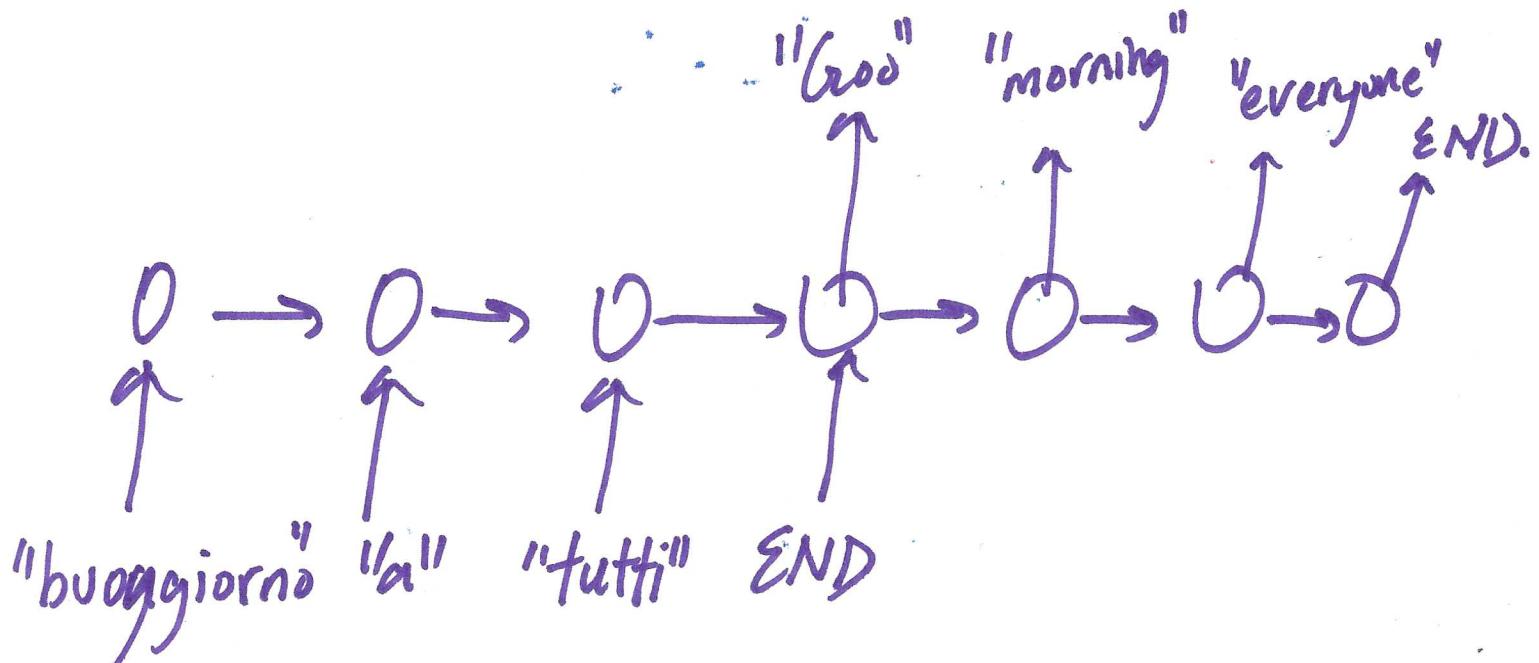


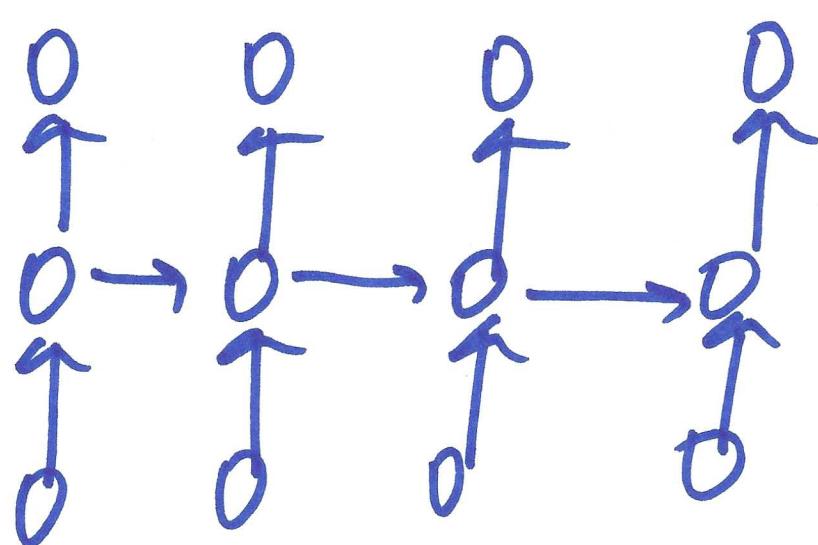
# Many-to-many architecture

6

## Example: Machine translation



## Synch. many-to-many architecture



e.g. frame-by-frame video classification.  
etc.

many other variants...

