

# Natural Language Processing with Deep Learning

## Word Embeddings



Navid Rekab-Saz

[navid.rekabsaz@jku.at](mailto:navid.rekabsaz@jku.at)

# Agenda

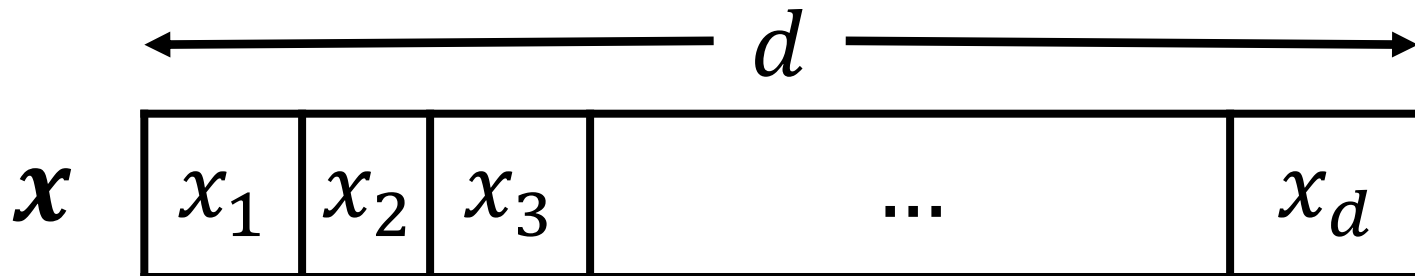
- Introduction
- Count-based word representation
- Prediction-based word embedding

# Agenda

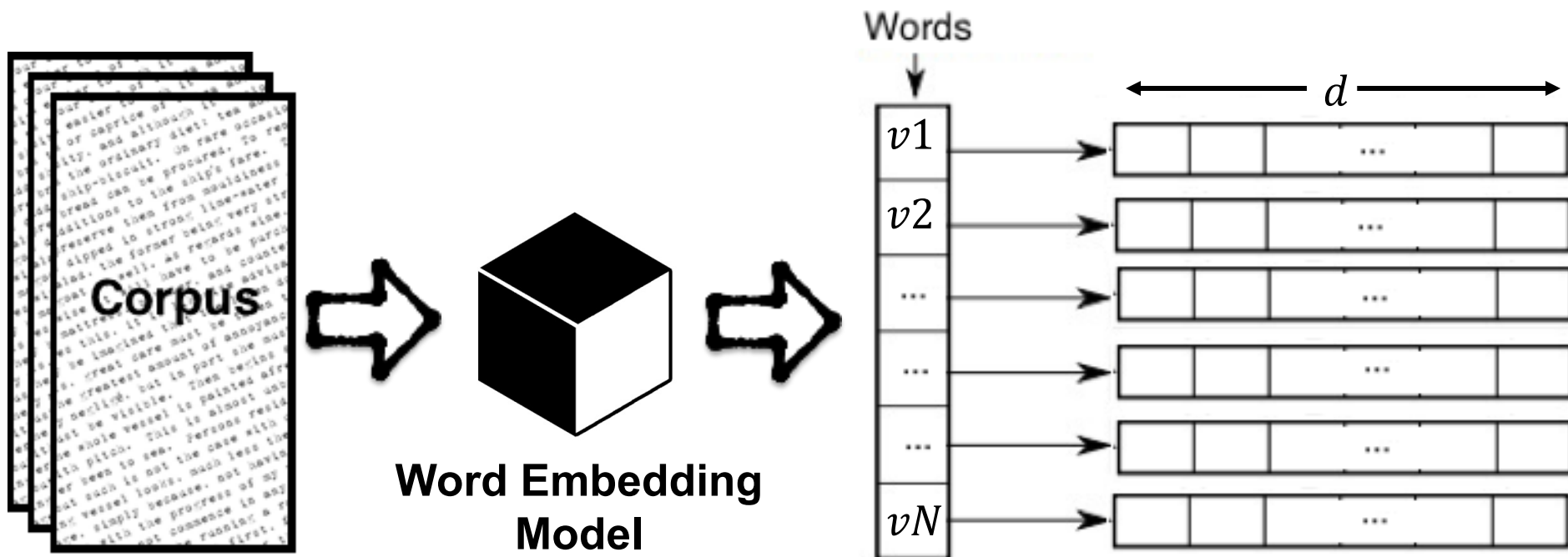
- **Introduction**
- Count-based word representation
- Prediction-based word embedding

# Distributional Representation

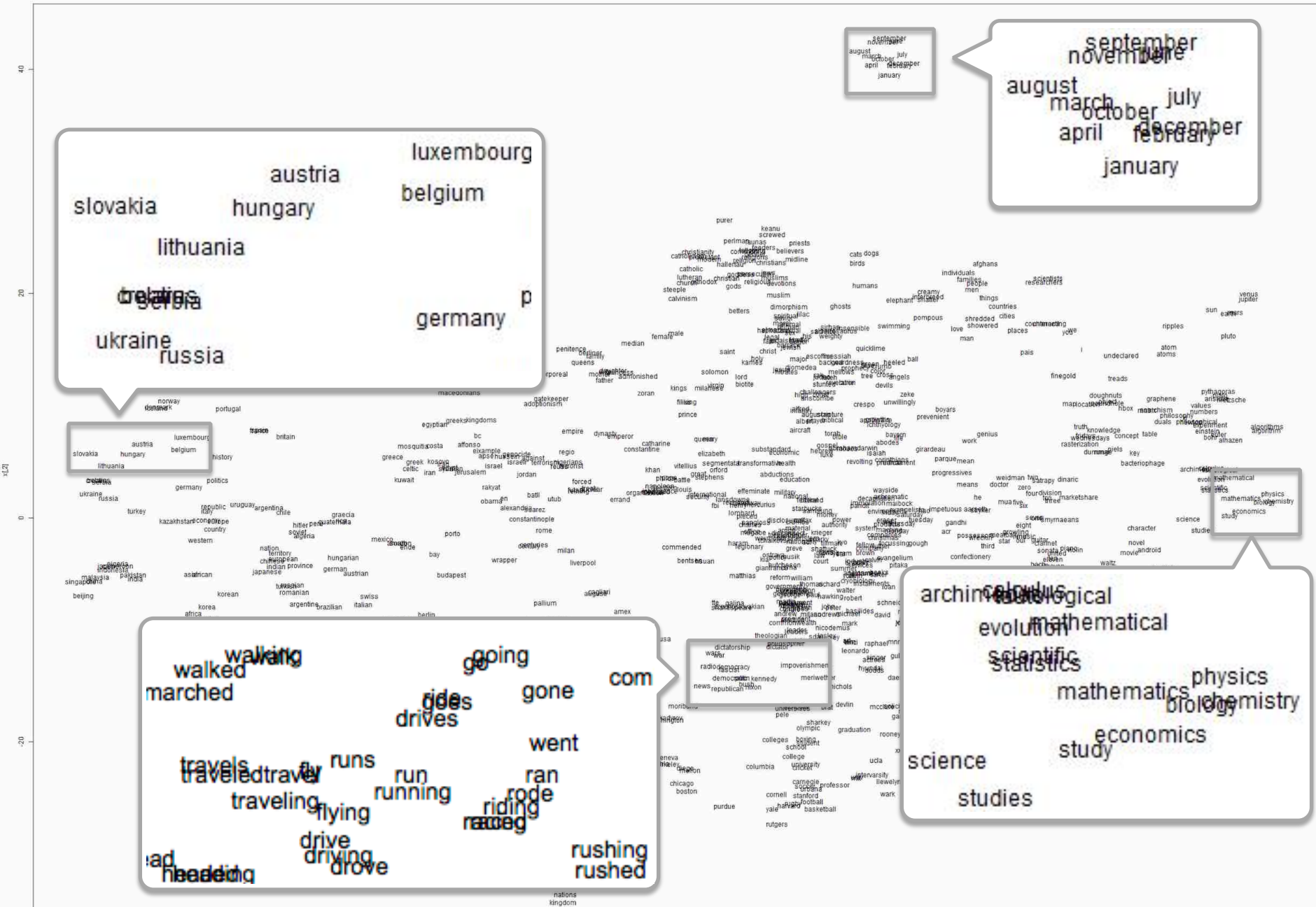
- An entity is represented with a **vector of  $d$  dimensions**
- **Distributed Representations**
  - Each dimension (units) is a **feature** of the entity
  - Units in a layer are not mutually exclusive
  - Two units can be “active” at the same time



# Word Embedding Model



When vector representations are dense, they are often called **embedding** e.g. word embedding



