

## Lecture 24

[Generating images by inverting  
CNNs]

A trained CNN yields  $p(y|\tilde{x})$

$\uparrow$   $\uparrow$   
 class label image

This can be used to produce a  
generative model:

$$p(\tilde{x}|y) \sim p(\tilde{x}) \underbrace{p(y|\tilde{x})}_{\text{prior for images}}$$

Need to sample the energy function:

$$E_c(\tilde{x}) = \log p(y=c|\tilde{x}) + \log p(\tilde{x}).$$

Can use Langevin dynamics:

$$(1) \quad \tilde{x}_{t+1} = \tilde{x}_t + \frac{\epsilon}{2} \underbrace{\nabla E_c(\tilde{x}_t)}_{-\text{force}} + \mathcal{N}(\vec{0}, \epsilon' \mathbb{I})$$

or even

$$(2) \quad \tilde{x}_{t+1} = \tilde{x}_t + \epsilon_1 \frac{\partial \log p(\tilde{x}_t)}{\partial \tilde{x}_t} + \epsilon_2 \underbrace{\frac{\partial \log p(y=c|\tilde{x}_t)}{\partial \tilde{x}_t}}_{\text{Jacobian}} + \mathcal{N}(\vec{0}, \epsilon_3 \mathbb{I})$$

If  $\epsilon_3 = 0$ , the method generates  
the 'most likely' image for class c.

## Image priors

Gaussian prior:  $p(\bar{x}) = \mathcal{N}(\bar{x} | \bar{0}, \mathbb{I})$

↑  
for centered  
image pixels

$$\nabla_{\bar{x}} \log p(\bar{x}) = \nabla_{\bar{x}} \left[ -\frac{|\bar{x}|^2}{2} \right] = -\bar{x}.$$

Using  $\epsilon_2 = 1$ ,  $\epsilon_3 = 0$ , we obtain:  
(Eq. (2))

$$\bar{x}_{t+1} = (1 - \epsilon_1) \bar{x}_t + \frac{\partial \log p(y=c | \bar{x}_t)}{\partial \bar{x}_t}$$

Total variation (TV) prior:

$$TV(\bar{x}) = \sum_k \sum_{ij} \left[ (x_{ijk} - x_{i+1,j,k})^2 + (x_{ijk} - x_{i,j+1,k})^2 \right]$$

channel row  $i$ ,  
column  $j$

One can use  $p(\bar{x}) \sim e^{-TV(\bar{x})}$

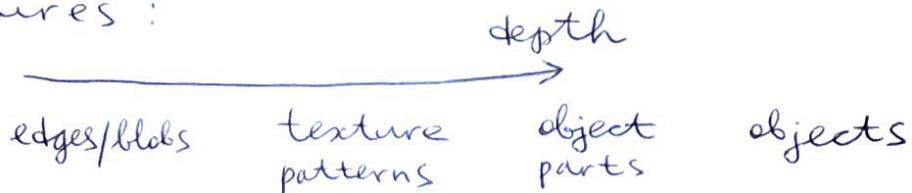
as the prior; penalizes high-freq.  
artifacts

Visualizing CNN features

Start with random image and optimize it to maximize the activation of a particular neuron  $\Rightarrow$  activation maximization

Goal: discover neuron functions

When applied to the AlexNet CNN, this procedure identifies a hierarchy of features:



Alternatively, one can simply search for images in a DB that maximally activate a given neuron  $\Rightarrow$  yields more 'relatable' images.

## Deep Dream

Goal: generate versions of an input image that emphasize/exaggerate certain features

Patterns of activity of neurons in a given layer = different features in the image

Then define  $\langle y_e \rangle = \frac{1}{HWC} \sum_{h,w,c} y_{e,hwc} (\vec{x})$

Use  $\sum_{\ell \in L} \langle f_\ell(\bar{x}) \rangle$  as the loss function,

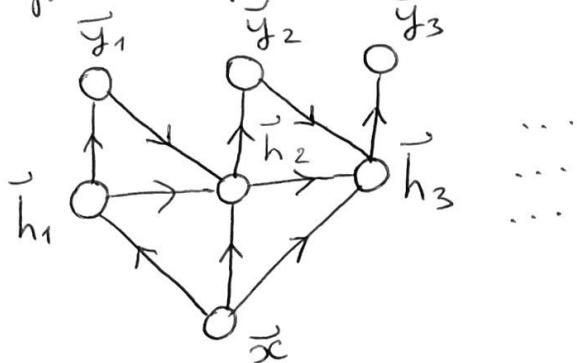
to be optimized via grad descent

This procedure results in a combination of the original image with "hallucinations".

# [ Recurrent neural networks (RNNs) ]

vec2seq. models:  $R^D \rightarrow R^{\frac{T}{C}}$   $\xrightarrow{C}$   $T$  vectors of size  $C$ ;  $T$  determined dynamically  
fixed size of the input vector

Output sequence:  $\tilde{y}_{1:T} = \{\tilde{y}_1, \dots, \tilde{y}_T\}$



computes  
 $p(\tilde{y}_{1:T} | \tilde{x})$  via  
hidden states  $\tilde{h}_{1:T}$

$$p(\tilde{y}_{1:T} | \tilde{x}) = \sum_{\tilde{h}_{1:T}} p(\tilde{y}_{1:T}, \tilde{h}_{1:T} | \tilde{x}) =$$

$$= \sum_{\tilde{h}_{1:T}} \prod_{t=1}^T p(\tilde{y}_t | \tilde{h}_t) p(\tilde{h}_t | \tilde{h}_{t-1}, \tilde{y}_{t-1}, \tilde{x})$$

Init. cond.:  $p(\tilde{h}_1 | \tilde{h}_0, \tilde{y}_0, \tilde{x}) = f(\tilde{h}_1 | \tilde{x})$   
initial hidden state distribution

Output distributions:

$$p(\tilde{y}_t | \tilde{h}_t) = \text{Cat}(\tilde{y}_t | \text{softmax}(\underbrace{W_{hy} \tilde{h}_t + \tilde{b}_y}_{\text{hidden} \rightarrow \text{output weights}}))$$

hidden  $\rightarrow$  output weights      output biases

or

$$p(\tilde{y}_t | \tilde{h}_t) = \mathcal{N}(\tilde{y}_t | W_{hy} \tilde{h}_t + \tilde{b}_y, \sigma^2 \mathbb{I})$$

$p(\tilde{h}_t | \tilde{h}_{t-1}, \tilde{y}_{t-1}, \tilde{x})$  is computed deterministically:

$$\tilde{h}_t = f(\tilde{h}_{t-1}, \tilde{y}_{t-1}, \tilde{x}) \text{ or,}$$

more explicitly,

$$\tilde{h}_t = g(w_1 \tilde{x} + w_2 \tilde{y}_{t-1} + w_3 \tilde{h}_{t-1} + b_h)$$

→ RNN can be viewed as a generalization of a standard Markov model

→ Training RNN: use labeled data as usual

→ Using RNN in a generative mode:

sample from  $p(\tilde{y}_t | \tilde{h}_t)$  to generate

$\tilde{y}_t$ , compute  $\tilde{h}_{t+1} = f(\tilde{h}_t, \tilde{y}_t, \tilde{x})$ ,

sample from  $p(\tilde{y}_{t+1} | \tilde{h}_{t+1})$ , etc.

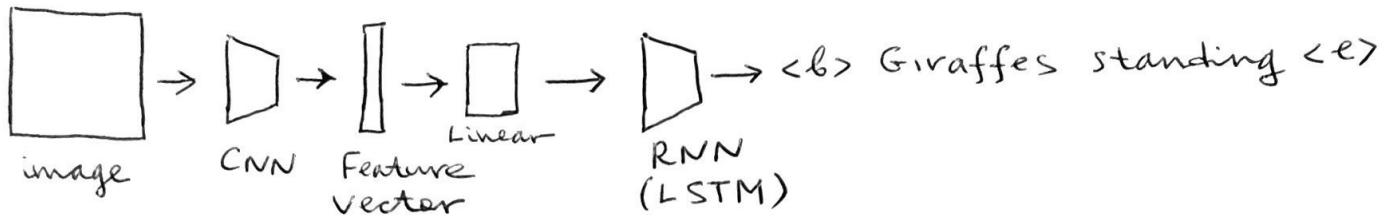
[Noise comes from the stochasticity in the outputs.]

Applications:

- ① Language modeling:  $\tilde{x} = \tilde{o}$ , learn  
 $p(y_1, \dots, y_T) = \text{sentence}$   
↑  
tokens

- ② Image captioning:  $\tilde{x} = \text{image embedding}$  produced by CNN

$p(y_1, \dots, y_T) = \text{image caption}$

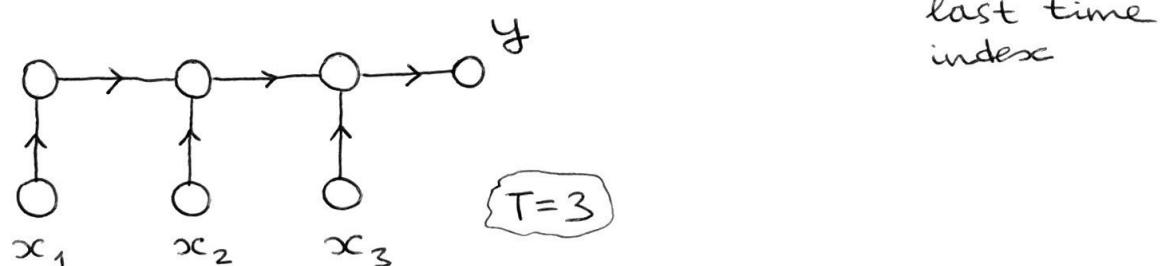


Seq2vec models:  $R^{\text{TD}} \xrightarrow{\downarrow} R^C$   
 variable-size sequence of vectors → fixed-size output

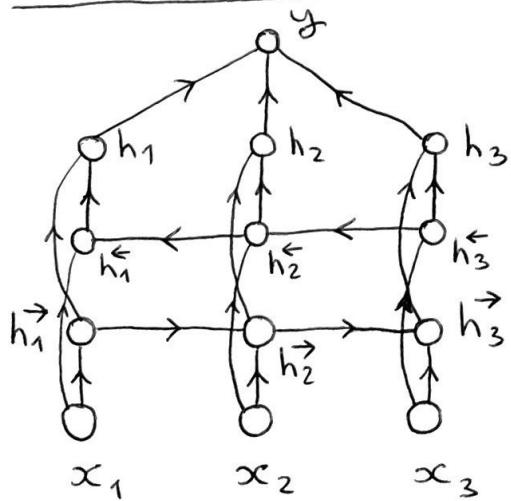
Focus on the class-label output:

$$y \in \{1, \dots, C\}$$

$$p(y | \tilde{x}_{1:T}) = \text{Cat}(y | \text{softmax}(W_{hy} \tilde{h}_T + \tilde{b}_y))$$



Bi-directional RNN:



$$\tilde{h}_t^> = \varphi(W_{xh}^> \tilde{x}_t + W_{hh}^> \tilde{h}_{t-1} + \tilde{b}_h^>),$$

$$\tilde{h}_t^< = \varphi(W_{xh}^< \tilde{x}_t + W_{hh}^< \tilde{h}_{t-1} + \tilde{b}_h^<),$$

$$\tilde{h}_t = [\tilde{h}_t^>, \tilde{h}_t^<]$$

$$\langle \tilde{h} \rangle = \underbrace{\frac{1}{T} \sum_{t=1}^T \tilde{h}_t}_{\text{concatenation}} \quad \text{average}$$

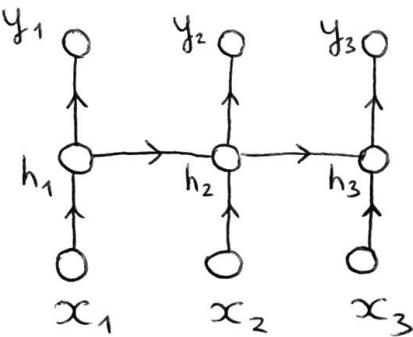
Finally,  $p(y|\tilde{x}_{1:T}) = \text{Cat}(y|\text{softmax}(w^T \tilde{h} + \tilde{b}_y))$

Applications: text to sentiment

Seq2Seq models:  $R^{TD} \rightarrow R^{TC}$

E.g.,  $T=T'$   $\Rightarrow$  aligned input/output  
 $T \neq T'$   $\Rightarrow$  unaligned  $-||- / -||-$

① aligned case: dense sequence labeling  
 (one label per location)

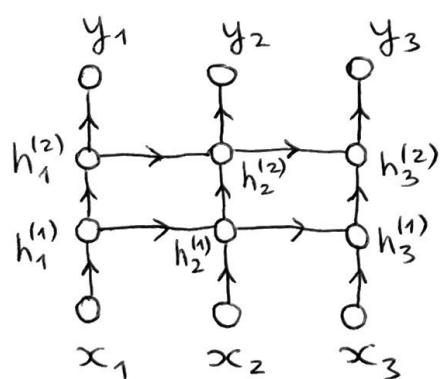


uni-directional  
 RNN (bi-directional RNN  
 possible as well)

$$p(\tilde{y}_{1:T} | \tilde{x}_{1:T}) = \sum_{\tilde{h}_{1:T}} \prod_{t=1}^T p(\tilde{y}_t | \tilde{h}_t) \\ f(\tilde{h}_{t-1}, \tilde{x}_t)$$

$$\tilde{h}_1 = f_o(\tilde{x}_1) \quad \text{initial condition}$$

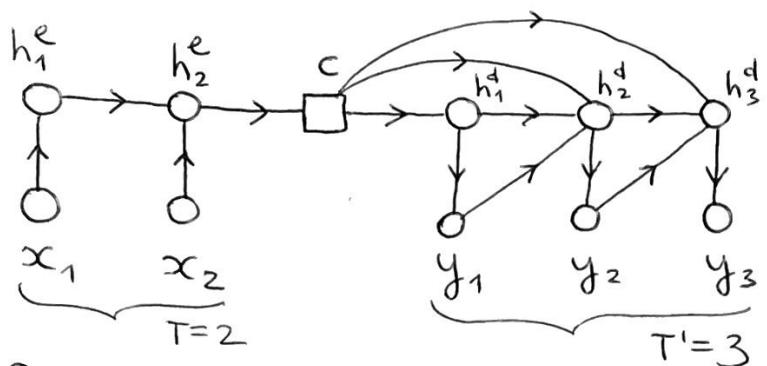
Deep RNN: multiple layers of hidden nodes



$$\tilde{h}_t^{(l)} = g_l(W_{xh}^l \tilde{h}_t^{(l-1)} + W_{hh}^l \tilde{h}_{t-1}^{(l)} + \tilde{b}_h^{(l)})$$

] layers      ↑  
 hidden layer updates

Unaligned case: encoder-decoder architecture



Context vector  $\tilde{c} = f_e(\tilde{x}_{1:T})$

RNN decoder  $\tilde{y}_{1:T} = f_d(\tilde{c})$

