

# Text Classification



Credit of the slides goes to Lena Voita and Stanford CS

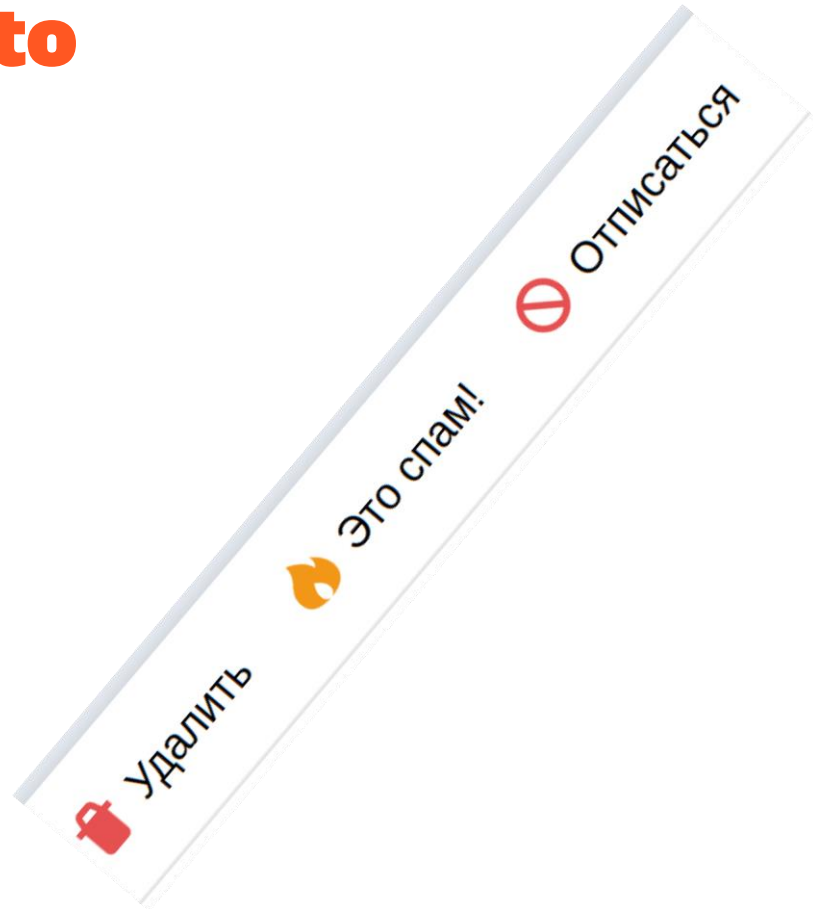
Why do we need to classify texts?

# Why do we need to classify texts?

- As a self-sufficient task:

# Why do we need to classify texts?

- As a self-sufficient task:
  - Spam filtering



# Why do we need to classify texts?

- As a self-sufficient task:
  - Spam filtering
  - Sentiment analysis

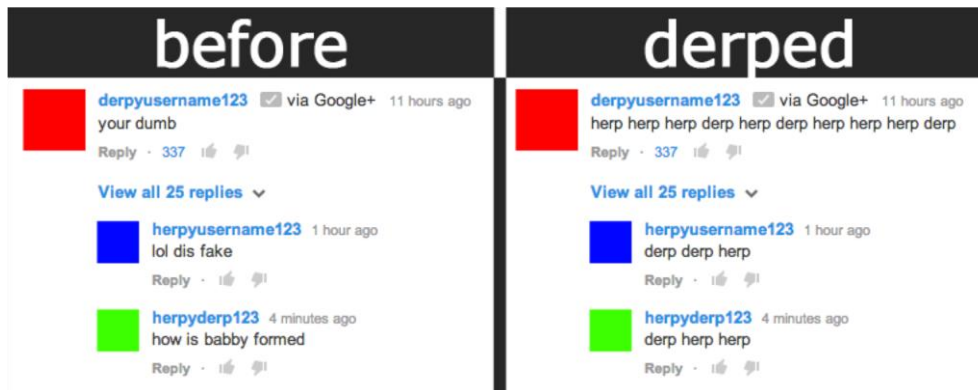


## Is Twitter better at predicting elections than opinion polls?



# Why do we need to classify texts?

- As a self-sufficient task:
  - Spam filtering
  - Sentiment analysis
  - Fake news/clickbait detection
  - Troll/bot protection



# Why do we need to classify texts?

- As a self-sufficient task:
  - Spam filtering
  - Sentiment analysis
  - Fake news/clickbait detection
  - Troll/bot protection
- As a part of more complicated NLP tasks



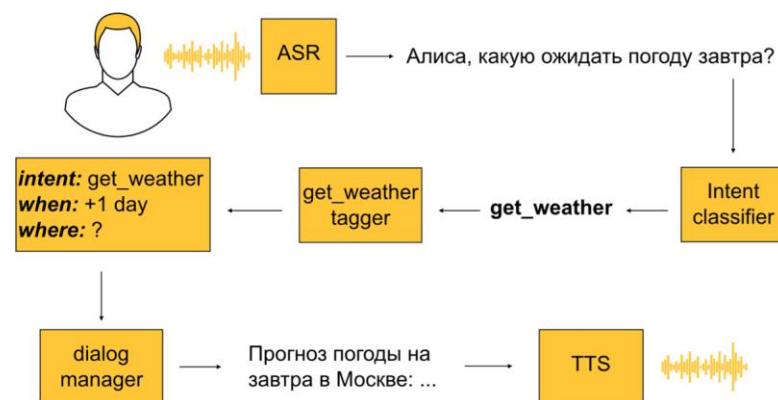
# Why do we need to classify texts?

- As a self-sufficient task:
  - Spam filtering
  - Sentiment analysis
  - Fake news/clickbait detection
  - Troll/bot protection
- As a part of more complicated NLP tasks
  - Data filtering



# Why do we need to classify texts?

- As a self-sufficient task:
  - Spam filtering
  - Sentiment analysis
  - Fake news/clickbait detection
  - Troll/bot protection
- As a part of more complicated NLP tasks
  - Data filtering
  - Intent classification in dialog systems



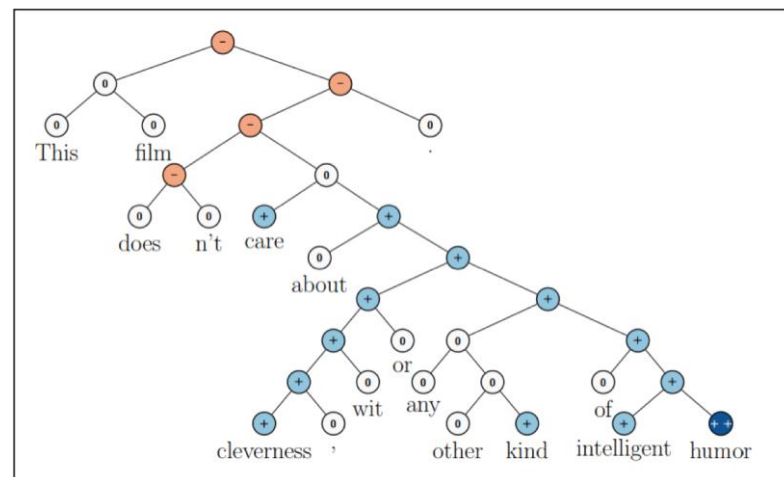
# Why do we need to classify texts?

- As a self-sufficient task:
  - Spam filtering
  - Sentiment analysis
  - Fake news/clickbait detection
  - Troll/bot protection
- As a part of more complicated NLP tasks
  - Data filtering
  - Intent classification in dialog systems
  - Hybrid machine translation systems :-)



# Popular Benchmarks

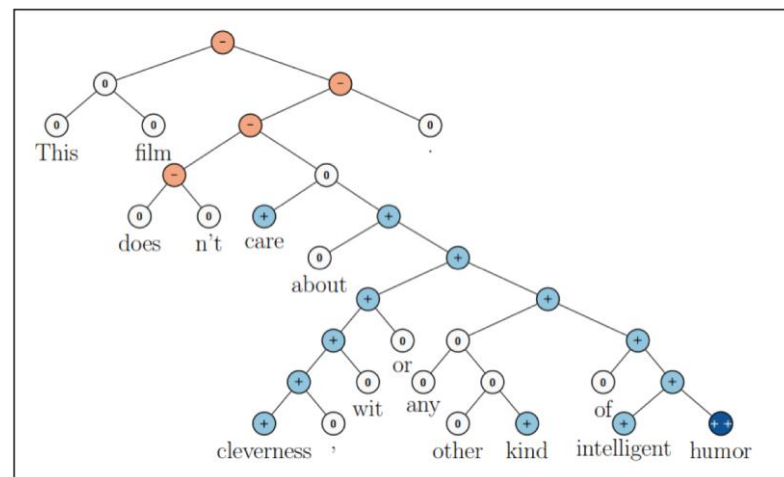
Dataset	Size
<a href="#">Question classification (TREC)</a>	6K
<a href="#">MPQA Opinion corpus (SUBJ)</a>	8K
<a href="#">Movie Reviews (MR)</a>	10K
<a href="#">Reuters-21578</a>	21.5K
<a href="#">IMDB Reviews</a>	25K
<a href="#">Stanford Sentiment Treebank (SST)</a>	9.5K
<a href="#">Sogou News</a>	0.5M
<a href="#">AG News</a>	1M
<a href="#">Yelp Dataset</a>	5.2M
<a href="#">Amazon Reviews</a>	35M
<a href="#">Flickr 100m</a>	100M



Subjective unigram	Objective unigram
amazing, beautiful, cheap, decent, effective, fantastic, good, happy, impress, jittery, light, madly, nice, outstanding, perfect, quick, responsive, sharp, terrible, ultimate, wonderful.	access, because, chance, default, entire, few, go, half, inside, job, keep, know, last, matter, new, only, past, quality, read, several, text, use, version, was, young.

# Popular Benchmarks

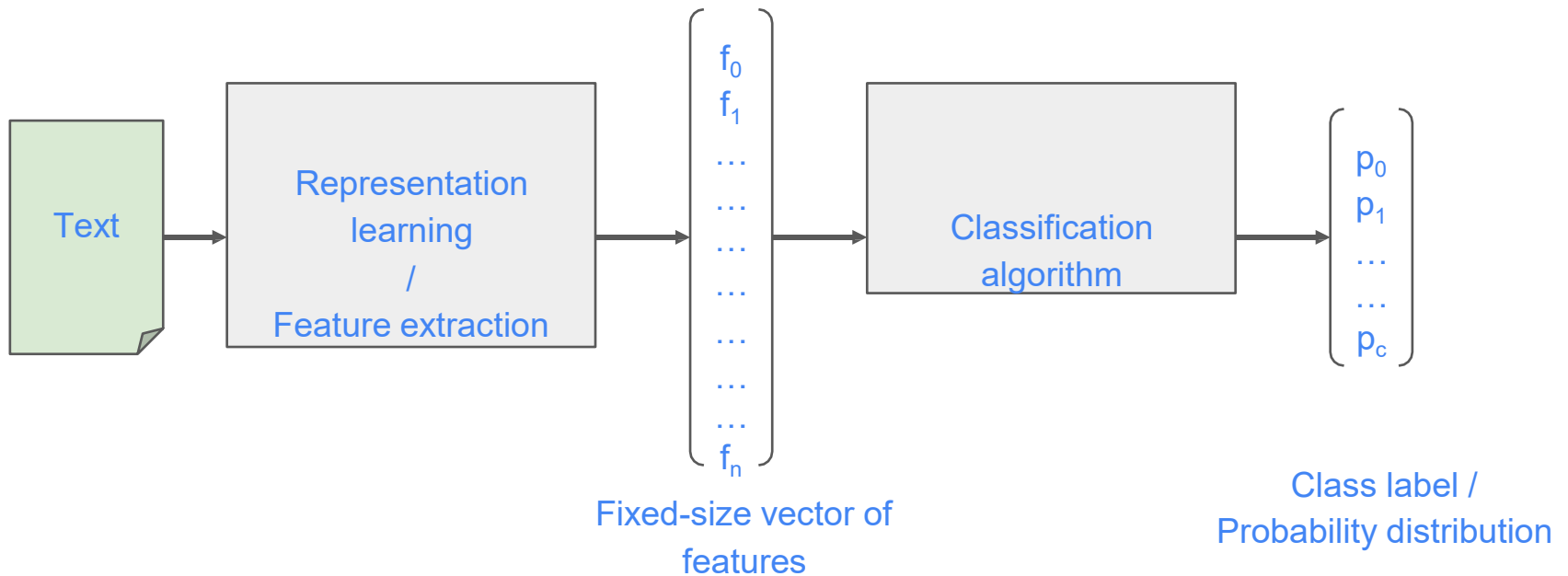
Dataset	Size
<a href="#">Question classification (TREC)</a>	6K
<a href="#">MPQA Opinion corpus (SUBJ)</a>	8K
<a href="#">Movie Reviews (MR)</a>	10K
<a href="#">Reuters-21578</a>	21.5K
<a href="#">IMDB Reviews</a>	25K
<del><a href="#">Stanford Sentiment Treebank (SST)</a></del>	<del>9.5K</del>
<a href="#">Sogou News</a>	0.5M
<a href="#">AG News</a>	1M
<a href="#">Yelp Dataset</a>	5.2M
<a href="#">Amazon Reviews</a>	35M
<a href="#">Flickr 100m</a>	100M



