

Plan

Recap

Reminders

Markov and n-gram

Recurrent Networks

↳ Architecture

↳ Loss

↳ Optimization

Improvements

Recap

Object Detection

1) grad CAM

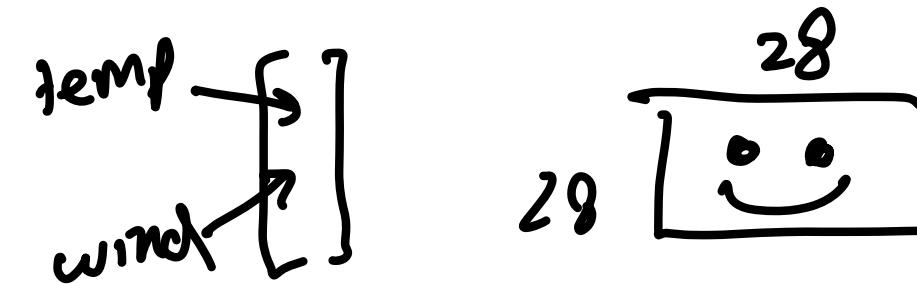
2) bounding boxes \leftarrow most popular

3) segmentation mask

Reminders

- Form 20/26
- Homework
 - ↳ 13.5 and above 😊
 - ↳ 12 to 13 😐
 - ↳ below 😞 talk to me!
 - ↳ Future submissions in LaTeX
- Project Proposal due Monday

Text Applications



- document retrieval
- convert audio to text
- translate between languages
- next word prediction

1) How do we represent words as numbers?
↳ assign numbers to words/letters

problems: distance

↳ one hot encoding

[?] How $\begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix} \leftarrow$ do

2) Varying length and long range dependency

"I am from France. ..."

I am going back home to —."

Assume we've solved 1)
and we have $w \rightarrow$ meaningful
vector encoding

$\lambda = (w^1, w^2, w^3, \dots, w^T)$ length
of sequence

what is $P(s)$?

$$= P(w^1) \cdot P(w^2 | w^1) \cdot P(w^3 | w^2, w^1) \cdot \dots$$

$\frac{P(w^2, w^1)}{P(w^2)}$ 2-gram
(word 1, word 2)
 $\sqrt{2}$

$$P(w^T | w^{T-1}, w^{T-2}, \dots, w^1)$$

n -gram
(word 1, word 2, ..., word n)
 \sqrt{n}

$P(1)$

$$= P(w^1) \cdot P(w^2|w^1) \cdot P(w^3|w^1, w^2) \dots$$

\sim Markov

$$\sim P(w^1) \cdot P(w^2|w^1) \cdot P(w^3|w^2) \dots$$

2-gram

Predicting next word

When Gamaliel Painter died, he was Middlebury's pride,
 A sturdy pioneer without a stain;
 And he left his all by will, to the college on the hill,
 And included his codicil cane.

Oh, its rap **rap** **rap**, and it's tap **tap** **tap**,
 If you listen you can hear it sounding plain;
 For a helper true and tried, as the generations glide,
 There is nothing like **Gamaliel Painter's** cane. [5][6]

What comes after died?
 w'

$$P(w|w') = \frac{P(w' \text{ and } w)}{P(w')}$$

2 grams with he

$$(he, was) \quad P(was/he) = \frac{P(\overset{1}{he} \overset{2}{was})}{P(he)}$$

(rap, rap)

