

# CSE440: Natural Language Processing II

Farig Sadeque  
Assistant Professor  
Department of Computer Science and Engineering  
BRAC University

# Lecture 6: Neural Nets and RNN

# Outline

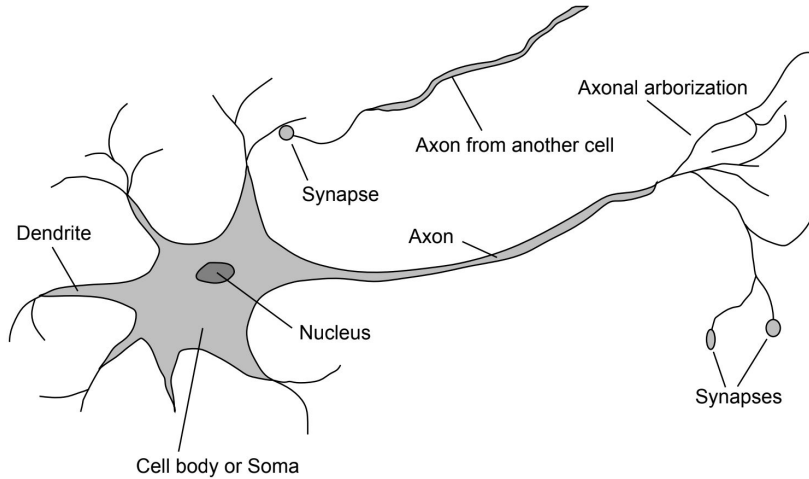
- Neural Networks (SLP 7 and lecture)
- Recurrent neural networks (SLP 9 and lecture)

# Before starting learning sequence ....

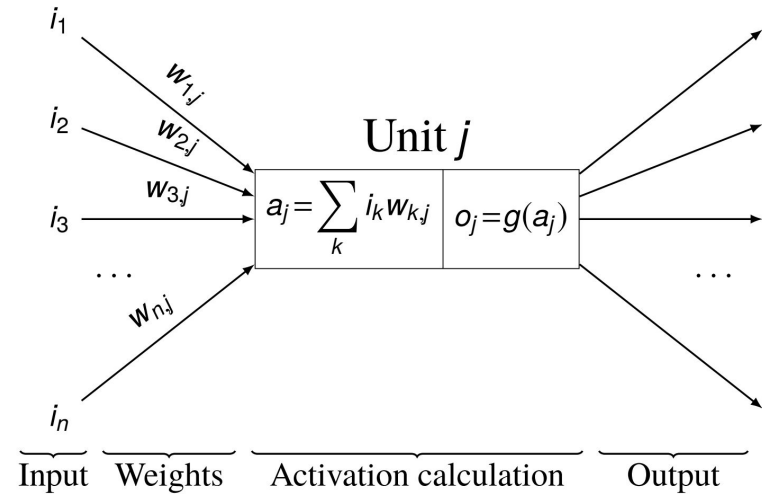
We need to remember some neural network basics.

