

INF721

2024/2



Deep Learning

L15: Word Embeddings



INF721

2024/2



Deep Learning

L15: Word Embeddings

Logistics

Announcements

- ▶ PA4 is out and due on Wednesday (13/11), 11:59pm

Last Lecture

- ▶ Language Models
- ▶ Implementing RNNs
- ▶ Vanishing/Exploding Gradients
- ▶ LSTM and GRUs

Lecture Outline

- ▶ Problems with one-hot encoding
- ▶ Word Embeddings
 - ▶ Featurized Representation
 - ▶ Visualization
 - ▶ Properties
 - ▶ Applications
- ▶ Word2Vec and Negative Sampling
- ▶ GloVe

Problem with one-hot encoding

A problem of **one-hot encoding** is that it represents each word as an independent category


Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \dots \\ 1 \\ \dots \\ 0 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \dots \\ 1 \\ \dots \\ 0 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ \dots \\ 1 \\ \dots \\ 0 \\ 0 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \dots \\ 1 \\ \dots \\ 0 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ \dots \\ 1 \\ \dots \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \dots \\ 1 \\ \dots \\ 0 \\ 0 \end{bmatrix}$
$\mathbf{0}_{5391}$	$\mathbf{0}_{5391}$	$\mathbf{0}_{5391}$	$\mathbf{0}_{5391}$	$\mathbf{0}_{5391}$	$\mathbf{0}_{5391}$

It doesn't let a model easily generalize across words!

- For example, consider a language model that gives high probability to the sentence:
*"I want a glass of **orange** juice"*
- Now, consider sampling from this model with context:
*"I want a glass of **apple** _____"*
- To the model, the relationship between **apple** and **orange** is the same as **apple** and **man**, or queen.
- This is because the distance between any two words is the same!

Word Embedding: Featurized Representation

Ideally, we would like to have featurized representation for words, where words with similar meaning to have a similar representation:

apple and orange are close in the embedding space 

	Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
Gender	-1	1	-0.95	0.97	0.00	0.01
Royal	0.01	0.02	0.93	0.95	-0.01	0.00
Age	0.03	0.02	0.7	0.69	0.03	-0.02
Food	0.04	0.01	0.02	0.01	0.95	0.97
...						
	e_{5391}	e_{9853}	e_{4914}	e_{7157}	e_{456}	e_{6257}

Word embeddings are learned featurized representations:

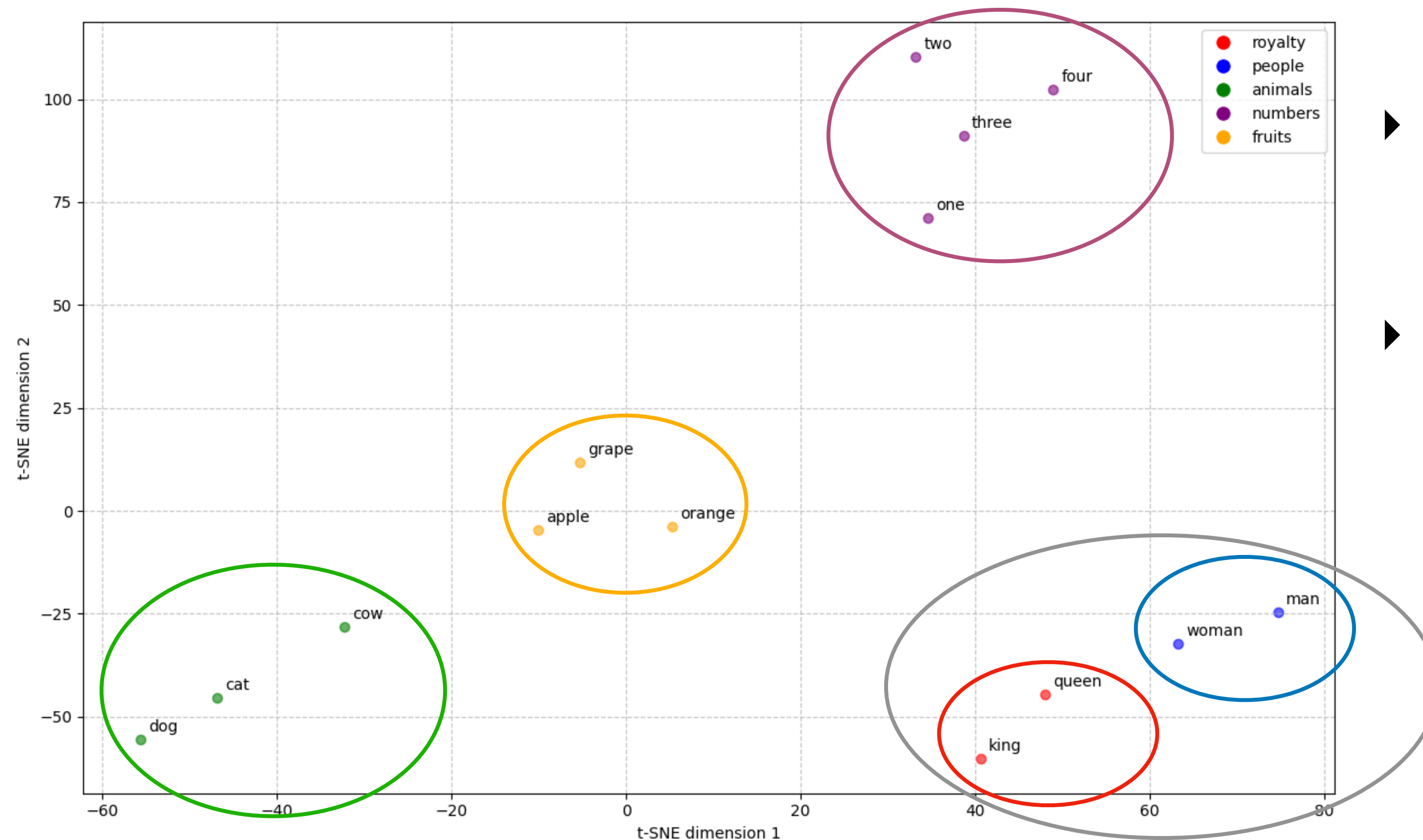
- ▶ We can define the number of features d to be learned (e.g., 300), but **not** what they represent (e.g., gender)
- ▶ Words with similar meaning to have a similar representation:
- ▶ They allow a model to generalize across words more easily:

*"I want a glass of **orange** juice"*

*"I want a glass of **apple** juice "*

Visualizing Word Embeddings (with t-SNE)

We can visualize high dimensional word embeddings (e.g., 300 features) in 2D using the t-SNE algorithm:



- ▶ t-SNE learns a non-linear mapping between the original representation and a 2D space
- ▶ Word embeddings tend to group words with similar meaning together in the embedding space.

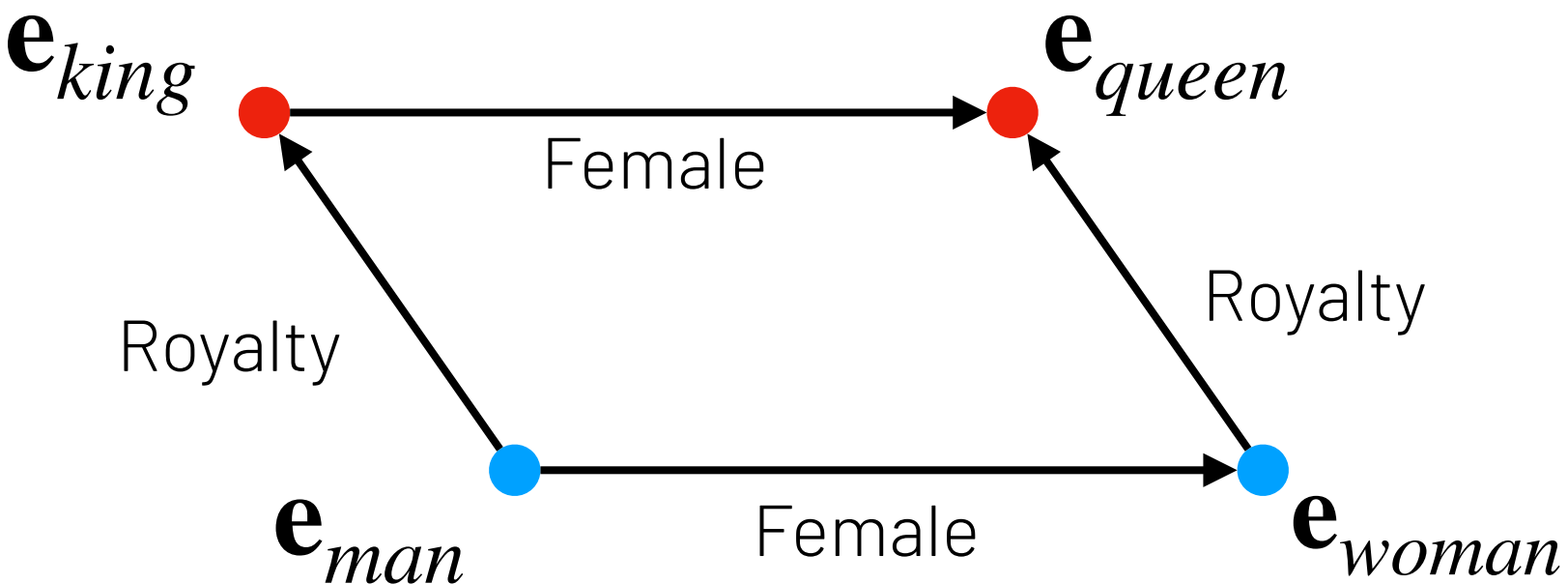
Properties of word embeddings

One interesting property of word embeddings is that they allow **word analogies** to be solved with vector arithmetic.

	Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
Gender	-1	1	-0.95	0.97	0.00	0.01
Royal	0.01	0.02	0.93	0.95	-0.01	0.00
Age	0.03	0.02	0.7	0.69	0.03	-0.02
Food	0.04	0.01	0.02	0.01	0.95	0.97
...						
	\mathbf{e}_{man}	\mathbf{e}_{woman}	\mathbf{e}_{king}	\mathbf{e}_{queen}	\mathbf{e}_{apple}	\mathbf{e}_{orange}

Man is to Woman as King is to _____?

$$\mathbf{e}_{king} - \mathbf{e}_{queen} \approx \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad \mathbf{e}_{man} - \mathbf{e}_{woman} \approx \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$



$$\mathbf{e}_{man} - \mathbf{e}_{woman} \approx \mathbf{e}_{king} - \mathbf{e}_w$$

$$\mathbf{e}_w = \underset{w}{\operatorname{argmax}} \underbrace{\operatorname{sim}(\mathbf{e}_w, \mathbf{e}_{king} - \mathbf{e}_{man} + \mathbf{e}_{woman})}_{\text{Similarity}}$$

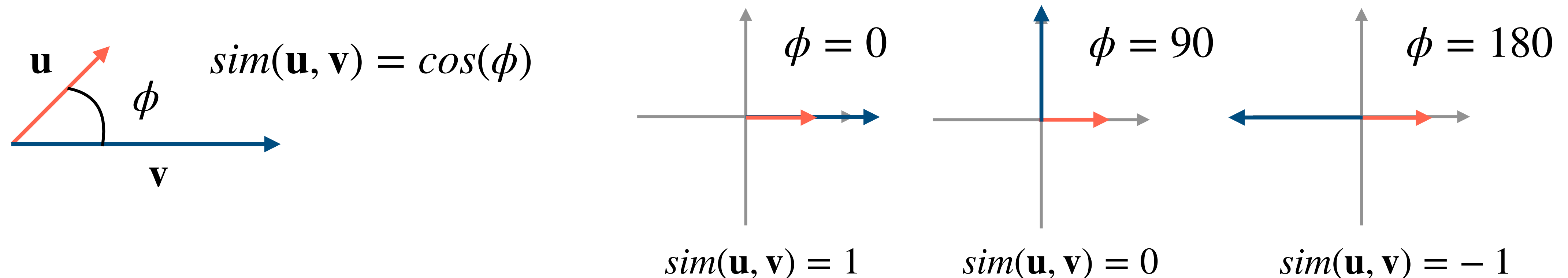
Similarity Between Vectors

Measuring similarity between vectors is very important in deep learning and one of the most popular functions to do that is the **cosine similarity**:

The **cosine similarity** between two vectors \mathbf{u} and \mathbf{v} is defined as:

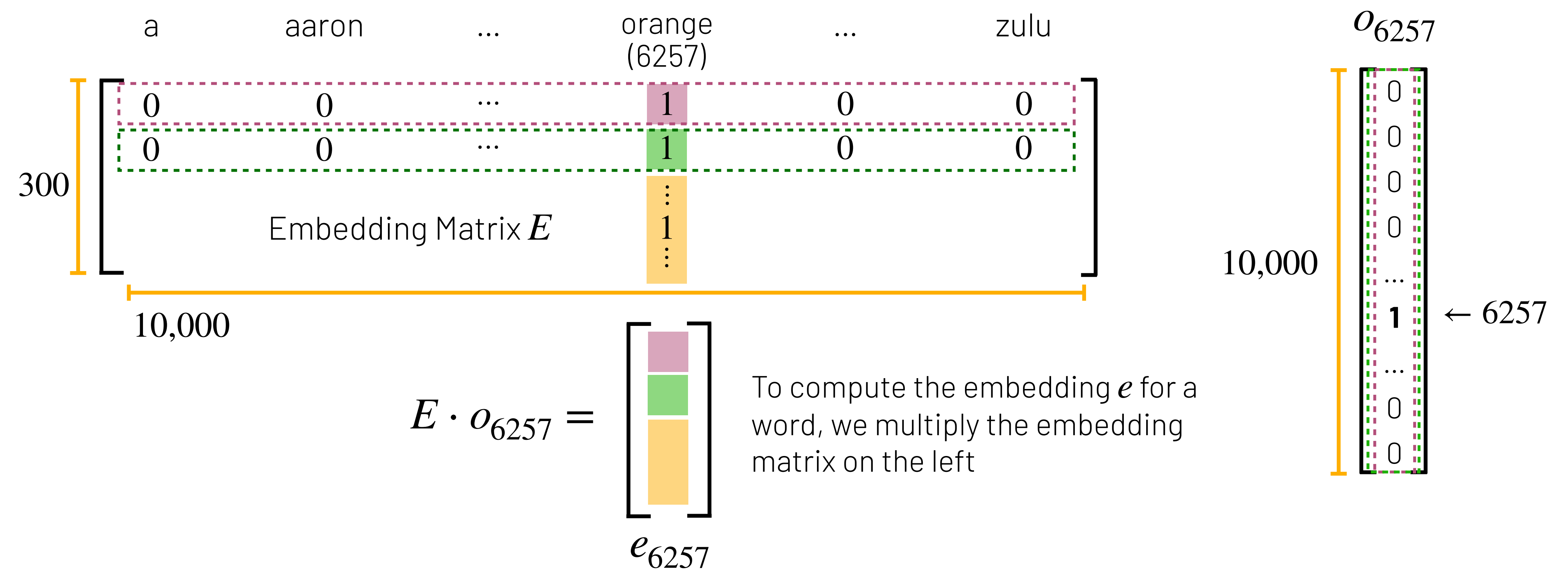
$$\text{sim}(\mathbf{u}, \mathbf{v}) = \frac{\mathbf{u} \cdot \mathbf{v}}{\|\mathbf{u}\|_2 \|\mathbf{v}\|_2}$$

This formula is equal to the cosine of the angle between \mathbf{u} and \mathbf{v} :



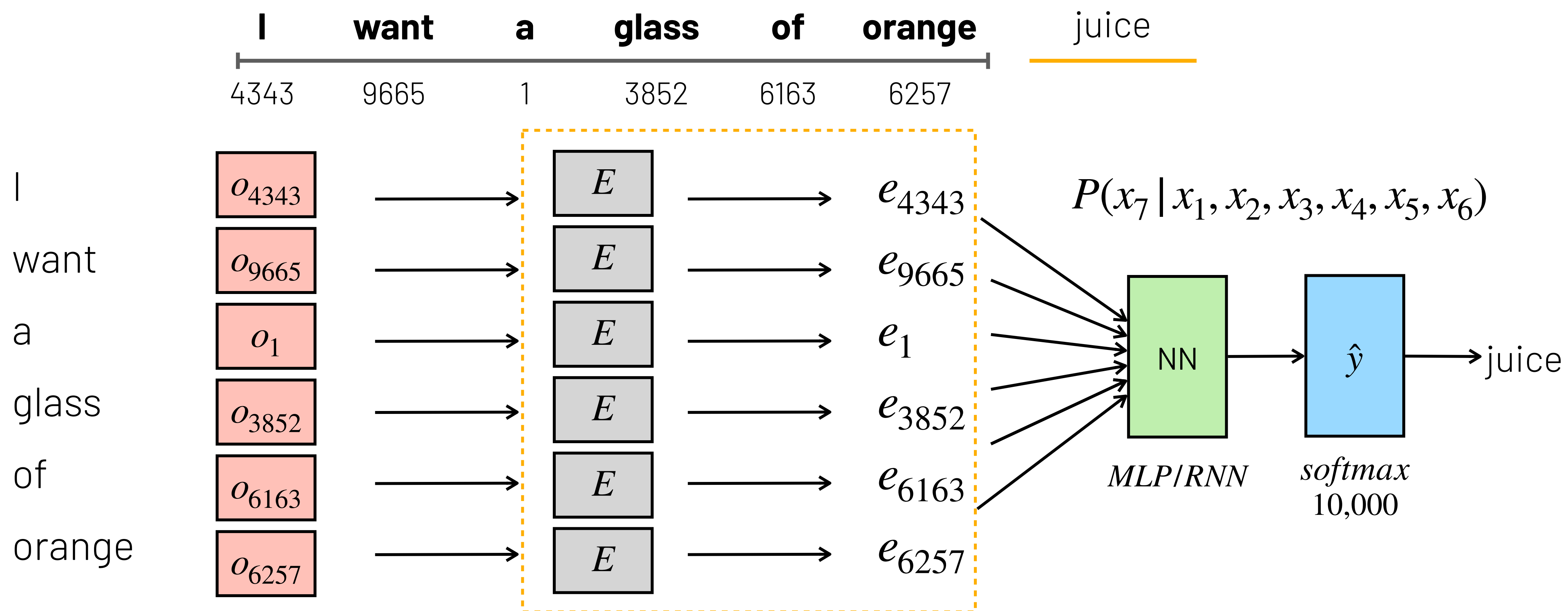
Embedding Matrix

Word embeddings with dimensionality d (e.g. 300) are learned by adding an extra weight matrix, called **embeddings matrix** E , to the model:



Learning Word Embeddings

There are many ways to learn word embeddings, but a simple and popular one is to train a **language model** with an **embedding layer** before the model:



Learning Word Embeddings

If your goal is not to learn a language model, but just the embeddings, we don't need to be limited to the previous words as **context**. We can use all the surrounding words:

Last words:
(Language Model)

I want a glass of orange juice to drink with my cereal

Context Target

Surrounding words:
(Word2Vec)

I want a glass of orange juice to drink with my cereal

Context Target Context

Word2Vec

Word2Vec is a model to learn word embeddings using the surrounding words as context

► Examples are produced by:

1. Randomly sampling a **target word** and;
2. Extracting the *m* (e.g, 3) words to left and to the right as **context**

I want a glass of orange juice to drink with my cereal

A horizontal sequence of words: "I", "want", "a", "glass", "of", "orange", "juice", "to", "drink", "with", "my", "cereal". The word "juice" is underlined in orange. Two horizontal brackets are positioned below the words. The first bracket starts under "glass" and ends under "orange". The second bracket starts under "to" and ends under "with".

[Mikolov et. al., 2013] evaluated two different types of models:

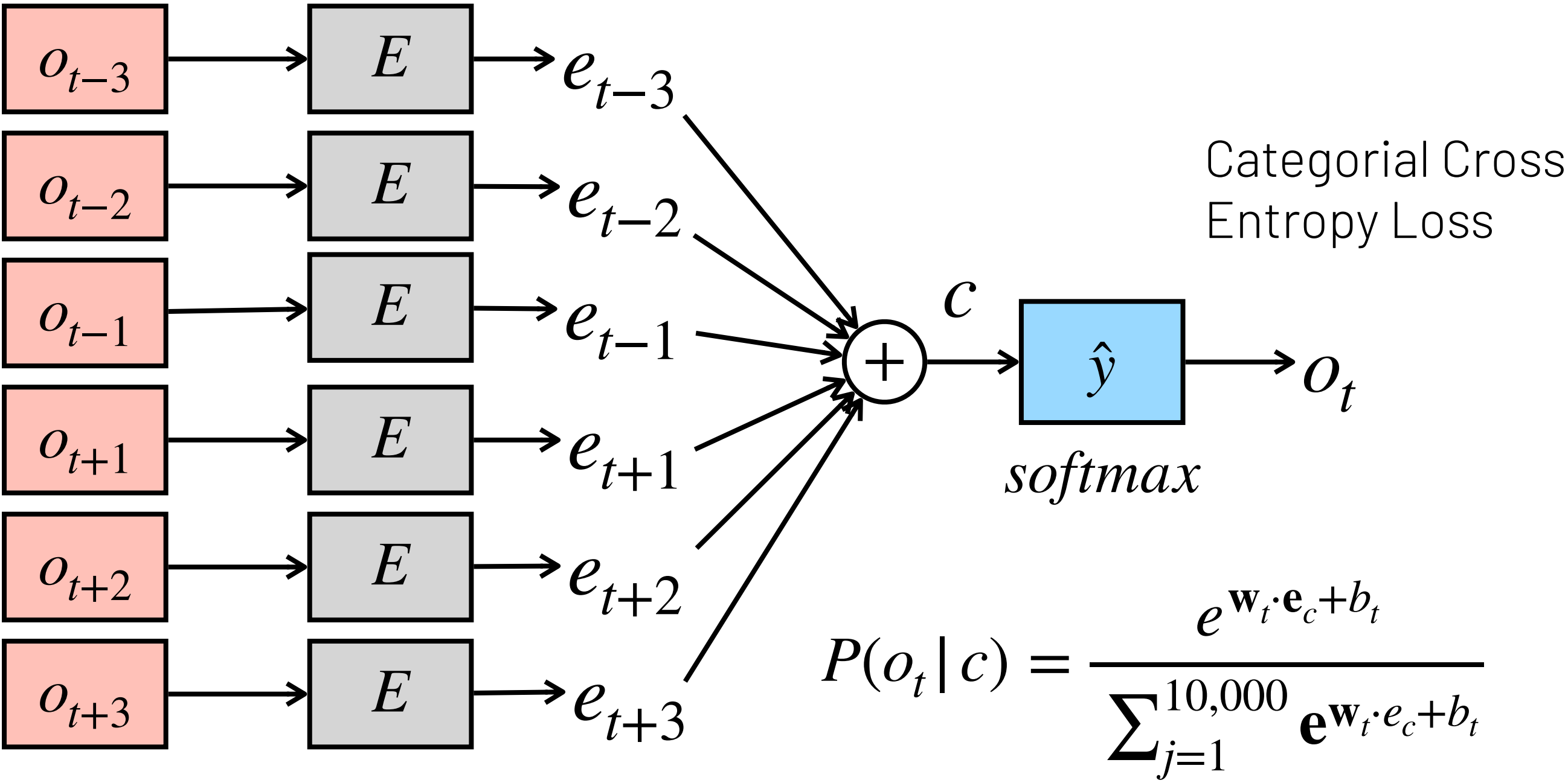
- **Continuous Bag of Words (CBOW)**: The goal is to predict the **target word** from the **context**
- **Skip-gram**: The goal is to predict the **context words** from the **target word**!

