



University of Colorado **Boulder**

Department of Computer Science

CSCI 5622: Machine Learning

Chenhao Tan

Lecture 12: Neural Networks II

Slides adapted from Chris Ketelsen, Jordan Boyd-Graber,
and Noah Smith

Administrivia

- HW2 solution is out
- HW2 feedback
- Google educational credit information in email

Learning Objectives

- Understand common components in neural architectures

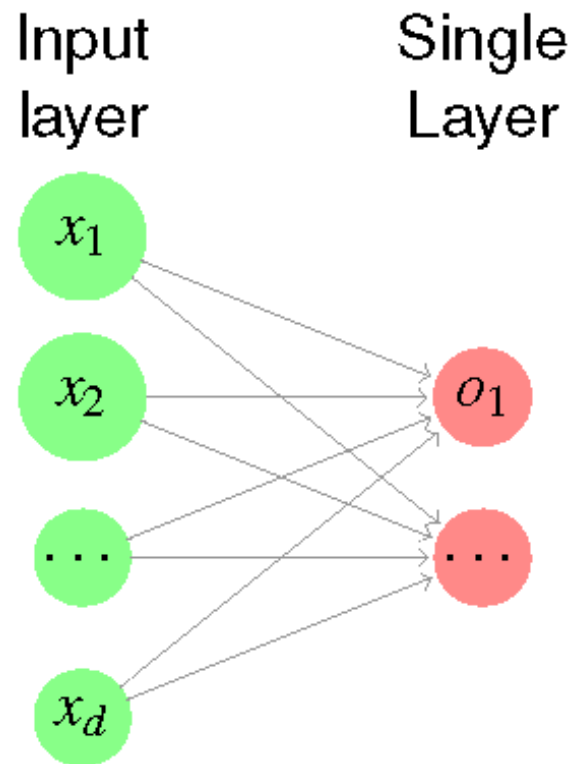
Outline

- Recap
- Layers for structured data

Outline

- Recap
- Layers for structured data

Logistic Regression as a Single-layer Neural Network

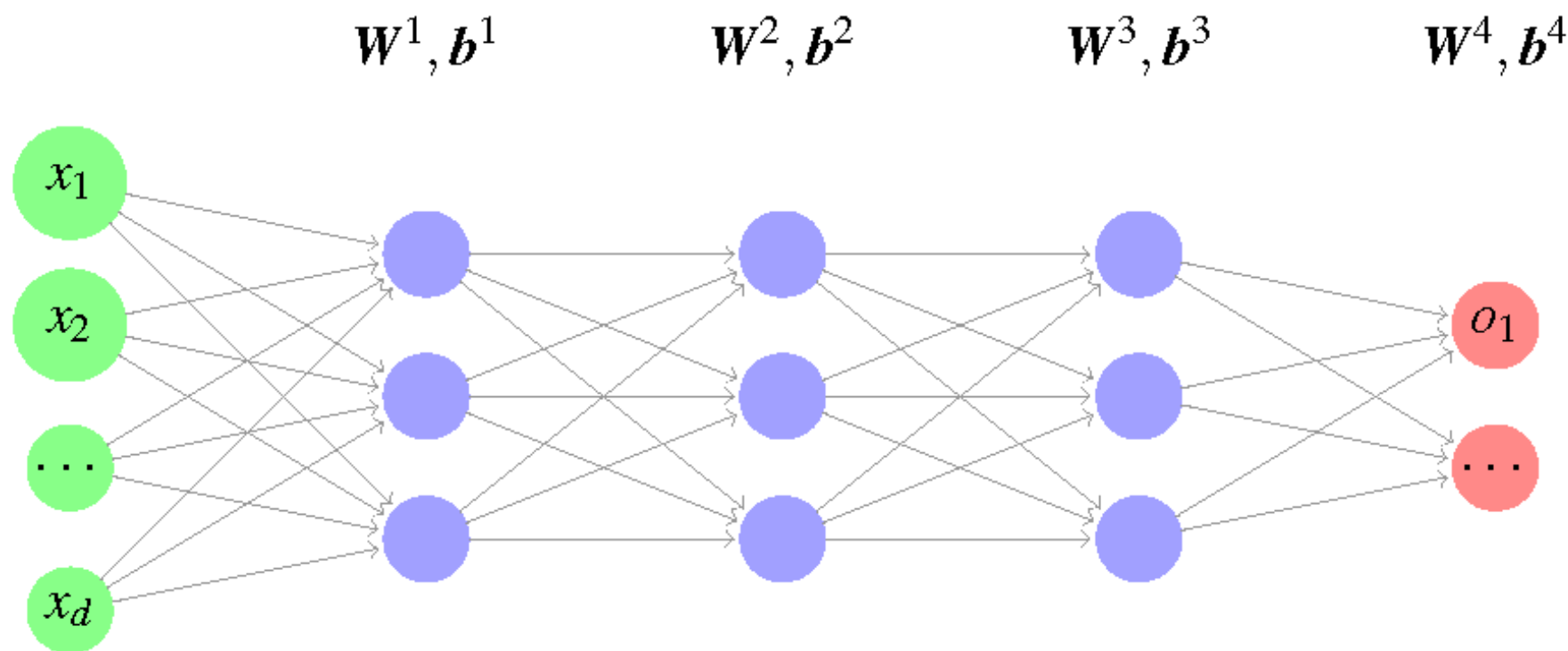


- What is the activation used in logistic regression? *softmax*
- What is the objective function used in logistic regression? *cross entropy*

Forward propagation algorithm

How do we make predictions based on a multi-layer neural network?

Store the biases for layer l in b^l , weight matrix in W^l



Forward propagation algorithm

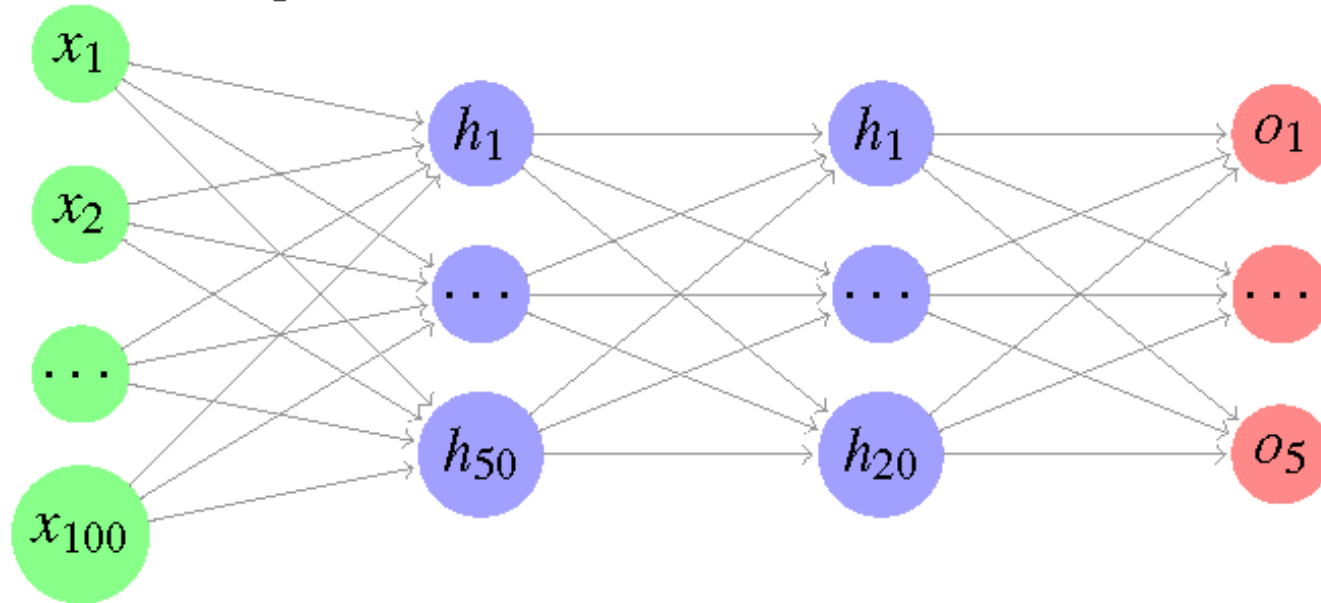
Suppose your network has L layers

Make prediction for an instance x

- 1: Initialize $a^0 = x$
- 2: **for** $l = 1$ to L **do**
- 3: $z^l = W^l a^{l-1} + b^l$
- 4: $a^l = g(z^l)$ // g represents the nonlinear activation
- 5: **end for**
- 6: The prediction \hat{y} is simply $\underline{a^L}$

Quiz 1

How many parameters are there in the following feed-forward neural networks (assuming no biases)?



- A. $100 * 50 + 50 * 20 + 20 * 5$ = $5000 + 1000 + 100$ ≈ 6100
- B. $100 * 50 + 50 + 50 * 20 + 20 + 20 * 5 + 5$

Neural networks in a nutshell

- Training data $S_{\text{train}} = \{(\mathbf{x}, y)\}$
- Network architecture (model)
- Loss function (objective function)

$$\hat{y} = f_w(\mathbf{x})$$

Handwritten red annotations: a red 'a' with a superscript 'L' above it, and a red 'T' below the function f_w .

$$\underline{\mathcal{L}(y, \hat{y})}$$

- Learning (next week)

Quiz 2

Which of the following statements is true? (Suppose that training data is large.)

- A. In training, K-nearest neighbors takes shorter time than neural networks. ✓
- B. In training, K-nearest neighbors takes longer time than neural networks.
- C. In testing, K-nearest neighbors takes shorter time than neural networks.
- D. In testing, K-nearest neighbors takes longer time than neural networks. ✓

Outline

- Recap
- Layers for structured data

Structured data

Spatial information



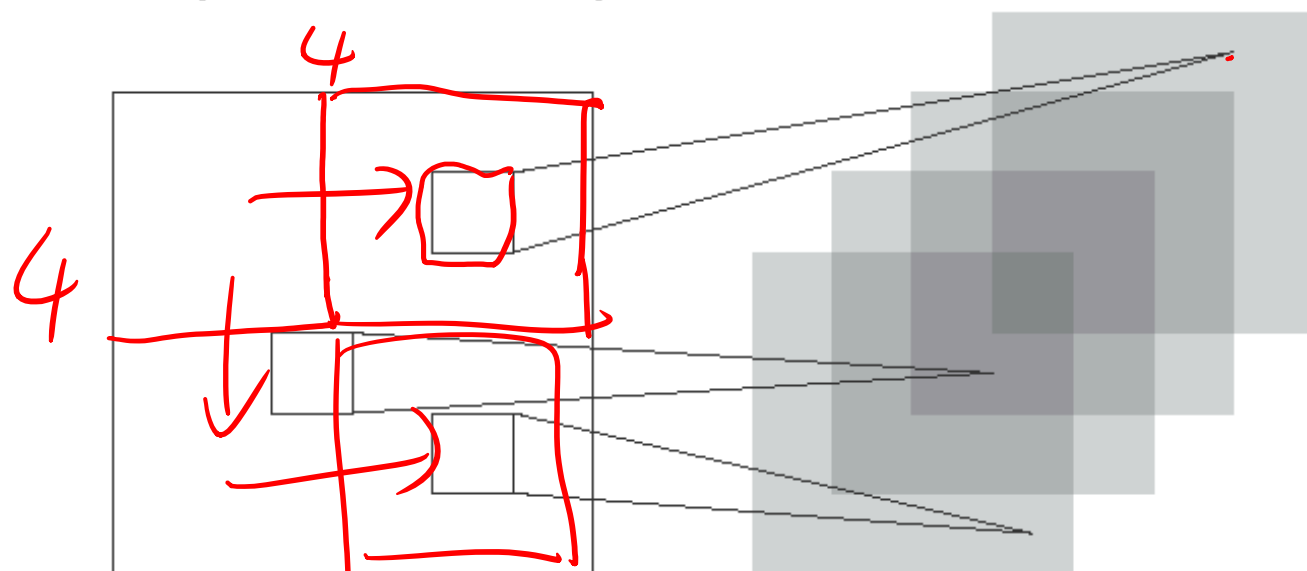
https://www.reddit.com/r/aww/comments/6ip21a/before_and_after_she_was_told_she_was_a_good_girl/

Convolutional Layers

Sharing parameters across patches

input image
or input feature map

output feature maps

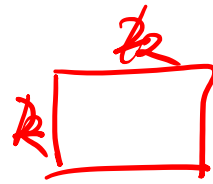


$$a_{i'j'} = \sum_{i=1}^k \sum_{j=1}^k w_{ij} x_{ij}$$

- Number of filters
- Filter shape $k \times k$
- Stride size

<https://github.com/davidstutz/latex-resources/blob/master/tikz-convolutional-layer/convolutional-layer.tex>

Quiz 3



$$\max_{1 \leq i \leq 2} x_{ij}$$
$$1 \leq j \leq 2$$

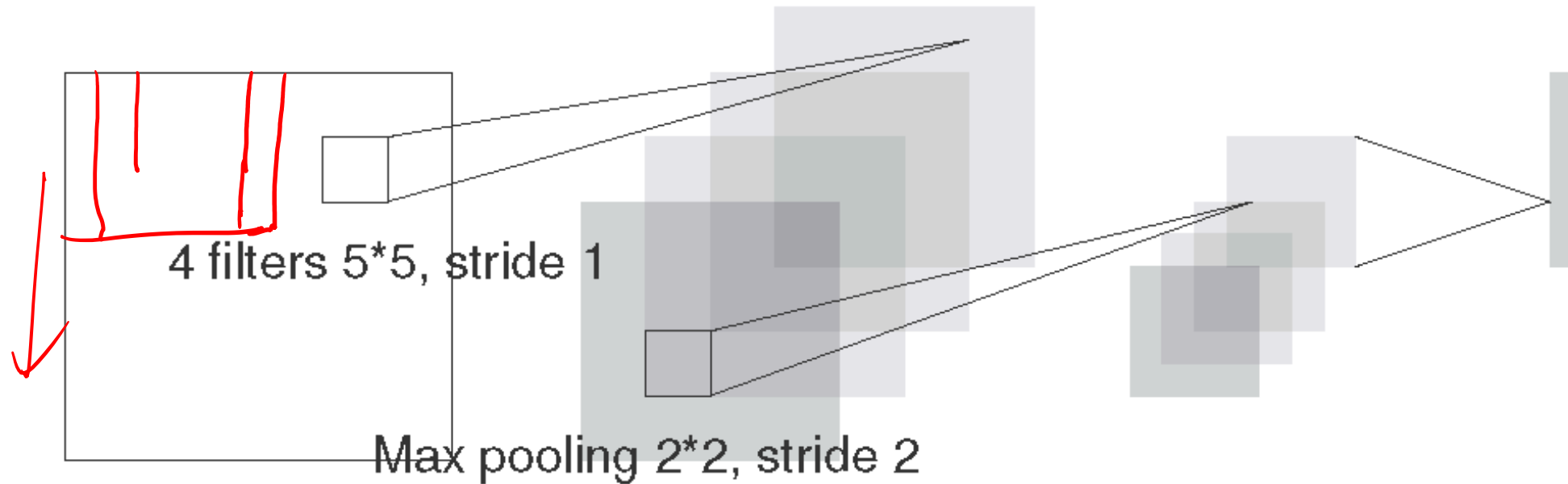
A concrete example (assuming no biases, convolution with 4 filters, ReLU, max pooling, and finally a fully-connected softmax layer)

input image (10*10)

4@6*6

4@3*3

5*1



Quiz 3

$\max(0, x)$

$\max_{ij} x_{ij}$

100
180

280

3800

How many parameters are there in the following convolutional neural networks?
(assuming no biases, convolution with 4 filters, ReLU, max pooling, and finally a
fully-connected softmax layer)