Word embeddings projected to a two-dimensional space

# Word Embeddings – Nearest neighbors

frog

---

frogs

---

toad

---

litoria

---

leptodactylidae

---

rana

---



book

---

books

---

foreword

---

author

---

published

---

preface

---

asthma

---

bronchitis

---

allergy

---

allergies

---

arthritis

---

diabetes

---

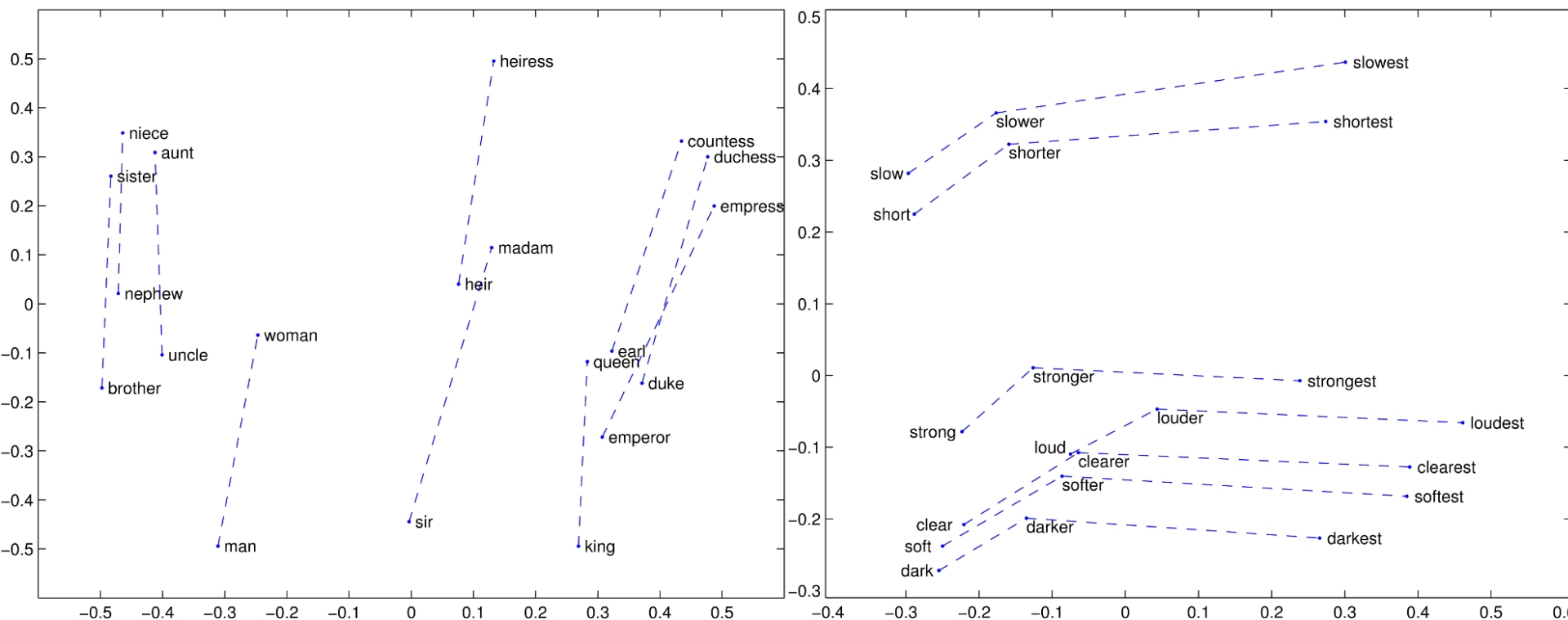
# Word Embeddings – Linear substructures

## ■ Analogy task:

- *man to woman is like king to ? (queen)*

$$\mathbf{x}_{\text{woman}} - \mathbf{x}_{\text{man}} + \mathbf{x}_{\text{king}} = \mathbf{x}^*$$

$$\mathbf{x}^* \approx \mathbf{x}_{\text{queen}}$$





# Agenda

- Introduction
- **Count-based word representation**
- Prediction-based word embedding

# Intuition for Computational Semantics



“You shall know a word  
by the company it  
keeps!”

*J. R. Firth, A synopsis of  
linguistic theory 1930–1955  
(1957)*

drink  
sacred  
alcoholic  
**Tesgüino**  
beverage  
out of corn  
fermented  
Mexico  
bottle

fermentation

bottle

grain

medieval

brew

**Ale**

pale

bar

drink

alcoholic



# Tesgüino



# Ale



**Algorithmic intuition:**

Two words are **related** when they have common **context words**

## Word-Document Matrix – recap

- $\mathbb{D}$  is a set of documents (plays of Shakespeare)  
 $\mathbb{D} = [d1, d2, \dots, dM]$
- $\mathbb{V}$  is the set of words (vocabularies) in dictionary  
 $\mathbb{V} = [v1, v2, \dots, vN]$
- Words as rows and documents as columns
- Values: term count  $tc_{v,d}$
- Matrix size  $N \times M$

	<i>d1</i>	<i>d2</i>	<i>d3</i>	<i>d4</i>
	As You Like It	Twelfth Night	Julius Caesar	Henry V
battle	1	1	8	15
soldier	2	2	12	36
fool	37	58	1	5
clown	6	117	0	0
...	...	...	...	...

# Cosine

- Cosine is the normalized dot product of two vectors
  - Its result is between -1 and +1

$$\cos(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x}}{\|\mathbf{x}\|_2} \cdot \frac{\mathbf{y}}{\|\mathbf{y}\|_2} = \frac{\sum_i x_i y_i}{\sqrt{\sum_i x_i^2} \sqrt{\sum_i y_i^2}}$$

- $\mathbf{x} = [1 \quad 1 \quad 0] \quad \mathbf{y} = [4 \quad 5 \quad 6]$

$$\cos(\mathbf{x}, \mathbf{y}) = \frac{1 * 4 + 1 * 5 + 0 * 6}{\sqrt{1^2 + 1^2 + 0^2} \sqrt{4^2 + 5^2 + 6^2}} = \frac{9}{\sim 12.4}$$

# Word-Document Matrix

	<i>d1</i>	<i>d2</i>	<i>d3</i>	<i>d4</i>
	As You Like It	Twelfth Night	Julius Caesar	Henry V
battle	1	1	8	15
soldier	2	2	12	36
fool	37	58	1	5
clown	6	117	0	0
...	...	...	...	...

- **Similarity** between two words:

$$\text{similarity}(\text{soldier}, \text{clown}) = \cos(\mathbf{x}_{\text{soldier}}, \mathbf{x}_{\text{clown}})$$



# Context

- Context can be
  - Document
  - Paragraph, tweet
  - **Window of (2-10) context words on each side of the word**
- Word-Context matrix
  - Every word as a unit (dimension)  
 $\mathbb{C} = [c_1, c_2, \dots, c_L]$
  - Matrix size:  $N \times L$
  - Usually  $\mathbb{C} = \mathbb{V}$  and therefore  $L = N$

# Word-Context Matrix

- Window context of 7 words

sugar, a sliced lemon, a tablespoonful of **apricot** preserve or jam, a pinch each of,  
their enjoyment. Cautiously she sampled her first **pineapple** and another fruit whose taste she likened  
well suited to programming on the digital **computer.** In finding the optimal R-stage policy from  
for the purpose of gathering data and **information** necessary for the study authorized in the

	<i>c1</i>	<i>c2</i>	<i>c3</i>	<i>c4</i>	<i>c5</i>	<i>c6</i>
	aardvark	computer	data	pinch	result	sugar
<i>v1</i> apricot	0	0	0	1	0	1
<i>v2</i> pineapple	0	0	0	1	0	1
<i>v3</i> digital	0	2	1	0	1	0
<i>v4</i> information	0	1	6	0	4	0

# Word-to-Word Relations

	<i>c1</i>	<i>c2</i>	<i>c3</i>	<i>c4</i>	<i>c5</i>	<i>c6</i>
	aardvark	computer	data	pinch	result	sugar
<i>v1</i> apricot	0	0	0	1	0	1
<i>v2</i> pineapple	0	0	0	1	0	1
<i>v3</i> digital	0	2	1	0	1	0
<i>v4</i> information	0	1	6	0	4	0

- First-order co-occurrence relation
  - Each cell of the word-context matrix
  - Words that appear in the proximity of each other
  - Like **drink** to **beer**, and **drink** and **wine**
- Second-order similarity relation
  - Cosine similarity between the representation vectors
  - Words that appear in similar contexts
  - Like **beer** to **wine**, **tesgüino** to **ale**, and **frog** to **toad**

# Point Mutual Information

- Problem with raw counting methods
  - Biased towards high frequent words (“and”, “the”) although they don’t contain much of information
- Point Mutual Information (PMI)
  - Rooted in information theory
  - A better measure for the **first-order relation** in word-context matrix in order to reflect the **informativeness** of co-occurrences
  - Joint probability of two events (random variables) divided by their marginal probabilities

$$\text{PMI}(X, Y) = \log \frac{p(X, Y)}{p(X)p(Y)}$$

## Point Mutual Information

$$\text{PMI}(v, c) = \log \frac{p(v, c)}{p(v)p(c)}$$

$$p(v, c) = \frac{\# \langle v, c \rangle}{S}$$

$$p(v) = \frac{\sum_{j=1}^L \# \langle v, c_j \rangle}{S} \quad p(c) = \frac{\sum_{i=1}^N \# \langle v_i, c \rangle}{S}$$

$$S = \sum_{i=1}^N \sum_{j=1}^L \# \langle v_i, c_j \rangle$$

- Positive Point Mutual Information (**PPMI**)

$$\text{PPMI}(t, c) = \max(\text{PMI}, 0)$$

# Point Mutual Information

	<i>c1</i>	<i>c2</i>	<i>c3</i>	<i>c4</i>	<i>c5</i>	<i>c6</i>
	aardvark	computer	data	pinch	result	sugar
<i>v1</i> apricot	0	0	0	1	0	1
<i>v2</i> pineapple	0	0	0	1	0	1
<i>v3</i> digital	0	2	1	0	1	0
<i>v4</i> information	0	1	6	0	4	0

$$P(v = \text{information}, c = \text{data}) = 6/19 = .32$$

$$P(v = \text{information}) = 11/19 = .58$$

$$P(c = \text{data}) = 7/19 = .37$$

$$\text{PPMI}(v = \text{information}, c = \text{data}) = \max(0, \log \frac{.32}{.58 * .37}) = .39$$

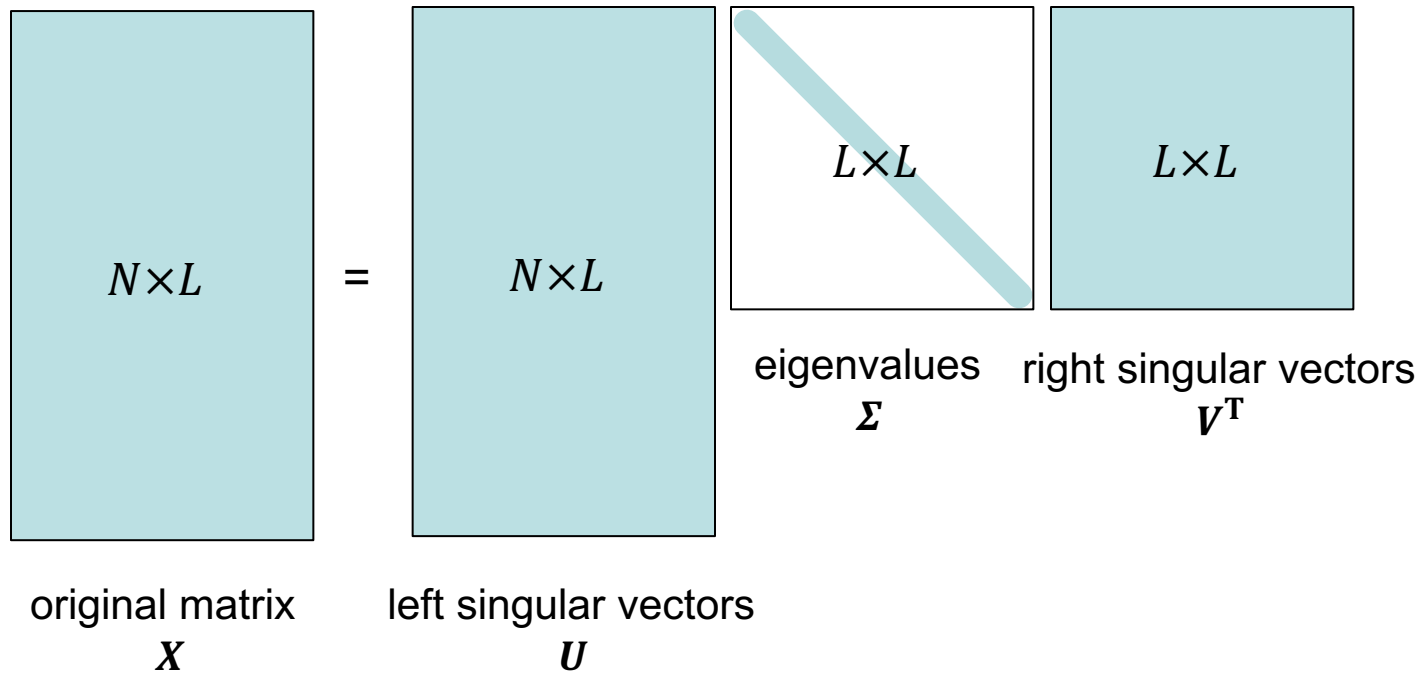
# Singular Value Decomposition - Recap

- An  $N \times L$  matrix  $X$  can be factorized to three matrices:

$$X = U\Sigma V^T$$

- $U$  (left singular vectors) is an  $N \times L$  unitary matrix
- $\Sigma$  is an  $L \times L$  diagonal matrix, diagonal entries
  - are eigenvalues,
  - show the importance of corresponding  $L$  dimensions in  $X$
  - are all positive and sorted from large to small values
- $V^T$  (right singular vectors) is an  $L \times L$  unitary matrix

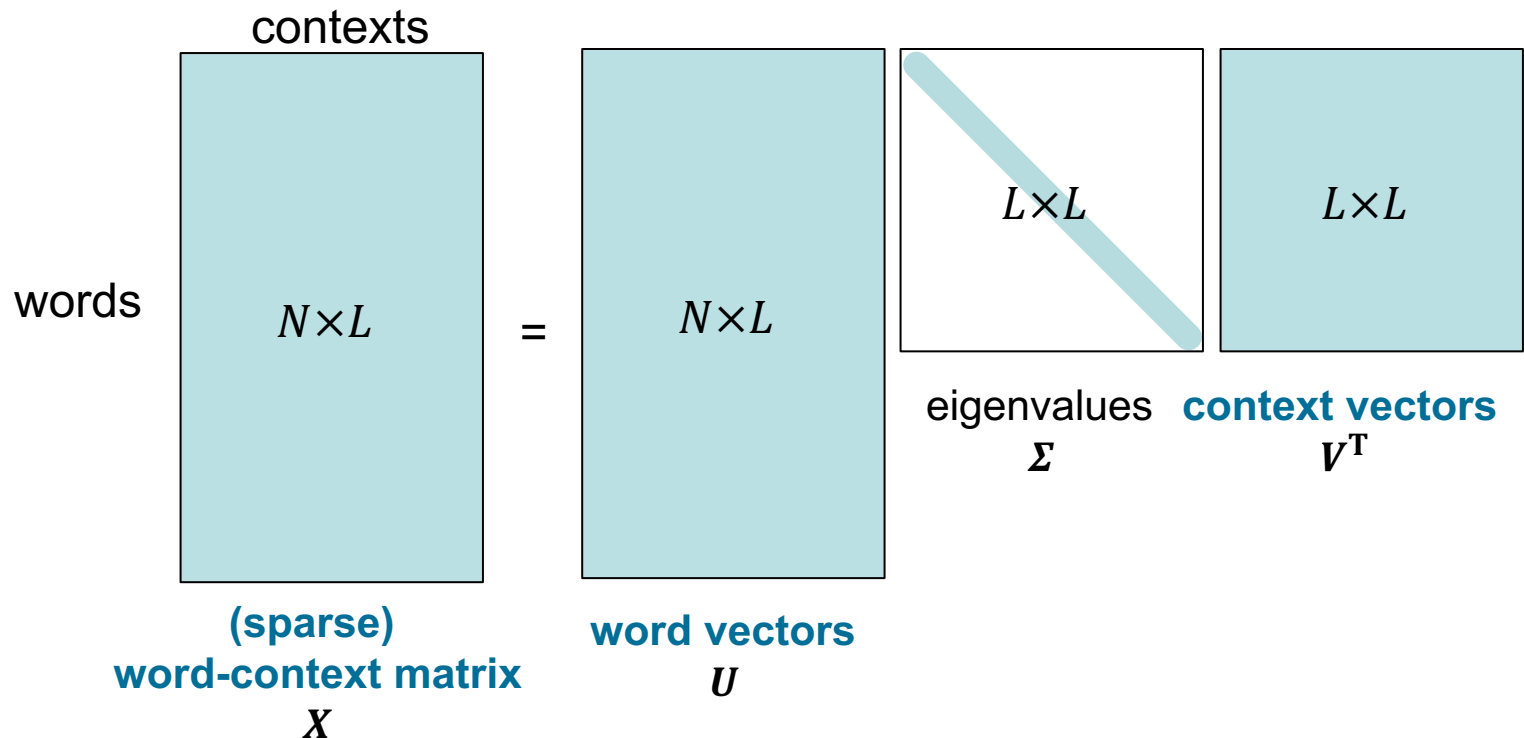
# Singular Value Decomposition – Recap





# Applying SVD to Word-Context Matrix

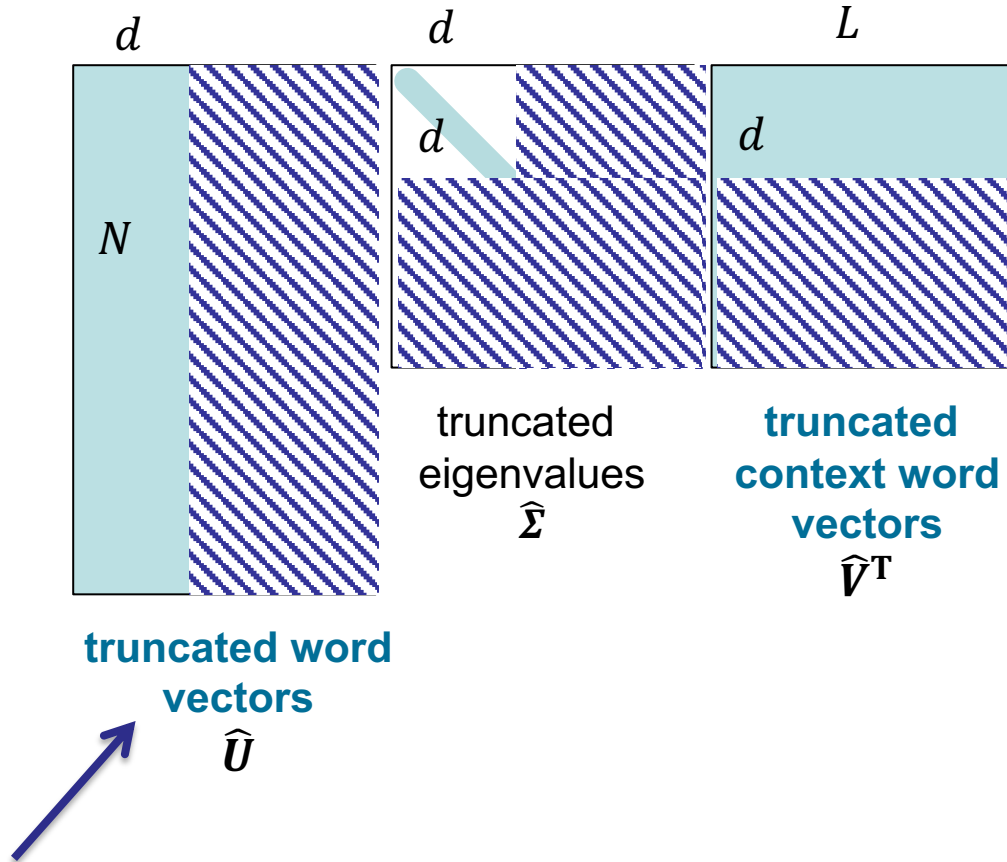
- Step 1: create a sparse PPMI matrix of the size  $N \times L$ ,
- Apply SVD



## Applying SVD to Term-Context Matrix

- Step 2: keep only top  $d$  eigenvalues in  $\Sigma$  and set the rest to zero
- Truncate the  $U$  and  $V^T$  matrices based on the changes in  $\Sigma$ , called  $\hat{U}$  and  $\hat{V}^T$  respectively

# Applying SVD to Term-Context Matrix



- $\hat{U}$  matrix is the dense low-dimensional word vectors



# Agenda

- Introduction
- Count-based word representation
- **Prediction-based word embedding**

# Word Embedding with Neural Networks

Recipe for creating (dense) word embedding with neural networks

- Design a neural network architecture!
- Loop over training data  $(v, c)$  for some epochs
  - Pass the word  $v$  as input and execute forward passing
  - Calculate the probability of observing context word  $c$  at output:

$$p(c|v)$$

- Optimize the network to maximize this likelihood probability

Details come next!

# Training Data $\mathcal{D}$

Window of size 2

Source Text

The quick brown fox jumps over the lazy dog. →

The quick brown fox jumps over the lazy dog. →

The quick brown fox jumps over the lazy dog. →

The quick brown fox jumps over the lazy dog. →

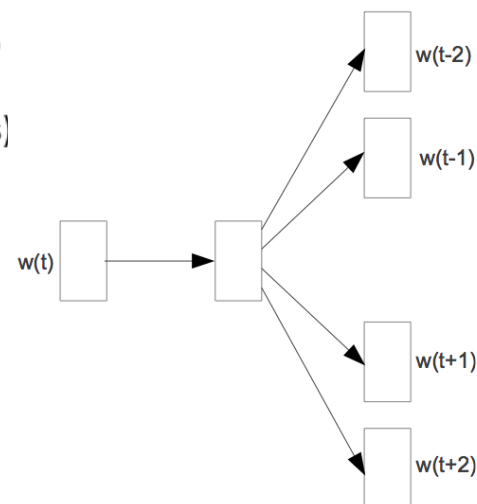
Training  
Samples

(the, quick)  
(the, brown)

(quick, the)  
(quick, brown)  
(quick, fox)

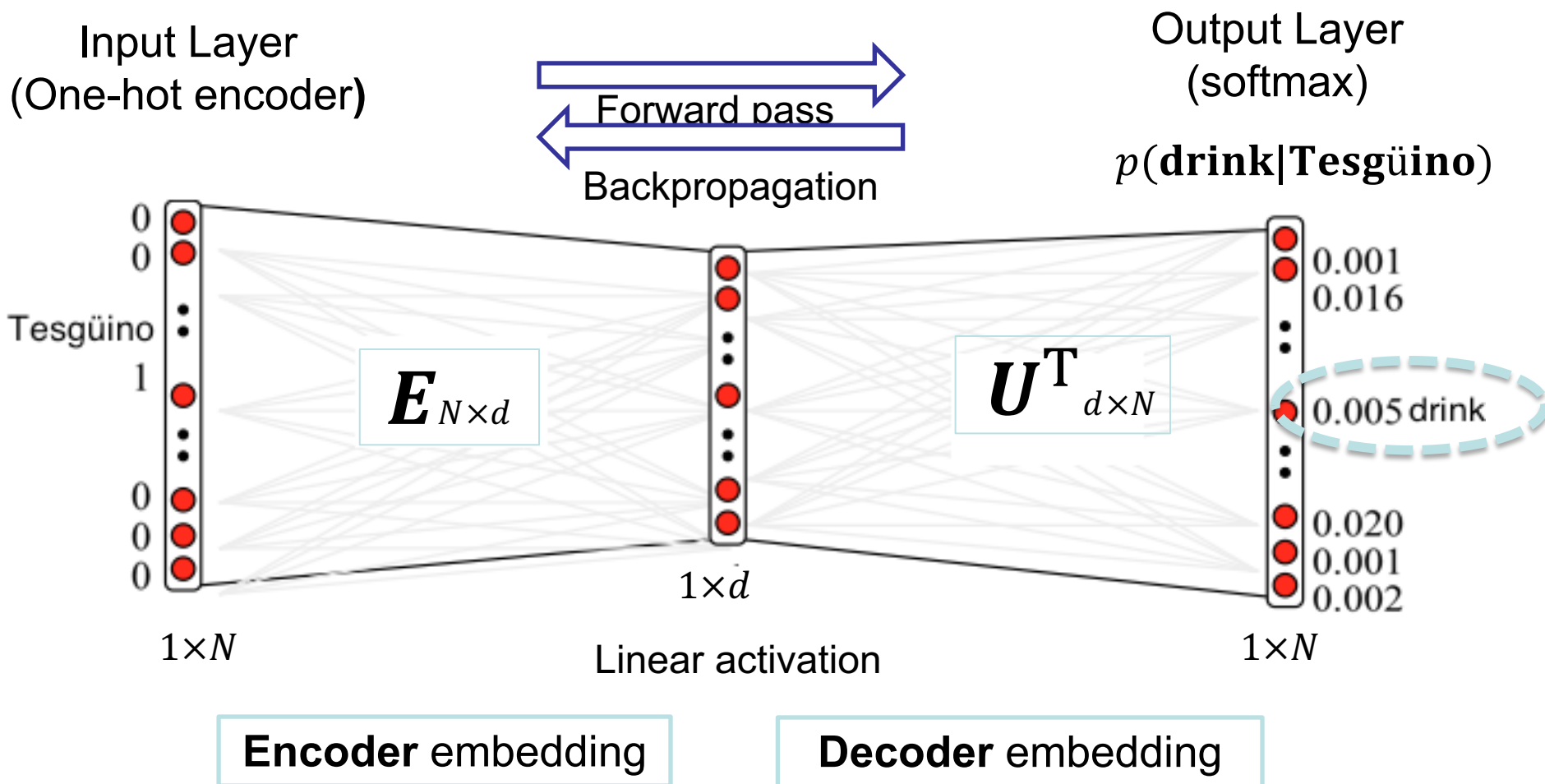
(brown, the)  
(brown, quick)  
(brown, fox)  
(brown, jumps)

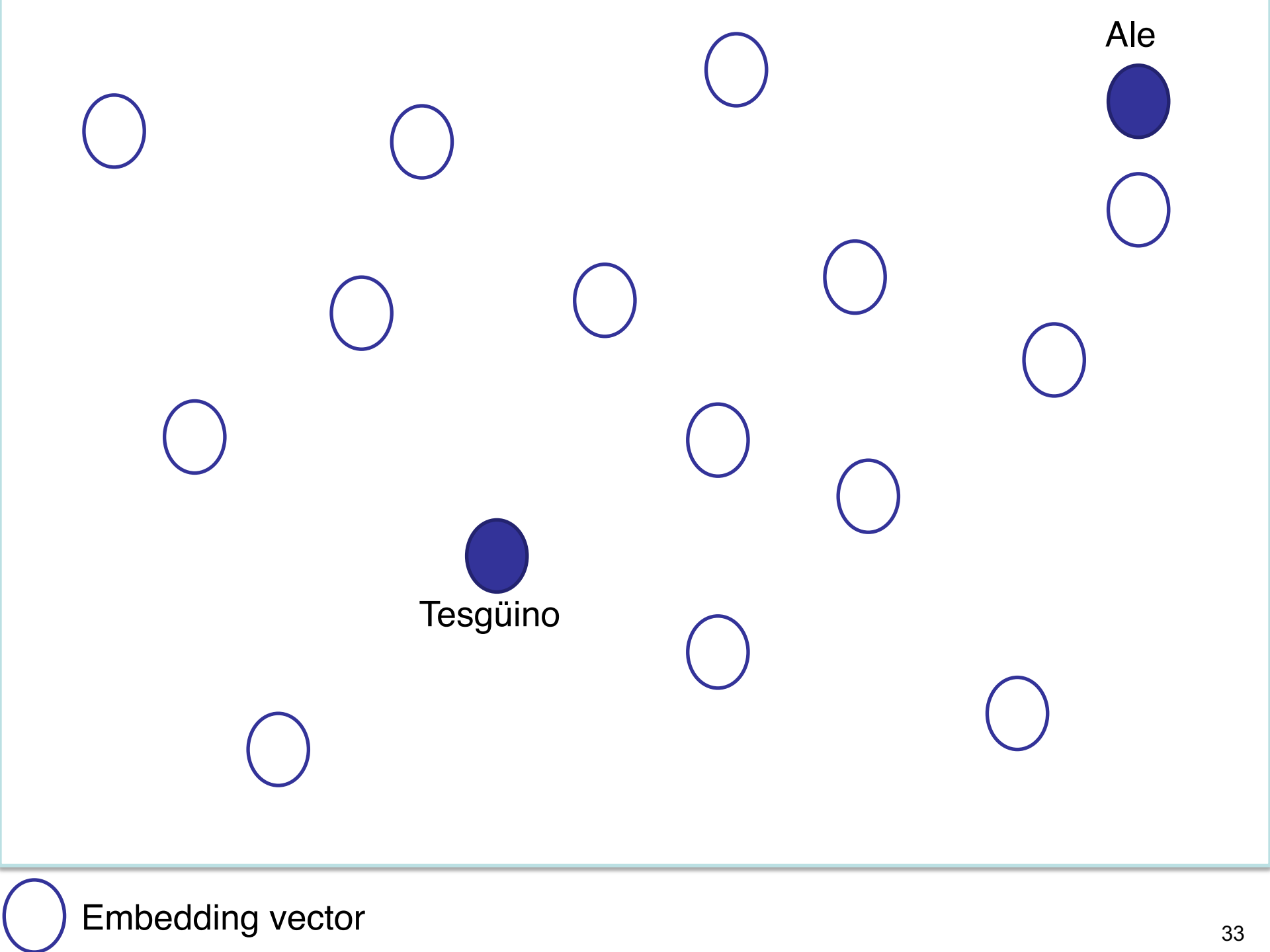
(fox, quick)  
(fox, brown)  
(fox, jumps)  
(fox, over)

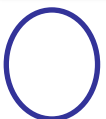


# Neural embeddings – architecture

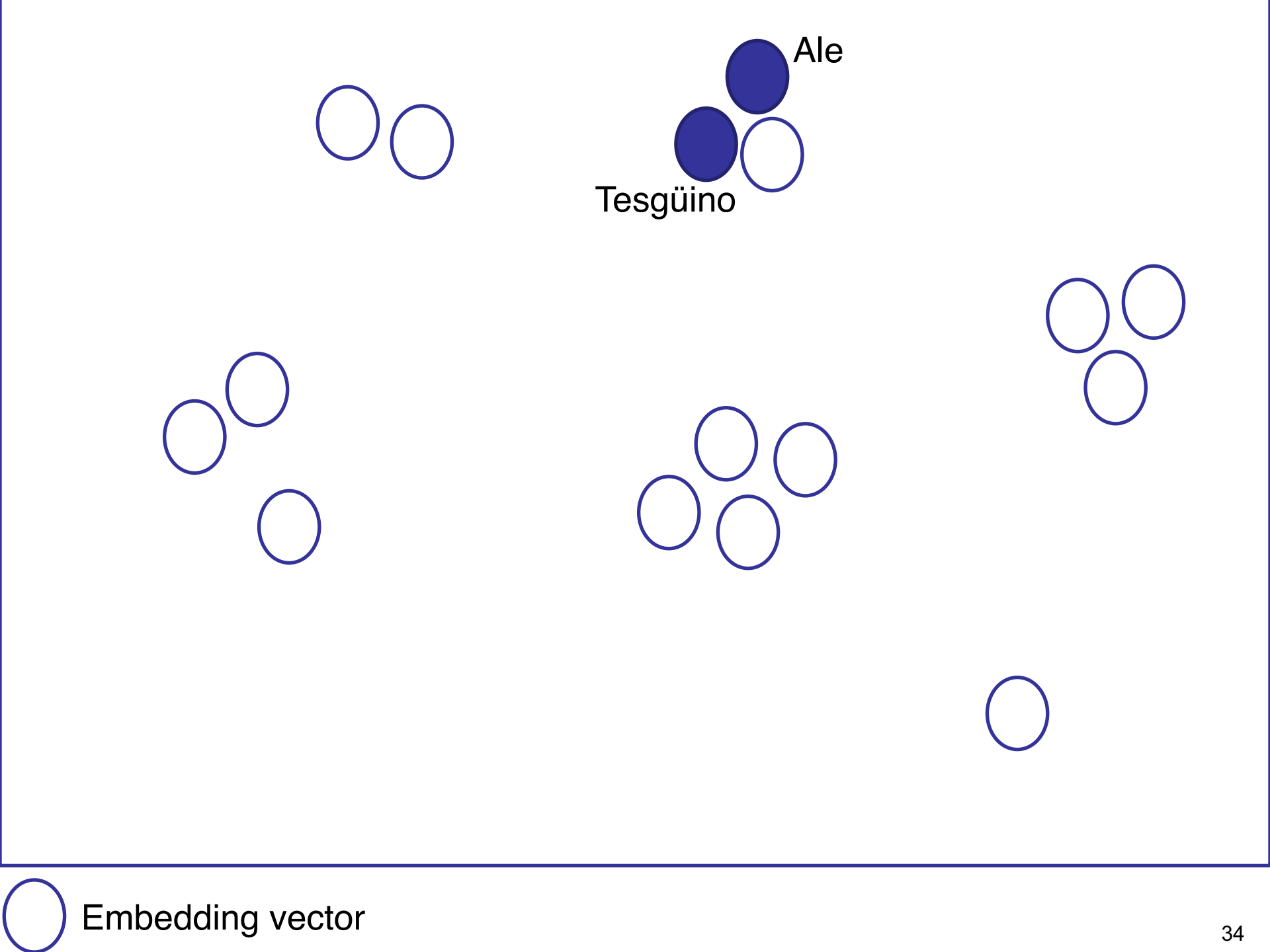
Train sample: (*Tesgüino*, *drink*)




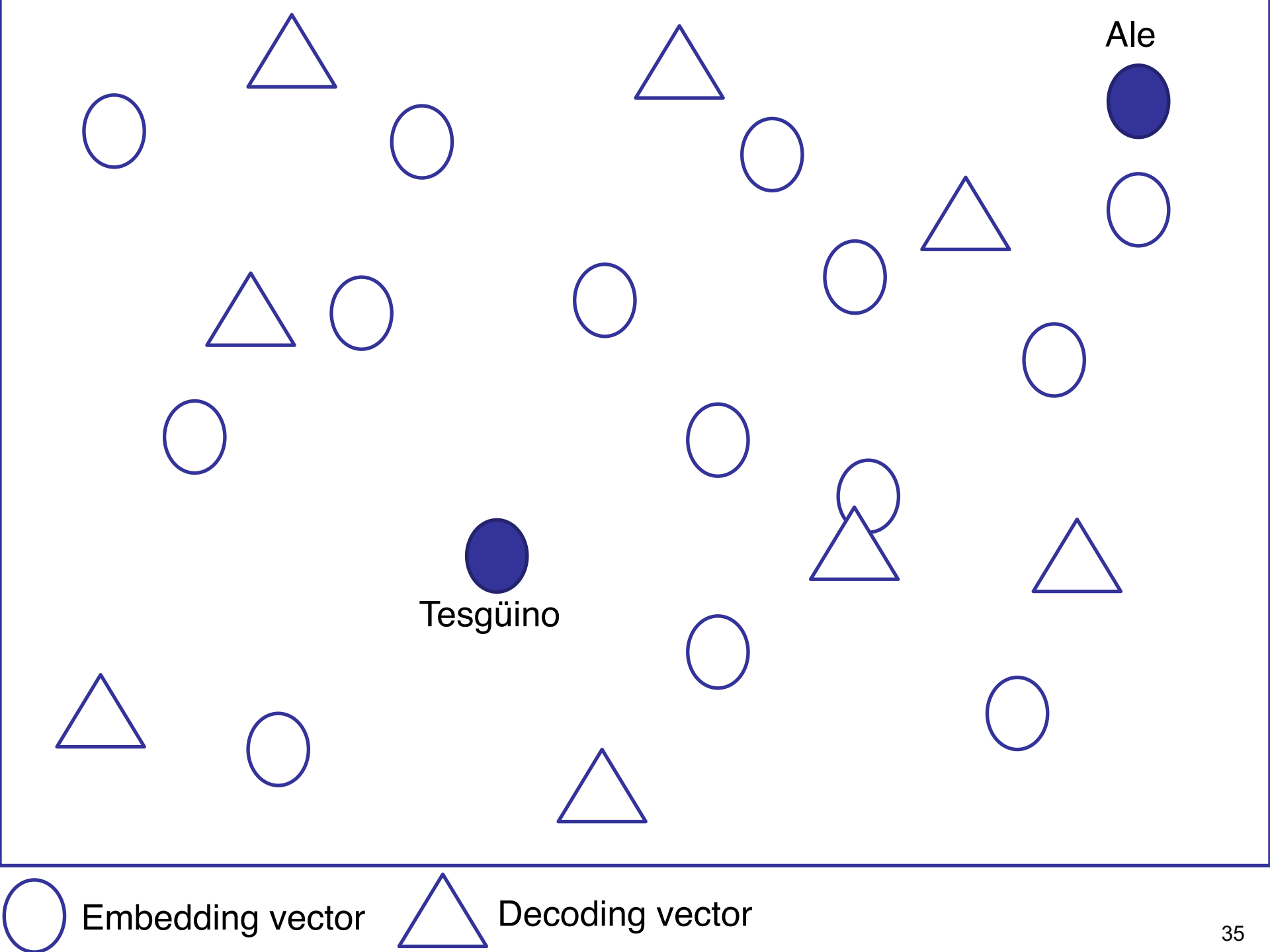


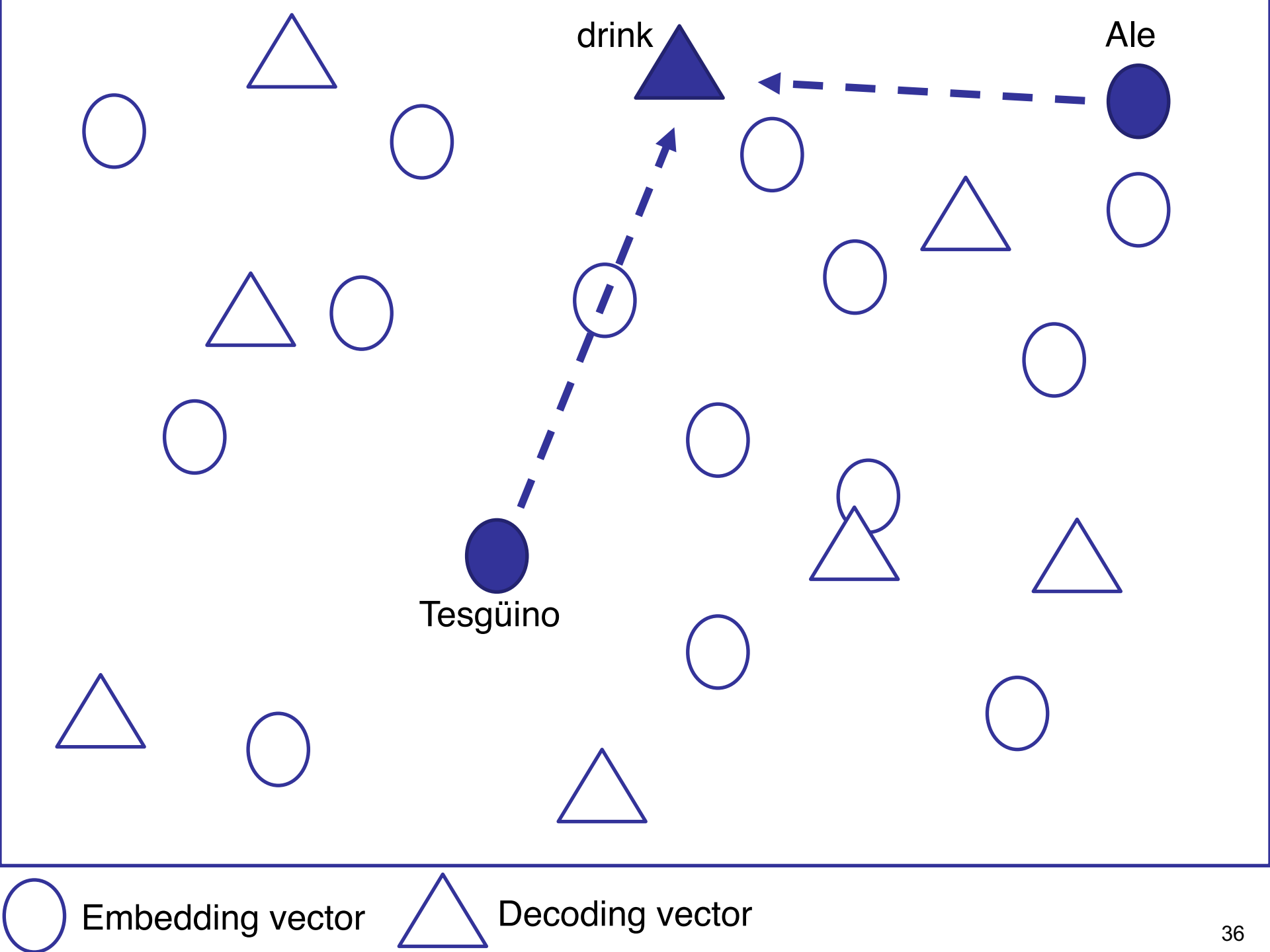
 Embedding vector

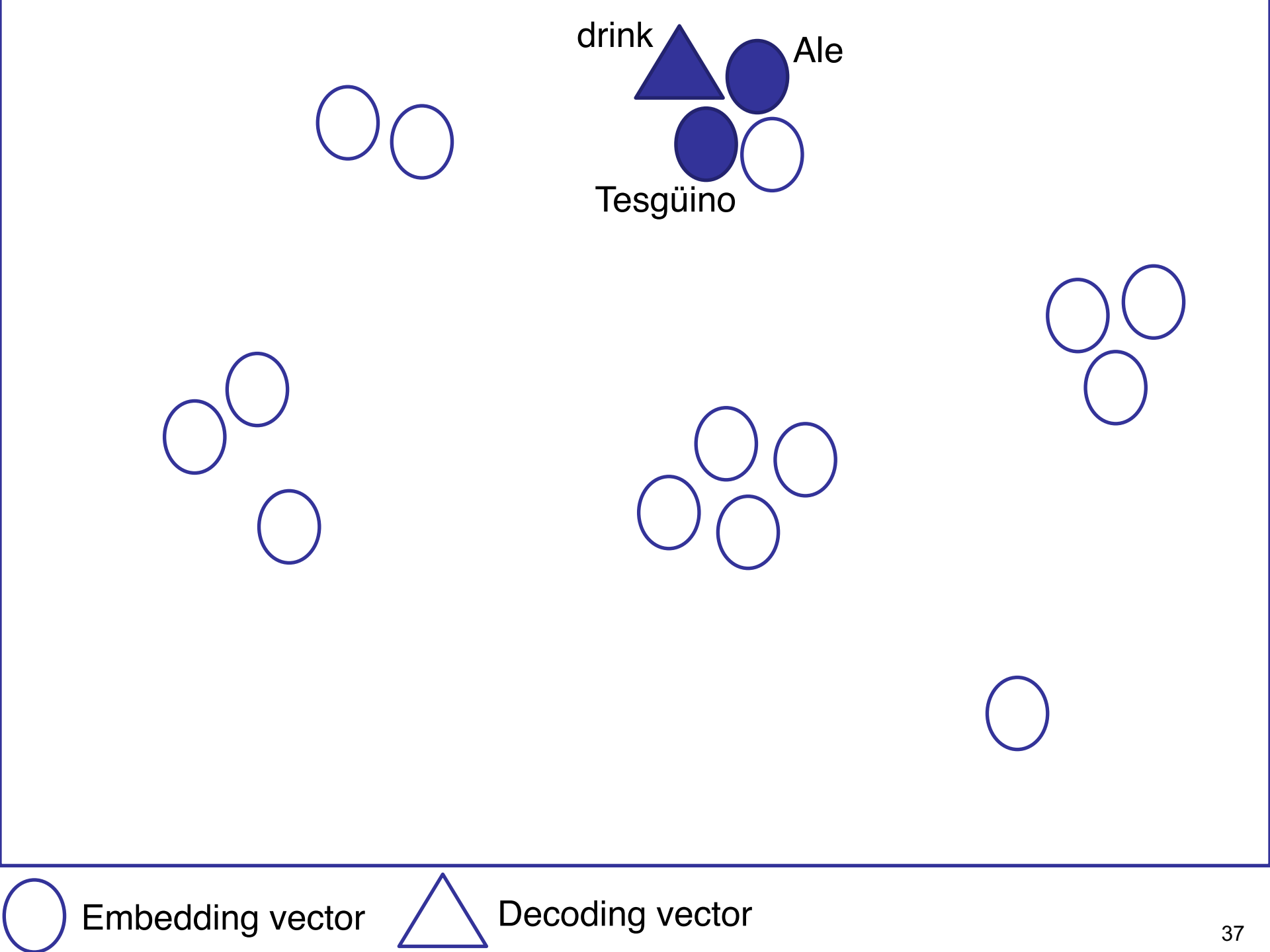


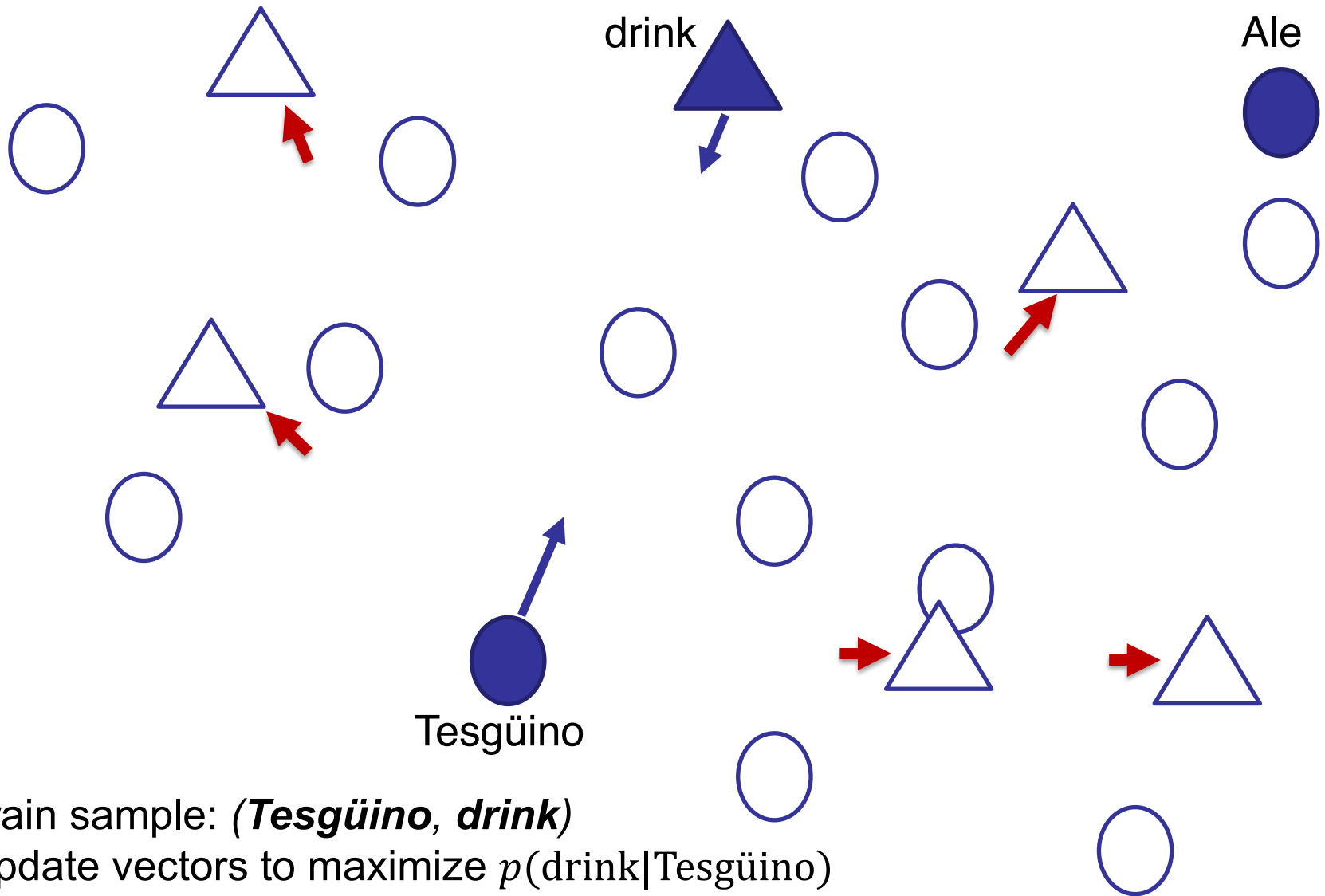


 Embedding vector









- Train sample: (***Tesgüino***, ***drink***)
- Update vectors to maximize  $p(\text{drink}|\text{Tesgüino})$

## Neural embeddings – prediction

- Since hidden layer is linear, output vector is:

$$\mathbf{z} = \mathbf{e}_v \mathbf{U}^T$$

- Probability distribution of output words using **softmax**:

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

- In this example:

$$p(\text{drink}|\text{Tesgüino}) = \frac{\exp(\mathbf{e}_{\text{Tesgüino}} \mathbf{u}_{\text{drink}}^T)}{\sum_{\tilde{c} \in \mathbb{V}} \exp(\mathbf{e}_{\text{Tesgüino}} \mathbf{u}_{\tilde{c}}^T)}$$



**Denominator (normalization) is expensive!**

## Neural embeddings – loss

- Loss function with **NLL** over all training samples  $\mathcal{D}$

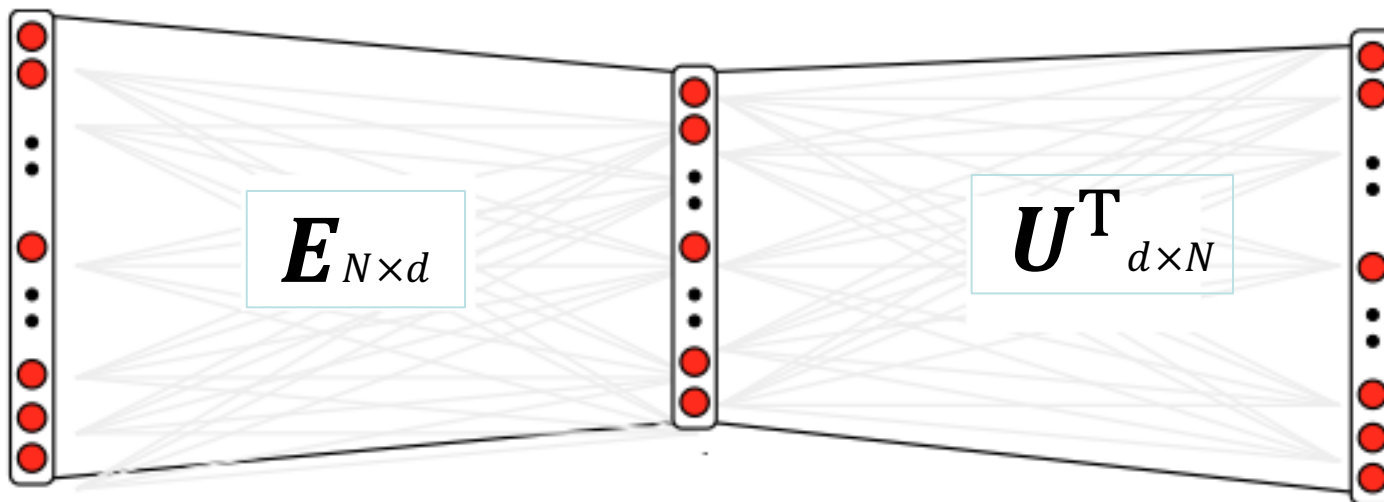
$$\mathcal{L} = -\mathbb{E}_{(v,c) \sim \mathcal{D}} \log p(c|v)$$

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

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

## Neural embeddings – word embedding

- word2vec uses the encoding matrix  $E$  as the final **word embeddings**
- Other possibilities
  - Mean of the vectors of  $E$  and  $U$  for each word
  - Concatenation of the vectors in  $E$  and  $U$  for each word
    - This results in vectors with  $2d$  dimensions





## word2vec (Skip-Gram) with Negative Sampling

- Word2vec is an **efficient** and **effective** algorithm that proposes **Negative Sampling** method to define loss
- In **Negative Sampling** instead of  $p(c|v)$ , the network estimates  $p(y = 1|v, c)$ : the probability that the **co-occurrence** of  $(v, c)$  comes from a **genuine** distribution, namely from training corpus
- $p(y = 1|v, c)$  is defined using sigmoid  $\sigma$ :

$$p(y = 1|v, c) = \sigma(\mathbf{e}_v \cdot \mathbf{u}_c^T)$$

## word2vec (Skip-Gram) with Negative Sampling

- When two words  $(v, c)$  appear in the training data (genuine distribution), it is a **positive sample**
- Negative Sampling aims to distinguish between the co-occurrence probability of  $v$  in a **positive sample**

$$p(y = 1|v, c)$$

and the co-occurrences in  $k$  **negative samples** with context words  $\tilde{c}$

$$p(y = 1|v, \tilde{c})$$

## word2vec (Skip-Gram) with Negative Sampling

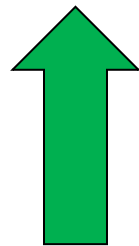
- **Negative samples** are drawn from the noisy distribution  $\tilde{\mathcal{D}}$ 
  - $\tilde{\mathcal{D}}$  represents a randomly created corpus  $\rightarrow$  words co-occur randomly!
  - Why a random co-occurrence can be assumed as a negative sample?
- The noisy distribution  $\tilde{\mathcal{D}}$  is defined using a *smooth* unigram distribution of words in the corpus,
  - In word2vec,  $\tilde{\mathcal{D}}$  is smoothed by raising unigram counts to the power of  $\alpha = 0.75 \rightarrow$  **Context Distribution Smoothing**
- Number of negative samples  $k$  is usually between 2 to 20

# word2vec with Negative Sampling – Objective Function

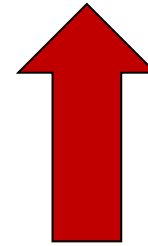
- The objective function
  - **increases** the probability for the **positive sample**  $(v, c)$
  - **decreases** the probability for the  $k$  **negative samples**  $(v, \tilde{c})$
- Loss function:

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

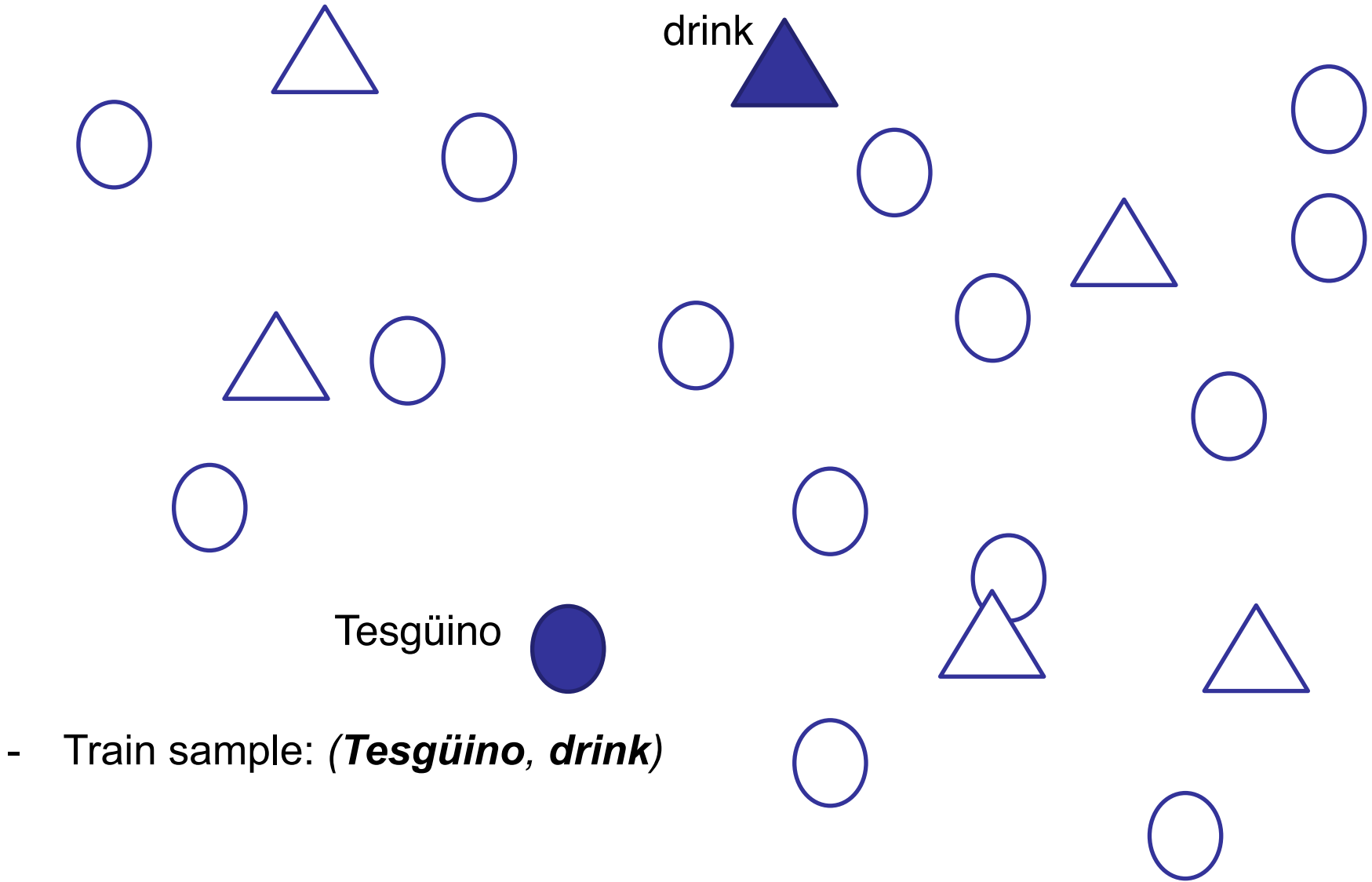
$$\mathcal{L} = -\mathbb{E}_{(v,c) \sim \mathcal{D}} \left[ \log \sigma(\mathbf{e}_v \mathbf{u}_c^T) - \sum_{\substack{\tilde{c} \sim \tilde{\mathcal{D}} \\ k \text{ times}}} \log \sigma(\mathbf{e}_v \mathbf{u}_{\tilde{c}}^T) \right]$$



positive samples



negative samples



- Train sample: (***Tesgüino***, ***drink***)

drink

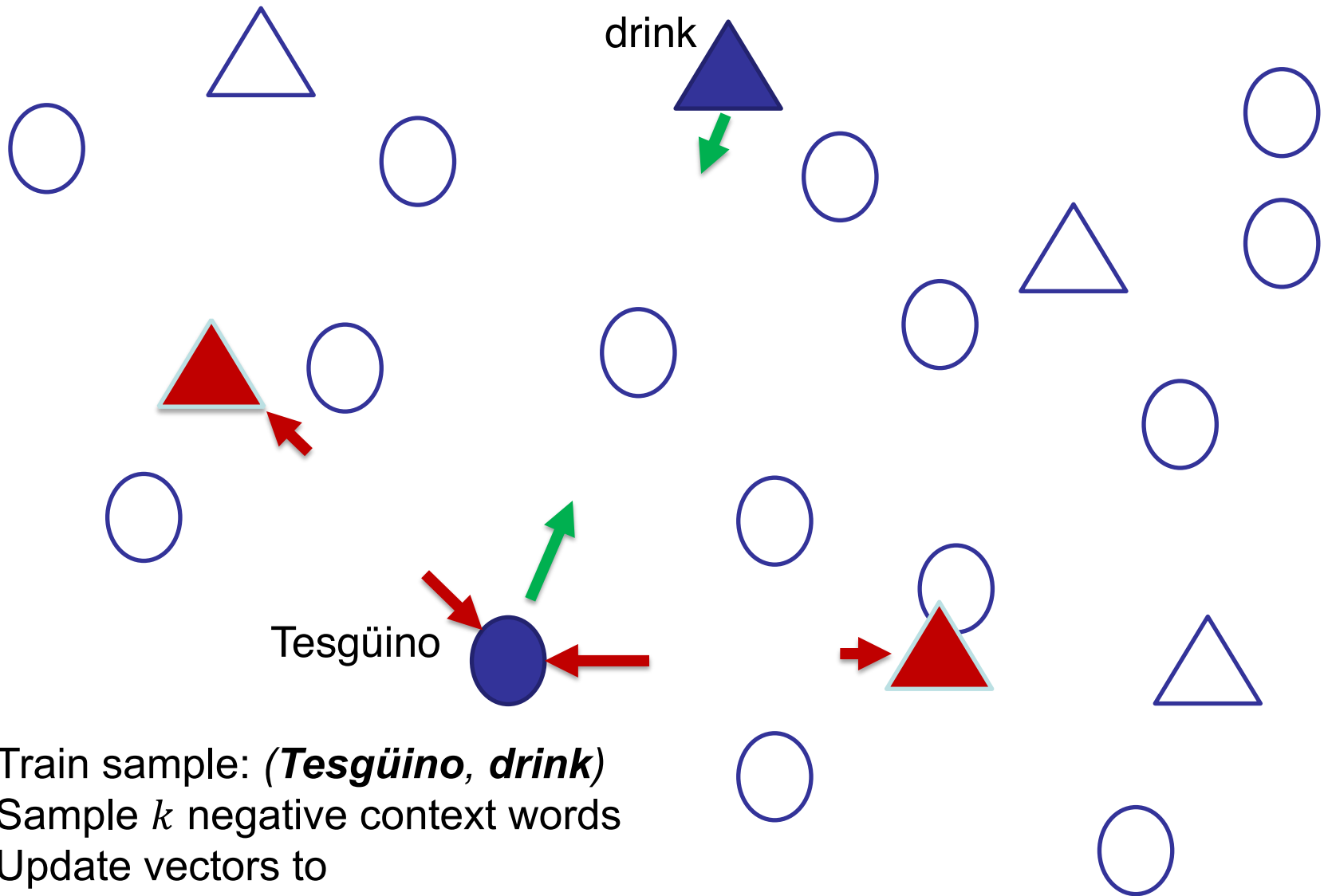


Tesgüino



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

 Embedding vector     Decoding vector



- Train sample: (***Tesgüino***, ***drink***)
- Sample  $k$  negative context words
- Update vectors to
  - Maximize  $p(y = 1 | \text{Tesgüino}, \text{drink})$
  - Minimize  $p(y = 1 | \text{Tesgüino}, \tilde{c})$

## Negative Sampling – final words!

- Negative Sampling turns the problem from multi-class classification to binary classification
- Negative Sampling is a biased approximation of softmax
  - **Noisy Contrastive Estimation** (the parent of Negative Sampling) is an unbiased approximation of softmax
- Softmax is a good choice for training language models, namely, to estimate  $p(c|v)$
- word2vec and Negative Sampling aim to train good embeddings