$$P(rap/rap) = \frac{1}{2}$$

• v^2 small

• long

range dependency

small

n -gram

big

• more info

• v^n exponentially big

$$P(w^t | w^{t-1}, w^{t-2}, \dots, w^1)$$

$$\approx f(w^{t-1}, h^{t-1})$$

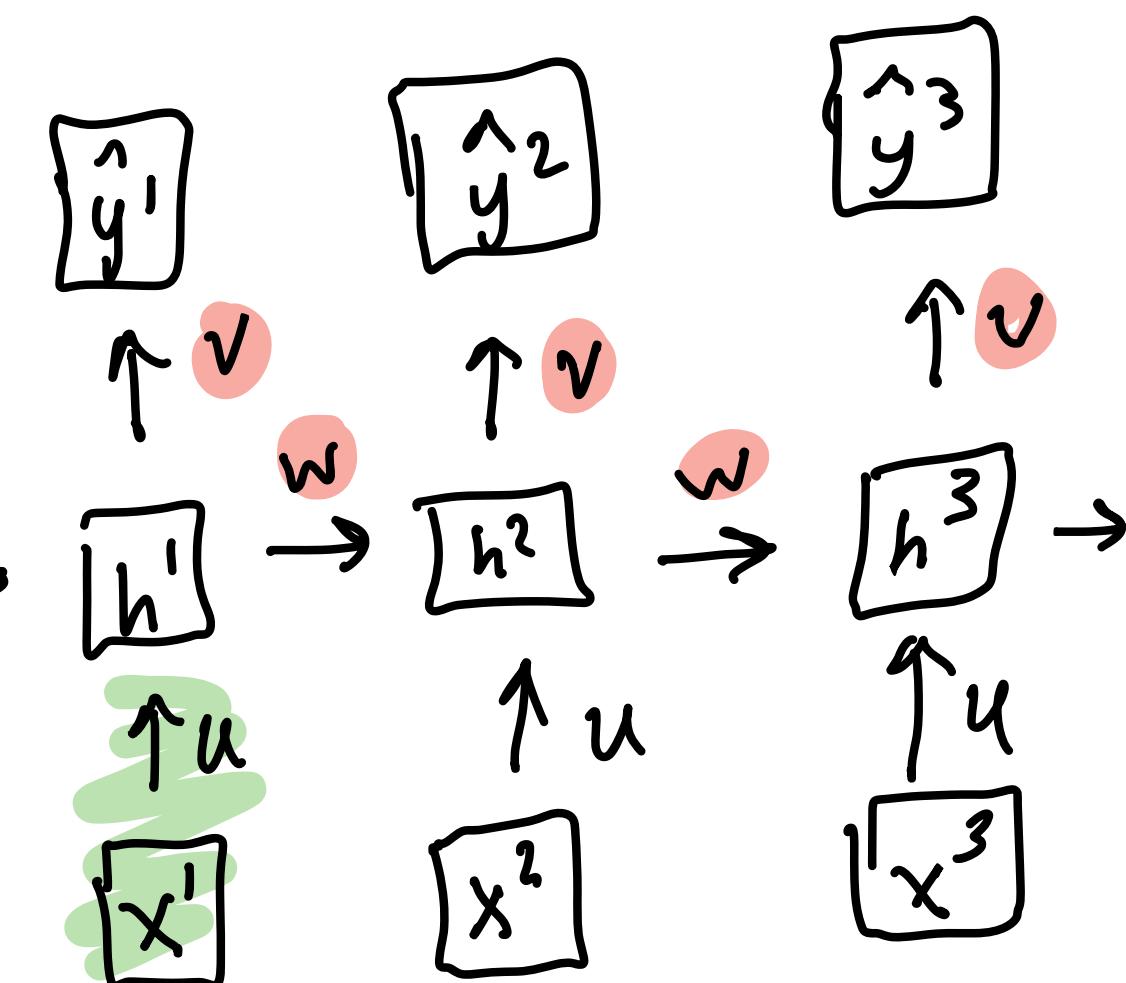
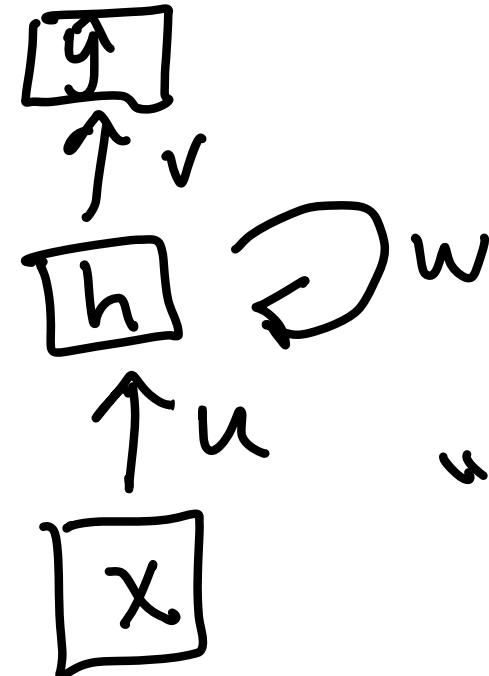
\downarrow encodes
information
 w^{t-1}, \dots, w^1

$$h \in \mathbb{R}^d$$

$$h^t = \sigma(Ux^t + Wh^{t-1})$$

$d \times l \quad l \times 1 \quad d \times d \quad d \times 1$

$$\hat{y}^t = \text{softmax}(vh^t)$$

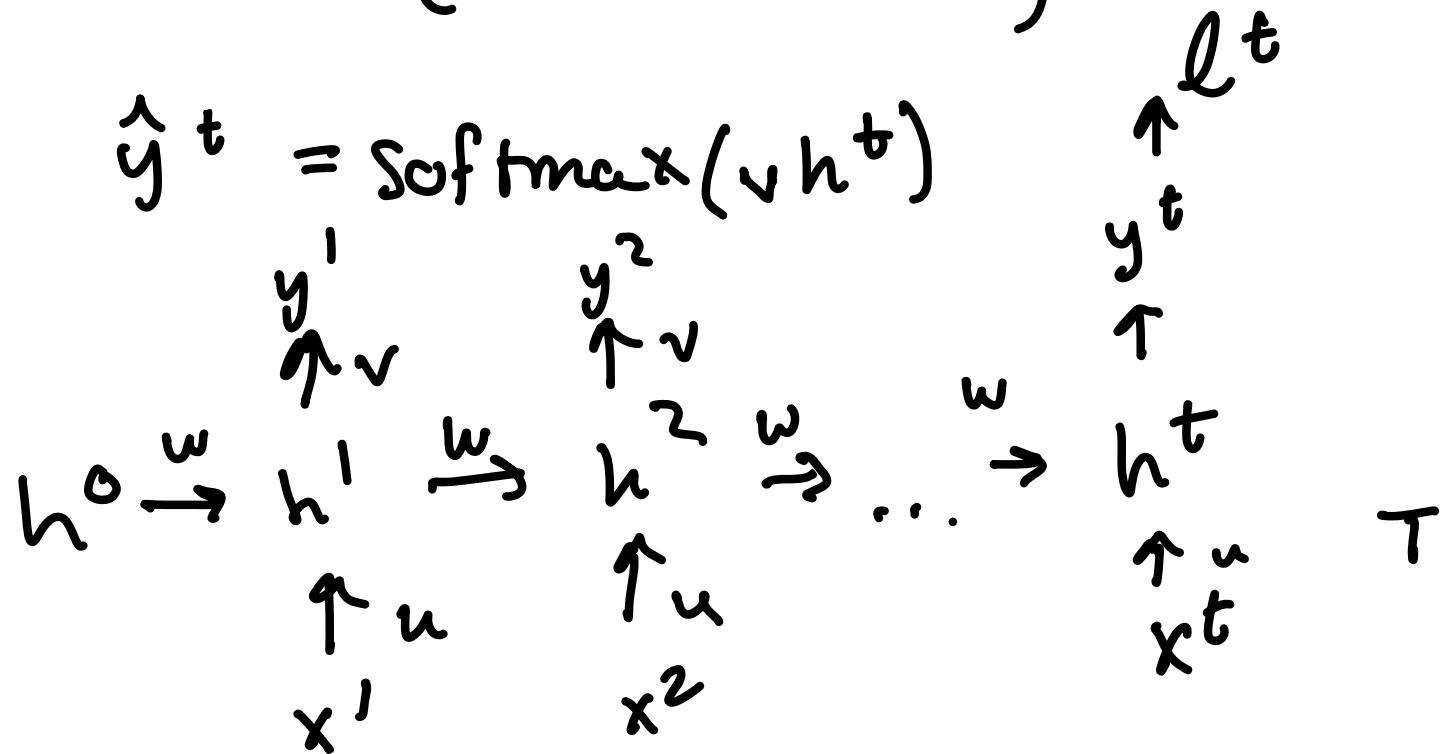


① weight sharing

② variable length

$$h^t = \sigma(u x^t + w h^{t-1})$$

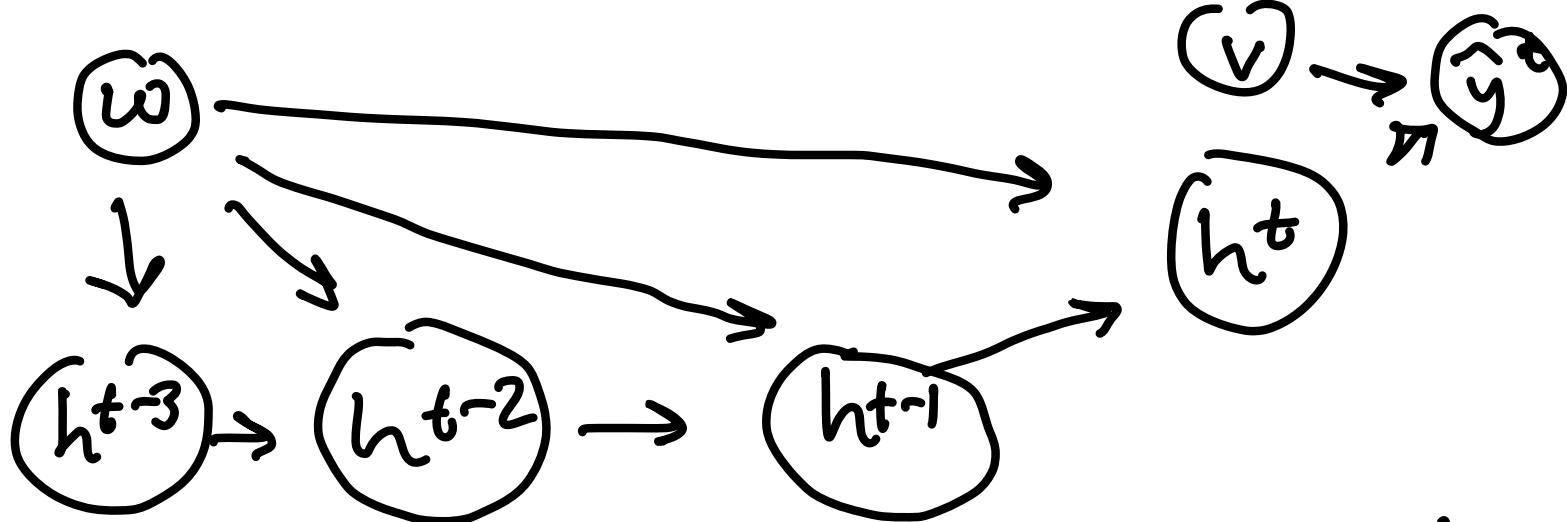
$$\hat{y}^t = \text{softmax}(v h^t)$$



$$\frac{\partial \mathcal{L}}{\partial w} = \frac{\partial}{\partial w} \left(\frac{1}{T} \sum_{t=1}^T l^t \right)$$

$$= \frac{1}{T} \sum_{t=1}^T \frac{\partial l^t}{\partial w}$$

$$\underbrace{\frac{\partial l^t}{\partial y^t}}_{\frac{\partial l^t}{\partial y_t}} \cdot \underbrace{\frac{\partial y^t}{\partial h^t}}_{\frac{\partial y_t}{\partial h^t}} \cdot \underbrace{\frac{\partial h^t}{\partial w}}$$



$$= \alpha > 1$$

$$\begin{aligned} \frac{\partial h^t}{\partial w} &= \frac{\partial h^t}{\partial w} + \frac{\partial h^t}{\partial h^{t-1}} \cdot \frac{\partial h^{t-1}}{\partial w} \\ &\quad + \frac{\partial h^t}{\partial h^{t-1}} \cdot \frac{\partial h^{t-1}}{\partial h^{t-2}} \cdot \frac{\partial h^{t-2}}{\partial w} + \dots \end{aligned}$$

Solutions

- truncate (only last n partial derivatives)

- $\nabla' \mathcal{L}(w) \leftarrow \frac{\nabla \mathcal{L}(w)}{\|\nabla \mathcal{L}(w)\|}$ scaling
- more complicated architectures

castle in

\hat{y}^1
 \hat{y}^2

$$h^0 \xrightarrow{w} h^1 \xrightarrow{w} h^2 \xrightarrow{w} \dots \xrightarrow{w} h^t$$

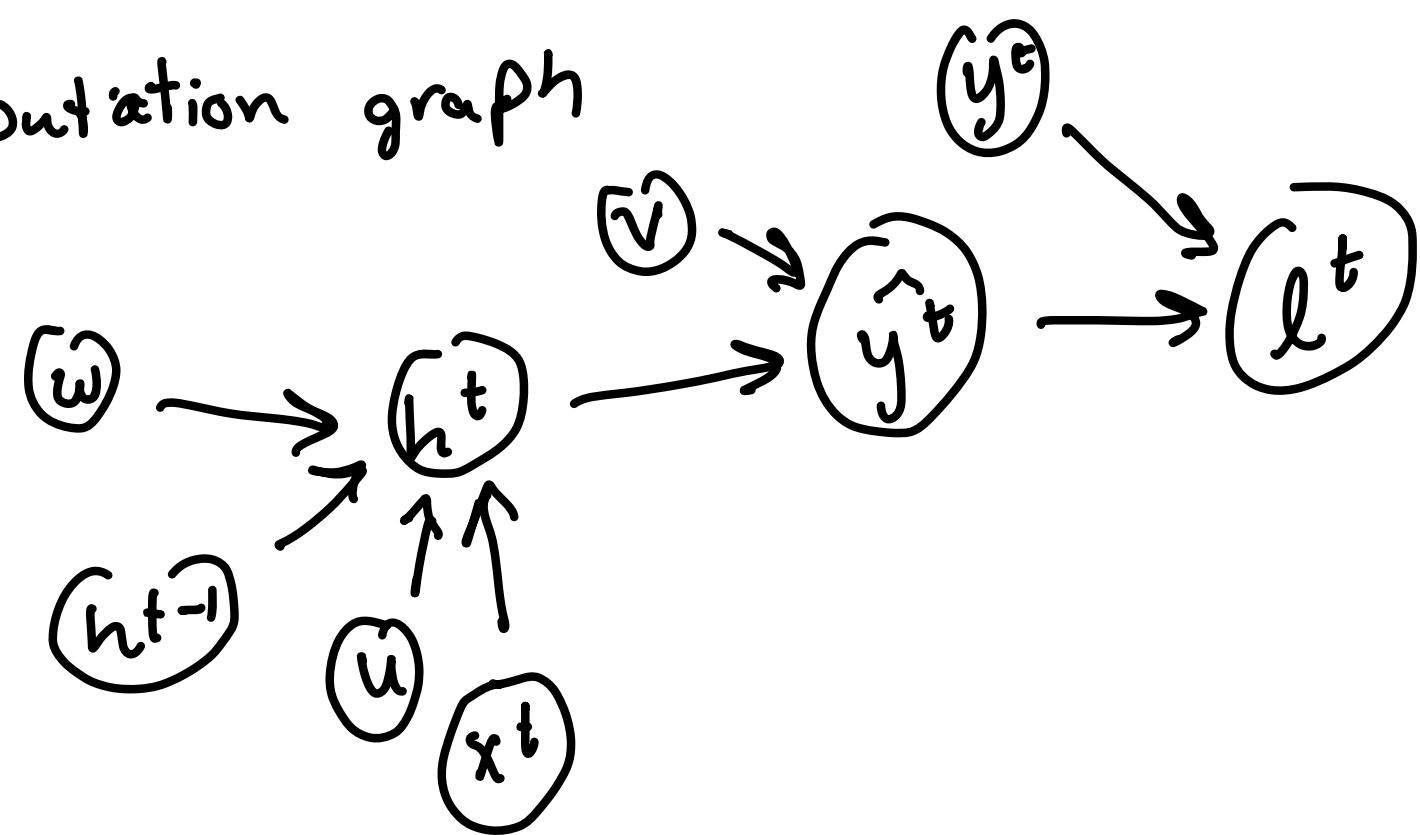
$\uparrow u$ $\uparrow u$ $\uparrow u$

The cat in the hat

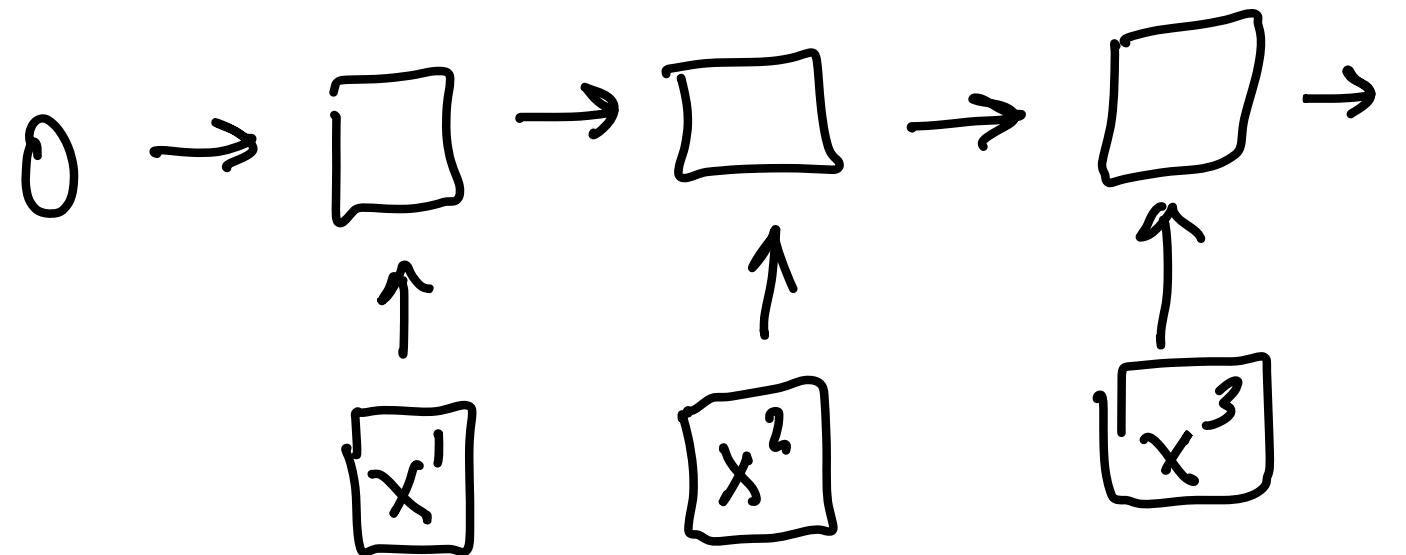
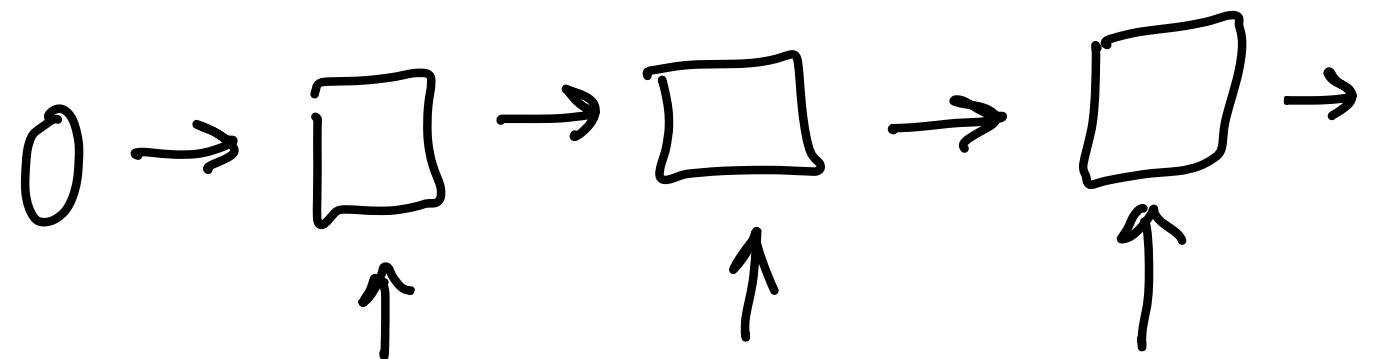
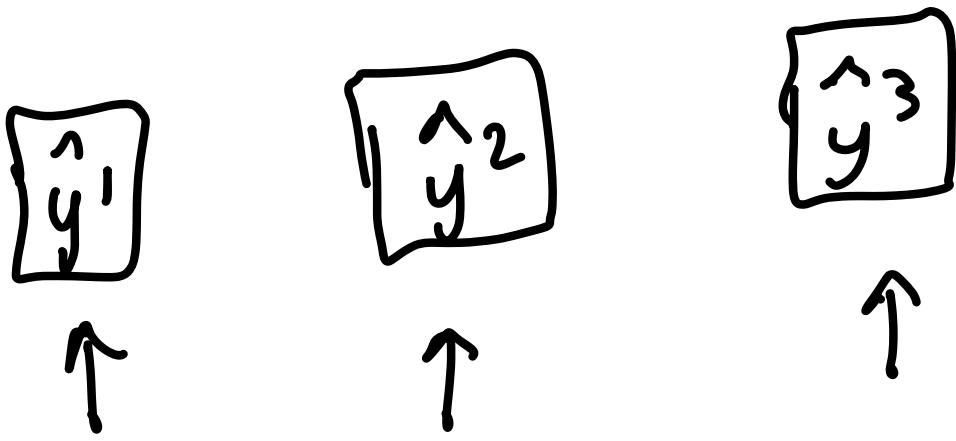
$$h^t = \sigma(u x^t + w h^{t-1})$$

$$\hat{y}^t = \text{Softmax}(v h^t)$$

computation graph



deep recurrent neural net



Long Short Term Memory

