

344.075 KV: Natural Language Processing

Language Modeling & Neural Word Embedding



Navid Rekab-Saz

navid.rekabsaz@jku.at

Agenda

- n -gram language models
- Language modeling with neural networks
- word2vec

Agenda

- ***n*-gram language models**
- Language modeling with neural networks
- word2vec

Language Modeling

- **Language Modeling** is the task of predicting a word (or a subword or character) given a context:

$$P(v|\text{context})$$

- A Language Model can answer the questions like



$$P(v|the\ students\ opened\ their)$$

Language Modeling – formal definition

- Given a sequence of words $x^{(1)}, x^{(2)}, \dots, x^{(t)}$, a **language model** calculates the **probability distribution** of next word $x^{(t+1)}$ over all words in vocabulary

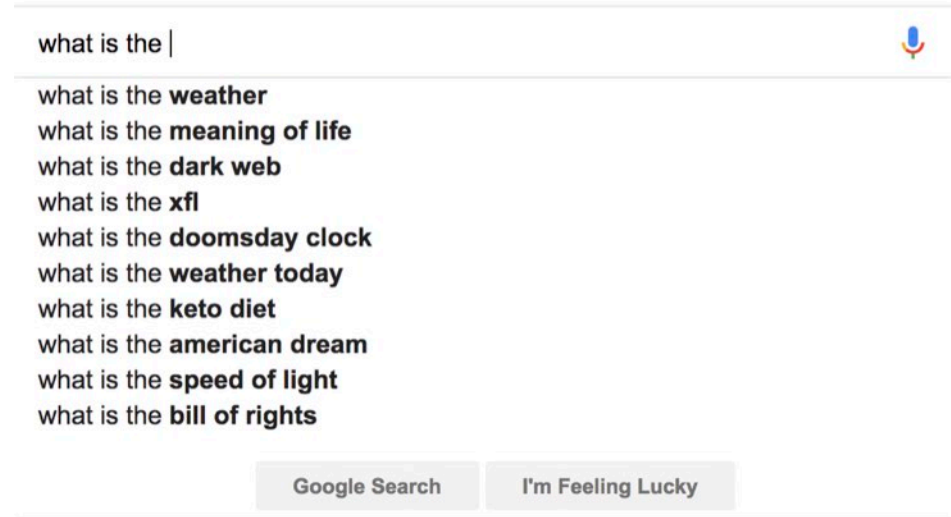
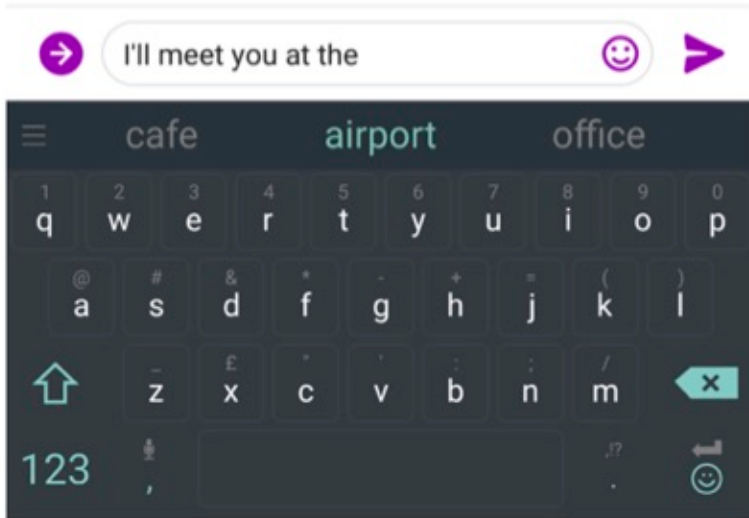
$$P(x^{(t+1)} | x^{(t)}, x^{(t-1)}, \dots, x^{(1)})$$

x is any word in the vocabulary $\mathbb{V} = \{v_1, v_2, \dots, v_N\}$

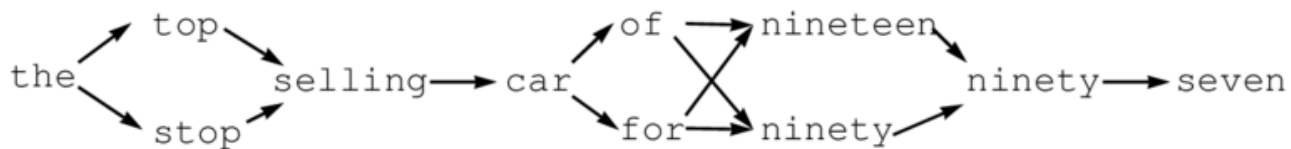
Why Language Modeling?

- Language Modeling is a **benchmark task** that helps us measure our progress on **understanding language**
- Language Models is a **subcomponent** of many NLP tasks, especially those involving **generating text** or estimating the **probability of text**:
 - Predictive typing
 - Spelling/grammar correction
 - Speech recognition
 - Handwriting recognition
 - Machine translation
 - Summarization
 - Dialogue /chatbots
 - etc.

Language Modeling



Language Modeling



Input Lattice



Lattice Rescoring
Tool



the top selling car of nineteen ninety seven
the top selling car of ninety ninety seven
the top selling car for ninety ninety seven
the stop selling car for nineteen ninety seven

N-Best List

n-gram Language Model

- Recall: a *n*-gram is a chunk of *n* consecutive words.

the students opened their _____

- unigrams: “*the*”, “*students*”, “*opened*”, “*their*”
 - bigrams: “*the students*”, “*students opened*”, “*opened their*”
 - trigrams: “*the students opened*”, “*students opened their*”
 - 4-grams: “*the students opened their*”
-
- A *n*-gram Language Model collects frequency statistics of different *n*-grams in a corpus, and use these to calculate probabilities

n-gram Language Model



- **Markov assumption**: decision at time t depends only on the current state
- In *n*-gram Language Model: predicting $x^{(t+1)}$ depends on preceding *n-1* words
- Without Markovian assumption:

$$P(x^{(t+1)} | x^{(t)}, x^{(t-1)}, \dots, x^{(1)})$$

- *n*-gram Language Model:

$$P(x^{(t+1)} | x^{(t)}, x^{(t-1)}, \dots, x^{(t-n+2)})$$


n-1 words

n-gram Language Model

- Based on definition of conditional probability:

$$P(x^{(t+1)} | x^{(t)}, \dots, x^{(t-n+2)}) = \frac{P(x^{(t+1)}, x^{(t)}, \dots, x^{(t-n+2)})}{P(x^{(t)}, \dots, x^{(t-n+2)})}$$

- The *n*-gram probability is calculated by counting *n*-grams and [*n*−1]-grams in a large corpus of text:

$$P(x^{(t+1)} | x^{(t)}, \dots, x^{(t-n+2)}) \approx \frac{\text{count}(x^{(t+1)}, x^{(t)}, \dots, x^{(t-n+2)})}{\text{count}(x^{(t)}, \dots, x^{(t-n+2)})}$$

n-gram Language Model

- Example: learning a 4-gram Language Model

~~as the exam clerk started the clock, the students opened their~~ _____

condition on this

$$P(v | \text{students opened their}) = \frac{P(\text{students opened their } v)}{P(\text{students opened their})}$$

- For example, suppose that in the corpus:
 - “students opened their” occurred 1000 times
 - “students opened their books” occurred 400 times
 - $P(\text{books} | \text{students opened their}) = 0.4$
 - “students opened their exams” occurred 100 times
 - $P(\text{exams} | \text{students opened their}) = 0.1$

n-gram Language Models – generating text

- A trigram Language Model trained on Reuters corpus (1.7 M words)

today the _____

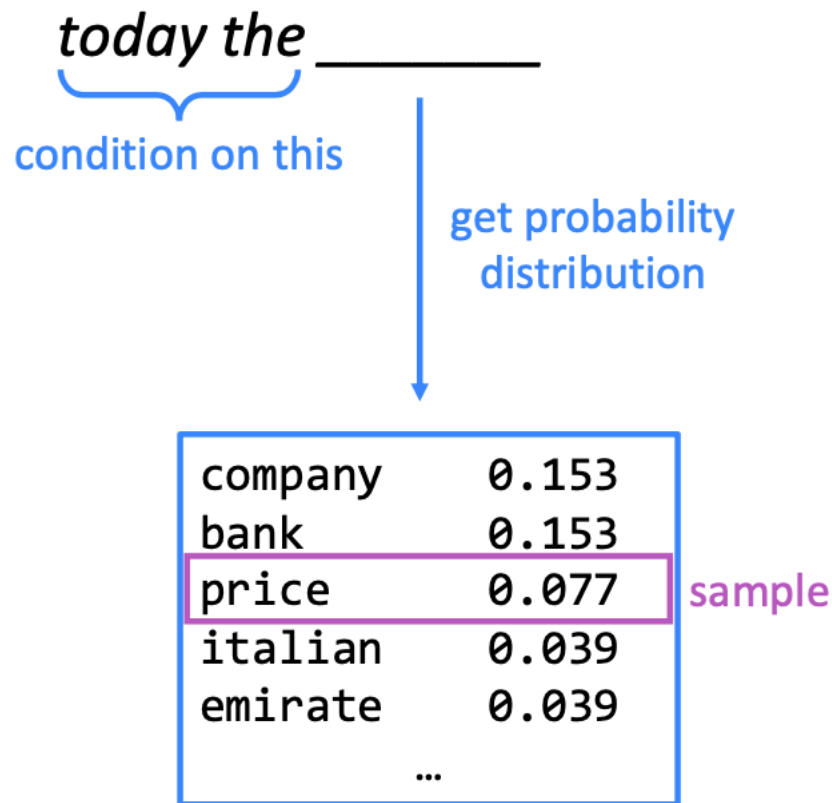
get probability
distribution

company	0.153
bank	0.153
price	0.077
italian	0.039
emirate	0.039
...	

