

FEATURE SPACES II

Prof. Alexander Huth

10/17/2017

REMINDER

- * Homework 1 due TODAY! Please email it to huth@cs.utexas.edu by midnight

SYSTEM IDENTIFICATION

$$Y = f(X)$$

* What kind of a function is f ?

SYSTEM IDENTIFICATION

* Linearized model

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



Let's invent some L's

LINEARIZING TRANSFORMATIONS

- * `L` should be a function that:
 - * `ingests stimuli`
 - * `emits feature vectors`

LINEARIZING TRANSFORMATIONS

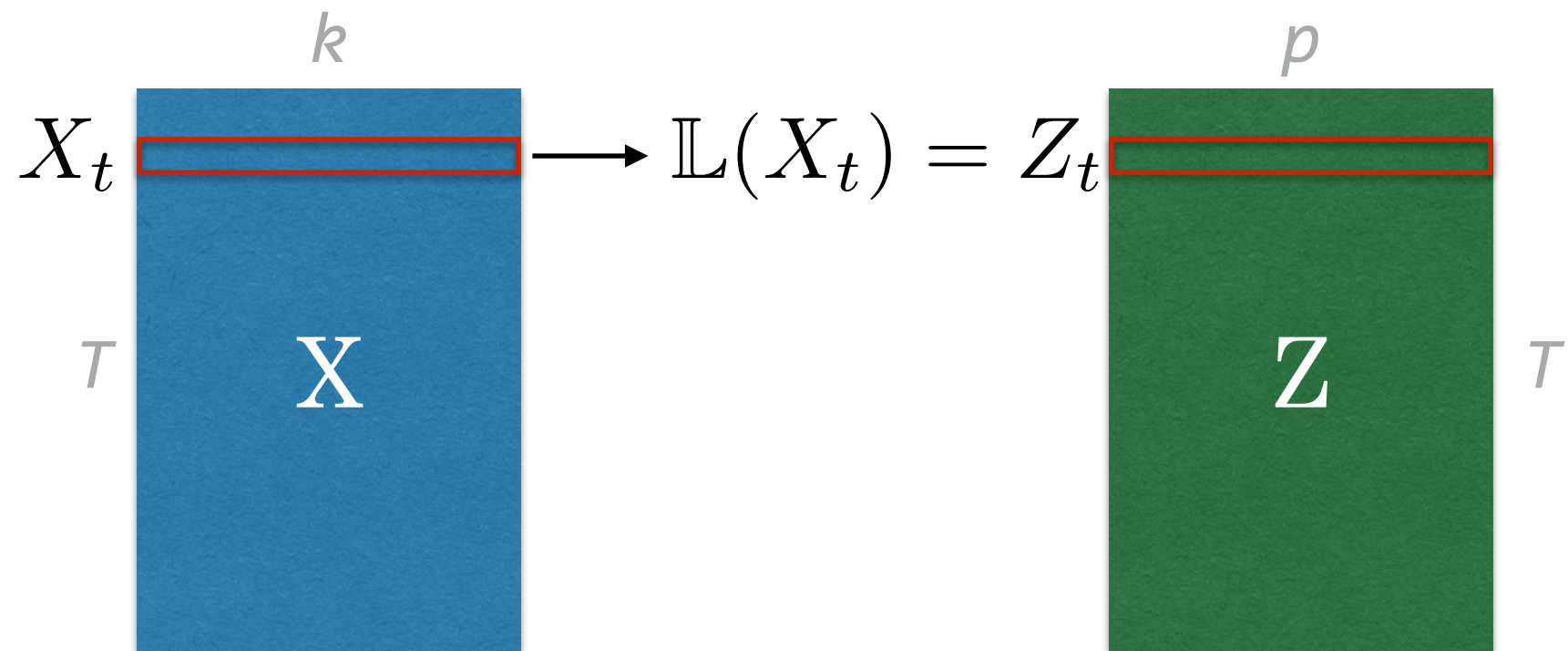
* Simplest version: time-invariant \mathbf{L}

$$X_t \in \mathbb{R}^k \quad k\text{-dim vector}$$

$$\mathbf{L}(X_t) = Z_t \in \mathbb{R}^p \quad p\text{-dim vector}$$

LINEARIZING TRANSFORMATIONS

* Simplest version: time-invariant \mathbb{L}



LINEARIZING TRANSFORMATIONS

- * Simplest version: time-invariant L
- * **Example:** X is an image with k pixels, Z is a binary vector saying which of p object categories are present in the image

LINEARIZING TRANSFORMATIONS

- * Simplest version: time-invariant L
- * Only suitable for some situations

LINEARIZING TRANSFORMATIONS

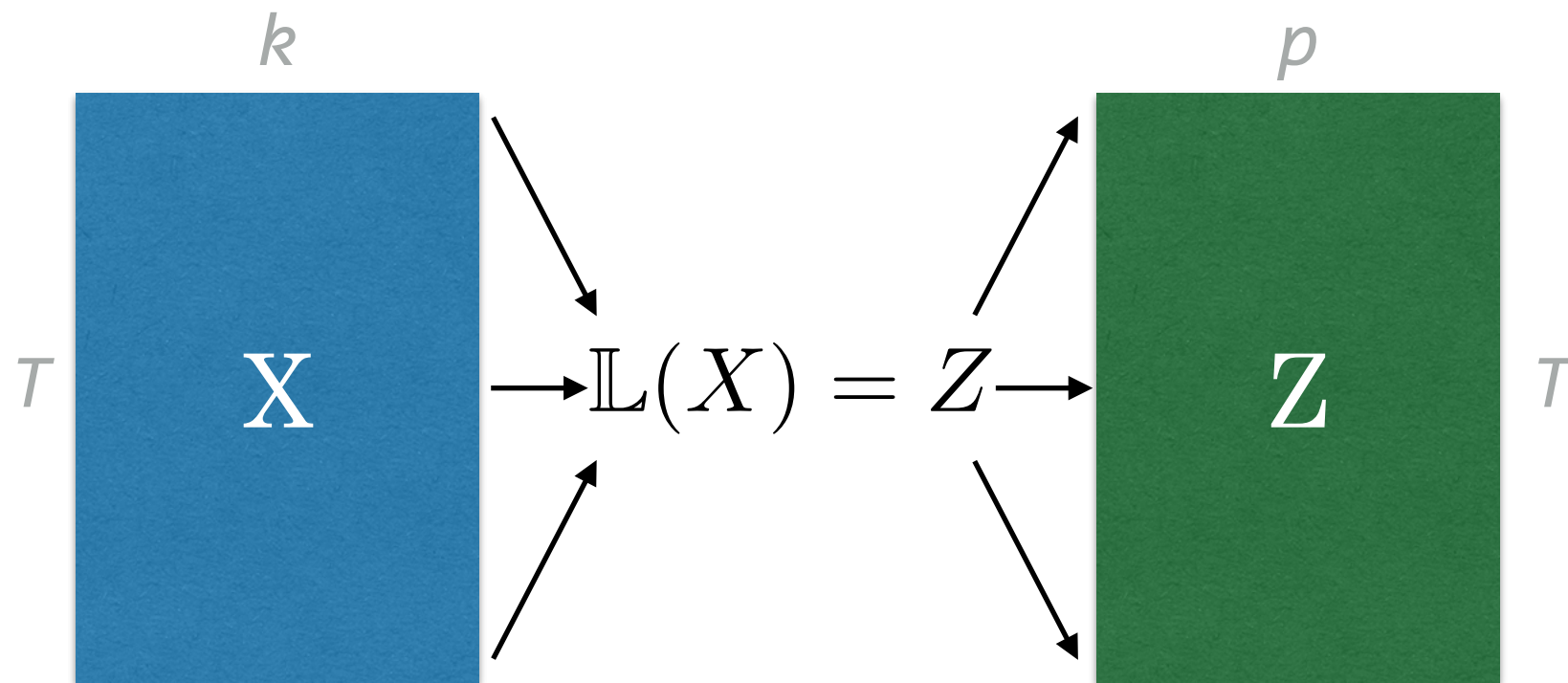
* More complex version: time-dependent \mathbf{L}

$$X \in \mathbb{R}^{T \times k} \quad T \times k \text{ matrix}$$

$$\mathbb{L}(X) = Z \in \mathbb{R}^{T \times p} \quad T \times p \text{ matrix}$$

LINEARIZING TRANSFORMATIONS

* More complex version: time-dependent L



LINEARIZING TRANSFORMATIONS

- * More complex version: time-dependent L
- * **Example:** X_t represents which word was presented as a 1-hot vector. Z_t is the part-of-speech for that word. Time-dependence is necessary because PoS depends on word context.

e.g. “I refuse” vs. “Refuse bin”

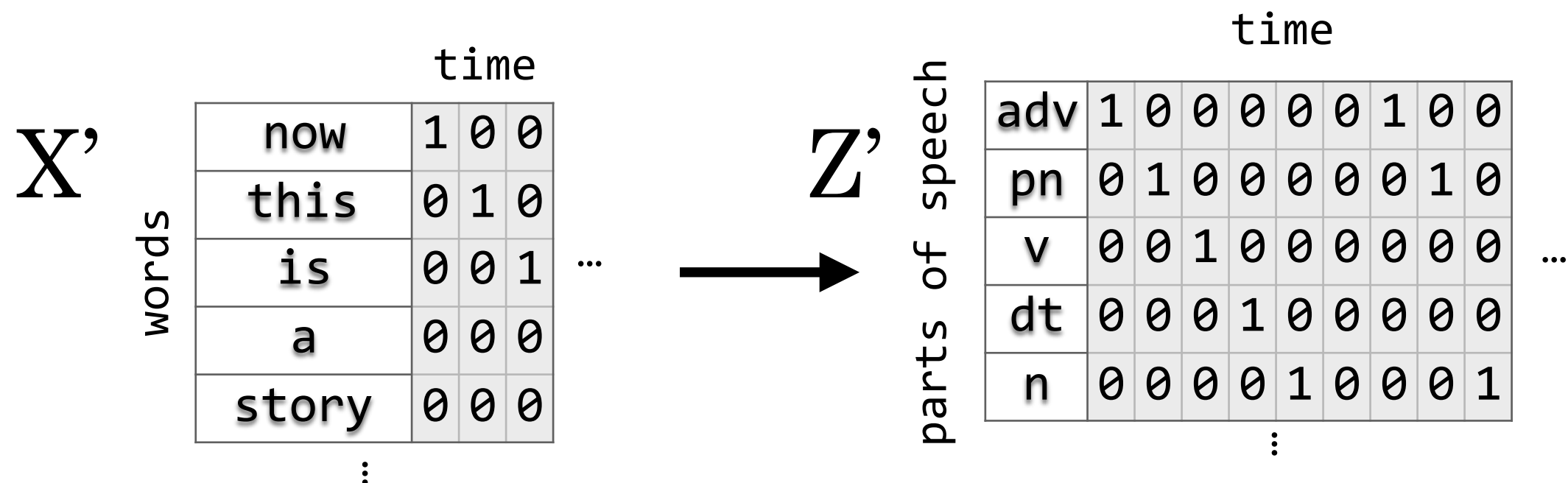
SYNTAX - PART OF SPEECH

“Now this is a story all about how my life

{adv} {pn} {v} {dt} {n} {adj} {prep} {adv} {pn} {n}

got flipped-turned upside down...”

{v-p} {v-p} {v-p} {prep phrase}



SYNTAX - HMM

- * **Hidden Markov model (HMM)**
- * **Example of a language model:** a function that, given the previous words, returns a probability distribution over the next word.