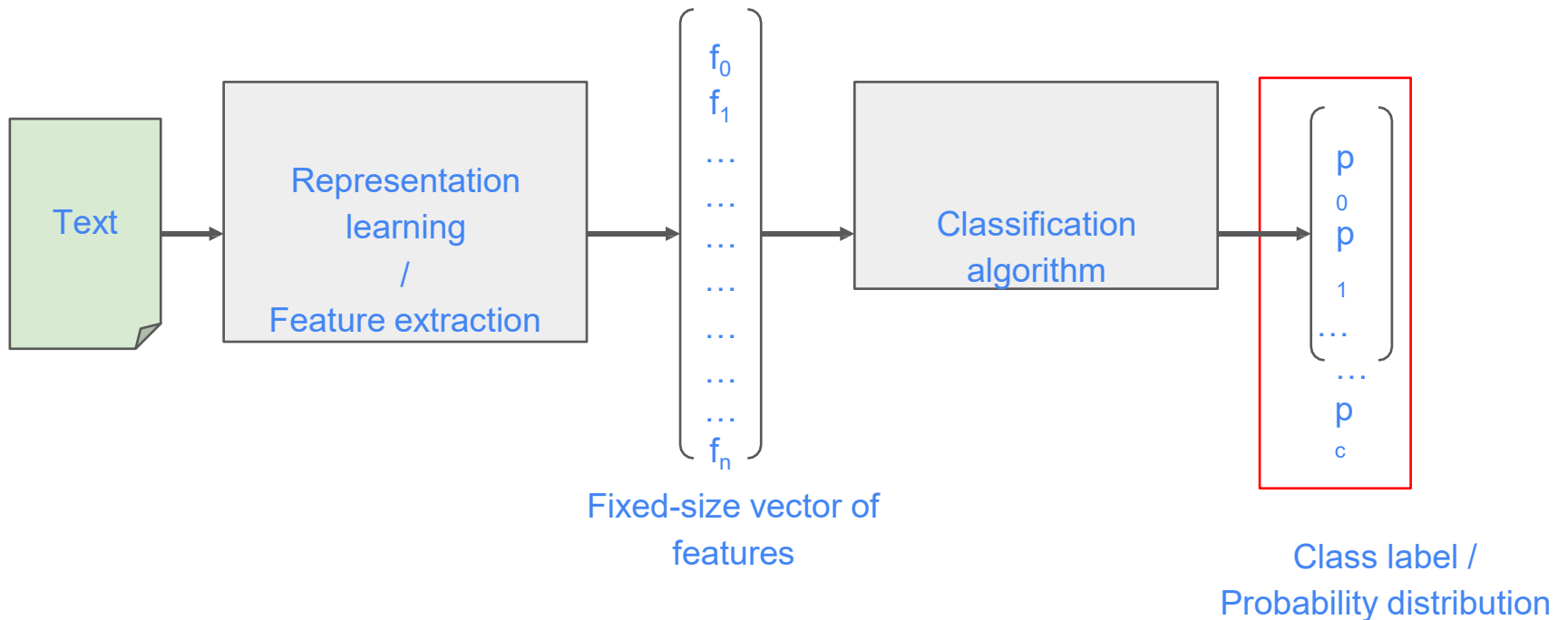
Subjective unigram	Objective unigram
amazing, beautiful, cheap, decent, effective, fantastic, good, happy, impress, jittery, light, madly, nice, outstanding, perfect, quick, responsive, sharp, terrible, ultimate, wonderful.	access, because, chance, default, entire, few, go, half, inside, job, keep, know, last, matter, new, only, past, quality, read, several, text, use, version, was, young.

~2012

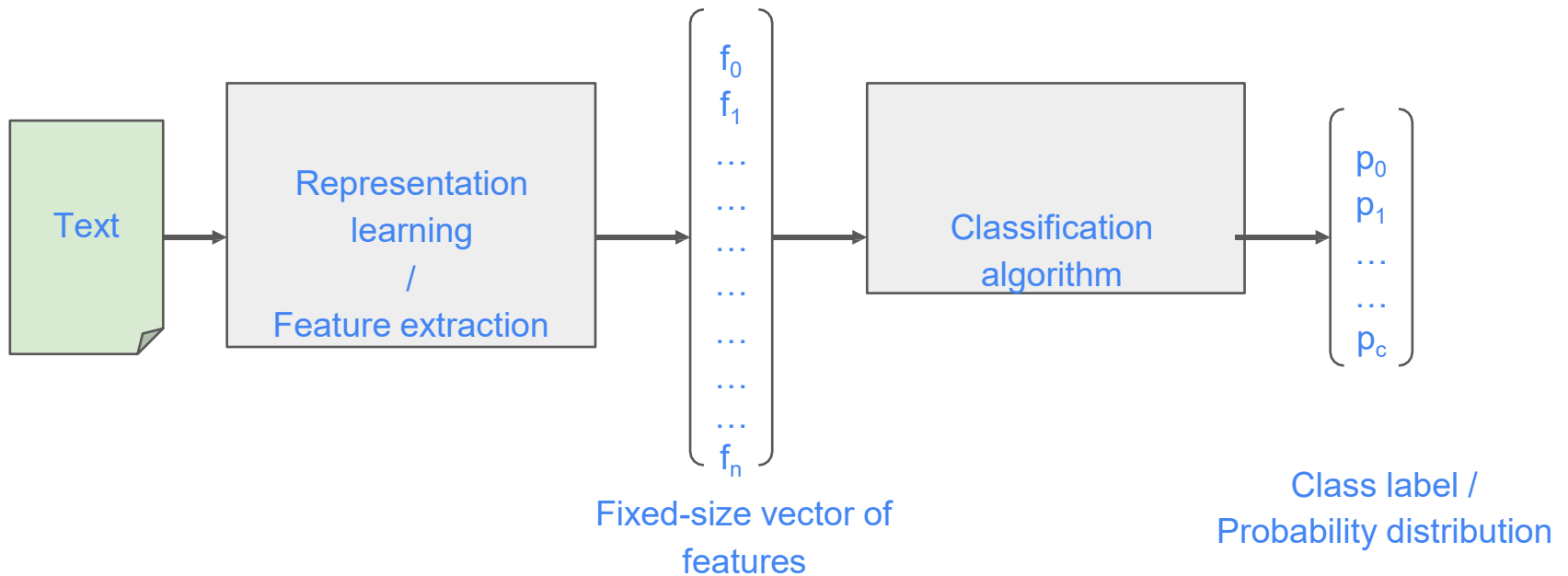
# Text classification in general



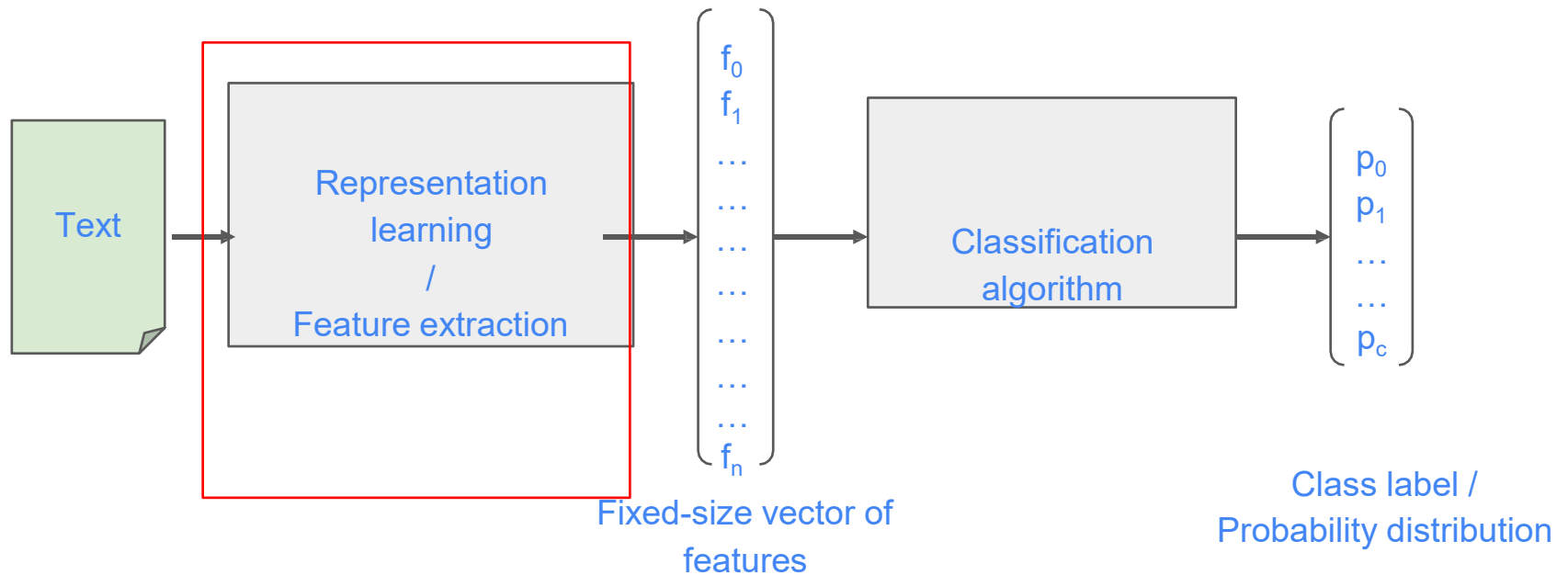
# Text classification in general



# Text classification in general



# Text classification in general





# Text representation: feature engineering

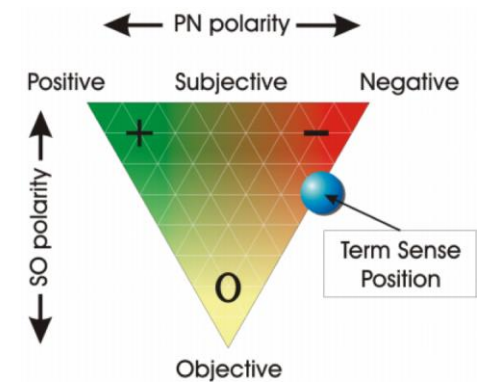
As for many ML tasks, it is possible to generate useful features by hands.

Like what?

# Text representation: feature engineering

As for many ML tasks, it is possible to generate useful features by hands.

- General statistics: text length, text length variance,...
- Scores from tagged word lists:
  - Sentiment dictionaries: [SentiWordNet](#), [SentiWords](#), ...
  - Subjectivity/objectivity dictionaries: [MPQA](#)
  - ...
- Syntactic features:
  - POS tags
- Ad-hoc features: e.g. number of emojis ( 🤮 or 😡 )



# Sparse text representation: BOW

The quick brown fox jumps over the dog .

# Sparse text representation: BOW

$$\begin{pmatrix} 1 \\ 0 \\ \dots \\ 0 \\ \dots \\ 0 \end{pmatrix}$$

The

$$\begin{pmatrix} 0 \\ 1 \\ \dots \\ 0 \\ \dots \\ 0 \end{pmatrix}$$

quick

$$\begin{pmatrix} 0 \\ 0 \\ \dots \\ \dots \\ \dots \\ 0 \end{pmatrix}$$

brown

fox

jumps

over

$$\begin{pmatrix} 1 \\ 0 \\ \dots \\ \dots \\ \dots \\ 0 \end{pmatrix}$$

the

dog

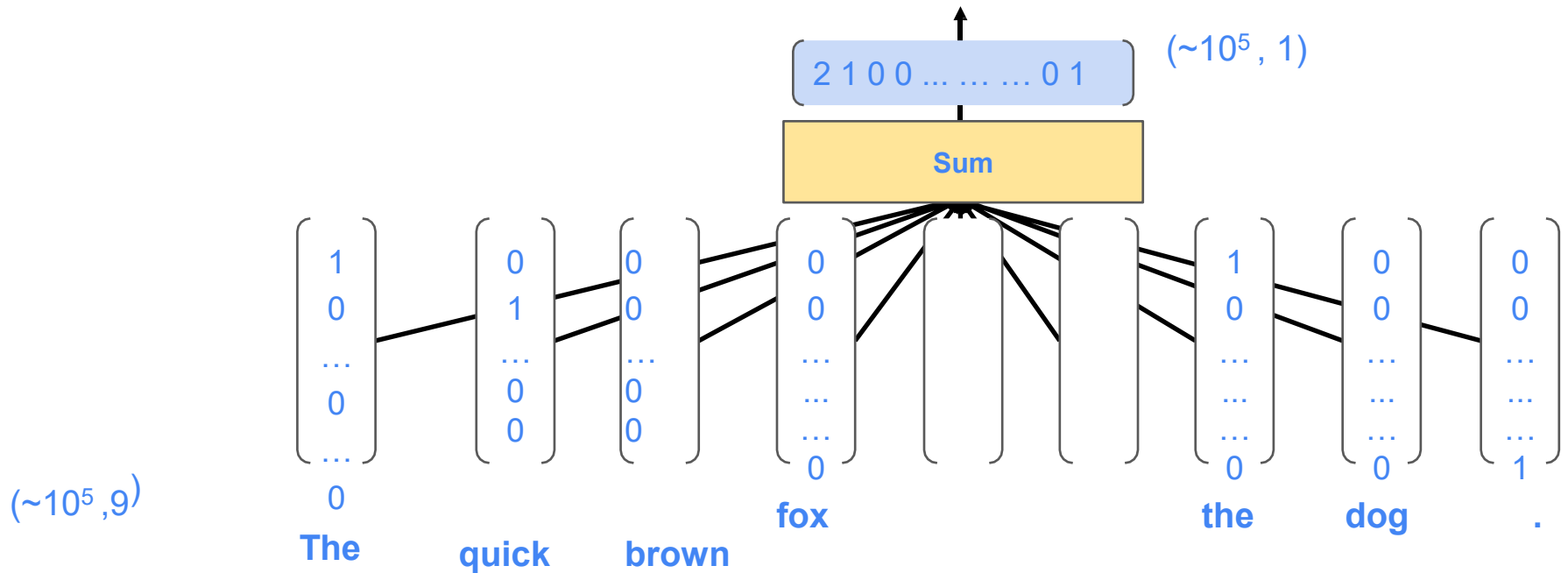
.

# Sparse text representation: BOW

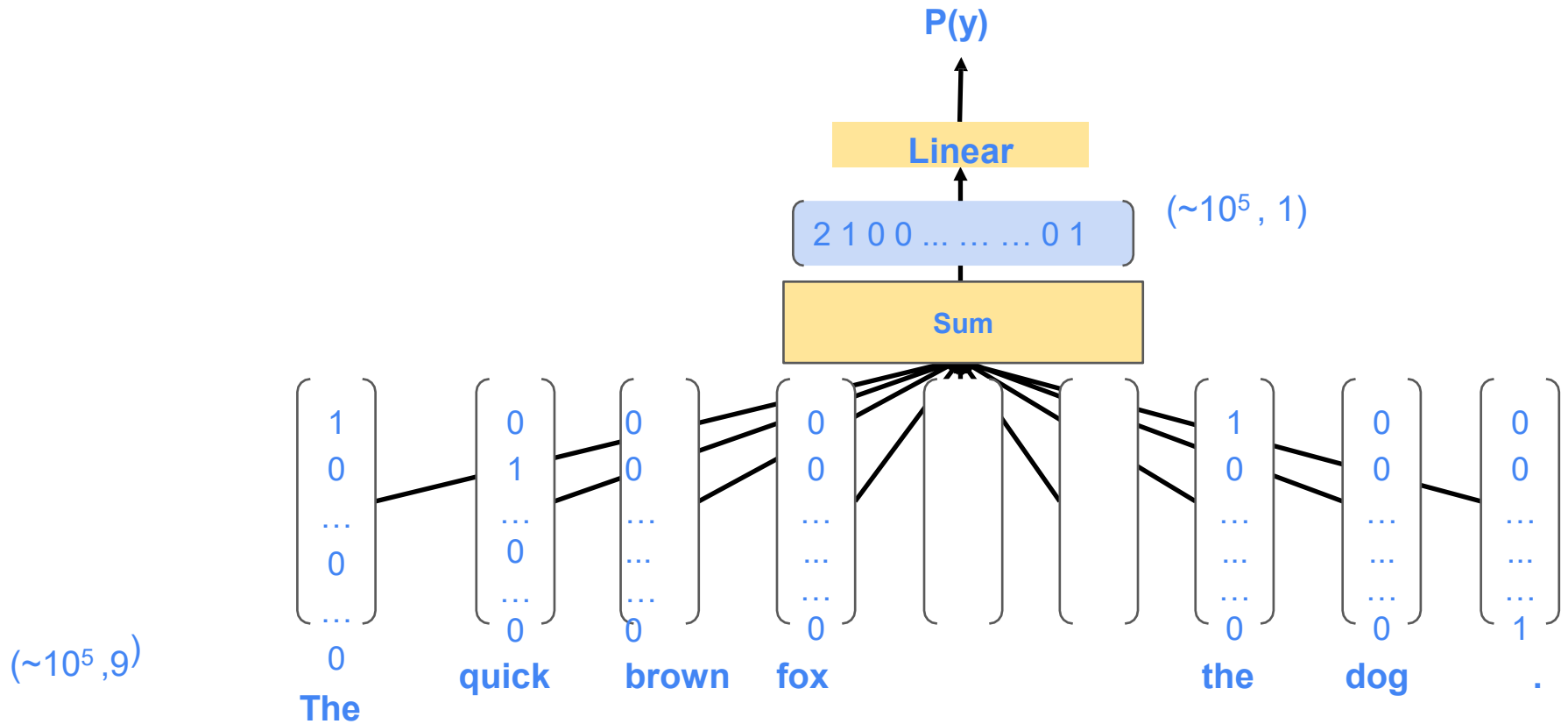
( $\sim 10^5, 9$ )

$\begin{pmatrix} 1 \\ 0 \\ \dots \\ 0 \\ \dots \end{pmatrix}$	$\begin{pmatrix} 0 \\ 1 \\ \dots \\ 0 \\ \dots \\ 0 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \\ \dots \\ \dots \\ \dots \\ 0 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \\ \dots \\ \dots \\ \dots \\ 0 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \\ \dots \\ \dots \\ \dots \\ 0 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \\ \dots \\ \dots \\ \dots \\ 0 \end{pmatrix}$	$\begin{pmatrix} 1 \\ 0 \\ \dots \\ \dots \\ \dots \\ 0 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \\ \dots \\ \dots \\ \dots \\ 0 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \\ \dots \\ \dots \\ \dots \\ 1 \end{pmatrix}$
The	quick	brown	fox	jumps	over	the	dog	.

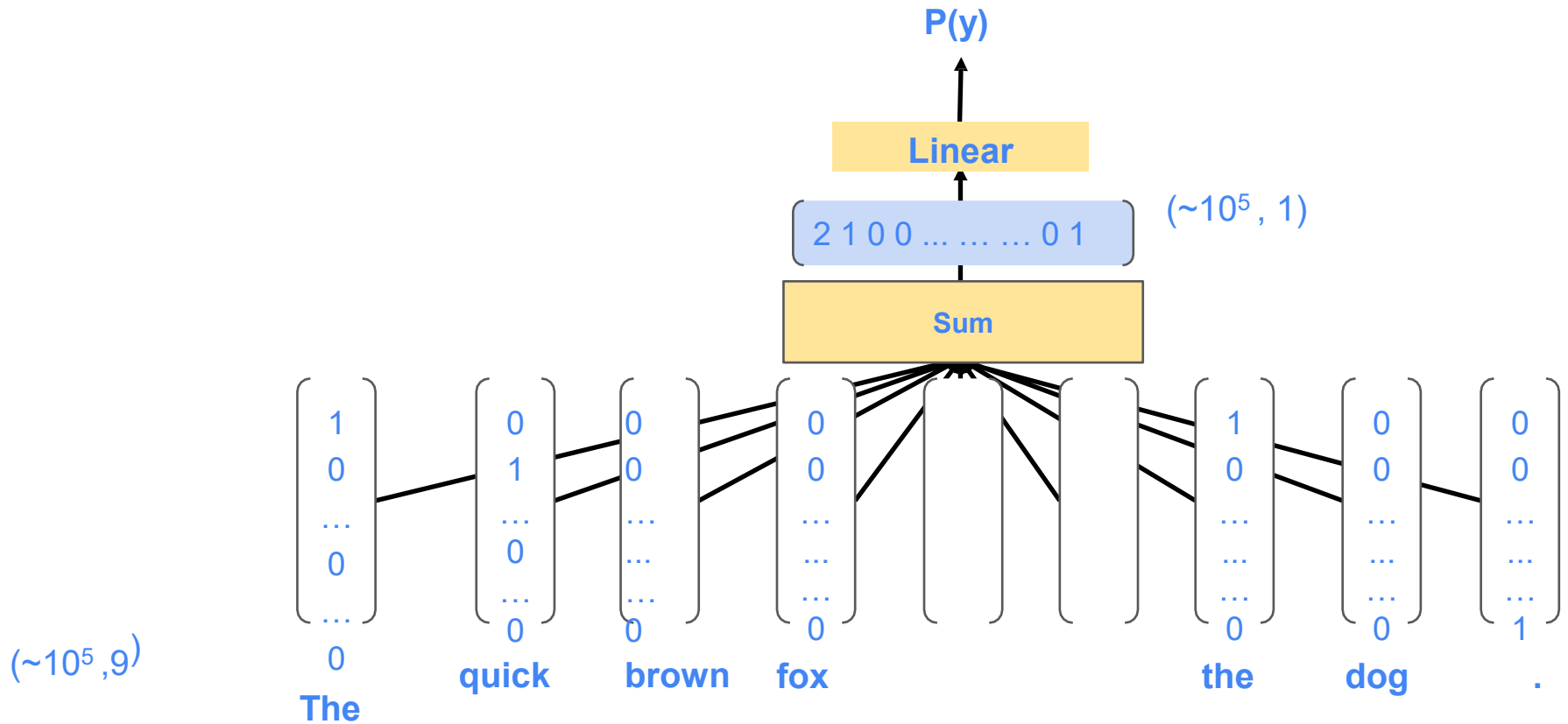
# Sparse text representation: BOW



# Sparse text representation: BOW



# Sparse text representation: BOW

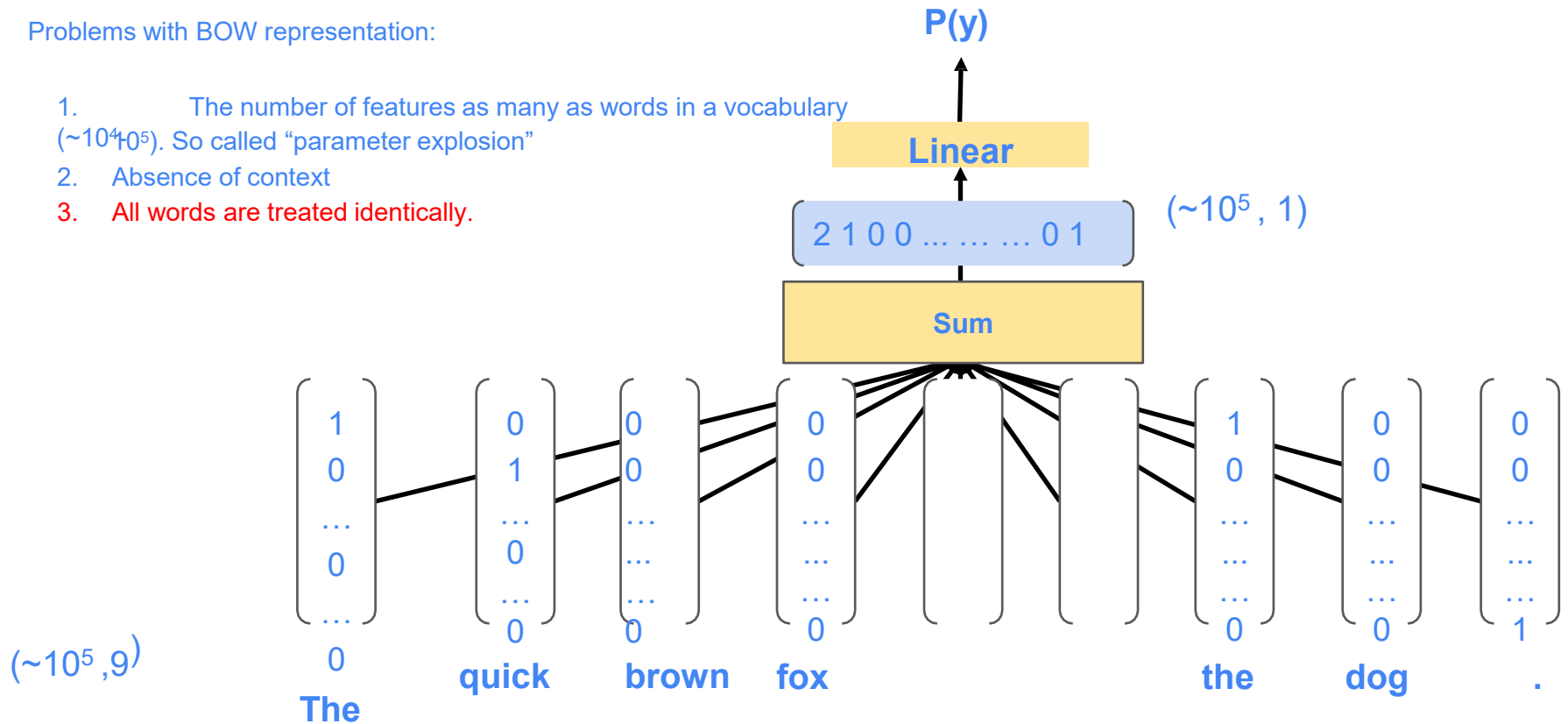




# Sparse text representation: BOW

Problems with BOW representation:

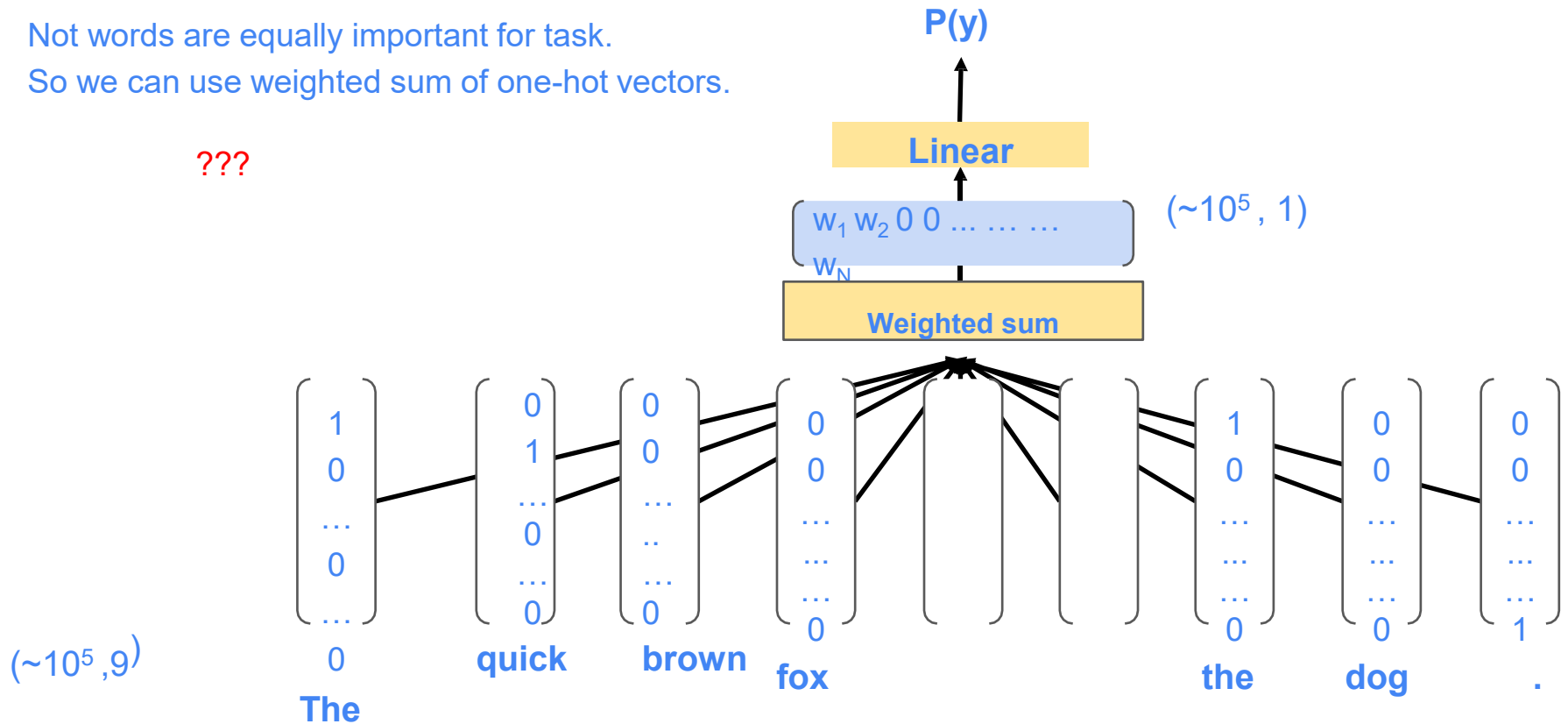
1. The number of features as many as words in a vocabulary ( $\sim 10^4$  to  $10^5$ ). So called “parameter explosion”
2. Absence of context
3. All words are treated identically.



# Weighting techniques for BOW

Not words are equally important for task.  
So we can use weighted sum of one-hot vectors.

???