Sparsity problem:
not much granularity
in the probability
distribution

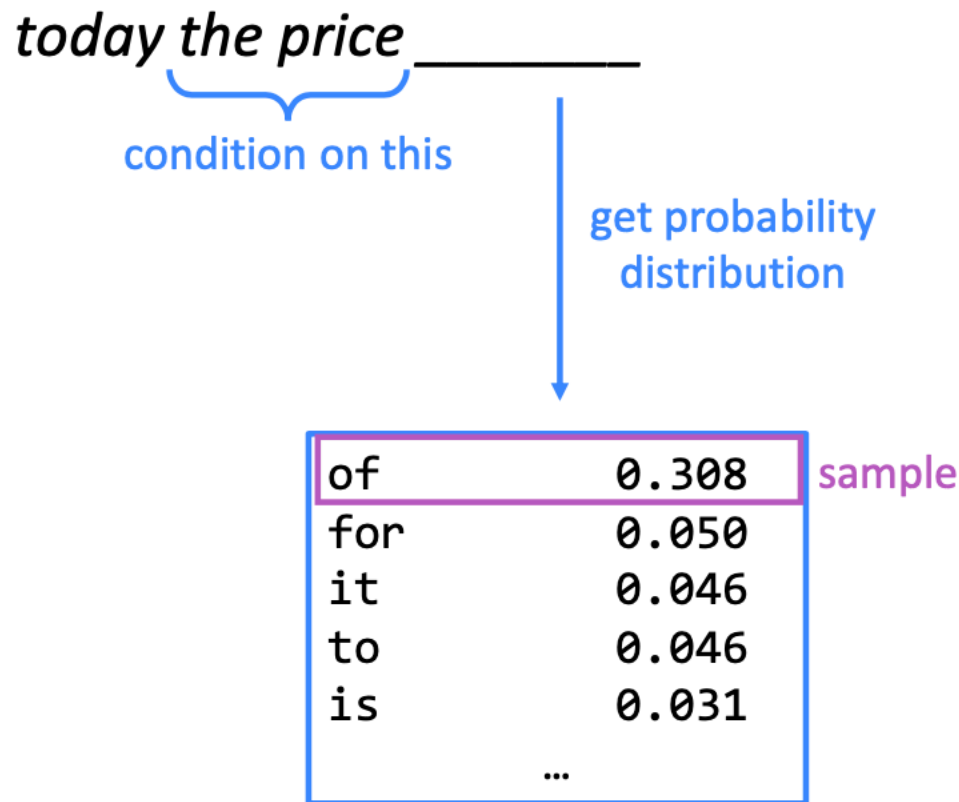
n-gram Language Models – generating text

- Generating text by sampling from the probability distributions



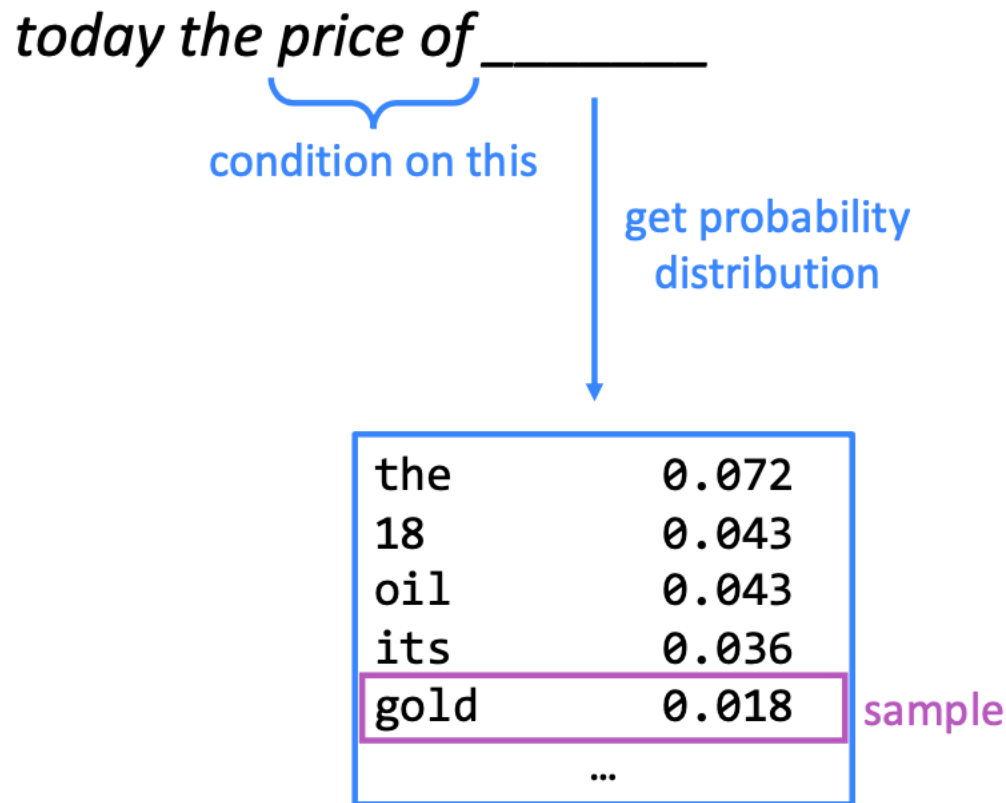
n-gram Language Models – generating text

- Generating text by sampling from the probability distributions



n-gram Language Models – generating text

- Generating text by sampling from the probability distributions



n -gram Language Models – generating text

- Generating text by sampling from the probability distributions

today the price of gold per ton , while production of shoe lasts and shoe industry , the bank intervened just after it considered and rejected an imf demand to rebuild depleted european stocks , sept 30 end primary 76 cts a share .

- Very good in **syntax** ... but **incoherent**!
- Increasing n makes the text more coherent but also intensifies the discussed issues

n-gram Language Model – issues

■ Sparsity

- What if the nominator „*students opened their v*“ never occurred in corpus?
 - **Smoothing**: add small hyper-parameter δ to all words
- What if the denominator „*students opened their*“ never occurred in corpus?
 - **Backoff**: condition on “*students opened*” instead
- Increasing n makes the sparsity problem worse!

■ Storage

- The model needs to store all n -grams (from unigram to n -gram), observed in the corpus
- Increasing n radically worsens the storage problem!

Agenda

- n -gram language models
- **Language modeling with neural networks**
- word2vec

n-gram Language Model with neural networks

Recall

- The aim of a *n*-gram Language Model is to calculate:

$$P(x^{(t+1)} | x^{(t)}, \dots, x^{(t-n+2)})$$

- We can use a feed-forward neural network to calculate this probability
 - Smooth probability estimation
 - Benefitting from the semantic space of word embeddings

Feed-Forward Language Model – training data

- Creating training data for a 4-gram Feed-Forward Language Model in the form of (context , next word), namely $(x^{(t-2)}x^{(t-1)}x^{(t)}, x^{(t+1)})$:

- Text corpus:

a fluffy cat sunbathes on the bank of river ...

- Resulting training data items:

(<bos> <bos> <bos>, a)

(<bos> <bos> a , fluffy)

(<bos> a fluffy , cat)

(a fluffy cat , sunbathes)

(fluffy cat sunbathes , on)

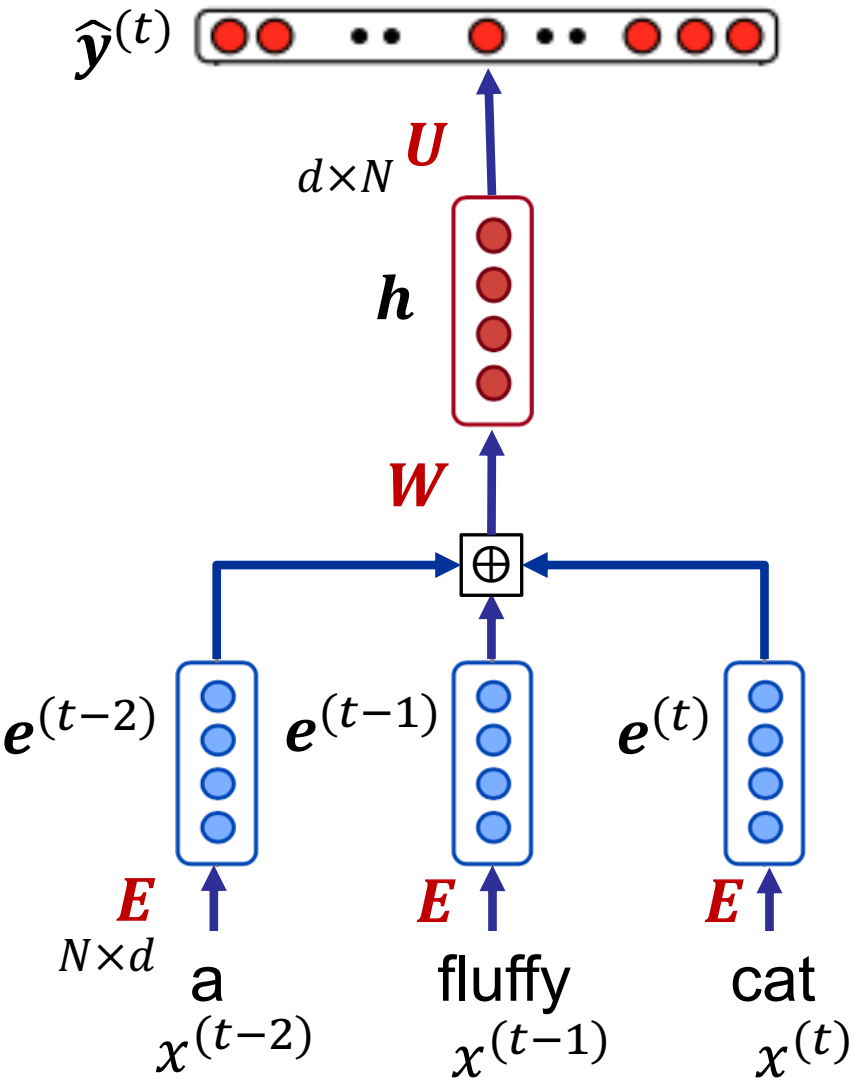
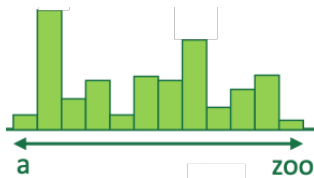
(cat sunbathes on , the)

(sunbathes on the , bank)

...

Feed-Forward Language Model

$P(\text{sunbathes} | \text{a fluffy cat})$



Training data item:
 $(\text{a fluffy cat}, \text{sunbathes})$

Feed-Forward Language Model – formulation

- Encoder

- One-hot vector of word $x^{(t)} \rightarrow \mathbf{x}^{(t)} \in \mathbb{R}^N$
- Fetching word embedding $\rightarrow \mathbf{e}^{(t)} = \mathbf{x}^{(t)} \mathbf{E}$

- Concatenation of word embeddings

$$\mathbf{e} = [\mathbf{e}^{(t-2)}, \mathbf{e}^{(t-1)}, \mathbf{e}^{(t)}]$$

- Hidden layer

$$\mathbf{h} = \tanh(\mathbf{W}\mathbf{e} + \mathbf{b})$$

- Decoder

- Predicted probability distribution:

$$\hat{\mathbf{y}}^{(t)} = \text{softmax}(\mathbf{U}\mathbf{h} + \mathbf{b}) \in \mathbb{R}^N$$