# Neural Networks

## Neuron in a human brain

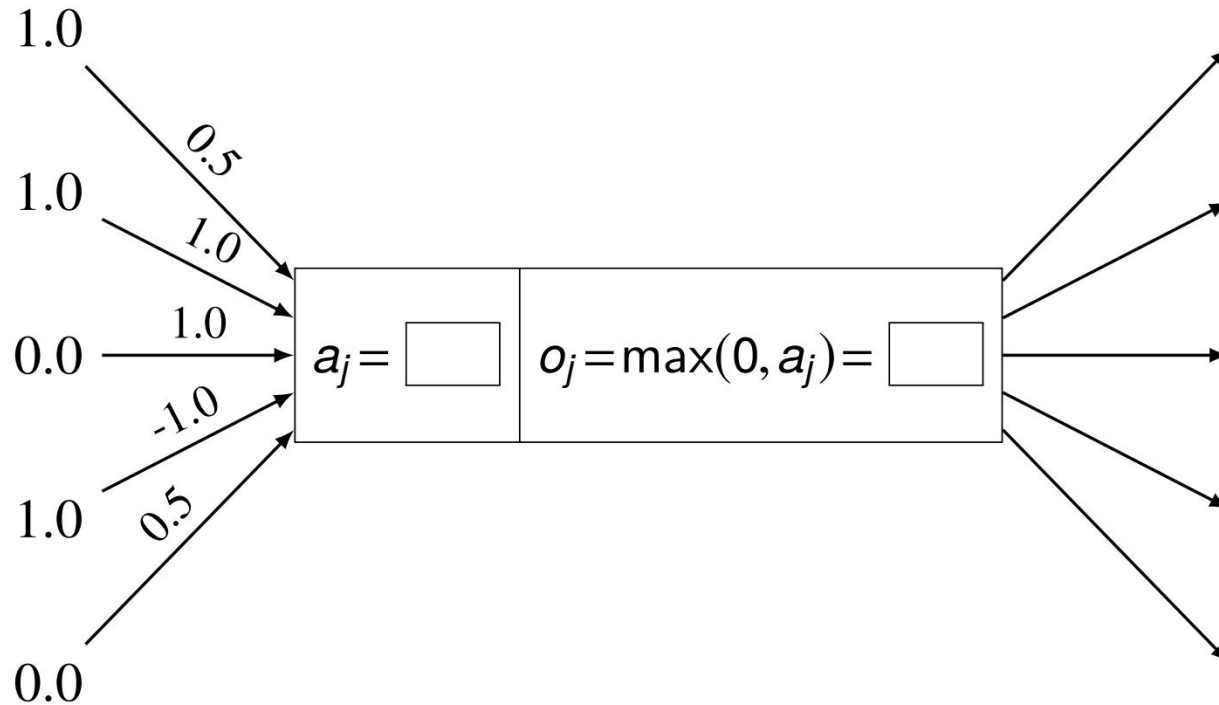


## Neuron in an ML model



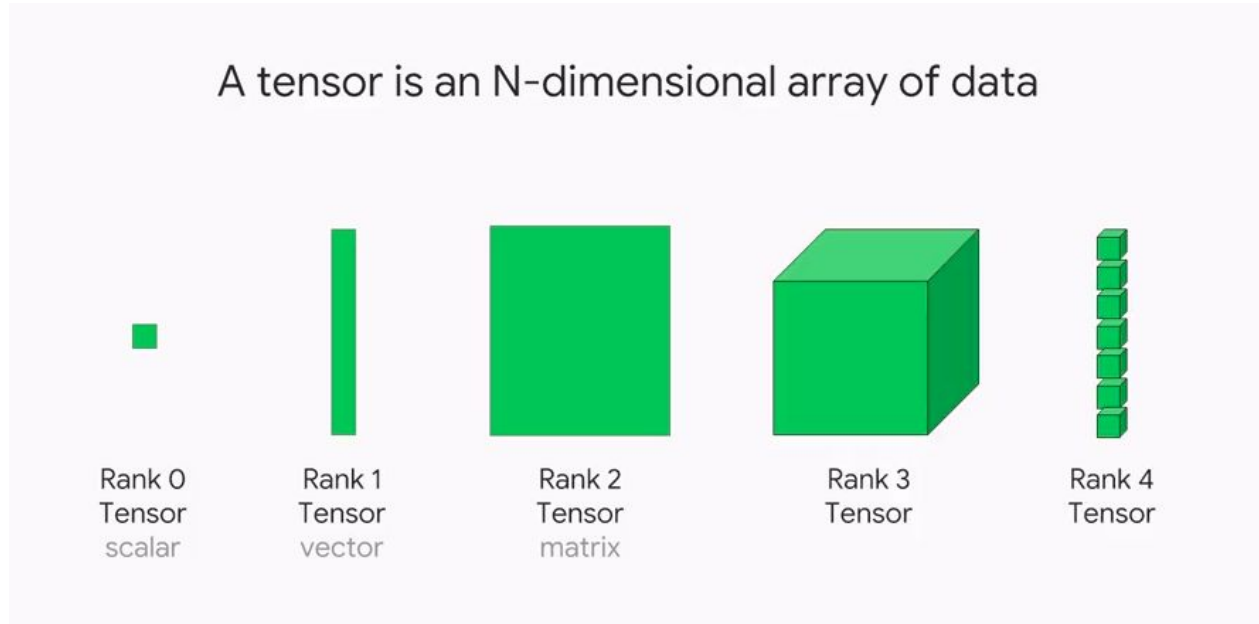
# Class work

Calculate the output of this unit:



# Unit calculations as tensor arithmetic

What is a tensor?



# Unit calculations as tensor arithmetic

Summation for a single unit:

$$o_j = g \left( \sum_k i_k w_{k,j} \right)$$

Vector arithmetic for a single unit:

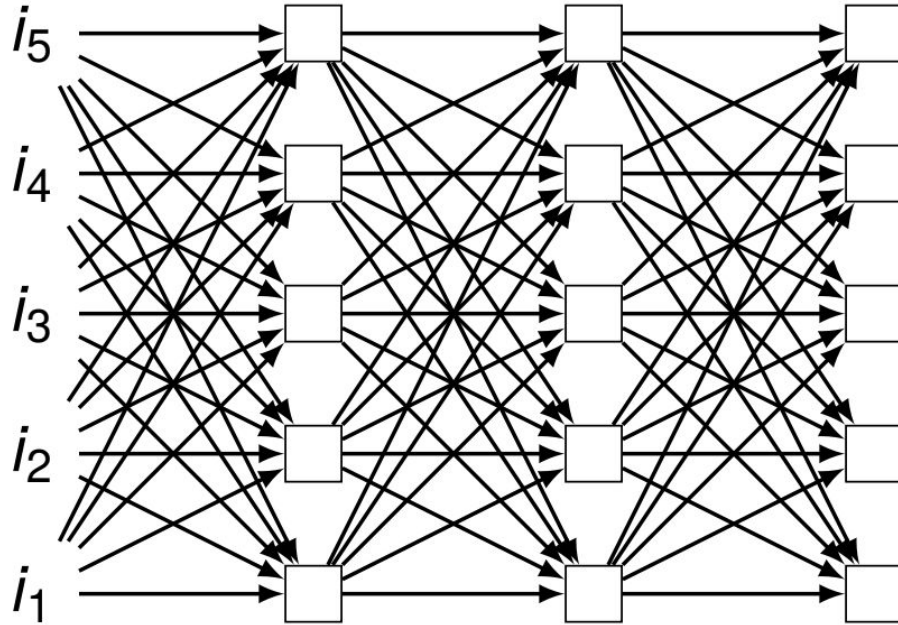
$$o_j = g \left( \begin{bmatrix} i_1 & i_2 & \dots & i_n \end{bmatrix} \begin{bmatrix} w_{1,j} \\ w_{2,j} \\ \dots \\ w_{n,j} \end{bmatrix} \right) = g (\mathbf{i} \times \mathbf{w}_{*,j})$$

Matrix arithmetic for multiple units:

$$\mathbf{o} = g (\mathbf{i} \times \mathbf{W})$$



# A feedforward network as composition



$$\mathbf{o} = g \left( g \left( g \left( \mathbf{i} \times \mathbf{W}^1 \right) \times \mathbf{W}^2 \right) \times \mathbf{W}^3 \right)$$

# Activation functions

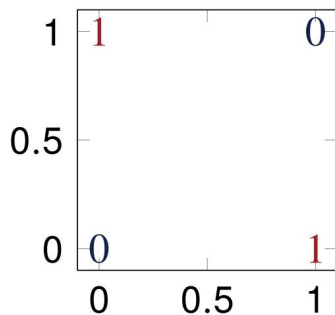
- Why do we need activation functions?
- What types of activation functions do we have?

# Learning the XOR

$$\mathbf{X} = \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix}$$

$$\mathbf{y} = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}$$

No such weights exist!



Can you solve it with linear regression?

$$\mathbf{y} = \mathbf{XW} + \mathbf{b}$$

$$\begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} \mathbf{W} + \mathbf{b}$$

$$\begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} ? \\ ? \end{bmatrix} + ?$$

# Solving XOR with NNs

$$f(\mathbf{X}; \mathbf{W}, \mathbf{c}, \mathbf{w}, b) = \max(0, \mathbf{XW} + \mathbf{c}) \mathbf{w} + b$$

$$\begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix} = \max \left( 0, \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} \mathbf{W} + \mathbf{c} \right) \mathbf{w} + b$$

$$\begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix} = \max \left( 0, \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} ? & ? \\ ? & ? \end{bmatrix} + \begin{bmatrix} ? & ? \end{bmatrix} \right) \begin{bmatrix} ? \\ ? \end{bmatrix} + ?$$

## Solving XOR with NNs

$$\begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix} = \max \left( 0, \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} + \begin{bmatrix} 0 & -1 \end{bmatrix} \right) \begin{bmatrix} 1 \\ -2 \end{bmatrix} + 0$$

$$\begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ -2 \end{bmatrix} + 0$$

# Common activation functions in hidden units

We have: affine transformation of input  $\mathbf{x}$ , followed by nonlinear activation function

$$\mathbf{h} = g(\mathbf{W}^\top \mathbf{x} + \mathbf{b})$$

$g$  could be just about anything! It can be linear, but a linear function is not preferred. Why?

Considerations:

- What specific behavior is needed?
- How will the gradients behave?

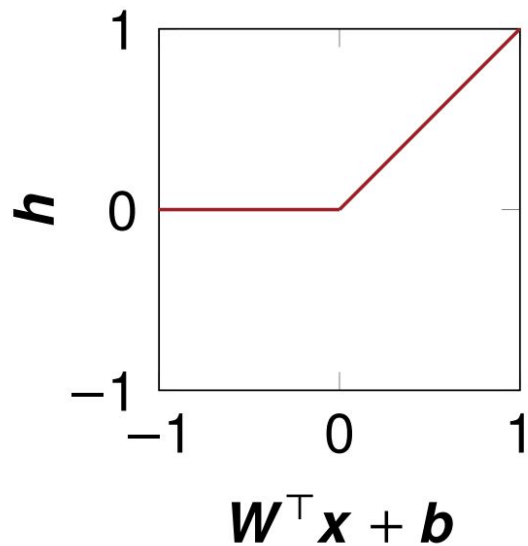
# Why linear functions are not preferred?

Because of the considerations.

- We need complex mappings between the inputs and the outputs
- All linear layers will translate the input linearly to output– that is, multiple linear transformation is basically one giant linear transformation
- Cannot use backpropagation as the derivative is constant

# ReLU

$$\mathbf{h} = \max(0, \mathbf{W}^\top \mathbf{x} + \mathbf{b})$$



Behavior?

- Active only when input is positive

Gradients?

- 1 when positive
- 0 when negative



# ReLU is non-differentiable

ReLU at  $z = 0$ :

- left derivative = 0
- right derivative = 1

So ReLU is not differentiable at  $z = 0$ !

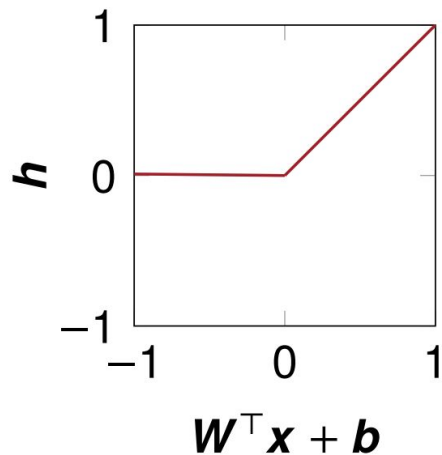
A few non-differentiable points are not a problem:

- Training rarely reaches a point with gradient 0 anyway
- Software simply returns either left or right derivative

# Leaky ReLU

$$f(x) = \max(0.1x, x)$$

$$\mathbf{h} = \max(0, \mathbf{W}^\top \mathbf{x} + \mathbf{b}) + 0.01 \min(0, \mathbf{W}^\top \mathbf{x} + \mathbf{b})$$



ReLU has a “dead neuron” problem!

Behavior?

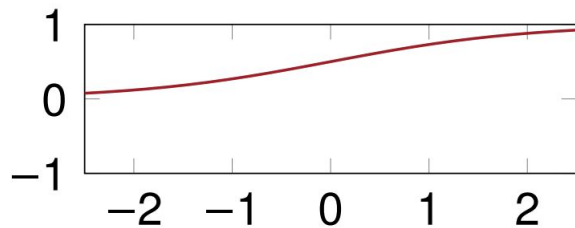
- Strong positive activation when positive
- Very weak negative activation when negative

Gradients?

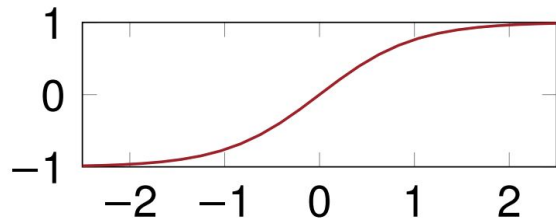
- 1 when positive
- 0.1 when negative

# Sigmoid and tanh

Sigmoid:  $f(x) = \frac{1}{1 + e^{-x}}$



Tanh:  $f(x) = \frac{(e^x - e^{-x})}{(e^x + e^{-x})}$



Behavior?

- Sigmoid: 0/1 switch
- Tanh: -1/+1 switch

Gradients?

- Saturate across most of their domain
- Tanh optimizes slightly better since it is similar to the identity function near 0

# Hyperparameters

- Network parameters are the ones that are being learned throughout the training process (e.g. weights)
- There are parameters that we can control to facilitate this learning: they are called hyperparameters
- Activation function is one such hyperparameter
- We have others ...

# Optimizers

- Algorithms used to update the learnable parameters in order to reduce loss
- Common optimizers:
  - Gradient descent
  - Stochastic gradient descent
  - Mini-batch gradient descent
  - Momentum
  - Adagrad
  - RMSProp
  - AdaDelta
  - Adaptive moment estimation
- GD methods maintain a single learning rate (with/without decay for all parameters and are known as first order optimizers (works with the first order derivative))
- Adaptive methods like Adagrad provide learning rates for each parameter, thus improving the learnability, but are computationally expensive
- RMSProp not only provides LR for each parameter, it adapts based on the mean of recent magnitudes of the gradients for the weight → first order. AdaDelta is similar but works with squared gradients (second order)

# Adaptive moment estimation

- Popularly known as Adam (not ADAM)
- Came out of OpenAI and University of Toronto
- Most likely the highest cited [paper](#) in recent history
- Why is it so popular?
  - Straightforward to implement
  - Little memory requirements
  - Well suited for problems that are large in terms of data and/or parameters
  - Needs little to no manual tuning

# Adam

- Adam works with momentums of first and **second** order
- Instead of adapting the parameter learning rates based on the average first moment (the mean),  $m_t$  as in RMSProp, Adam also makes use of the average of the second moments of the gradients (the uncentered variance)  $v_t$

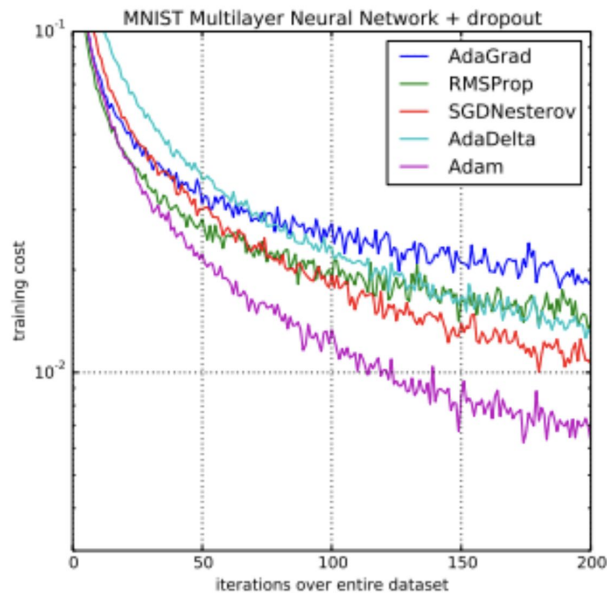
$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}.$$

$$\hat{v}_t = \frac{v_t}{1 - \beta_2^t}.$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \hat{m}_t.$$

- The values for  $\beta_1$  is 0.9 , 0.999 for  $\beta_2$  and epsilon is an extremely small number to avoid zero division

# Adam's popularity



Sebastian Ruder: “... RMSprop, Adadelata, and Adam are very similar algorithms that do well in similar circumstances. .... Adam might be the best overall choice.”

Andrej Karpathy: “In practice Adam is currently recommended as the default algorithm to use, and often works slightly better than RMSProp.”