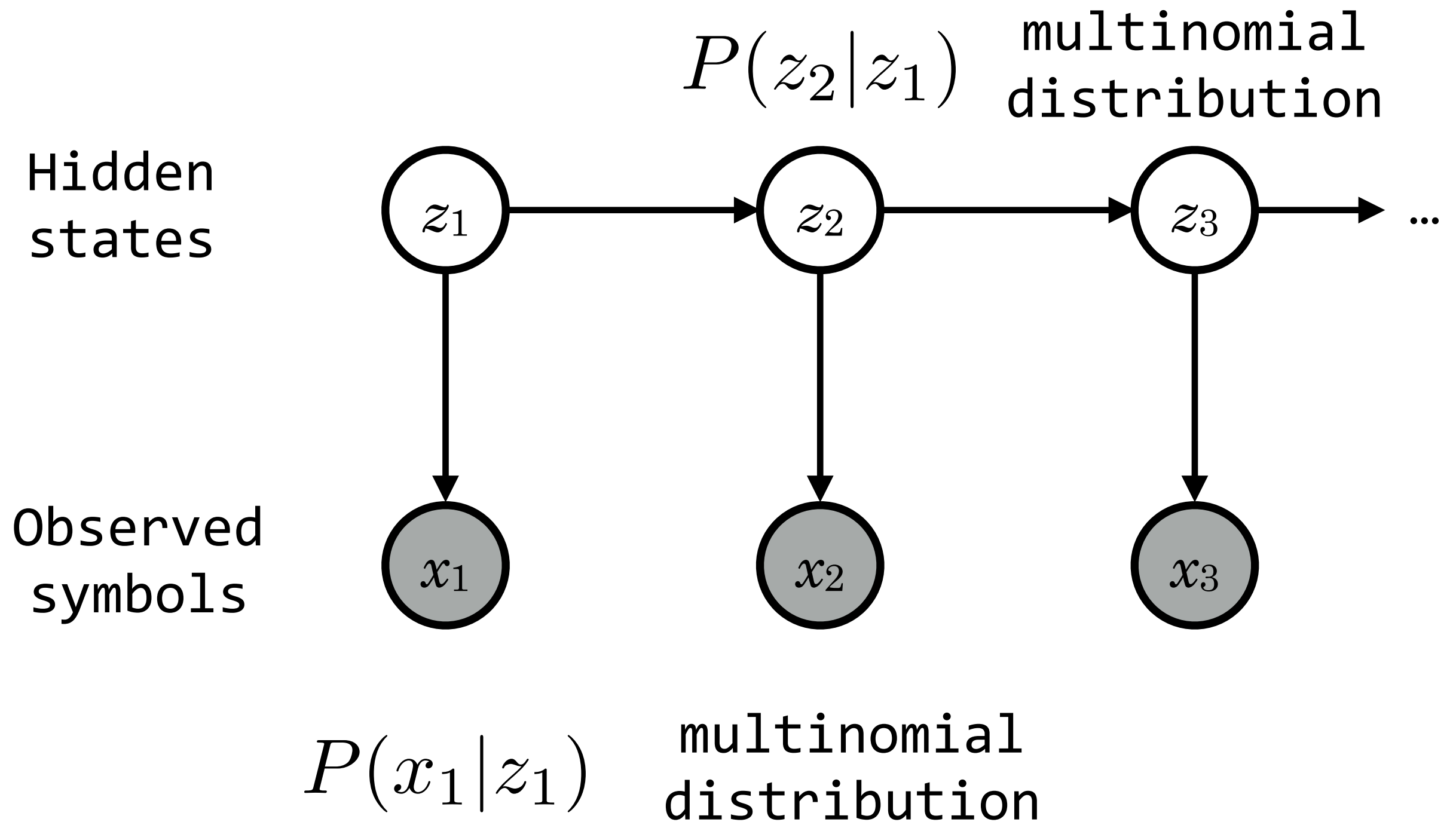
$$P(w_i | w_1 \dots w_{i-1}; \theta, \phi)$$

next word *previous words* *HMM parameters*

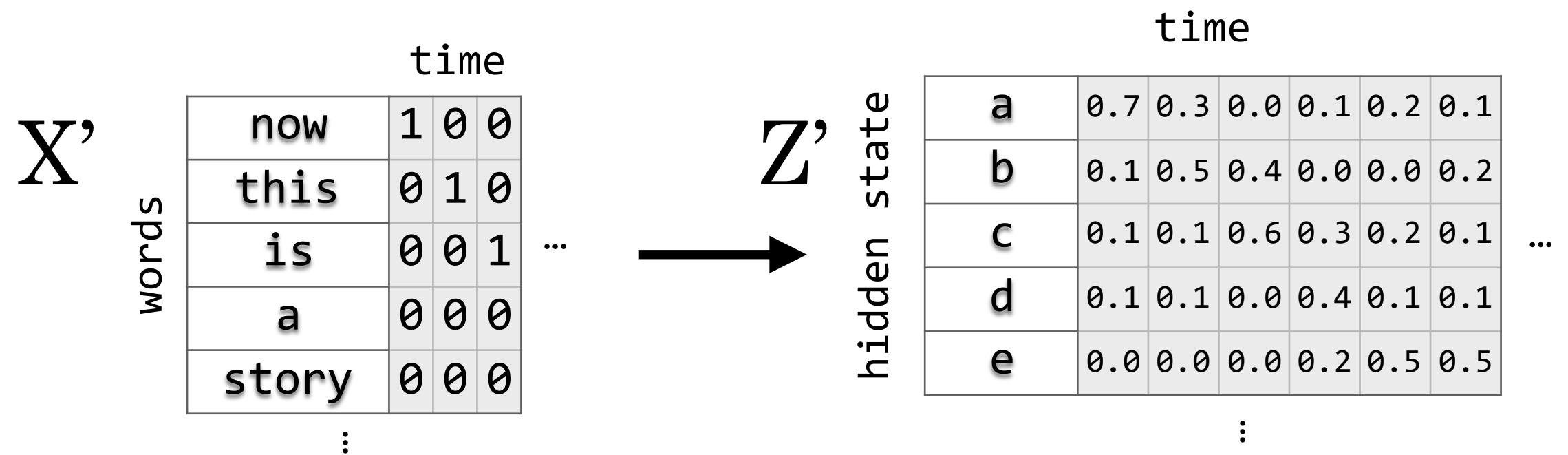
SYNTAX – HMM

- * *We get:* a sequence of observed symbols
[“now”, “this”, “is”, “a”, “story”, ...]
- * *We think:* there are hidden, underlying
states
[a, b, c, d, e, ...]