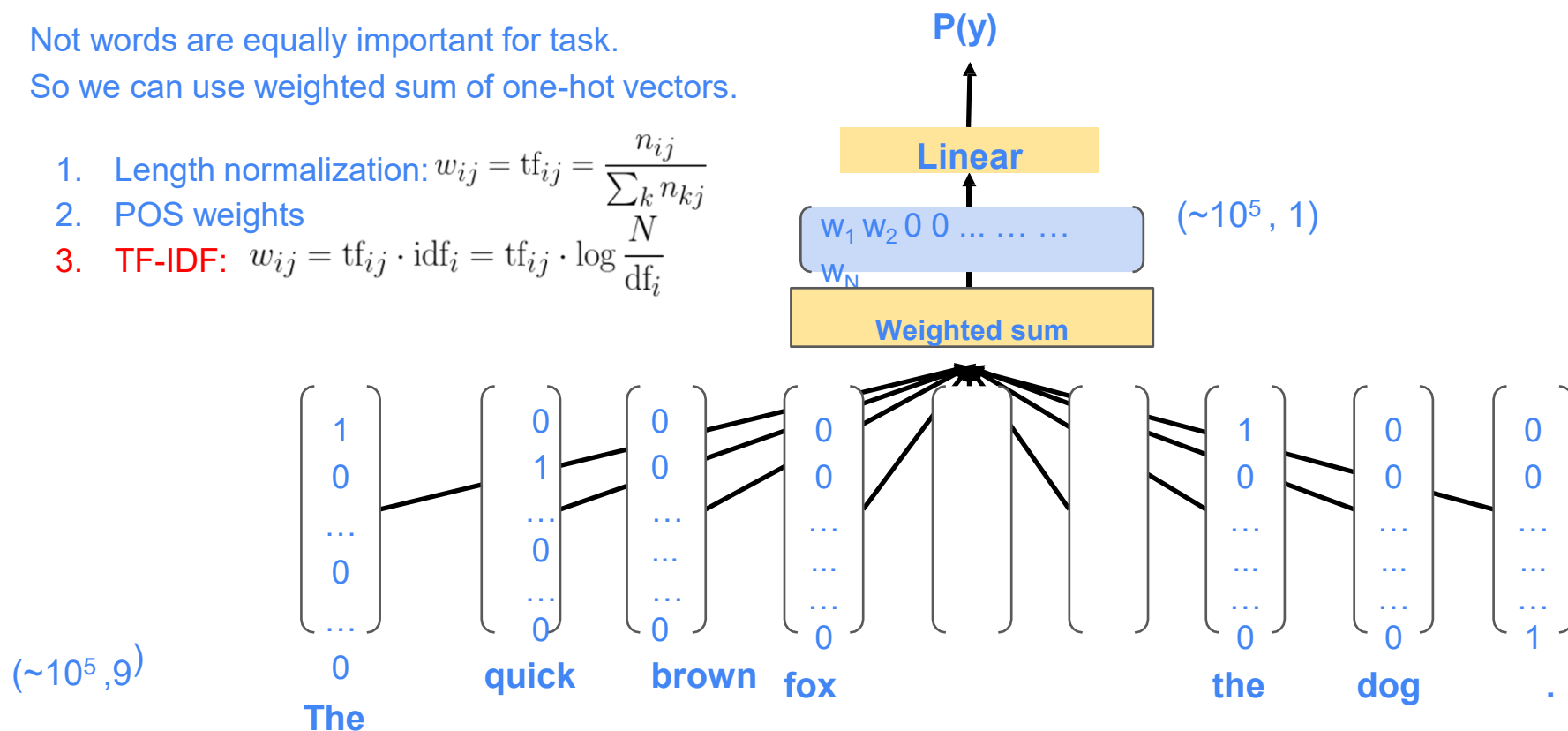


# Weighting techniques for BOW

Not words are equally important for task.

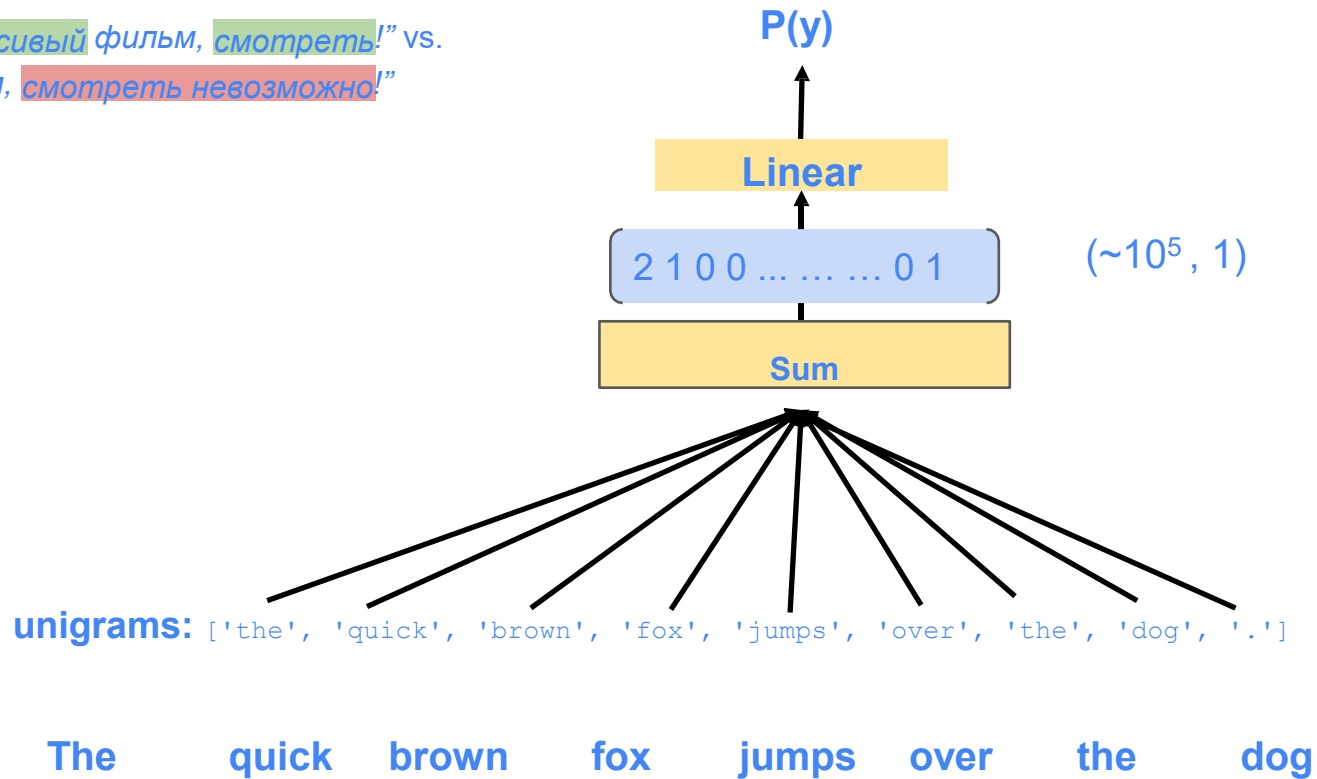
So we can use weighted sum of one-hot vectors.

1. Length normalization:  $w_{ij} = \text{tf}_{ij} = \frac{n_{ij}}{\sum_k n_{kj}}$
2. POS weights
3. **TF-IDF**:  $w_{ij} = \text{tf}_{ij} \cdot \text{idf}_i = \text{tf}_{ij} \cdot \log \frac{N}{\text{df}_i}$



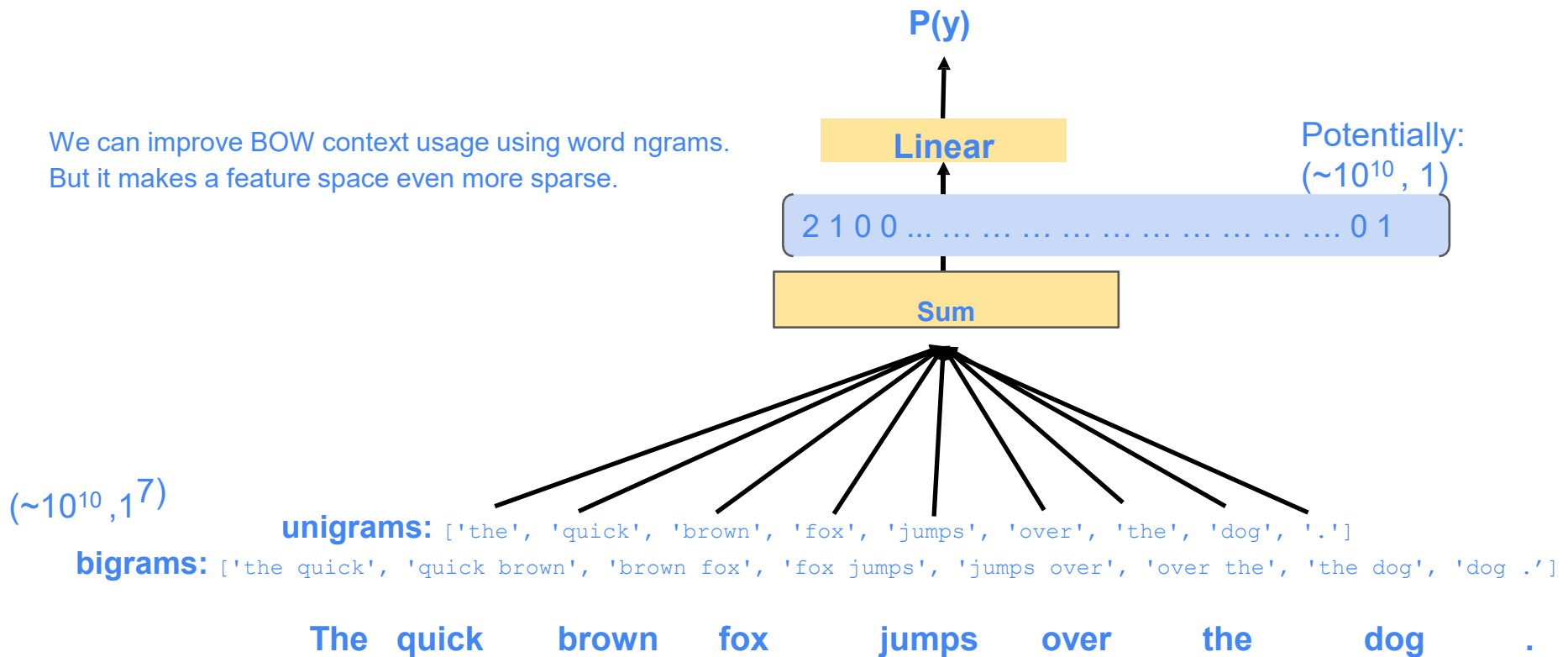
# Context importance

“Невозможно красивый фильм, смотреть!” vs.  
“Красивый фильм, смотреть невозможно!”



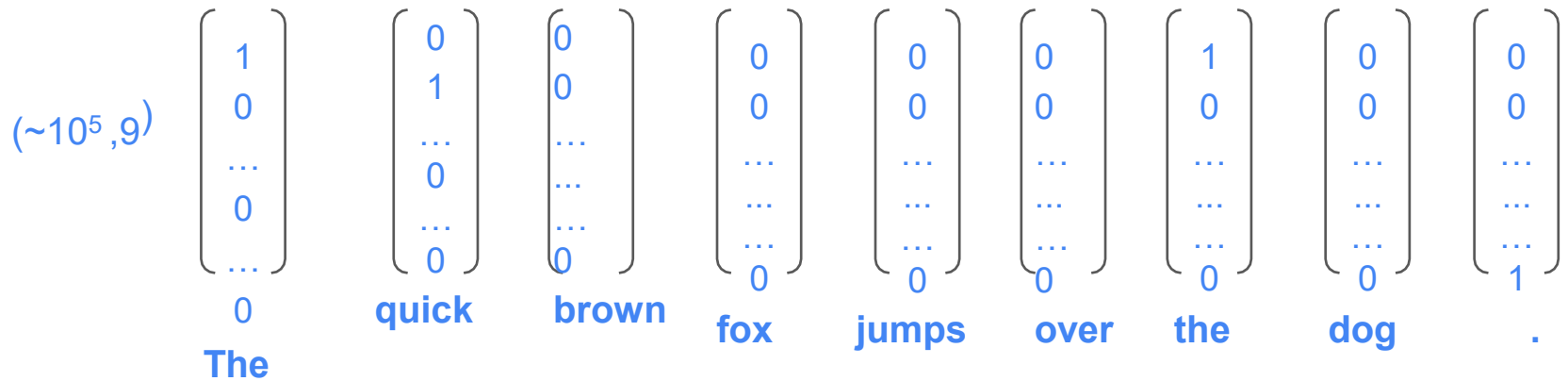
# Context importance

We can improve BOW context usage using word ngrams.  
But it makes a feature space even more sparse.



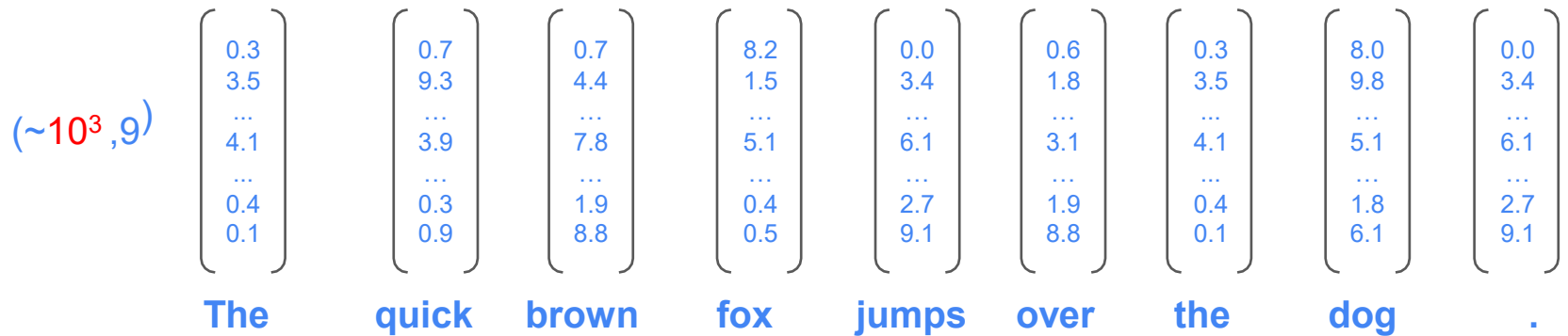
# Dense text representation: NBOW

Instead of sparse one-hot encoding we can use pre-trained word embeddings. It is so-called “neural bag-of- words”.



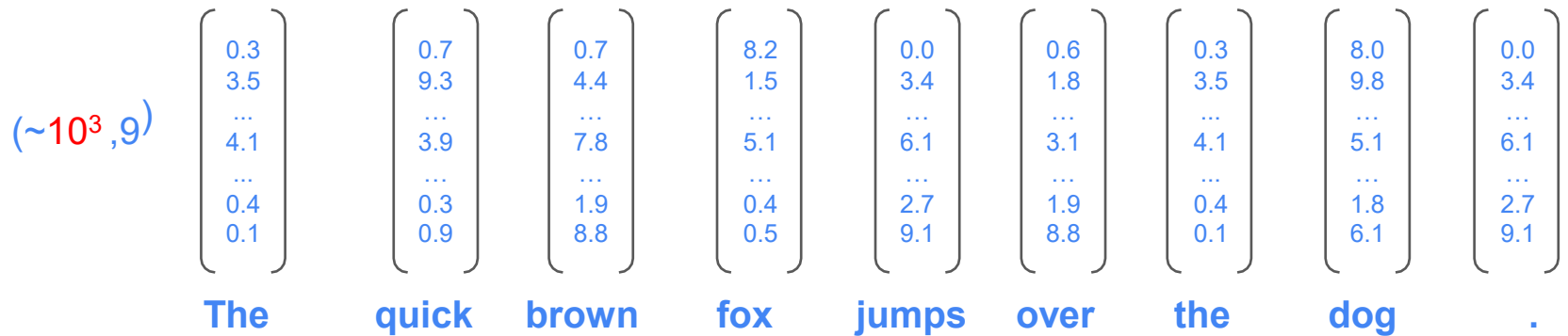
# Dense text representation: NBOW

Instead of sparse one-hot encoding we can use pre-trained word embeddings. It is so-called “neural bag-of-words”.



# Dense text representation: NBOW

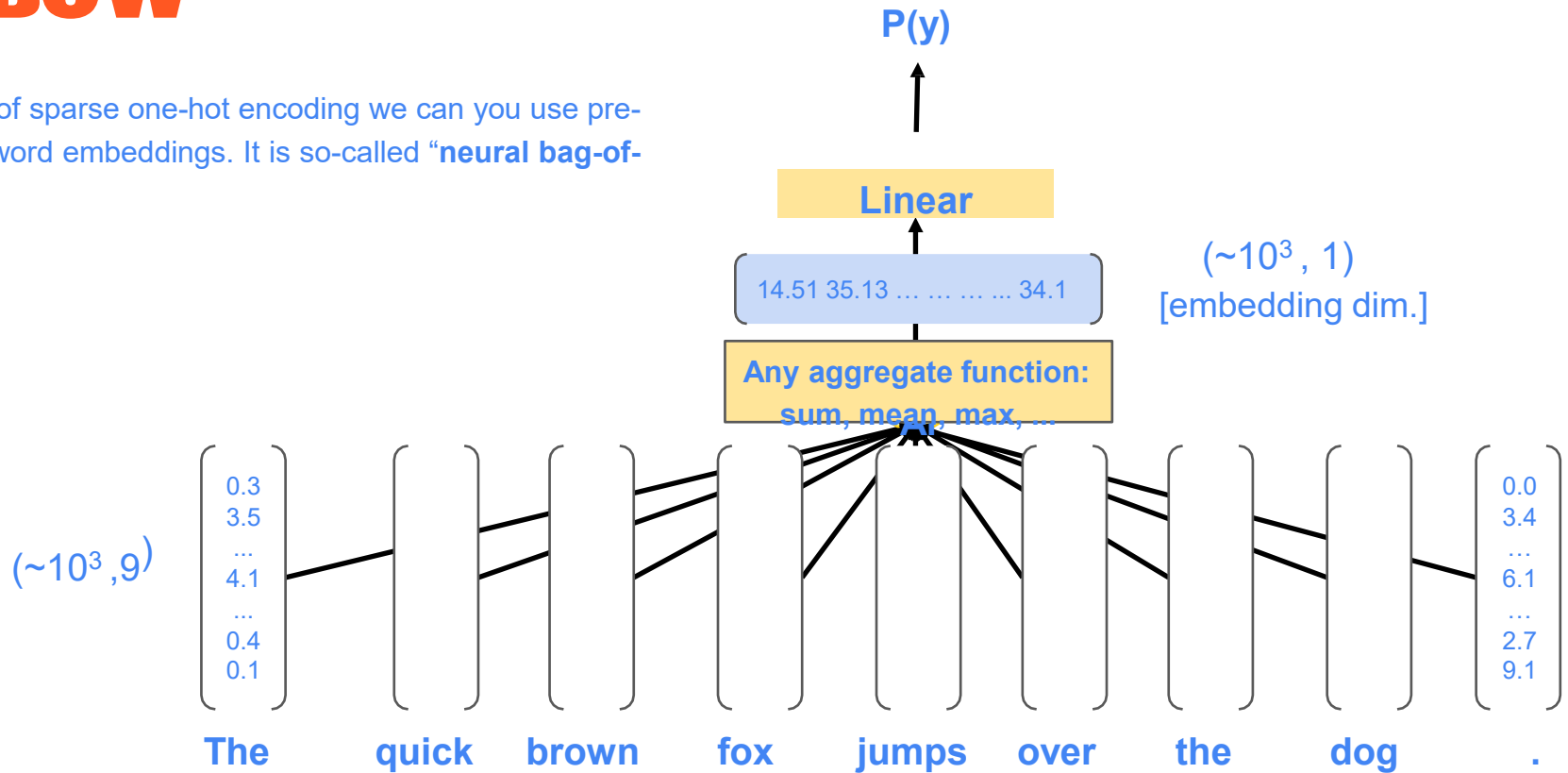
Instead of sparse one-hot encoding we can use pre-trained word embeddings. It is so-called “neural bag-of-words”.





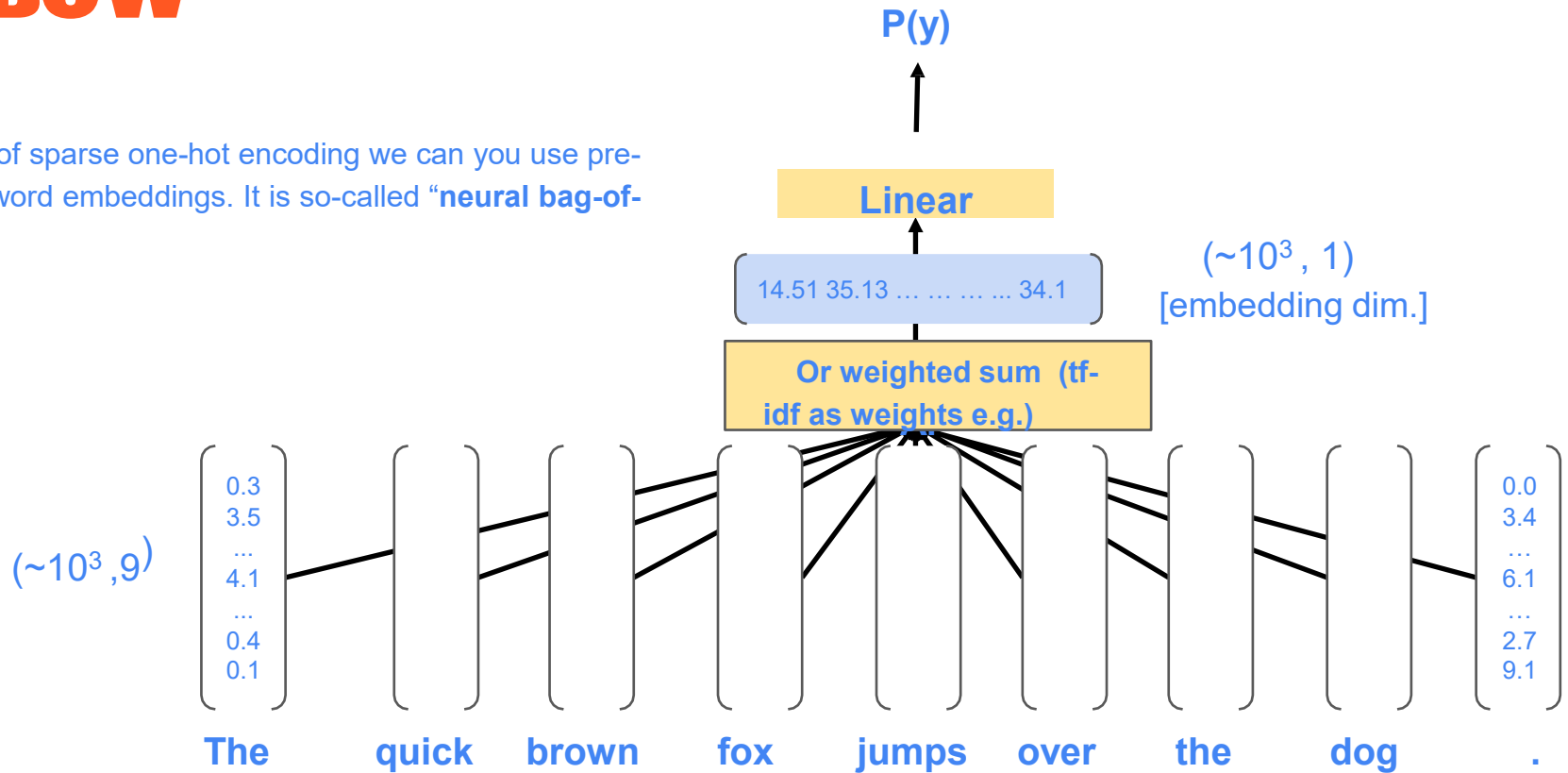
# Dense text representation: NBOW

Instead of sparse one-hot encoding we can use pre-trained word embeddings. It is so-called “**neural bag-of-words**”.



# Dense text representation: NBOW

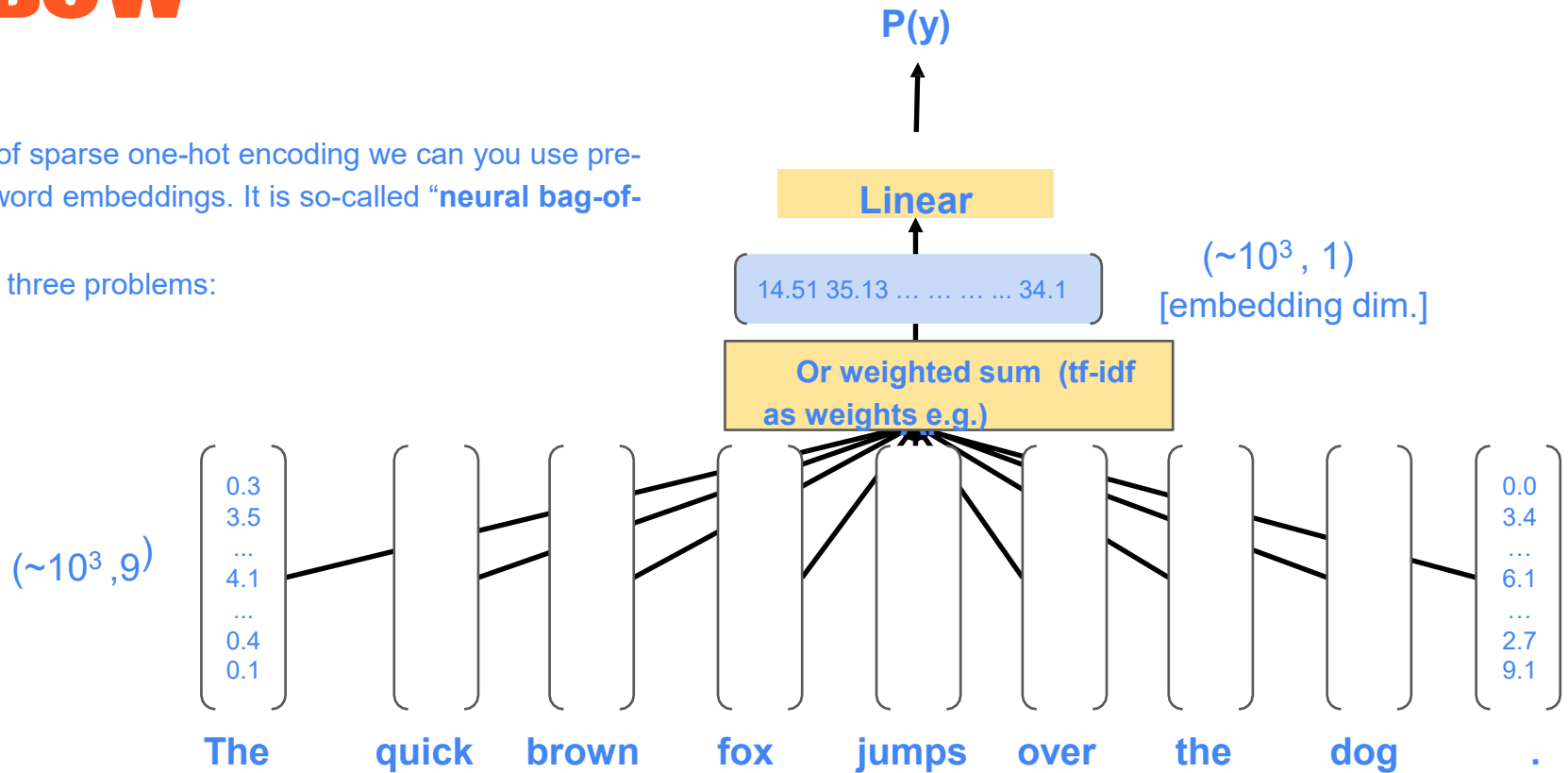
Instead of sparse one-hot encoding we can use pre-trained word embeddings. It is so-called “**neural bag-of-words**”.



# Dense text representation: NBOW

Instead of sparse one-hot encoding we can use pre-trained word embeddings. It is so-called “**neural bag-of-words**”.

