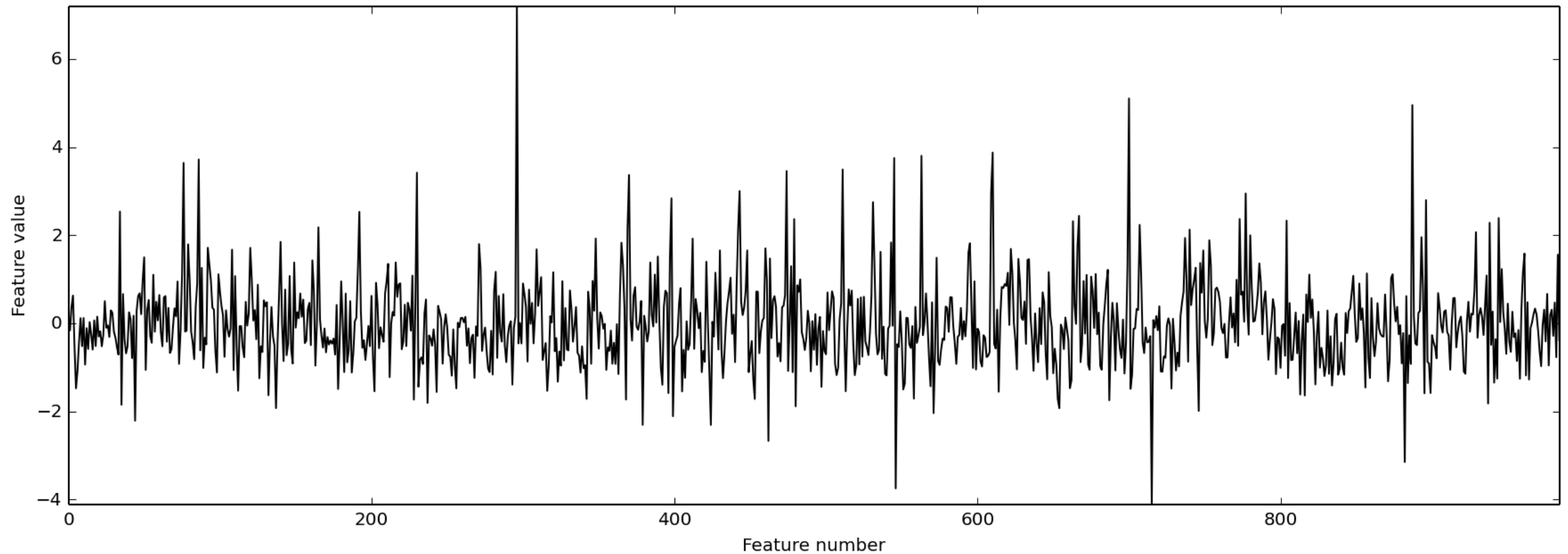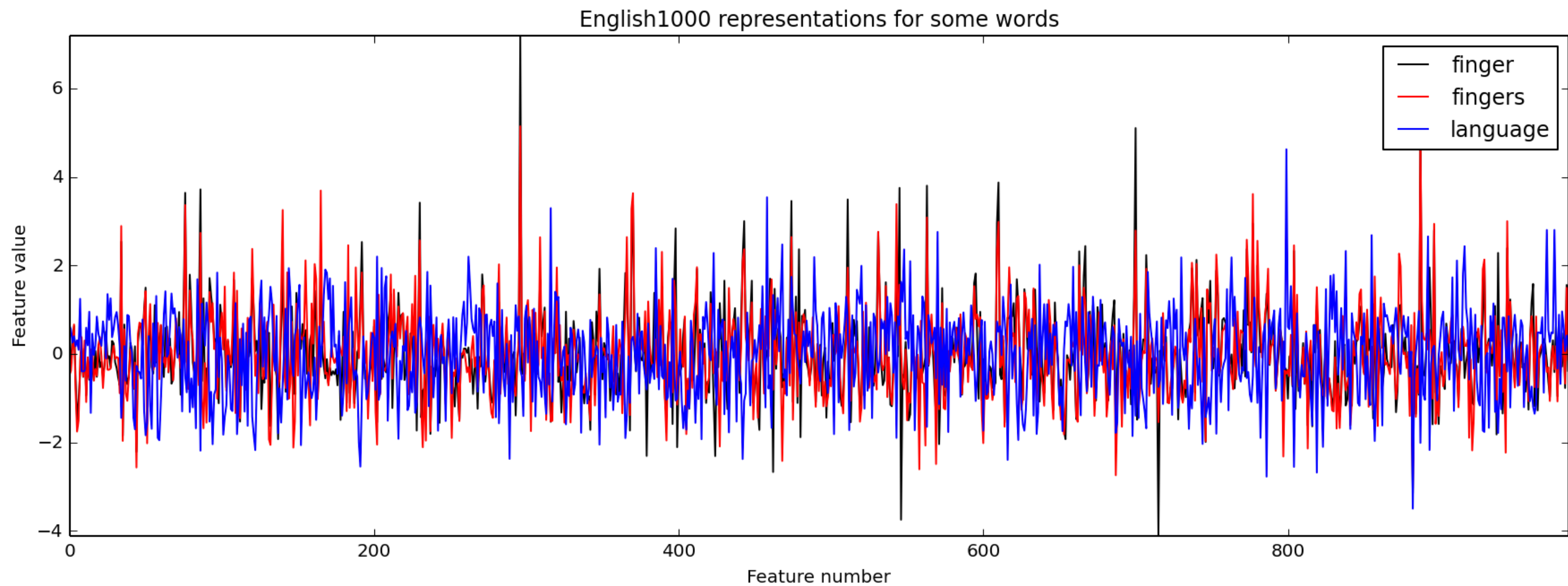
# SEMANTICS - ENGLISH-1000

* The corpus was used to build a 985 x 10,470 matrix **M**

  * $M_{i,j}$ is the number of times target $i$ occurs within 15 words of word $j$

* Then log-transform:  $M^*_{i,j} = Log(M_{i,j}+1)$

* Then $z$-score each row, then each column

* ... yielding 985-D vector representation of each word in the lexicon

# SEMANTICS - ENGLISH-1000
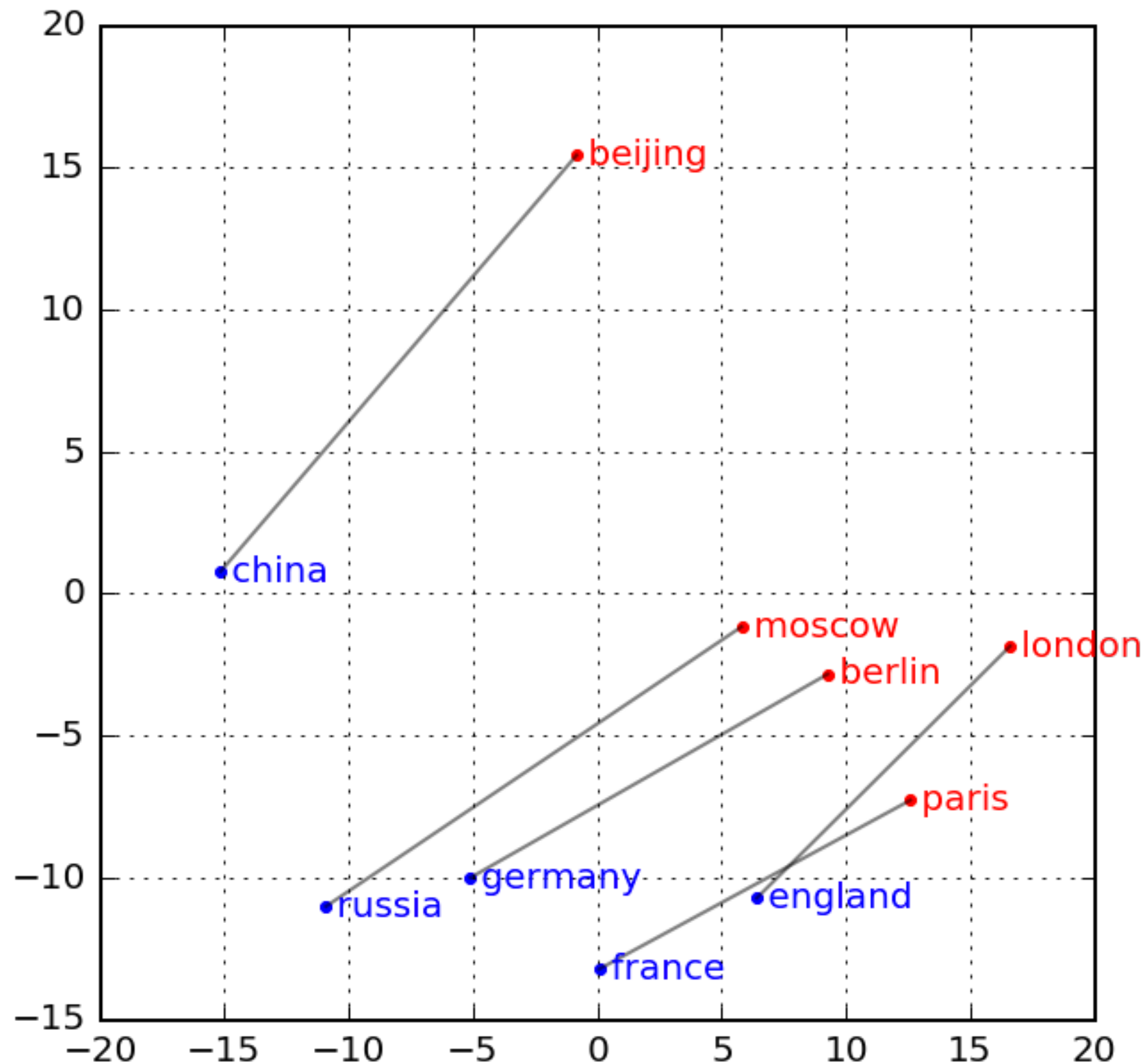


"finger" in english-1000

# SEMANTICS - ENGLISH-1000



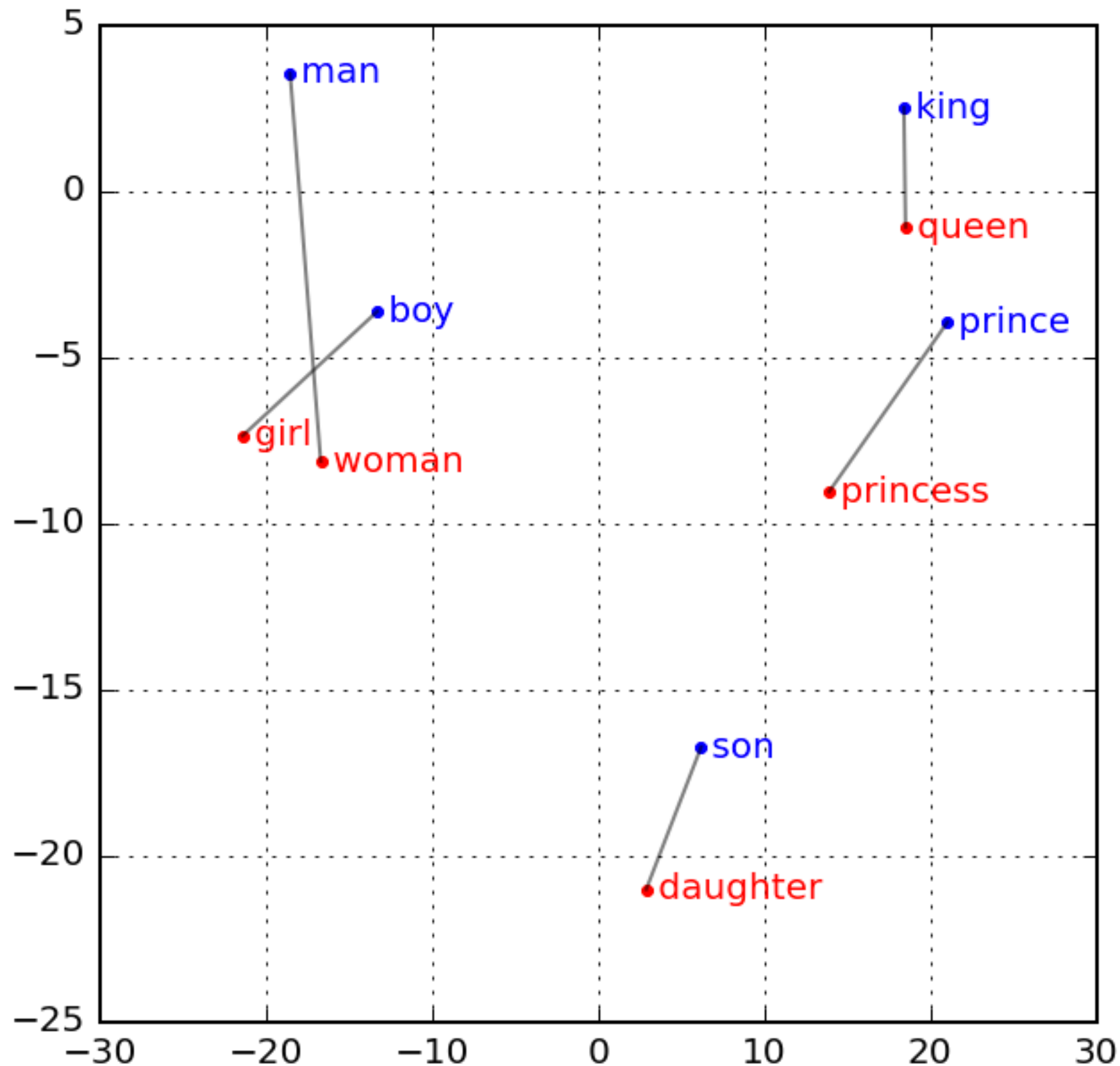English1000 representations for some words

# SEMANTICS - ENGLISH-1000

correlation
with "finger"

word

```
1.00, 'finger'
0.81, 'fingers'
0.67, 'hand'
0.67, 'nose'
0.66, 'arm'
0.64, 'mouth'
0.64, 'stick'
0.63, 'neck'
0.63, 'forehead'
0.62, 'tongue'
```
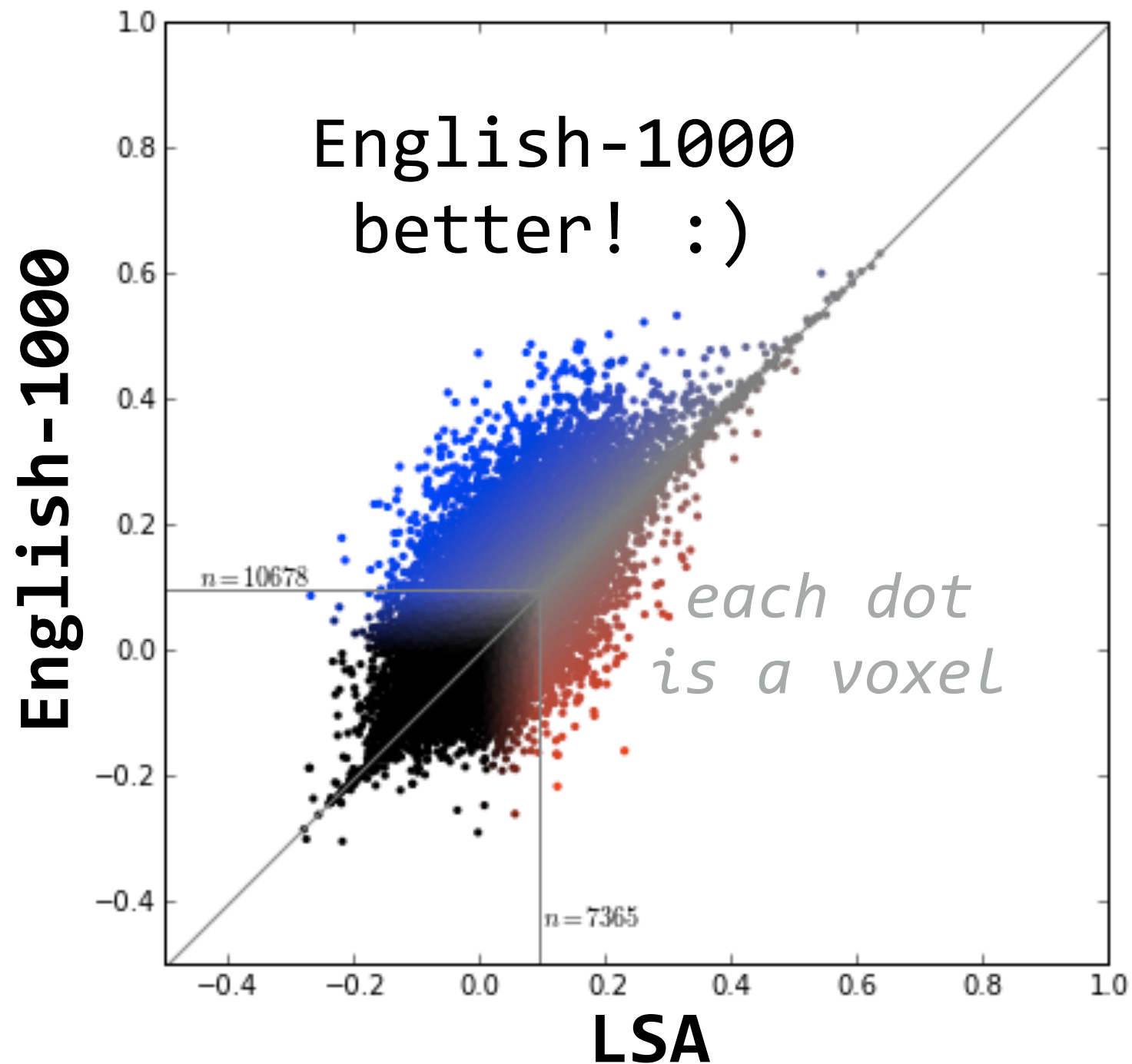