- Probability of any next word v at step t :

$$P(v | x^{(t)}, \dots, x^{(t-n+2)}) = \hat{y}_v^{(t)}$$

Feed-Forward Language Model – parameters

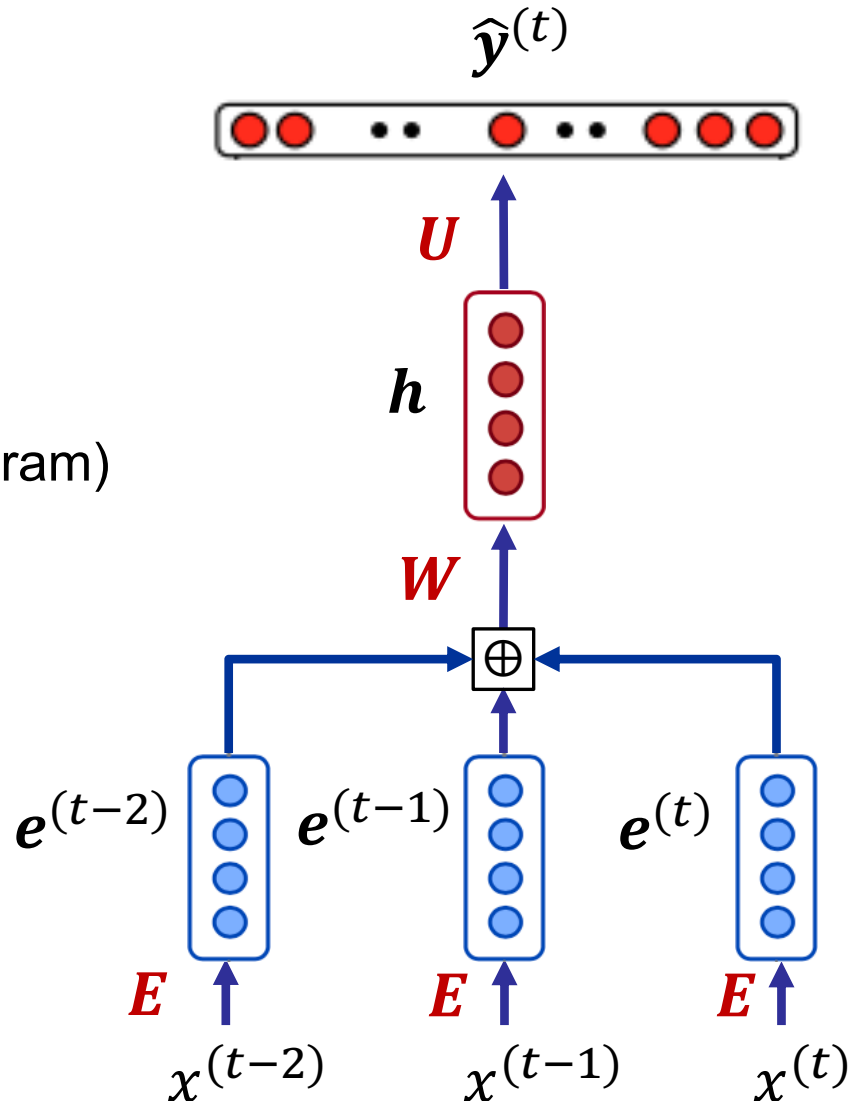
- $E \rightarrow N \times h$
- $W \rightarrow (n \times h) \times h$
- $U \rightarrow h \times N$

h embeddings dimension

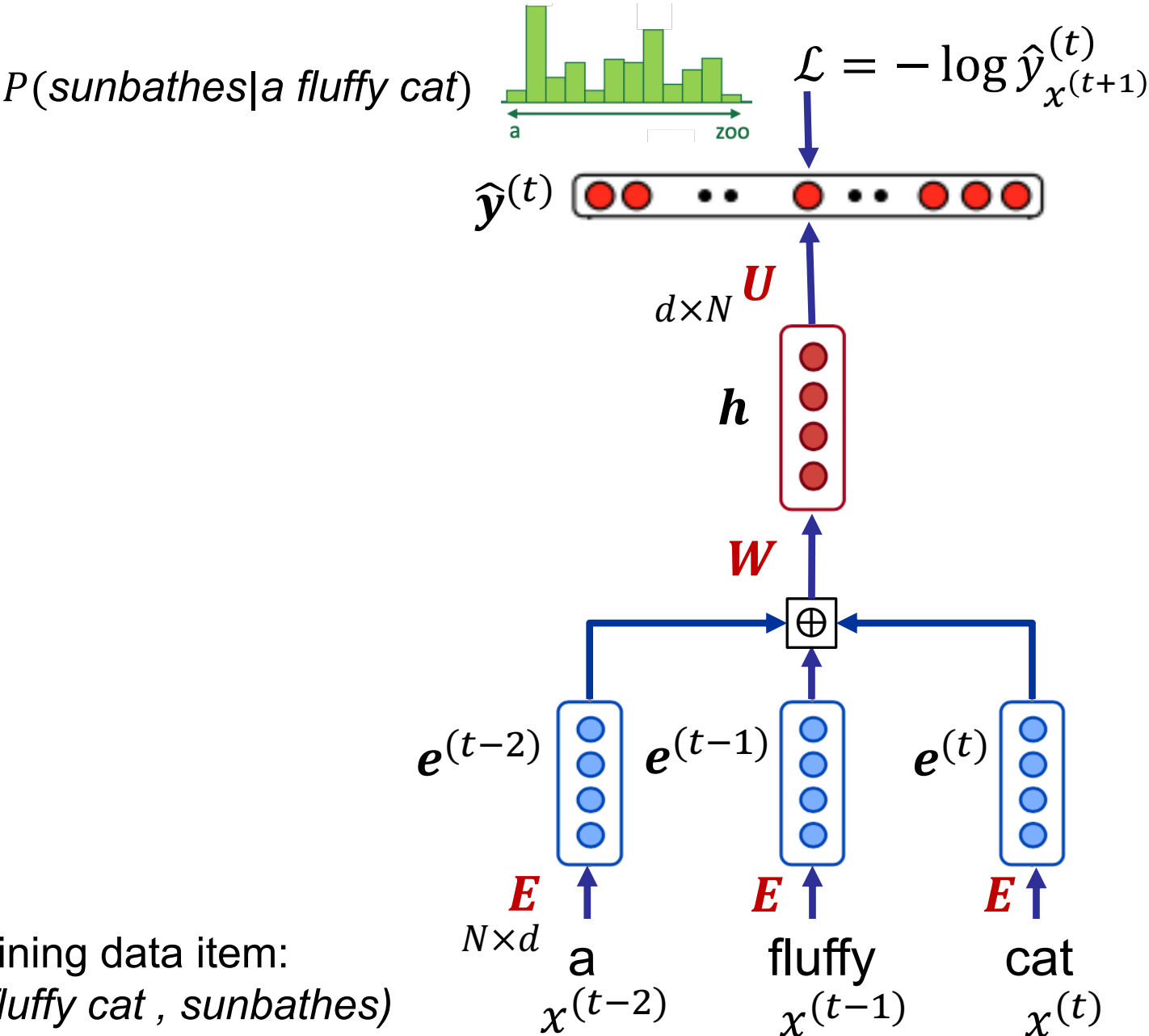
n number of preceding words (n -gram)

E is called **encoder embedding** or simply **word embedding**

U is called **decoder embedding** or **output projection**



Feed-Forward Language Model – loss



Training a Feed-Forward Language Model

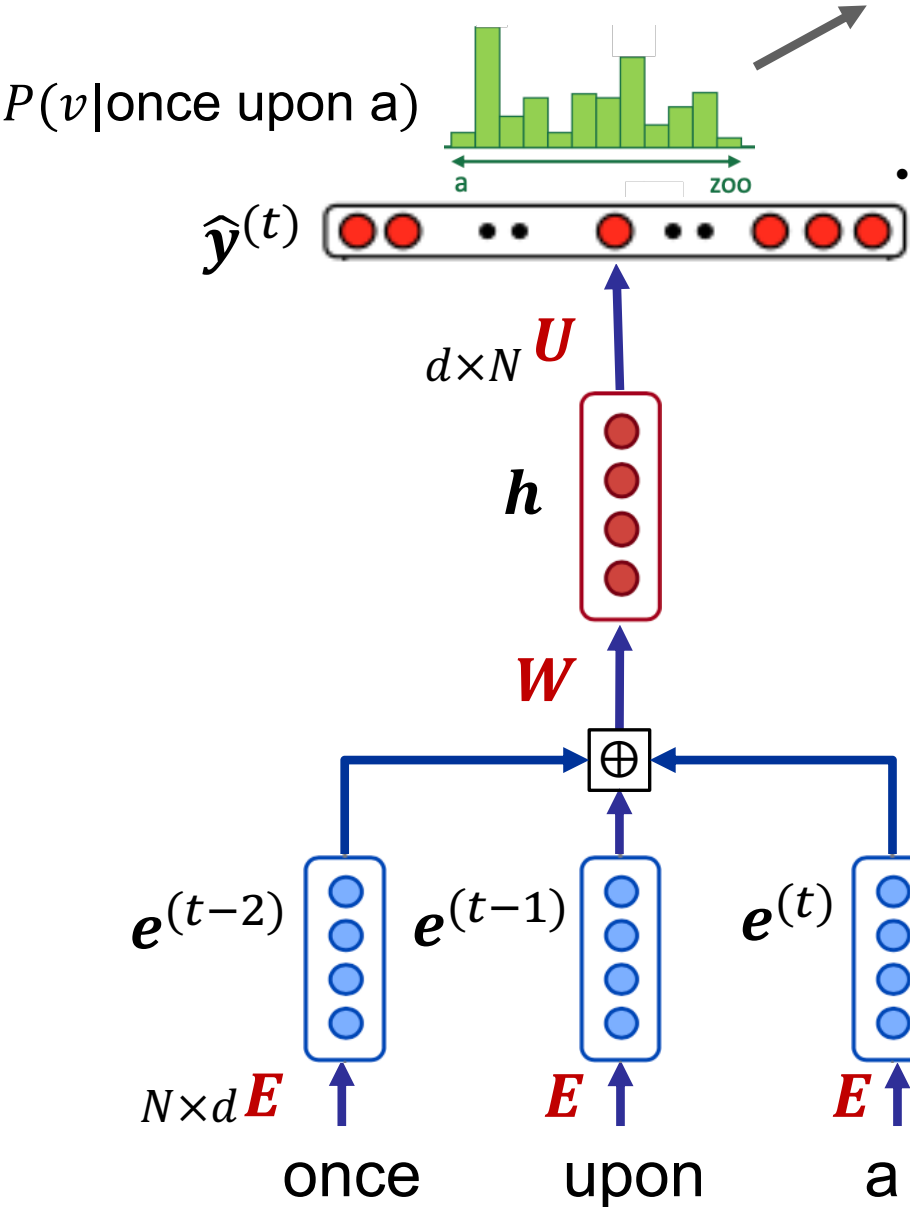
- Start with a large text corpus: $x^{(1)}, \dots, x^{(T)}$
- For every step t predict the output distribution $\hat{y}^{(t)}$ given n previous words
- Loss function at t is Negative Log Likelihood of the predicted probability of the word at $x^{(t+1)}$

$$\mathcal{L}^{(t)} = -\log \hat{y}_{x^{(t+1)}}^{(t)}$$

- Overall loss is the average over all time steps:

$$\mathcal{L} = \frac{1}{T} \sum_{t=1}^T \mathcal{L}^{(t)}$$

Feed-Forward Language Model – text generating



- Generate the next word by sampling from the probability distribution (e.g. next word: **time**)
- Use the generated word and continue generating next words

Generating text with Language Models

- Trained on Obama speeches



Jobs

The United States will step up to the cost of a new challenges of the American people that will share the fact that we created the problem. They were attacked and so that they have to say that all the task of the final days of war that I will not be able to get this done. The promise of the men and women who were still going to take out the fact that the American people have fought to make sure that they have to be able to protect our part. It was a chance to stand together to completely look for the commitment to borrow from the American people.