It solves three problems:

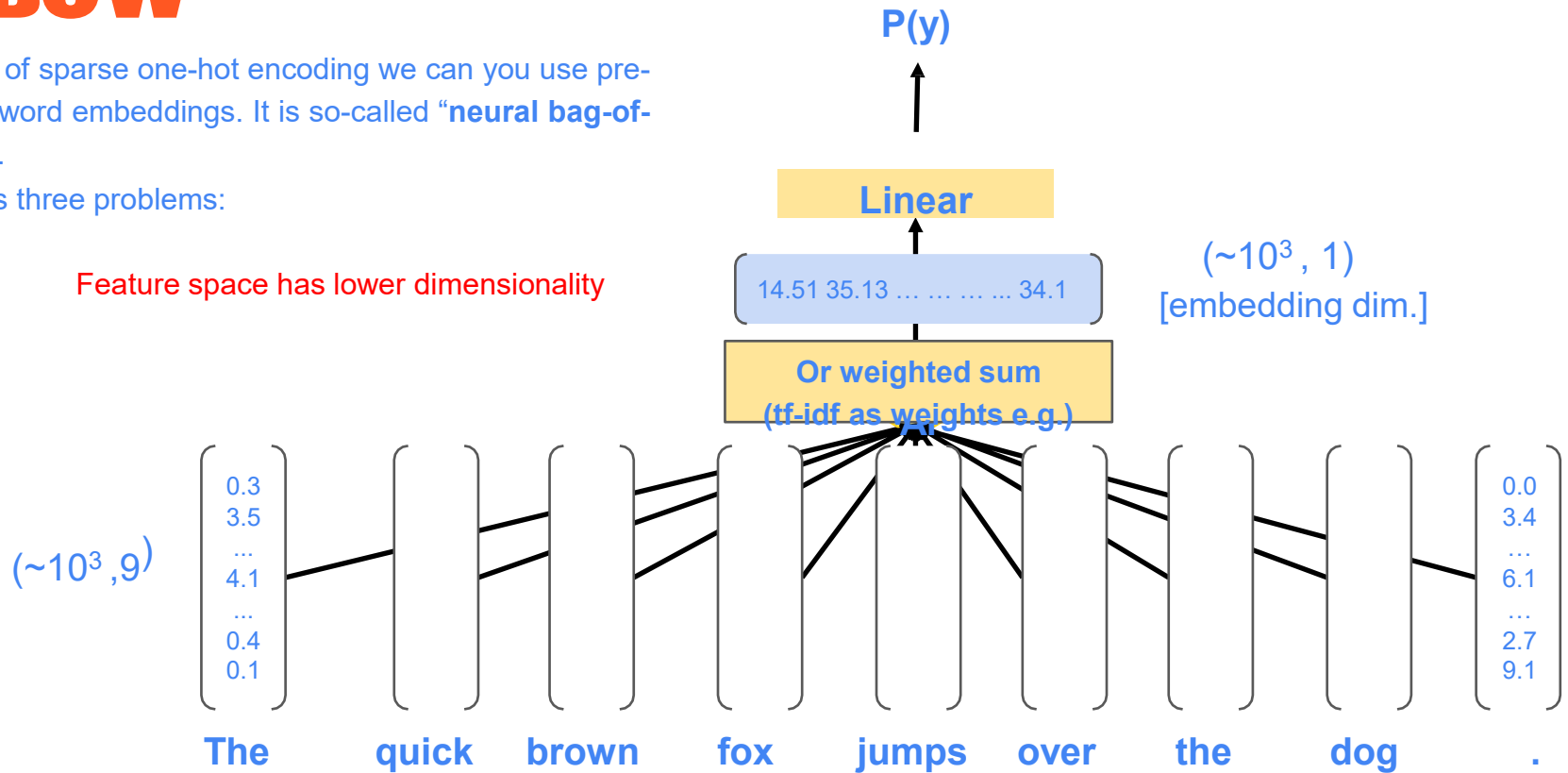


# Dense text representation: NBOW

Instead of sparse one-hot encoding we can use pre-trained word embeddings. It is so-called “**neural bag-of-words**”.

It solves three problems:

1. Feature space has lower dimensionality

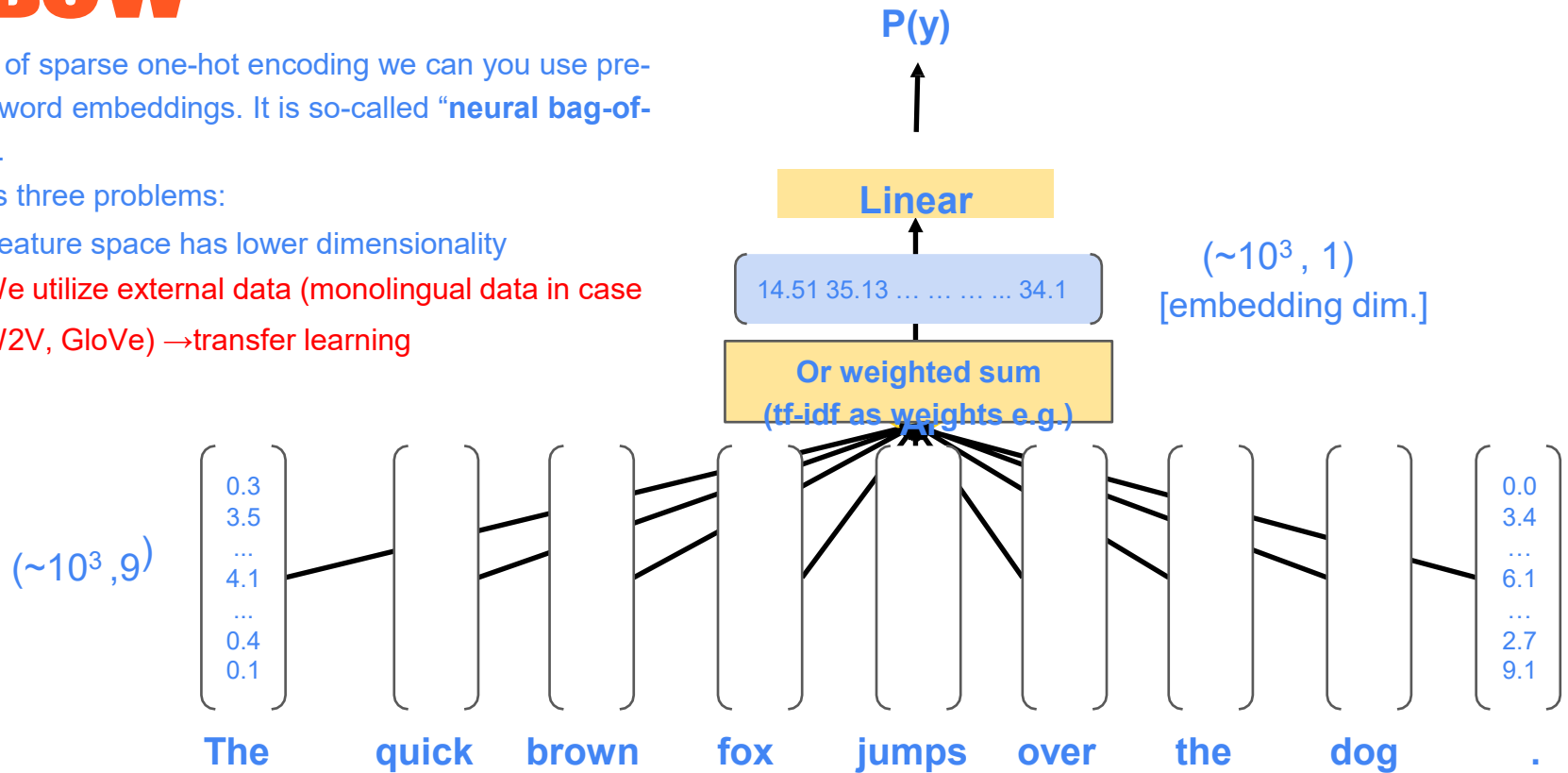


# Dense text representation: NBOW

Instead of sparse one-hot encoding we can use pre-trained word embeddings. It is so-called “**neural bag-of-words**”.

It solves three problems:

1. Feature space has lower dimensionality
2. We utilize external data (monolingual data in case W2V, GloVe) → transfer learning

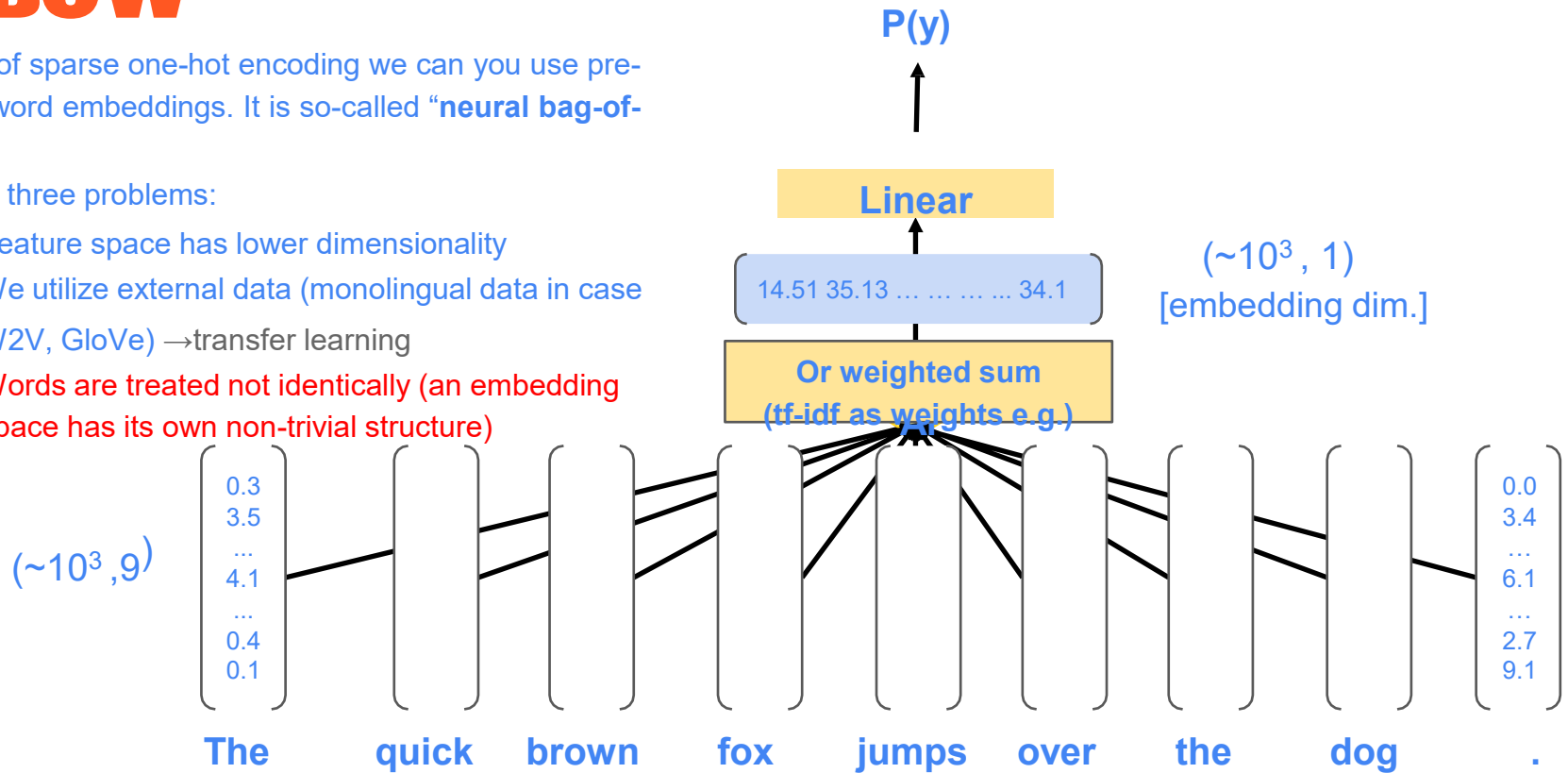


# Dense text representation: NBOW

Instead of sparse one-hot encoding we can use pre-trained word embeddings. It is so-called “**neural bag-of-words**”.

It solves three problems:

1. Feature space has lower dimensionality
2. We utilize external data (monolingual data in case W2V, GloVe) → transfer learning
3. Words are treated not identically (an embedding space has its own non-trivial structure)

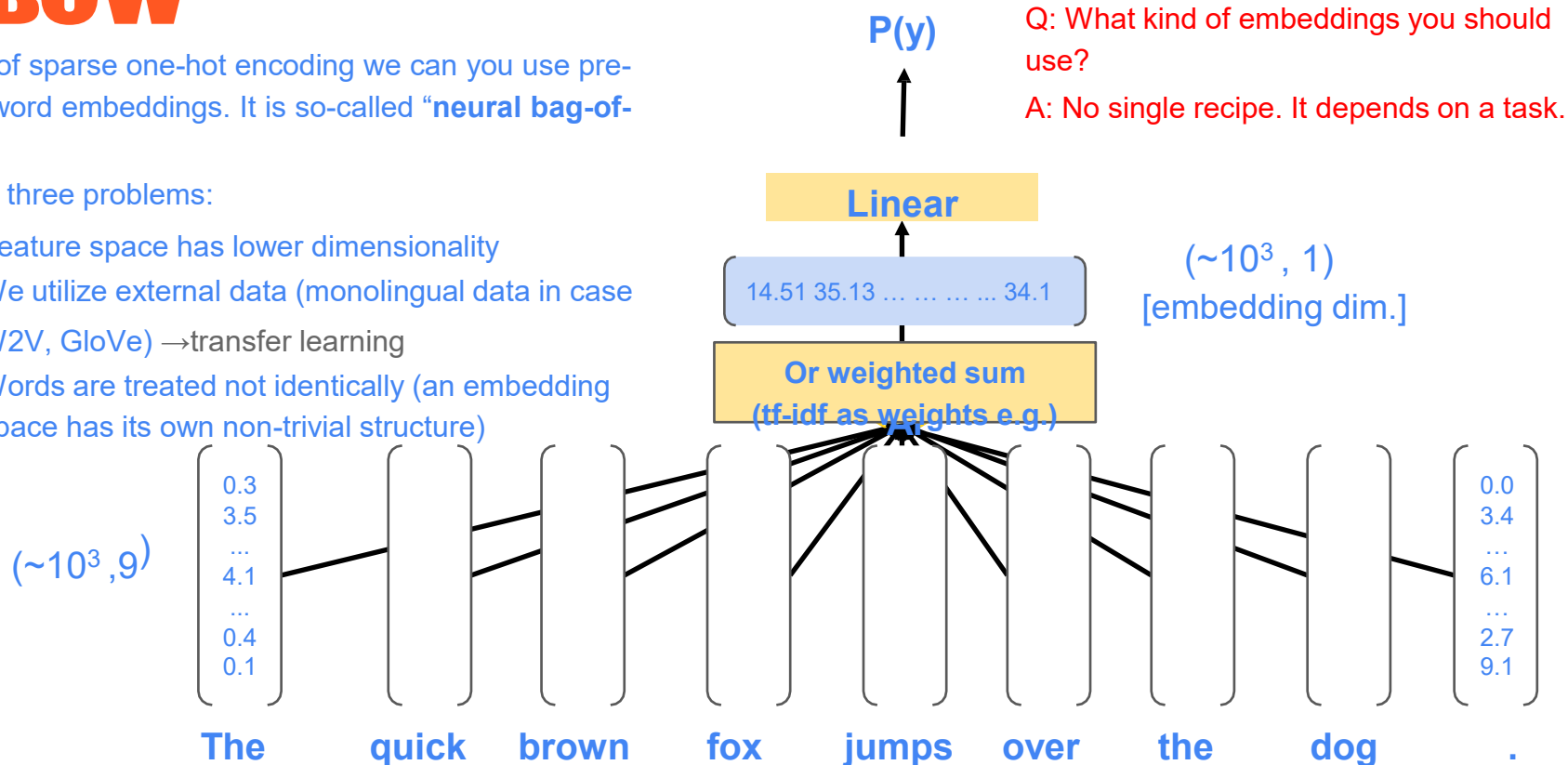


# Dense text representation: NBOW

Instead of sparse one-hot encoding we can use pre-trained word embeddings. It is so-called “**neural bag-of-words**”.

It solves three problems:

1. Feature space has lower dimensionality
2. We utilize external data (monolingual data in case W2V, GloVe) → transfer learning
3. Words are treated not identically (an embedding space has its own non-trivial structure)



# Dense text representation: NBOW

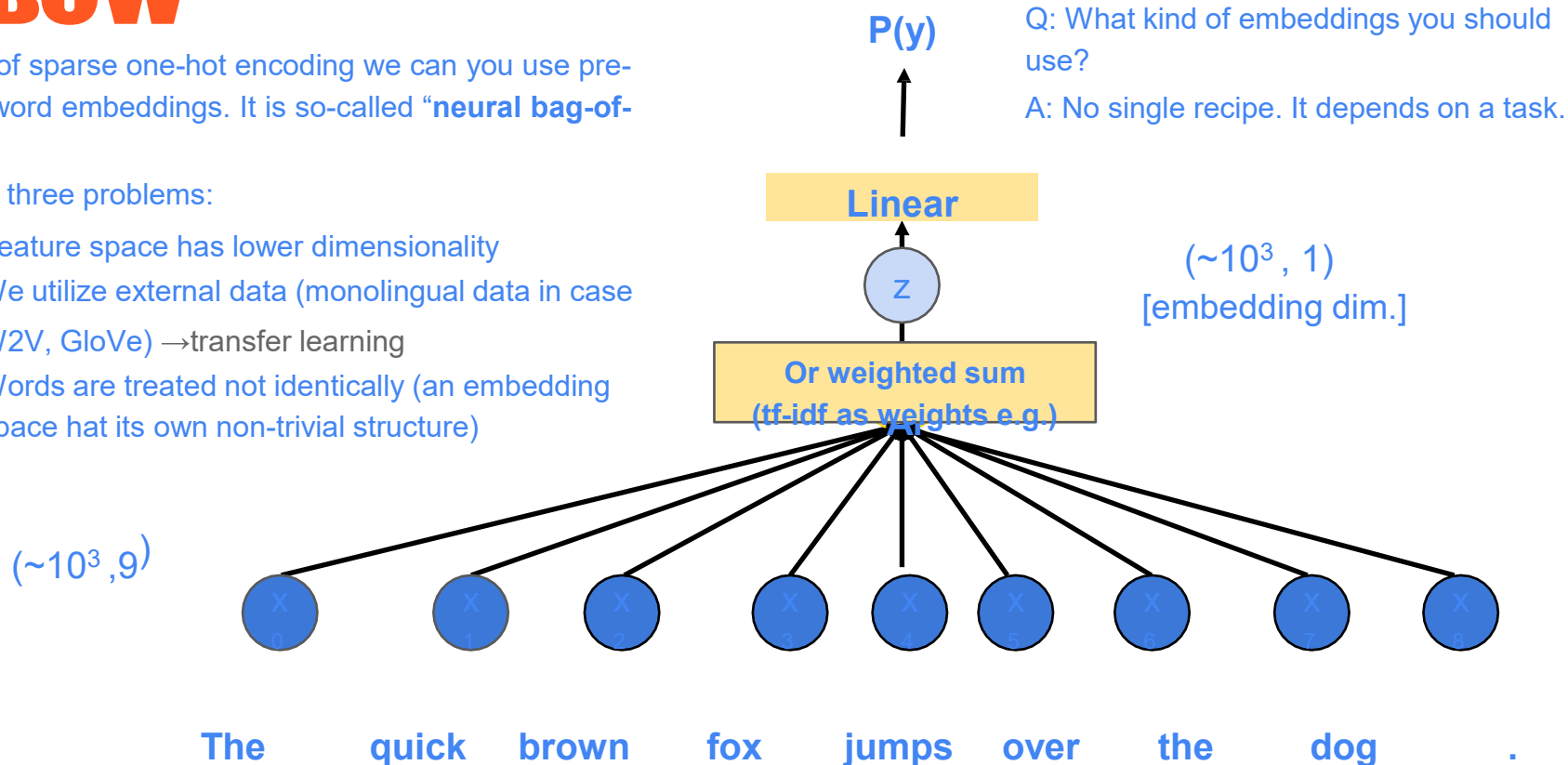
Instead of sparse one-hot encoding we can use pre-trained word embeddings. It is so-called “**neural bag-of-words**”.

It solves three problems:

1. Feature space has lower dimensionality
2. We utilize external data (monolingual data in case W2V, GloVe) → transfer learning
3. Words are treated not identically (an embedding space has its own non-trivial structure)

Q: What kind of embeddings you should use?

A: No single recipe. It depends on a task.





# BOW and NBOW: the shared problems

1. The importance weights for the word vectors aren't defined fully.
2. The only way to use context for these models is to utilize word ngrams.



Hmm...

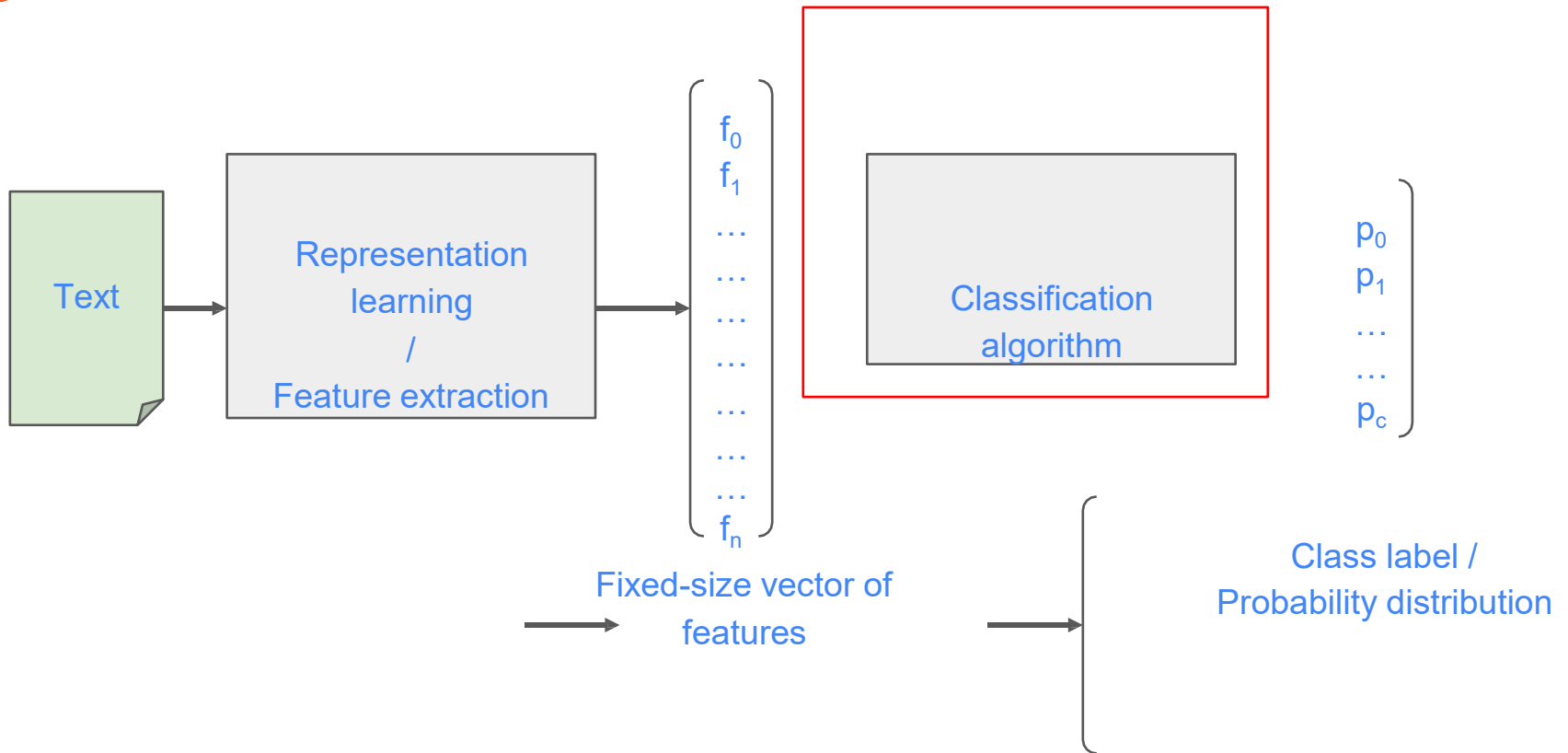
# BOW and NBOW: the shared problems

1. The importance weights for the word vectors aren't defined fully.
2. The only way to use context for these models is to utilize word ngrams.

We can use a learnable aggregation function to overcome the difficulties.  
The learnable function is a neural network (the universal approximator)



# Text classification in general



# Classification Algorithms

Let's explore how to use Naïves bayes to solve this task

# Naïve Bayes Intuition

Simple (“naïve”) classification method based on Bayes rule

Relies on very simple representation of document

Bag of words

# The bag of words representation

$\gamma$ (

I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet.

) = c



# The bag of words representation

$\gamma$ (

I **love** this movie! It's **sweet**, but with **satirical** humor. The dialogue is **great** and the adventure scenes are **fun**... It manages to be **whimsical** and **romantic** while **laughing** at the conventions of the fairy tale genre. I would **recommend** it to just about anyone. I've seen it **several** times, and I'm always **happy** to see it **again** whenever I have a friend who hasn't seen it yet.

) = C



## The bag of words representation: using a subset of words

$\gamma$ (

```
x love xxxxxxxxxxxxxxxxxxxx sweet
xxxxxxxx satirical xxxxxxxxxxxx
xxxxxxxxxxxx great xxxxxxxx
xxxxxxxxxxxxxxxxxxxxxxxxxxxx fun xxxx
xxxxxxxxxxxxxxxx whimsical xxxx
romantic xxxx laughing
xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx
xxxxxxxxxxxxxxxx recommend xxxxxx
xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx
xx several xxxxxxxxxxxxxxxxxxxx
xxxxx happy xxxxxxxx again
xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx
xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx
```

) = C







## The bag of words representation

$\gamma$  (

great	2
love	2
recommend	1
laugh	1
happy	1
...	...

) = C

# Bayes' Rule Applied to Documents and Classes

- For a document  $d$  and a class  $c$

$$P(c | d) = \frac{P(d | c)P(c)}{P(d)}$$

# Naïve Bayes Classifier (I)

$$c_{MAP} = \operatorname{argmax}_{c \in C} P(c | d)$$

MAP is “maximum a posteriori” = most likely class

$$= \operatorname{argmax}_{c \in C} \frac{P(d | c)P(c)}{P(d)}$$

Bayes Rule

$$= \operatorname{argmax}_{c \in C} P(d | c)P(c)$$

Dropping the denominator

Likelihood

Prior probability

# Naïve Bayes Classifier (II)

$$c_{MAP} = \operatorname{argmax}_{c \in C} P(d \mid c)P(c)$$

$$= \operatorname{argmax}_{c \in C} P(x_1, x_2, \dots, x_n \mid c)P(c)$$

Document  $d$   
represented as  
features  $x_1..x_n$

# Naïve Bayes Classifier (III)

$$c_{MAP} = \operatorname{argmax}_{c \in C} P(x_1, x_2, \dots, x_n | c) P(c)$$

$O(|X|^n \cdot |C|)$  parameters

How often does this class occur?

Could only be estimated if a very, very large number of training examples was available.

We can just count the relative frequencies in a corpus

# Multinomial Naïve Bayes Independence Assumptions

$$P(x_1, x_2, \dots, x_n \mid c)$$

**Bag of Words assumption:** Assume position doesn't matter

**Conditional Independence:** Assume the feature probabilities  $P(x_i \mid c_j)$  are independent given the class  $c$ .

$$P(x_1, \dots, x_n \mid c) = P(x_1 \mid c) \bullet P(x_2 \mid c) \bullet P(x_3 \mid c) \bullet \dots \bullet P(x_n \mid c)$$

# Multinomial Naïve Bayes Classifier

$$c_{MAP} = \operatorname{argmax}_{c \in C} P(x_1, x_2, \dots, x_n | c) P(c)$$

$$c_{NB} = \operatorname{argmax}_{c \in C} P(c_j) \prod_{x \in X} P(x | c)$$

# Applying Multinomial Naive Bayes Classifiers to Text Classification

positions  $\leftarrow$  all word positions in test document

$$c_{NB} = \operatorname{argmax}_{c_j \in C} P(c_j) \prod_{i \in \text{positions}} P(x_i | c_j)$$



# Problems with multiplying lots of probs

There's a problem with this:

$$c_{NB} = \operatorname{argmax}_{c_j \in C} P(c_j) \prod_{i \in \text{positions}} P(x_i | c_j)$$

Multiplying lots of probabilities can result in floating-point underflow!

.0006 \* .0007 \* .0009 \* .01 \* .5 \* .000008....

Idea: Use logs, because  $\log(ab) = \log(a) + \log(b)$

We'll sum logs of probabilities instead of multiplying probabilities!

# We actually do everything in log space

Instead of this:  $c_{NB} = \operatorname{argmax}_{c_j \in C} P(c_j) \prod_{i \in \text{positions}} P(x_i | c_j)$

This:  $c_{NB} = \operatorname{argmax}_{c_j \in C} \left[ \log P(c_j) + \sum_{i \in \text{positions}} \log P(x_i | c_j) \right]$

Notes:

1) Taking log doesn't change the ranking of classes!

The class with highest probability also has highest log probability!

2) It's a linear model:

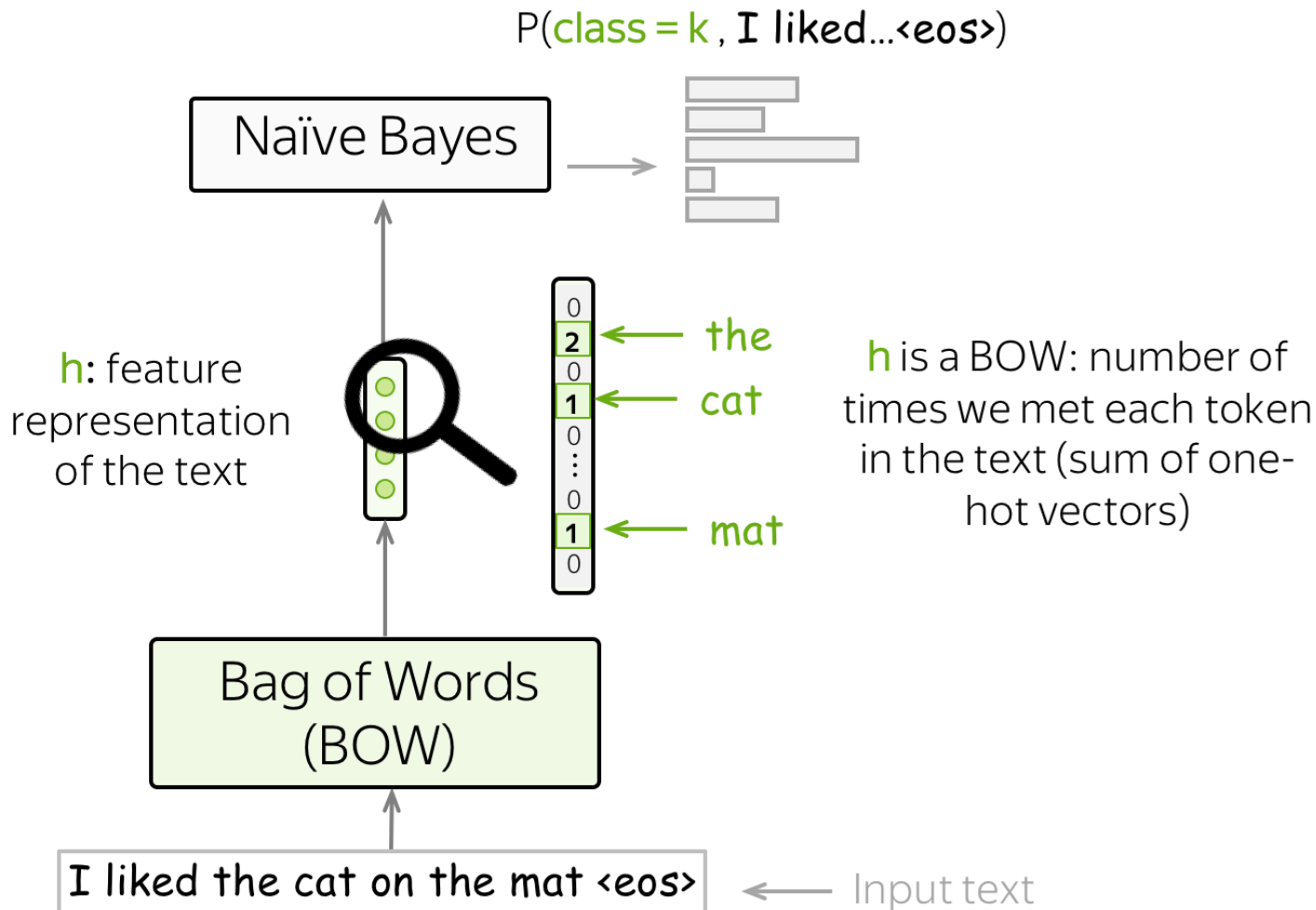
Just a max of a sum of weights: a **linear** function of the inputs

So naive bayes is a **linear classifier**

## Laplace (add-1) smoothing for Naïve Bayes

$$\begin{aligned}\hat{P}(w_i | c) &= \frac{\text{count}(w_i, c) + 1}{\sum_{w \in V} (\text{count}(w, c) + 1)} \\ &= \frac{\text{count}(w_i, c) + 1}{\left( \sum_{w \in V} \text{count}(w, c) \right) + |V|}\end{aligned}$$

# Final Framework



# A Worked Example

$$\hat{P}(c) = \frac{N_c}{N}$$

$$\hat{P}(w|c) = \frac{\text{count}(w,c)+1}{\text{count}(c)+|V|}$$

	Doc	Words	Class
Training	1	Chinese Beijing Chinese	c
	2	Chinese Chinese Shanghai	c
	3	Chinese Macao	c
	4	Tokyo Japan Chinese	j
Test	5	Chinese Chinese Chinese Tokyo Japan	?

**Priors:**

$$P(c) = \frac{3}{4}$$

$$P(j) = \frac{1}{4}$$

**Choosing a class:**

$$P(c|d5) \propto \frac{3}{4} * \left(\frac{3}{7}\right)^3 * \frac{1}{14} * \frac{1}{14} \approx 0.0003$$

**Conditional Probabilities:**

$$P(\text{Chinese}|c) = \frac{(5+1)}{(8+6)} = \frac{6}{14} = \frac{3}{7}$$

$$P(\text{Tokyo}|c) = \frac{(0+1)}{(8+6)} = \frac{1}{14}$$

$$P(\text{Japan}|c) = \frac{(0+1)}{(8+6)} = \frac{1}{14}$$

$$P(\text{Chinese}|j) = \frac{(1+1)}{(3+6)} = \frac{2}{9}$$

$$P(\text{Tokyo}|j) = \frac{(1+1)}{(3+6)} = \frac{2}{9}$$

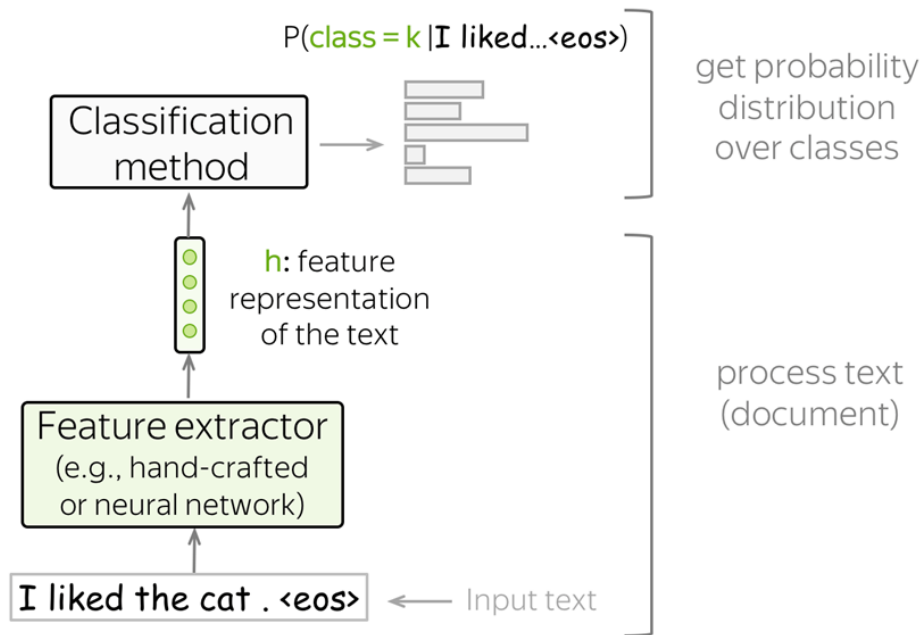
$$P(\text{Japan}|j) = \frac{(1+1)}{(3+6)} = \frac{2}{9}$$

$$P(j|d5) \propto \frac{1}{4} * \left(\frac{2}{9}\right)^3 * \frac{2}{9} * \frac{2}{9} \approx 0.0001$$

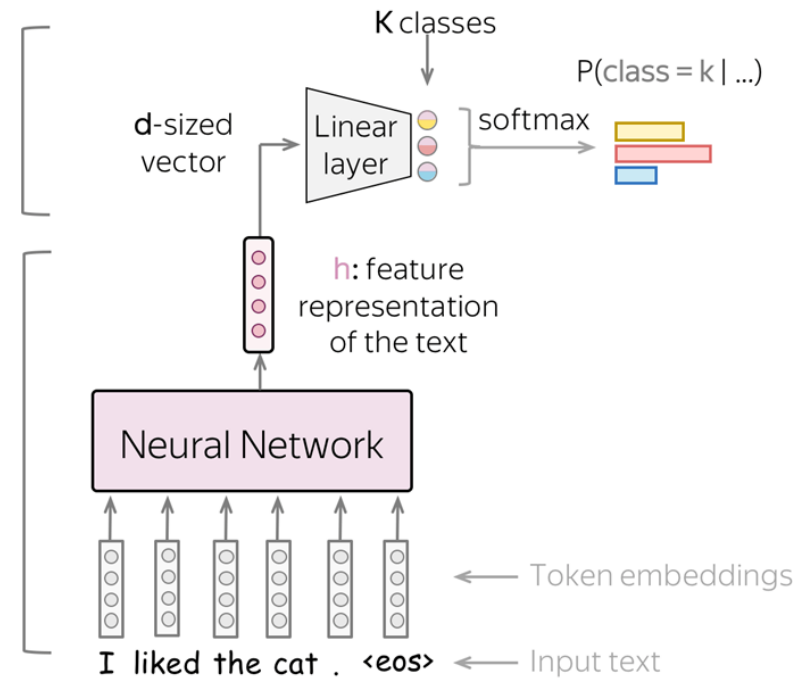
# Text Classification with Neural Networks

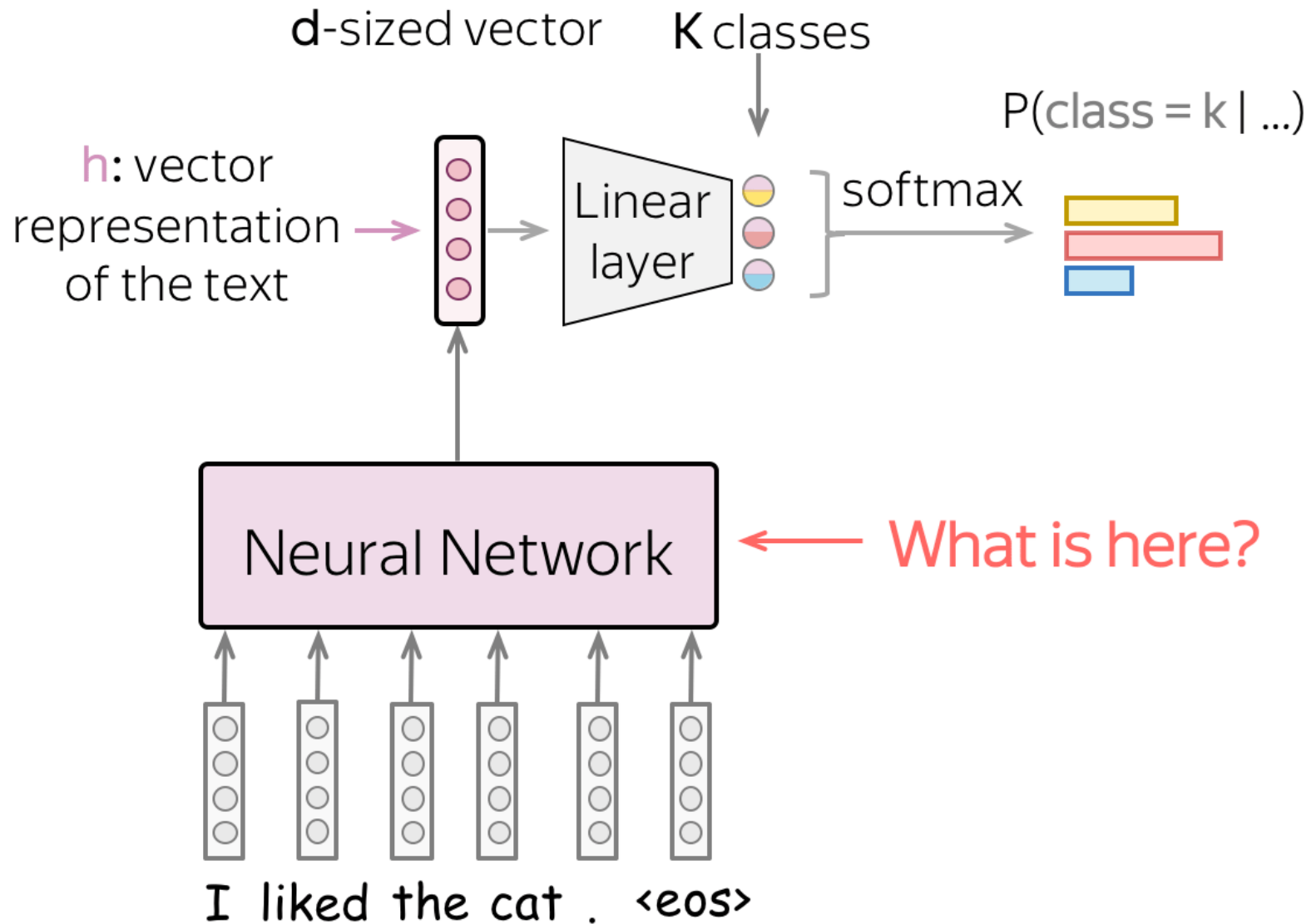
- Feature representation of the input text can be obtained using a neural network
- We feed the embeddings of the input tokens to a neural network, and this neural network gives us a vector representation of the input text.
- After that, this vector is used for classification.

## General Classification Pipeline



## Classification with Neural Networks







**Thank you**