# Simple RNN architecture 7

$x_t$ : Input at time  $t \in \mathbb{R}^{d \times t}$

$s_t$ : State at time  $t \in \mathbb{R}^{N \times 1}$

$y_t$ : Output at time  $t \in \mathbb{R}^{d \times 1}$ .

$$s_t = \tanh(u x_t + \underbrace{W s_{t-1}}_{\text{recurrence}})$$

$W, U \in \mathbb{R}^{N \times d}$   
are trainable weights.

$$y_t = V s_t$$

(one can also do  
sigmoid, hard  
decisions (step func.)  
etc.)

e.g.  $y_t = \sigma(V s_t)$   
where  $\sigma$  is the  
sigmoid act.  
function.

(8)

Ex: An RNN trained to generate words over the alphabet

$\{ h, e, l, o, \square \}$



to indicate end of word.

one hot encoding we used to represent each letter.

$$\text{E.g.: } x_1 = "e", \quad x_2 = "e"$$

(First two letters of the word are "ee"  
what will come after is likely "l",  
indicating eel (a certain kind of fish))

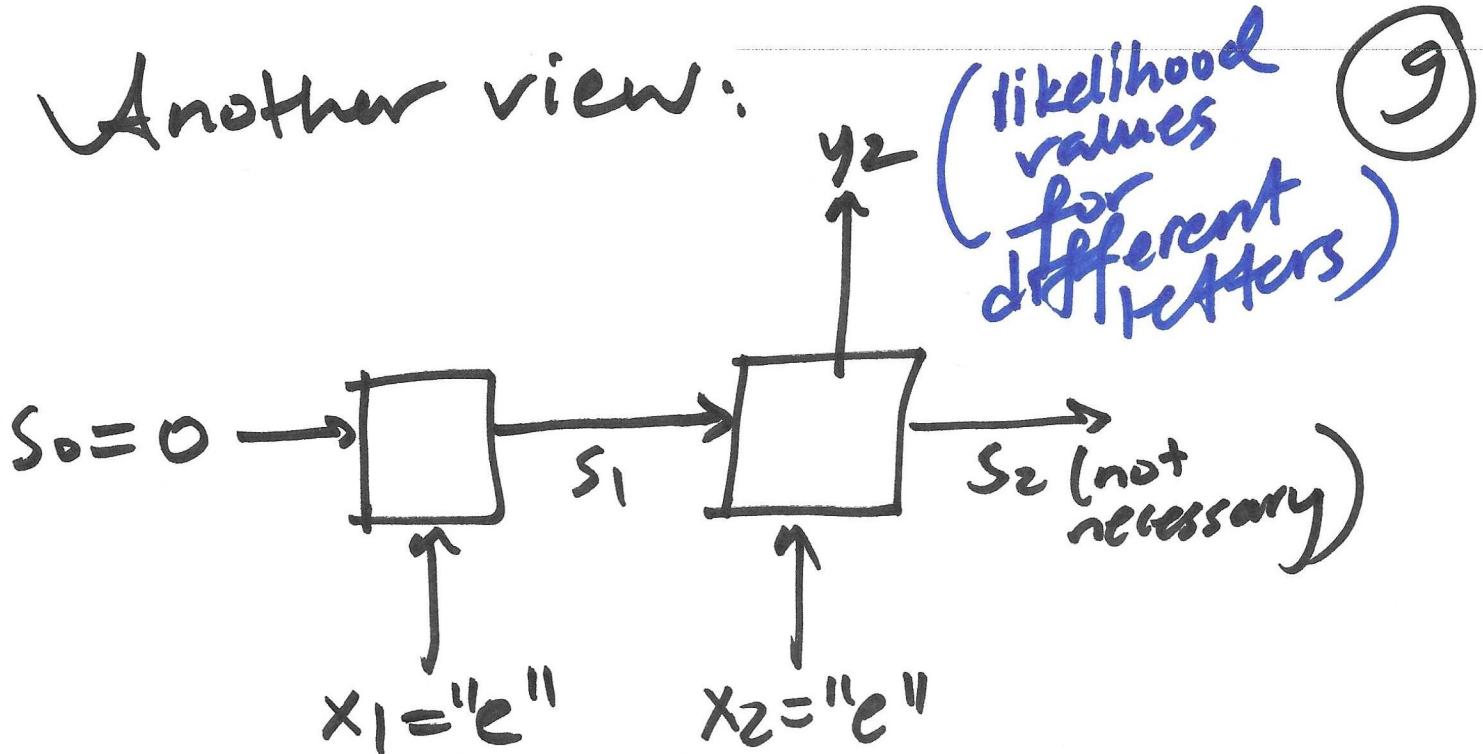
We are interested in predicting what comes after "ee", which we determine through  $y_2$ .