Generating text with Language Models

- Trained on Trump speeches



make the country rich. it was terrible. but saudi arabia, they make a billion dollars a day. i was the king. i was the king. i was the smartest person in yemen, we have to get to business. i have to say, but he was an early starter. and we have to get to business. i have to say, donald, i can't believe it. it's so important. but this is what they're saying, dad, you're going to be really pro, growth, blah, blah. it's disgusting what's disgusting., and it was a 19 set washer and to go to japan. did you hear that character. we are going to have to think about it. but you know, i've been nice to me.

Agenda

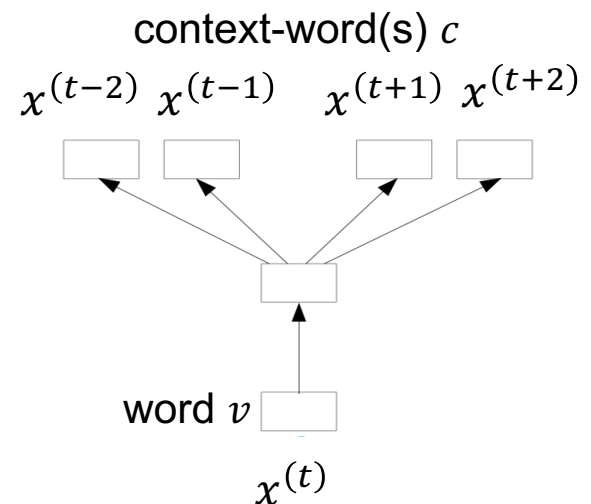
- n -gram language models
- Language modeling with neural networks
- **word2vec**

Word embedding with Neural Networks

- Word embeddings are a by-product of neural Language Models
- To create word embedding, we use a neural Language Model ...
 - but instead of predicting the next word, ...
 - ... the model predicts the probability of appearance of a context-word c in a window around the word v

$$P(c|v)$$

- This approach is known as skip-gram



drink sacred beer
ritual
Tesgüino
corn
fermented
Mexico Tarahumara people

$$P(\text{drink}|\text{Tesgüino}) = ?$$

Training data \mathcal{D}

- Creating training data with a window size of 2 in the form of (word , context- word), namely (v, c) :

Tarahumara	people	drink
------------	--------	-------

 Tesgüino while following rituals ...

(Tarahumara , people)

(Tarahumara , drink)

Tarahumara	people	drink	Tesgüino
------------	--------	-------	----------

 while following rituals ...

(people , Tarahumara)

(people , drink)

(people , Tesgüino)

...

Tarahumara	people	drink	Tesgüino	while	following
------------	--------	-------	----------	-------	-----------

 rituals ...

(Tesgüino , people)

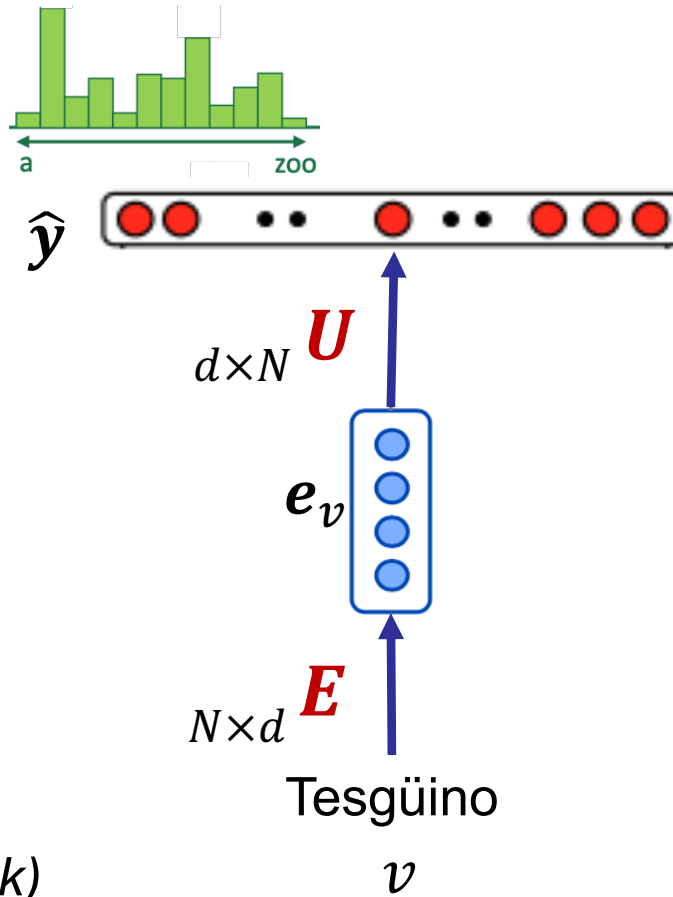
(Tesgüino , drink)

(Tesgüino , while)

(Tesgüino , following)

Neural word embedding – architecture

$$P(c|v) = P(\text{drink}|\text{Tesgüino})$$



Training data: (*Tesgüino* , *drink*)

Neural word embedding – formulation

■ Encoder

- One-hot vector of word $v \rightarrow \mathbf{v} \in \mathbb{R}^N$
- Fetching word embedding $\rightarrow \mathbf{e}_v = \mathbf{v}\mathbf{E}$

■ Decoder

- Predicted probability distribution:

$$\hat{\mathbf{y}} = \text{softmax}(\mathbf{U}\mathbf{e}_v) \in \mathbb{R}^N$$

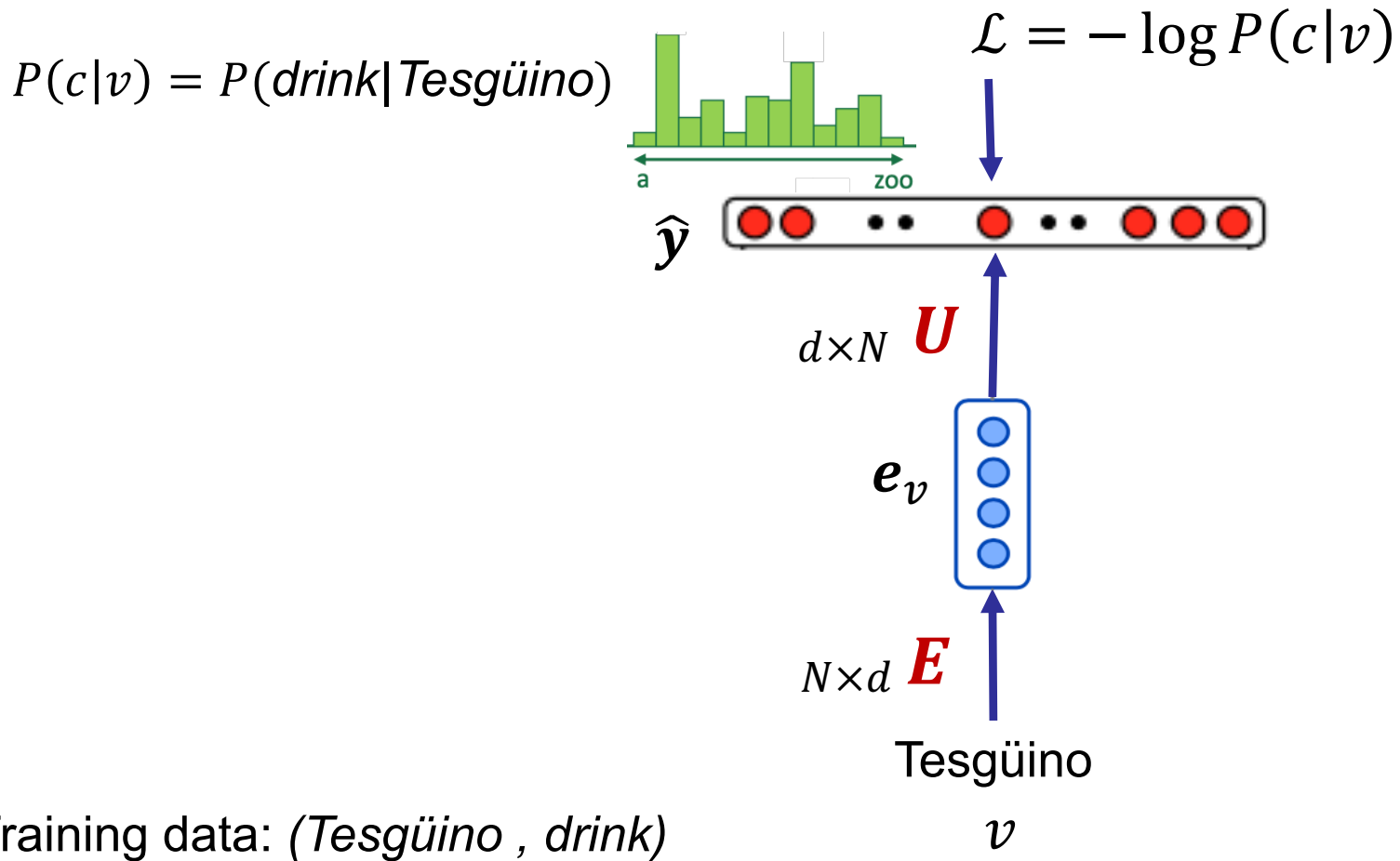
- Probability of an arbitrary context-word c given the word v :

$$P(c|v) = \hat{y}_c$$

■ Putting all together:

$$P(c|v) = \text{softmax}(\mathbf{U}\mathbf{e}_v)_c = \frac{\exp(\mathbf{e}_v \mathbf{u}_c)}{\sum_{\tilde{c} \in \mathbb{V}} \exp(\mathbf{e}_v \mathbf{u}_{\tilde{c}})}$$

Neural word embedding – loss



Neural word embedding – all together

- Probability distribution of output words:

$$P(c|v) = \frac{\exp(\mathbf{e}_v \mathbf{u}_c)}{\sum_{\tilde{c} \in \mathbb{V}} \exp(\mathbf{e}_v \mathbf{u}_{\tilde{c}})}$$

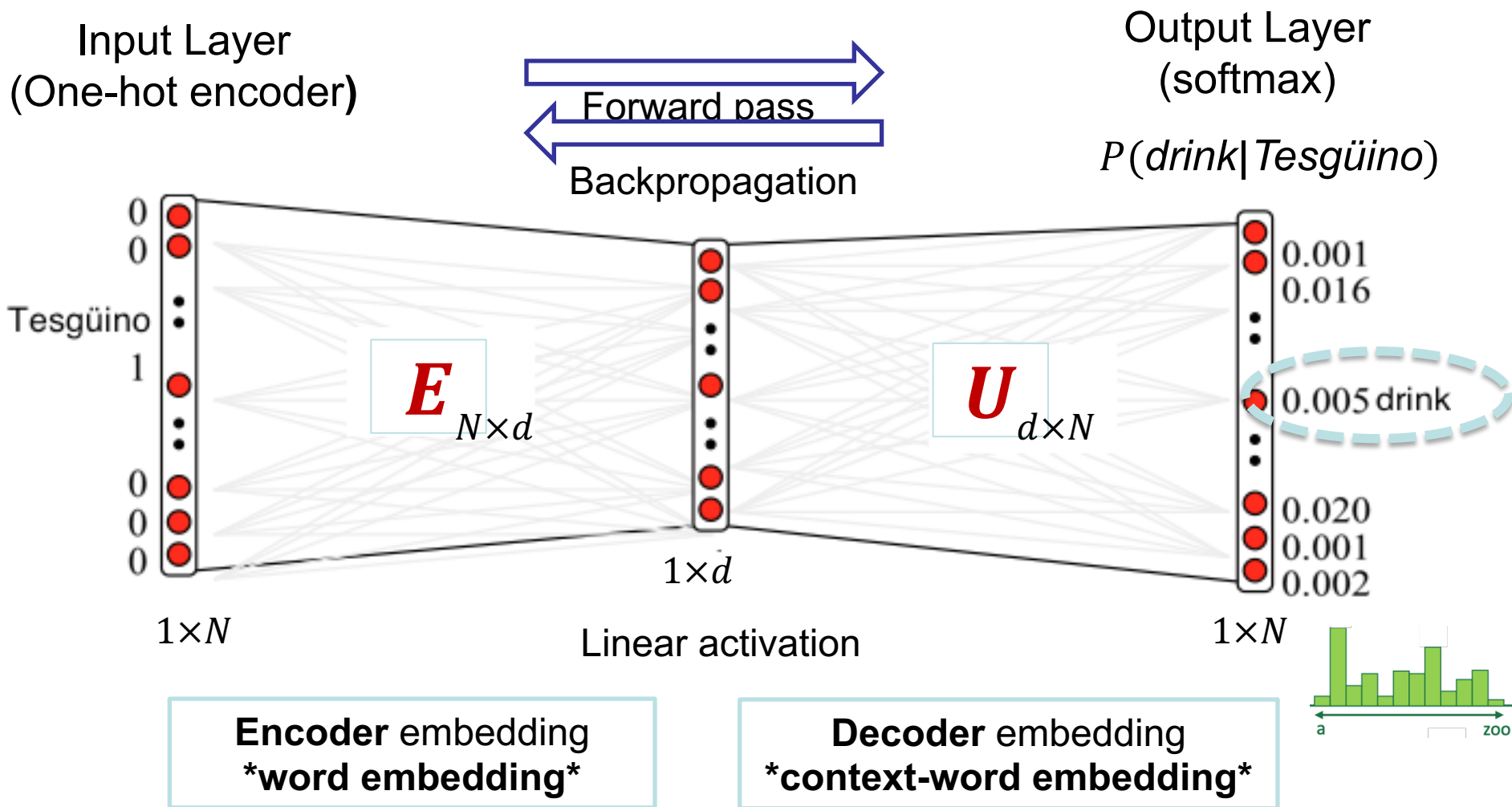
- In the example: $P(\text{drink}|\text{Tesgüino}) = \frac{\exp(\mathbf{e}_{\text{Tesgüino}} \mathbf{u}_{\text{drink}})}{\sum_{\tilde{c} \in \mathbb{V}} \exp(\mathbf{e}_{\text{Tesgüino}} \mathbf{u}_{\tilde{c}})}$

- Loss is the **NLL** over all training data:

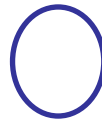
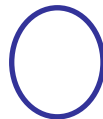
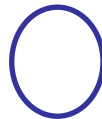
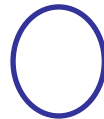
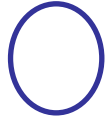
$$\mathcal{L} = -\mathbb{E}_{(v,c) \sim \mathcal{D}} \log P(c|v)$$

Neural word embedding – architecture (another view!)

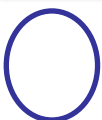
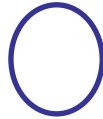
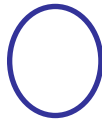
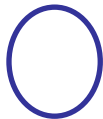
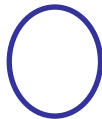
Training data: (*Tesgüino*, *drink*)