Word2Vec: CBOW

In **CBOW** model, the goal is to predict the **target word** from the **context**.

I want a glass of orange juice to drink with my cereal

Context	Target
...	...
glass of orange _____ to drink with	juice
...	...

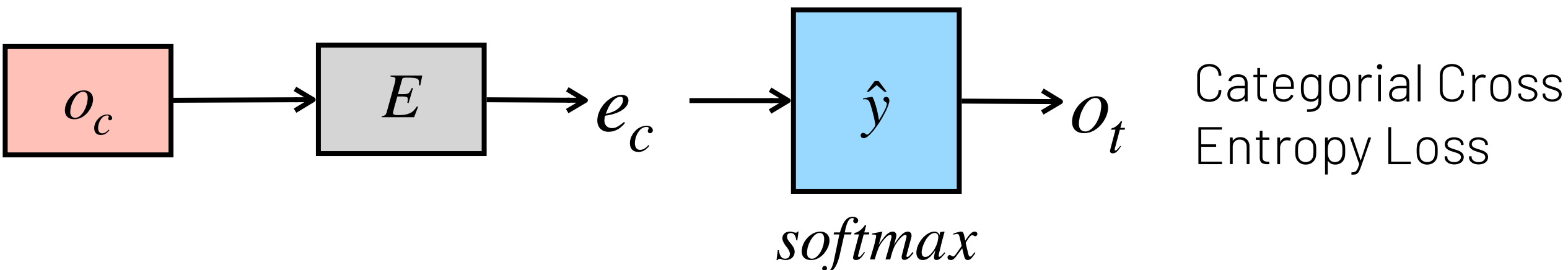


Word2Vec: Skip-Gram

In the **Skip-gram** model, the goal is to predict the **context words** from the **target word**!

I want a glass of orange juice to drink with my cereal

Target	Context
juice	glass
juice	of
juice	orange
juice	to
juice	drink
juice	with



$$P(o_t | c) = \frac{e^{\mathbf{w}_t \cdot \mathbf{e}_c + b_t}}{\sum_{j=1}^{10,000} e^{\mathbf{w}_t \cdot \mathbf{e}_c + b_t}}$$

Negative Sampling

Since the vocabulary size is typically large (e.g., 10,000), computing a probability with the softmax layer is expensive!

$$P(t | c) = \frac{e^{\mathbf{w}_t \cdot \mathbf{e}_c + b_t}}{\sum_{j=1}^{10,000} e^{\mathbf{w}_t \cdot \mathbf{e}_c + b_t}}$$

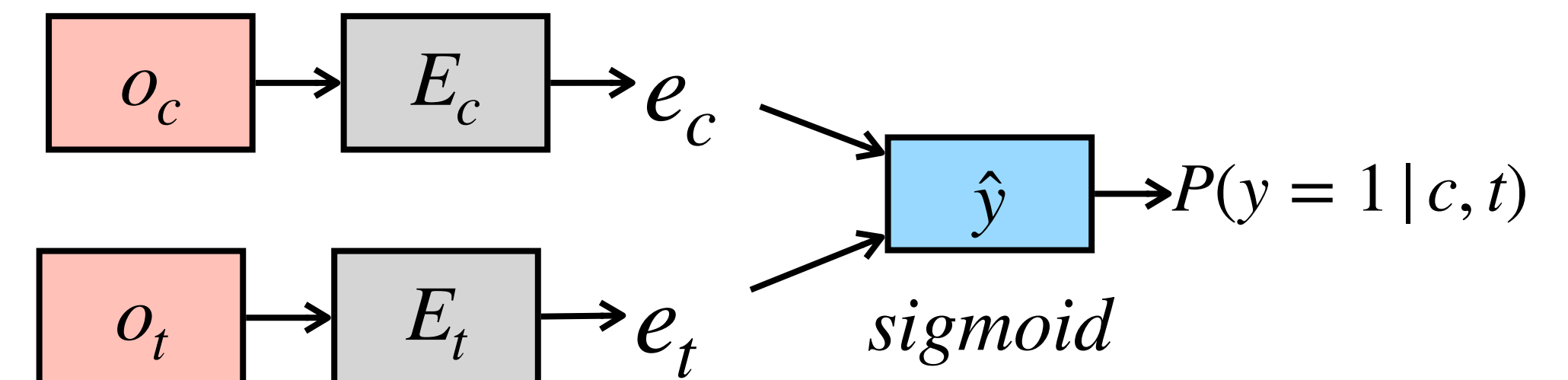
(we have to sum 10,000 terms)

Negative sampling changes the classification text to a binary classification problem by generating a dataset with positive and negative examples:

Context	Word	Target?
orange	juice	1
orange	king	0
orange	book	0
orange	the	0
orange	of	0

Sample with skip-gram

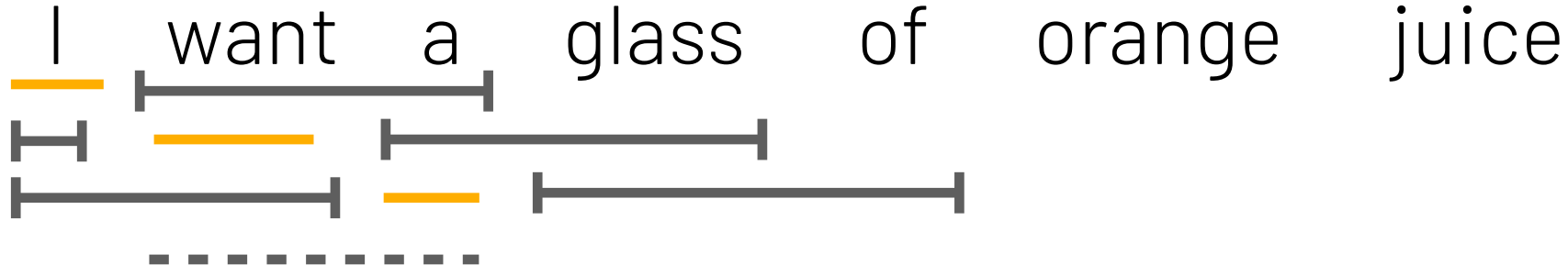
Sample k
random words
from the
dictionary



$$P(y = 1 | c, t) = \sigma(\mathbf{e}_t \cdot \mathbf{e}_c)$$

GloVe (Global Vectors for Word Representation)

In Word2Vec, the global information is not preserved. To address that, GloVe predicts the number of times a target word t appears in the context of the context word c

- 1. Build a word co-occurrence matrix X

- 2. Assign a weight vector per context \mathbf{w}_i and target word $\tilde{\mathbf{w}}_j$
- 3. Predict log co-occurrence from \mathbf{w}_i and $\tilde{\mathbf{w}}_j$ using weighted MSE

	$\tilde{\mathbf{w}}_1$	$\tilde{\mathbf{w}}_2$	$\tilde{\mathbf{w}}_3$	$\tilde{\mathbf{w}}_4$	$\tilde{\mathbf{w}}_5$	$\tilde{\mathbf{w}}_6$	$\tilde{\mathbf{w}}_7$
	I	want	a	glass	of	orange	juice
\mathbf{w}_1 I	0	1.0	0.5	0	0	0	0
\mathbf{w}_2 want	1.0	0	1.0	0.5	0	0	0
\mathbf{w}_3 a	0.5	1.0	0	1.0	0.5	0	0
\mathbf{w}_4 glass	0	0.5	1.0	0	1.0	0.5	0
\mathbf{w}_5 of	0	0	0.5	1.0	0	1.0	0.5
\mathbf{w}_6 orange	0	0	0	0.5	1.0	0	1.0
\mathbf{w}_7 juice	0	0	0	0	0.5	1.0	0

$X_{ij} = \sum_j^C \frac{1}{dist(i,j)}$
where C is context window size (e.g, 2)

$$L = \sum_{i=1}^{10,000} \sum_{j=1}^{10,000} f(X_{ij}) (\mathbf{w}_i \tilde{\mathbf{w}}_j + b_i + b_j - \log(X_{ij}))^2$$

