# Lecture 6: Recurrent & Graph Neural Networks
## Efstratios Gavves

# Lecture overview

o Sequential data

o Recurrent Neural Networks

o Backpropagation through time

o Exploding and vanishing gradients

o LSTMs and variants

o Encoder-Decoder Architectures

o Graph Neural Networks

Sequence data

Sequence applications

# Example of sequential data

- Videos
- Other?

# Example of sequential data

o Videos

o Other?

o Time series data
   ◦ Stock exchange
   ◦ Biological measurements
   ◦ Climate measurements
   ◦ Market analysis

o Speech/Music

o User behavior in websites

o .....

# Applications

o Machine translation

o Image captioning

o Question answering

o Video generation

o Speech synthesis

o Speech recognition

# A sequence of probabilities

o Sequence → Chain rule of probabilities

$$p(x) = \prod_i p(x_i | x_1, \ldots, x_{i-1})$$

o For instance, let's model that "This is the best course!"

$p(\texttt{This is the best course!}) =$
$= p(\texttt{This}) \cdot$
$p(\texttt{is}|\texttt{This}) \cdot$
$p(\texttt{the}|\texttt{This is}) \cdot \ldots \cdot$
$p(\texttt{!}|\texttt{This is the best course})$

# What is the problem with sequences?

- ???

# What is the problem with sequences?

o Sequences might be of arbitrary or even infinite lengths

o Infinite parameters?

# What is the problem with sequences?

o Sequences might be of arbitrary or even infinite lengths

o Infinite parameters?

o No, better share and reuse parameters

o RecurrentModel(`I think, therefore, I am. |` $\boldsymbol{\theta}$`)`

    can be reused also for

    RecurrentModel(`Everything is repeated in circles. History is a Master because it teaches that it doesn't exist. It is the permutations that matter|` $\boldsymbol{\theta}$`)`

o For a ConvNet that is not straightforward

o Why?

# What is the problem with sequences?

o Sequences might be of arbitrary or even infinite lengths

o Infinite parameters?

o No, better share and reuse parameters

o RecurrentModel(`I think, therefore, I am. |` $\boldsymbol{\theta}$`)`

can be reused also for

RecurrentModel(`Everything is repeated in circles. History is a Master because it teaches that it doesn't exist. It is the permutations that matter|` $\boldsymbol{\theta}$`)`

o For a ConvNet that is not straightforward

o Why? Fixed dimensionalities

# Some properties of sequences?

# Some properties of sequences

o Data inside a sequence are non identically, independently distributed (IID)
  ◦ The next "word" depends on the previous "words"
  ◦ Ideally on all of them

o We need context, and we need memory!

o **Big question:** How to model context and memory ?

I    am    Bond    ,    James    Bond    ⎰ McGuire

Bond

tired

am

!

# Properties of sequences

○ Data inside a sequence are non identically, independently distributed (IID)
  ◦ The next "word" depends on the previous "words"
  ◦ Ideally on all of them

○ We need context, and we need memory!

○ **Big question:** How to model context and memory ?

McGuire

| Bond |

I    am    Bond    ,    James    Bond    tired

am

!

# One-hot vectors

o A vector with all zeros except for the active dimension

o 12 words in a sequence → 12 One-hot vectors

o After the one-hot vectors apply an embedding
  ◦ Word2Vec, GloVE

| Vocabulary | | One-hot vectors | | | | | | |
|---|---|---|---|---|---|---|---|---|
| I | I | 1 | I | 0 | I | 0 | I | 0 |
| am | am | 0 | am | 1 | am | 0 | am | 0 |
| Bond | Bond | 0 | Bond | 0 | Bond | 1 | Bond | 0 |
| James | James | 0 | James | 0 | James | 0 | James | 1 |
| tired | tired | 0 | tired | 0 | tired | 0 | tired | 0 |
| , | , | 0 | , | 0 | , | 0 | , | 0 |
| McGuire | McGuire | 0 | McGuire | 0 | McGuire | 0 | McGuire | 0 |
| ! | ! | 0 | ! | 0 | ! | 0 | ! | 0 |

# Why not indices instead of one-hot vectors?

One-hot representation            OR?            Index representation

$$x_{t=1,2,3,4} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

Column headers: I, am, James, McGuire

I am James McGuire

$$x_{"I"} = 1$$
$$x_{"am"} = 2$$
$$x_{"James"} = 4$$
$$x_{"McGuire"} = 7$$

# Why not indices instead of one-hot vectors?

One-hot representation          OR?          Index representation

$$x_{t=1,2,3,4} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

Column headers: I, am, James, McGuire

I am James McGuire

$$x_{"I"} = 1$$
$$x_{"am"} = 2$$
$$x_{"James"} = 4$$
$$x_{"McGuire"} = 7$$

$$\ell_2(x_{am}, x_{McQuire}) = \sqrt{2}$$

$$=$$

$$\ell_2(x_I, x_{am}) = \sqrt{2}$$

$$\ell_2(x_{am}, x_{McQuire}) = (7-2)^2 = 5$$

$$\neq$$

$$\ell_2(x_I, x_{am}) = (2-1)^2 = 1$$

# Recurrent Neural Networks

# Backprop through time

Output

NN Cell State

Recurrent connections

Input

# Memory

o Memory is a mechanism that learns a representation of the past

o At timestep $t$ project all previous information $1, \dots, t$ onto a latent space $c_t$
  ◦ Memory controlled by a neural network $h_\theta$ with shared parameters $\theta$

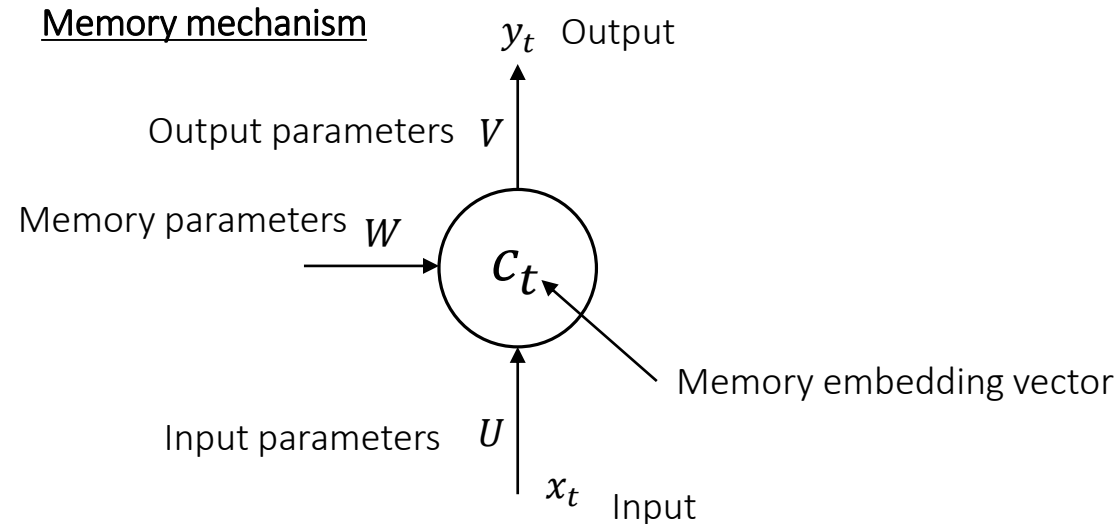o Then, at timestep $t + 1$ re-use the parameters $\theta$ and the previous $c_t$

$$c_{t+1} = h_\theta(x_{t+1}, c_t)$$
$$\dots$$
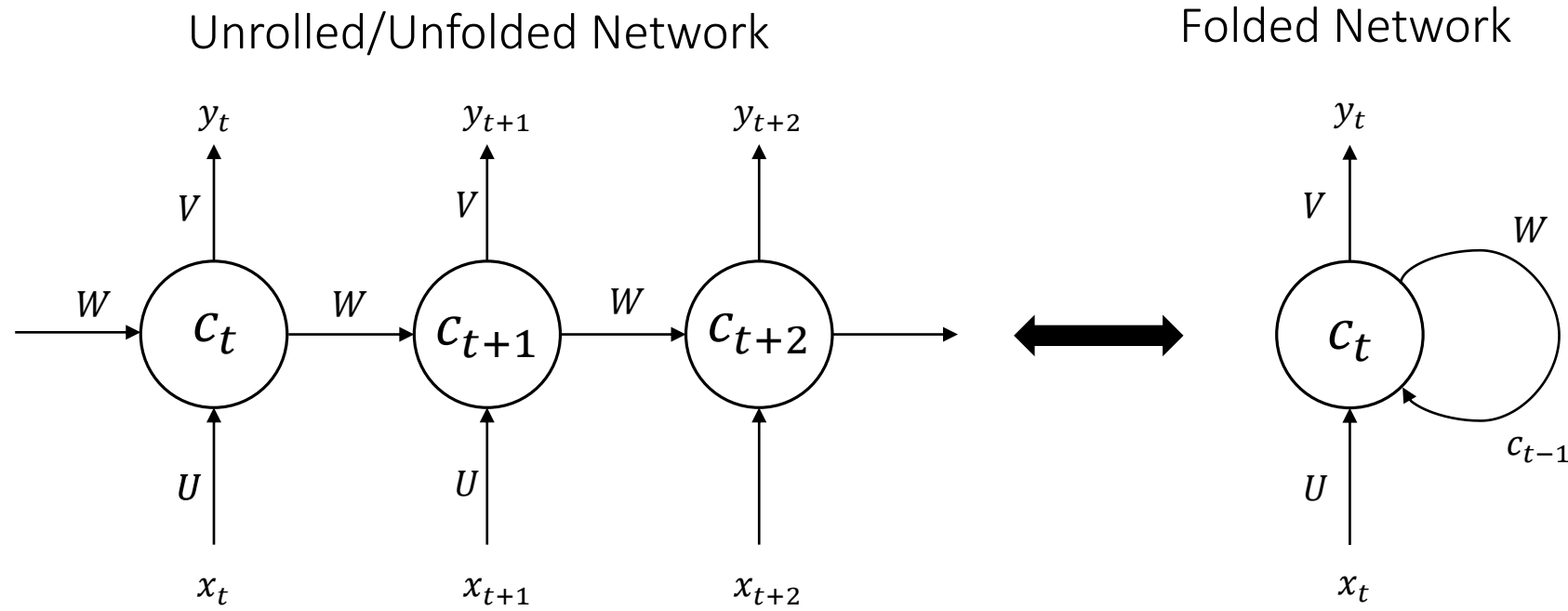$$c_{t+1} = h_\theta(x_{t+1}, h_\theta(x_t, h_\theta(x_{t-1}, \dots h_\theta(x_1, c_0))))$$

# A graphical representation of memory

o In the simplest case, what are the Inputs/Outputs of our system

o Sequence inputs ➔ we model them with parameters $U$

o Sequence outputs ➔ we model them with parameters $V$

o Memory I/O ➔ we model it with parameters $W$

Memory mechanism

$y_t$  Output

Output parameters  $V$

Memory parameters  $W$

$c_t$

Memory embedding vector

Input parameters  $U$

$x_t$  Input

# A graphical representation of memory

o In the simplest case, what are the Inputs/Outputs of our system

o Sequence inputs → we model them with parameters $U$

o Sequence outputs → we model them with parameters $V$

o Memory I/O → we model it with parameters $W$

# Folding the memory

Unrolled/Unfolded Network

Folded Network

# Recurrent Neural Networks - RNNs

o Basically, two equations

$$c_t = \tanh(U\,x_t + W\,c_{t-1})$$
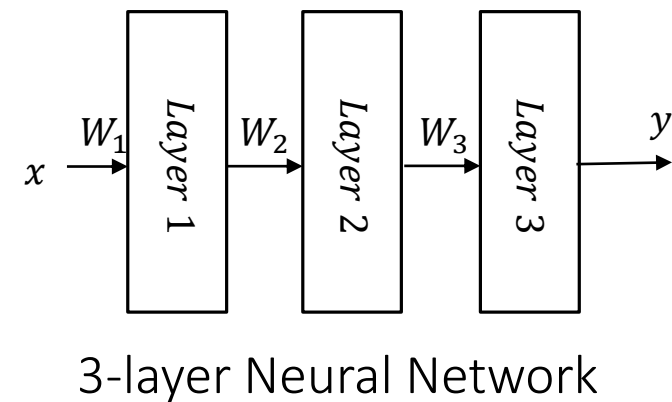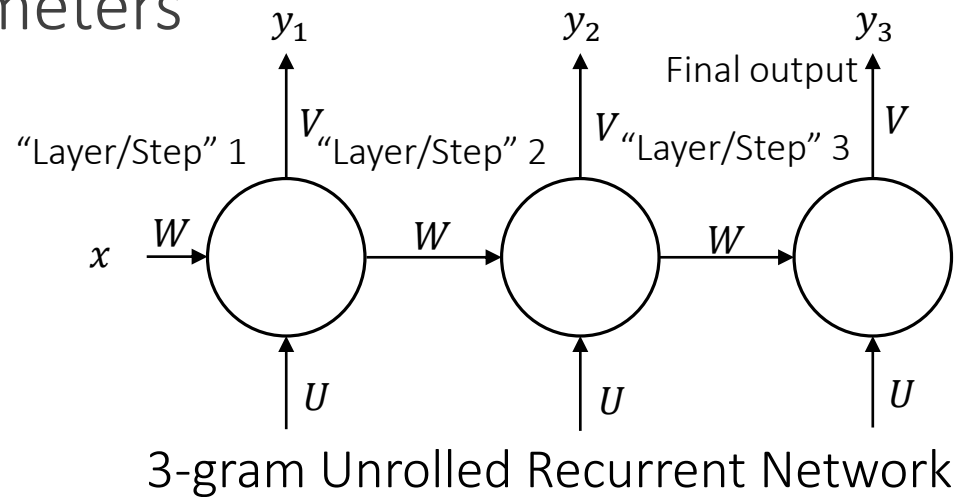$$y_t = \mathrm{softmax}(V\,c_t)$$

o And a loss function

$$\mathcal{L} = \sum_t \mathcal{L}_t(y_t, y_t^*)$$
$$= \sum_t y_t^* \log y_t$$

assuming the cross-entropy loss function

# RNNs vs MLPs

o Is there a big difference?

o Instead of layers ➔ Steps

o Outputs at every step ➔ MLP outputs in every layer possible

o Main difference: Instead of layer-specific parameters ➔ Layer-shared parameters



3-gram Unrolled Recurrent Network

3-layer Neural Network

# Hmm, layers share parameters ?!?

o How is the training done? Does Backprop remain the same?

# Hmm, layers share parameters ?!?

o <span style="color:red">How is the training done? Does Backprop remain the same?</span>

o Basically, chain rule

◦ So, again the same concept

o Yet, a bit more tricky this time, as the gradients survive over time

# Backpropagation through time

$$c_t = \tanh(U\, x_t + W\, c_{t-1})$$
$$y_t = \text{softmax}(V\, c_t)$$
$$\mathcal{L} = \sum_t y_t^* \log y_t$$

o Let's say we focus on the third timestep loss

$$\frac{\partial \mathcal{L}}{\partial V} = \cdots$$
$$\frac{\partial \mathcal{L}}{\partial W} = \cdots$$
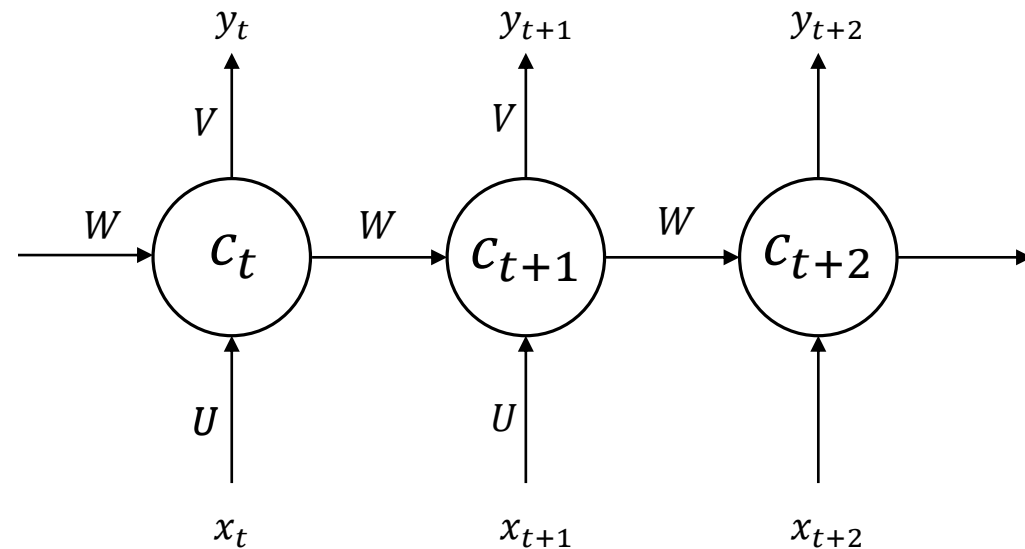$$\frac{\partial \mathcal{L}}{\partial U} = \cdots$$

# Backpropagation through time: $\partial \mathcal{L}_t / \partial V$

- Expanding the chain rule

$$\frac{\partial \mathcal{L}_t}{\partial V} = \frac{\partial \mathcal{L}_t}{\partial y_{t_k}} \frac{\partial y_{t_k}}{\partial c_{t_l}} \frac{\partial c_{t_l}}{\partial V_{ij}} = \cdots$$
$$= \cdots = (y_t - y_t^*) \otimes c_t$$

- All terms depend only on the current timestep $t$

- Then, we should sum up all the gradients for all time steps

$$\frac{\partial \mathcal{L}}{\partial V} = \sum_t \frac{\partial \mathcal{L}_t}{\partial V}$$

# Backpropagation through time: $\partial \mathcal{L}_t / \partial W$

$$c_t = \tanh(U\,x_t + W c_{t-1})$$
$$y_t = \mathrm{softmax}(V\,c_t)$$

o Expanding with the chain rule

$$\frac{\partial \mathcal{L}_t}{\partial W} = \frac{\partial \mathcal{L}_t}{\partial y_t} \frac{\partial y_t}{\partial c_t} \frac{\partial c_t}{\partial W}$$

o However, $c_t$ itself depends on $c_{t-1} \rightarrow \frac{\partial c_t}{\partial W}$ depends also on $c_{t-1} \rightarrow$
  The current dependency of $c_t$ to $W$ is recurrent
  ◦ And continuing till we reach $c_{-1} = [0]$
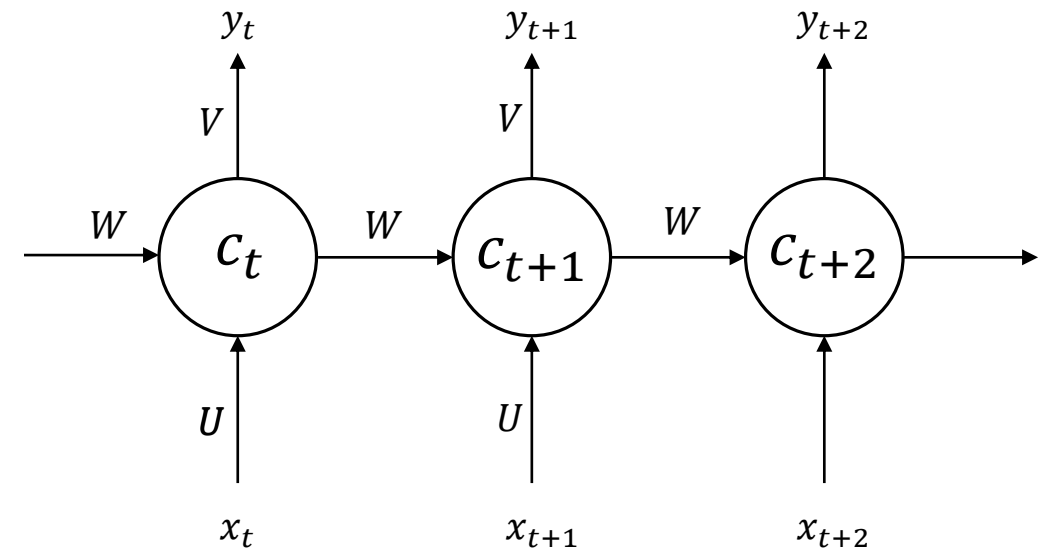
o So, in the end we have

$$\frac{\partial \mathcal{L}_t}{\partial W} = \sum_{k=0}^{t} \frac{\partial \mathcal{L}_t}{\partial y_t} \frac{\partial y_t}{\partial c_t} \frac{\partial c_t}{\partial c_k} \frac{\partial c_k}{\partial W}$$

o The gradient $\frac{\partial c_t}{\partial c_k}$ itself is subject to the chain rule

$$\frac{\partial c_t}{\partial c_k} = \frac{\partial c_t}{\partial c_{t-1}} \frac{\partial c_{t-1}}{\partial c_{t-2}} \cdots \frac{\partial c_{k+1}}{\partial c_k} = \prod_{j=k+1}^{t} \frac{\partial c_j}{\partial c_{j-1}}$$

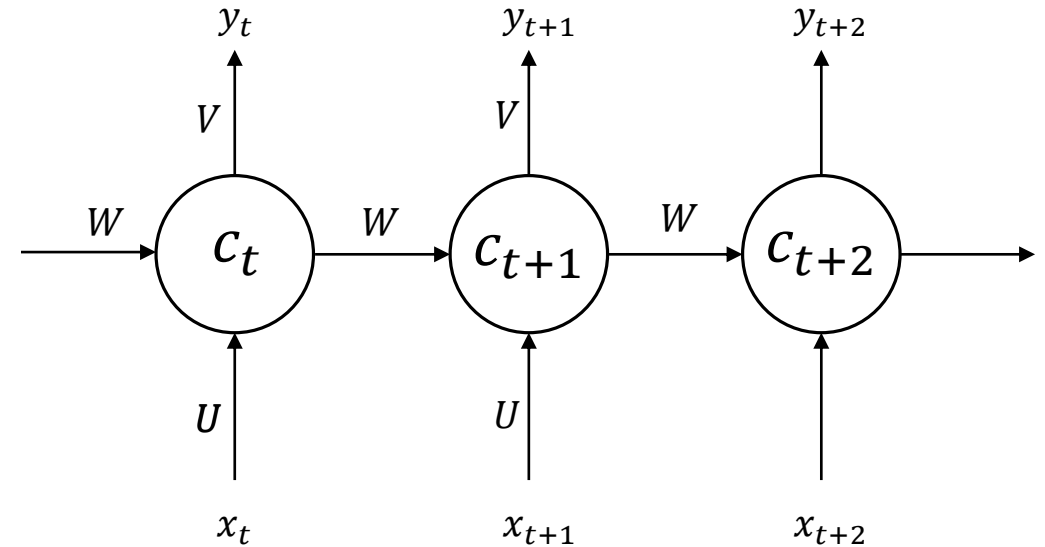o Then, we should sum up all the gradients for all time steps

# Backpropagation through time: $\partial \mathcal{L}_t / \partial U$

o For parameter matrix $U$ a similar process

$$\frac{\partial \mathcal{L}_t}{\partial U} = \sum_{k=0}^{t} \frac{\partial \mathcal{L}_t}{\partial y_t} \frac{\partial y_t}{\partial c_t} \frac{\partial c_t}{\partial c_k} \frac{\partial c_k}{\partial U}$$
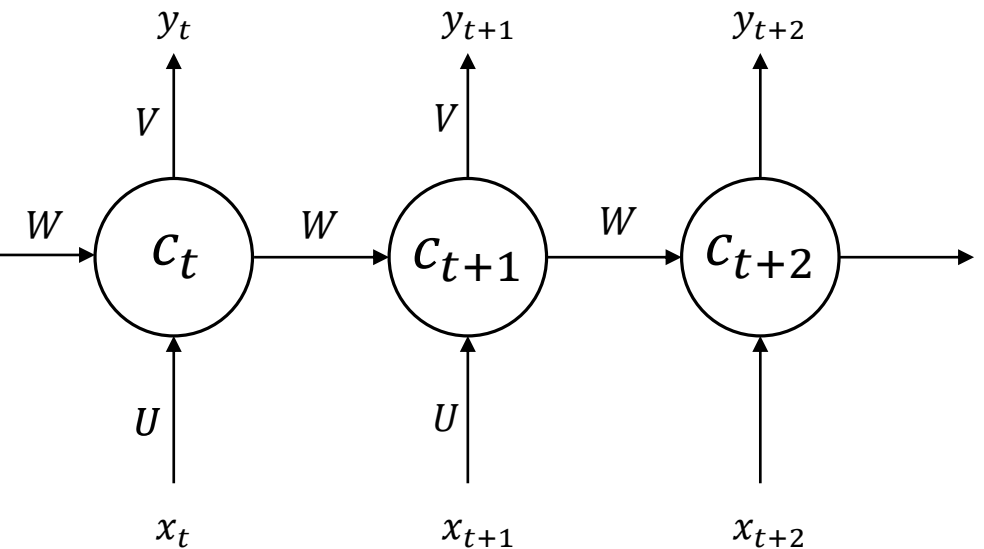
$$c_t = \tanh(U\, x_t + W c_{t-1})$$
$$y_t = \mathrm{softmax}(V\, c_t)$$

# Trading off Weight Update Frequency & Gradient Accuracy

o At time $t$ we use current weights $w_t$ to compute states $\mathbf{c}_t$ and outputs $y_t$

o Then, we use the states and outputs to backprop and get $w_{t+1}$

o Then, at $t+1$ we use $w_{t+1}$ and the current state $c_t$ to $y_{t+1}$ and $c_{t+1}$

o Then we update the weights again with $y_{t+1}$.
  ◦ The problem is $y_{t+1}$ was computed with $c_t$ in mind, which in turns depends on the old weights $w_t$, not the current ones $w_{t+1}$. So, the new gradients are only an estimate
  ◦ Getting worse and worse, the more we backprop through time

$$c_t = \tanh(U\,x_t + W c_{t-1})$$
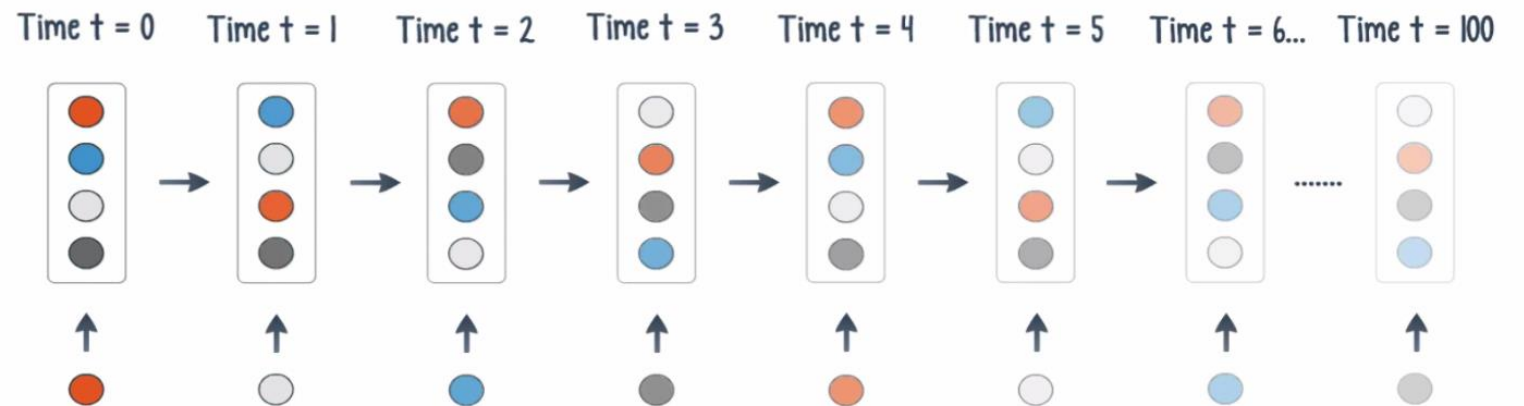$$y_t = \mathrm{softmax}(V\,c_t)$$

# Potential solutions

o Do fewer updates
- That might slow down training

o We can also make sure we do not backprop through more steps than our frequency of updates
- But then we do not compute the full gradients
- Bias again → not really gaining much

Vanishing gradients
Exploding gradients
Truncated backprop



Decay of information through time

Time t = 0   Time t = 1   Time t = 2   Time t = 3   Time t = 4   Time t = 5   Time t = 6...   Time t = 100

# An alternative formulation of an RNN

o Easier for mathematical analysis, and doesn't change the mechanics of the recurrent neural network

$$c_t = W \cdot \tanh(c_{t-1}) + U \cdot x_t + b$$

$$\mathcal{L} = \sum_t \mathcal{L}_t(c_t)$$

$$\theta = \{W, U, b\}$$

# What is the problem

o As we just saw, the gradient $\frac{\partial c_t}{\partial c_k}$ itself is subject to the chain rule

$$\frac{\partial c_t}{\partial c_k} = \frac{\partial c_t}{\partial c_{t-1}} \frac{\partial c_{t-1}}{\partial c_{t-2}} \cdots \frac{\partial c_{k+1}}{\partial c_k} = \prod_{j=k+1}^{t} \frac{\partial c_j}{\partial c_{j-1}}$$

o Product of ever expanding Jacobians
◦ Ever expanding because we multiply more and more for longer dependencies

# Let's look again the gradients

o Minimize the total loss over all time steps

$$\arg\min_{\theta} \sum_{t} \mathcal{L}_t(c_{t,\theta})$$

$$\frac{\partial \mathcal{L}_t}{\partial W} = \cdots$$

# Let's look again the gradients

o Minimize the total loss over all time steps

$$\arg\min_\theta \sum_t \mathcal{L}_t(c_{t,\theta})$$

$$\frac{\partial \mathcal{L}_t}{\partial W} = \sum_{\tau=1}^{t} \frac{\partial \mathcal{L}_t}{\partial c_t} \frac{\partial c_t}{\partial c_\tau} \frac{\partial c_\tau}{\partial W}$$

$$\frac{\partial \mathcal{L}}{\partial c_t} \frac{\partial c_t}{\partial c_\tau} = \frac{\partial \mathcal{L}}{\partial c_t} \cdot \underbrace{\frac{\partial c_t}{\partial c_{t-1}} \cdot \frac{\partial c_{t-1}}{\partial c_{t-2}}}_{t \ll \tau \; \rightarrow \; short\text{-}term \; factors} \cdot \ldots \cdot \underbrace{\frac{\partial c_{\tau+1}}{\partial c_\tau}}_{t \gg \tau \; \rightarrow \; long\text{-}term \; factors}$$

# Let's look again the gradients

o Minimize the total loss over all time steps

$$\arg \min_{\theta} \sum_{t} \mathcal{L}_t(c_{t,\theta})$$

$$\frac{\partial \mathcal{L}_t}{\partial W} = \sum_{\tau=1}^{t} \frac{\partial \mathcal{L}_t}{\partial c_t} \frac{\partial c_t}{\partial c_\tau} \frac{\partial c_\tau}{\partial W}$$

$$\frac{\partial \mathcal{L}}{\partial c_t} \frac{\partial c_t}{\partial c_\tau} = \frac{\partial \mathcal{L}}{\partial c_t} \cdot \frac{\partial c_t}{\partial c_{t-1}} \cdot \frac{\partial c_{t-1}}{\partial c_{t-2}} \cdot \ldots \cdot \frac{\partial c_{\tau+1}}{\partial c_\tau}$$

$$\left\| \frac{\partial c_{t+1}}{\partial c_t} \right\| \leq \|W^T\| \cdot \|diag(\sigma'(c_t))\|$$

# Let's look again the gradients

$$\left\|\frac{\partial c_{t+1}}{\partial c_t}\right\| \leq \|W^T\| \cdot \|diag(\sigma'(c_t))\|$$

o If we assume that the norm of the weight $W$ is bounded

　◦ Spectral radius (max eigenvalue) is smaller than an arbitrary small number $\lambda_1 < \frac{1}{\gamma}$

o And if we assume that the non linearity is bounded

$$\|diag(\sigma'(c_t))\| < \gamma$$
$$\Rightarrow$$
$$\left\|\frac{\partial c_{t+1}}{\partial c_t}\right\| < \frac{1}{\gamma}\gamma < 1$$

# Let's look again the gradients

- Minimize the total loss over all time steps

$$\arg \min_{\theta} \sum_{t} \mathcal{L}_t(c_{t,\theta})$$

$$\frac{\partial \mathcal{L}_t}{\partial W} = \sum_{\tau=1}^{t} \frac{\partial \mathcal{L}_t}{\partial c_t} \frac{\partial c_t}{\partial c_\tau} \frac{\partial c_\tau}{\partial W}$$

$$\frac{\partial \mathcal{L}}{\partial c_t} \frac{\partial c_t}{\partial c_\tau} = \frac{\partial \mathcal{L}}{\partial c_t} \cdot \underbrace{\frac{\partial c_t}{\partial c_{t-1}} \cdot \frac{\partial c_{t-1}}{\partial c_{t-2}}}_{t \ll \tau \;\to\; short\text{-}term\; factors} \cdot \ldots \cdot \underbrace{\frac{\partial c_{\tau+1}}{\partial c_\tau}}_{t \gg \tau \;\to\; long\text{-}term\; factors} \leq \eta^{t-\tau} \frac{\partial \mathcal{L}_t}{\partial c_t}$$

- RNN gradients expanding product of $\frac{\partial c_t}{\partial c_{t-1}}$

- With $\eta < 1$ long-term factors $\to 0$ exponentially fast

Pascanu, Mikolov, Bengio, On the difficulty of training recurrent neural networks, JMLR 2013

# Some cases

○ Let's assume we have 10 time steps and $\frac{\partial c_t}{\partial c_{t-1}} > 1$, e.g. $\frac{\partial c_t}{\partial c_{t-1}} = 1.5$

○ What would happen to the total $\frac{\partial \mathcal{L}_t}{\partial W}$?

# Some cases

○ Let's assume we have 100 time steps and $\frac{\partial c_t}{\partial c_{t-1}} > 1$, e.g. $\frac{\partial c_t}{\partial c_{t-1}} = 1.5$

○ What would happen to the total $\frac{\partial \mathcal{L}_t}{\partial W}$?

$$\frac{\partial \mathcal{L}}{\partial c_t} \frac{\partial c_t}{\partial c_\tau} \propto 1.5^{10} = 4.06 \cdot 10^{17}$$
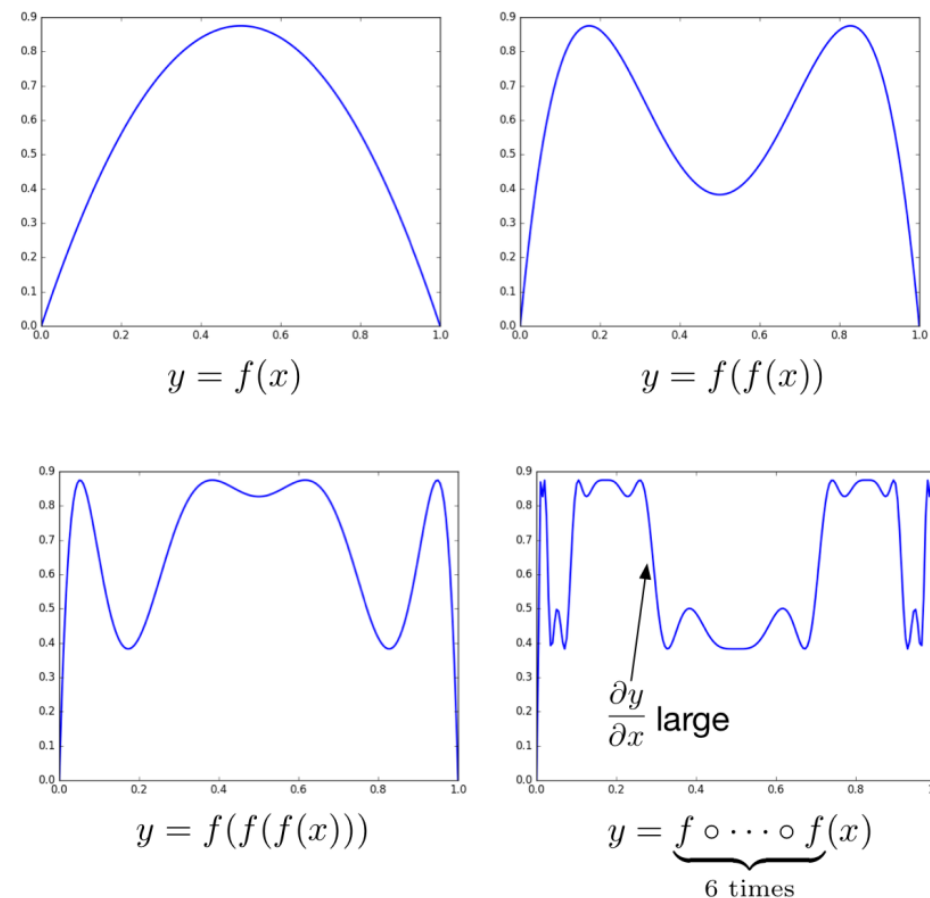
# Some cases

○ Let's assume now that $\frac{\partial c_t}{\partial c_{t-1}} < 1$, e.g. $\frac{\partial c_t}{\partial c_{t-1}} = 0.5$

○ What would happen to the total $\frac{\partial \mathcal{L}_t}{\partial W}$?

# Some cases

○ Let's assume now that $\frac{\partial c_t}{\partial c_{t-1}} < 1$, e.g. $\frac{\partial c_t}{\partial c_{t-1}} = 0.5$

○ What would happen to the total $\frac{\partial \mathcal{L}_t}{\partial W}$?

$$\frac{\partial \mathcal{L}}{\partial c_t} \frac{\partial c_t}{\partial c_\tau} \propto 0.5^{10} = 9.7 \cdot 10^{-5}$$

○ Do you think our optimizers like these kind of gradients?

# Some cases

o Let's assume now that $\frac{\partial c_t}{\partial c_{t-1}} < 1$, e.g. $\frac{\partial c_t}{\partial c_{t-1}} = 0.5$

o What would happen to the total $\frac{\partial \mathcal{L}_t}{\partial W}$?

$$\frac{\partial \mathcal{L}}{\partial c_t} \frac{\partial c_t}{\partial c_\tau} \propto 0.5^{10} = 9.7 \cdot 10^{-5}$$

o Do you think our optimizers like these kind of gradients?

o Too large ➔ unstable training, oscillations, divergence

o Too small ➔ very slow training, has it converged?

# A visual example

o Recurrent networks as iterated functions

Credit: R. Grosse



Figure 2: Iterations of the function $f(x) = 3.5\,x\,(1-x)$.

# Vanishing & Exploding Gradients

○ In recurrent networks, and in very deep networks in general (an RNN is not very different from an MLP), gradients are much affected by depth

$$\frac{\partial \mathcal{L}}{\partial c_t} = \frac{\partial \mathcal{L}}{\partial c_T} \cdot \frac{\partial c_T}{\partial c_{T-1}} \cdot \frac{\partial c_{T-1}}{\partial c_{T-2}} \cdot \ldots \cdot \frac{\partial c_{t+1}}{\partial c_{c_t}} \text{ and } \frac{\partial c_{t+1}}{\partial c_t} < 1 \Rightarrow \frac{\partial \mathcal{L}}{\partial W} \ll 1 \Rightarrow \text{Vanishing gradient}$$

$$\frac{\partial \mathcal{L}}{\partial c_t} = \frac{\partial \mathcal{L}}{\partial c_T} \cdot \frac{\partial c_T}{\partial c_{T-1}} \cdot \frac{\partial c_{T-1}}{\partial c_{T-2}} \cdot \ldots \cdot \frac{\partial c_{t+1}}{\partial c_{c_t}} \text{ and } \frac{\partial c_{t+1}}{\partial c_t} > 1 \Rightarrow \frac{\partial \mathcal{L}}{\partial W} \gg 1 \Rightarrow \text{Exploding gradient}$$

# Vanishing gradients & long memory

o Vanishing gradients are particularly a problem for long sequences

o Why?

# Vanishing gradients & long memory

o Vanishing gradients are particularly a problem for long sequences

o Why?

o Exponential decay

$$\frac{\partial \mathcal{L}}{\partial c_t} = \prod_{t \geq k \geq \tau} \frac{\partial c_k}{\partial c_{k-1}} = \prod_{t \geq k \geq \tau} W \cdot \partial \tanh(c_{k-1})$$

o The further back we look (long-term dependencies), the smaller the weights automatically become
  ◦ exponentially smaller weights

# Why are vanishing gradients bad?

- Weight updates focus on early time steps
- Updates for longer time steps become exponentially smaller
- Bad learning, even if we train the model exponentially longer. Why?
- Weights quickly learn (prefer) to "model" short-term transitions
  - And ignore long-term transitions
- At best, even after longer training, they will try "fine-tune" the whatever bad "modelling" of long-term transitions
  - After the short-term transitions are learned, the weights are set for them and are likely suboptimal for long-term
- Eventually, as the short-term transitions are inherently more prevalent, they will dominate the learning and gradients

$$\frac{\partial \mathcal{L}_1}{\partial W}$$

$$\frac{\partial \mathcal{L}_2}{\partial W}$$

$$\frac{\partial \mathcal{L}_3}{\partial W}$$

$$\frac{\partial \mathcal{L}_4}{\partial W}$$

$$\frac{\partial \mathcal{L}_5}{\partial W}$$

$$\frac{\partial \mathcal{L}}{\partial W} = \frac{\partial \mathcal{L}_1}{\partial W} + \frac{\partial \mathcal{L}_2}{\partial W} + \frac{\partial \mathcal{L}_3}{\partial w} + \frac{\partial \mathcal{L}_4}{\partial W} + \frac{\partial \mathcal{L}_5}{\partial W}$$

# Quick fix for exploding gradients: Rescaling!

○ First, get the gradient $\mathbf{g} \leftarrow \dfrac{\partial \mathcal{L}}{\partial w}$

○ Check if the norm is larger than a threshold $\theta_0$

○ If it is, rescale it to have same direction and threshold norm

$$\mathbf{g} \leftarrow \frac{\theta_0}{\|g\|}\, g$$

○ Simple, but works!

# An illustration



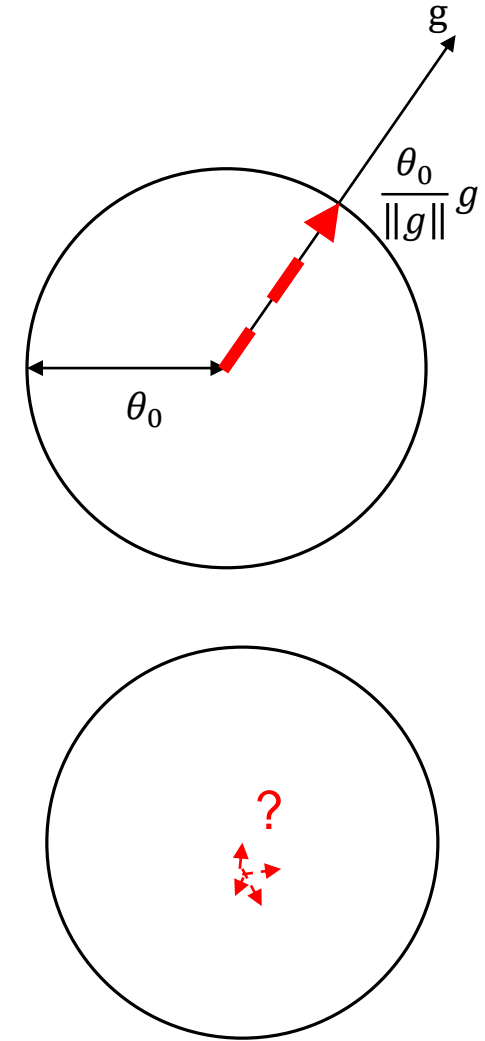Without clipping

With clipping

$J(w,b)$

$J(w,b)$

$w$

$b$

$w$

$b$

— Goodfellow et al., *Deep Learning*

# Can we rescale gradients also for vanishing gradients? No!

o The nature of the problem is different

o Exploding gradients → you might have bouncing and unstable optimization

o Vanishing gradients → you simply do not have a gradient to begin with
  ◦ Rescaling of what exactly?

o Unclear how would you rescale in a principled way, without affecting the rest of the time-steps

o In any case, even with re-scaling we would still focus on the short-term gradients
  ◦ Long-term dependencies would still be ignored

# Building intuition

o **With exploding gradients, the gradient is sort of good just too large**

  ◦ That is, the direction of the gradient is good, but the magnitude is too much

  ◦ Problem with optimization➔ bouncing, oscillation, etc.

o **With vanishing gradients, the gradient is not good in the first place**

  ◦ Neither the direction because of numerical instabilities, nor the magnitude are good

  ◦ Even if we rescale, are we sure we are going to change weights in the right direction? We cannot be sure.

# Biased gradients?

o Backpropagating all the way till infinity is unrealistic
  ◦ We would backprop forever (or simply it would be computationally very expensive)
  ◦ And in case, the gradients would be inaccurate because of intermediate updates

o What about truncating backprop to the last K steps

$$\tilde{g}_{t+1} \propto \left.\frac{\partial \mathcal{L}}{\partial w}\right|_{t=0}^{t=k}$$

o Unfortunately, this leads to biased gradients

$$g_{t+1} = \left.\frac{\partial \mathcal{L}}{\partial w}\right|_{t=0}^{t=\infty} \neq \tilde{g}_{t+1}$$

o Other algorithms exist but they are not as successful
  ◦ Maybe we will visit them later

# LSTM and variants

# How to fix the vanishing gradients?

o Error signal over time must have not too large, not too small norm

o Let's have a look at the loss function

$$\frac{\partial \mathcal{L}_t}{\partial W} = \sum_{\tau=1}^{t} \frac{\partial \mathcal{L}_r}{\partial y_t} \frac{\partial y_t}{\partial c_t} \textcolor{red}{\frac{\partial c_t}{\partial c_\tau}} \frac{\partial c_\tau}{\partial W}$$

$$\textcolor{red}{\frac{\partial c_t}{\partial c_\tau}} = \prod_{t \geq k \geq \tau} \frac{\partial c_k}{\partial c_{k-1}}$$

o How to make the product roughly the same no matter the length?

# How to fix the vanishing gradients?

o Error signal over time must have not too large, not too small norm

o Let's have a look at the loss function

$$\frac{\partial \mathcal{L}_t}{\partial W} = \sum_{\tau=1}^{t} \frac{\partial \mathcal{L}_r}{\partial y_t} \frac{\partial y_t}{\partial c_t} \textcolor{red}{\frac{\partial c_t}{\partial c_\tau}} \frac{\partial c_\tau}{\partial W}$$

$$\textcolor{red}{\frac{\partial c_t}{\partial c_\tau}} = \prod_{t \geq k \geq \tau} \frac{\partial c_k}{\partial c_{k-1}}$$

o How to make the product roughly the same no matter the length?

o Use the identity function with gradient of 1

# Main idea of LSTMs

o Over time the state change is $c_{t+1} = c_t + \Delta c_{t+1}$

o This constant over-writing over long time steps leads to chaotic behavior

o Input weight conflict
  ◦ Are all inputs important enough to write them down?

o Output conflict
  ◦ Are all outputs important enough to be read?

o Forget conflict
  ◦ Is all information important enough to be remembered over time?

# LSTMs

○ RNNs

$$c_t = W \cdot \tanh(c_{t-1}) + U \cdot x_t + b$$

○ LSTMs

$$i = \sigma\big(x_t U^{(i)} + m_{t-1} W^{(i)}\big)$$
$$f = \sigma\big(x_t U^{(f)} + m_{t-1} W^{(f)}\big)$$
$$o = \sigma\big(x_t U^{(o)} + m_{t-1} W^{(o)}\big)$$
$$\widetilde{c}_t = \tanh(x_t U^{(g)} + m_{t-1} W^{(g)})$$
$$c_t = c_{t-1} \odot f + \widetilde{c}_t \odot i$$
$$m_t = \tanh(c_t) \odot o$$

# LSTMs: A marking difference

- RNNs

$$c_t = W \cdot \tanh(c_{t-1}) + U \cdot x_t + b$$

- LSTMs

$$i = \sigma\left(x_t U^{(i)} + m_{t-1} W^{(i)}\right)$$
$$f = \sigma\left(x_t U^{(f)} + m_{t-1} W^{(f)}\right)$$
$$o = \sigma\left(x_t U^{(o)} + m_{t-1} W^{(o)}\right)$$
$$\tilde{c}_t = \tanh\left(x_t U^{(g)} + m_{t-1} W^{(g)}\right)$$
$$c_t = c_{t-1} \odot f + \tilde{c}_t \odot i$$
$$m_t = \tanh(c_t) \odot o$$

- The previous state $c_{t-1}$ and the next state $c_t$ are also connected by addition
  - It is also connected by the **tanh**, but at least there is the addition to make sure of good gradients

Additivity leads to strong gradients
Bounded by sigmoidal $f$



Nice tutorial: http://colah.github.io/posts/2015-08-Understanding-LSTMs/

# Cell state

$$i = \sigma\left(x_t U^{(i)} + m_{t-1} W^{(i)}\right)$$
$$f = \sigma\left(x_t U^{(f)} + m_{t-1} W^{(f)}\right)$$
$$o = \sigma\left(x_t U^{(o)} + m_{t-1} W^{(o)}\right)$$
$$\widetilde{c}_t = \tanh\left(x_t U^{(g)} + m_{t-1} W^{(g)}\right)$$
$$c_t = c_{t-1} \odot f + \widetilde{c}_t \odot i$$
$$m_t = \tanh(c_t) \odot o$$



Cell state line

# LSTM nonlinearities

$i = \sigma(x_t U^{(i)} + m_{t-1} W^{(i)})$

$f = \sigma(x_t U^{(f)} + m_{t-1} W^{(f)})$

$o = \sigma(x_t U^{(o)} + m_{t-1} W^{(o)})$

$\widetilde{c}_t = \tanh(x_t U^{(g)} + m_{t-1} W^{(g)})$

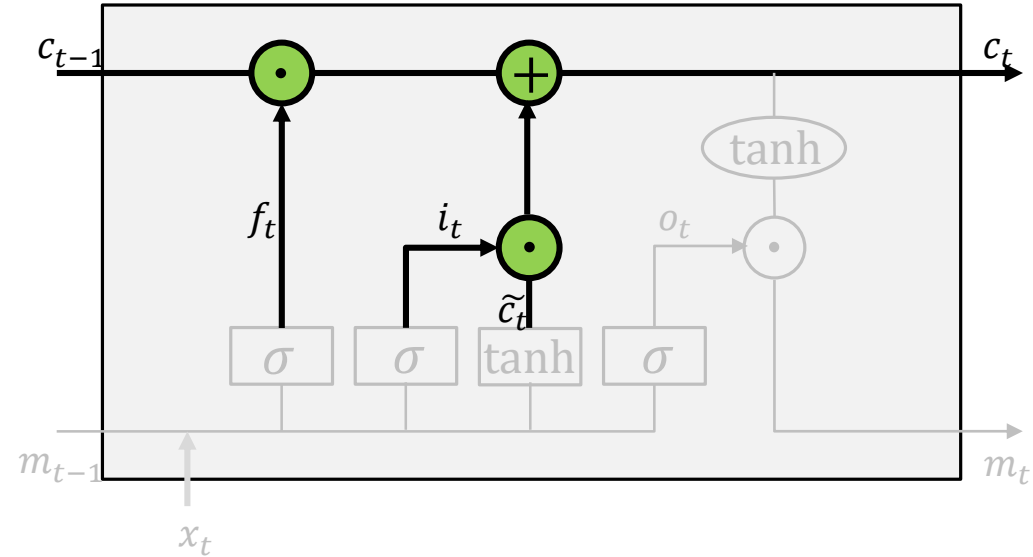$c_t = c_{t-1} \odot f + \widetilde{c}_t \odot i$

$m_t = \tanh(c_t) \odot o$



○ $\sigma \in (0, 1)$: control gate – something like a switch

○ $\tanh \in (-1, 1)$: recurrent nonlinearity

# LSTM Step by Step #1

$i = \sigma\left(x_t U^{(i)} + m_{t-1} W^{(i)}\right)$

$f = \sigma\left(x_t U^{(f)} + m_{t-1} W^{(f)}\right)$

$o = \sigma\left(x_t U^{(o)} + m_{t-1} W^{(o)}\right)$

$\widetilde{c}_t = \tanh\left(x_t U^{(g)} + m_{t-1} W^{(g)}\right)$

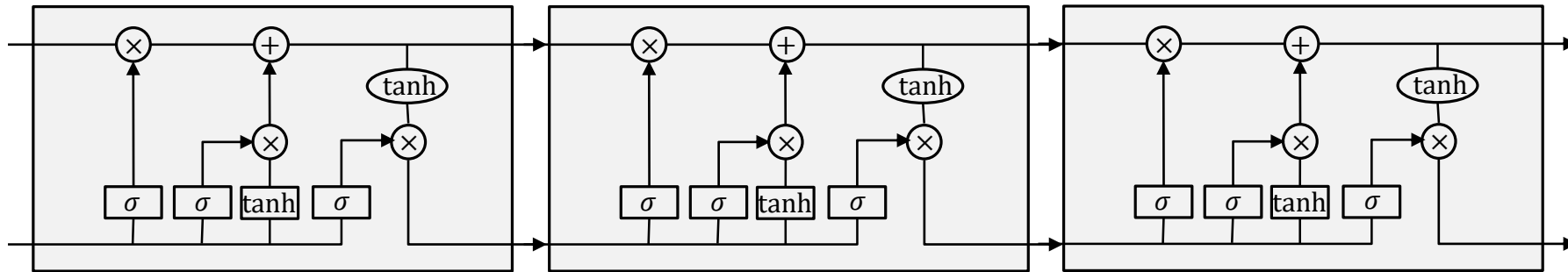$c_t = c_{t-1} \odot f + \widetilde{c}_t \odot i$

$m_t = \tanh(c_t) \odot o$

$$i = \sigma\big(x_t U^{(i)} + m_{t-1} W^{(i)}\big)$$
$$f = \sigma\big(x_t U^{(f)} + m_{t-1} W^{(f)}\big)$$
$$o = \sigma\big(x_t U^{(o)} + m_{t-1} W^{(o)}\big)$$
$$\widetilde{c}_t = \tanh(x_t U^{(g)} + m_{t-1} W^{(g)})$$
$$c_t = c_{t-1} \odot f + \widetilde{c}_t \odot i$$
$$m_t = \tanh(c_t) \odot o$$



o Decide what new information is relevant from the new input and should be added to the new memory
  ◦ Modulate the input $i_t$
  ◦ Generate candidate memories $\widetilde{c}_t$

$$i = \sigma\big(x_t U^{(i)} + m_{t-1} W^{(i)}\big)$$
$$f = \sigma\big(x_t U^{(f)} + m_{t-1} W^{(f)}\big)$$
$$o = \sigma\big(x_t U^{(o)} + m_{t-1} W^{(o)}\big)$$
$$\widetilde{c}_t = \tanh(x_t U^{(g)} + m_{t-1} W^{(g)})$$
$$c_t = c_{t-1} \odot f + \widetilde{c}_t \odot i$$
$$m_t = \tanh(c_t) \odot o$$



- Compute and update the current cell state $c_t$
  - Depends on the previous cell state
  - What we decide to forget
  - What inputs we allow
  - The candidate memories

$$i = \sigma\left(x_t U^{(i)} + m_{t-1} W^{(i)}\right)$$
$$f = \sigma\left(x_t U^{(f)} + m_{t-1} W^{(f)}\right)$$
$$o = \sigma\left(x_t U^{(o)} + m_{t-1} W^{(o)}\right)$$
$$\tilde{c}_t = \tanh\left(x_t U^{(g)} + m_{t-1} W^{(g)}\right)$$
$$c_t = c_{t-1} \odot f + \tilde{c}_t \odot i$$
$$m_t = \tanh(c_t) \odot o$$



o Modulate the output
  ◦ Does the new cell state relevant? → Sigmoid 1
  ◦ If not → Sigmoid 0
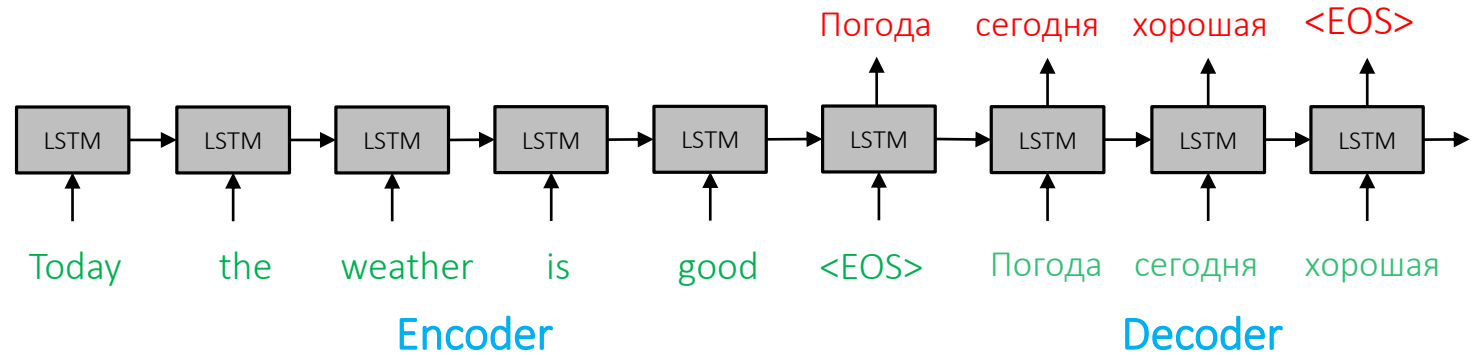
o Generate the new memory

# Unrolling the LSTMs

o Just the same like for RNNs

o The engine is a bit different (more complicated)
  ◦ Because of their gates LSTMs capture long and short term dependencies

# LSTM variants

○ LSTM with peephole connections

○ Gates have access also to the previous cell states $c_{t-1}$ (not only memories)

○ Bi-directional recurrent networks

○ Gated Recurrent Units (GRU)

○ Phased LSTMs

○ Skip LSTMs

○ And many more …

# Encoder-Decoder Architectures

Погода   сегодня   хорошая   <EOS>

| LSTM | LSTM | LSTM | LSTM | LSTM | LSTM | LSTM | LSTM | LSTM |

Today    the    weather    is    good    <EOS>    Погода    сегодня    хорошая

Encoder                                                                Decoder

# Machine translation

o The phrase in the source language is one sequence
  ◦ "Today the weather is good"

o It is captured by an Encoder LSTM

o The phrase in the target language is also a sequence
  ◦ "Погода сегодня хорошая"

o It is captured by a Decoder LSTM

# Image captioning

o Similar to image translation

o The only difference is that the Encoder LSTM is an image ConvNet
  ◦ VGG, ResNet, …

o Keep decoder the same

# Image captioning demo

NeuralTalk and Walk, recognition, text description of the image while walking
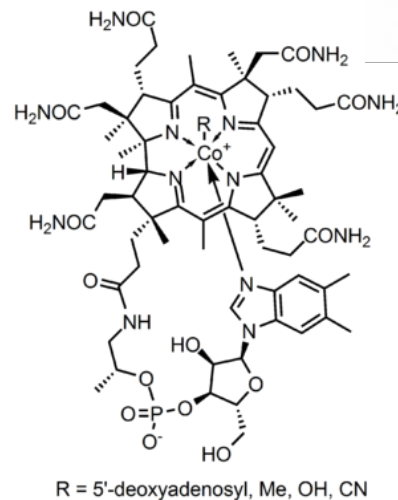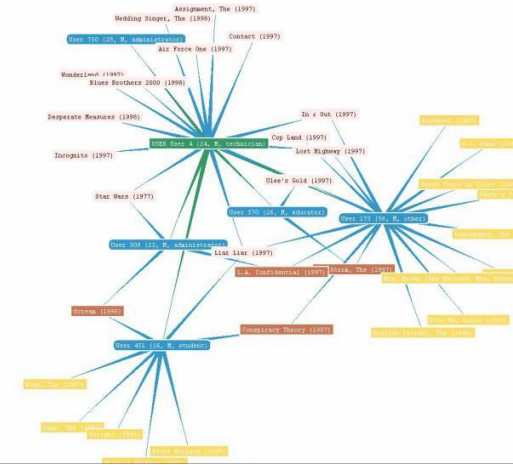
# Graph Neural Networks

# Why Graphs?

o Many domains & data have graph structure

o Examples?

# Why Graphs?

o Many domains & data have graph structure

o Social networks

o Knowledge graphs

o Recommender systems

o Chemical compounds

o And more
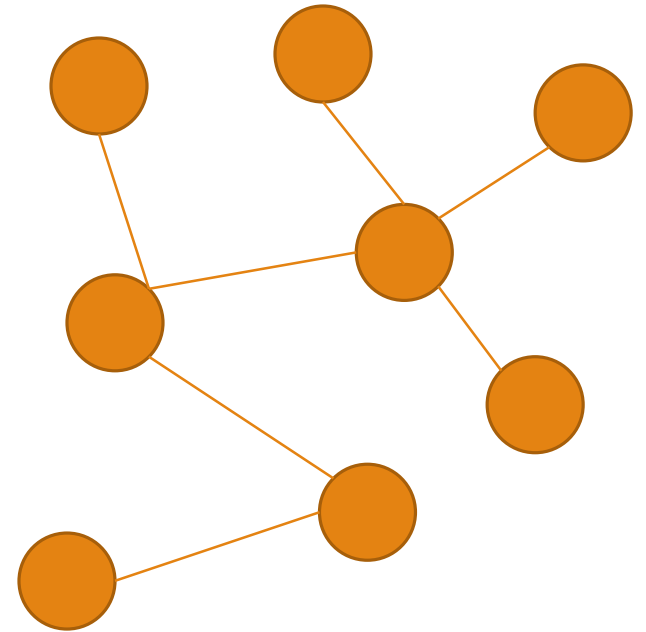
# Predictions tasks on graphs?

# Predictions tasks on graphs?

o Node classification

o Filling out missing edges

o Filling out missing nodes

o Novel graph generation

# DeepWalk

Algorithm

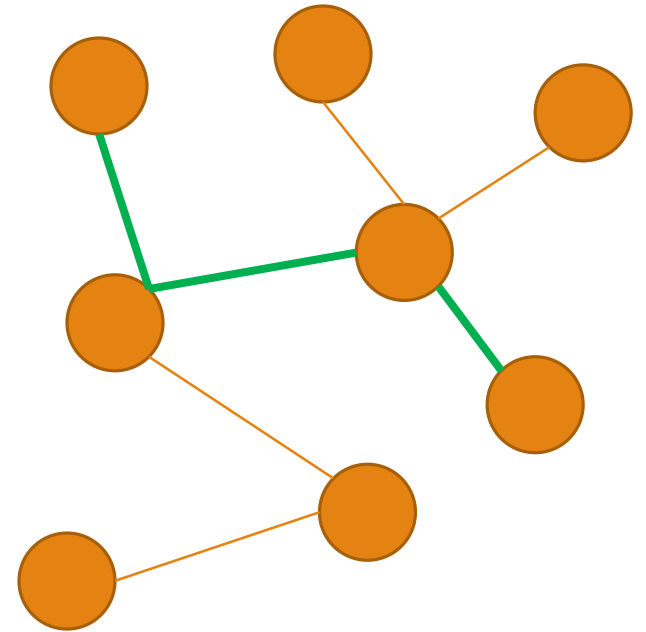1. Perform random walks on the graph to generate node sequences



DeepWalk: Online Learning of Social Representations, Perozzi et al, 2014

# DeepWalk

Algorithm

1. Perform random walks on the graph to generate node sequences
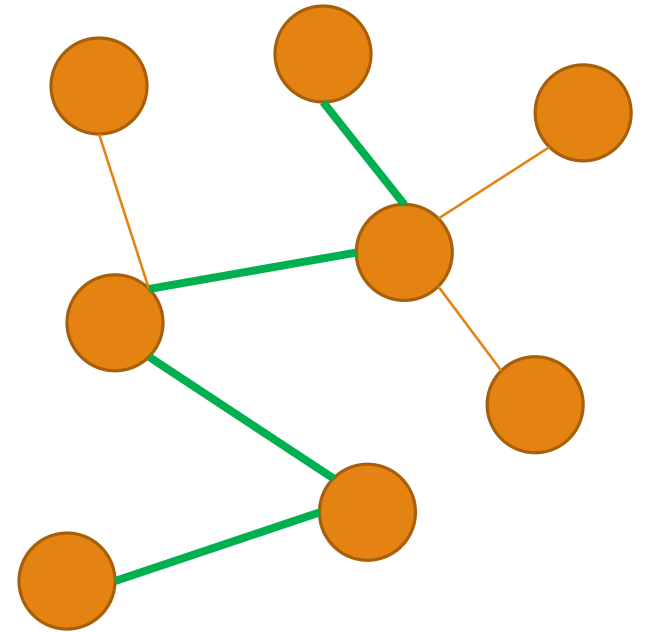2. Run skip-gram to learn the node embedding

DeepWalk: Online Learning of Social Representations, Perozzi et al, 2014

# DeepWalk

Algorithm

1. Perform random walks on the graph to generate node sequences



DeepWalk: Online Learning of Social Representations, Perozzi et al, 2014

# DeepWalk

Algorithm

1. Perform random walks on the graph to generate node sequences
2. Run skip-gram to learn node embeddings



DeepWalk: Online Learning of Social Representations, Perozzi et al, 2014

# DeepWalk

Algorithm

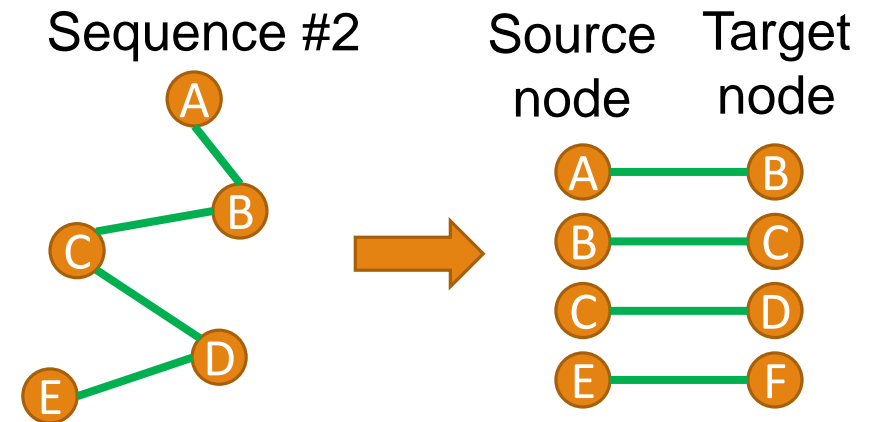1. Perform random walks on the graph to generate node sequences
2. Run skip-gram to learn node embeddings

Sequence #2    Source node    Target node

DeepWalk: Online Learning of Social Representations, Perozzi et al, 2014

# DeepWalk

Algorithm

1. Perform random walks on the graph to generate node sequences
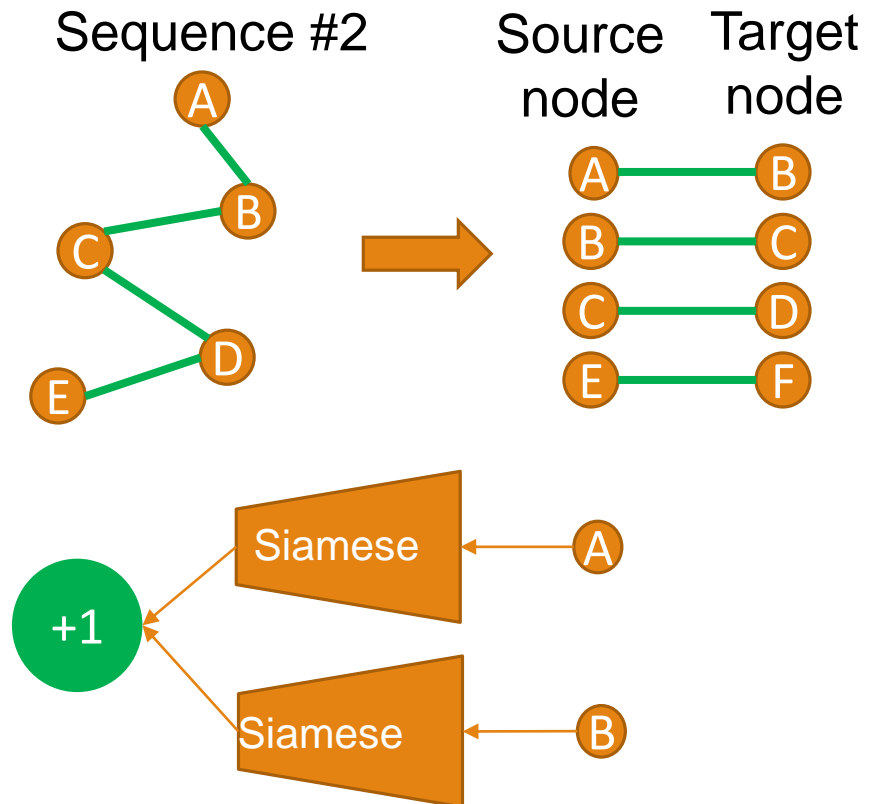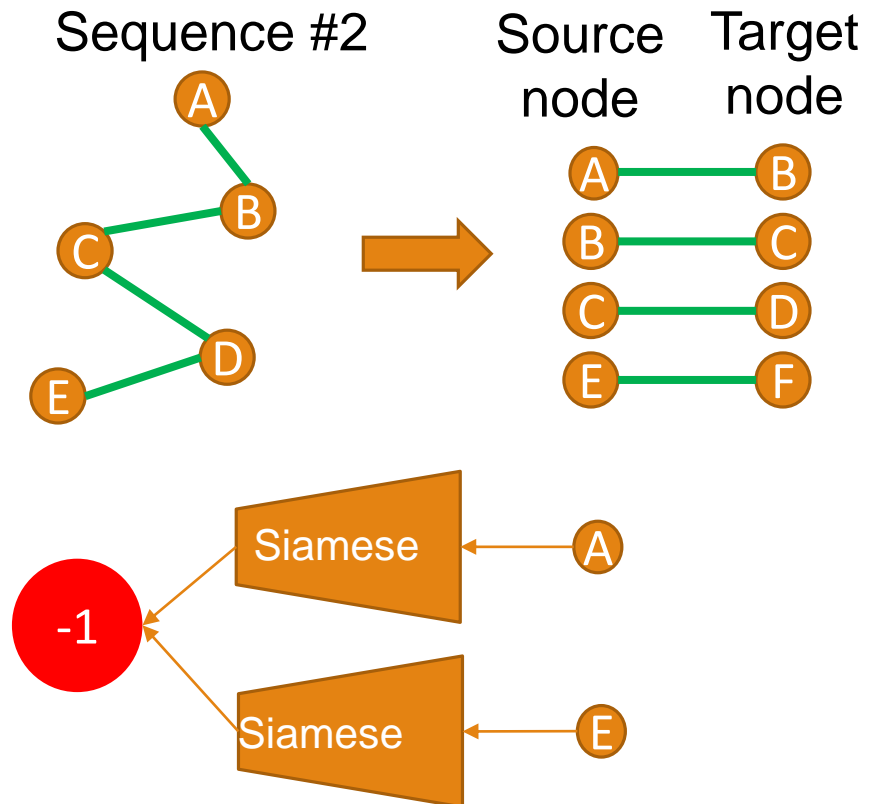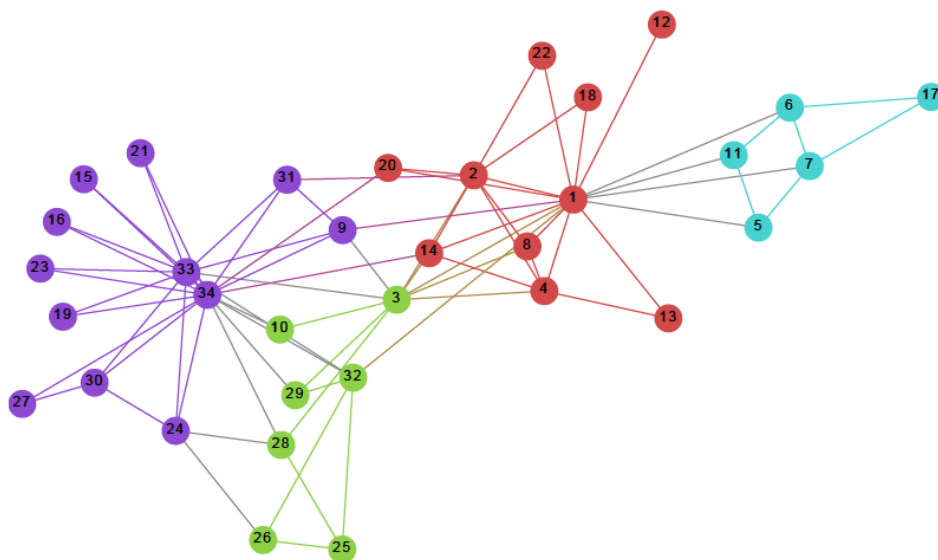2. Run skip-gram to learn node embeddings

DeepWalk: Online Learning of Social Representations, Perozzi et al, 2014

# DeepWalk: Results



(a) Input: Karate Graph    (b) Output: Representation

DeepWalk: Online Learning of Social Representations, Perozzi et al, 2014

# DeepWalk: A problem

o The method is transductive

o Whenever a new node is added to the graph, the model must be retrained

o This is not useful for dynamic graphs

DeepWalk: Online Learning of Social Representations, Perozzi et al, 2014

# GraphSage

**Algorithm 1:** GraphSAGE embedding generation (i.e., forward propagation) algorithm

**Input** : Graph $\mathcal{G}(\mathcal{V}, \mathcal{E})$; input features $\{\mathbf{x}_v, \forall v \in \mathcal{V}\}$; depth $K$; weight matrices $\mathbf{W}^k, \forall k \in \{1, ..., K\}$; non-linearity $\sigma$; differentiable aggregator functions $\text{AGGREGATE}_k, \forall k \in \{1, ..., K\}$; neighborhood function $\mathcal{N} : v \rightarrow 2^{\mathcal{V}}$

**Output :** Vector representations $\mathbf{z}_v$ for all $v \in \mathcal{V}$

1 $\mathbf{h}_v^0 \leftarrow \mathbf{x}_v, \forall v \in \mathcal{V}$ ;
2 **for** $k = 1...K$ **do**
3      **for** $v \in \mathcal{V}$ **do**
4          $\mathbf{h}_{\mathcal{N}(v)}^k \leftarrow \text{AGGREGATE}_k(\{\mathbf{h}_u^{k-1}, \forall u \in \mathcal{N}(v)\});$
5          $\mathbf{h}_v^k \leftarrow \sigma\left(\mathbf{W}^k \cdot \text{CONCAT}(\mathbf{h}_v^{k-1}, \mathbf{h}_{\mathcal{N}(v)}^k)\right)$
6      **end**
7      $\mathbf{h}_v^k \leftarrow \mathbf{h}_v^k / \|\mathbf{h}_v^k\|_2, \forall v \in \mathcal{V}$
8 **end**
9 $\mathbf{z}_v \leftarrow \mathbf{h}_v^K, \forall v \in \mathcal{V}$

GraphSage: Inductive Representation Learning on Large Graphs, Hamilton et al., 2017
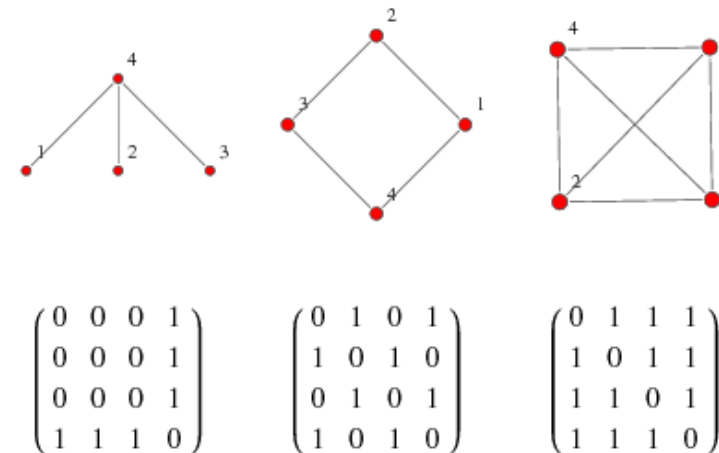
# GraphSage: How to aggregate?

- Mean aggregation $\quad \mathbf{h}_v^k \leftarrow \sigma(\mathbf{W} \cdot \text{MEAN}(\{\mathbf{h}_v^{k-1}\} \cup \{\mathbf{h}_u^{k-1}, \forall u \in \mathcal{N}(v)\})$

- LSTM aggregation

- Pooling aggregation $\quad \text{AGGREGATE}_k^{\text{pool}} = \max(\{\sigma\left(\mathbf{W}_{\text{pool}}\mathbf{h}_{u_i}^k + \mathbf{b}\right), \forall u_i \in \mathcal{N}(v)\})$

- Loss $\quad J_{\mathcal{G}}(\mathbf{z}_u) = -\log\left(\sigma(\mathbf{z}_u^\top \mathbf{z}_v)\right) - Q \cdot \mathbb{E}_{v_n \sim P_n(v)} \log\left(\sigma(-\mathbf{z}_u^\top \mathbf{z}_{v_n})\right)$

# Graph Convolutional Networks

o Assuming a graph $G = (\mathcal{V}, \mathcal{E})$

o A node has a description $x_i$, all stored in a $N \times D$ matrix $X = [\ldots, x_i, \ldots]$

o The graph structure is encoded by the adjacency matrix $A$

o A neural network on this graph then is
$$H^{(l+1)} = h(H^{(l)}, A)$$

Graph Convolutional Networks, Kipf and Welling, 2016

$$\begin{pmatrix} 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{pmatrix} \quad \begin{pmatrix} 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 \end{pmatrix} \quad \begin{pmatrix} 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{pmatrix}$$

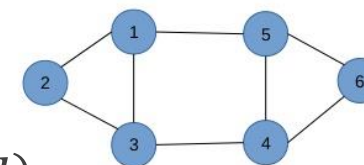# Graph Convolutional Networks: A simple example

○ $h(H^{(l)}, A) = \sigma(AH^{(l)}W^{(l)})$

○ Two problems

   ◦ Given a node, the adjacency matrix $A$ considers neighboring nodes but not the node itself → Aggregation does not use the node itself

   ◦ A node might have different numbers of neighbors and change the scale of the multiplication

○ Add the identity matrix to $A$

**Degree matrix**

$$D_{ij} = \begin{cases} d(i) & \text{if } i = j \\ 0 & \text{otherwise} \end{cases}$$

○ Left multiply by $D^{-1}A$: $D$ is the degree matrix

○ Combining all, we have the following module

$$h(H^{(l)}, A) = \sigma(D^{-\frac{1}{2}} \hat{A} D^{-\frac{1}{2}} H^{(l)} W^{(l)})$$
$$\hat{A} = A + I$$

Graph Convolutional Networks, Kipf and Welling, 2016

$$D = \begin{bmatrix} 3 & 0 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 0 & 0 \\ 0 & 0 & 3 & 0 & 0 & 0 \\ 0 & 0 & 0 & 3 & 0 & 0 \\ 0 & 0 & 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 0 & 0 & 2 \end{bmatrix}$$

# Summary

o Sequential data

o Recurrent Neural Networks

o Backpropagation through time

o Exploding and vanishing gradients

o LSTMs and variants

o Encoder-Decoder Architectures

o Graph Neural Networks