$4 \times 5 \times 5$
100

0

0

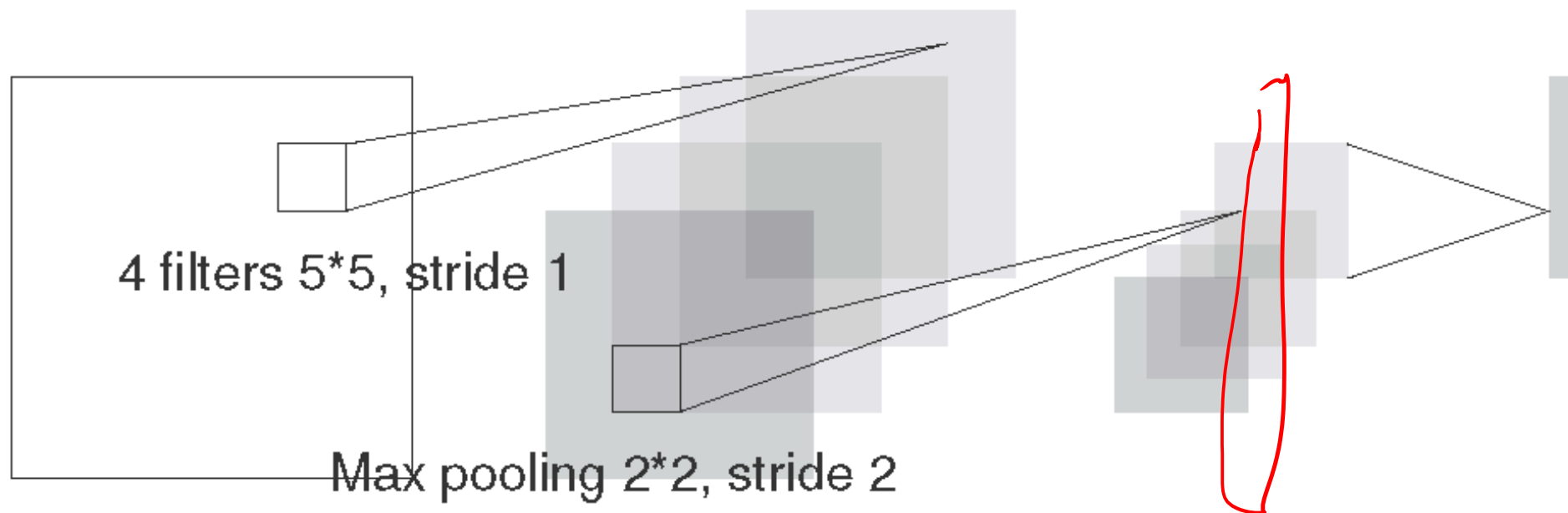
$(4 \times 3 \times 3 \times 5)$

input image (10*10)

4@6*6

4@3*3

5*1 (80)



Quiz 3

100
144
180

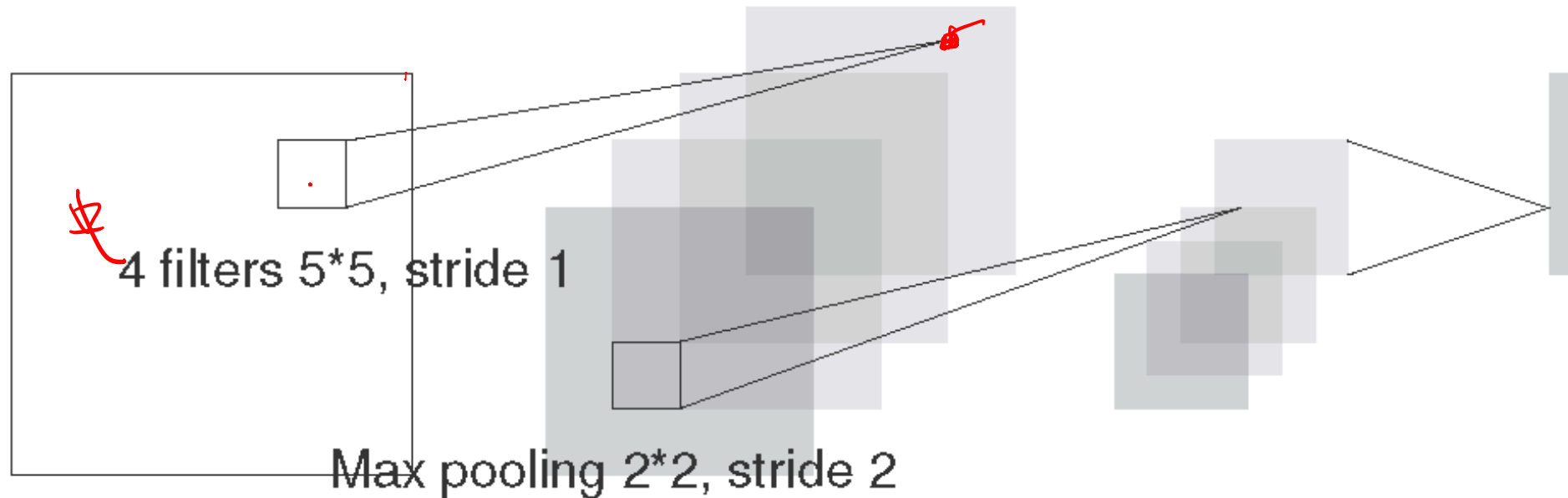
How many ReLU operations are performed on the forward pass? (assuming no biases, convolution with 4 filters, ReLU, max pooling, and finally a fully-connected softmax layer)

input image (10*10)

4@6*6

4@3*3

5*1



Structured data

Sequential information

“My words fly up, my thoughts remain below: Words without thoughts never to heaven go.”

—Hamlet

Structured data

Sequential information

“My words fly up, my thoughts remain below: Words without thoughts never to heaven go.”

—Hamlet

- language
- activity history

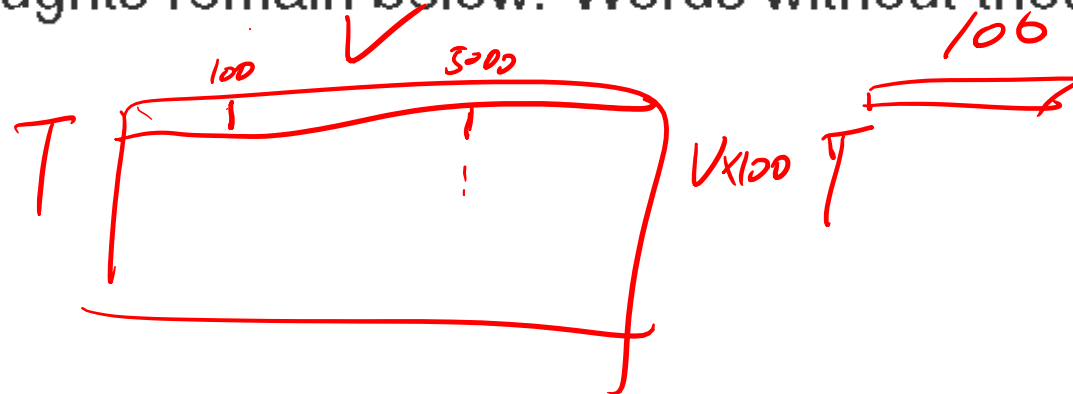
Structured data

(V) 50000 x 100



Sequential information

“My words fly up, my thoughts remain below: Words without thoughts never to heaven go.”



—Hamlet

- language
- activity history

$$\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_T)$$



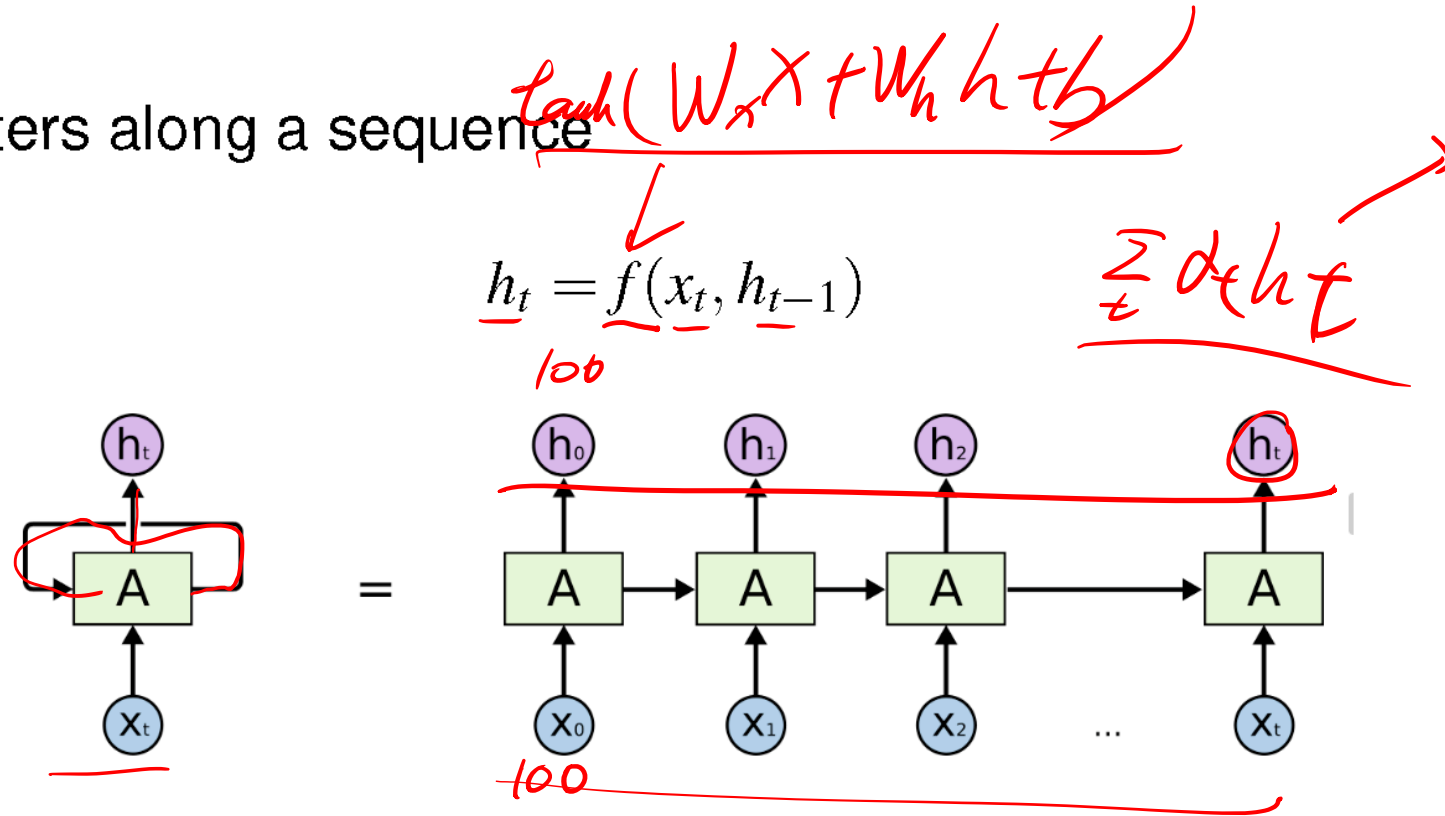
Recurrent Layers

Sharing parameters along a sequence

$$h_t = f(\underbrace{x_t}_{\uparrow}, \underbrace{h_{t-1}})$$

Recurrent Layers

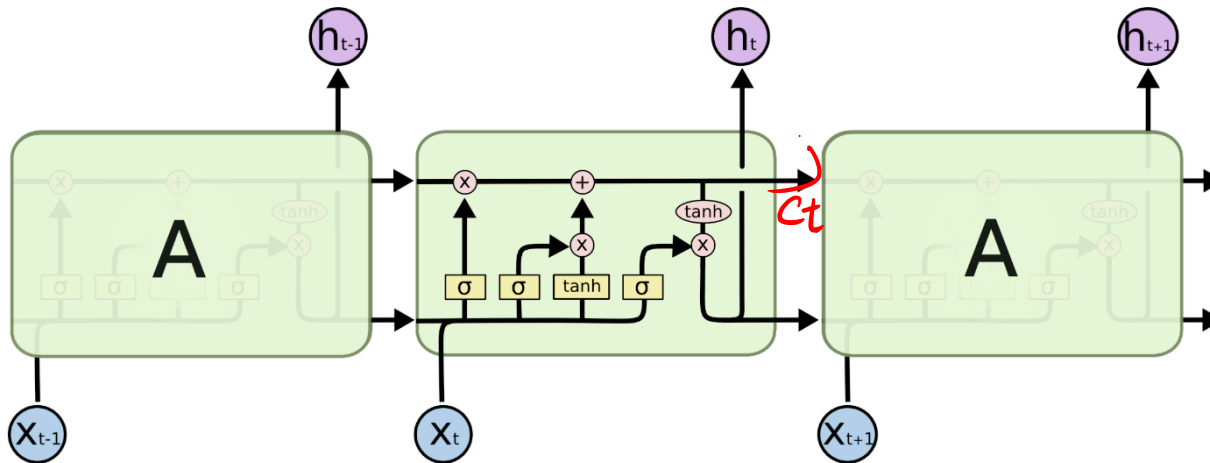
Sharing parameters along a sequence



<https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

Long short-term memory

A commonly used recurrent neural network in natural language processing



$$\begin{aligned} f_t &= \sigma(W_f \cdot [h_{t-1}, x_t] + b_f) \\ i_t &= \sigma(W_i \cdot [h_{t-1}, x_t] + b_i) \\ o_t &= \sigma(W_o \cdot [h_{t-1}, x_t] + b_o) \\ \tilde{C}_t &= \tanh(W_C \cdot [h_{t-1}, x_t] + b_C) \\ C_t &= f_t * C_{t-1} + i_t * \tilde{C}_t \\ h_t &= o_t * \tanh(C_t) \end{aligned}$$

Wrap up

Neural networks

- Network architecture (a lot more than fully connected layers)
 - Convolutional layer
 - Recurrent layer
- Loss function

What is missing?

- How to find good weights?
- How to make the model work (regularization, even more architecture, etc)?