# Regularization

- A set of strategies used in Machine Learning to reduce the generalization error
- Why?
- Bias-variance tradeoff
- Bias: error from wrong assumptions in the learning algorithm
  - High bias can cause an algorithm to miss the relevant relations between features and target outputs → Observations don't matter
  - Underfitting
- Variance: error from sensitivity to small fluctuations in the training set
  - High variance may result in modeling the random noise in the training data → focusing too much on observations
  - Overfitting

# Regularization

- A good model needs to balance bias and variance
- Regularization helps us do that

# Regularization techniques

- Introduce regularization term to the loss function
  - Introduces a small amount of bias to counter variance → reduce overfitting
  - Loss function: negative log likelihood or binary cross entropy
  - Most common terms:
    - L2 regularizer: Ridge regression → adds the “squared magnitude” of the coefficient as the penalty term to the loss function
    - L1 regularizer: Lasso regression → adds the “absolute value of magnitude” of the coefficient as a penalty term to the loss function
    - Uses  $\lambda$  that controls the sensitivity of the model to the input → less sensitive, less likely to overfit
    - L2 focuses on larger weights, so higher  $\lambda$  means L2 penalizes higher value weights more than lower ones → no feature essentially goes away
    - L1 has equal focus, so, higher  $\lambda$  → low weight features go away → feature selection

# Regularization techniques

- Dropout: Drops out (ignore) a layer's output with a probability  $p$ 
  - Choice of  $p$  depends on the architecture
- Early stopping: Stops when performance gets saturated
  - Stops model from being overfit
  - Returns the best current model
- Data augmentation
  - Introduce new training data with variation
  - Injects noise

# Other common hyperparameters

- Learning rate: how quickly a network updates its parameters
  - Low learning rate slows down the learning process but converges smoothly
  - Larger learning rate speeds up the learning but may not converge
  - Decaying learning rate is preferred
- Epochs: How many times the entire training dataset has passed through the model
- Batch size: (Mini) batch size refers to a subset of the training data. Weights are updated after each mini batch training
  - Small mini batch → too many updates
  - Large mini batch → too few updates, too many epochs to converge

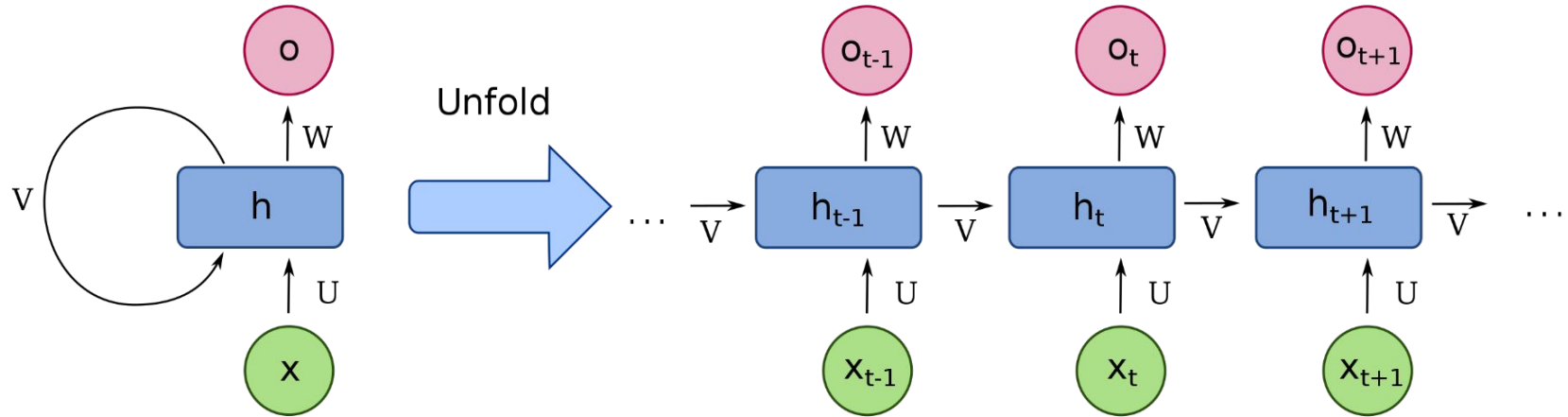
# Recurrent Neural Networks (RNN)

- Simple recurrent networks
- Bidirectional and gated recurrent networks
- Recurrent architectures
- Seq2Seq models
- Attention

# Short history of RNN

- 1986: RNNs are Introduced by David Rumelhart
- 1995: LSTMs are introduced by Sepp Hochreiter and Jürgen Schmidhuber based on Hochreiter's 1991 research on vanishing gradient problem
- 2001: Gers and Schmidhuber trained LSTMs to learn language models (unlearnable by HMMs)
- 2009: Graves et al. won ICDAR handwriting recognition competition using LSTMs
- 2013: Hinton and his team destroyed previous record for speech recognition using LSTM
- 2014: GRU is introduced by Cho et al.
- 2015: Widespread use in both academia and industry due to Google's adaptation of LSTM in their Google Voice speech recognition system

