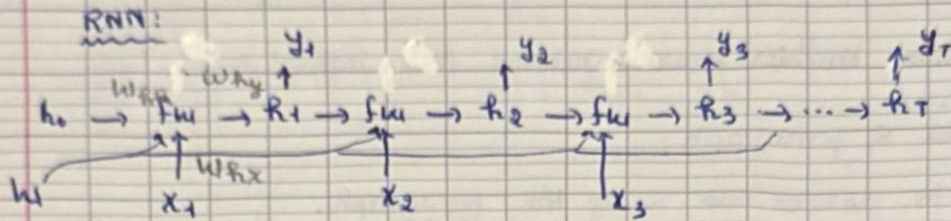


$$\frac{e^3 + e^{-3}}{e^3 - e^{-3}} = \frac{e^3 + \frac{1}{e^3}}{e^3 - \frac{1}{e^3}} = \frac{\frac{e^{23} + 1}{e^3}}{\frac{e^{23} - 1}{e^3}} = \frac{e^{23} + 1}{e^{23} - 1}$$

HW 5: Image Captioning.



$$\begin{aligned} h_t &= f_w(h_{t-1}, x_t) \\ &= \theta(w_{hh} h_{t-1} + w_{hx} x_t) \\ &= \theta\left([w_{hh} \quad w_{hx}] \begin{bmatrix} h_{t-1} \\ x_t \end{bmatrix} + b_h\right) \end{aligned}$$

Many to Many:

$$y_t = \theta(w_{hy} h_t + b_y)$$

ENCODER: Encode input sequence into single vector;

MANY to ONE; $y = w_{hy} h_T + b_y$

LSTM:

$$\begin{bmatrix} i \\ f \\ o \\ g \end{bmatrix} = \begin{bmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{bmatrix} w \begin{bmatrix} h_{t-1} \\ x_t \end{bmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

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

Image captioning:

$$h_t = \tanh(w_{hx} x_t + w_{hh} h_{t-1} + b_h)$$

$$y_t = w_{hy} h_t + b_y$$

$$L = \frac{1}{2} \sum_t (y_t - \hat{y}_t)^2$$

EX:

$$T=2; x = [1, 2]; w_{hx} = \begin{bmatrix} 0.5 \\ -0.3 \end{bmatrix}; w_{hh} = \begin{bmatrix} 0.1 & 0.2 \\ 0.4 & 0.5 \end{bmatrix}$$

$$w_{hy} = [0.6 \quad 0.4]; b_h = \begin{bmatrix} 0 \\ 0 \end{bmatrix}; b_y = 0$$

Find pass: $t=1$:

$$h_1 = \tanh\left(\begin{bmatrix} 0.5 \\ -0.3 \end{bmatrix} \cdot 1 + \begin{bmatrix} 0 \\ 0 \end{bmatrix}\right) = \begin{pmatrix} 0.462 \\ -0.291 \end{pmatrix}$$

$$y_1 = [0.6 \quad 0.4] \cdot h_1 = 0.277$$

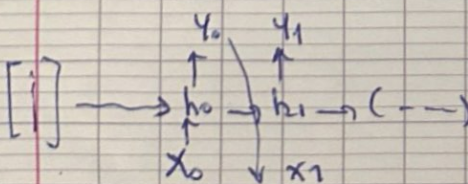
$$t=2: h_2 = \tanh\left(\begin{bmatrix} 0.5 \\ -0.3 \end{bmatrix} \cdot 2 + \begin{bmatrix} 0.1 & 0.2 \\ 0.4 & 0.5 \end{bmatrix} \begin{pmatrix} 0.462 \\ -0.291 \end{pmatrix}\right) = \begin{pmatrix} 0.744 \\ -0.545 \end{pmatrix}$$

Encoder: CNN to transform image into vector.

Decoder: RNN.

$$\text{Image} \rightarrow \begin{bmatrix} \end{bmatrix}$$

$$x_0 = \langle \text{START} \rangle$$



$$\frac{\partial L}{\partial w_{hy}} = \sum_t \frac{\partial L}{\partial y_t} \cdot \frac{\partial y_t}{\partial w_{hy}}$$

$$\frac{\partial L}{\partial y_2} = y_2 - \hat{y}_2 = -1.754$$

$$\frac{\partial L}{\partial y_2} = h_2$$

SOLO

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$\frac{d \tanh}{dx} =$$

HW 5:

train 2014 - Vgg16 - pca } Extracted Features & PCA applied.
val 2014 - Vgg16 - pca

caption - encoding. Integer ID \rightarrow word.

model - layers.py: contains all fundamental rnn layer implem.

mn - model.py: bring layers together, & create Image Captioning.

forward - rnn - step: process incoming data & updates internal state memory.

① Forward layer: feed pass through linear dense.

(1) inp: Input data
weights matrix.
Bias Vector.

(2) output = inp * weights + bias.

$$\begin{bmatrix} \cdot & \cdot \\ \cdot & \cdot \end{bmatrix} \cdot \begin{bmatrix} w_1 & w_2 \\ w_3 & w_4 \end{bmatrix}$$

② Backward layer:

inp, weights, bias = cache_data

$\partial \text{inp} = \nabla \cdot W^T$ = Propagates error backward

$$\frac{\partial L}{\partial W_{hx}} = \frac{\partial L}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial h_t} \cdot \frac{\partial h_t}{\partial W_{hx}}$$

$$\begin{cases} z_{hx} = W_{hx} x_t + b_h \\ z_{hh} = W_{hh} h_{t-1} \end{cases} \quad \begin{cases} z = W_{hx} x_t + W_{hh} h_{t-1} + b_h \\ h_t = \tanh(z) \end{cases}$$

Given:

d_inp = grad. weights.
d_weights = inp^T grad
d_bias = Σ grad.

$$\frac{\partial h_t}{\partial z} = 1 - \tanh^2(z) \cdot \frac{\partial L}{\partial z} = \frac{\partial L}{\partial h_t} \times \frac{\partial h_t}{\partial z}$$

we're given $\frac{\partial L}{\partial h_t}$

$$\begin{aligned} z &= W_{hx} x_t + W_{hh} h_{t-1} + b_h \\ h_t &= \tanh(z) \end{aligned}$$

$$\text{I have } \frac{\partial L}{\partial h_t} \quad \frac{\partial L}{\partial z} = \frac{\partial L}{\partial h_t} \times \frac{\partial h_t}{\partial z} = d_{\text{next-state}} \times (1 - \tanh^2(z))$$

$$\frac{\partial L}{\partial W_{hx}} = \frac{\partial L}{\partial z} \times \frac{\partial z}{\partial W_{hx}} = d^g \times x_t^T$$

SOLO

$$z = W_{hx} x_t + W_{hh} h_{t-1} + b_h$$

$$\frac{\partial L}{\partial W_{hx}} = \frac{\partial L}{\partial z} \times \frac{\partial z}{\partial W_{hx}} = dz \times (x_t^T) = 1$$

$$\frac{\partial L}{\partial x_t} = \frac{\partial L}{\partial z} \times \frac{\partial z}{\partial x_t} = dz \times W_{hx}^T$$

$$\frac{\partial L}{\partial b_h} = dz \times 1$$

$$\frac{\partial L}{\partial W_{hh}} = \frac{\partial L}{\partial z} \times \frac{\partial z}{\partial W_{hh}} = dz \times h_{t-1}$$

$$\frac{\partial L}{\partial h_{t-1}} = dz \times W_{hh}$$

From Example

Final steps

inp is: $\begin{bmatrix} - & - & - & - & - \\ - & - & - & - & - \\ - & - & - & - & - \end{bmatrix}$

$\Rightarrow 3 \times 10$ matrix

• Row \rightarrow one embedding.

• Col \rightarrow size/features of word.

prev-hidden: $\begin{bmatrix} [] \\ [] \\ [] \end{bmatrix}$

$\Rightarrow 3 \times 4$ matrix;

word weights: $\begin{bmatrix} [] \\ [] \\ [] \\ [] \end{bmatrix}$

vocab size:

$10 \times 4 \rightarrow 4$ hidden dim.

memory-weights: $\begin{bmatrix} [] \\ [] \\ [] \\ [] \end{bmatrix}$

$$z_{hx} = W_{hx} \cdot x + b_h$$

$$z_{hh} = W_{hh} \cdot h_{t-1}$$

$$z = z_{hx} + z_{hh}$$

$$h_t = \tanh(z)$$

time step here represents no. of words in a caption.

SOLO

word sequence:

word 1 2 samples
word 2 3 words
word 3 4 emb. dim

Init. cont:

word 1
word 2
word 3

2x5

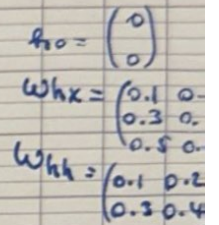
1 capture } [] } → word 1
 [] } → word 2
 [] } → word 3
[] }
[] }
 2 x 5

$\cdot h_i = \begin{bmatrix} h_{i1} \\ h_{i2} \\ h_{i3} \end{bmatrix}$; Initial context
= initial hidden state

$\cdot h_i = \begin{bmatrix} h_{i1} \\ h_{i2} \\ h_{i3} \end{bmatrix}$; Initial context
= initial hidden state

$$dh_3 \approx$$

$\frac{1}{2} \times \frac{1}{2} = \frac{1}{4}$

$$-h_0 \rightarrow f_w \rightarrow$$


[- - -]
[- - -]
[- - -]

$$x_t = x_t w_{hx} + h_{t-1} w_{hh} + b$$

$\varphi_0 = (a \ b); \varphi_1 = (c \ d)$ hidden states
 $\begin{pmatrix} x & x \\ x & x \end{pmatrix}$ er
 $\begin{pmatrix} - & - \end{pmatrix}$ feature vector

$$dh_{next} = dh_{states}[1] = (- \quad -).$$
$$d_3 = 1 - \frac{1}{2}$$
$$dw_{-np} \equiv -$$

Low-4-

$\frac{1}{\sqrt{2}}$

[illegible]


--	--	--	--	--

--	--	--	--	--	--

--	--	--	--	--

$$\begin{bmatrix} 0 & - \\ - & - \\ - & - \end{bmatrix}$$
[illegible]

next-hidden-state =

Scanned with
 CamScanner™

LSTM:

R

$$\begin{bmatrix} i \\ f \\ o \\ g \end{bmatrix} = \begin{bmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{bmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

dec

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

LSTM. step: given: $x_t, h_{t-1}, c_{t-1}, W_{hx}, W_{hh}, b$

At time "t":

$$x_t \in \mathbb{R}^D; h_{t-1} \in \mathbb{R}^H; c_{t-1} \in \mathbb{R}^H$$

$$\text{Need to learn: } W_{hx} \in \mathbb{R}^{H \times D}$$

$$W_{hh} \in \mathbb{R}^{H \times H}$$

$$b \in \mathbb{R}^H$$

$$(1): a = W_{hx} x_t + W_{hh} h_{t-1} + b$$

$$(2): \text{Split "a" into 4 parts } v_i \in \mathbb{R}^H:$$

np.split

$$a_i = \text{first } H \text{ elem.}$$

$$a_o = \text{---}$$

$$a_f = \text{2nd } H \text{ elem.}$$

$$a_g = \text{---}$$

$$i = \sigma(a_i) \quad o = \sigma(a_o)$$

$$f = \sigma(a_f) \quad g = \tanh(a_g)$$

!! DO:

$$a = x_t W_{hx} + h_{t-1} W_{hh}$$

backward:

$$\text{input: } (2, 3, 10)$$

$$W_{hx} = (10, 5)$$

$$h_0 = (2, 5)$$

$$W_{hh} = (5, 1)$$

$$b_{hx} = (5, 1)$$

$$h_t = (2, 3, 5)$$

$$\text{prev } h = (2, 5)$$

$$h_0 =$$

$$t=0:$$

$$x_t = 1 \times 10$$

$$\text{next } h_t = ((1, 10); (1, 5); (10, 5); (5, 5); (5, 1))$$

$$z = 1 \times 10 \times 10 \times 5 + (1, 5) \times (5, 1) = (1, 5) + (1, 5) = (2, 5)$$

$$\text{next } h_t = (1, 5)$$

$$\text{hidden states } (h_t) = \text{next hidden.}$$

backward pass:

$$x; \text{inp}; W_{hx}; W_{hh}; h_0; b$$

$$\text{arr} =$$

$$(1, 5); (2, 3, 10); (10, 5); (5, 5); (2, 5); (5, 1)$$

$$\text{dinp} =$$

SOLO

$$dh_t = (2, 5); \quad \hat{z} = (2, 5); \quad inp = (2, 3, 10); \quad W_{hx} = (10, 5); \quad W_{hh} = (5, 5); \quad b =$$

$$d\hat{z} = (2, 5).$$

$$dinp = \hat{z} \cdot d\hat{z} \times W_{hx}^T \Rightarrow dinp = (2, 10); \quad dW_{hh} = (5, 5);$$

$$dW_{hx} = inp^T \cdot d\hat{z}$$

$$dW_{hx} = (10, 5).$$

$$dh_{t-1} = \dots$$

$$dh_{t-1} = (10, 5)$$

backward pass:

$$dinp = (2, 3, 10); \quad dW_{hx} = (10, 5) \quad dh_t = (2, 5)$$

$$dho = (2, 5); \quad dW_{hh} = (5, 5)$$

$$i = 0, 1, 2; \rightarrow dh_{t-i} = [- \dots -]$$

$$t = 3, 2, 1$$

$$dh_{t+1} = dh_t(i) + dh_{t-1}$$

$$dx, dh_{t-1}, dW_{hx}, dW_{hh}, db = \text{backwardstep}(dh_{t+1}, \text{cache})$$

$$dx = dh_{t+1}$$

one for loop:

$$t = 3, 2, 1$$

$t = \text{reverse}$:

$$dh_t = dh_t$$

$$dh_t(i) = dh_{states}[t, i]; \Rightarrow \text{access "t" for all samples.}$$

LOGIC:

$$\begin{bmatrix} i \\ 0 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} f \\ f \\ f \\ \tanh \end{bmatrix} W$$

$$c_t = f * c_{t-1} + i \odot g$$

$$h_t = 0 \odot \tanh(c_t)$$

$$da = \begin{bmatrix} da_t \\ dg_t \\ dho \\ d\hat{z} \end{bmatrix} \quad \text{cache}$$

backward pass

$$\begin{pmatrix} n+1 \\ x_t \end{pmatrix}$$

Given: dc_t & dh_t

$$dinp = da \times W_{hx}^T$$

$$dh_{t-1} = da \times W_{hh}^T$$

$$dc_{t-1} = dc_t * f$$

$$dW_{hx} = inp^T * da$$

$$dW_{hh} = h_{t-1}^T * da$$

$$dbias =$$

$$dt = dc_t * g$$

$$do = dh_t \times \frac{\partial \tanh}{\partial c} = dh_t * \tanh'(c_t)$$

$$dg = dc_t * i$$

$$dc_t = dc_t + dh_t * 0$$

$$dc_{t-1} = dc_t * c_{t-1}$$

RNN model:

Input: $V_{D \times 1}$

Process with RNN: $H_{H \times 1}$

Output: W_{W} ; lengths of seq.

(1) $V_{D \times 1} \rightarrow \text{Init. } h_0 \Rightarrow$
(batch, D) \rightarrow (batch, H)

(2) W.E: Take indices of captions \rightarrow Feed-W.E.

3) Seq Processing

SOLO