Prediction

Target

weighting function that helps balance the importance of different co-occurrences

Using log helps compress the range of values

4. Embedding:

$$e_i = \frac{w_i + \tilde{w}_i}{2}$$

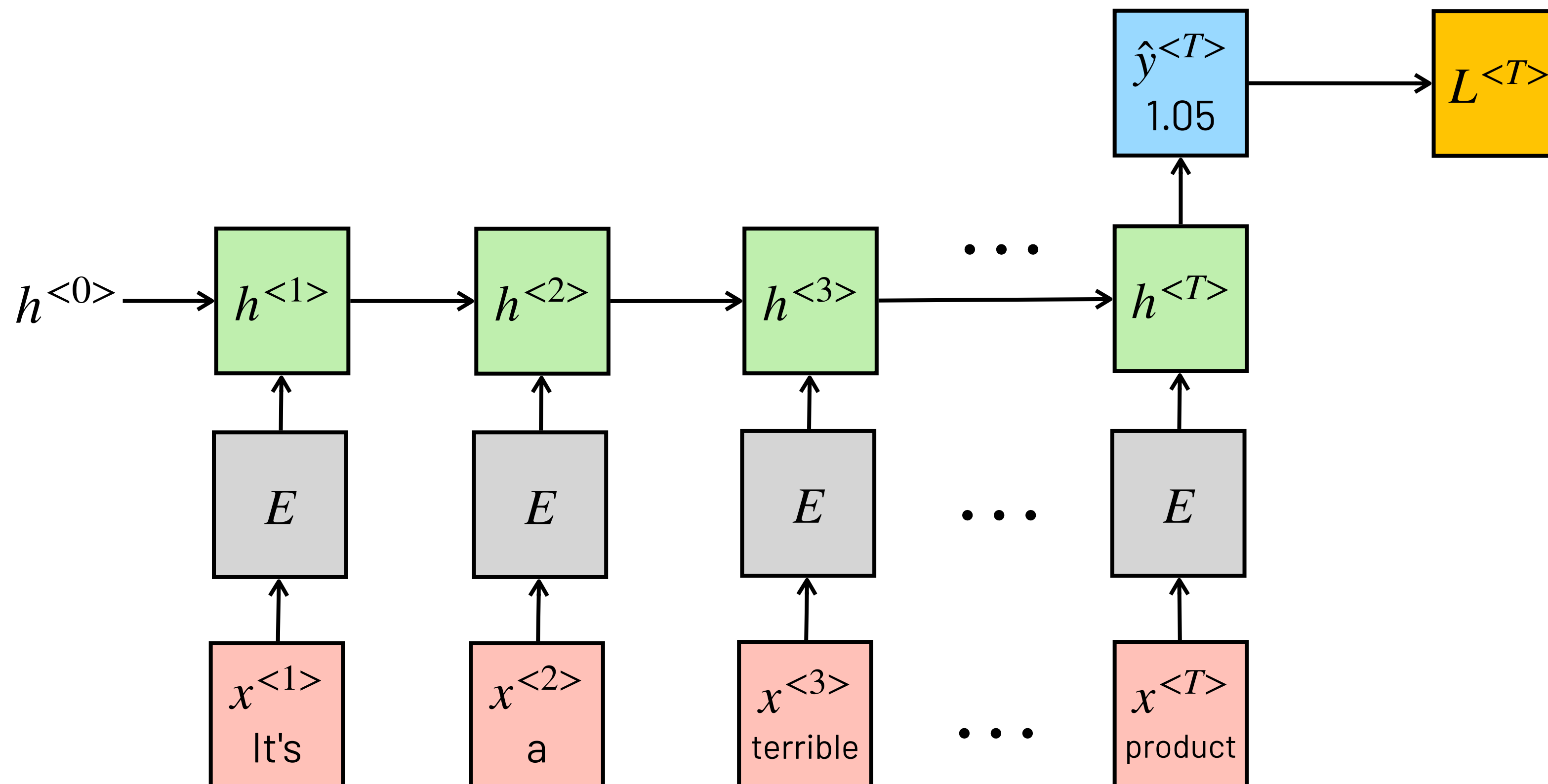
Using Word Embeddings

Sentiment Analysis

"It's a terrible product."



When working with Natural Language Processing problems, you can learn word embeddings by just adding an embedding layer to your network:



Using Word Embeddings

Sentiment
Analysis

"It's a terrible product." ★☆☆☆☆

In PyTorch, you can use the Embedding layer to add an Embedding Matrix E to your model:

```
class RNNSentiment(nn.Module):
    def __init__(self, vocab_size, embedding_dim, hidden_dim, output_dim, padding_idx=None):
        super().__init__()

        # Embedding layer
        self.embedding = nn.Embedding(
            num_embeddings=vocab_size, embedding_dim=embedding_dim, padding_idx=padding_idx
        )

        # RNN layer
        self.rnn = nn.RNN(input_size=embedding_dim, hidden_size=hidden_dim, num_layers=1,
                           batch_first=True)

        # Output layer
        self.fc = nn.Linear(hidden_dim, output_dim)
```

Next Lecture

L15: Attention Mechanisms

Machine Translation, Decoding Strategies, Attention Mechanisms in RNNs