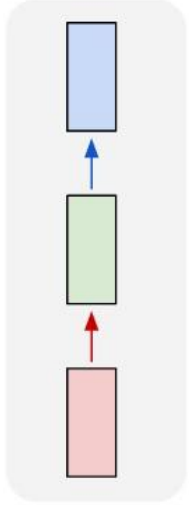
# Structure of an RNN



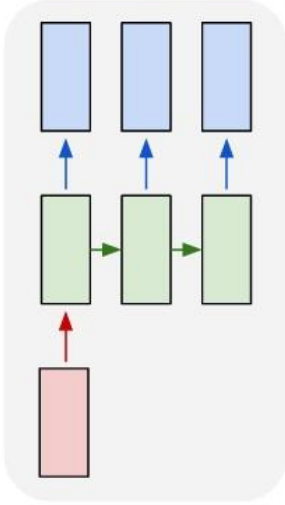


# Types of RNNs

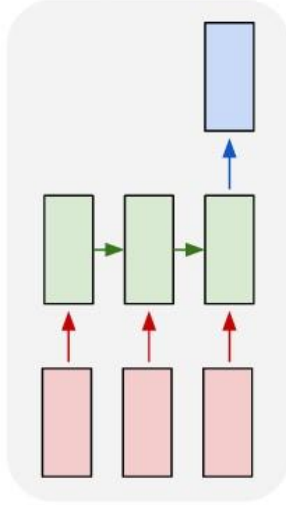
one to one



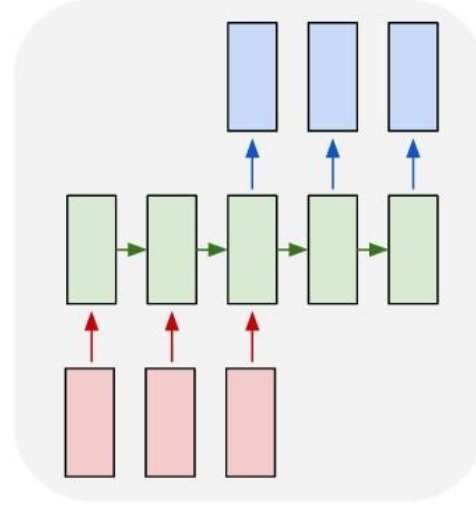
one to many



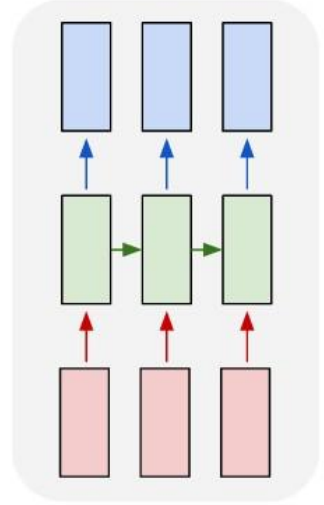
many to one



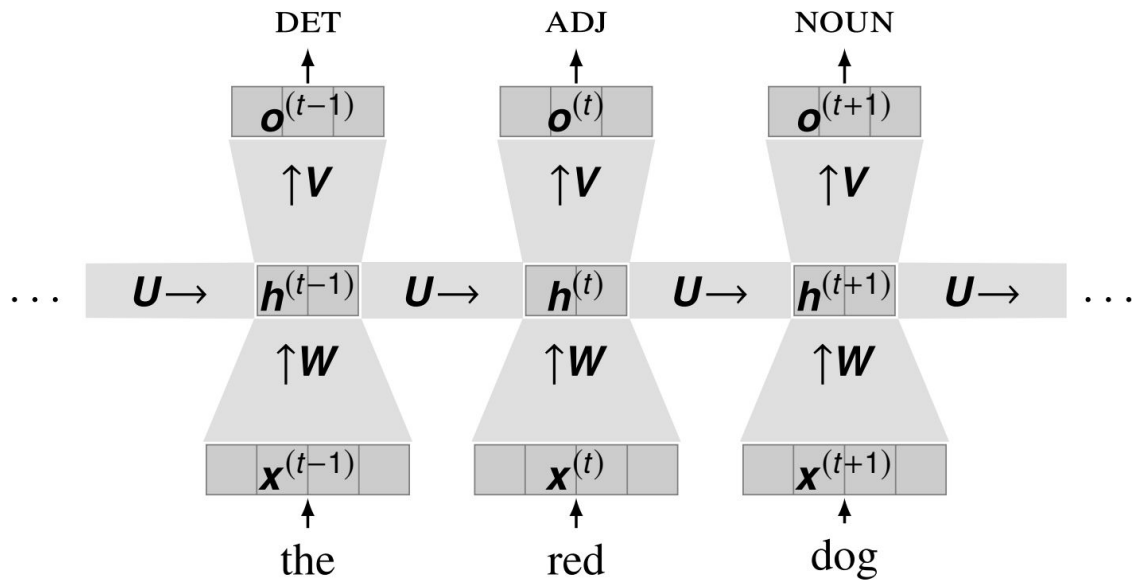
many to many



many to many



# A completely unrolled Simple RNN



$$\mathbf{o}^{(t)} = \text{softmax}(\mathbf{V}\mathbf{h}^{(t)})$$

$$\mathbf{h}^{(t)} = \tanh(\mathbf{U}\mathbf{h}^{(t-1)} + \mathbf{W}\mathbf{x}^{(t)})$$

# Simple RNN

Intuitions:

- Each step combines the current input with the history
- Each prediction is made based on this combination

Observations:

- The input and hidden state change at each time step
- The parameters  $W$ ,  $U$ ,  $V$  are the same at each step

Equations:

$$\mathbf{h}^{(t)} = \tanh(\mathbf{U}\mathbf{h}^{(t-1)} + \mathbf{W}\mathbf{x}^{(t)})$$

$$\mathbf{o}^{(t)} = \text{softmax}(\mathbf{V}\mathbf{h}^{(t)})$$

# Classwork

Consider an RNN that predicts as:

$$\mathbf{o}^{(t)} = \text{softmax}(\mathbf{V}\mathbf{h}^{(t)})$$

$$\mathbf{h}^{(t)} = \mathbf{U}\mathbf{h}^{(t-1)} + \mathbf{W}\mathbf{x}^{(t)}$$

whose parameters have been set to:

$$\mathbf{U} = \begin{bmatrix} 2 & 0 \\ 1 & 2 \end{bmatrix} \quad \mathbf{W} = \begin{bmatrix} 2 & 0 & -1 \\ 1 & -1 & 1 \end{bmatrix} \quad \mathbf{V} = \begin{bmatrix} 1 & 0 \\ -1 & 1 \\ 0 & 3 \end{bmatrix}$$

If you are given the following input:

$$\mathbf{h}^{(0)} = \begin{bmatrix} 1 \\ -1 \end{bmatrix} \quad \mathbf{x}^{(1)} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \quad \mathbf{x}^{(2)} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

Which label will be predicted for each word if in the final softmax, index 0=ADJ, index 1=DET, and index 2=NOUN?

# Drawbacks of simple RNN

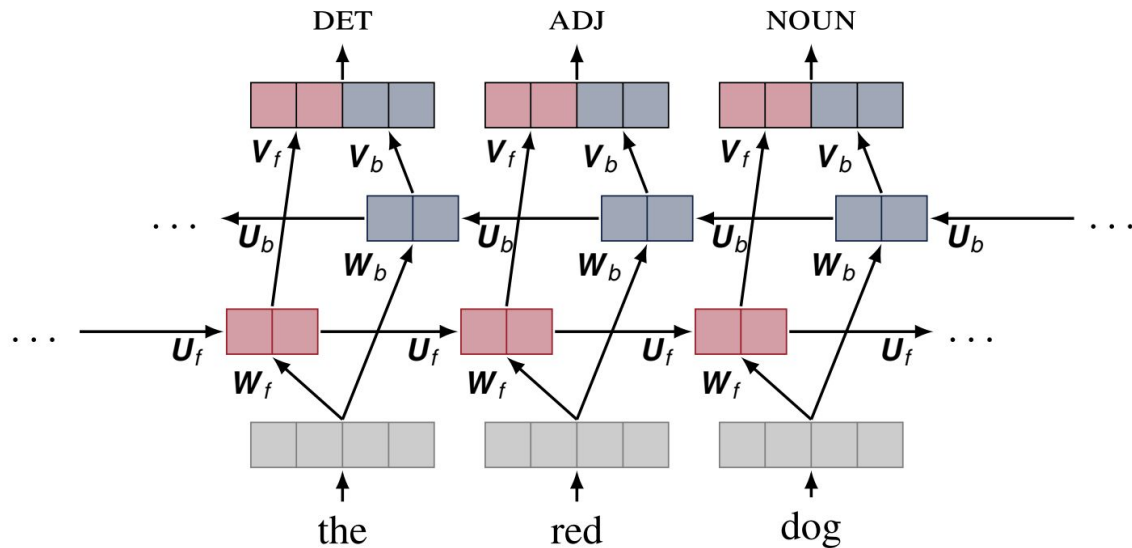
- They can only see the past, not the future
- They must forget the same amount of history at each time step
  - Theoretically, they don't have to
  - Maintaining long term memory is difficult