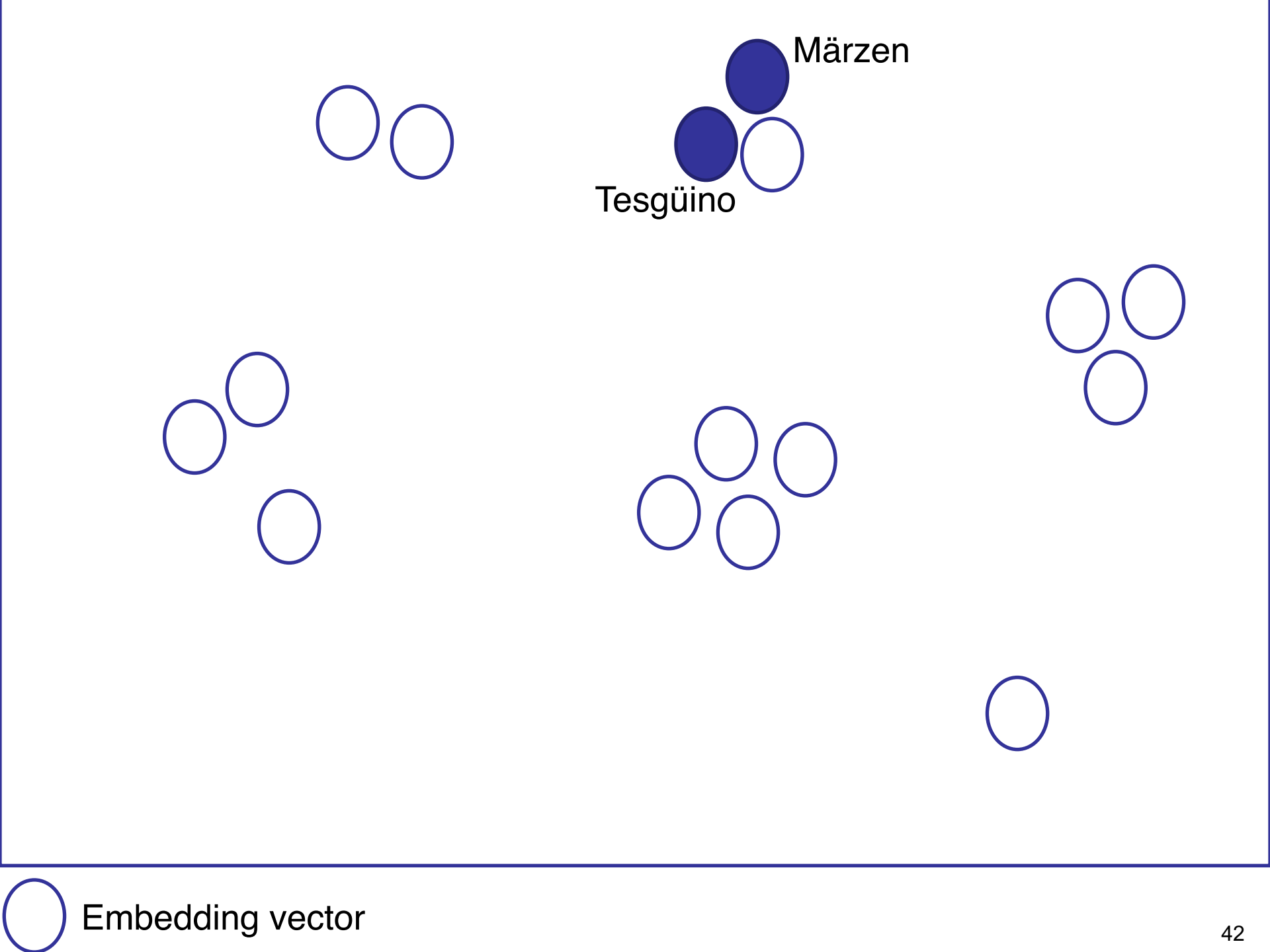
Märzen

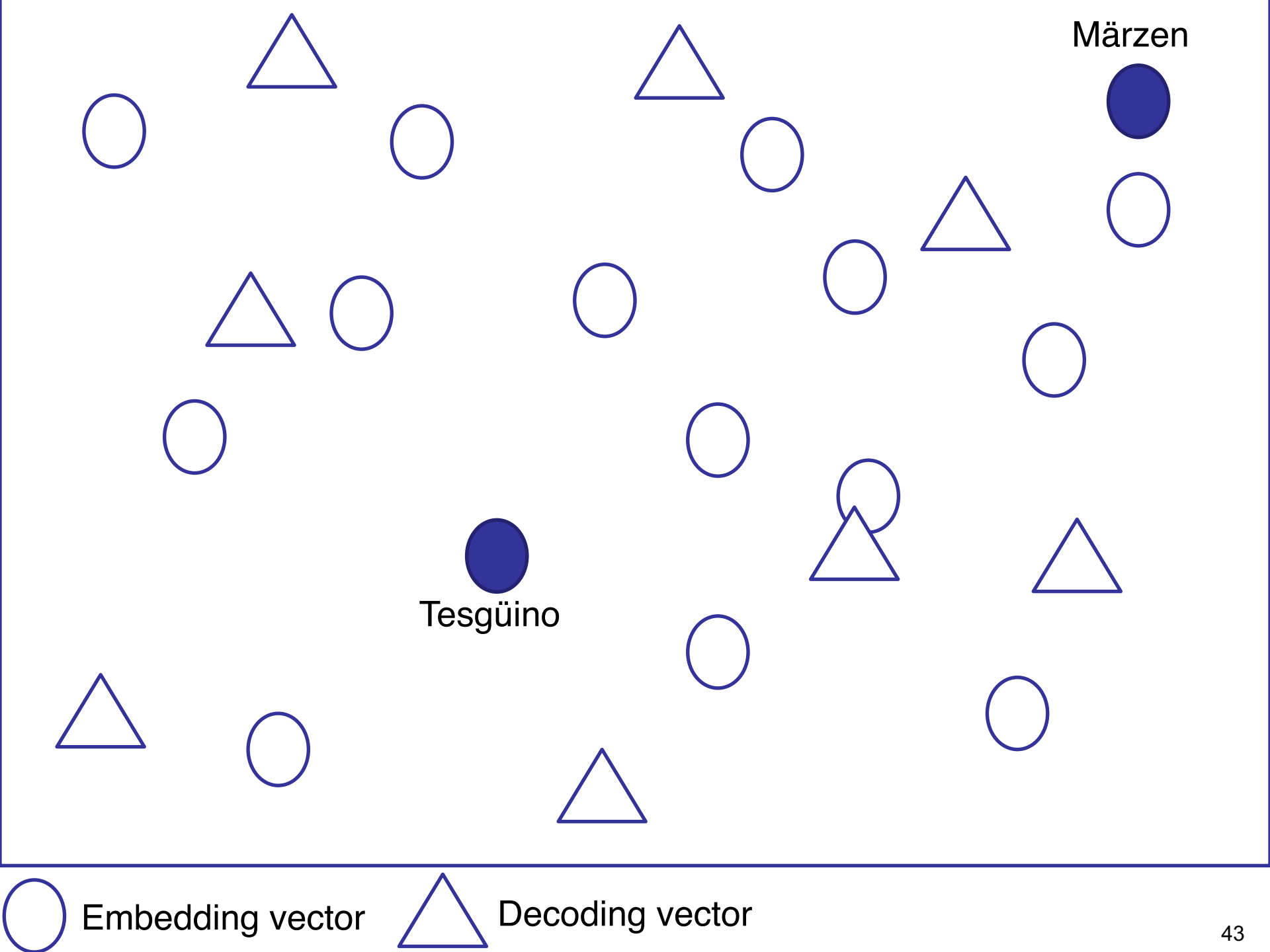


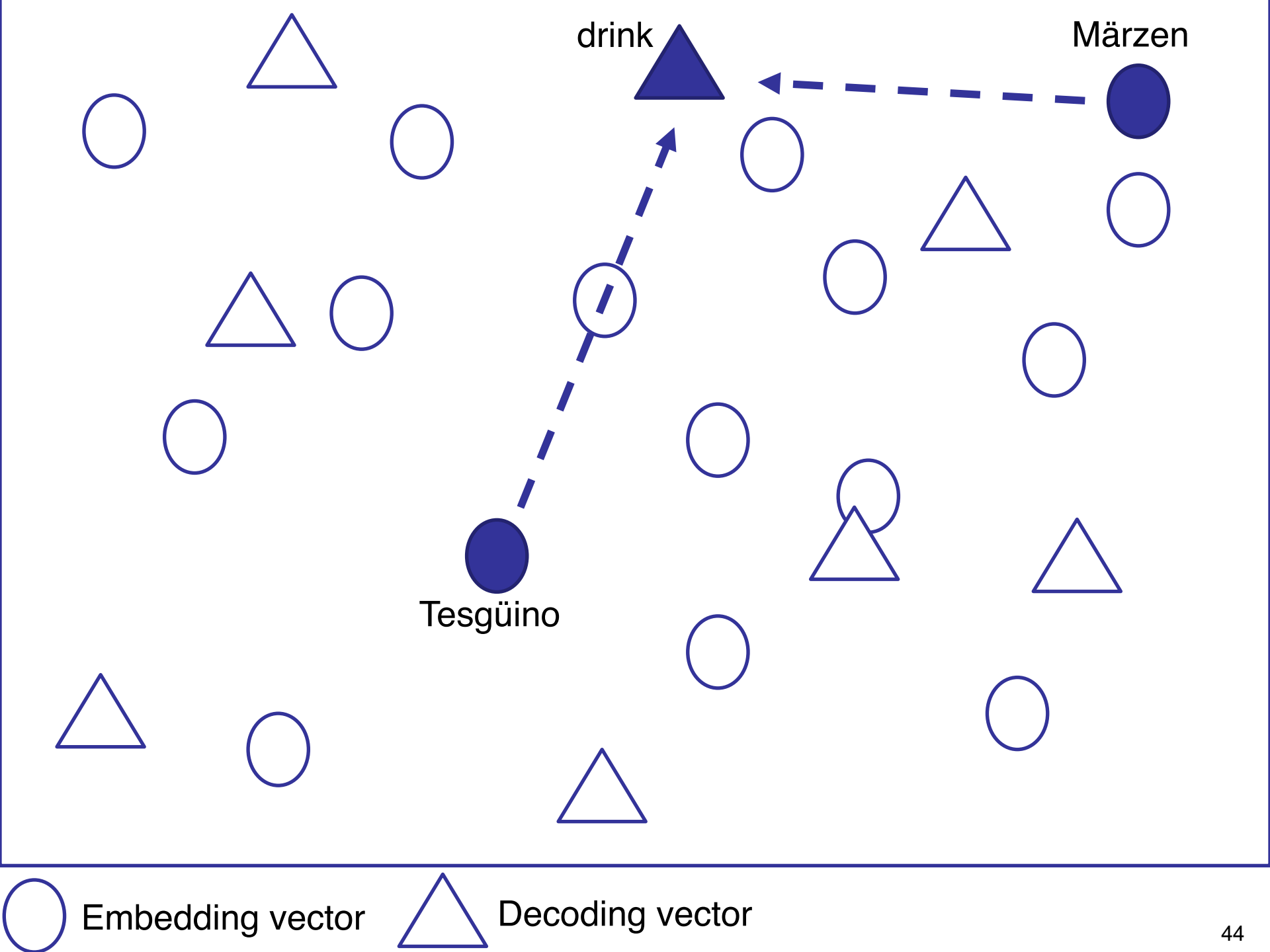
Tesgüino

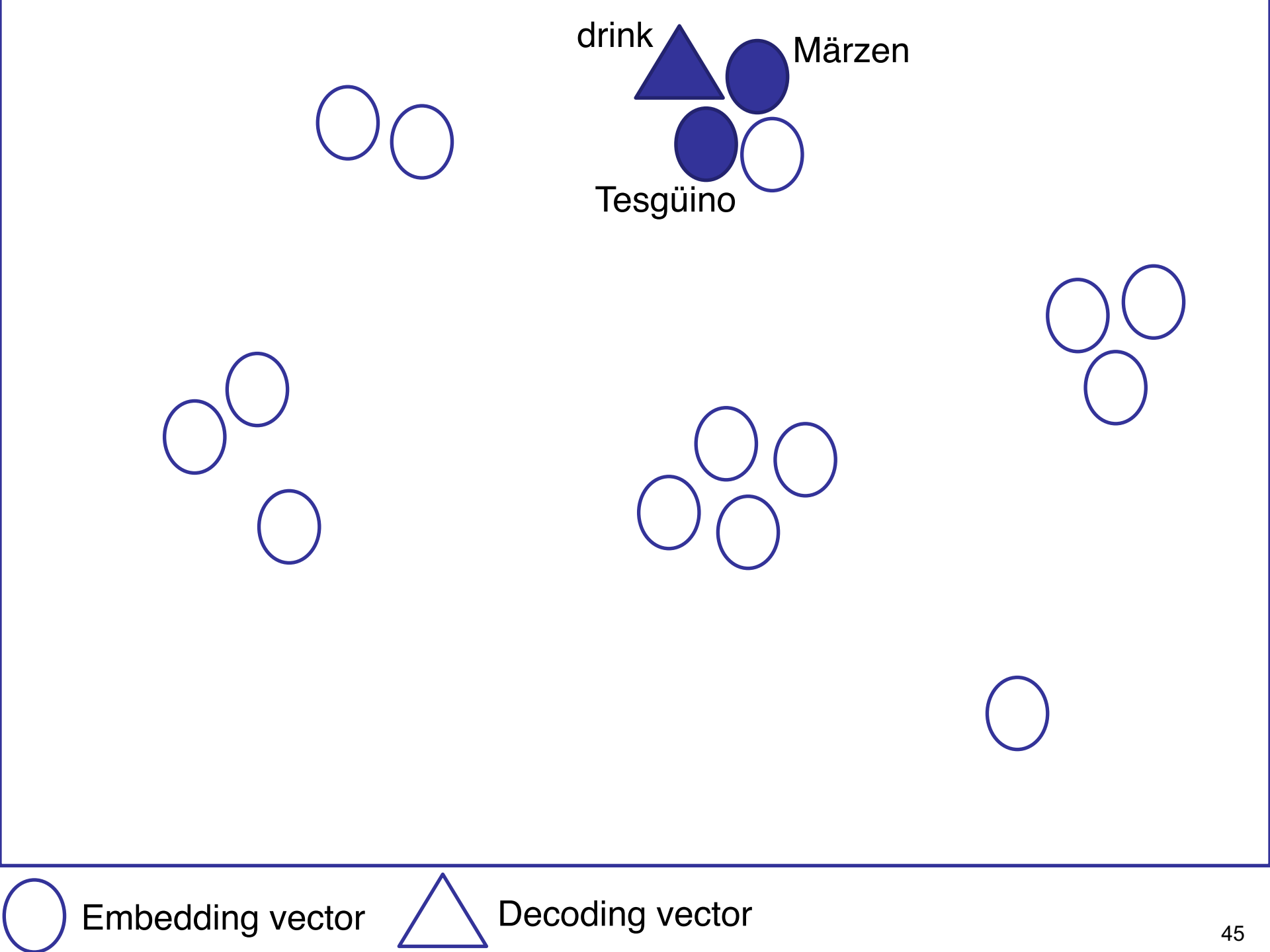


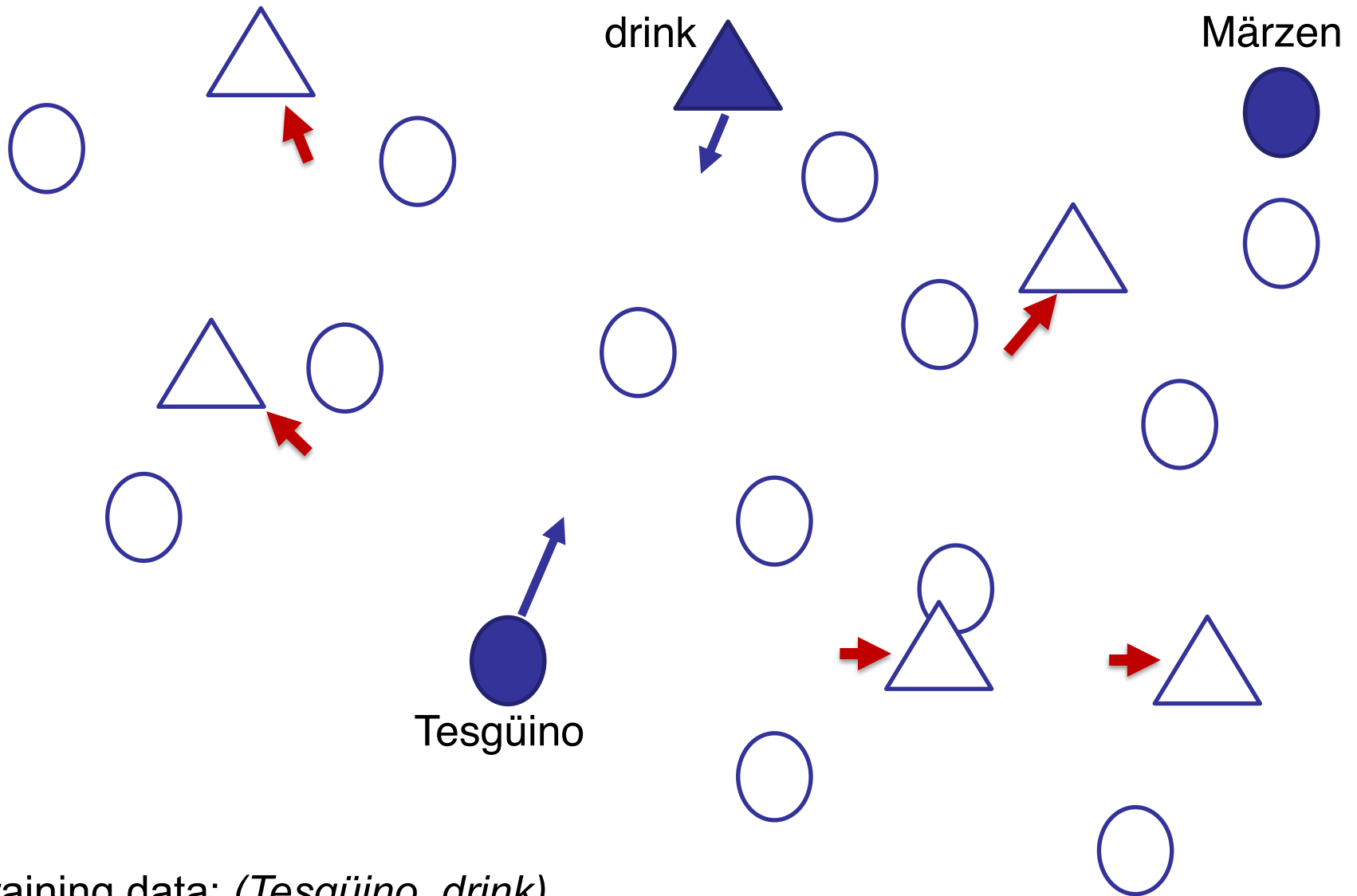
Embedding vector











- Training data: $(Tesgüino, drink)$
- Update vectors to maximize $P(drink|Tesgüino)$

Neural word embedding – NLL + softmax

$$P(c|v) = \frac{\exp(\mathbf{e}_v \mathbf{u}_c)}{\sum_{\tilde{c} \in \mathbb{V}} \exp(\mathbf{e}_v \mathbf{u}_{\tilde{c}})}$$

$$\mathcal{L} = -\mathbb{E}_{(v,c) \sim \mathcal{D}} \log P(c|v)$$

$$\mathcal{L} = -\mathbb{E}_{(v,c) \sim \mathcal{D}} \left[\log \frac{\exp(\mathbf{e}_v \mathbf{u}_c)}{\sum_{\tilde{c} \in \mathbb{V}} \exp(\mathbf{e}_v \mathbf{u}_{\tilde{c}})} \right]$$

$$\mathcal{L} = -\mathbb{E}_{(v,c) \sim \mathcal{D}} \left[\mathbf{e}_v \mathbf{u}_c - \log \sum_{\tilde{c} \in \mathbb{V}} \exp(\mathbf{e}_v \mathbf{u}_{\tilde{c}}) \right]$$



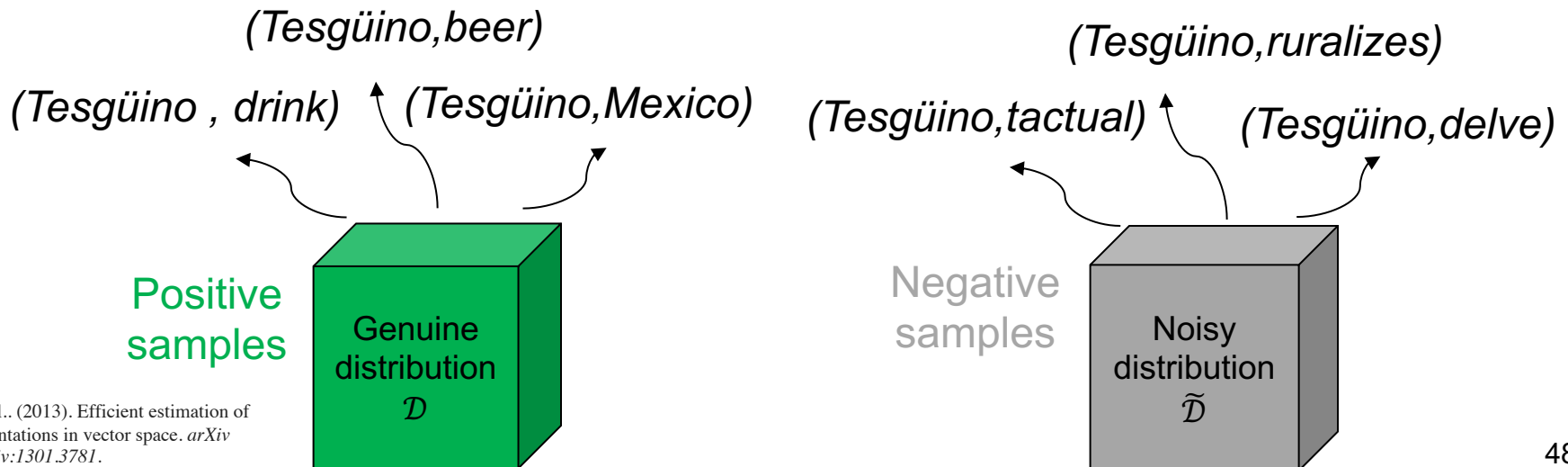
calculating this normalization term is expensive!

word2vec Skip-Gram with Negative Sampling

- Word2vec is an **efficient** and **effective** algorithm that proposes **Negative Sampling** method to define loss

Central idea of Negative Sampling:

- Consider two data distributions that generate [word, context-word] pairs:
 - The **genuine** one that comes from the training data $\rightarrow \mathcal{D}$
 - The **noisy** one that generates (almost) random pairs $\rightarrow \tilde{\mathcal{D}}$
- Aim: given a pair, the model should be able to distinguish whether it comes from the genuine distribution (binary classification)



Negative Sampling

- Instead of $P(c|v)$ as in neural word embedding, Negative Sampling calculates ...
- $P(y = 1|v, c) \rightarrow$ probability that the pair (v, c) comes from the **genuine** data distribution
- **Positive sample**: pair (v, c) that comes from the genuine data distribution \mathcal{D} (training data)
- **Negative sample**: pair (v, \tilde{c}) drawn from the noisy distribution $\tilde{\mathcal{D}}$
 - $\tilde{\mathcal{D}}$ generates pairs almost randomly!
 - Why random pairs can be considered as negative samples?
 - $\tilde{\mathcal{D}}$ in word2vec is a *smoothed* unigram distribution of words in corpus
 - Implementation details: In word2vec, $\tilde{\mathcal{D}}$ is smoothed by raising unigram counts to the power of $\alpha = 0.75$

Negative Sampling

- Negative Sampling loss aims to
 - increase the probability of **positive samples**, $P(y = 1|v, c)$ and ...
 - decrease the probabilities of k **negative samples**, $P(y = 1|v, \tilde{c})$
 - k is usually between 2 to 20
- $P(y = 1|v, c)$ is defined using sigmoid σ :

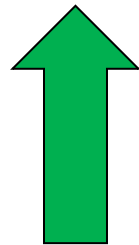
$$P(y = 1|v, c) = \sigma(\mathbf{e}_v \mathbf{u}_c)$$

word2vec – loss

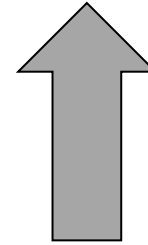
- Negative Sampling loss
 - **increases** the probability for the **positive sample** (v, c)
 - **decreases** the probability for the k **negative samples** (v, \tilde{c})
- Loss function:

$$\mathcal{L} = -\mathbb{E}_{(v,c) \sim \mathcal{D}} \left[\log P(y = 1 | v, c) - \sum_{\substack{\tilde{c} \sim \tilde{\mathcal{D}} \\ k \text{ times}}} \log P(y = 1 | v, \tilde{c}) \right]$$

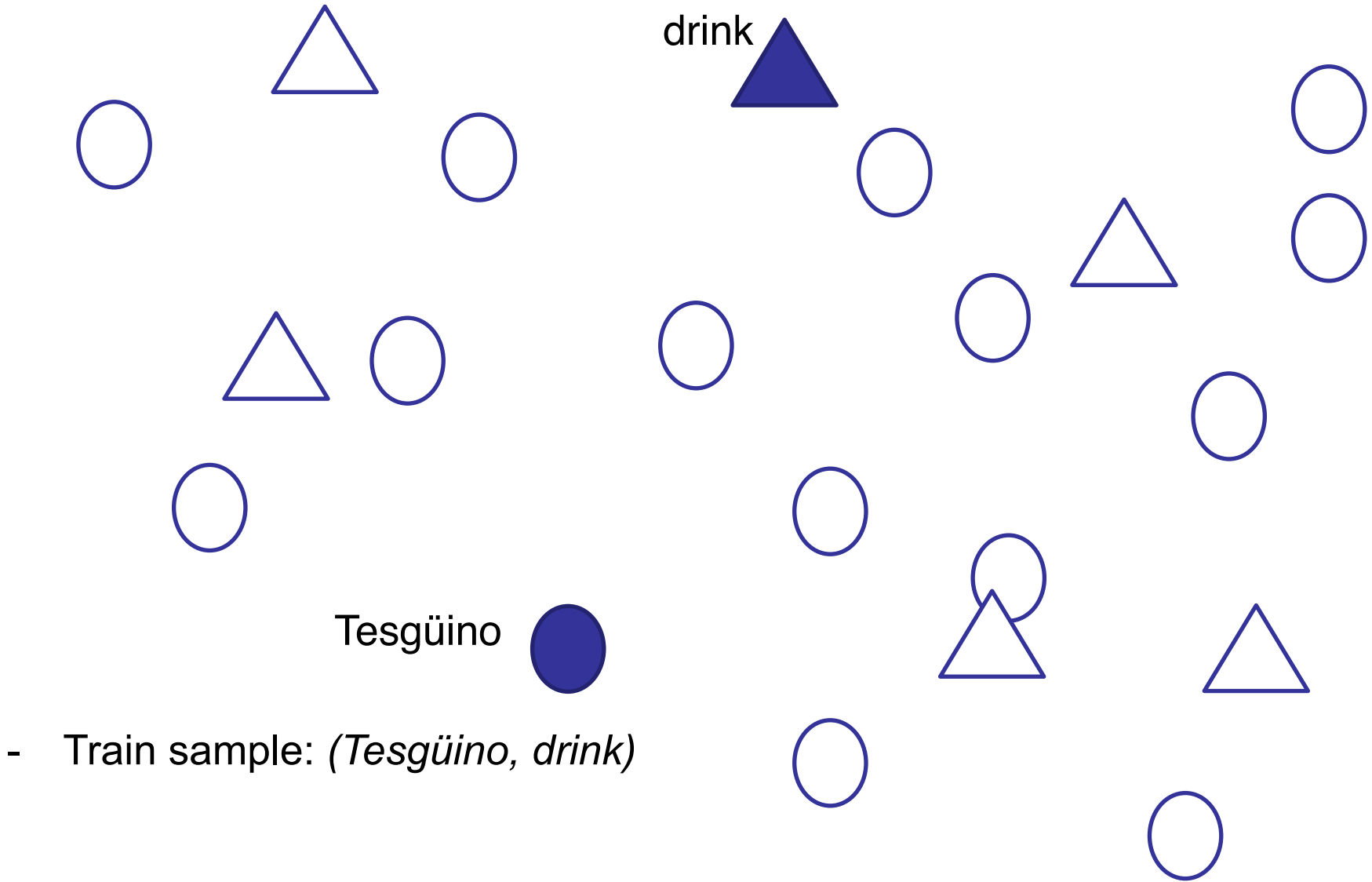
$$\mathcal{L} = -\mathbb{E}_{(v,c) \sim \mathcal{D}} \left[\log \sigma(\mathbf{e}_v \mathbf{u}_c) - \sum_{\substack{\tilde{c} \sim \tilde{\mathcal{D}} \\ k \text{ times}}} \log \sigma(\mathbf{e}_v \mathbf{u}_{\tilde{c}}) \right]$$



positive samples



negative samples



- Train sample: (*Tesgüino*, *drink*)

