



Natural Language Processing: Classification

HSE Faculty of Computer Science
Machine Learning and Data-Intensive Systems

Murat Khazhgeriev



Table of Content

- **General view on classification with text**
- Pre-deep learning approaches
- Deep learning approaches
- Practical Tips

Extract features, apply model, compare the distribution to the target

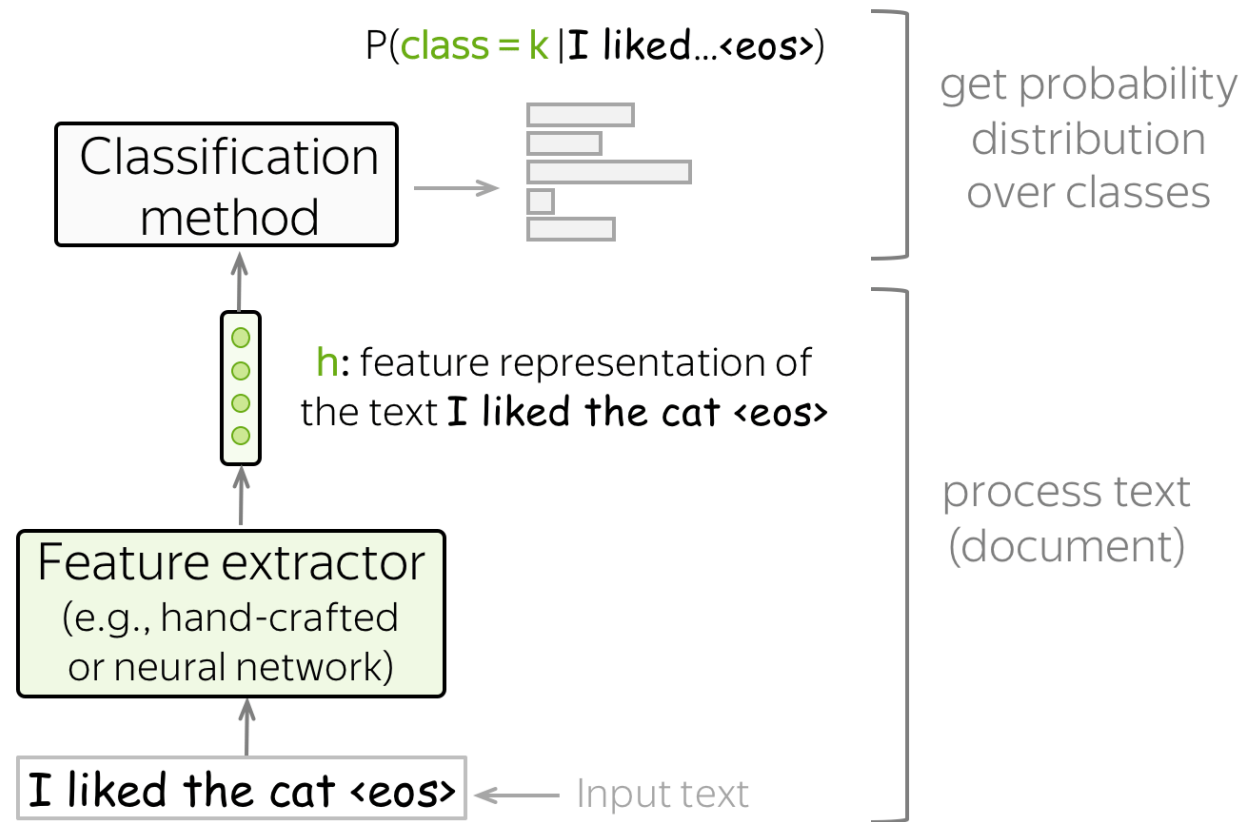




Table of Content

- General view on classification with text
- **Pre-deep learning approaches: Naïve Bayes**
- Deep learning approaches
- Practical Tips

Naively count the conditional probability

Bayes' rule
(hence Naïve Bayes)

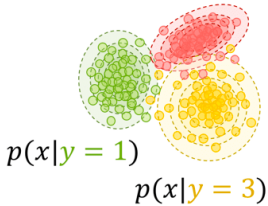
Ignore $P(x)$ – it does not
influence the argmax

$$y^* = \arg \max_k P(y = k|x) = \arg \max_k \frac{P(x|y = k) \cdot P(y = k)}{P(x)} = \arg \max_k \underbrace{P(x|y = k)}_{\text{need to define this}} \cdot \underbrace{P(y = k)}_{\text{need to define this}}$$

Naively count the conditional probability

The Naive Bayes assumptions are

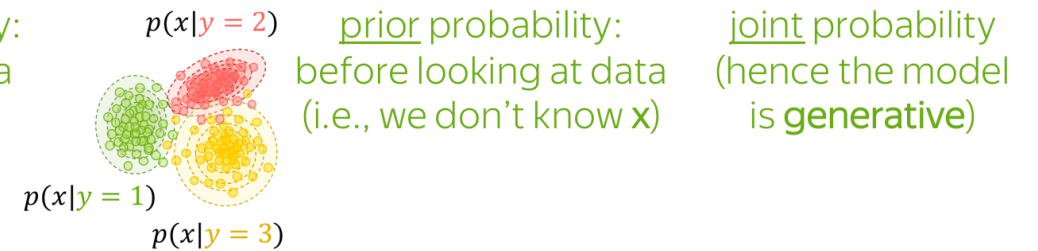
- Bag of Words assumption: word order does not matter
- Conditional Independence assumption: features (words) are independent given the class

$$y^* = \arg \max_k \underbrace{P(y = k|x)}_{\text{posterior probability: after looking at data (i.e., we know } \mathbf{x})} = \arg \max_k \underbrace{P(x|y = k)}_{\text{prior probability: before looking at data (i.e., we don't know } \mathbf{x})} \cdot \underbrace{P(y = k)}_{\text{joint probability (hence the model is generative)}} = \arg \max_k \underbrace{P(x, y = k)}$$


How to define $P(x|y=k)$ and $P(y=k)$

$$y^* = \arg \max_k \underbrace{P(y = k|x)}_{\text{posterior probability: after looking at data (i.e., we know } x)} = \arg \max_k \underbrace{P(x|y = k)}_{\text{prior probability: before looking at data (i.e., we don't know } x)} \cdot \underbrace{P(y = k)}_{\text{joint probability (hence the model is generative)}} = \arg \max_k \underbrace{P(x, y = k)}_{\text{joint probability (hence the model is generative)}}$$

$$P(y = k) = \frac{N(y = k)}{\sum_i N(y = i)}$$



$$P(x|y = k) = P(x_1, \dots, x_n|y = k) = \prod_{t=1}^n P(x_t|y = k)$$

$$P(x_i|y = k) = \frac{N(x_i, y = k)}{\sum_{t=1}^{|V|} N(x_t, y = k)}$$

How to define $P(x | y=k)$ and $P(y=k)$?

The Naive Bayes assumptions are

- Bag of Words assumption: word order does not matter,
- Conditional Independence assumption: features (words) are independent given the class.

$$P(y = k) = \frac{N(y = k)}{\sum_i N(y = i)}$$

$$P(x|y = k) = P(x_1, \dots, x_n|y = k) = \prod_{t=1}^n P(x_t|y = k)$$

What if $N(x_i, y = k) = 0$? Need to avoid this!

nulls out token prob. \Rightarrow nulls out document probability \Rightarrow Bad!

$$\underbrace{N(x_i, y = k)}_{\substack{\uparrow \\ \text{In training data, haven't seen} \\ \text{token } x_i \text{ in documents of class } k}} \Rightarrow P(x_i|y = k) = \frac{N(x_i, y = k)}{\sum_{t=1}^{|V|} N(x_t, y = k)} = 0 \Rightarrow P(x|y = k) = \prod_{i=1}^n P(x_i|y = k) = 0$$

$$P(x_i|y = k) = \frac{\delta + N(x_i, y = k)}{\sum_{t=1}^{|V|} (\delta + N(x_t, y = k))} = \frac{\delta + N(x_i, y = k)}{\delta \cdot |V| + \sum_{t=1}^{|V|} N(x_t, y = k)}$$

Making a prediction

Data: $x =$ This film is awesome !
 $x_1 \quad x_2 \quad x_3 \quad x_4 \quad x_5$

Compute joint probability of data and class

Positive class

$$\begin{aligned} P(x, y = +) &= P(y = +) \cdot P(x|y = +) \\ &= P(y = +) \cdot \\ &\quad \cdot P(\text{This}|y = +) \\ &\quad \cdot P(\text{film}|y = +) \\ &\quad \cdot P(\text{is}|y = +) \\ &\quad \cdot P(\text{awesome}|y = +) \\ &\quad \cdot P(!|y = +) \end{aligned}$$

Prior class probability (often 0.5)

Neutral words – not much difference
in probabilities for classes

This is where we expect the difference:

$$P(\text{awesome}|y = +) \gg P(\text{awesome}|y = -)$$

Negative class

$$\begin{aligned} P(x, y = -) &= P(y = -) \cdot P(x|y = -) \\ &= P(y = -) \cdot \\ &\quad \cdot P(\text{This}|y = -) \\ &\quad \cdot P(\text{film}|y = -) \\ &\quad \cdot P(\text{is}|y = -) \\ &\quad \cdot P(\text{awesome}|y = -) \\ &\quad \cdot P(!|y = -) \end{aligned}$$

$$P(x, y = +) > P(x, y = -) \Rightarrow y = +$$

View in the general framework

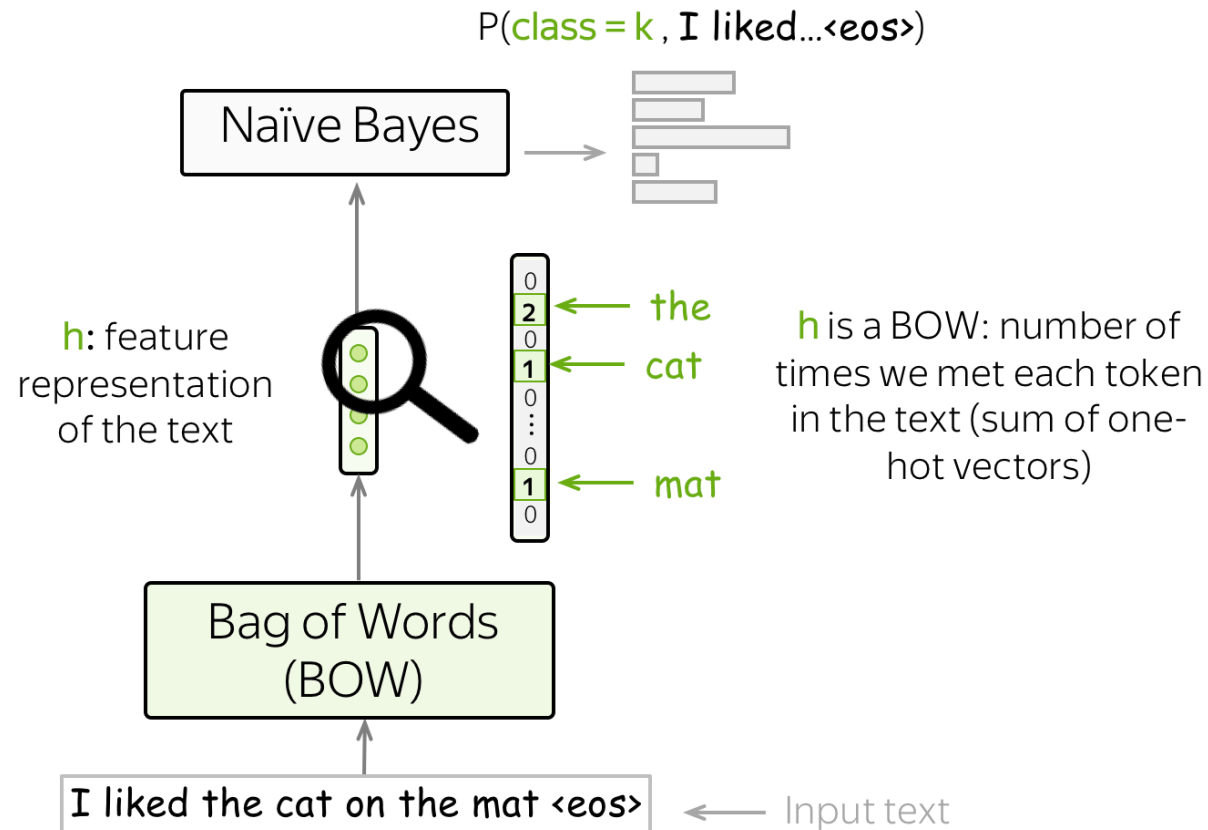
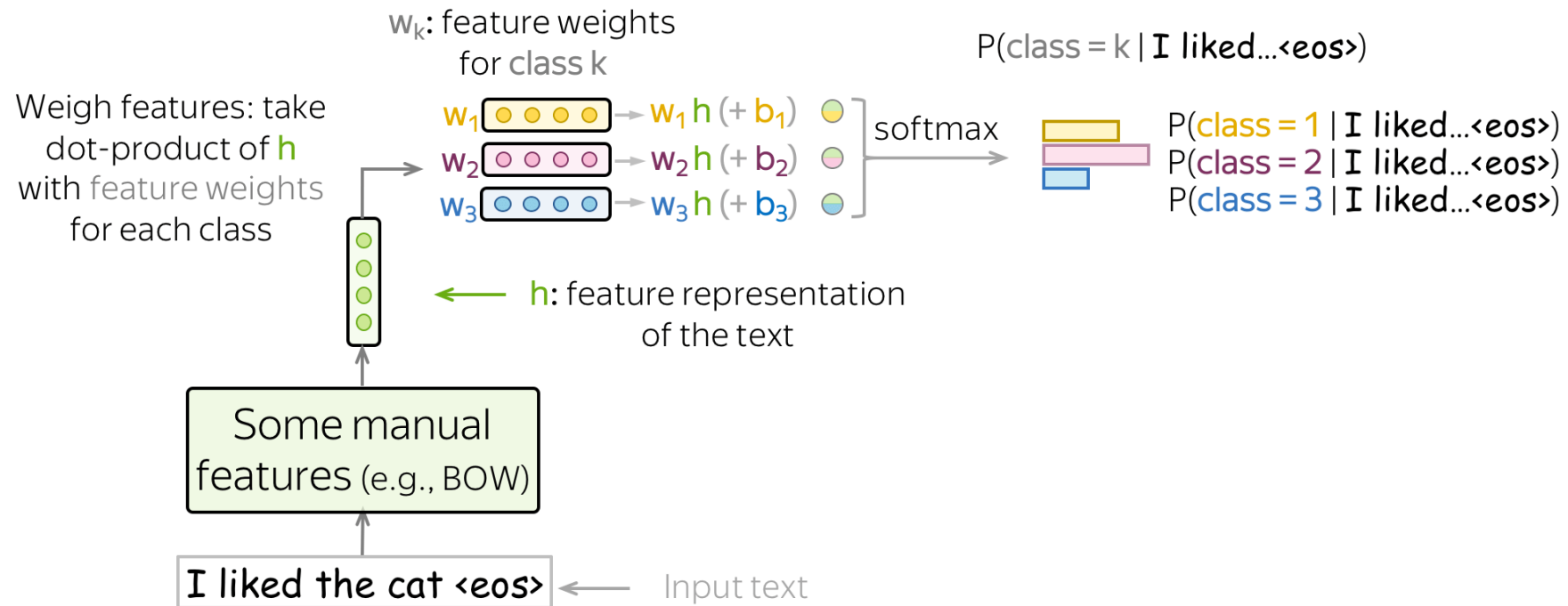




Table of Content

- General view on classification with text
- **Pre-deep learning approaches: Logistic Regression**
- Deep learning approaches
- Practical Tips

Train a logistic regression



Loss function: log-loss

Given training examples x^1, \dots, x^N with labels y^1, \dots, y^N , $y^i \in \{1, \dots, K\}$, we pick those weights $w^{(k)}$, $k = 1..K$ which maximize the probability of the training data:

$$w^* = \arg \max_w \sum_{i=1}^N \log P(y = y^i | x^i).$$

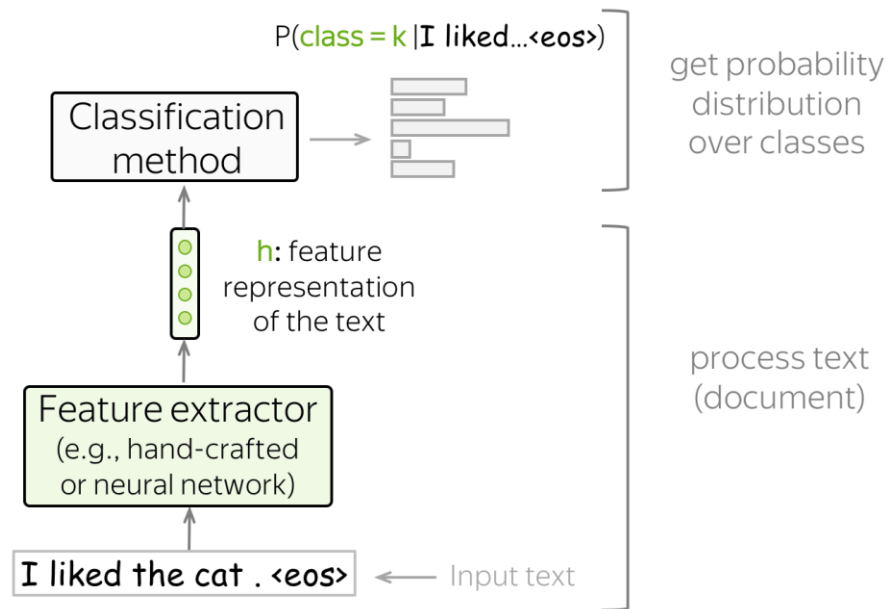


Table of Content

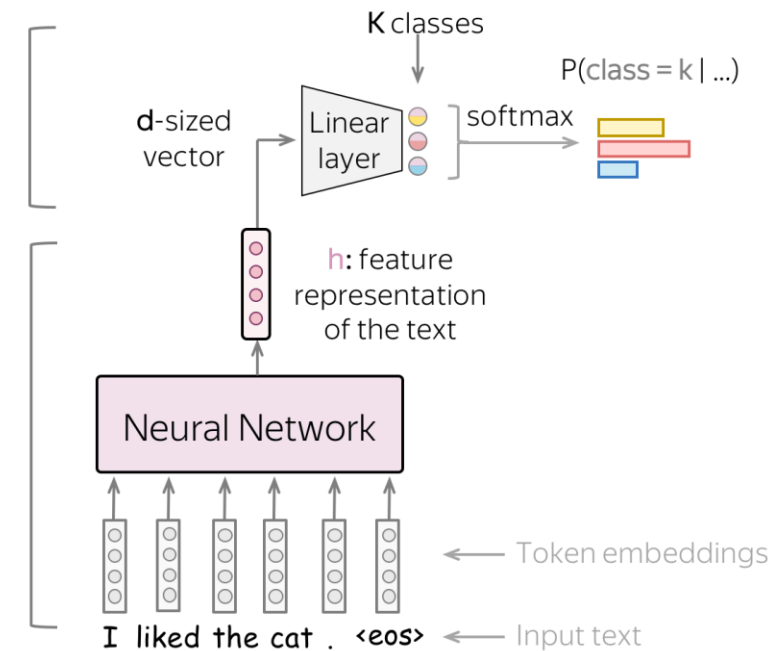
- General view on classification with text
- Pre-deep learning approaches
- **Deep learning approaches: Overview**
- Practical Tips

Overview of DL-based approaches

General Classification Pipeline



Classification with Neural Networks



Overview of DL-based approaches

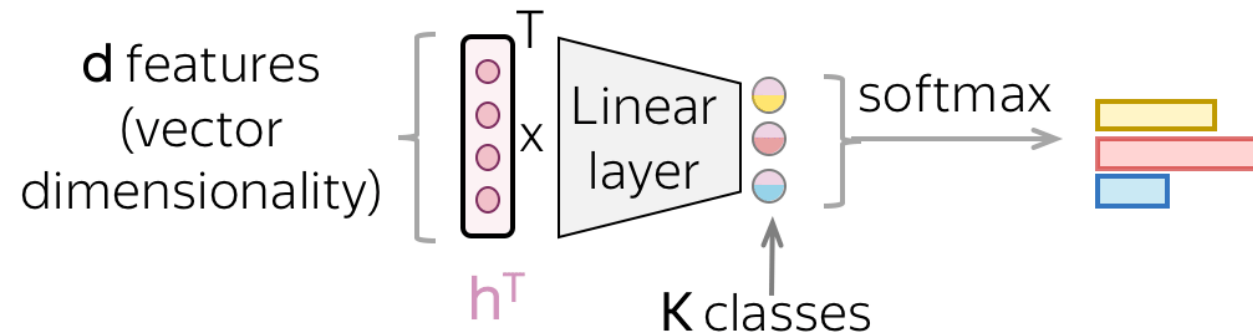
We have:

- h - vector of size d

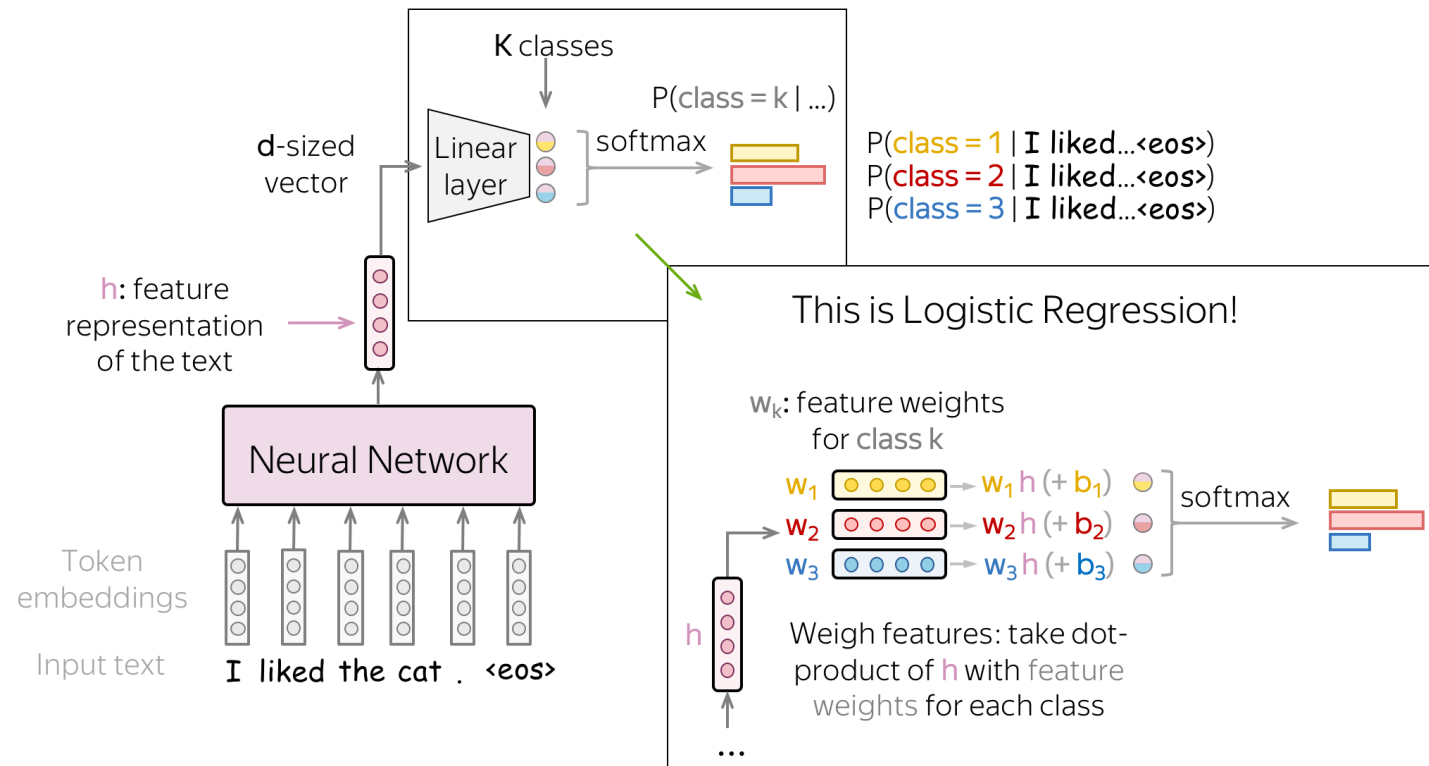
We need:

- vector of size K – probabilities for K classes

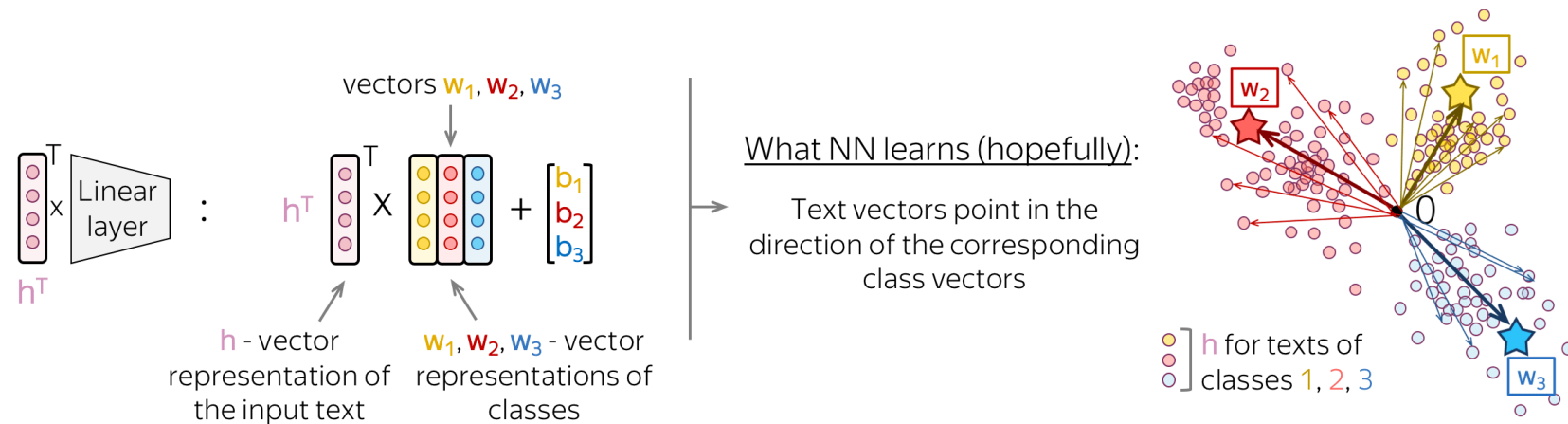
Transform linearly
from size d to size K



Overview of DL-based approaches



Overview of DL-based approaches



Overview of DL-based approaches

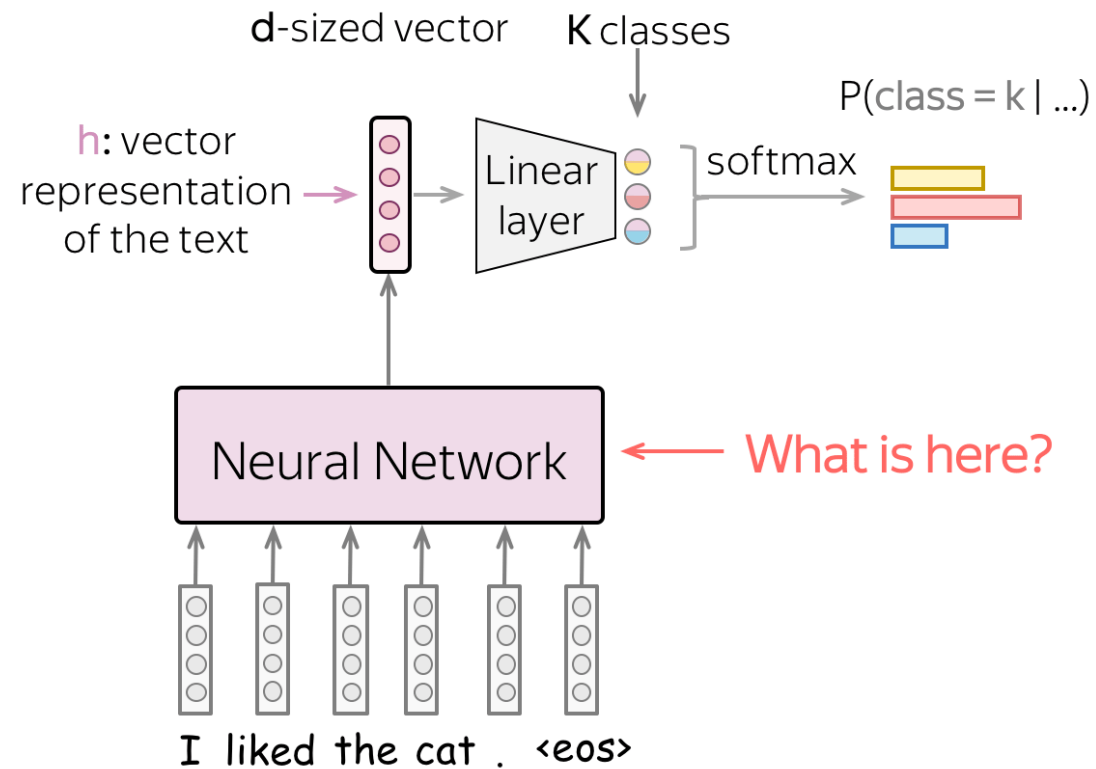




Table of Content

- General view on classification with text
- Pre-deep learning approaches
- **Deep learning approaches: Basics**
- Practical Tips

Vanilla embeddings

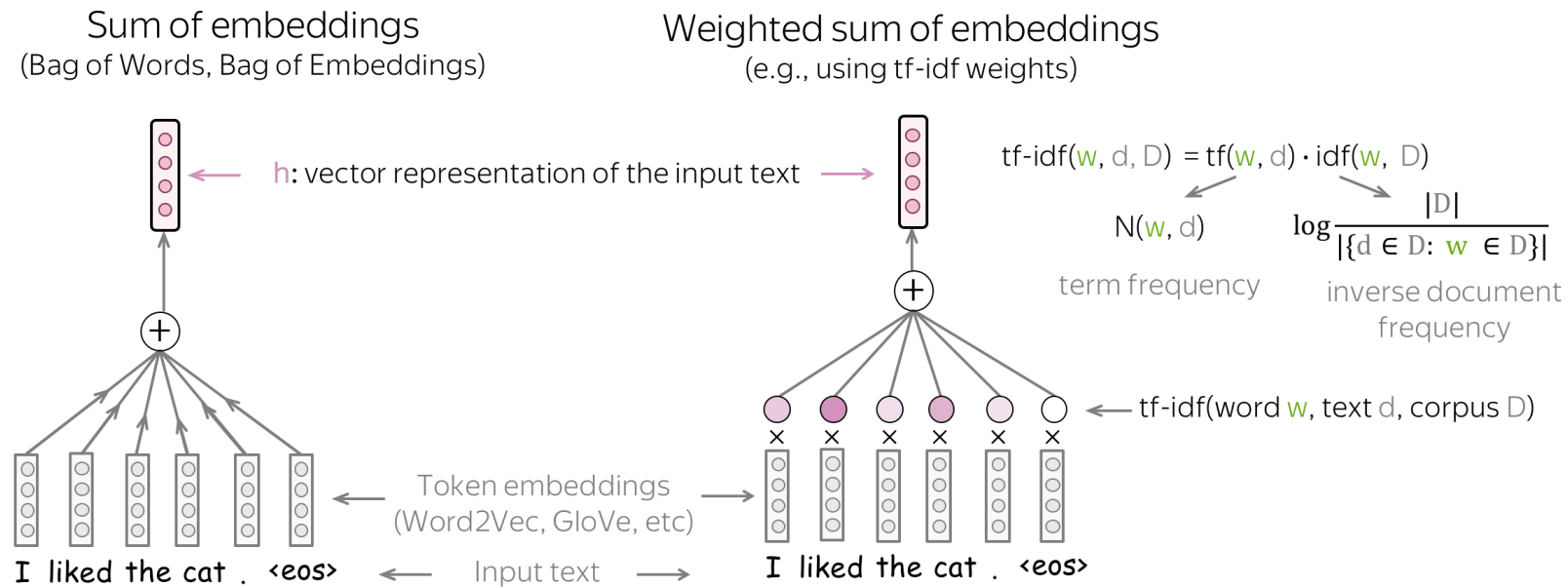
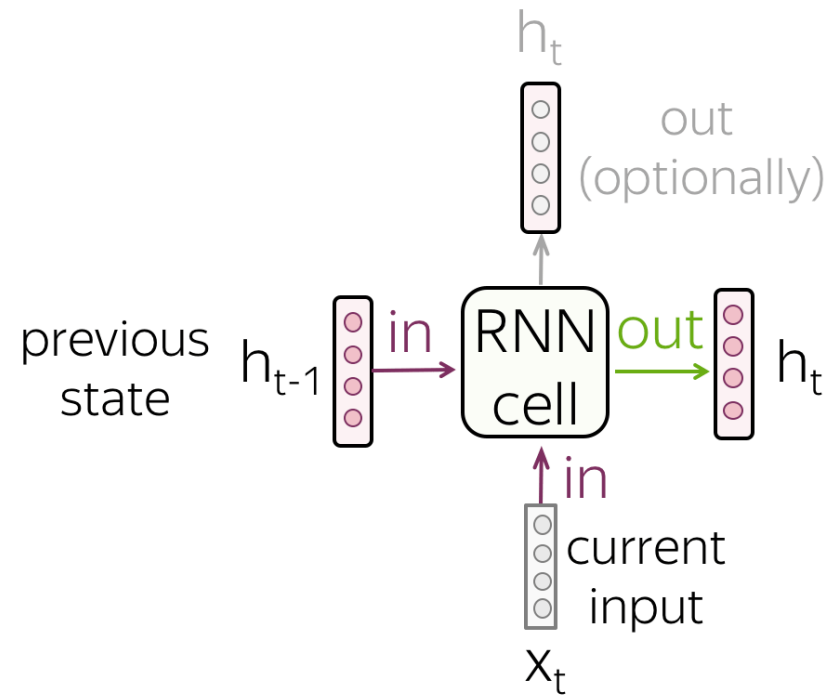




Table of Content

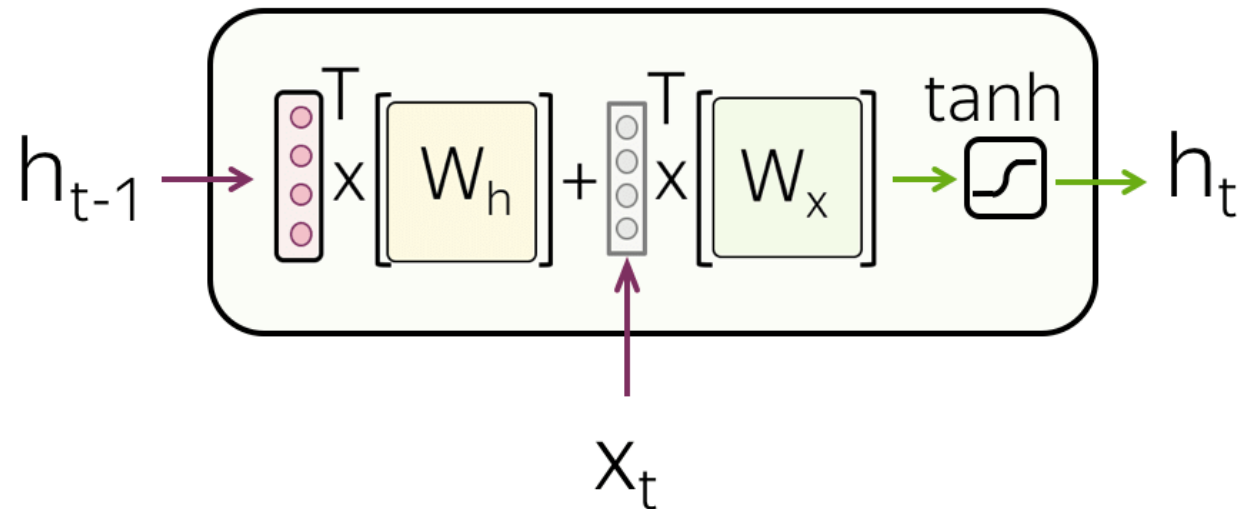
- General view on classification with text
- Pre-deep learning approaches
- **Deep learning approaches: RNN**
- Practical Tips

Vanilla embeddings

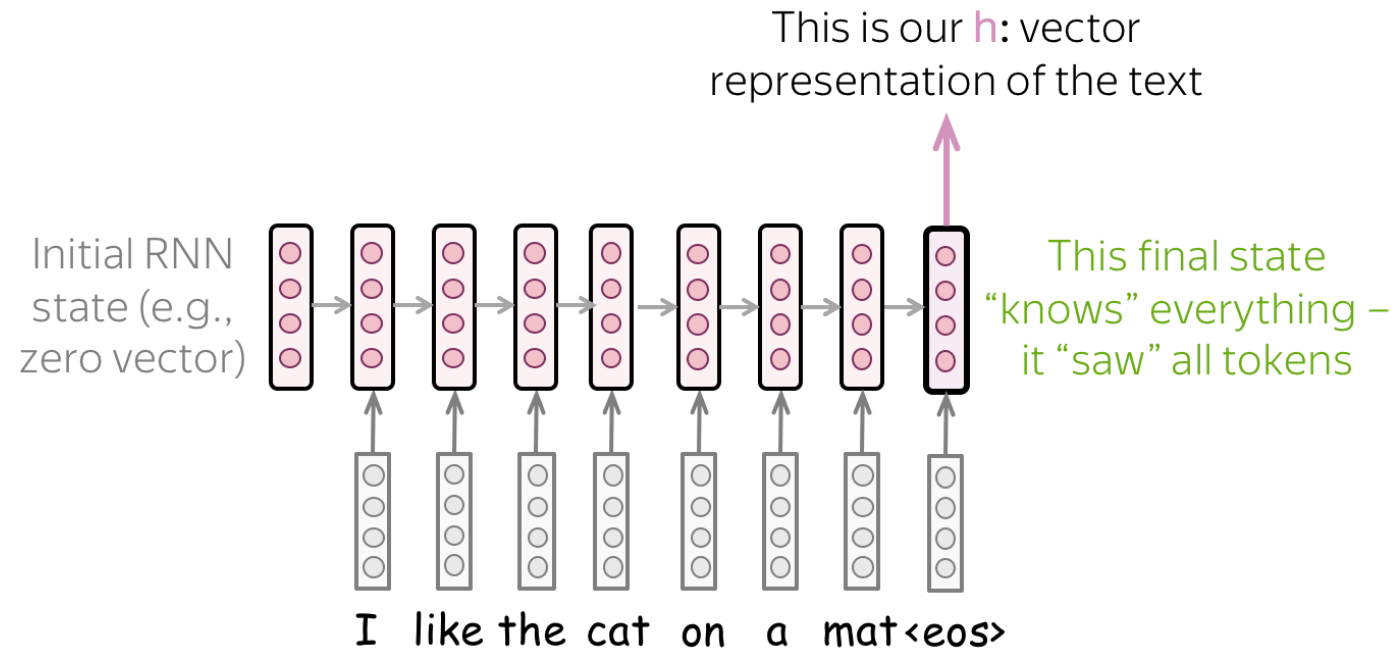


Vanilla RNN

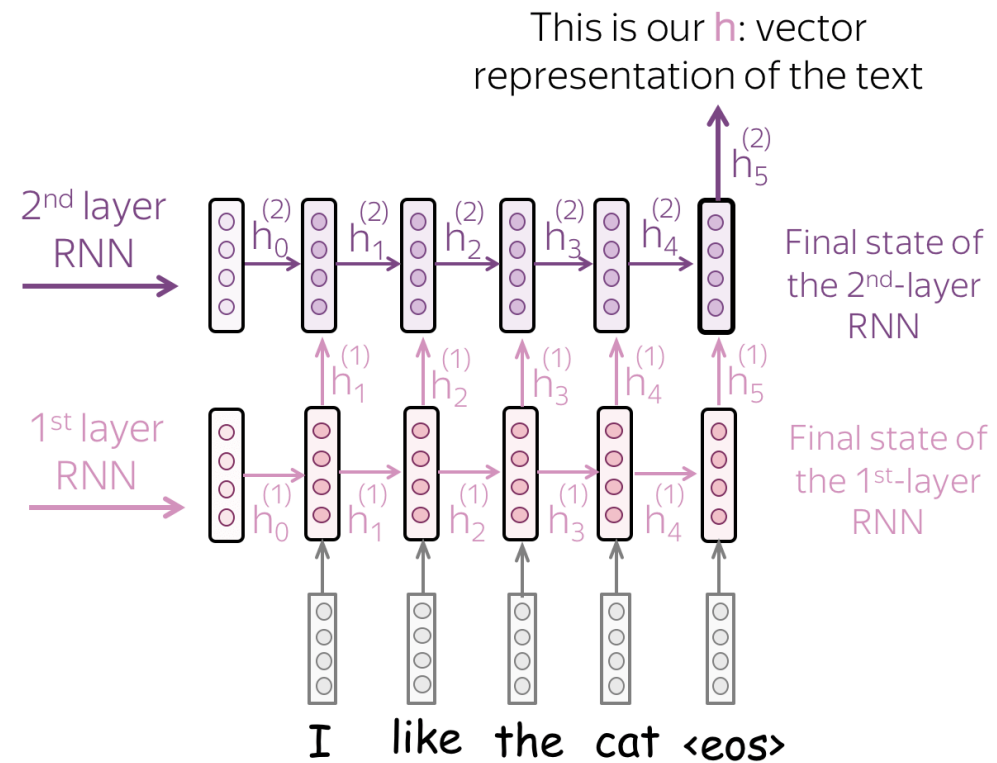
$$h_t = \tanh(h_{t-1}W_h + x_tW_x)$$



Vanilla RNN



Stacked RNN



Bidirectional RNN

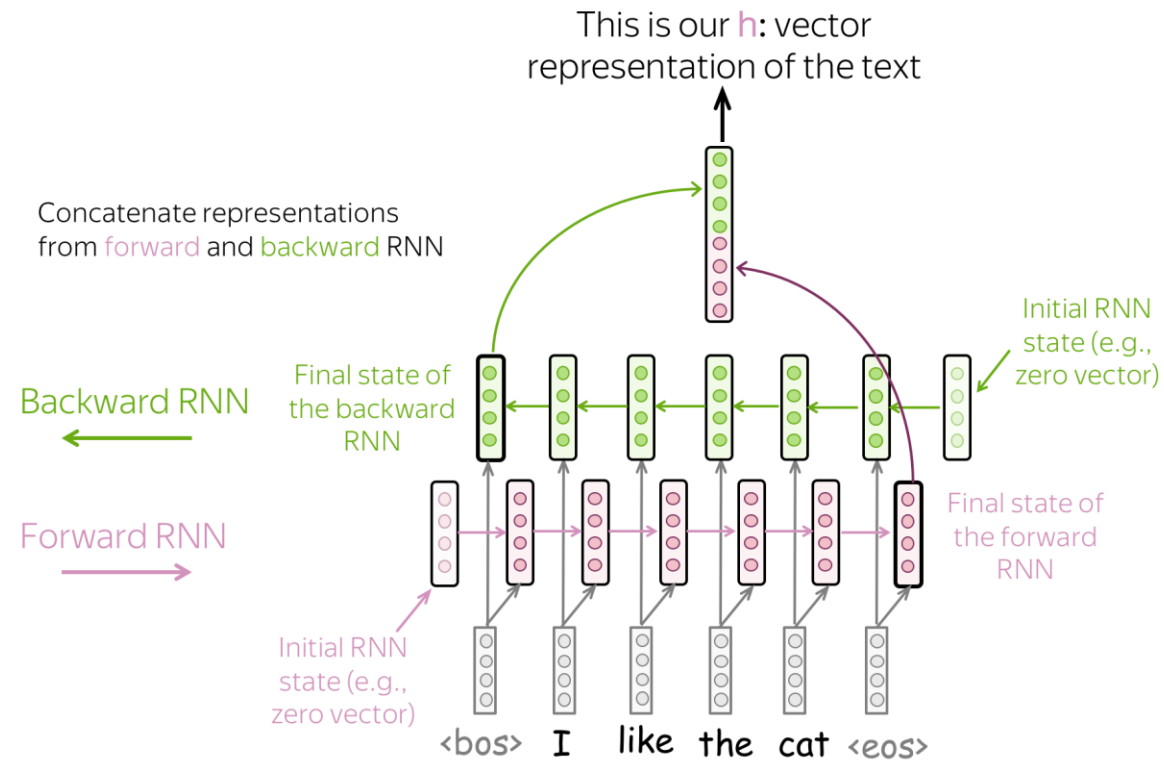
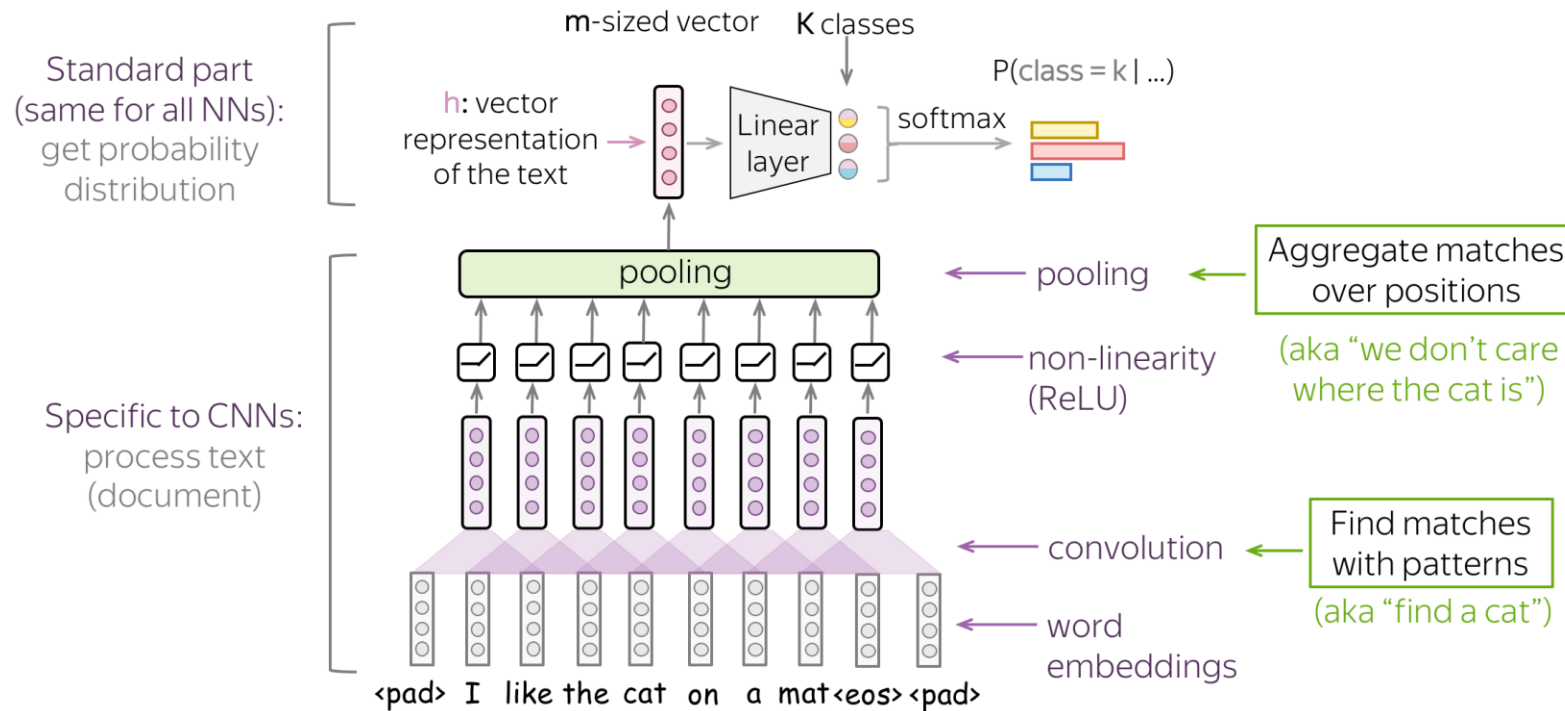




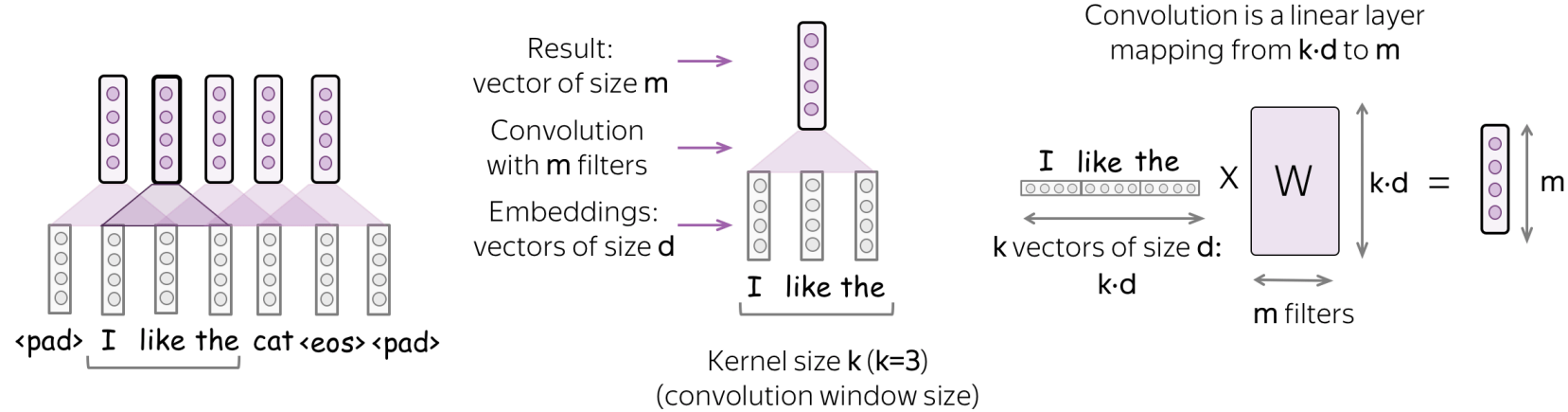
Table of Content

- General view on classification with text
- Pre-deep learning approaches
- **Deep learning approaches: CNN**
- Practical Tips

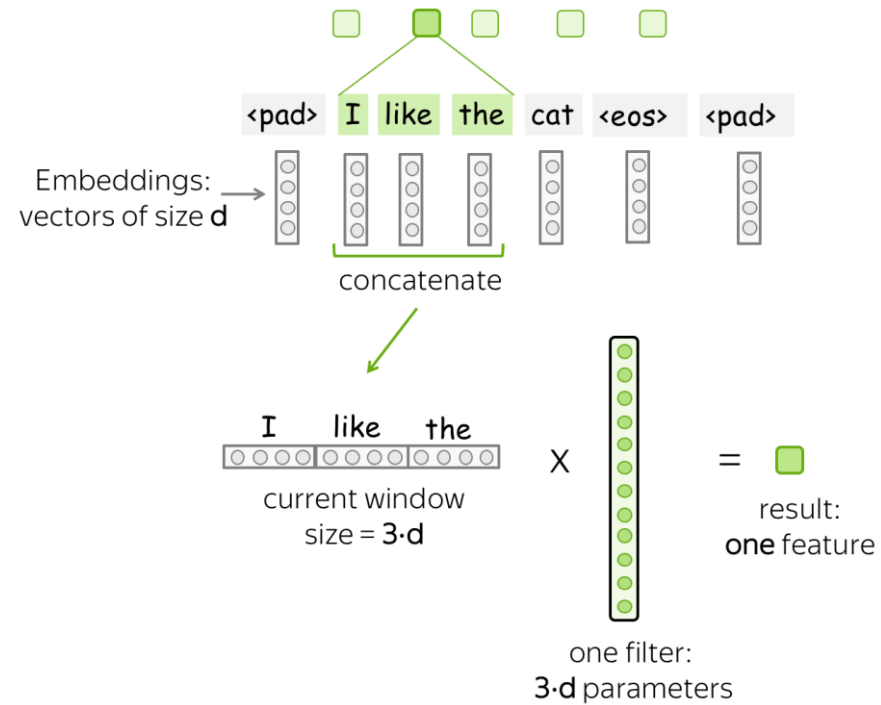
CNN overview



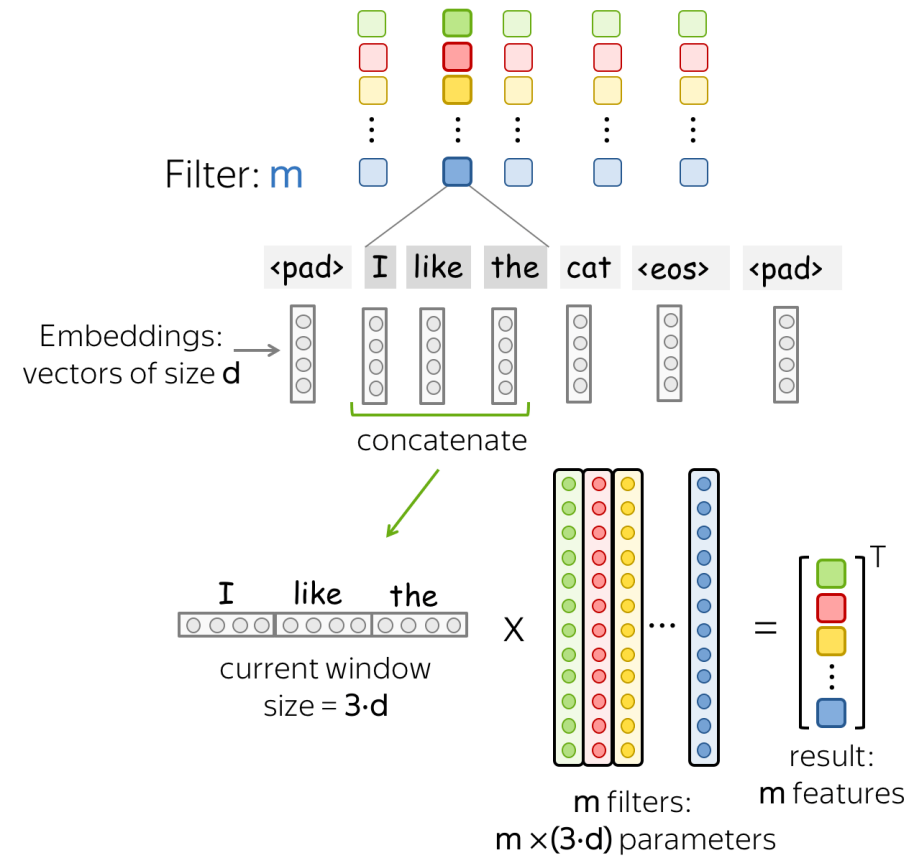
Convolution is a Linear Operation Applied to Each Window



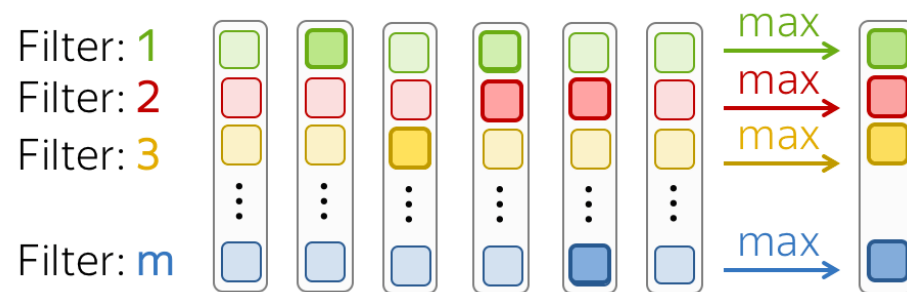
CNN intuition



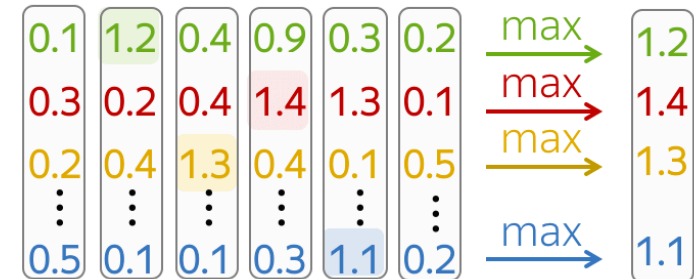
CNN intuition



Pooling



Max pooling:
maximum for each
dimension (feature)



Larger kernel size for a broader context

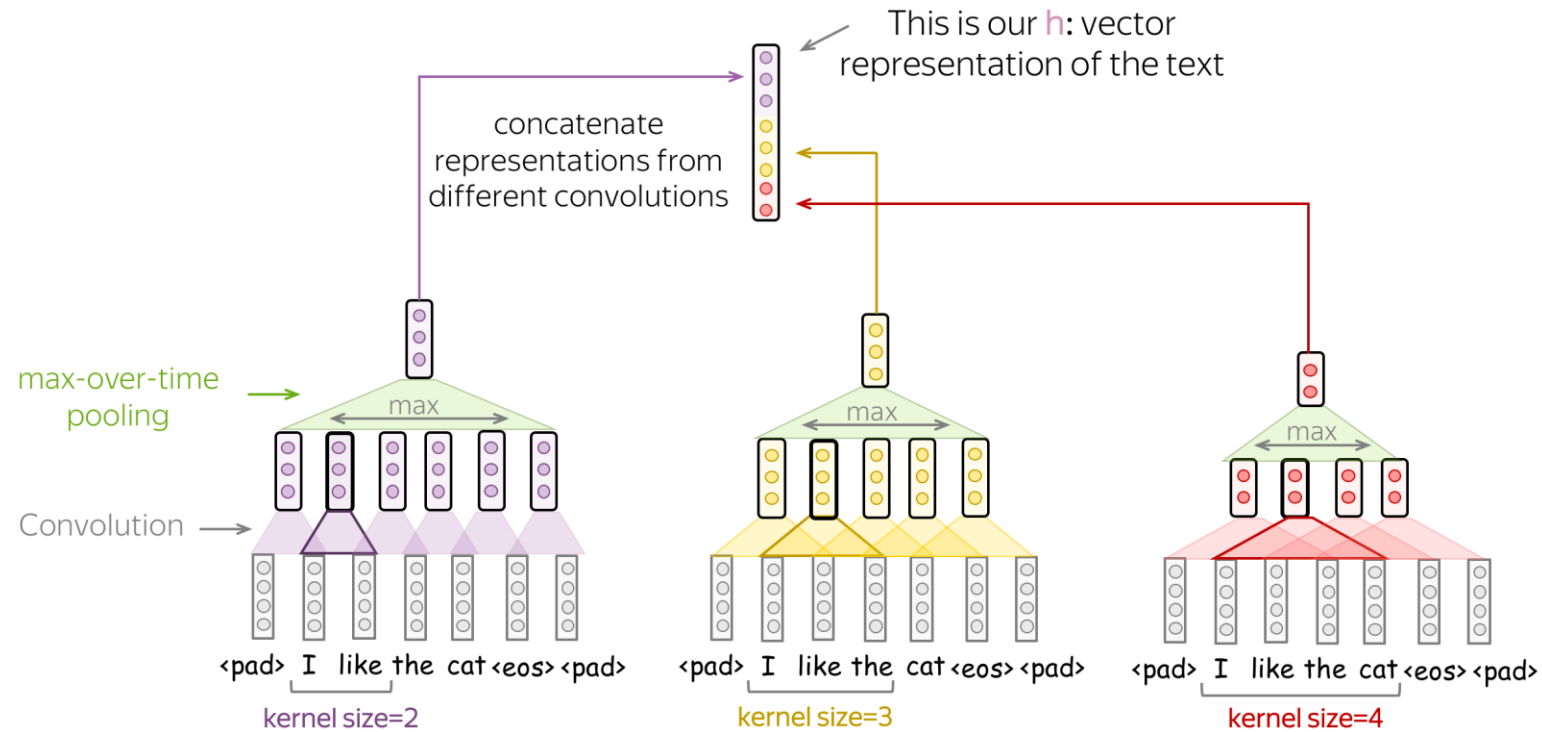
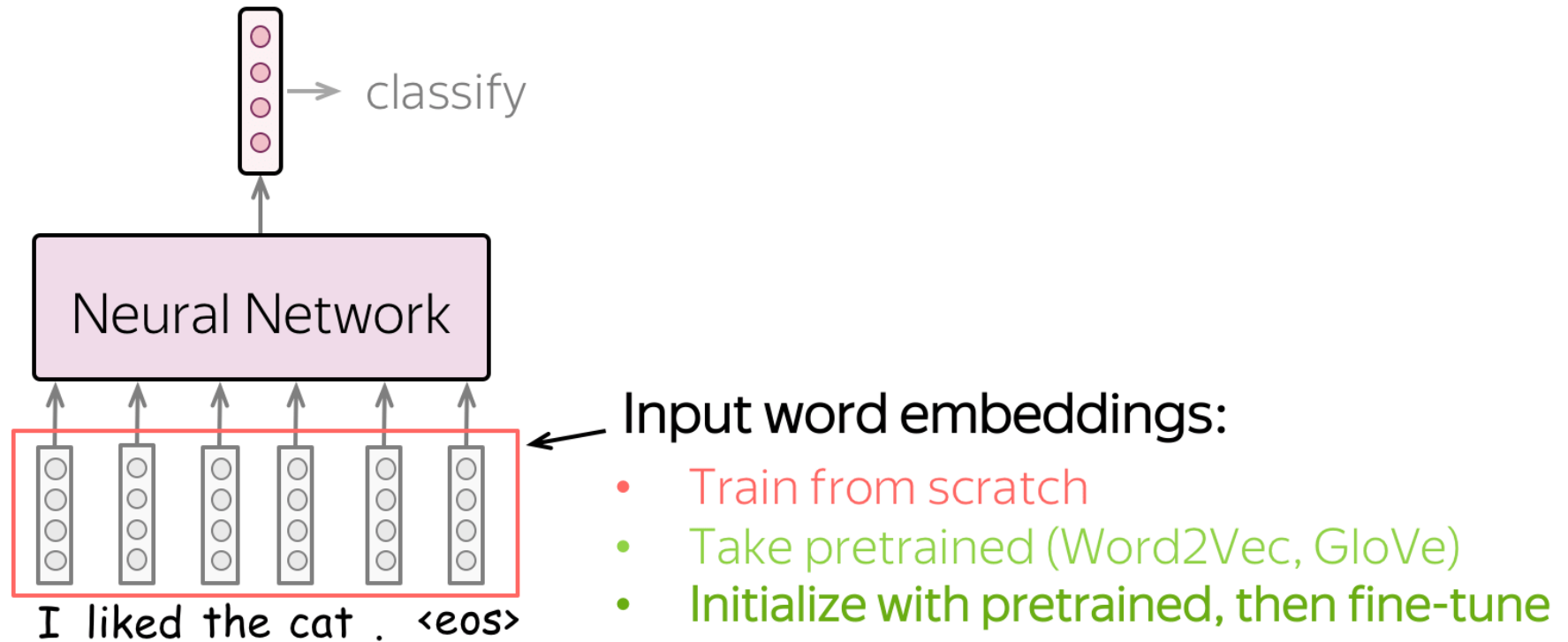




Table of Content

- General view on classification with text
- Pre-deep learning approaches
- Deep learning approaches
- **Practical Tips**

Larger kernel size for a broader context



Data augmentation