# SEMANTICS - ENGLISH-1000
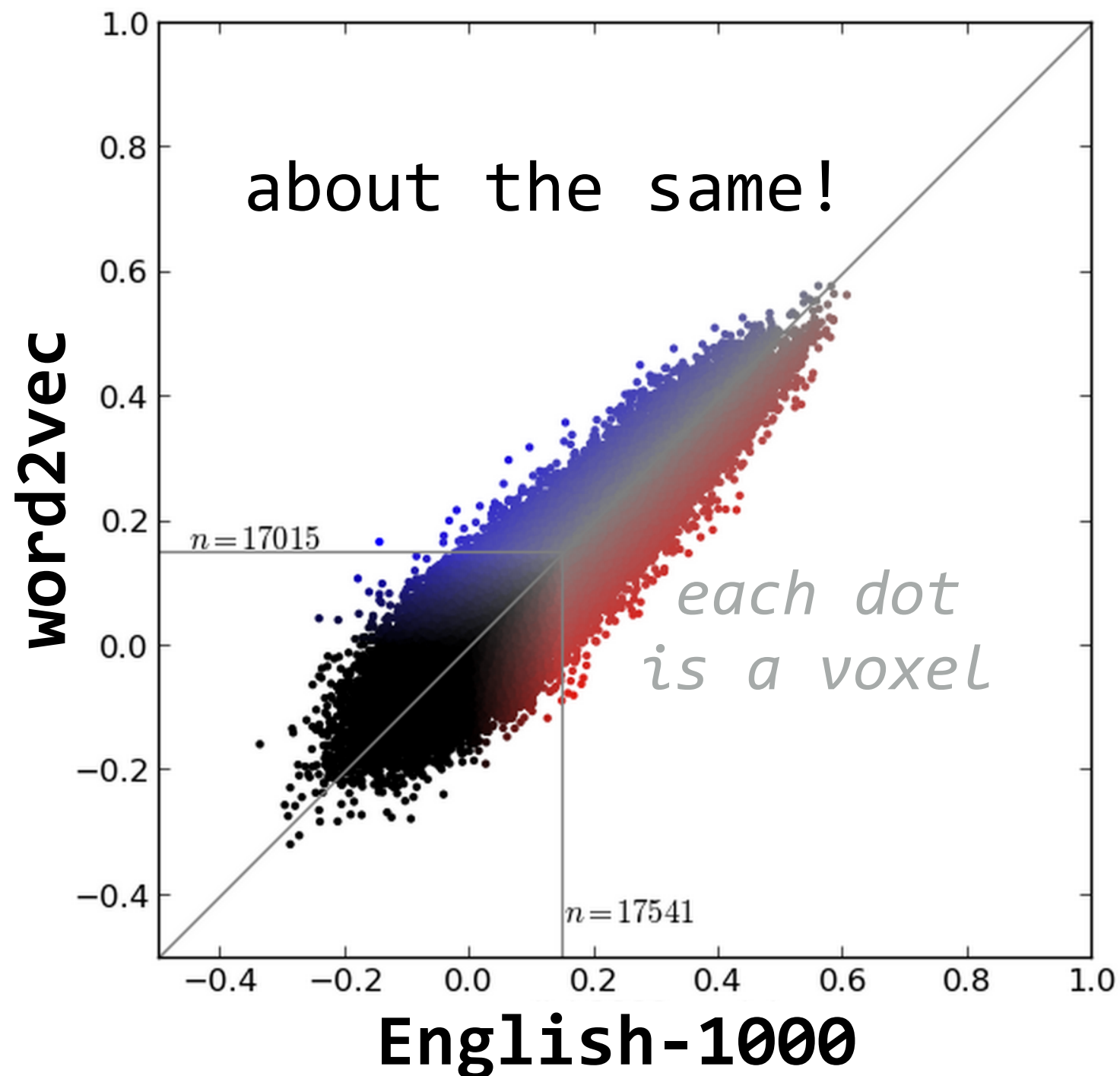
# SEMANTICS - ENGLISH-1000

# ENGLISH-1000 VS LSA

## model performance on held-out data

# ENGLISH-1000 VS WORD2VEC

model performance on held-out data

# NEXT TIME

* Model comparison

* Variance partitioning