SYNTAX - HMM



SYNTAX – HMM



SYNTAX - HMM

- * **Learning time:** we know x , what are theta and phi?
- * **Objective:** find theta and phi that *maximize probability* of observed x
- * **Hyperparameters:** number of hidden states (number of possible Z values), priors on theta and phi

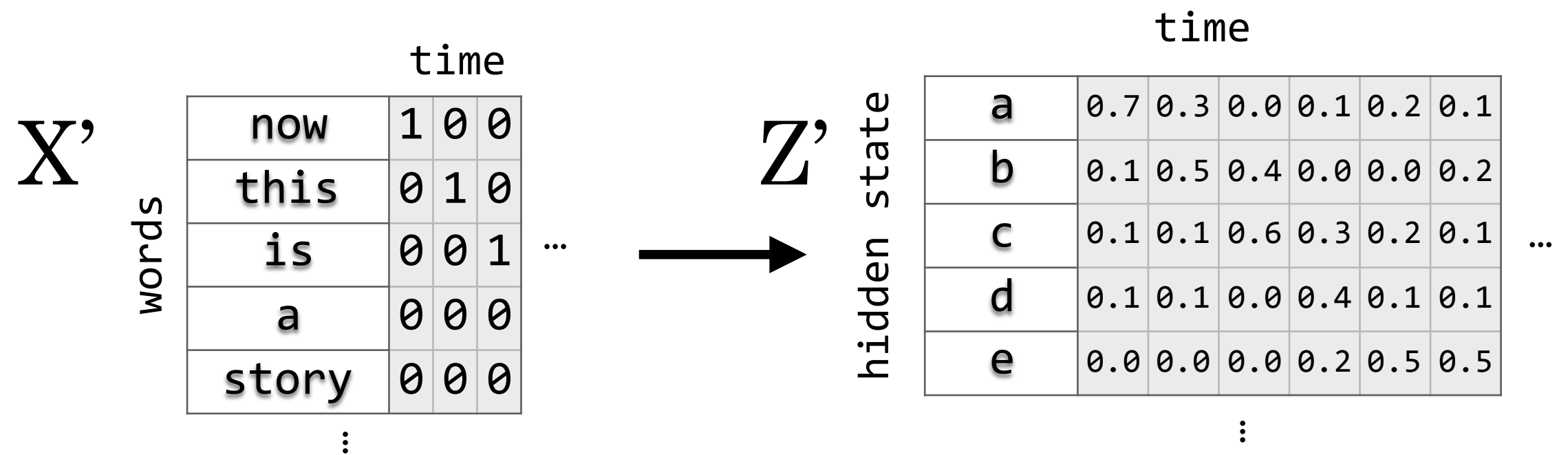
SYNTAX - HMM

- * **Learning time:** (The easy way) Markov chain Monte Carlo (MCMC) w/ *Gibbs sampling*
- * (Init.) Guess random Z
- * Update θ & ϕ based on current X , Z
- * For each Z_t , update based on Z_{t-1} , X_t , ϕ , & θ
- * Repeat, repeat, repeat, repeat, repeat...

SYNTAX - HMM

- * **Inference time:** we know x , we know θ , we know ϕ ; what is $P(z|x)$?
- * Finally, use inferred state probabilities as features in a linearized model!

SYNTAX - HMM



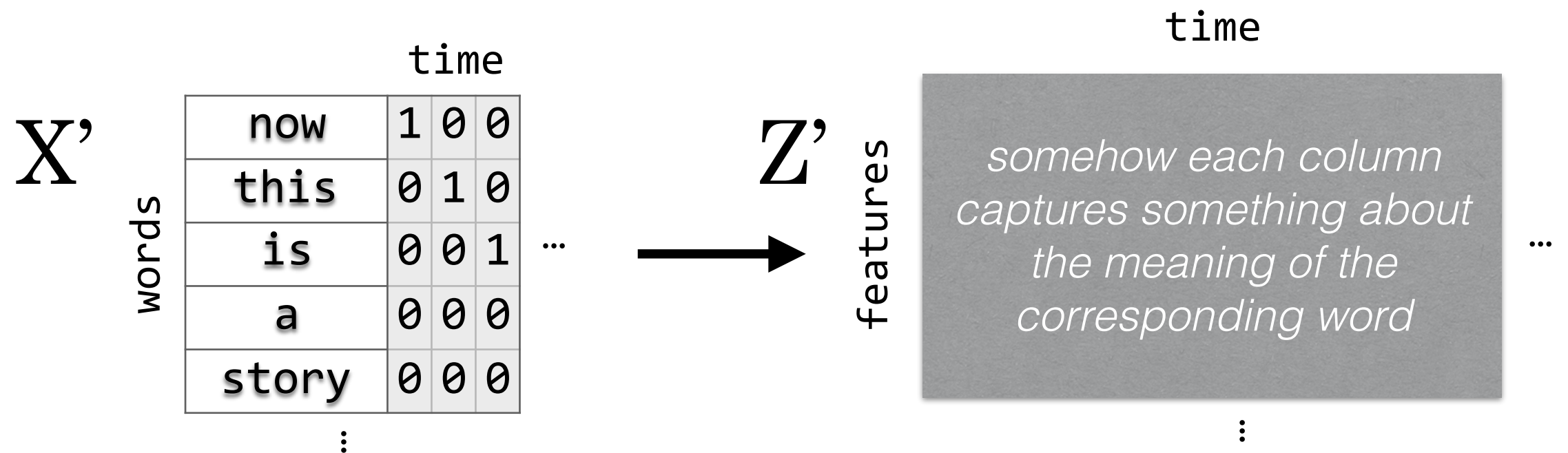
LEXICAL SEMANTICS

- * Let's create an L that captures word-level semantic information
- * Unlike the `awful` syntax models, this model will be *time-invariant*