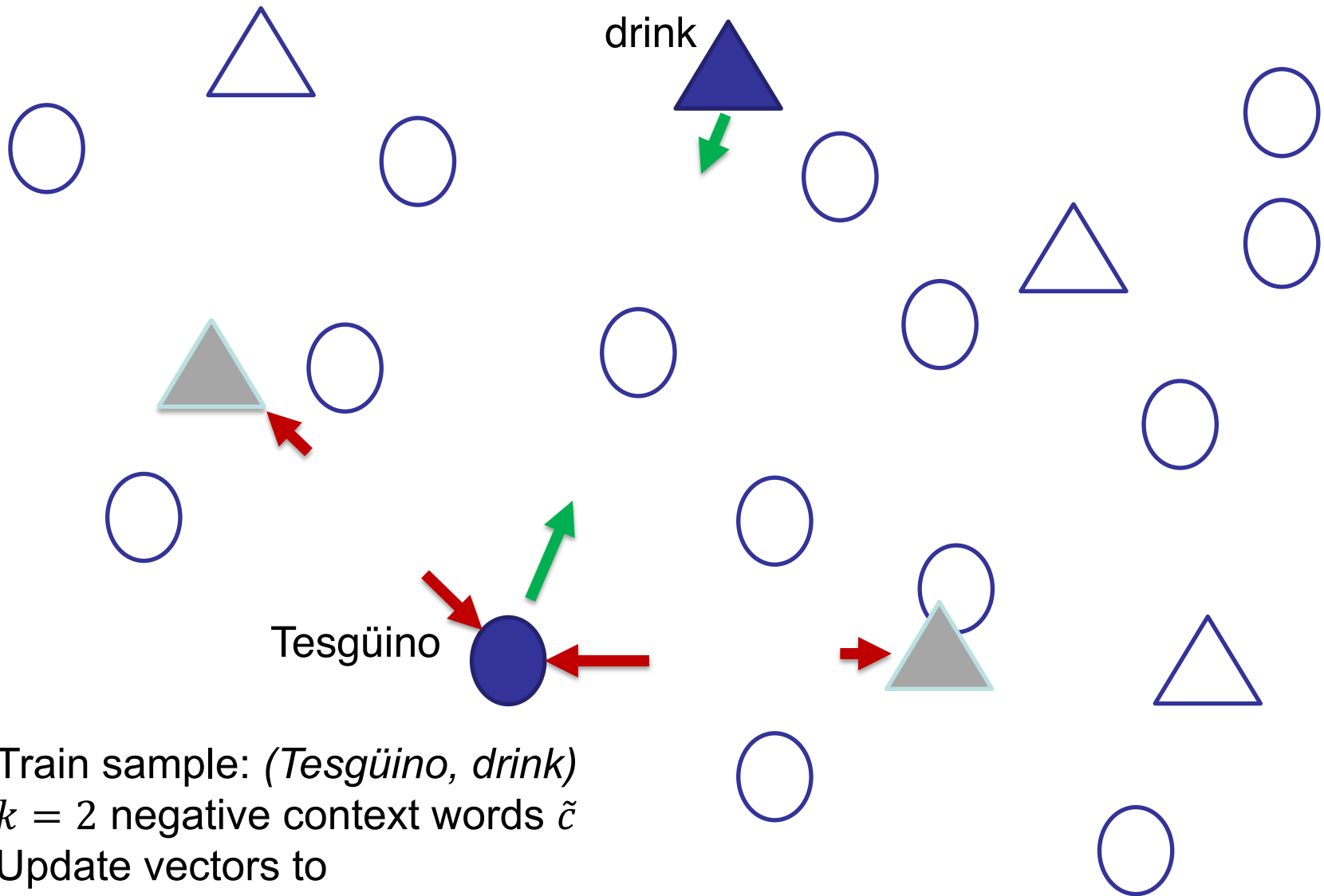


drink

Tesgüino

- Train sample: (*Tesgüino*, *drink*)
- $k = 2$ negative context words

 Embedding vector  Decoding vector



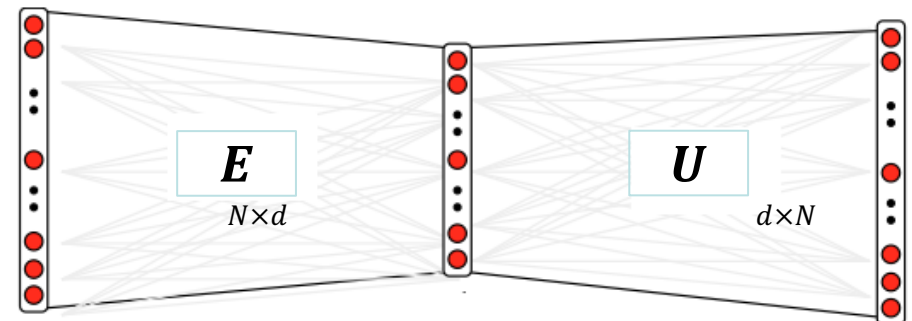
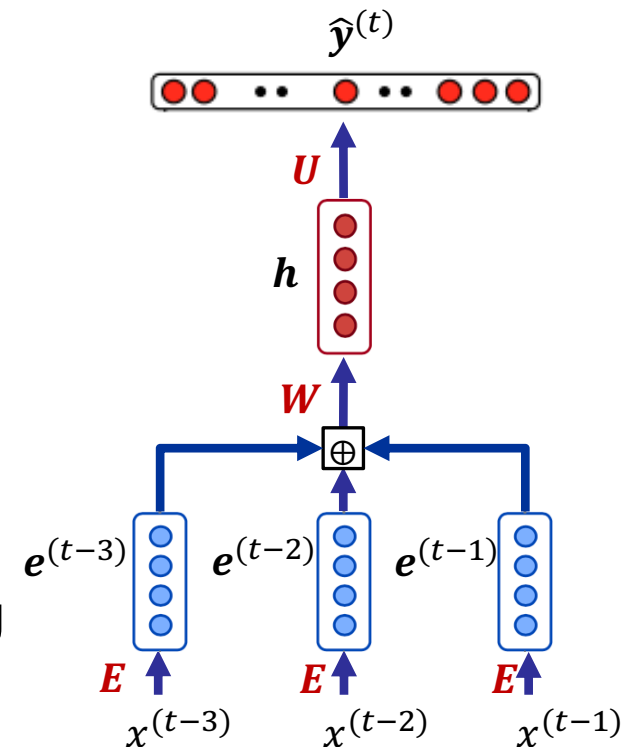
- Train sample: $(\text{Tesgüino}, \text{drink})$
- $k = 2$ negative context words \tilde{c}
- Update vectors to
 - Increase $P(y = 1 | \text{Tesgüino}, \text{drink})$
 - Decrease $P(y = 1 | \text{Tesgüino}, \tilde{c})$

Negative Sampling – final words!

- Negative Sampling turns the problem from multi-class classification to binary classification
- Softmax is a good choice for training Language Models, namely to estimate $P(v|\text{context})$
- Negative Sampling is shown to be effective for training good embeddings
- Negative Sampling is a biased approximation of softmax
 - **Noisy Contrastive Estimation** (the parent of Negative Sampling) is an unbiased approximation of softmax

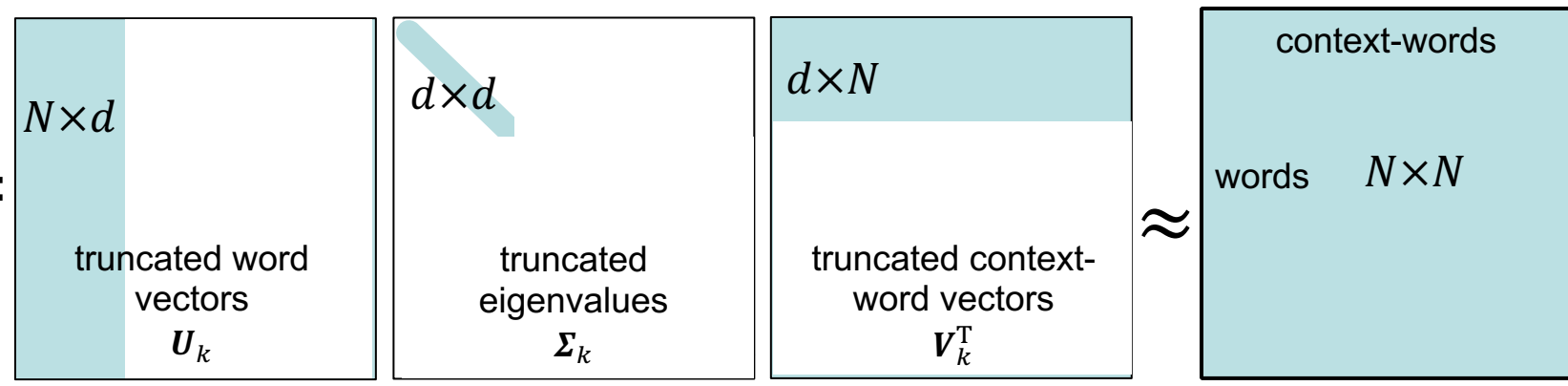
Summary

- n -gram Language Models count co-occurrences
- Neural n -gram Language Models predict co-occurrence probabilities
- word2vec creates word embedding by ...
 - following a skip-gram Language Modeling objective and ...
 - exploiting Negative Sampling loss

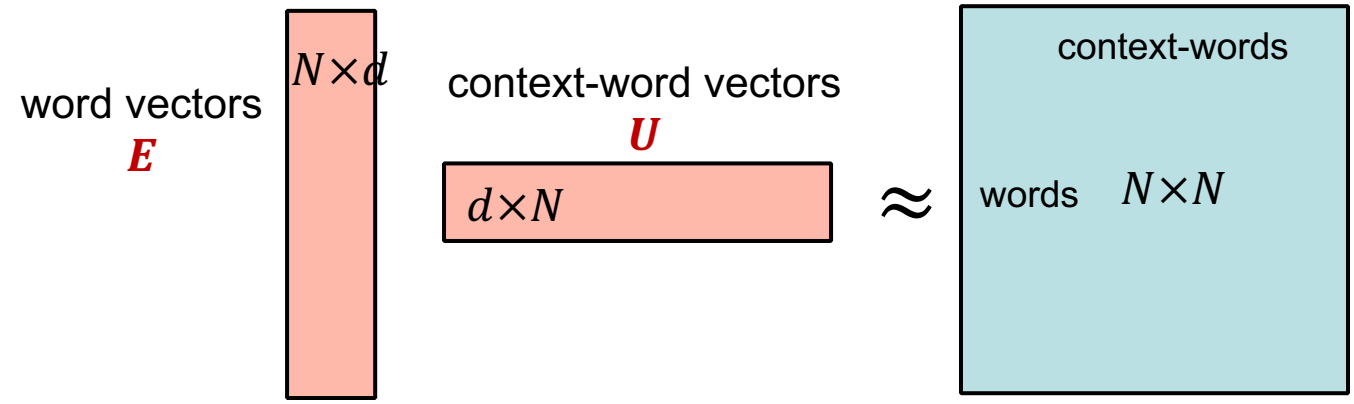


Word embedding models – in one frame!

PPMI+SVD:



GloVe:



word2vec:

