

Go to

join.iClicker.com

HXFY



COS 484

Natural Language Processing

L8: Neural networks for NLP

Spring 2025

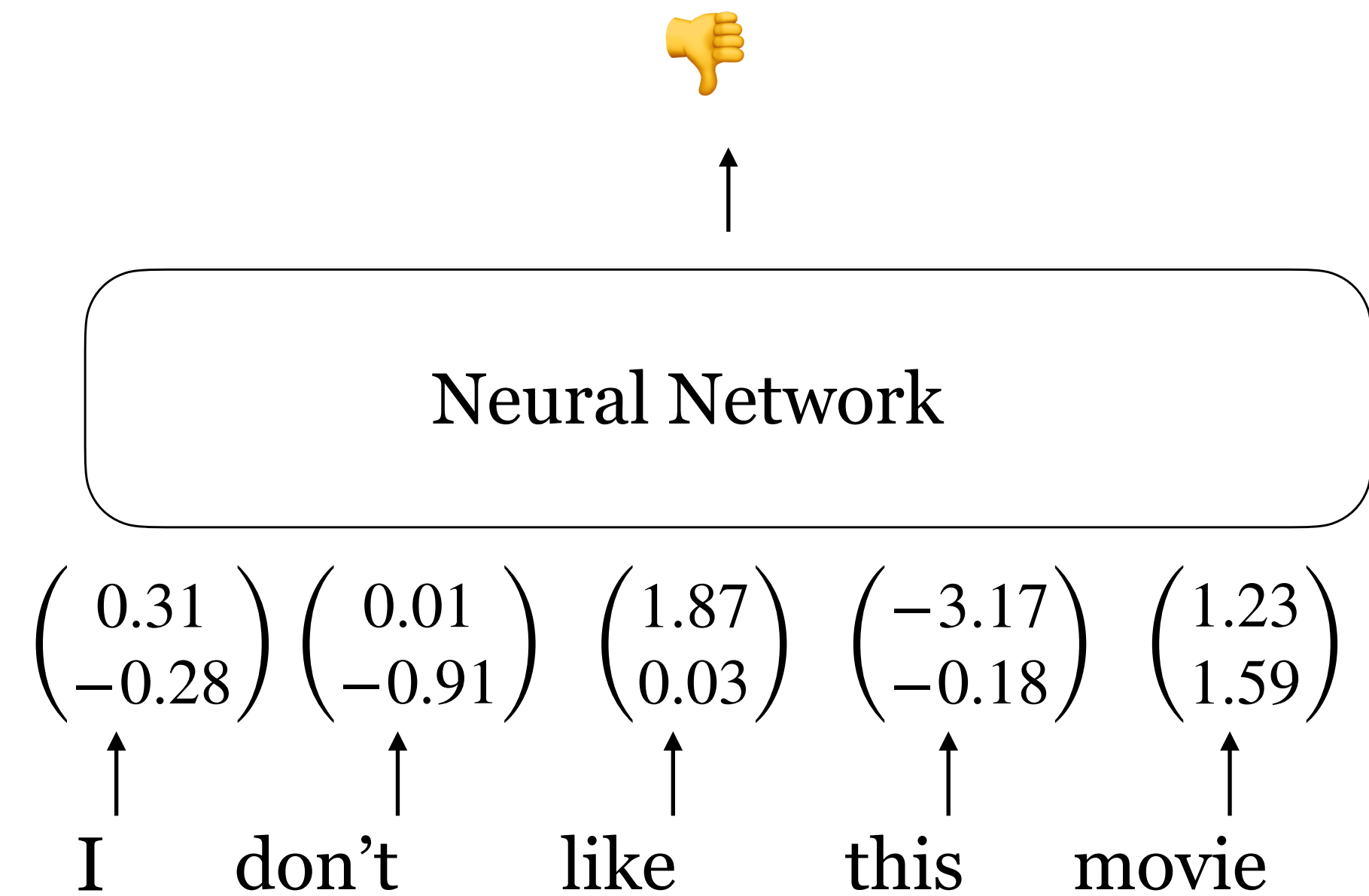
From word embeddings to neural networks

$$v_{\text{cat}} = \begin{pmatrix} -0.224 \\ 0.130 \\ -0.290 \\ 0.276 \end{pmatrix}$$

$$v_{\text{dog}} = \begin{pmatrix} -0.124 \\ 0.430 \\ -0.200 \\ 0.329 \end{pmatrix}$$

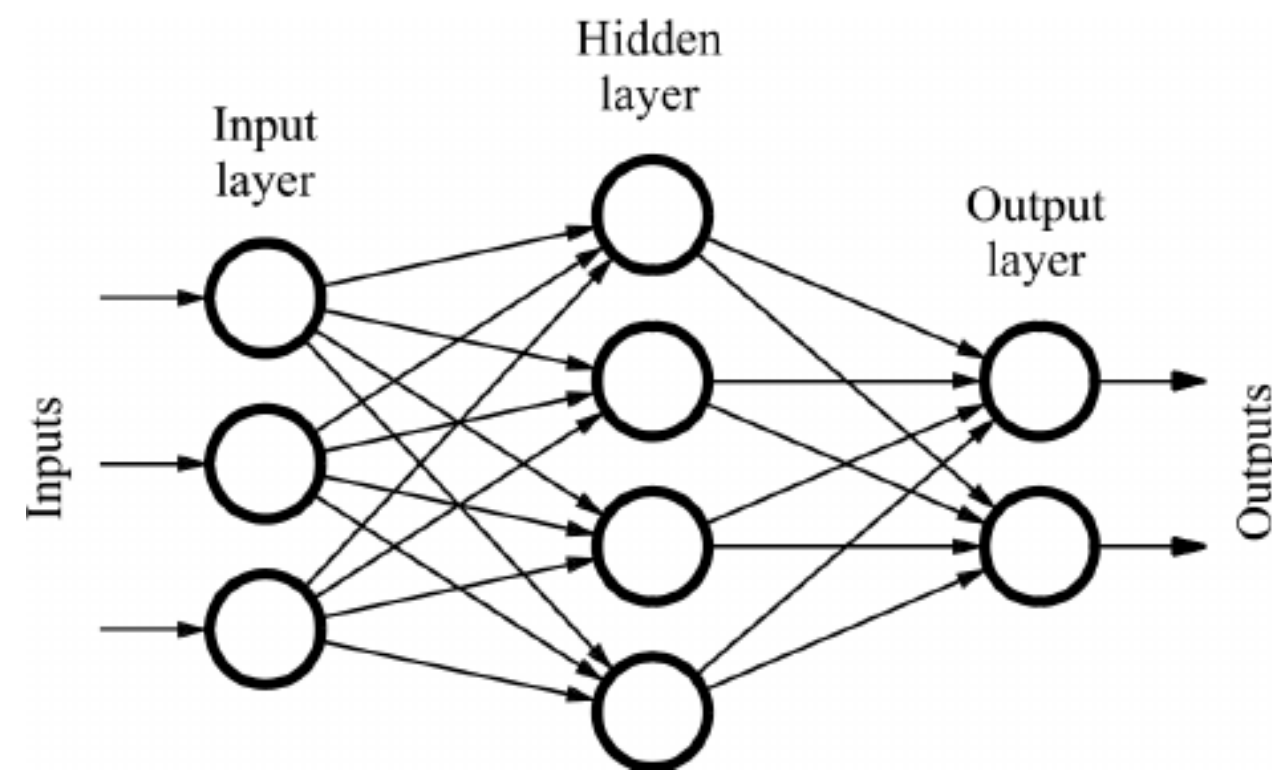
$$v_{\text{the}} = \begin{pmatrix} 0.234 \\ 0.266 \\ 0.239 \\ -0.199 \end{pmatrix}$$

$$v_{\text{language}} = \begin{pmatrix} 0.290 \\ -0.441 \\ 0.762 \\ 0.982 \end{pmatrix}$$

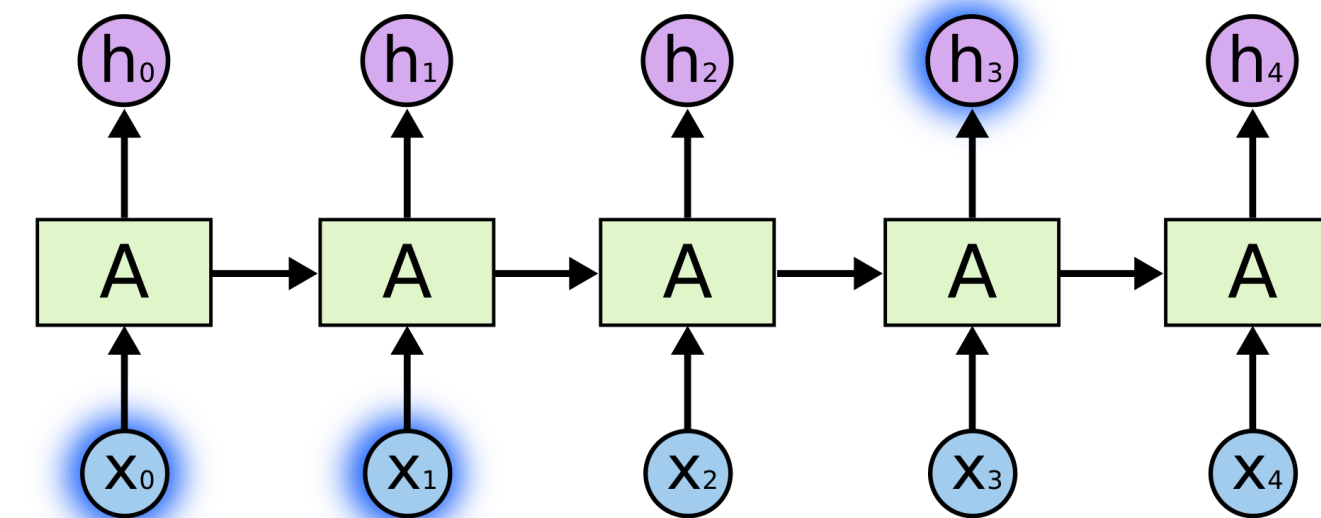


Neural networks in NLP

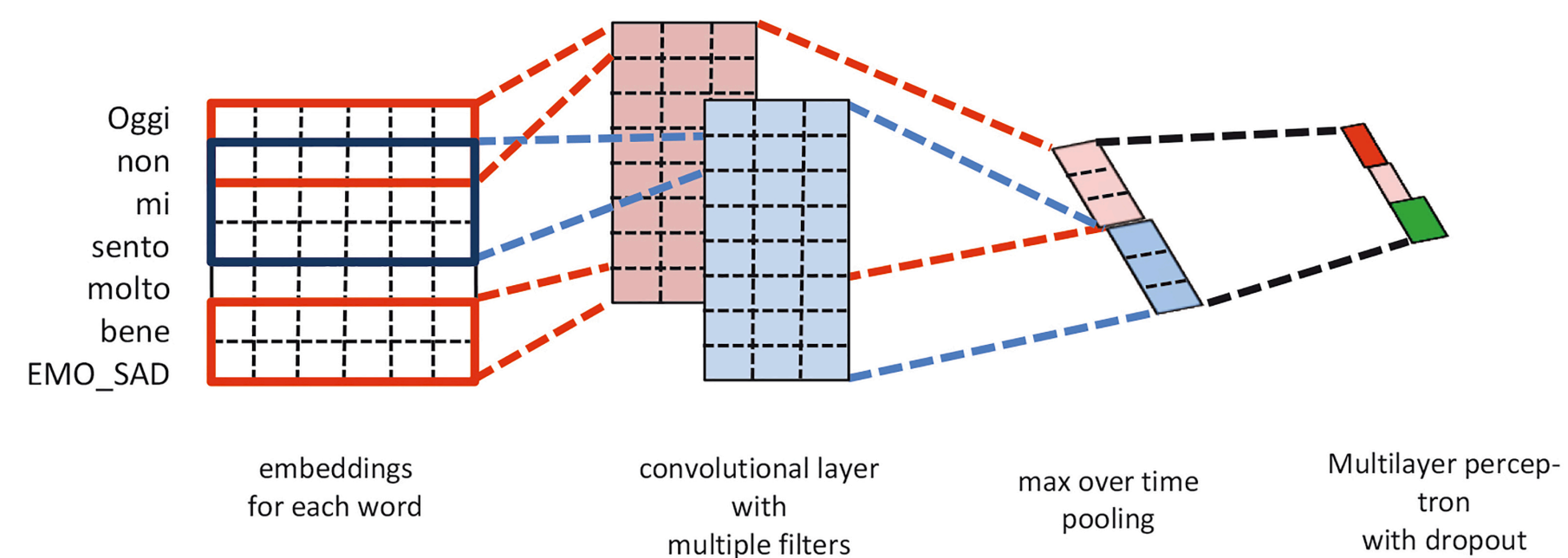
Feed-forward NNs



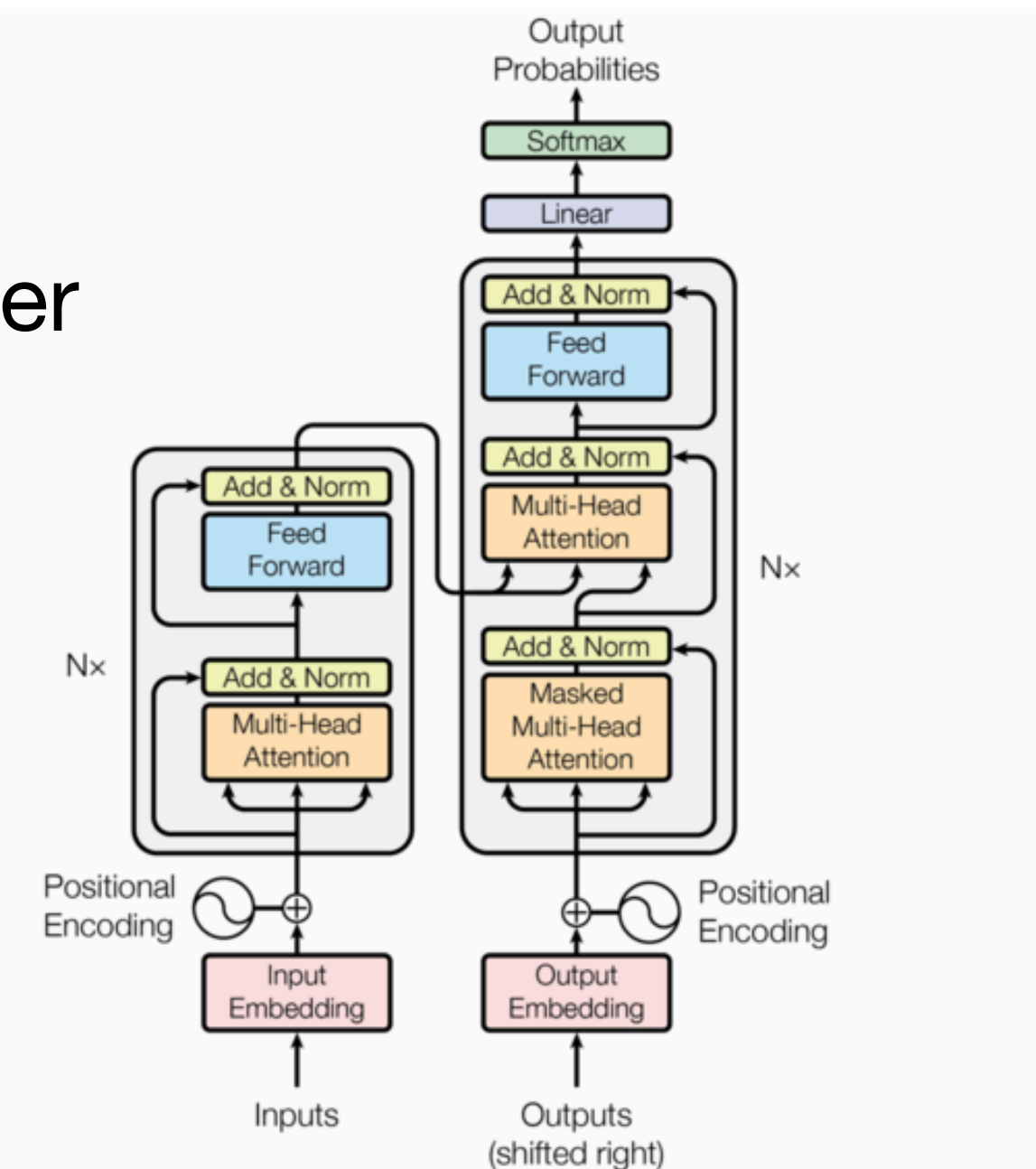
Recurrent NNs



Convolutional NNs



Transformer



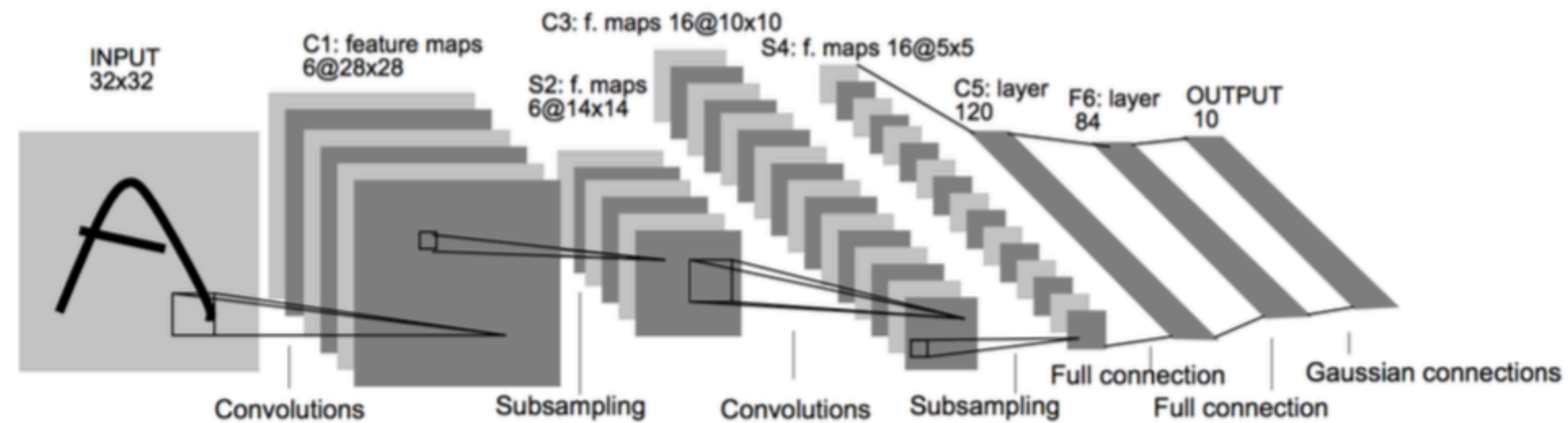
Lecture plan

- Neural networks for NLP: a brief history
- Recap: feedforward neural networks
- Main focus: two simple applications of NNs for NLP
 - Neural bag-of-words model for text classification
 - Feedforward neural language models

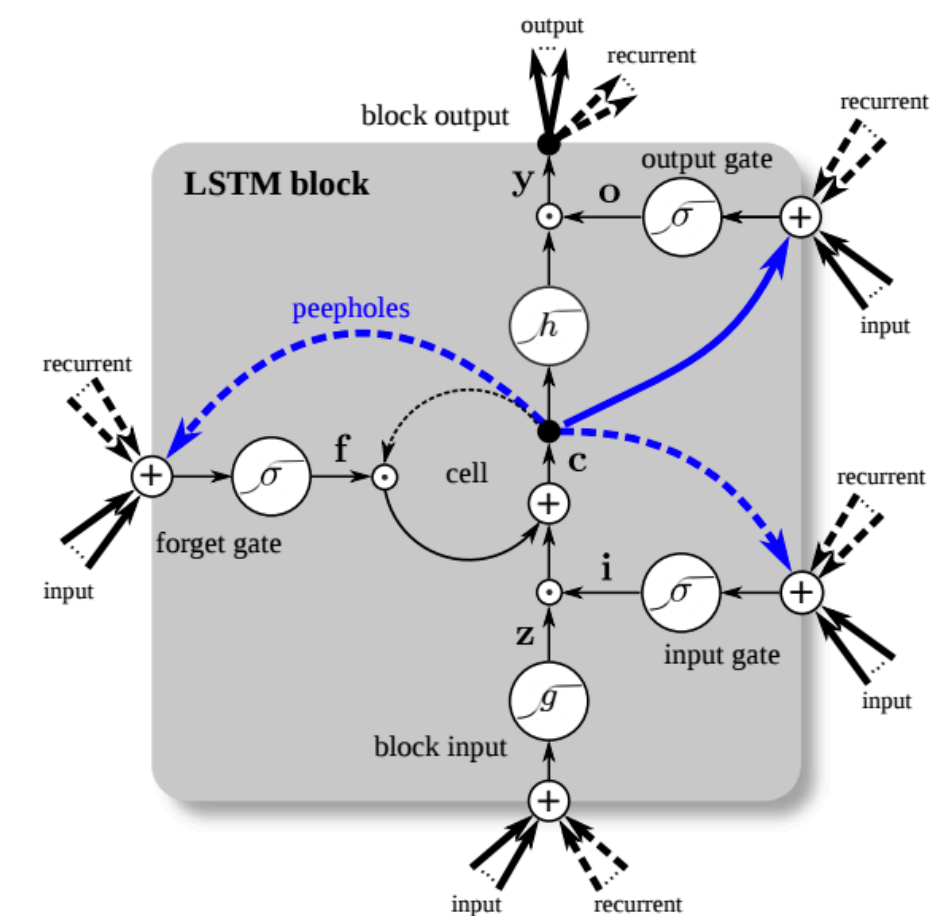
Neural networks for NLP: History

NN “dark ages”

- Neural network algorithms date from the 80s
- ConvNets: applied to MNIST by LeCun in 1998



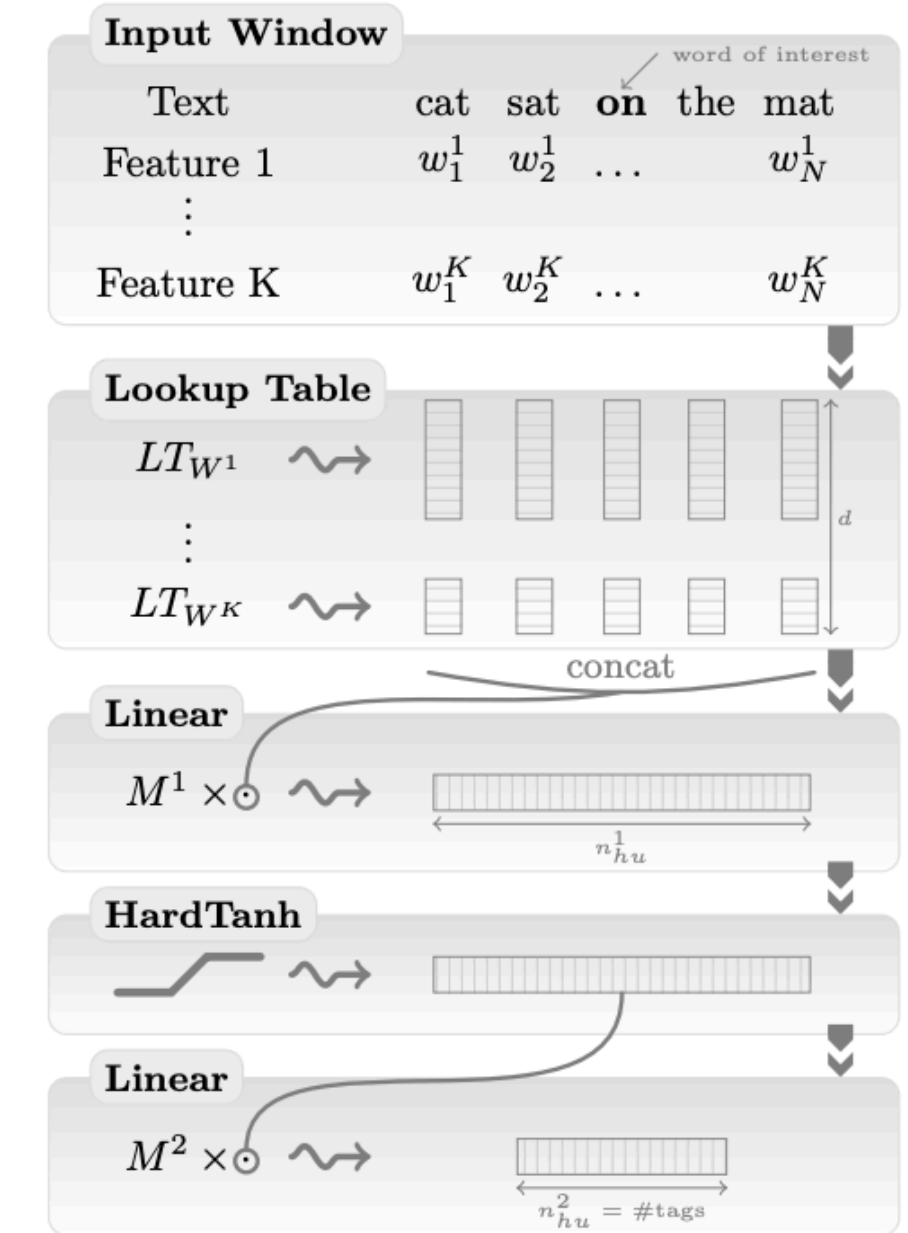
- Long Short-term Memory Networks (LSTMs): Hochreiter and Schmidhuber 1997
- Henderson 2003: neural shift-reduce parser, not SOTA



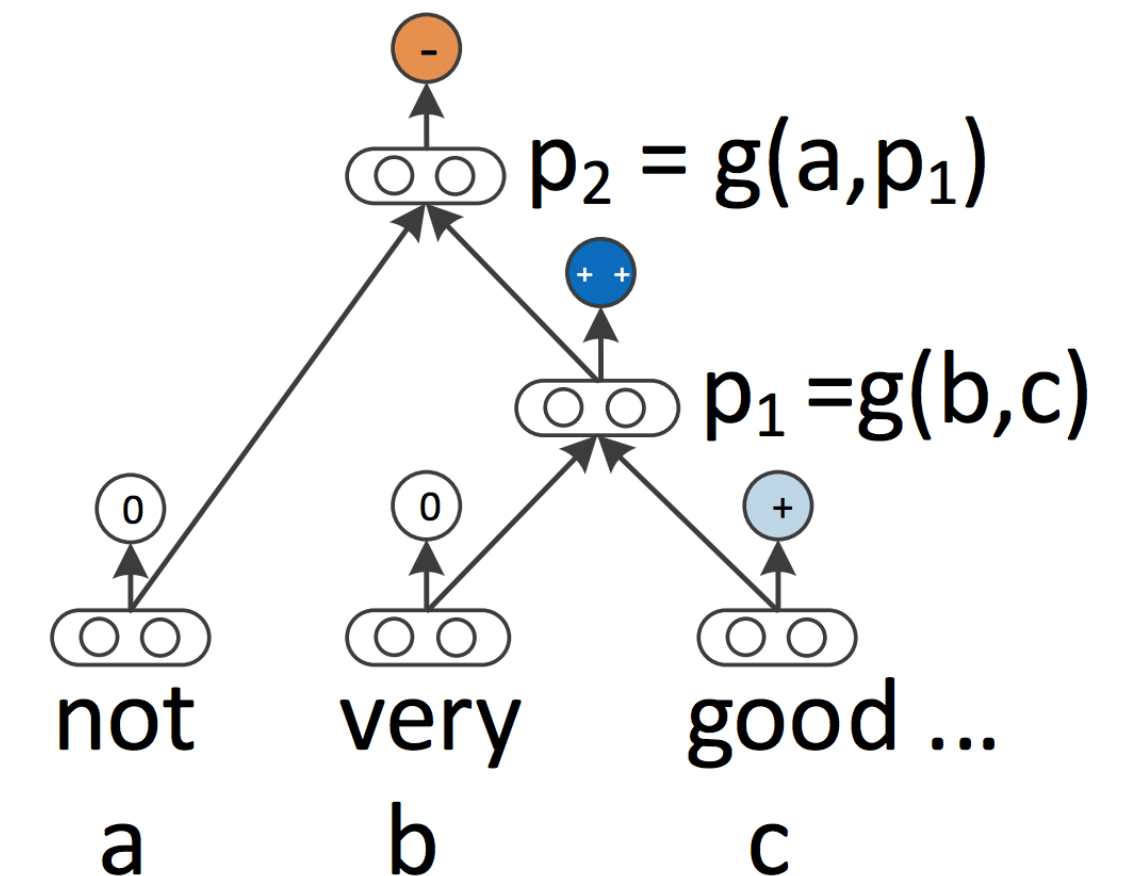
Slide credit: Greg Durrett

2008-2013: A glimmer of light

- Collobert and Weston 2011: “**NLP (almost) from Scratch**”
- Feedforward NNs can replace “feature engineering”
- 2008 version was marred by bad experiments, claimed SOTA but wasn’t, 2011 version tied SOTA



- Krizhevsky et al, 2012: AlexNet for ImageNet Classification
- Socher 2011-2014: tree-structured RNNs working okay



2014: Stuff starts working

- Kim (2014) + Kalchbrenner et al, 2014: sentence classification
 - ConvNets work for NLP!
- Sutskever et al, 2014: sequence-to-sequence for neural MT
 - LSTMs work for NLP!
- Chen and Manning 2014: dependency parsing
 - Even feedforward networks work well for NLP!
- 2015: explosion of neural networks for everything under the sun
- 2018-2019: NLP has entered the era of pre-trained models (ELMo, GPT, BERT)
- 2020+: the emergence of large language models (GPT-3, ChatGPT, OpenAI o1/DeepSeek R1)

Convolutional Neural Networks for Sentence Classification

Yoon Kim
New York University
yhk255@nyu.edu

A Fast and Accurate Dependency Parser using Neural Networks

Danqi Chen
Computer Science Department
Stanford University
danqi@cs.stanford.edu

Christopher D. Manning
Computer Science Department
Stanford University
manning@stanford.edu

Why didn't they work before?

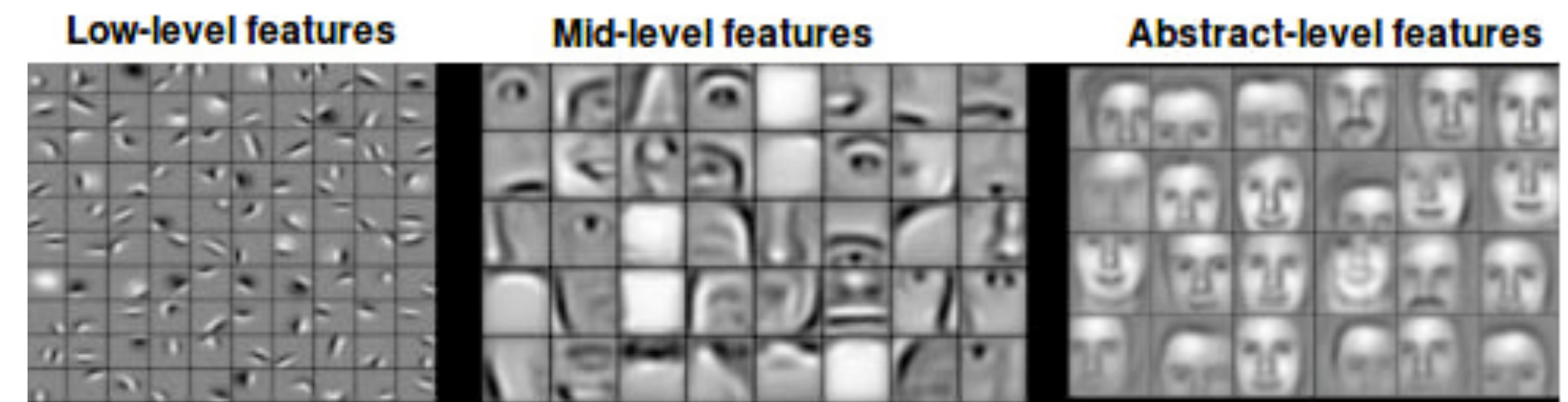
- **Datasets too small:** for machine translation, not really better until you have 1M+ parallel sentences (and really need a lot more)
- **Optimization not well understood:** good initialization, per-feature scaling + momentum (Adagrad/Adam) work best out-of-the-box
 - Regularization: dropout is pretty helpful
 - Computers not big enough: can't run for enough iterations
- Inputs: need **word embeddings** to have the right continuous semantics

The “promise” of deep learning

- Most NLP works in the past focused on human-designed representations and input features

Var	Definition	Value in Fig. 5.2
x_1	count(positive lexicon) \in doc)	3
x_2	count(negative lexicon) \in doc)	2
x_3	$\begin{cases} 1 & \text{if "no"} \in \text{doc} \\ 0 & \text{otherwise} \end{cases}$	1
x_4	count(1st and 2nd pronouns \in doc)	3
x_5	$\begin{cases} 1 & \text{if "!"} \in \text{doc} \\ 0 & \text{otherwise} \end{cases}$	0
x_6	log(word count of doc)	$\ln(64) = 4.15$

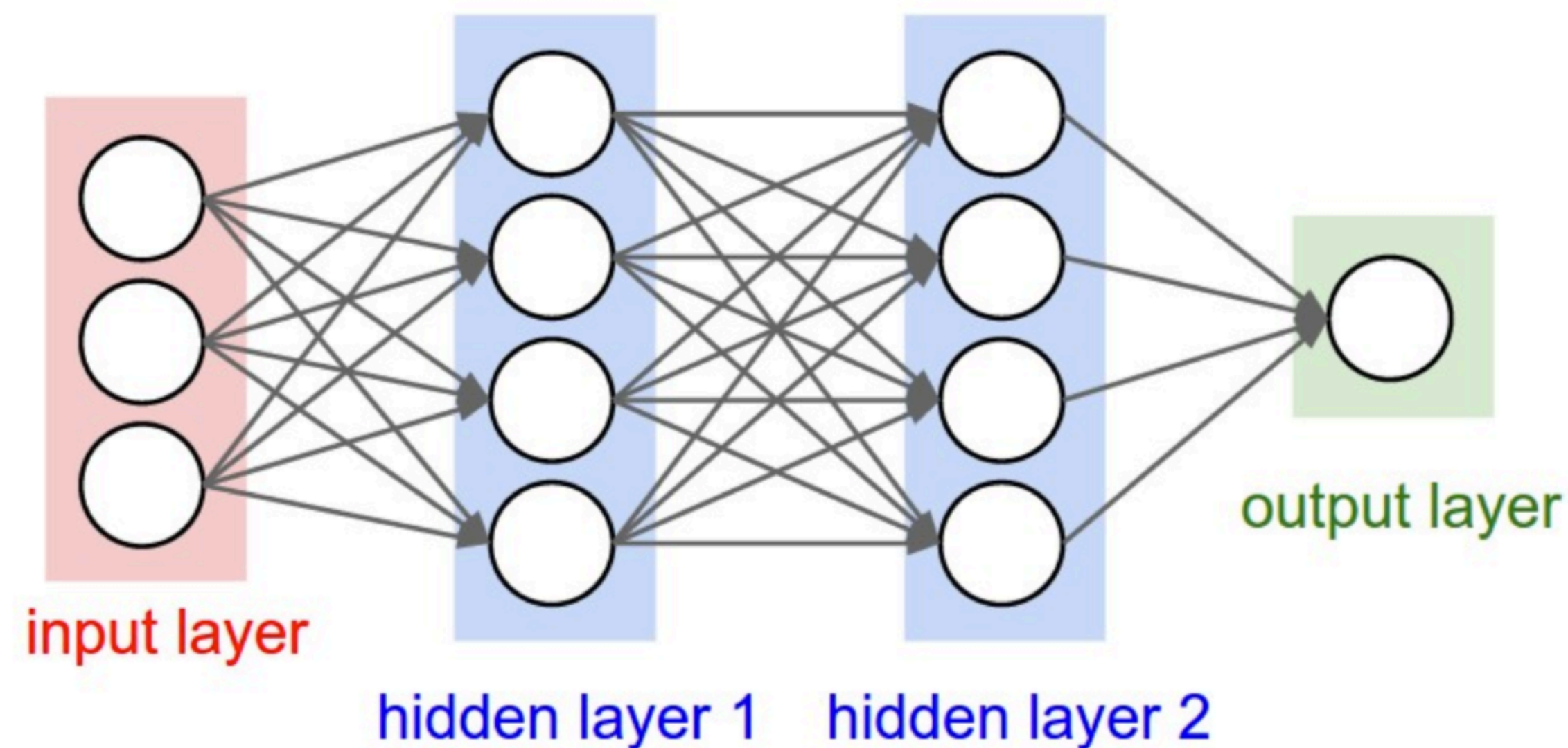
- **Representation learning** attempts to automatically learn good features and representations
- **Deep learning** attempts to learn multiple levels of representations on increasing complexity/abstraction



Review: Feedforward neural networks

Feed-forward NNs

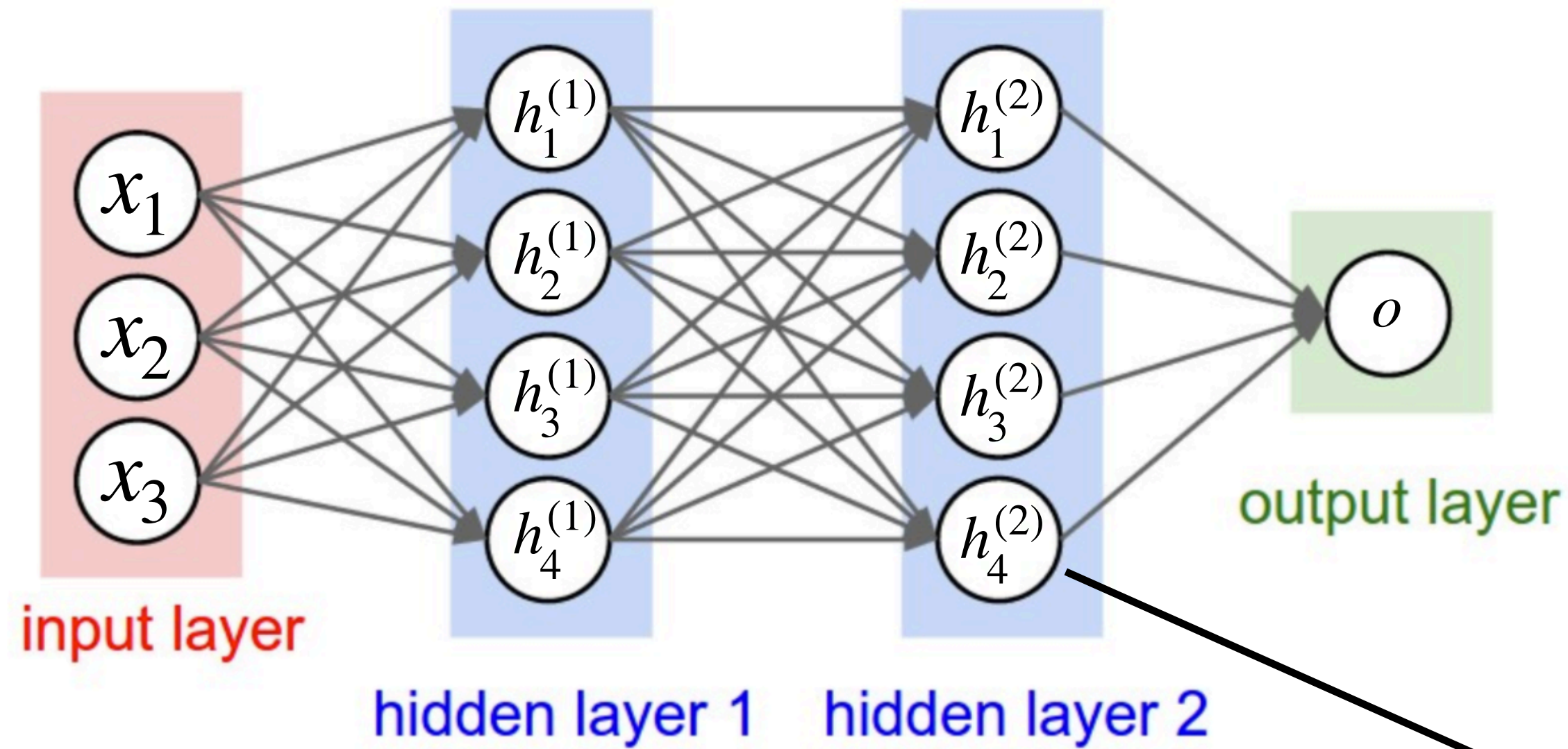
- The units are connected with no cycles
- The outputs from units in each layer are passed to units in the next higher layer
- No outputs are passed back to lower layers



Fully-connected (FC) layers:

All the units from one layer are fully connected to every unit of the next layer.

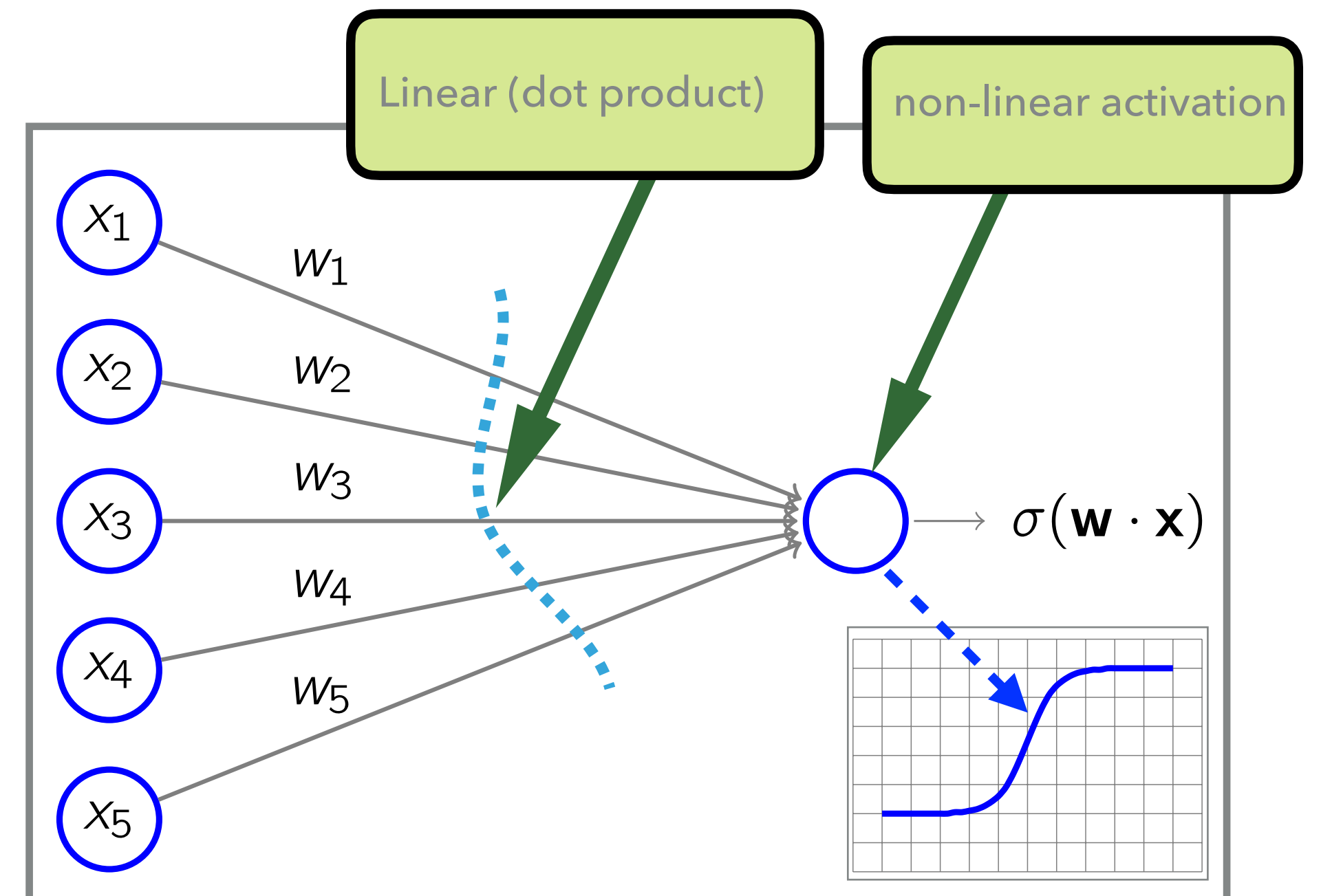
Feed-forward NNs



$$h_1^{(1)} = f(w_{1,1}^{(1)}x_1 + w_{1,2}^{(1)}x_2 + w_{1,3}^{(1)}x_3)$$

$$h_3^{(2)} = f(w_{3,1}^{(2)}h_1^{(1)} + w_{3,2}^{(2)}h_2^{(1)} + w_{3,3}^{(2)}h_3^{(1)} + w_{3,4}^{(2)}h_4^{(1)})$$

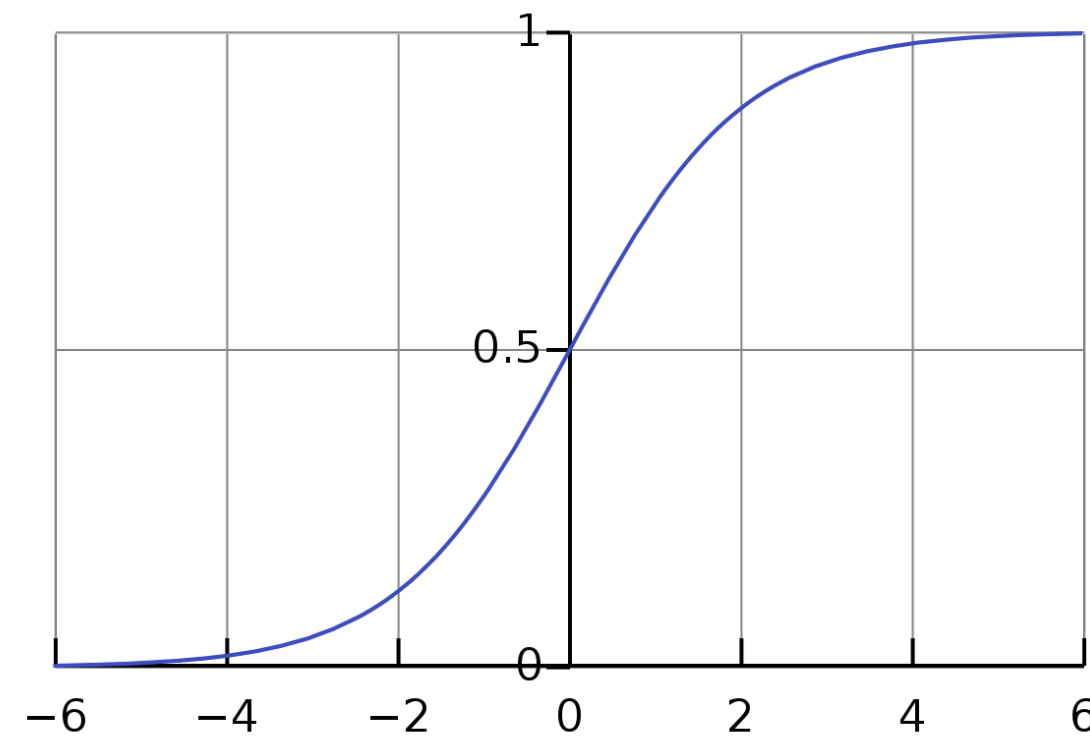
non-linearity f : σ , tanh or ReLU.



Activation functions

sigmoid

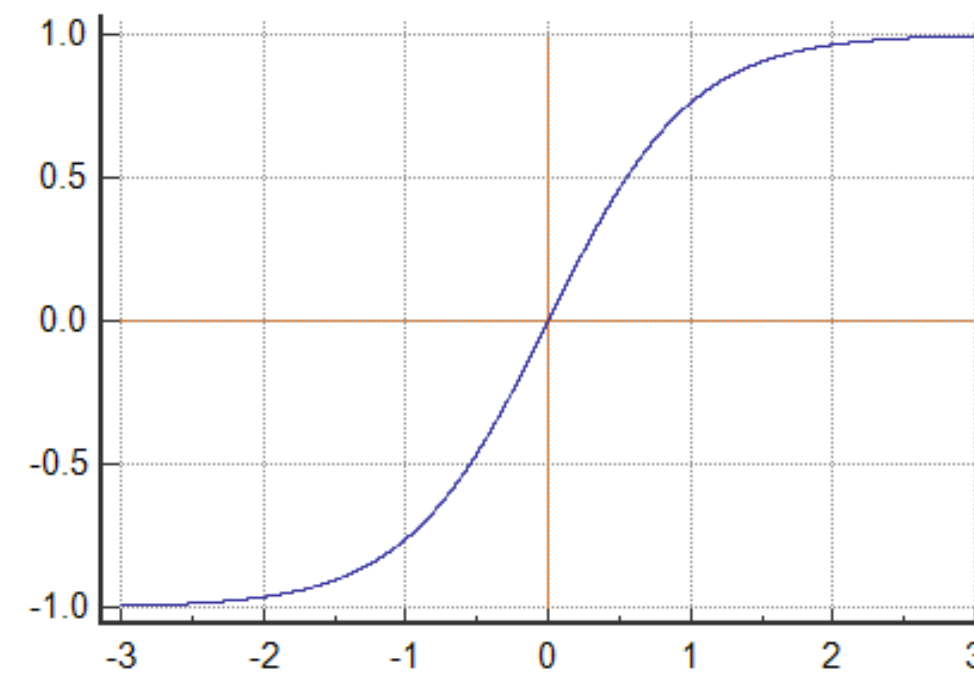
$$f(z) = \frac{1}{1 + e^{-z}}$$



$$f'(z) = f(z) \times (1 - f(z))$$

tanh

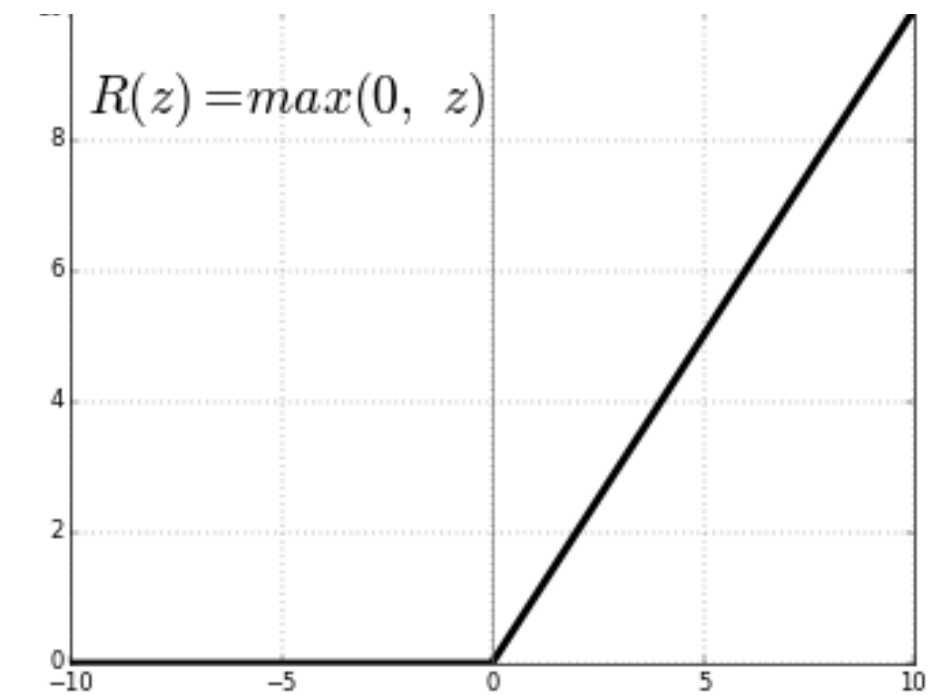
$$f(z) = \frac{e^{2z} - 1}{e^{2z} + 1}$$



$$f'(z) = 1 - f(z)^2$$

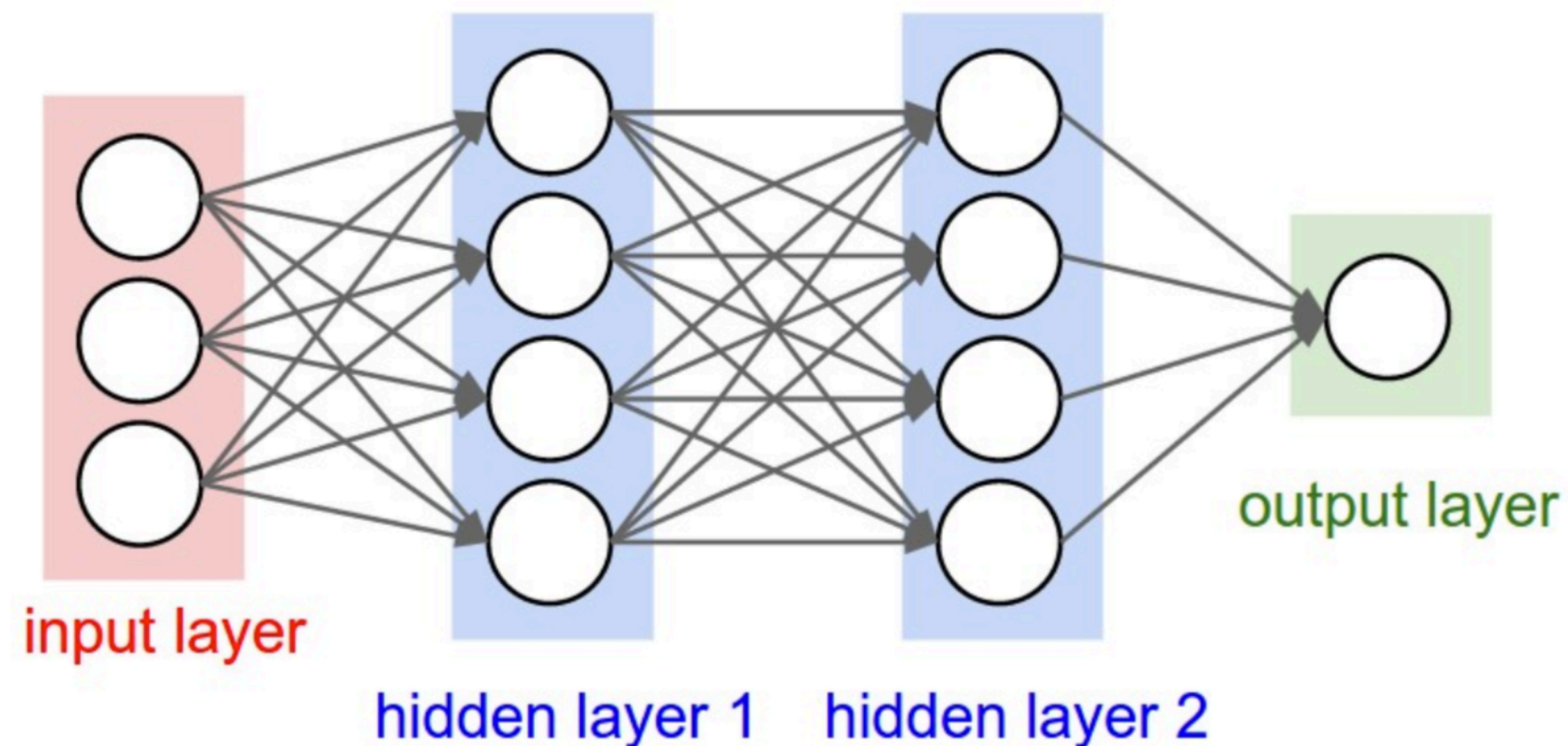
ReLU
(rectified linear unit)

$$f(z) = \max(0, z)$$



$$f'(z) = \begin{cases} 1 & z > 0 \\ 0 & z < 0 \end{cases}$$

Matrix notations



*: f is applied element-wise

$$f([z_1, z_2, z_3]) = [f(z_1), f(z_2), f(z_3)]$$

C : number of classes

d : input dimension, d_1, d_2 : hidden dimensions

- Input layer: $\mathbf{x} \in \mathbb{R}^d$

- Hidden layer 1:

$$\mathbf{h}_1 = f(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) \in \mathbb{R}^{d_1}$$

$$\mathbf{W}^{(1)} \in \mathbb{R}^{d_1 \times d}, \mathbf{b}^{(1)} \in \mathbb{R}^{d_1}$$

- Hidden layer 2:

$$\mathbf{h}_2 = f(\mathbf{W}^{(2)}\mathbf{h}_1 + \mathbf{b}^{(2)}) \in \mathbb{R}^{d_2}$$

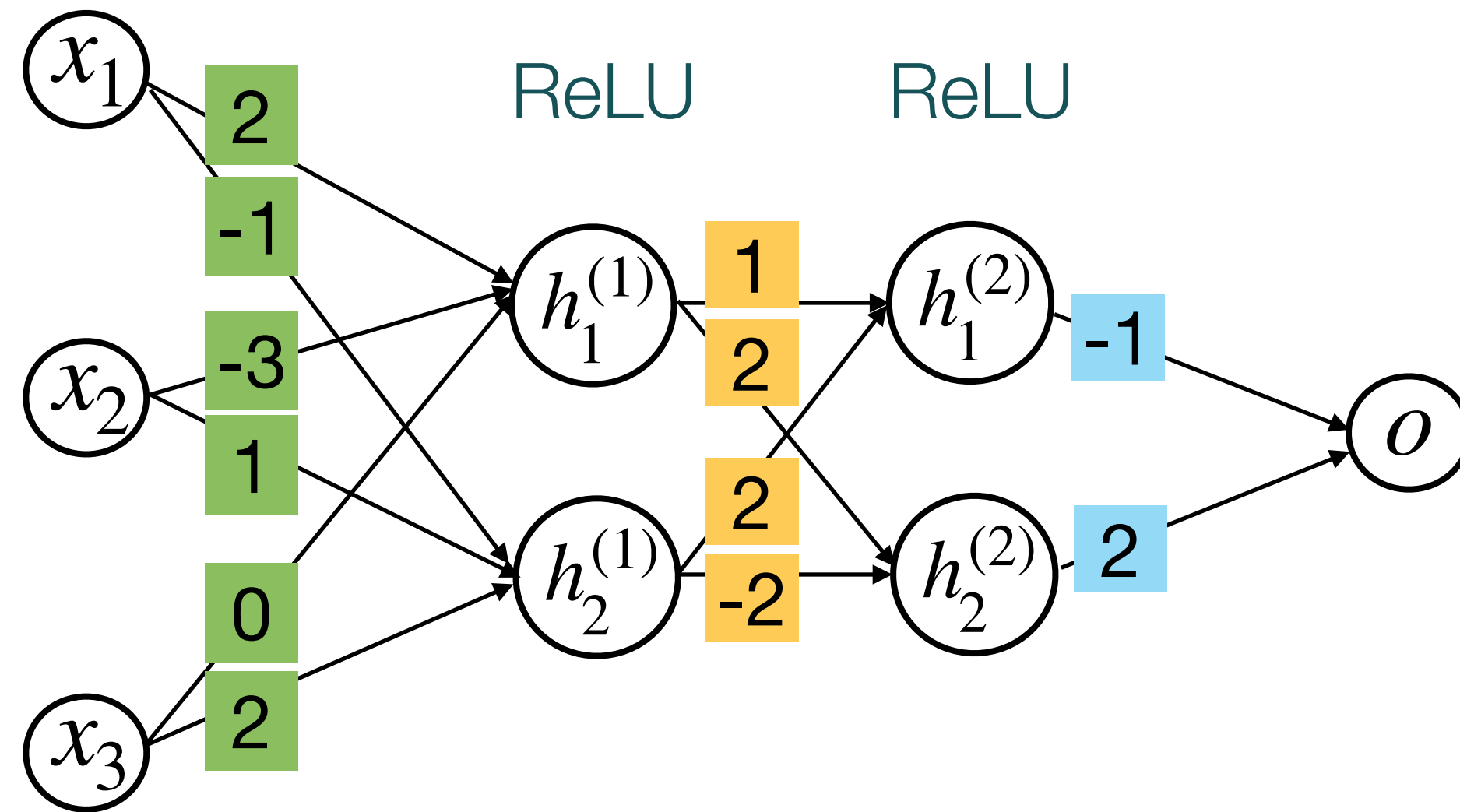
$$\mathbf{W}^{(2)} \in \mathbb{R}^{d_2 \times d_1}, \mathbf{b}^{(2)} \in \mathbb{R}^{d_2}$$

- Output layer:

$$\mathbf{y} = \mathbf{W}^{(o)}\mathbf{h}_2, \mathbf{W}^{(o)} \in \mathbb{R}^{C \times d_2}$$



Feedforward NNs



(Bias terms omitted in the next few slides)

For $x_1 = x_2 = x_3 = 1$, what is the value of $h_1^{(1)}$?

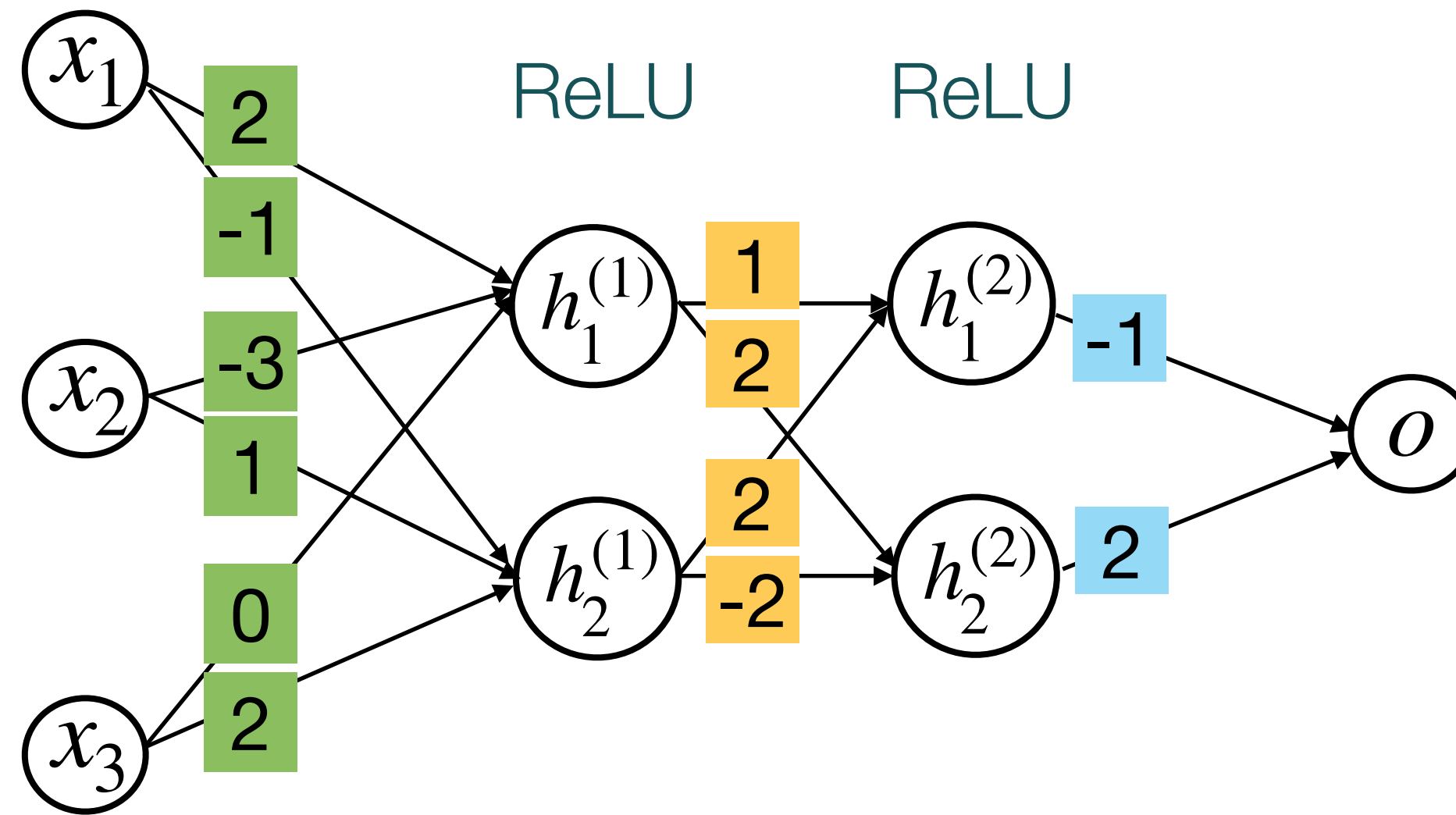
- (a) 0 (b) -1 (c) 1 (d) 2

Correct: (a), because of the ReLU:

$$\max(2 \times 1 + (-3) \times 1 + 0 \times 1, 0) = \max(-1, 0) = 0$$



Feedforward NNs



(Bias terms omitted in the next few slides)

$\mathbf{h}^{(1)} = \text{ReLU}(\mathbf{W}^{(1)}\mathbf{x})$ What is the matrix $\mathbf{W}^{(1)}$?

(a) $\begin{bmatrix} 2 & -1 \\ -3 & 1 \\ 0 & 2 \end{bmatrix}$

(b) $\begin{bmatrix} 2 & -3 & 0 \\ -1 & 1 & 2 \end{bmatrix}$

(c) $\begin{bmatrix} 1 & 2 \\ 2 & -2 \end{bmatrix}$

Correct: (b). $\mathbf{W}^{(1)}$ is a 2 x 3 matrix.

Feedforward NNs for multi-class classification

$$\mathbf{y} = \mathbf{W}^{(o)} \mathbf{h}_2, \mathbf{W}^{(o)} \in \mathbb{R}^{C \times d_2}$$

$$\hat{\mathbf{y}} = \text{softmax}(\mathbf{y}) \quad \text{softmax}(\mathbf{y})_k = \frac{\exp(y_k)}{\sum_{j=1}^C \exp(y_j)} \quad \mathbf{y} = [y_1, y_2, \dots, y_C]$$

Training loss:

$$\min_{\mathbf{W}^{(1)}, \mathbf{W}^{(2)}, \mathbf{W}^{(o)}} - \sum_{(\mathbf{x}, y) \in D} \log \hat{y}_y$$

Training feedforward NNs:
stochastic gradient descent!

Neural networks are difficult to optimize.
SGD can only converge to local minimum.
Initializations and optimizers matter a lot!

Back-propagation

Forward propagation:
from input to output layer

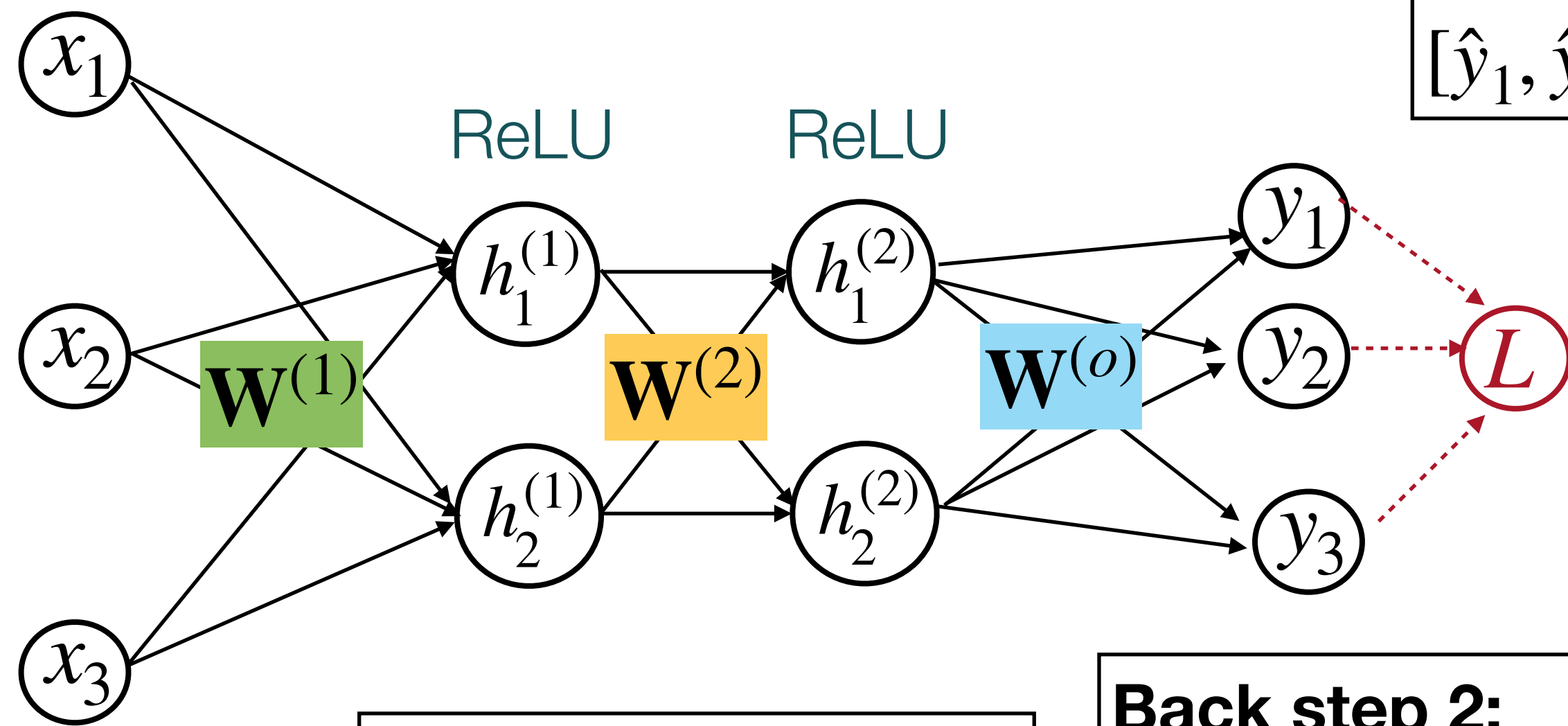
Given: x_1, x_2, x_3
and the class
label y
(a single training
example)

Forward step 1:
Compute $h_1^{(1)}, h_2^{(1)}$

Forward step 2:
Compute $h_1^{(2)}, h_2^{(2)}$

Forward step 3:
Compute y_1, y_2, y_3 and
 $[\hat{y}_1, \hat{y}_2, \hat{y}_3] = \text{softmax}[y_1, y_2, y_3]$

Forward step 4:
Compute loss
 $L = -\log \hat{y}_y$



Goal:
 $\frac{\partial L}{\partial W^{(1)}}$,
 $\frac{\partial L}{\partial W^{(2)}}$,
 $\frac{\partial L}{\partial W^{(o)}}$

Back step 4:
Compute
 $\frac{\partial L}{\partial W^{(1)}}$

Back step 3:
Compute
 $\frac{\partial L}{\partial h_1^{(1)}}, \frac{\partial L}{\partial h_2^{(1)}}, \frac{\partial L}{\partial W^{(2)}}$

Back step 2:
Compute
 $\frac{\partial L}{\partial h_1^{(2)}}, \frac{\partial L}{\partial h_2^{(2)}}, \frac{\partial L}{\partial W^{(o)}}$

Back step 1:
Compute
 $\frac{\partial L}{\partial y_1}, \frac{\partial L}{\partial y_2}, \frac{\partial L}{\partial y_3}$


Back propagation:
from output to input layer

Back-propagation in PyTorch

```
1 import torch.nn as nn
2 import torch.nn.functional as F
3
4 class Net(nn.Module):
5     def __init__(self):
6         super().__init__()
7         self.fc1 = nn.Linear(784, 128)
8         self.fc2 = nn.Linear(128, 64)
9         self.fc3 = nn.Linear(64, 10)
10
11     def forward(self, x):
12         x = F.relu(self.fc1(x))
13         x = F.relu(self.fc2(x))
14         x = self.fc3(x)
15         return x
16
```

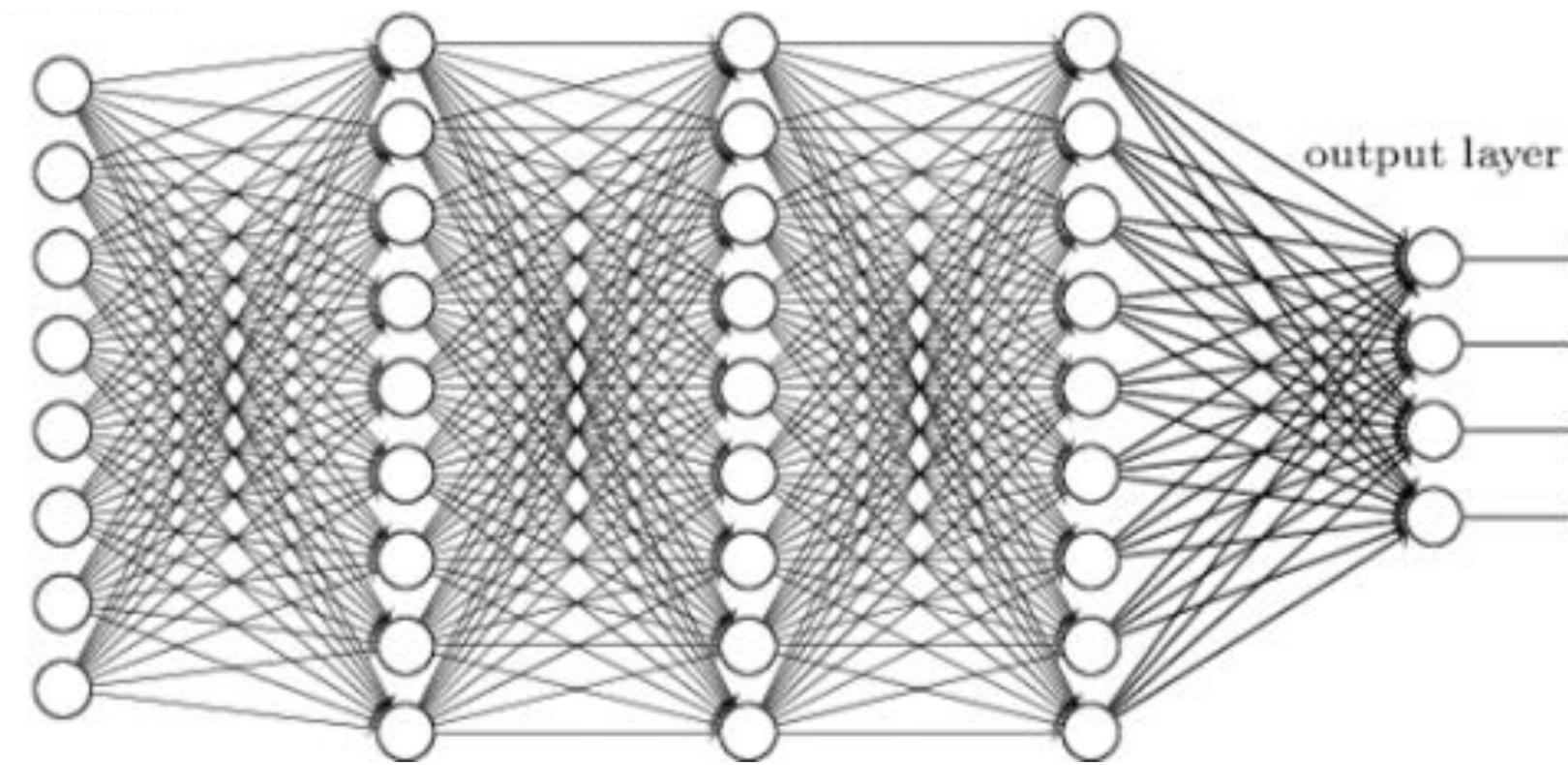
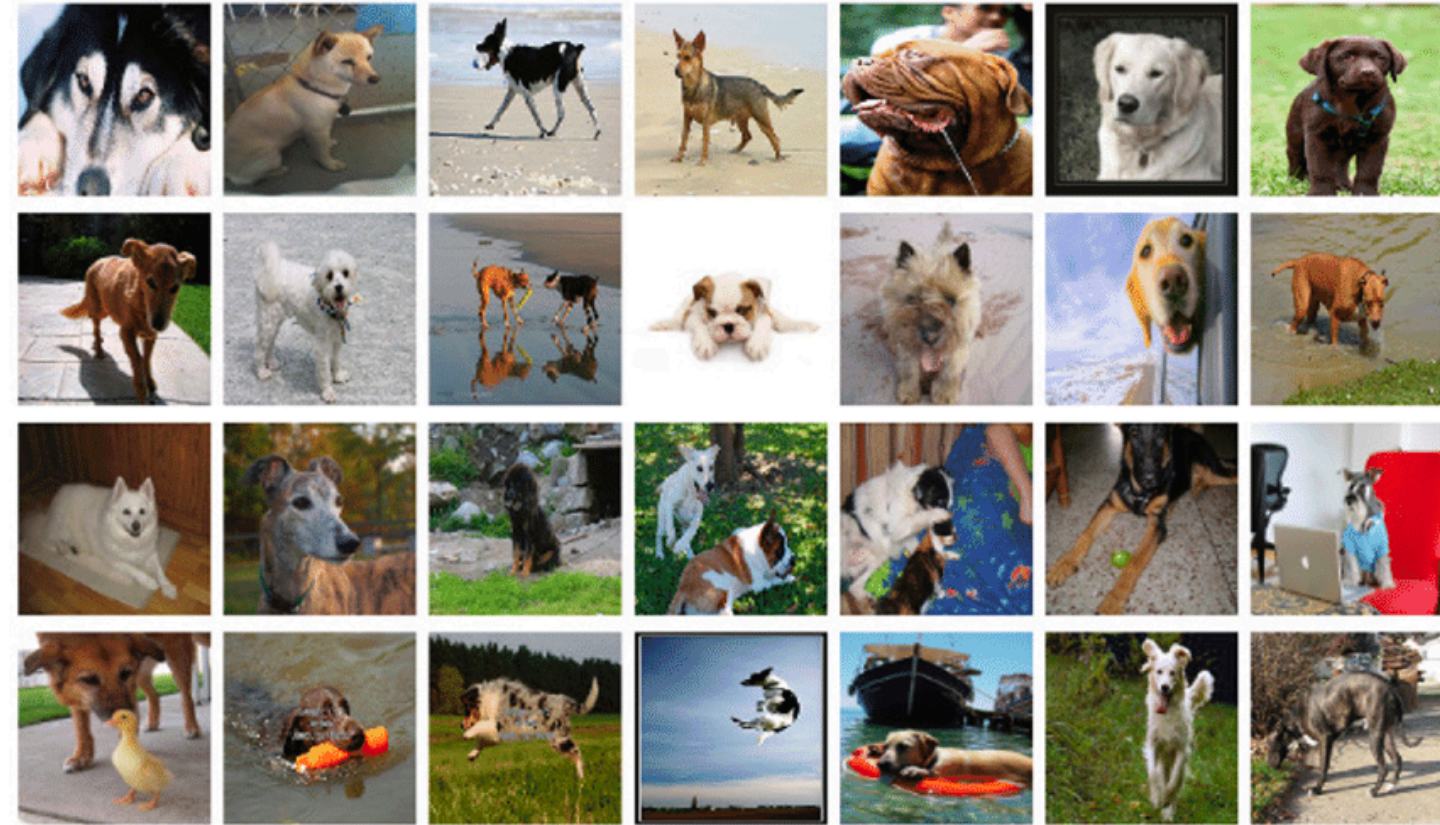
```
1 import torch.optim as optim
2
3 net = Net()
4 criterion = nn.CrossEntropyLoss()
5 optimizer = optim.SGD(net.parameters(), lr=0.001, momentum=0.9)
```

```
1 outputs = net(inputs)
2 loss = criterion(outputs, labels)
3 loss.backward()
4 optimizer.step()
```



PyTorch did back-propagation for you in this one line of code!

Comparison: image vs text inputs



label = "dog"

label = positive

a sometimes tedious film
i had to look away - this was god awful .
a gorgeous , witty , seductive movie .

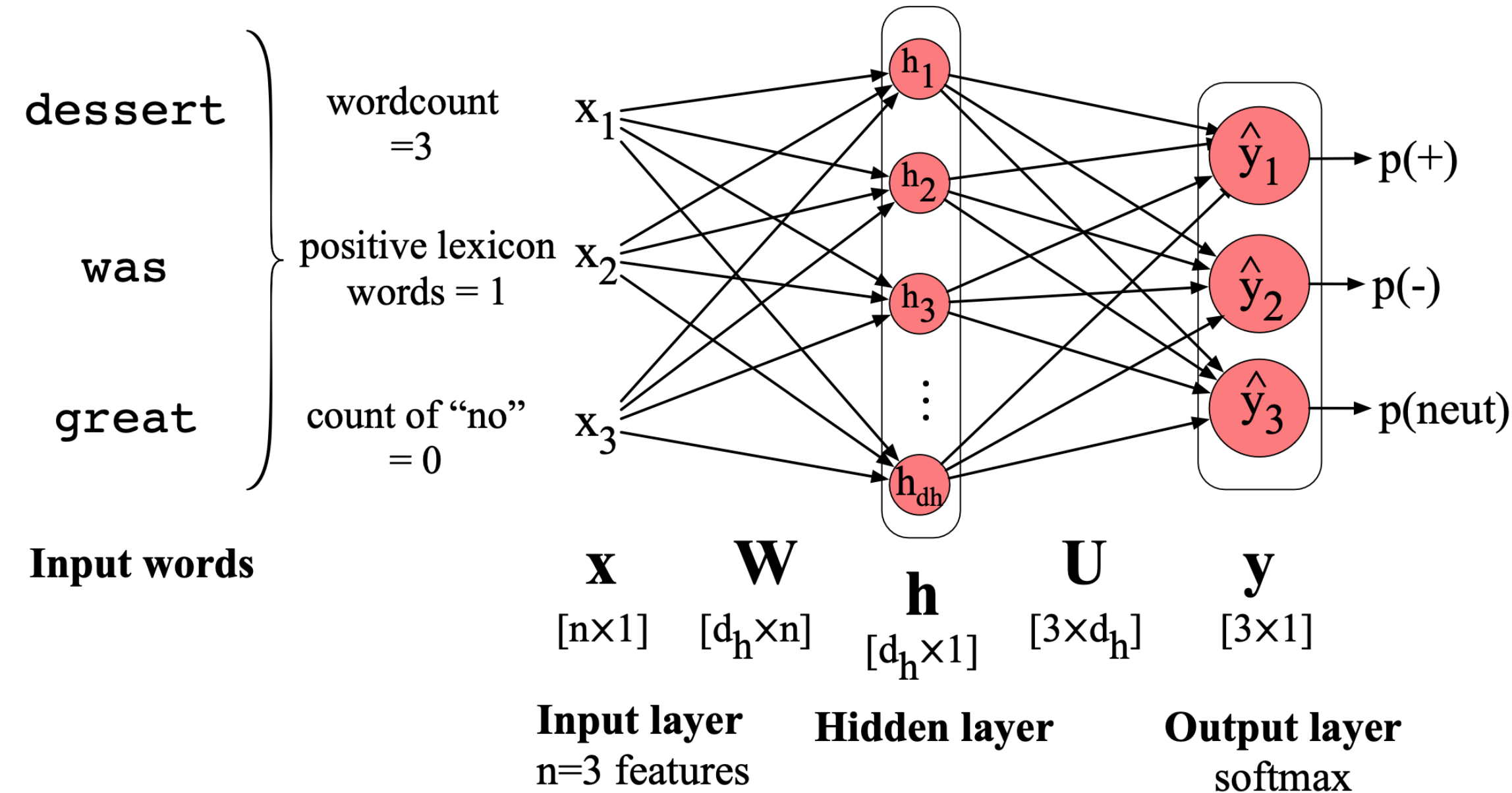
- Images: fixed-size input, continuous values
- Text: **variable-length** input, discrete words
 - need to convert into vectors - word embeddings!

Neural “bag-of-words” models for text classification

Neural networks for text classification

- Input: $w_1, w_2, \dots, w_K \in V$
- Output: $y \in C$
- Input: dessert was great
- Output: positive $C = \{\text{positive, negative, neutral}\}$

Solution #1: You can construct a feature vector \mathbf{x} from the input and simply feed the vector to a **neural network**, instead of a **logistic regression classifier**!



(each \mathbf{x}_i is a hand-designed feature)

- $\mathbf{x} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]$
- $\mathbf{h} = \text{ReLU}(\mathbf{W}\mathbf{x} + \mathbf{b})$
- $\mathbf{y} = \mathbf{U}\mathbf{h}$
- $\hat{\mathbf{y}} = \text{softmax}(\mathbf{y})$

Deep learning has the promise to learn good features automatically..

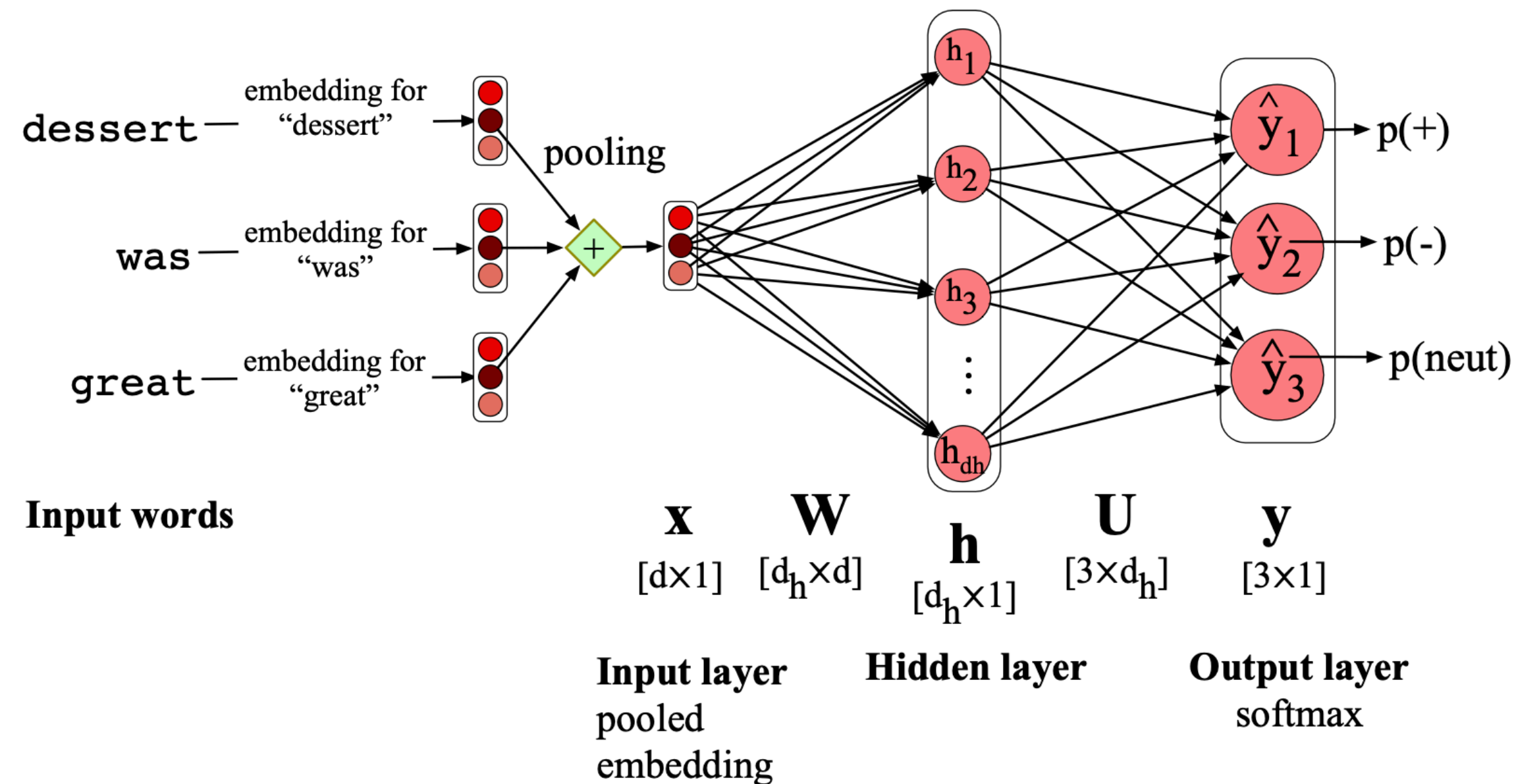
Neural networks for text classification

- How can we feed a **variable-length** input to a neural network classifier? $w_1, w_2, \dots, w_K \in V$

Solution #2: Let's take the all the word embeddings of these words and aggregate them into a vector through some **pooling** function!

$$\mathbf{x}_{\text{mean}} = \frac{1}{K} \sum_{i=1}^K e(w_i)$$

pooling: sum, mean or max



- $\mathbf{x} = \frac{1}{K} \sum_{i=1}^K e(w_i)$
- $\mathbf{h} = \text{ReLU}(\mathbf{W}\mathbf{x} + \mathbf{b})$
- $\mathbf{y} = \mathbf{U}\mathbf{h}$
- $\hat{\mathbf{y}} = \text{softmax}(\mathbf{y})$

Important note: each input has a different K

Neural networks for text classification

- (+): This provides a simple and flexible way to handle variable-length input
- (+): It learns feature representations automatically from the data
- (+): It can generalize to similar inputs through word embeddings
- (-): The model throws away any sequential information of the text

neural bag-of-words model (NBOW)

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!

15



it	6
I	5
the	4
to	3
and	3
seen	2
yet	1
would	1
whimsical	1
times	1
sweet	1
satirical	1
adventure	1
genre	1
fairy	1
humor	1
have	1
great	1
...	...

How to train this model?

- $\mathbf{x} = \frac{1}{K} \sum_{i=1}^K e(w_i)$

- $\mathbf{h} = \text{ReLU}(\mathbf{W}\mathbf{x} + \mathbf{b})$

- $\mathbf{y} = \mathbf{U}\mathbf{h}$

- $\hat{\mathbf{y}} = \text{softmax}(\mathbf{y})$

- Training data: $\{(d^{(1)}, y^{(1)}), \dots, (d^{(m)}, y^{(m)})\}$

- Parameters: $\{\mathbf{W}, \mathbf{b}, \mathbf{U}\}$

- Optimize these parameters using stochastic gradient descent!

- Word embeddings can be treated as parameters too!

$$\mathbf{E} \in \mathbb{R}^{|V| \times d}$$

How to train this model?

- Common practice: initialize \mathbf{E} using word embeddings (e.g. word2vec), and optimize them using SGD!
- When the training data is small, don't treat \mathbf{E} as parameters!
- When the training data is very large (e.g., language modeling), initialization doesn't matter much either (= can use random initialization)

Why? $\mathbf{v}(\text{good}) \approx \mathbf{v}(\text{bad})$

	Most Similar Words for	
	Static	Non-static
bad	<i>good</i> <i>terrible</i> <i>horrible</i> <i>lousy</i>	<i>terrible</i> <i>horrible</i> <i>lousy</i> <i>stupid</i>
good	<i>great</i> <i>bad</i> <i>terrific</i> <i>decent</i>	<i>nice</i> <i>decent</i> <i>solid</i> <i>terrific</i>

(Kim 2014)

Feedforward neural language models

N-gram vs neural language models

Language models: Given $x_1, x_2, \dots, x_n \in V$, the goal is to model:

$$P(x_1, x_2, \dots, x_n) = \prod_{i=1}^n P(x_i | x_1, \dots, x_{i-1})$$

Bigram: $P(x_1, x_2, \dots, x_n) = \prod_{i=1}^n P(x_i | x_{i-1})$

Maximum likelihood estimate:

Trigram: $P(x_1, x_2, \dots, x_n) = \prod_{i=1}^n P(x_i | x_{i-2}, x_{i-1})$

$$P(\text{sat} | \text{the cat}) = \frac{\text{count}(\text{the cat sat})}{\text{count}(\text{the cat})}$$

Limitations? Can't handle long histories!

As the proctor started the clock, the students opened their _____

The **keys** to the cabinet is/are

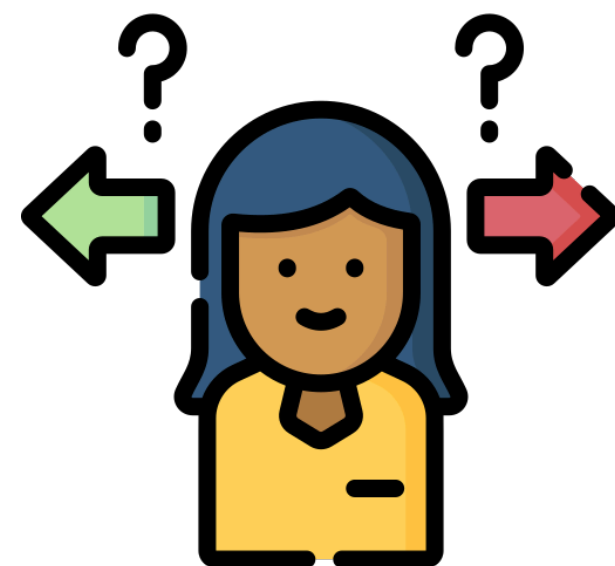
N-gram vs neural language models

- If we use a 4-gram, 5-gram, 6-gram language model, it will become too sparse to estimate the probabilities:

$$P(w \mid \text{students opened their}) = \frac{\text{count}(\text{students opened their } w)}{\text{count}(\text{students opened their})}$$

Dilemma:

- We need to model bigger context!
- The # of probabilities that we need to estimate grow exponentially with window size!



- A lot of contexts are similar and simply counting them won't generalize

I am a **good** _____ count(I am a good w)

I am a **great** _____ count(I am a great w)

$e(\text{good}) \approx e(\text{great})$

Can we estimate the probabilities better?

Feedforward neural language models

A Neural Probabilistic Language Model (Bengio et al., 2003)



Yoshua Bengio
Réjean Ducharme
Pascal Vincent
Christian Jauvin

BENGIOY@IRO.UMONTREAL.CA
DUCHARME@IRO.UMONTREAL.CA
VINCENTP@IRO.UMONTREAL.CA
JAUVINC@IRO.UMONTREAL.CA

Yoshua Bengio

Probabilistic models of sequences: In the 1990s, Bengio combined neural networks with probabilistic models of sequences, such as hidden Markov models. These ideas were incorporated into a system used by AT&T/NCR for reading handwritten checks, were considered a pinnacle of neural network research in the 1990s, and modern deep learning speech recognition systems are extending these concepts.

High-dimensional word embeddings and attention: In 2000, Bengio authored the landmark paper, "A Neural Probabilistic Language Model," that introduced high-dimension word embeddings as a representation of word meaning. Bengio's insights had a huge and lasting impact on natural language processing tasks including language translation, question answering, and visual question answering. His group also introduced a form of attention mechanism which led to breakthroughs in machine translation and form a key component of sequential processing with deep learning.

Generative adversarial networks: Since 2010, Bengio's papers on generative deep learning, in particular the Generative Adversarial Networks (GANs) developed with Ian Goodfellow, have spawned a revolution in computer vision and computer graphics. In one fascinating application of this work, computers can actually create original images, reminiscent of the creativity that is considered a hallmark of human intelligence.

<https://awards.acm.org/about/2018-turing>

Feedforward neural language models

A Neural Probabilistic Language Model (Bengio et al., 2003)



Yoshua Bengio
Réjean Ducharme
Pascal Vincent
Christian Jauvin

BENGIOY@IRO.UMONTREAL.CA
DUCHARME@IRO.UMONTREAL.CA
VINCENTP@IRO.UMONTREAL.CA
JAUVINC@IRO.UMONTREAL.CA

Key idea: Instead of estimating raw probabilities, let's use a **neural network** to fit the **probabilistic distribution of language!**

$$P(w \mid \text{I am a good})$$

$$P(w \mid \text{I am a great})$$

Key ingredient: word embeddings $\mathbf{e}(\text{good}) \approx \mathbf{e}(\text{great})$

Hope: this would give us similar distributions for similar contexts!

Feedforward neural language models

- Feedforward neural language models approximate the probability based on the previous m (e.g., 5) words - m is a hyper-parameter!

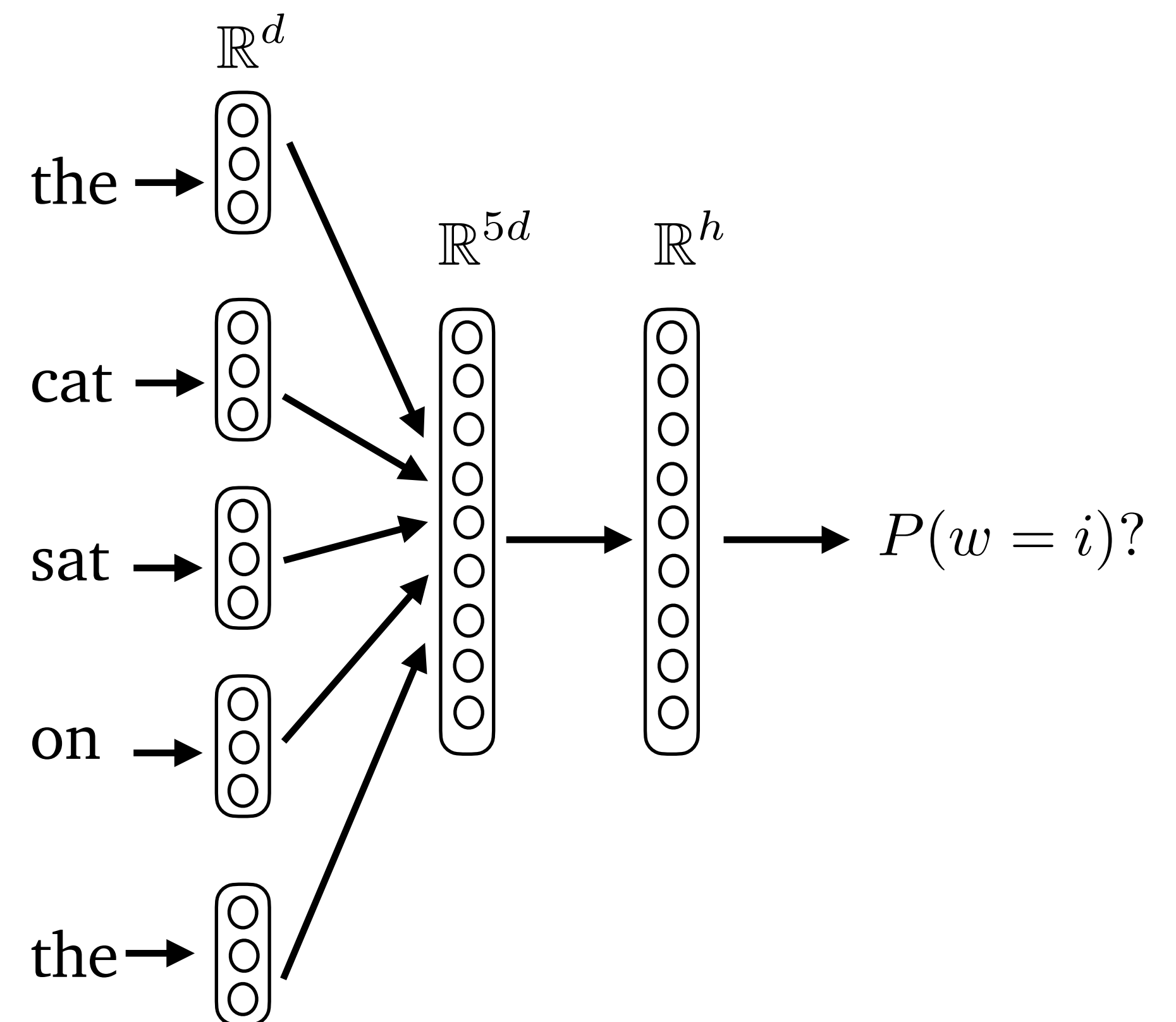
$$P(x_1, x_2, \dots, x_n) \approx \prod_{i=1}^n P(x_i | x_{i-m+1}, \dots, x_{i-1})$$

$P(\text{mat} | \text{the cat sat on the}) = ?$

d : word embedding size

h : hidden size

It is a $|V|$ -way classification problem!



Feedforward neural language models

$P(\text{mat} \mid \text{the cat sat on the}) = ?$ d : word embedding size h : hidden size

- Input layer ($m=5$): **Q: why concat instead of taking the average?**

$$\mathbf{x} = [e(\text{the}); e(\text{cat}); e(\text{sat}); e(\text{on}); e(\text{the})] \in \mathbb{R}^{md}$$

- Hidden layer:

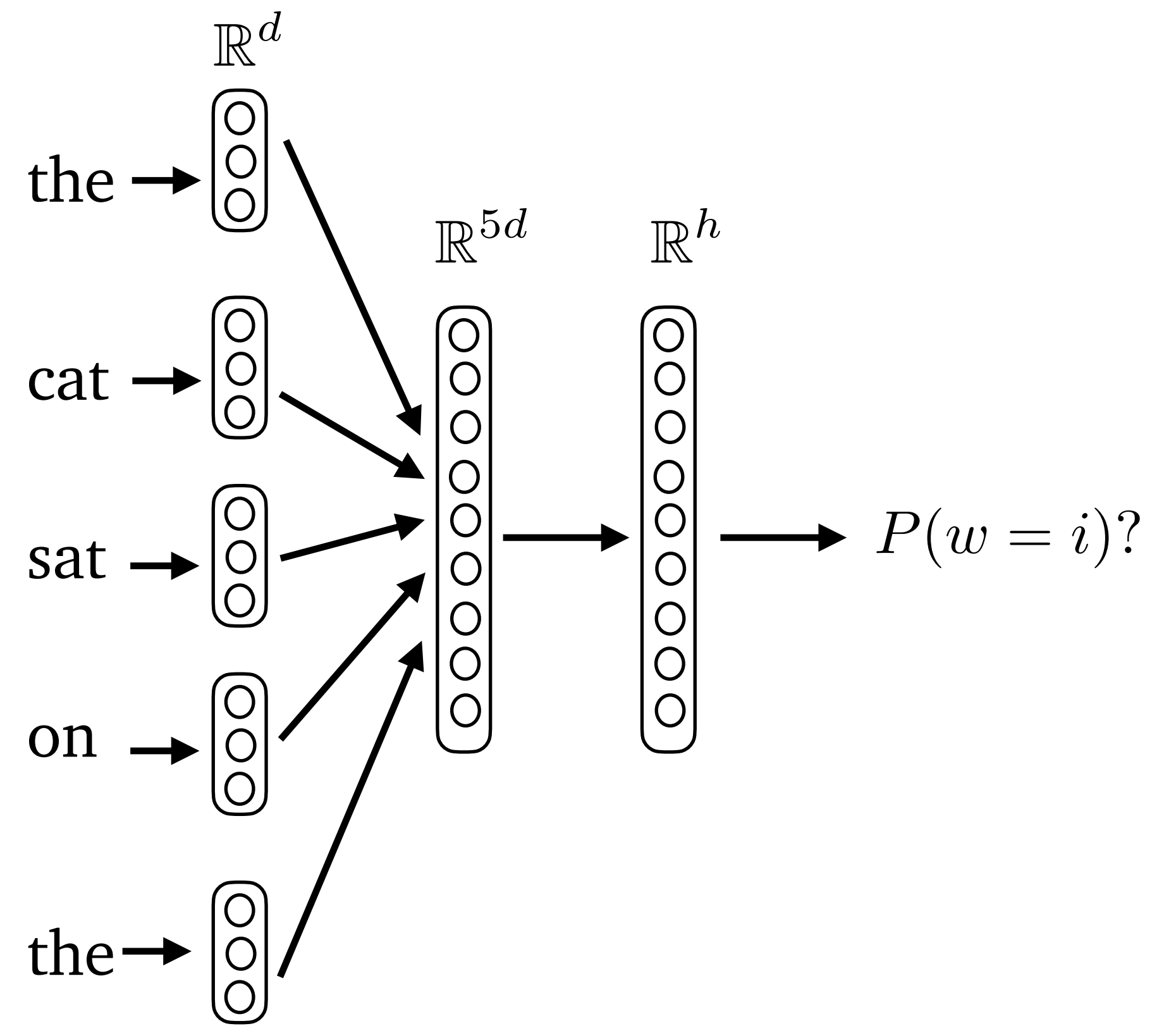
$$\mathbf{h} = \tanh(\mathbf{W}\mathbf{x} + \mathbf{b}) \in \mathbb{R}^h$$

- Output layer

$$\mathbf{z} = \mathbf{U}\mathbf{h} \in \mathbb{R}^{|V|}$$

$$P(w = i \mid \text{the cat sat on the})$$

$$= \text{softmax}_i(\mathbf{z}) = \frac{e^{z_i}}{\sum_k e^{z_k}}$$





What are the dimensions of W and U ?

d : word embedding size, h : hidden size

- (a) $W \in \mathbb{R}^{h \times d}, U \in \mathbb{R}^{|V| \times h}$
- (b) $W \in \mathbb{R}^{h \times 5d}, U \in \mathbb{R}^{|V| \times h}$
- (c) $W \in \mathbb{R}^{h \times 5d}, U \in \mathbb{R}^{|V| \times d}$
- (d) $W \in \mathbb{R}^{h \times d}, U \in \mathbb{R}^{d \times h}$

Correct: (b)

- Input layer ($m=5$):

$$\mathbf{x} = [e(\text{the}); e(\text{cat}); e(\text{sat}); e(\text{on}); e(\text{the})] \in \mathbb{R}^{md}$$

- Hidden layer:

$$\mathbf{h} = \tanh(\mathbf{W}\mathbf{x} + b) \in \mathbb{R}^h$$

- Output layer

$$\mathbf{z} = \mathbf{U}\mathbf{h} \in \mathbb{R}^{|V|}$$

$$P(w = i \mid \text{the cat sat on the})$$

$$= \text{softmax}_i(\mathbf{z}) = \frac{e^{z_i}}{\sum_k e^{z_k}}$$

Feedforward neural language models

- How to train this model? A: Use a lot of raw text to create training examples and run gradient-descent optimization!

The Fat Cat Sat on the Mat is a 1996 children's book by Nurit Karlin. Published by Harper Collins as part of the reading readiness program, the book stresses the ability to read words of specific structure, such as -at.



the fat cat sat on → the
fat cat sat on the → mat
cat sat on the mat → is
sat on the mat is → a
...

- Limitations?
 - **W linearly** scales with the context size m
 - The model learns separate patterns for different positions!
- Better solutions: recurrent NNs, Transformers..

the fat cat **sat on** → the
fat cat **sat on** the → mat
cat **sat on** the mat → is

“sat on” corresponds to different parameters in **W**