# Bidirectional RNNs

Intuition:

- Run one forward RNN
- Run one backward RNN
- Combine their outputs

# Bidirectional RNNs



$$\mathbf{o}^{(t)} = \text{softmax}(\mathbf{V}_f \mathbf{h}_f^{(t)} + \mathbf{V}_b \mathbf{h}_b^{(t)})$$

$$\mathbf{h}_b^{(t)} = \tanh(\mathbf{U}_b \mathbf{h}^{(t+1)} + \mathbf{W}_b \mathbf{x}^{(t)})$$

$$\mathbf{h}_f^{(t)} = \tanh(\mathbf{U}_f \mathbf{h}^{(t-1)} + \mathbf{W}_f \mathbf{x}^{(t)})$$

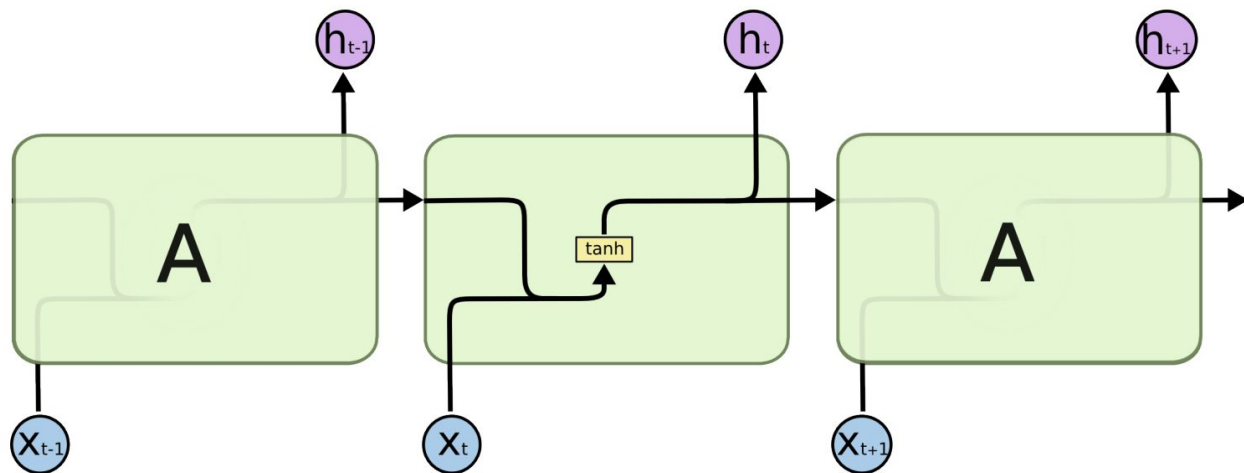
# Gated RNNs

Intuition:

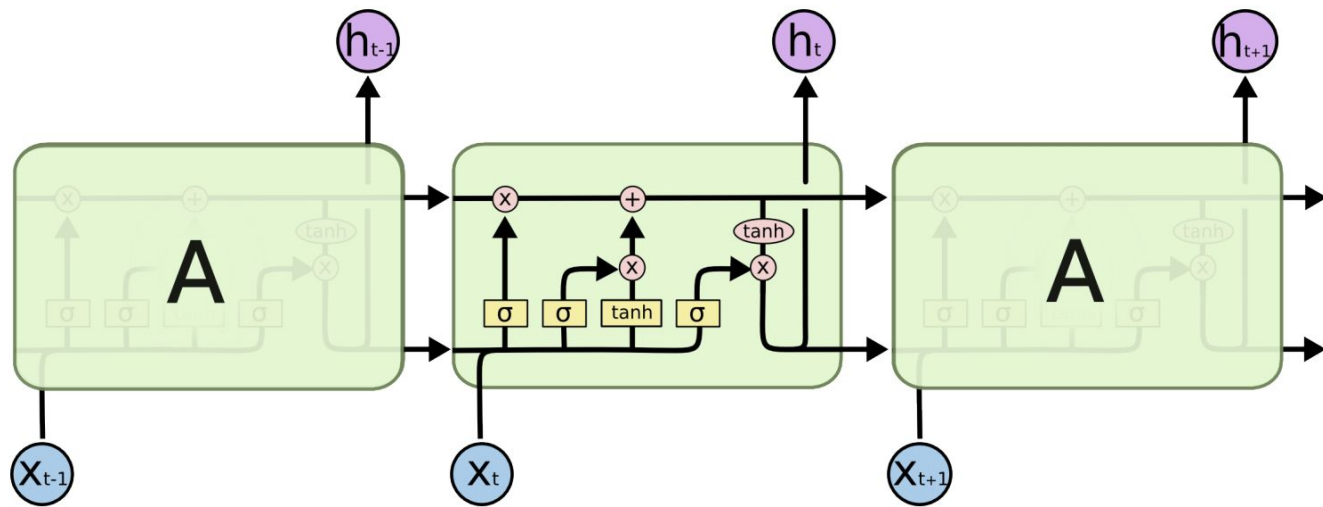
- Simple RNNs forget the same amount at each time step
- Look at the previous hidden state and the current input
- Decide how much to forget based on those
- Two gated RNNs
  - Long Short-Term Memory (LSTM)
  - Gated Recurrent Units (GRU)



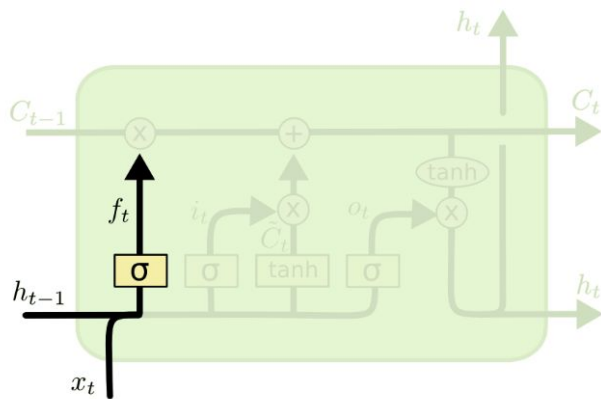
# Simple RNN



# LSTM

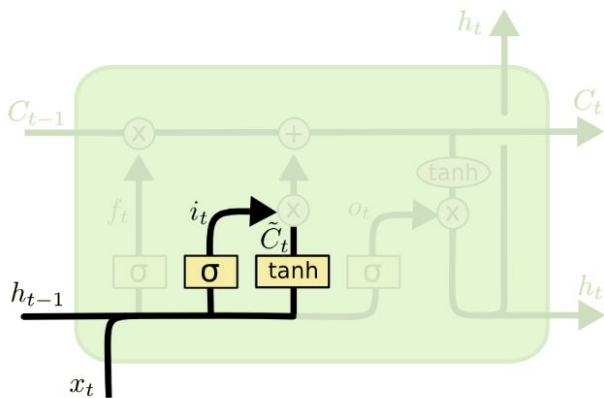


# Forget gate



$$f_t = \sigma (W_f \cdot [h_{t-1}, x_t] + b_f)$$

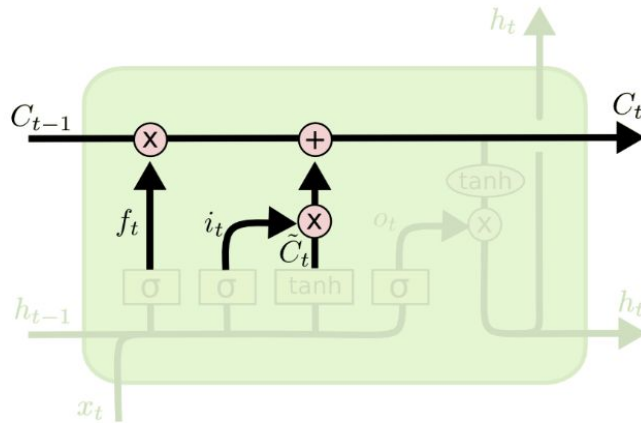
# Input gate



$$i_t = \sigma (W_i \cdot [h_{t-1}, x_t] + b_i)$$

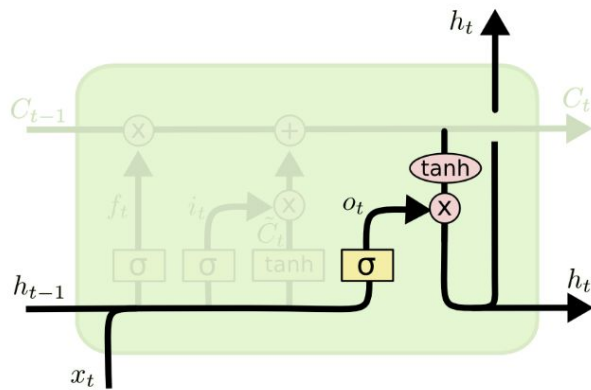
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

# Cell update



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

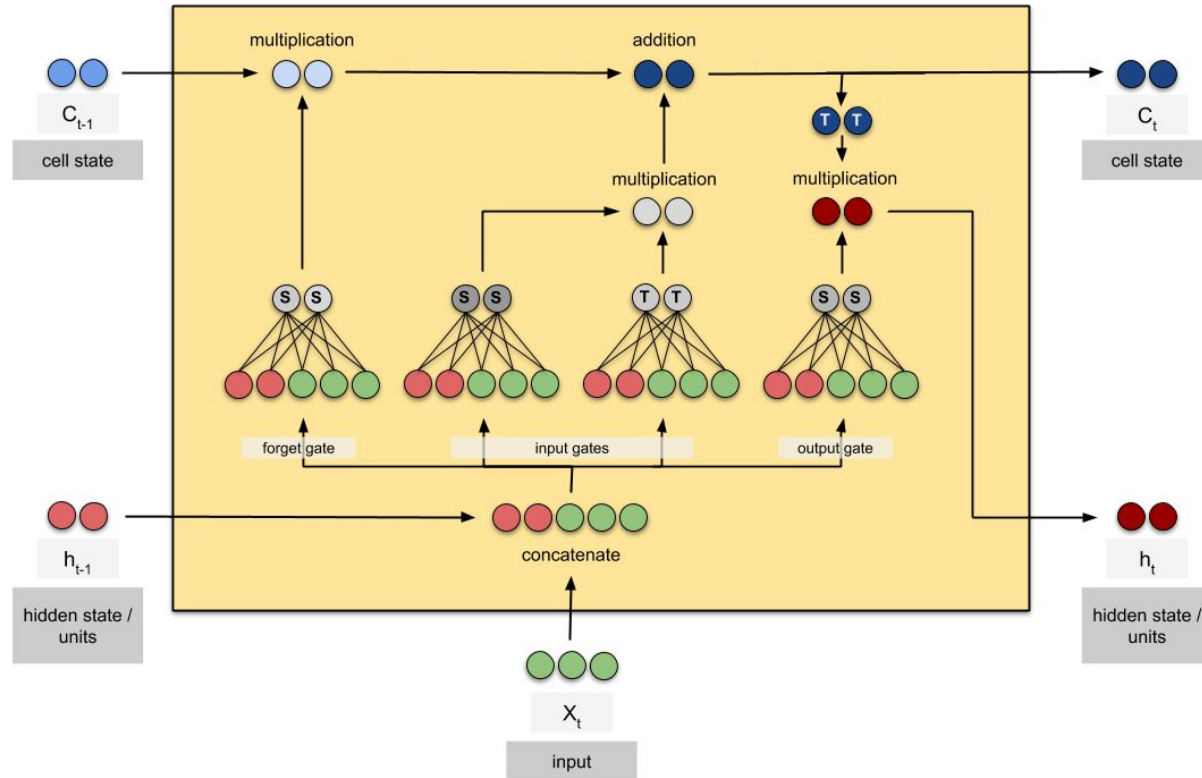
# Output gate



$$o_t = \sigma (W_o [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh (C_t)$$

# One LSTM cell



# LSTM

Step 1: calculate forget, input, and output gates

$$\mathbf{f}^{(t)} = \sigma(\mathbf{b}_f + \mathbf{W}_f \mathbf{h}^{(t-1)} + \mathbf{U}_f \mathbf{x}^{(t)})$$

$$\mathbf{i}^{(t)} = \sigma(\mathbf{b}_i + \mathbf{W}_i \mathbf{h}^{(t-1)} + \mathbf{U}_i \mathbf{x}^{(t)})$$

$$\mathbf{o}^{(t)} = \sigma(\mathbf{b}_o + \mathbf{W}_o \mathbf{h}^{(t-1)} + \mathbf{U}_o \mathbf{x}^{(t)})$$

Step 2: update “cell”: forget some old, add some new

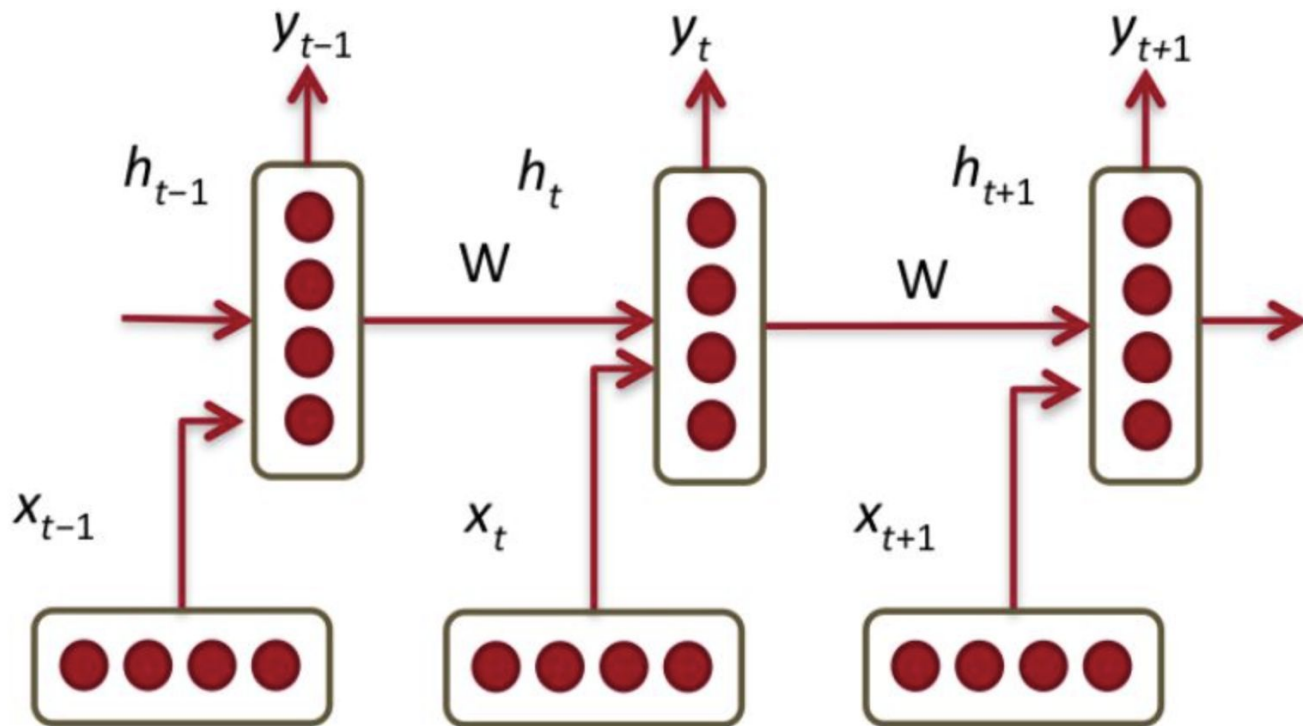
$$\mathbf{c}^{(t)} = \mathbf{f}^{(t)} \odot \mathbf{c}^{(t-1)} + \mathbf{i}^{(t)} \odot \tanh(\mathbf{b}_c + \mathbf{W}_c \mathbf{h}^{(t-1)} + \mathbf{U}_c \mathbf{x}^{(t)})$$

Step 3: update hidden state: output some of the cell

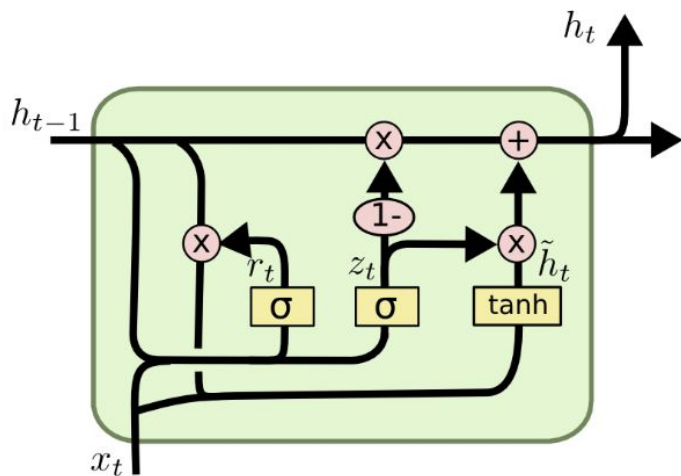
$$\mathbf{h}^{(t)} = \mathbf{o}^{(t)} \odot \tanh(\mathbf{c}^{(t)})$$



# LSTM Network



# GRU



$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

# GRU

Step 1: calculate update and reset gates

$$\mathbf{u}^{(t)} = \sigma(\mathbf{b}_u + \mathbf{W}_u \mathbf{h}^{(t-1)} + \mathbf{U}_u \mathbf{x}^{(t)})$$

$$\mathbf{r}^{(t)} = \sigma(\mathbf{b}_r + \mathbf{W}_r \mathbf{h}^{(t-1)} + \mathbf{U}_r \mathbf{x}^{(t)})$$

Step 2: update hidden: forget some old, add some new

$$\begin{aligned} \mathbf{h}^{(t)} = & (1 - \mathbf{u}^{(t)}) \odot \mathbf{h}^{(t-1)} \\ & + \mathbf{u}^{(t)} \odot \tanh(\mathbf{b}_h + \mathbf{W}_h (\mathbf{r}^{(t)} \odot \mathbf{h}^{(t-1)}) + \mathbf{U}_h \mathbf{x}^{(t)}) \end{aligned}$$

GRUs generally perform as well or better than LSTMs,  
with fewer parameters

# Gated RNNs

Properties:

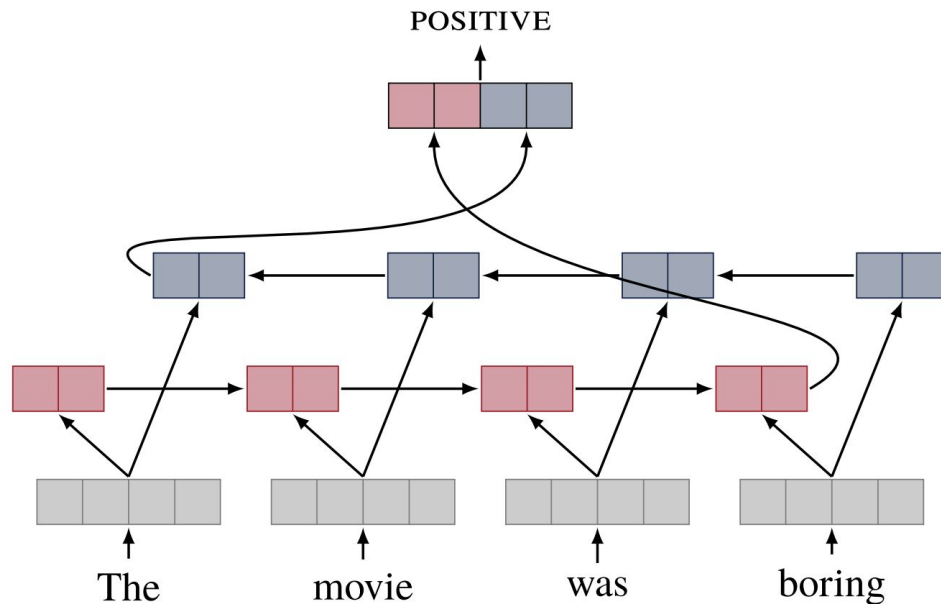
- Can forget different amounts at each time step
- Much better at using long distance information

A bidirectional GRU is a good starting point for many sequence tagging tasks

# Recurrent architectures for related tasks

## RNNs for text classification

- Last hidden state of the RNN represents the entire sentence.



# Recurrent architectures for related tasks

What other tasks can RNNs handle?

# Next to come

- Seq2seq models
  - Encoders and decoders
- Attention
- Autoregressive models
- Transformers
- A whole lotta variations
  - Transformer-encoders vs. transformer-decoders