LEXICAL SEMANTICS

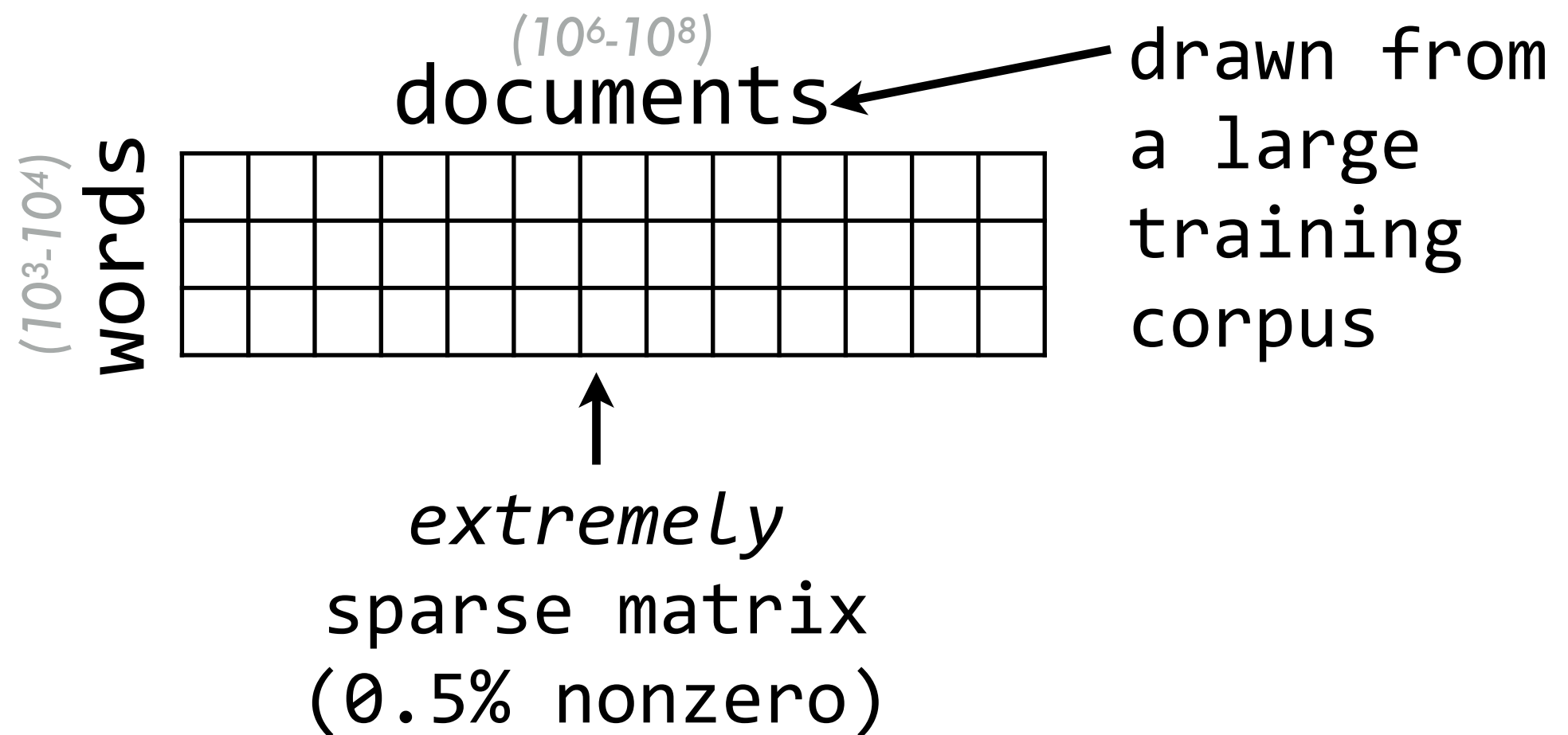
$$\mathbb{L}_{semantic}(X_t) = Z_t \in \mathbb{R}^p$$

1-hot vector ↓ *p-dim vector*



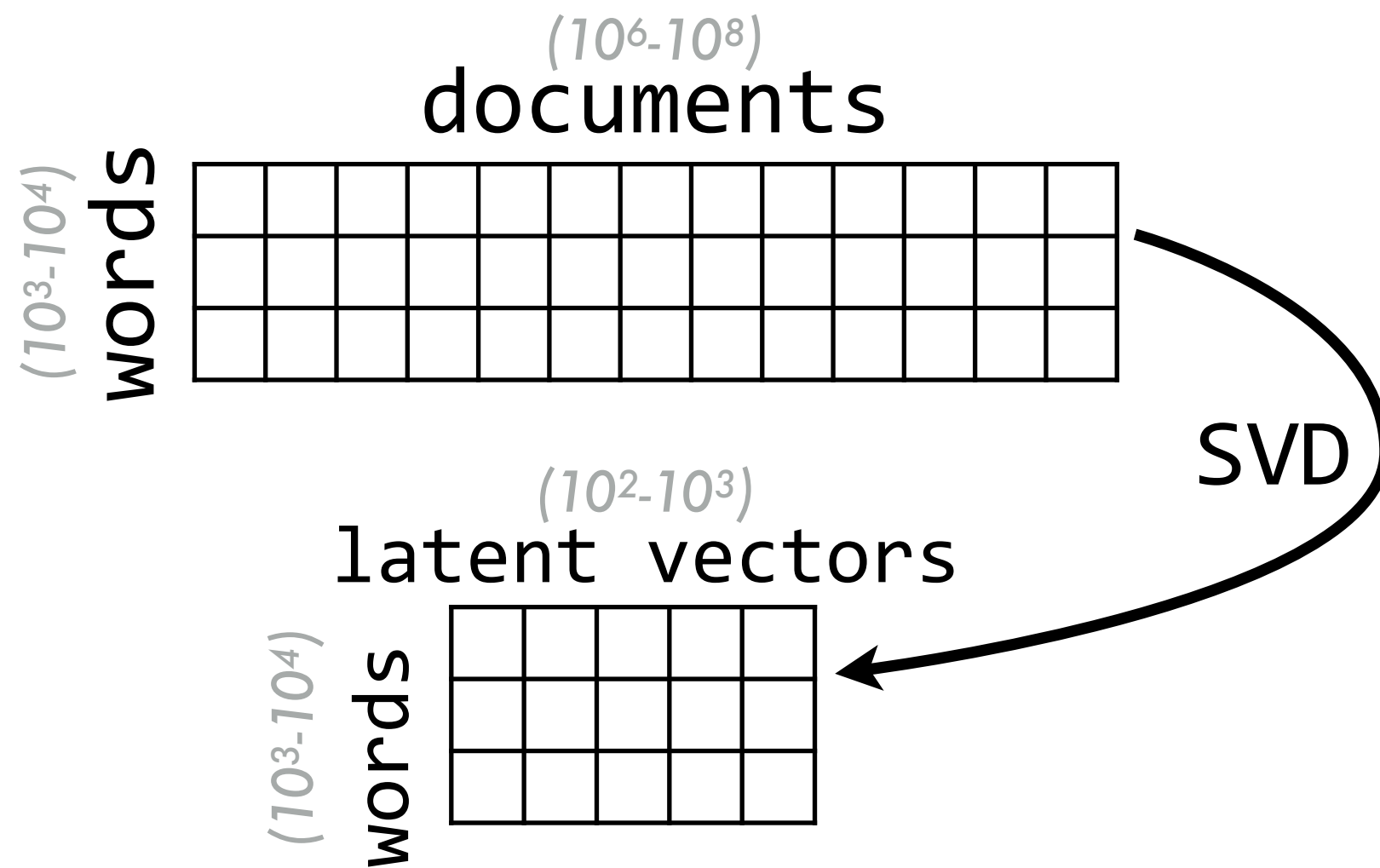
LEXICAL SEM. - LSA

- * Latent Semantic Analysis (LSA)

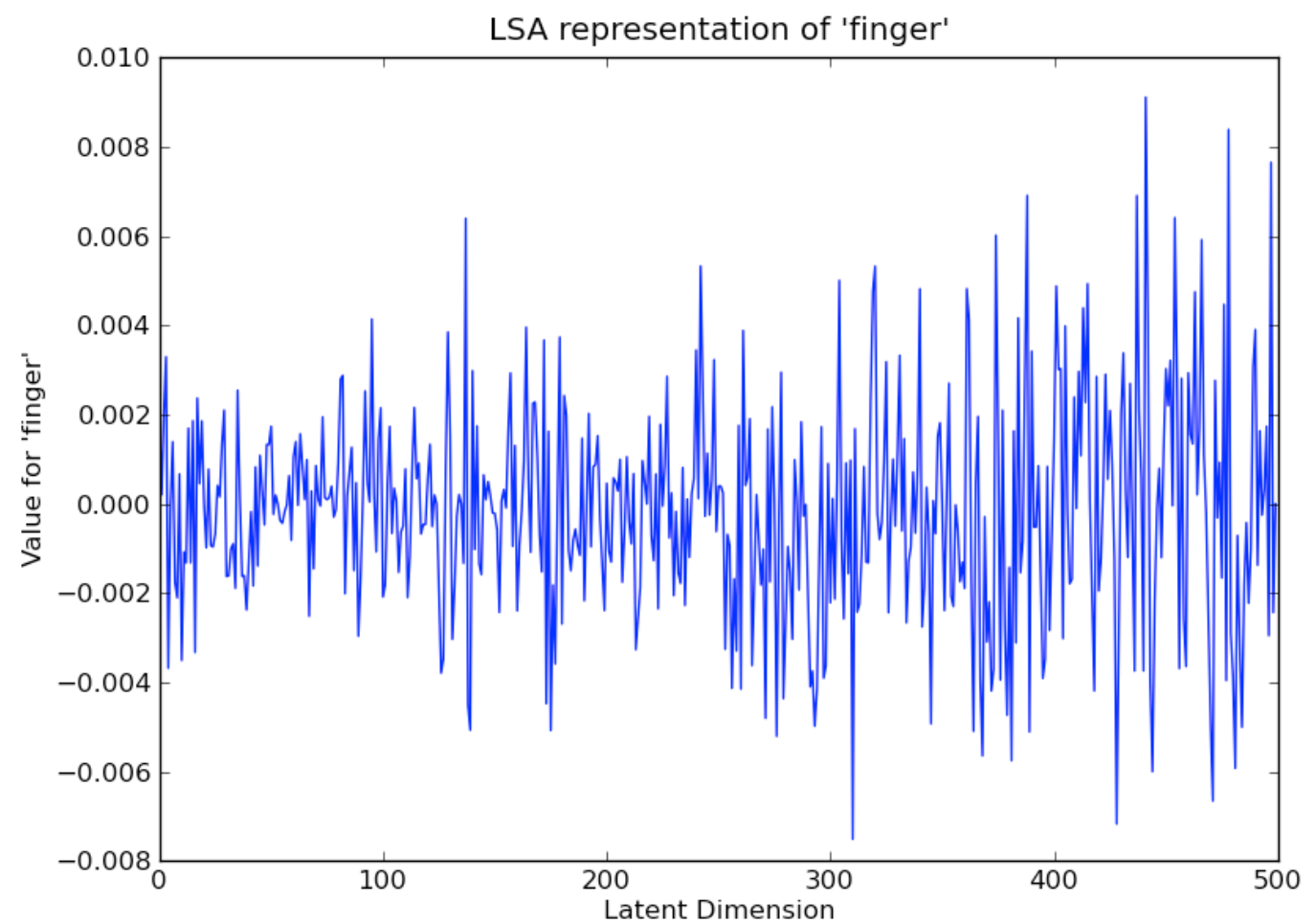
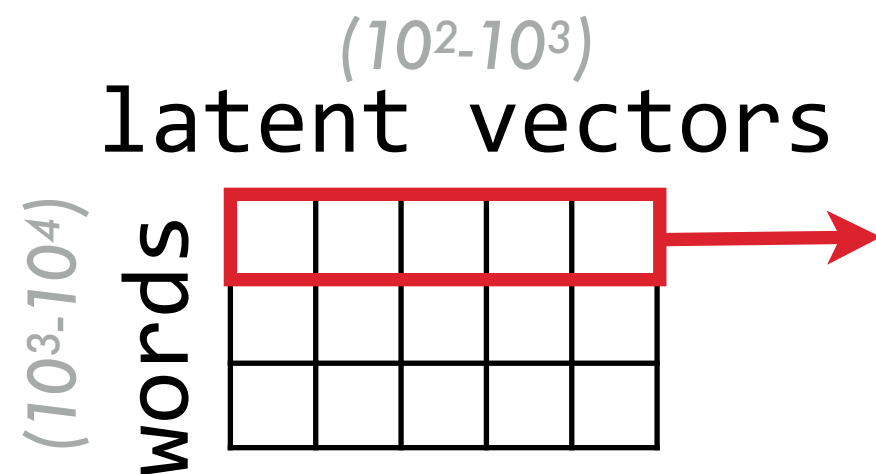


LEXICAL SEM. - LSA

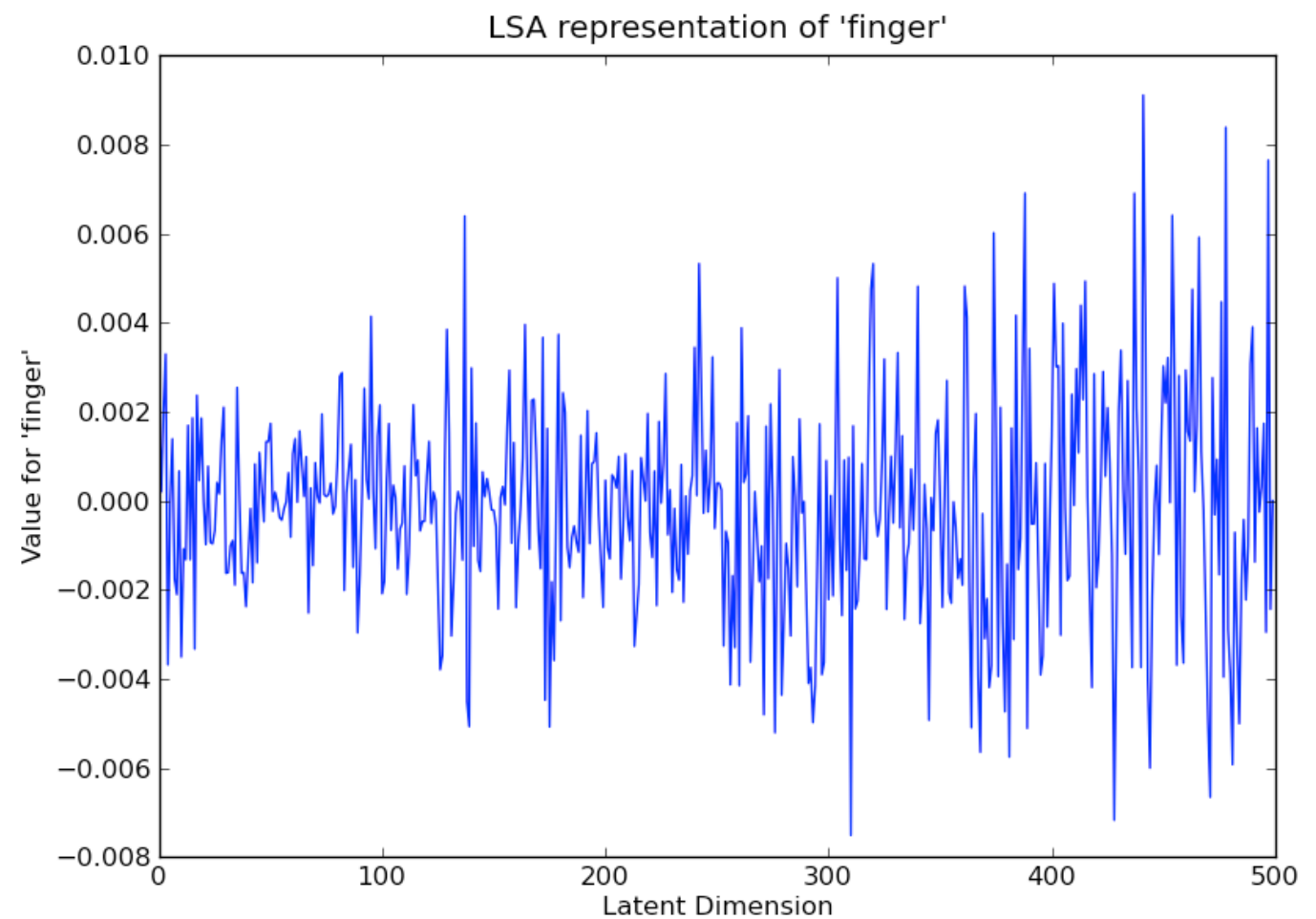
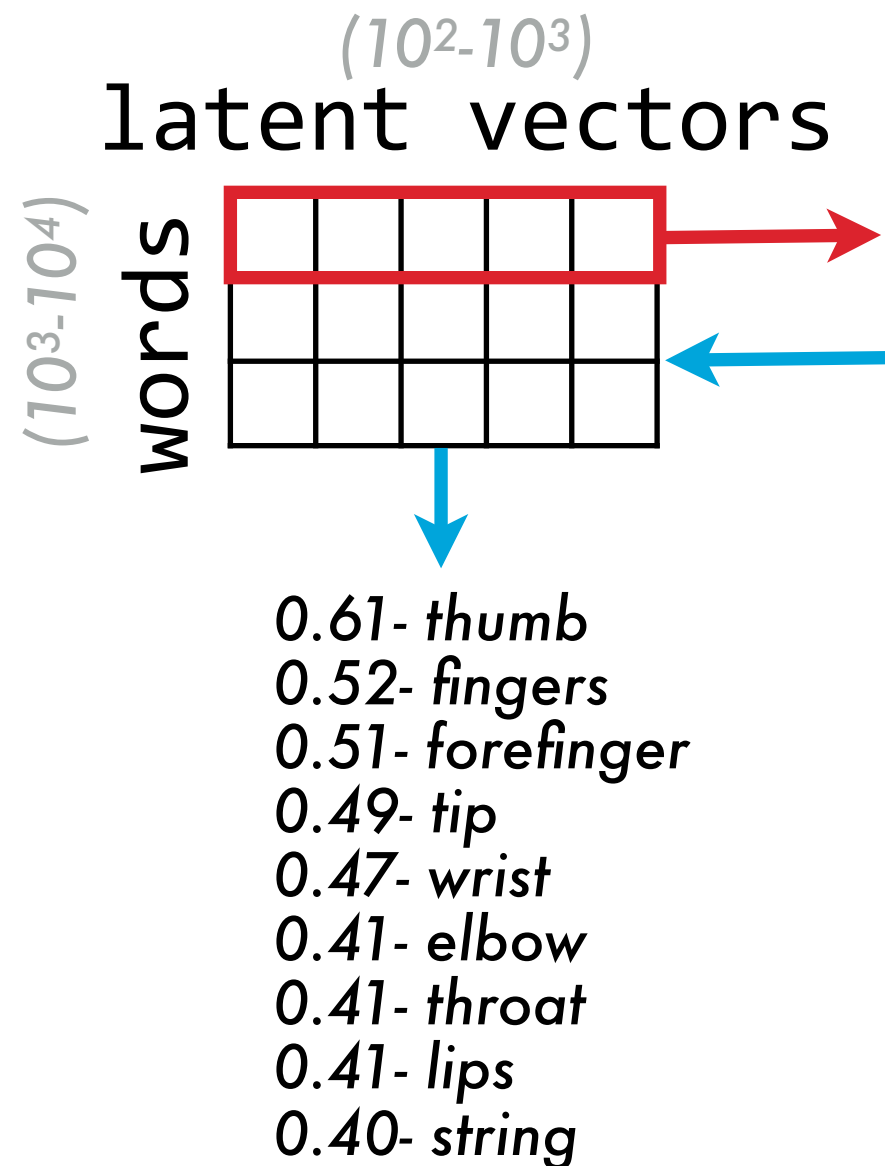
* Latent Semantic Analysis (LSA)



LEXICAL SEM. - LSA



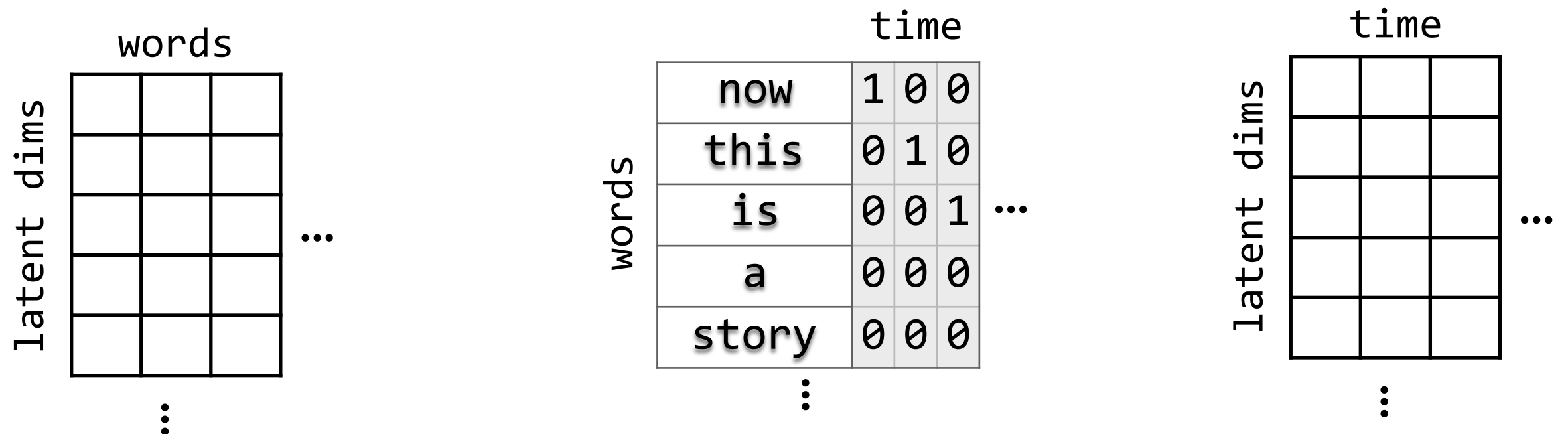
LEXICAL SEM. - LSA



LEXICAL SEM. - LSA

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

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



***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
COVARIANCEINVERSE OF
PENALTY
INNER PRODUCTEMBEDDING
INNER PRODUCT

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

NEXT TIME

* MORE SEMANTICS!!