$$y_2 = f(Vs_2) = f(V \tanh(Ux_2 + Ws_1))$$

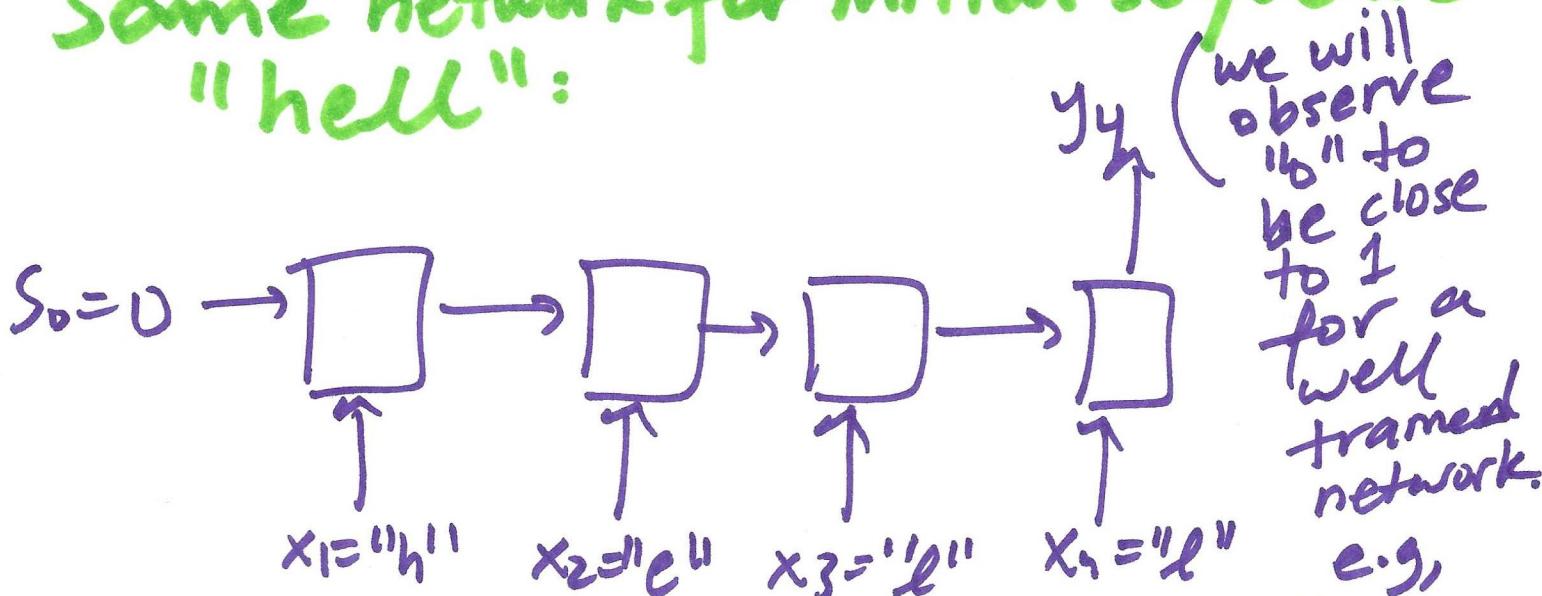
$$s_1 = \tanh(Ux_1 + Ws_0)$$

$s_0 = \text{some fixed 5-dim vector}$   
(e.g. all zero vector).

Another view:



Some network for initial sequence "hell":



Note that

$$"0" = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$$

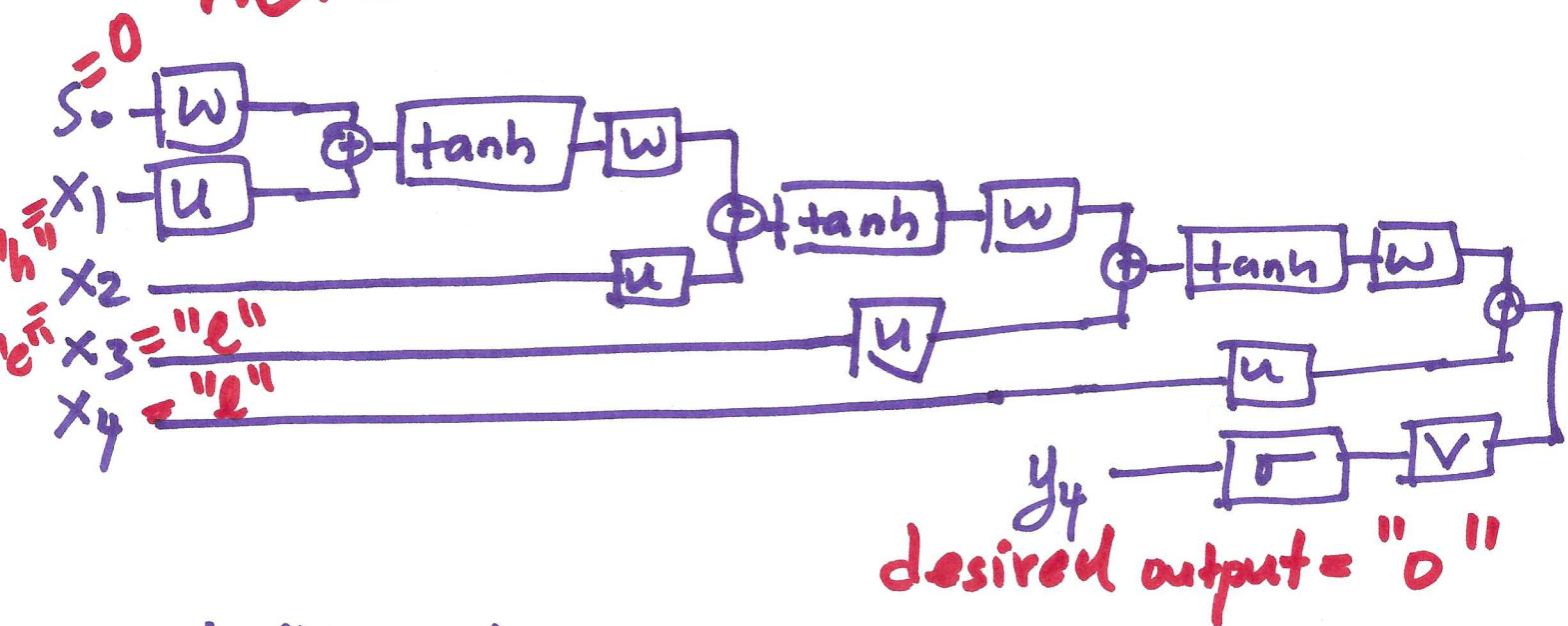
can be

$$y_4 = \begin{bmatrix} 0.05 \\ 0.01 \\ 0 \\ 0.95 \\ 0 \end{bmatrix}$$

# Training RNNs

10.

View each example above as one instance of an ordinary feedforward NN. For example, for the "hello" prediction, the network is:



Similarly for any such desired inputs - desired output pairs the RNN can be unrolled any number of times and trained using standard backpropagation.

The resulting algorithm is referred to as backpropagation through time (BPTT).

# Long-Short Term Memory

(11)

When we consider longer as longer sequences, we face the vanishing & exploding gradient problems in the standard RNN structure.

LSTMs alleviate (though not entirely solve) the vanishing & exploding  $\nabla$  problems of RNNs.

We simply replace the basic recurrence block with another block. Define

$x_t$ : Input at time  $t$

$y_t$ : Output at time  $t$

$s_t$ : State at time  $t$

As usual. LSTMs define several new "gates":

"Forget" Gate: (how much of the previous state to forget) 12

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

Input Gate: (how much of the new candidate state to keep)

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

Output Gate: (how much of the calculated output to expose)

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

New candidate state:

$$g = \tanh(U^g x_t + W^g y_{t-1} + b^g)$$

New state:  $\rightarrow$  Hadamard (element-wise) product

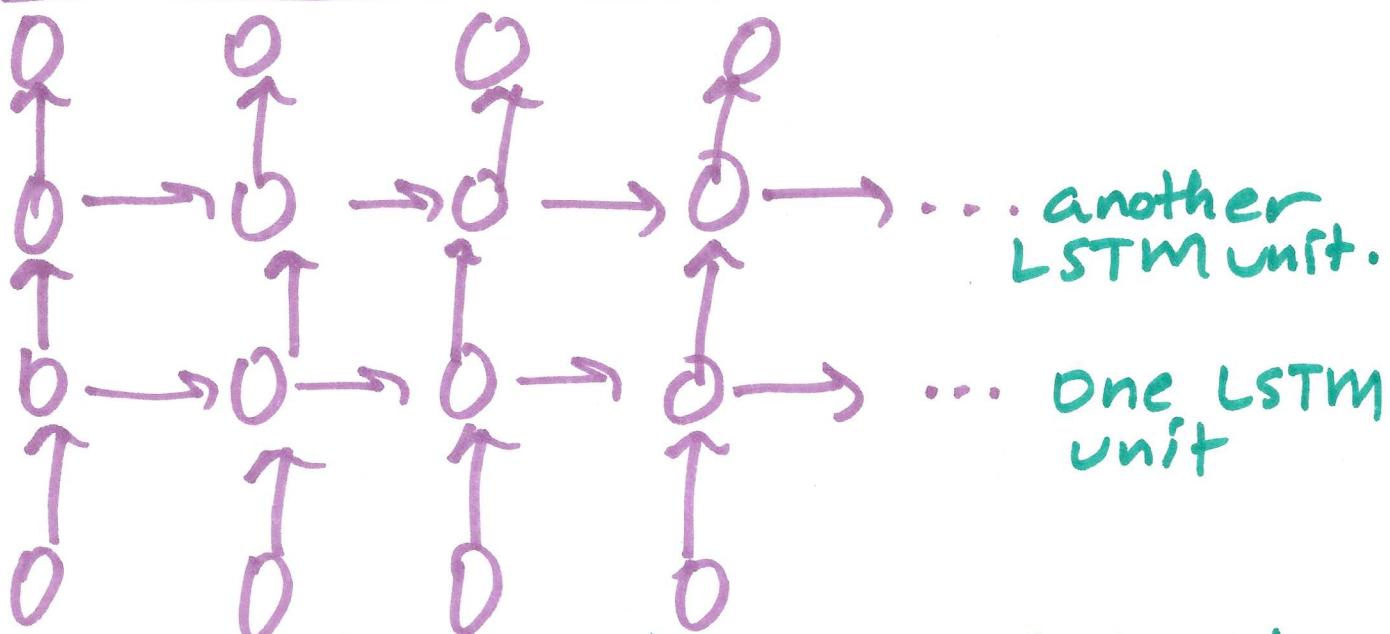
$$s_t = s_{t-1} \odot f_t + g \odot i_t$$

Output:

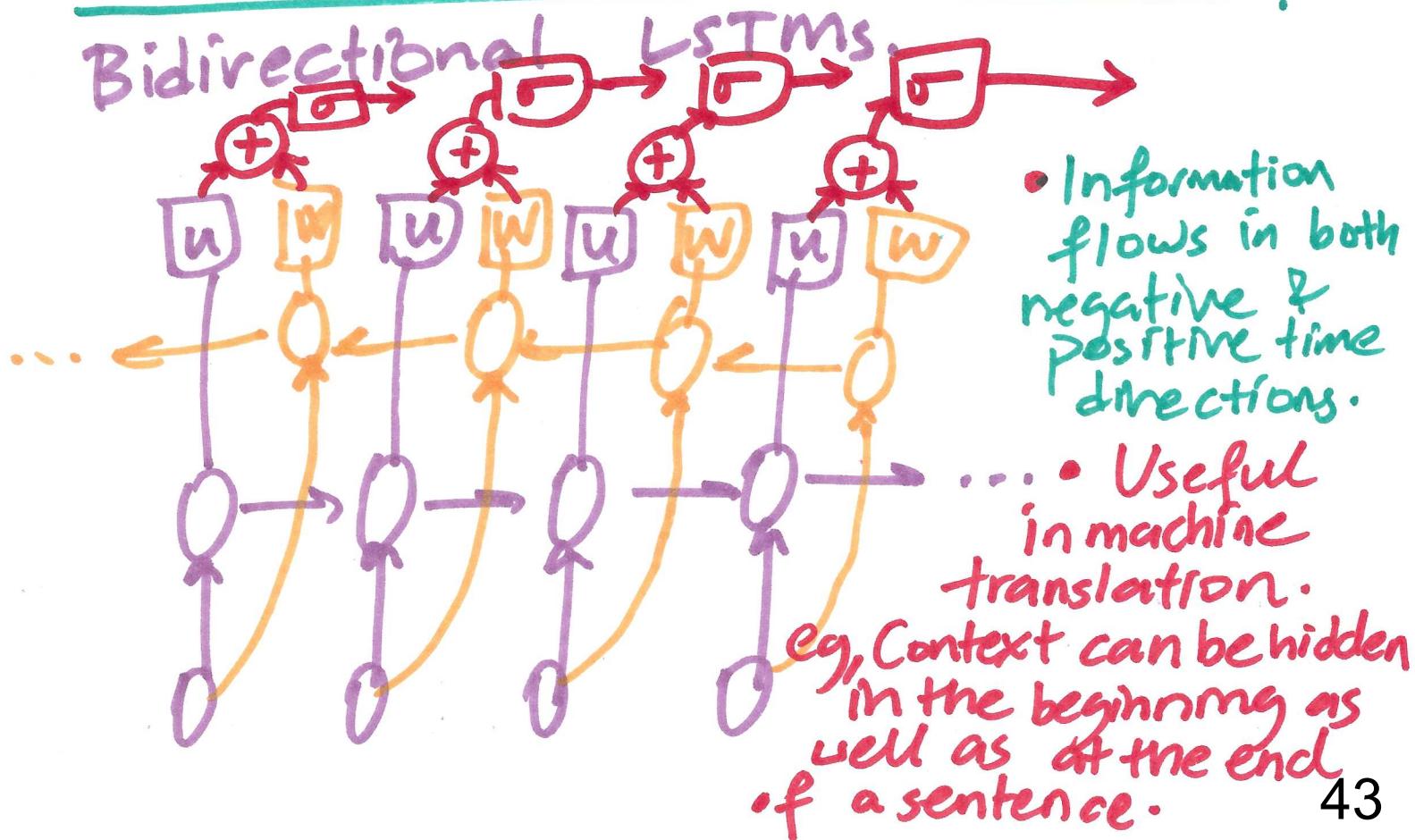
$$y_t = \tanh(s_t) \odot o_t$$

- All weights, biases are trainable.
- Peephole LSTMs include a trainable weight  $\times$  previous state to the induced local fields.

# Stacked LSTMs



- \* motivation: capturing more complex features (as in deep vs shallow NNs).
- \* Not much gains to have  $> 3$  stacks.



# L22&23

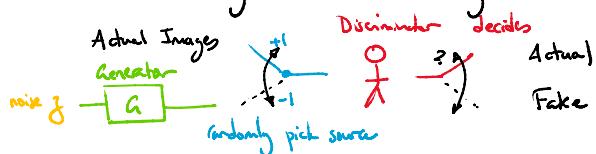
ECE/CS 559 Lecture 23 11/19/2024

Last time: Hierarchical Autoencoders

## (1) Generative-Adversarial Networks (GANs)

Motivation: In both density estimation and auto-encoders, we use likelihoods / ELBO losses.

- These are problematic in high dimensions, where the sparsity of the data makes it hard to be accurate.
- They evaluate each data point individually, instead of seeing whether the generation overall is good.



- We measure the performance of the discriminator via binary cross-entropy:  $\mathbb{E}_S [\mathbb{E}_{x \sim S} [-\log P_{\text{disc}}(S|x)]]$

- Empirically, generate as many  $z$ 's as  $x$ 's in Data, match each  $x$  to a  $z$ , then calculate:

$$R(w) = \frac{1}{|\text{Data}|} \sum_{x \in \text{Data}} -\frac{1}{2} \log \underbrace{P_{\text{disc}}(+1|x)}_{D(x)} - \frac{1}{2} \log \underbrace{P_{\text{disc}}(-1|G(z))}_{1 - D(G(z))}$$

- Goal: Discriminator tries to ↓ cross-entropy and generator tries to ↑. (get better at telling apart) (get better at fooling)

$$\arg \max_{P_{\text{gen}}} \min_{P_{\text{disc}}} R(w) \quad \text{↑ all weights (Disc & gen)}$$

## (2) Analyzing GANs Why do we think this would work?

- Let's assume infinite data + universality (can approximate all)

$$\arg \max_{P_{\text{gen}}} \min_{P_{\text{disc}}} \mathbb{E}_S [\mathbb{E}_{x \sim S} [-\log P_{\text{disc}}(S|x)]]$$

actual  $P(S|x) = \frac{P(S) \cdot P_{\text{disc}}(S)}{P(x) = \frac{1}{2} P_{\text{gen}} + \frac{1}{2} P_{\text{disc}}}$

$$\arg \max_{P_{\text{gen}}} \frac{1}{2} \mathbb{E}_{x \sim P_{\text{gen}}} [-\log P(+1|x)] + \frac{1}{2} \mathbb{E}_{x \sim P_{\text{gen}}} [-\log P(-1|x)]$$

$\frac{1}{2} P_{\text{gen}} + \frac{1}{2} P_{\text{disc}}$

$$\text{KL Divergence} \quad \text{KL}(P_{\text{gen}} \parallel \frac{1}{2}(P_{\text{gen}} + P_{\text{disc}})) + \text{KL}(P_{\text{gen}} \parallel \frac{1}{2}(P_{\text{gen}} + P_{\text{disc}}))$$

$$\arg \max_{P_{\text{gen}}} JSD(P_{\text{gen}}, P_{\text{gen}}) \Rightarrow P_{\text{gen}}^* = P_{\text{gen}}$$

Jensen-Shannon Divergence

• Intuition: If a very good **discriminator** cannot distinguish actual images from those made by a **generator**, then the **generator** is quite good. Note:

- No  $P_x$ , instead the discriminator judges images (perceptual)
- The discriminator judges the overall (not individual) performance

Model • Let  $S$  be a coin toss:  $P_S(s) = \begin{cases} \frac{1}{2} & s=1 \\ \frac{1}{2} & s=-1 \end{cases}$  (actual)

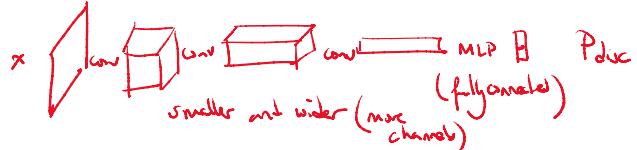
- The **discriminator** produces  $P_{\text{disc}}(s|x)$ , judging probability of actual/fake
- The **generator** samples  $x$  from  $P_{\text{gen}}(x)$ , by sampling  $z$  then  $x = G(z)$
- The discriminator sees  $x$  from data,  $P_{\text{obs}}$ , if  $s=1$ , or generator if  $s=-1$

$$P_X(x) = P_S(+1) \cdot P_{\text{obs}}(x) + P_S(-1) \cdot P_{\text{gen}}(x)$$

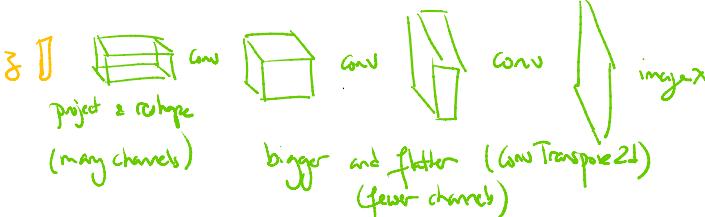
$P_{\text{obs}}(x|s=1) \quad P_{\text{obs}}(x|s=-1)$

- What do the discriminator and generator look like.

- Discriminator: CNNs, standard



- Generator: CNNs, in reverse



## (3) Conditional Generation

- All a form of translation, e.g., image-to-image style change  $x \rightarrow \begin{cases} \text{dog} \end{cases} \rightarrow \begin{cases} \text{cat} \end{cases} x$

- (a) With paired images  $(x, x')$  given  $x'$  generate  $x$

$$\arg \min_{\text{gen}} \mathbb{E}_{\text{gen}} [\|x - G(z, x')\|_1] + \lambda \max_{\text{disc}} \mathbb{E}_{\text{disc}} [-\log P_{\text{disc}}(S|x)]$$

- (b) With unpaired images  $\{x'\} \quad \{x\} \quad R_{\text{disc}}$

Cycle GAN: 2 generators & 2 discriminators (for  $x \neq x'$ )

$$\mathbb{E}[\|G_1(G_1(z_1, x')) - x'\|_1] \quad (\text{cycle loss}) + \lambda R_{\text{disc}} \\ + \mathbb{E}[\|G_1(G_2(z_2, x)) - x\|_1] \quad (\text{loss}) + \lambda R_{\text{disc}}$$

# L19&20

## ECE/CS 559 Lecture 19/20 Th 10/31 ; Th 11/7

Last time: Hebbian updates

(1) Generative Models

- All unsupervised learning is about understanding  $P(x)$
- Generative models attempt to create one  $x$ :
  - They can do this by modeling  $P(x)$  & sampling from it.
  - Or, they can directly model the sampling mechanism  
e.g., sample  $z$  (standard distribution), then let  $x = f(z) + \text{noise}$
  - e.g.,  $y = \text{caption}$   $x = \text{image}$  (diffusion models)
  - Conditional generative models attend to generation via prompts:  $P(x|y)$
  - Examples:  $y = \text{caption}$   $x = \text{image}$  (diffusion models)  
 $y = \text{question}$   $x = \text{answer}$  (transformers)

(2) Density Estimation

Given Data  $x$ , assume i.i.d. from  $P$ .

Parametric methods: Let  $P(x) = f(x; w)$

Maximize likelihood:  $\underset{w}{\operatorname{argmax}} P(\text{Data}) = \prod_{x \in \text{Data}} P(x)$

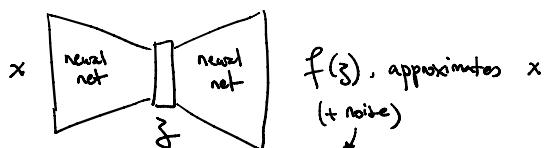
$$\Leftrightarrow \underset{w}{\operatorname{argmin}} \sum_{x \in \text{Data}} -\log P(x)$$

Non-parametric methods: Allow # of parameters  $w$  to grow with data size. Example, kernel density estimation  
Data  $\rightarrow$  density  $\rightarrow$  estimated density

Works well when  $x$  discrete or small dimensional.  
Otherwise, requires regularization, e.g. smooth densities.

(3) Autoencoders (Variational; VAEs)

- They try to capture  $P(x)$  by reducing the "essence" of  $x$  to a simple representation  $z$  (encoding) then generating  $x$  from a sample of  $z$  (decoding)



Globally:  $P_{enc}(z|x)$   $P_{prior}(z)$   $P_{dec}(x|z)$

- Why sampling  $z \sim P_{prior}(z)$  then letting  $x = f(z)$  work?  
How do we choose  $P_{enc}, P_{dec}$ ?

(4) Derived distributions:

- Say  $z, x \in \mathbb{R}^d$ . Let  $z \sim P(z)$  and let  $x = f(z)$ .
- What is  $P(x)$ ? Let's assume  $P(z)$  are densities,  $f$  monotonic
- $P(x) = P(X \leq x) = P(f(z) \leq x) = P(z \leq f^{-1}(x))$
- thus  $P(x) = \frac{dP}{dx} = \frac{d}{dx} F_z(f^{-1}(x)) = \frac{d}{dz} F_z(z) \cdot \frac{dz}{dx} = P(f(z)) \cdot \frac{df^{-1}}{dx}$

Example:  $P_z \sim \text{Uniform}[0, 1]$   $x = -\log z$   $f^{-1}(x) = e^{-x}$   
 $\Rightarrow P_x(x) = \begin{cases} 1 \cdot (-e^{-x}) & e^{-x} \in (0, 1) \\ 0 & \text{otherwise} \end{cases} = \begin{cases} e^{-x} & x < 0 \\ 0 & x \geq 0 \end{cases}$

- Can create any new random variable from standard random variables (e.g. uniform, Gaussian, etc.)

(b) From maximum likelihood to ELBO

How do we train an autoencoder?

Idea 1: Fix  $P_{prior}(z)$  and only train the decoder  $P_{dec}(x|z)$

$$\text{by minimizing } \sum_{x \in \text{Data}, z \sim P_{prior}} -\log P_{dec}(x|z)$$

Issue: Creates disconnect between  $x$  &  $z$  (close  $x$ 's  $\nRightarrow$  close  $z$ 's)

Idea 2: Fix  $x$   $P_{enc}(z|x)$  and train  $P_{prior}(z) \circ P_{dec}(x|z)$

$$\text{by minimizing } \sum_{x \in \text{Data}, z \sim P_{enc}(z|x)} -\log P_{prior}(z) \cdot P_{dec}(x|z)$$

$$\text{approximates } \mathbb{E}_x \mathbb{E}_{z \sim P_{enc}(z|x)} -\log P(z) P_{dec}(x|z) \Leftrightarrow$$

implies  $P(x) P(z|x)$

$$\mathbb{E}_x \left[ \log \frac{1}{P_{gen}(x)} \right] + \mathbb{E}_x \left[ \mathbb{E}_{z \sim P_{enc}(z|x)} \left[ \log \frac{1}{P_{gen}(f(z))} \right] \right] - \mathbb{E}_x \left[ \mathbb{E}_z \left[ \log \frac{1}{P_{enc}(z|x)} \right] \right]$$

min at  $P_{gen}(f(z)) = P_{enc}(z|x)$       min at  $P_{gen}(f(z)) = P_{enc}(z|x)$       - Entropy ( $P_{enc}$ )

Issue: If the neural net of the decoder cannot easily "invert" the choice of  $P_{enc}$ , it won't work well.

Idea 3: Train  $P_{enc}$  and  $P_{dec}$  jointly.

Issue:  $P_{enc}$  in the loss is just a sampling distribution.

Hard to optimize: bad local minima /  $P_{enc}$  can degenerate by following w/  $P_{dec}$

Idea 4: Include  $P_{enc}$  in the loss:  $\mathbb{E} \left[ -\log \frac{P_{prior}(z) P_{dec}(x|z)}{P_{enc}(z|x)} \right]$

"max-entropy" effect

Evidence Lower Bound (ELBO) Loss  $< \log P(x)$

③ Example: Gaussian Variational Autoencoder

$$P_{\text{prior}}(\mathbf{z}) \propto \exp\left(-\frac{\|\mathbf{z}\|^2}{2}\right) \quad \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$P_{\text{encoder}}(\mathbf{z}|\mathbf{x}) \propto \exp\left(-\frac{\|\mathbf{z} - g(\mathbf{x}; \mathbf{w})\|^2}{2}\right) \quad \mathcal{N}(g(\mathbf{x}; \mathbf{w}), \mathbf{I})$$

$$P_{\text{decoder}}(\mathbf{x}|\mathbf{z}) \propto \exp\left(-\frac{\|\mathbf{x} - f(\mathbf{z}; \mathbf{w})\|^2}{2}\right) \quad \mathcal{N}(f(\mathbf{z}; \mathbf{w}), \mathbf{I})$$

$$\text{ELBO has explicit form: } \mathbb{E}\left[\frac{\|\mathbf{z}\|^2}{2} - \frac{\|\mathbf{z} - g(\mathbf{x}; \mathbf{w})\|^2}{2}\right] = \frac{1}{2}\mathbb{E}\left[\|g(\mathbf{x}; \mathbf{w})\|^2\right]$$

(minimize representation norm)

$$+ \frac{1}{2}\mathbb{E}\left[\|\mathbf{x} - \underbrace{f(g(\mathbf{x}; \mathbf{w}) + \mathbf{N}; \mathbf{w})}_{\text{encode}}\|_2^2\right] \quad \mathbf{N} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

(minimize reconstruction error)

This looks almost like k-means and it's not a coincidence! The similarity is closer if we quantize  $\mathbf{z}$ .

④ Hierarchical VAEs

- Instead of one representation stage, have multiple:

$$\mathbf{x} \xrightarrow{\text{encode}} \mathbf{z}_1 \xrightarrow{\text{encode}} \mathbf{z}_2 \xrightarrow{\text{encode}} \mathbf{z}_3 \xrightarrow{\text{decode}} \mathbf{z}_2 \xrightarrow{\text{decode}} \mathbf{z}_1 \xrightarrow{\text{decode}} \mathbf{x}$$

- Why? Allows for gradual tweaks, making it easier to train.

Example: Diffusion models use encoders that just add noise.

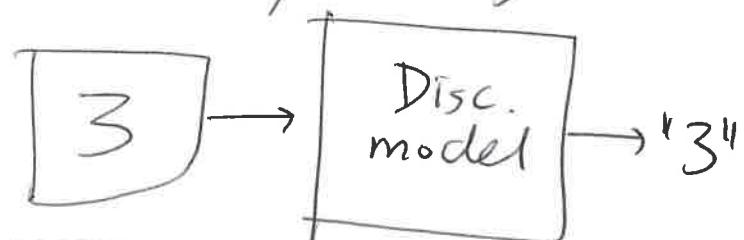
$$\text{ELBO: } \log \frac{P_{\text{prior}}(\mathbf{z}_1) P_{\text{dec}}(\mathbf{z}_2|\mathbf{z}_1) P_{\text{dec}}(\mathbf{z}_3|\mathbf{z}_2) P_{\text{dec}}(\mathbf{x}|\mathbf{z}_1)}{P_{\text{enc}}(\mathbf{z}_3|\mathbf{z}_2) P_{\text{enc}}(\mathbf{z}_2|\mathbf{z}_1) P_{\text{enc}}(\mathbf{z}_1|\mathbf{x})}$$

# ECE/CS 559 - Neural Networks.

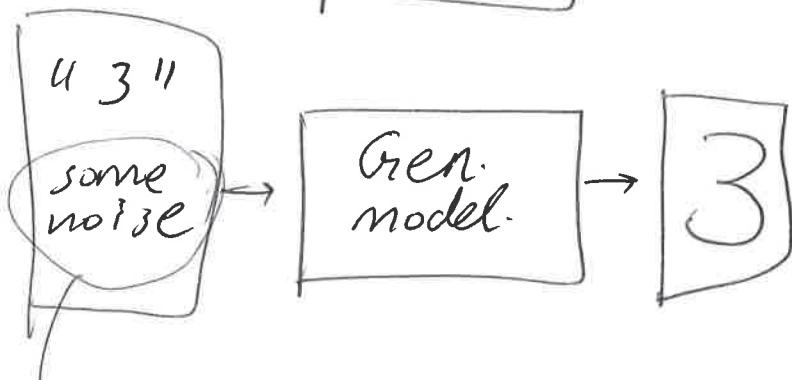
## Generative Models. ①

- \* Two main approaches in machine learning:  
generative vs. discriminative approach.
- \* Observation:  $X$ , Target:  $Y$   
(e.g. an input image) (e.g. a class label).
- \* A discriminative model can provide  $P(Y|X=x)$   
(e.g. given the input, what are the likelihoods  
of different class labels).
- \* Whereas in a generative model, we are  
interested in finding  $P(X|Y=y)$   
Given label, generate some samples ~~some~~  
belonging to that label (kind of an inverse  
problem):

Discriminative  
model :



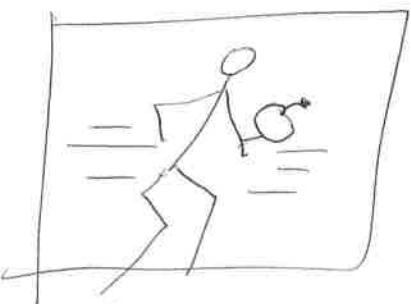
Generative  
model



→ Noise could be considered  
as a "seed". Different noises  
would generate different images of 3.

## Applications of generative models

- \* Low ~~to~~ high resolution image synthesis.  
Text to image translation  
"A guy running with an apple on his hand".



- \* Speech synthesis.
- \* Error correction

: