

RNNs

John Doe



Recurrent Neural Networks:

Sentence of m words

$$P(w_1, \dots, w_m) = \prod_{i=1}^m P(w_i | w_1, \dots, w_{i-1})$$

The probability of a sentence is the product of probabilities before it

e.g. "He went to buy chocolate" \Rightarrow probability of chocolate given "He went to buy some"
multiplied by the probability of "some" given "He went to buy" etc

Score of observing a sentence

\Rightarrow Scoring mechanism

predict next word, then generate new text

Generative model

1. Tokenize text

predictions on a word / word basis

NLTK's word tokenization and sent-tokenize

2. Remove infrequent words

Vocabulary-size \Rightarrow 800 most common

Rest of tokens are replaced by unknown_token

unknown_token becomes like any other word, to predict

\Rightarrow random sample after

\Rightarrow keep generating until no unknown tokens

3. Prepend start and end tokens to start and end of a sentence.

\Rightarrow Figure out which words would likely start the first word of a sentence

4. Training data matrices

words β indices

index \rightarrow word

word \rightarrow index

$x [0, 179, 3u1, 416]$
 $(179, 3u1, 416, 1)$

Goal is to predict the next word so
 y is just the x vector shifted by
one position with the last element
being the Sentence-End token.

Input x is a sequence of words, each x_t is a single word

Represent each word as a one hot vector of size vocabulary size

x is the matrix, and x_t is the word when each row represents a word

Output: Vector of Vocabulary size \Rightarrow each element represents the probability of the word being the next word in the sentence.

$$S_t = \tanh(u x_t + w s_{t-1})$$

$$\alpha_t = \text{softmax}(V S_t)$$

Vocab size $C = 8000$ eg.

Hidden layer size $H = 100$

$$x_t \in \mathbb{R}^{8000}$$

$$o_t \in \mathbb{R}^{8000}$$

$$s_t \in \mathbb{R}^{100}$$

$$u \in \mathbb{R}^{100 \times 8000}$$

$$v \in \mathbb{R}^{8000 \times 100}$$

$$w \in \mathbb{R}^{100 \times 100}$$

We want to learn u, v and w during training

We want to keep v (Vocab size)

Small as possible (Large mat mul)

Calculating Loss : find u, v, w that minimize loss function

Cross entropy :

$$L(y, o) = -\frac{1}{N} \sum_{n \in N} y_n \log o_n$$

Sums over training examples and adds loss based

o_n how far our predictions are.

SGD and Backpropagation through time

Calculate $\frac{\partial L}{\partial u}$ $\frac{\partial L}{\partial v}$ $\frac{\partial L}{\partial w}$

Gradient checking:

$$\frac{\partial L}{\partial \theta} \approx \lim_{h \rightarrow 0} \frac{J(\theta + h) - J(\theta - h)}{2h}$$

Compare to gradient calculated for back propagation to the gradient estimated by the formula

The vanilla RNN cannot generate text because it's unable to learn dependencies between words that are several steps apart.

Back propagation through time

Fancy algorithm on an unrolled RNN

s_t to \hat{y}_t change

$$s_t = \tanh(u x_t + u s_{t-1})$$

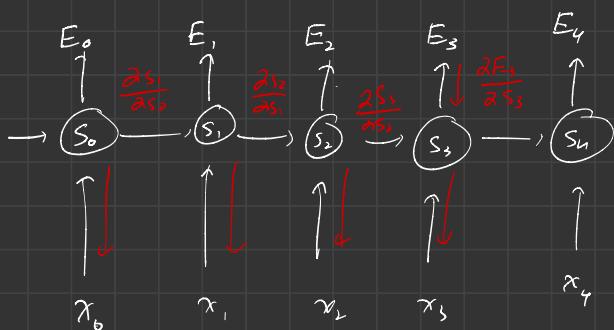
$$\hat{y}_t = \text{softmax}(V s_t)$$

$$E_t(y_t, \hat{y}_t) = -y_t \log \hat{y}_t$$

$$E(u_t, \hat{y}_t) = \sum_t E_t(y_t, \hat{y}_t)$$

$$= \sum_t y_t \log \hat{y}_t$$

\hat{y}_t our prediction, y_t correct word



$$\frac{\partial E}{\partial w} = \sum_t \frac{\partial E_t}{\partial w}$$

$$\frac{\partial E_3}{\partial v} = \frac{\partial E_3}{\partial \hat{y}_3} \frac{\partial \hat{y}_3}{\partial v}$$

$$= \frac{\partial E}{\partial \hat{y}_3} \frac{\partial \hat{y}_3}{\partial z_3} \frac{\partial z_3}{\partial v}$$

$$= (\hat{y}_3 - y_3) \odot s_3$$

$$\frac{\partial E_3}{\partial w} = \frac{\partial E_3}{\partial \hat{y}_3} \frac{\partial \hat{y}_3}{\partial s_3} \frac{\partial s_3}{\partial w}$$

$$\frac{\partial E_3}{\partial w} = \sum_{k=0}^3 \frac{\partial E_3}{\partial \hat{y}_3} \frac{\partial \hat{y}_3}{\partial s_3} \frac{\partial s_3}{\partial s_k} \frac{\partial s_k}{\partial w}$$

$$t=0 \rightarrow t=3$$

Define error as well:

$$\delta_2^3 = \frac{\partial E_3}{\partial z_2} = \frac{\partial E_3}{\partial s_3} \frac{\partial s_3}{\partial s_2} \frac{\partial s_2}{\partial w}$$

$$z_2 = w x_2 + b$$

Vanishing gradient problem:

Consider $\frac{\partial E_3}{\partial w} = \sum_{k=0}^3 \frac{\partial E_3}{\partial \hat{y}_k} \frac{\partial \hat{y}_k}{\partial s_k} \frac{\partial s_k}{\partial s_i} \frac{\partial s_i}{\partial w}$

↓
chain rule

If we take the derivative of a

vector function wrt to a vector, the output is a Jacobian Matrix

$$\frac{\partial s_i}{\partial s_k} = \left(\prod_{j=k+1}^3 \frac{\partial s_j}{\partial s_{j-1}} \right) \frac{\partial s_k}{\partial w}$$

Saturated neurons: neurons with zero gradient that drive other gradients in previous layers to zero
⇒ Gradient contributions from far away become zero

⇒ Initialize w matrix effectively
ReLU activation function.

LSTM networks:

uses a gating mechanism

$$i = \sigma(x_t u^i + s_{t-1} w^i)$$

$$f = \sigma(x_t u^f + s_{t-1} w^f)$$

$$o = \sigma(x_t u^o + s_{t-1} w^o)$$

$$g = \tanh(x_t u^g + s_{t-1} w^g)$$

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

$$s_t = \tanh(c_t) \circ o$$

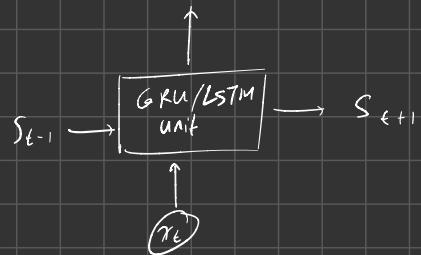
Treat LSTMs as a black box

Given current input and hidden state, they compute the next hidden layer in some way.

i, f, o are input, forget, and output gates

Sigmoid makes values between 0 and 1

Multiplying element wise \Rightarrow defines how much of a vector you want to let through.



Input gate defines how much newly computed state lct through
Forget gate defines how much of previous state you want to let through

Output gate defines how much you want to expose to the external network (hidden layers and next time step)

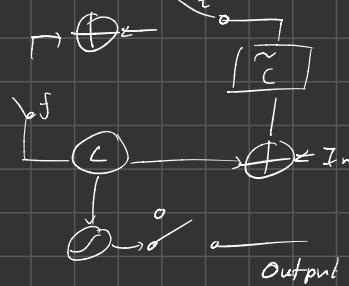
All gates have dimensions d_s , size of the hidden layer.

g is a candidate hidden state that is computed based on the previous input and previous hidden state.

c_t is the internal memory of the unit previous memory multiplied by the forget gate and newly computed hidden state g , then multiplied by the input gate.

"that we want to combine previous memory and the new input"

Given memory C_t , computing the output hidden state involves multiplying by the output gate.



Plain RNNs are a special case of LSTMs
(fix all input gates to 1)
(fix all forget states to 0)

GRUs

$$z = \sigma(x_t u^z + s_{t-1} w^z)$$

$$r = \sigma(x_t u^r + s_{t-1} w^r)$$

$$h = \tanh(x_t u^h + (s_{t-1} \circ r) w^h)$$

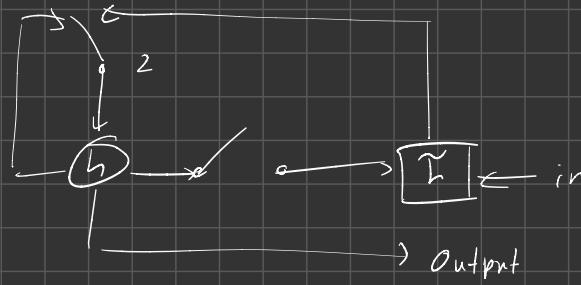
$$s_t = (1 - z) \circ h + z \circ s_{t-1}$$

A GRU has 2 gates, an LSTM has 3 gates

GRUs don't have internal memory C_t and output gate

Input and forget gates are coupled with update gate z , reset gate r , applied directly to the hidden state. No second non-linearity applied either.

GRU:



GRU vs LSTMs

\Rightarrow No clear winner

\Rightarrow Hyper tuning is more important

understanding LSTM Networks

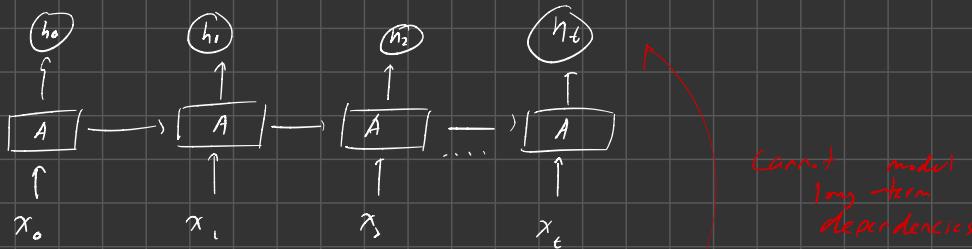
Recurrent neural networks allow information to persist



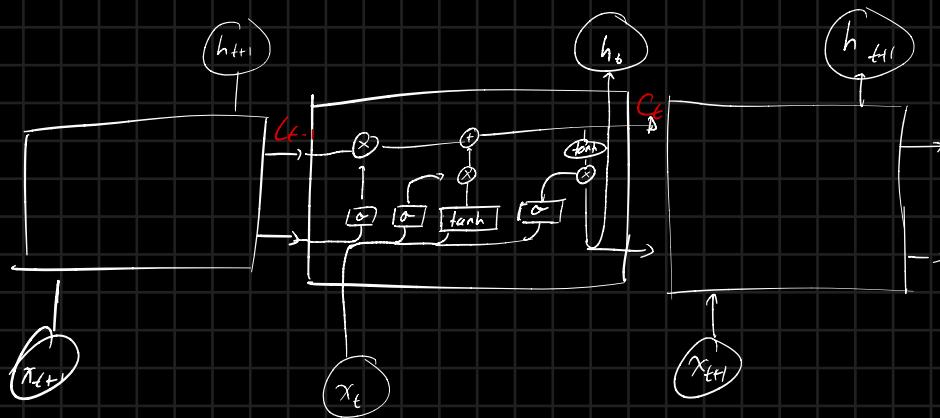
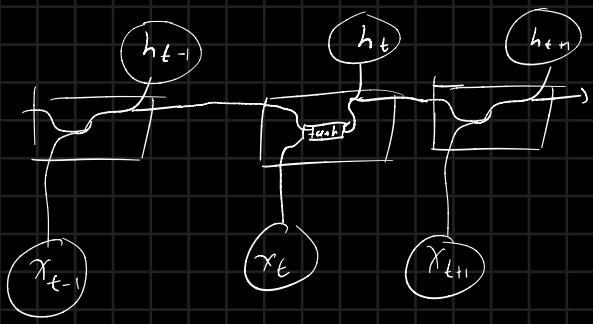
A looks at some input x_t and outputs a value h_t .

Information in a loop can be passed from one step of the network to the next

=> Multiple copies of the same network, each passing a message to its successor



Cannot model
long-term
dependencies



$\boxed{\quad}$ neural network layer

\circ pointwise operation

\longrightarrow vector transfer

\curvearrowright concatenation

\swarrow copy

Each line carries a vector from the output of one node to the input of others

Cell State \rightarrow flow of information

Forget (add into) to cell state

Gates: optimally let information through



Sigmoid (neural) net layer

\Rightarrow pointwise multiplication operation

between 0 and 1, denoting how much of each component should be let through

① Decide what information to throw away, from cell state

Sigmoid layer \rightarrow forget gate layer

Precision is made by another sigmoid layer called the forget layer.

Looks at h_{t-1} , and x_t , outputting a value between 0 and 1, for each number in the cell state.

$0 \rightarrow 1$
forget keep

$$f_t = \sigma(\omega_f \cdot [h_{t-1}, x_t] + b_f)$$

② Decide what information to store in the cell state

Sigmoid layer called the input gate layer decides what values to update

tanh layer creates a vector of new candidate, \tilde{C}_t that could be added to the state

$$i_t = \sigma(\omega_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(\omega_c \cdot [h_{t-1}, x_t] + b_c)$$

③ update the old cell state C_{t-1} into the new cell state C_t

\Rightarrow multiply old state by f_t , forgetting things we decided to forget

\Rightarrow Add $i_t * \tilde{C}_t$, adding candidate values \times how much we decided to update them by

$$\text{overall: } C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

④ prepare to output

Run a sigmoid layer, that decides what parts of the cell state to output

put cell state into tanh (-1, 1) multiply by output of the sigmoid gate
to output only parts we decide to.

$$o_t = \sigma(\omega_o \cdot [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(C_t)$$

Some LSTMs have peep holes

Some forget and update in one step

CPU → combines forget and update gate into a single "update gate"

CNN - Long short - term Memory Networks

CNNs for feature extraction

LSTMs → sequence prediction

CNN → trained on image classification
→ re purpose as a feature extractor

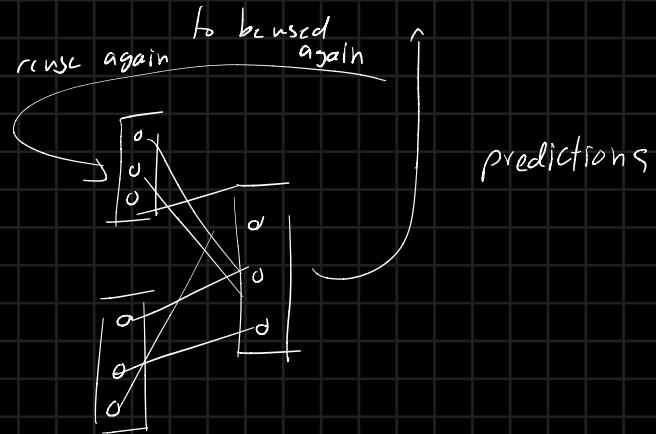
encoder

→ RNN decoder that generates text

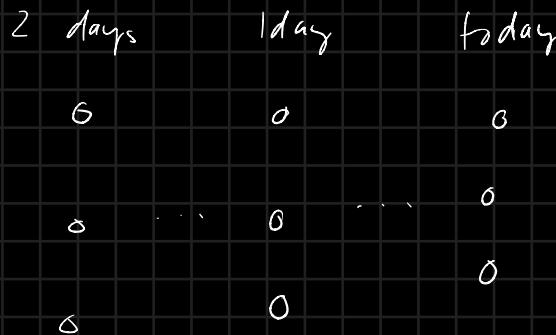
pass the output matrix to LSTM as a single time step

= Time Distributed wrapper
= applies layer multiple times

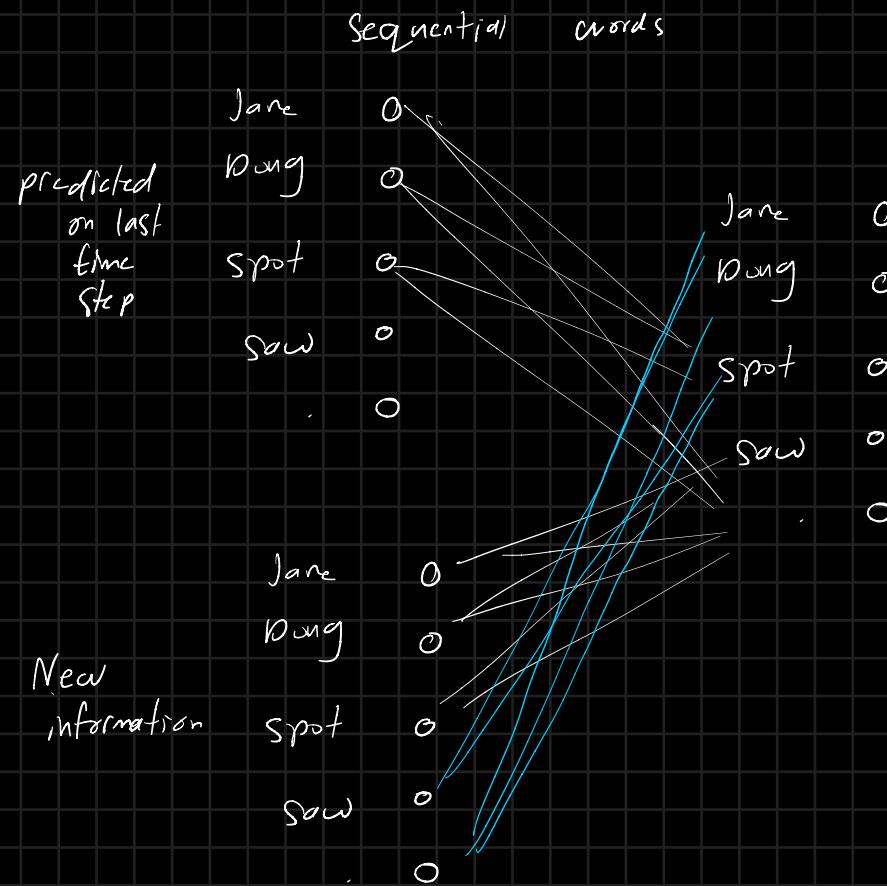
RNNs ≠ LSTMs



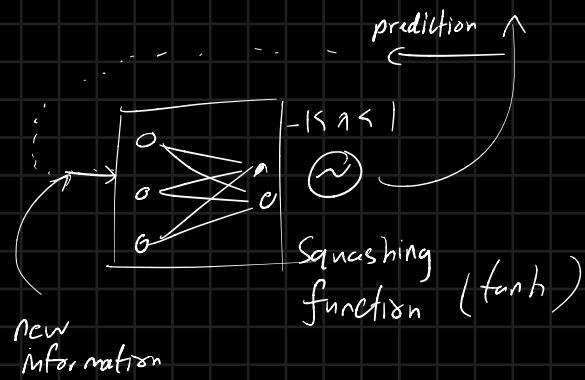
unwrapping this loop



Trace forward / backward



Recurrent neural networks

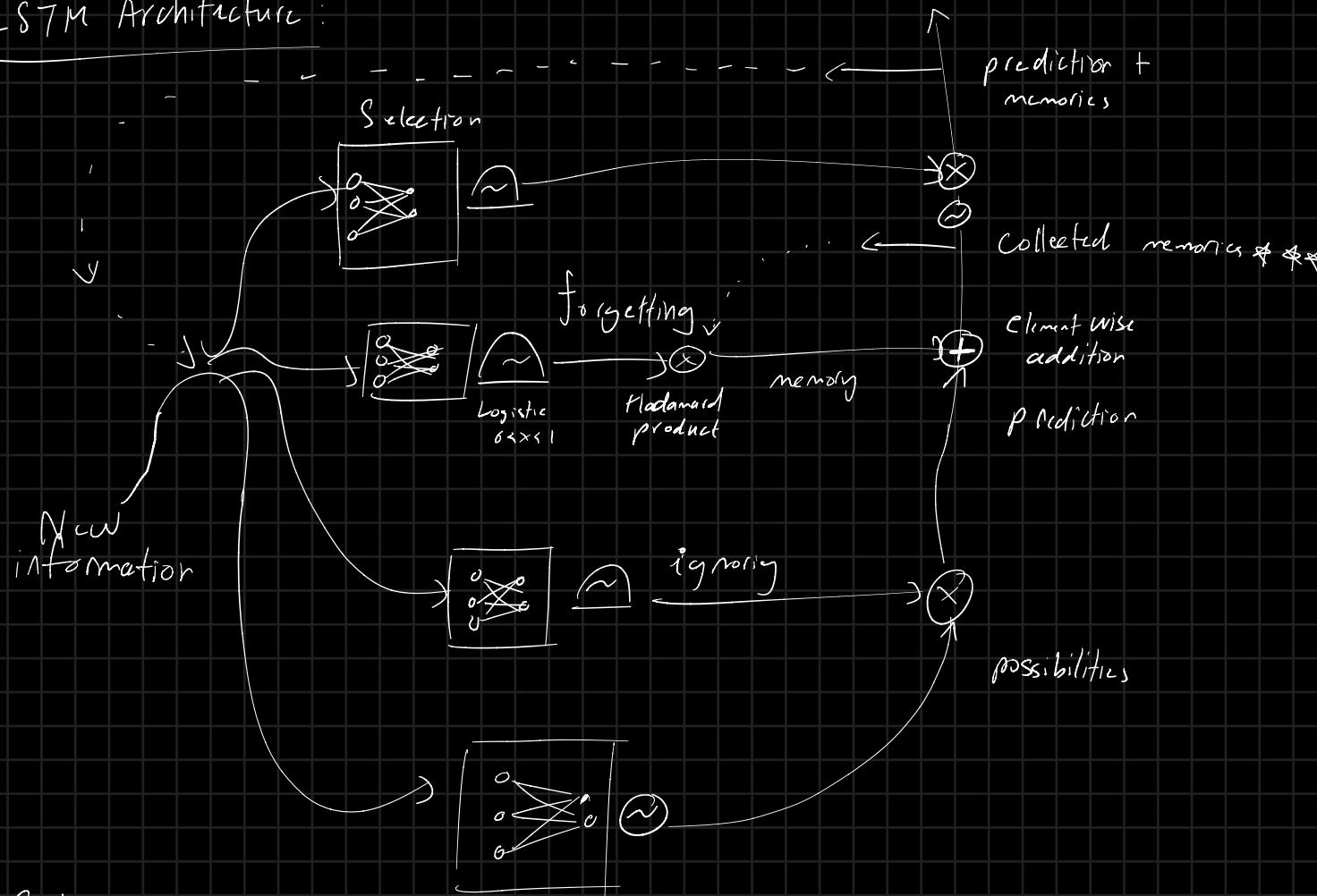


Mistakes an RNN makes:

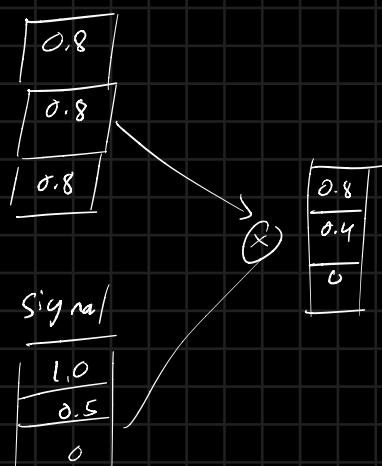
Dong Sow Dong. (possibly)

RNNs only look back at one time step

LSTM Architecture:



Gating:



Creates a gate that determines
Selection: what memories to let out as
a prediction

Collect memories over time

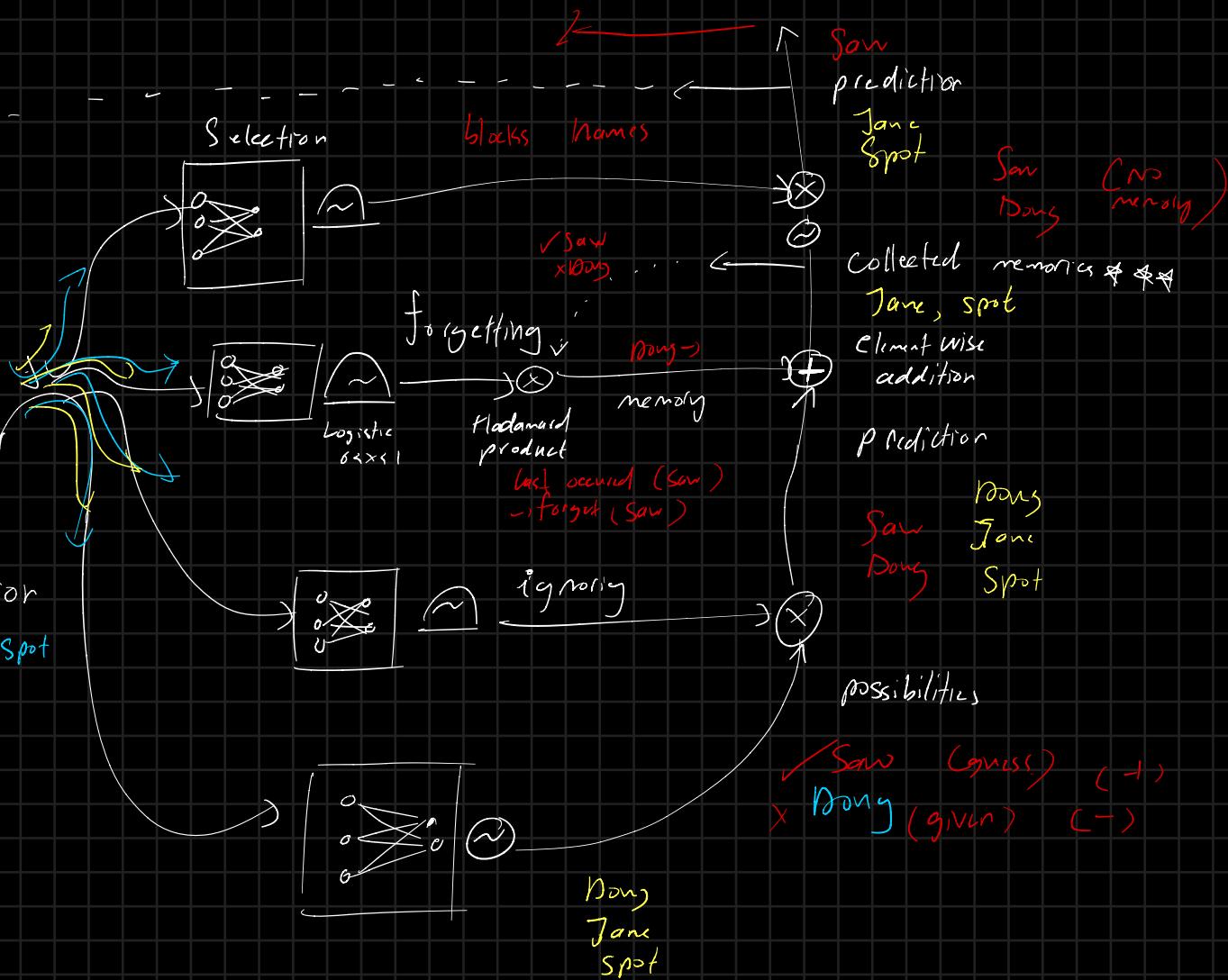
Ignoring: Creates a gate that determines what
predictions in memory go forward

Doug
Jane
Spot

Saw

New information

Jane saw Spot
Doug ...



LSTMs can look back several time steps

RNNs can only look back a few

Initial values $c_0 = 0$, $h_0 = 0$

$$f_t = \sigma_f(w_f x_t + u_f h_{t-1} + b_f)$$

$$i_t = \sigma_i(w_i x_t + u_i h_{t-1} + b_i)$$

$$\sigma_t = \sigma_g(c_0, x_t, u_t, h_{t-1}, b_o)$$

$$c_t = f_t \circ c_{t-1} + i_t \circ \sigma_c(w_c x_t + u_c h_{t-1} + b_c)$$

$$h_t = \sigma_h \circ \sigma_h(c_t)$$

x_t → input vector

h_t → output vector

c_t → Cell state vector

w, u, b → parameter and matrices vector

f_t, i_t, σ_t gate vectors

f_t : forget gate vector → old information

i_t : acquiring new information → input gate

σ_t : output gate vector

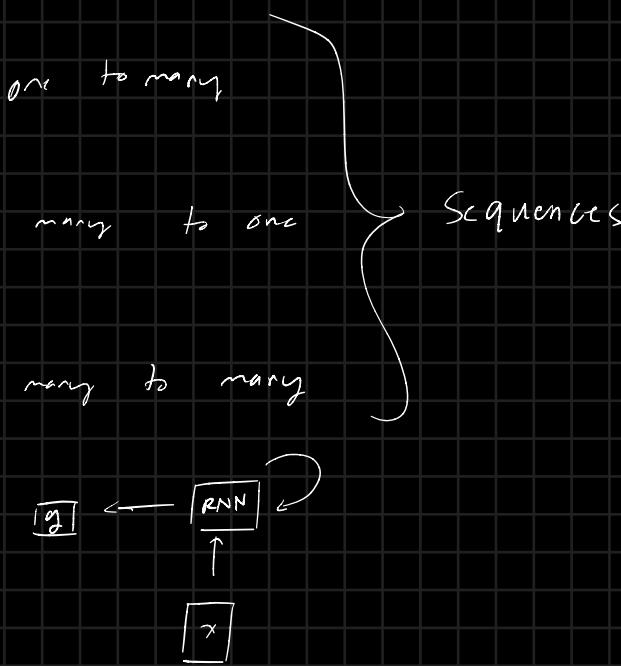
σ_g : sigmoid

σ_c : tanh

σ_h : tanh, parabolic LSTM → $\sigma_h(x) = x$

Karpathy RNNs Lecture

Flexibility



process a sequence of vectors x_t ,
applying a recurrence formula

$$h_t = f_w(h_{t-1}, x_t)$$

new state old state input vector at
state time t

Fixed function w

Vanilla RNN

$$h_t = f_w(h_{t-1}, x_t)$$

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

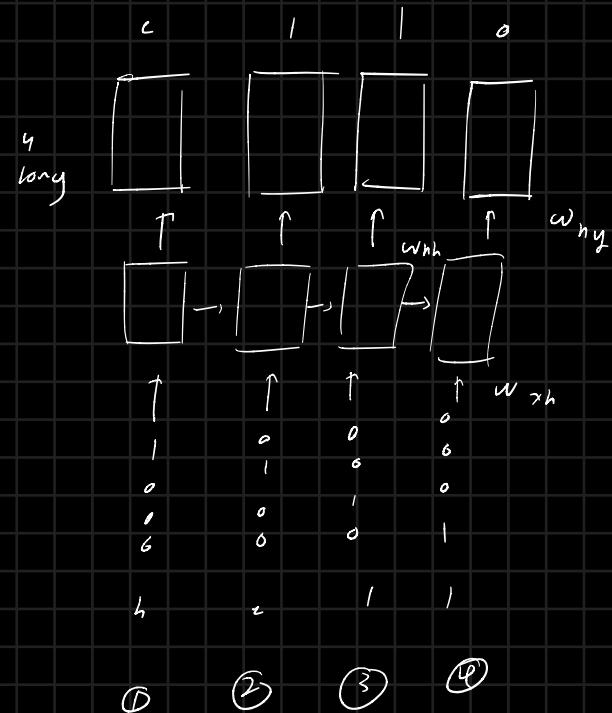
$$y_t = C_{hy} h_t$$

character level language models

feed a seq of chars

predict next char

hello one hot encode



Back propagate (all occurrences)

min-char-rnn.py

Data loading ... itos, choi

hidden_size

seq_length - 25 batch size

learning_rate

$w_{xh} = (\text{vocab}, \text{hidden})$

$w_{hh} = (\text{hidden}, \text{hidden})$

$why = (\text{vocab}, \text{hidden})$

$b_h = (\text{hidden}, 1)$

$b_y = (\text{vocab}, 1)$

Sample batch

Targets is next character into future

Sample, feed, sample, feed

Loss fn \rightarrow fn in a while True loop

Loss: $h_t = \tanh(w_{hh} h_{t-1} + w_{xh} x_t)$

$y_t = w_{hy} h_t$

(using dictionary) to keep track of time steps
softmax

↓

Loss at correct index

Back propagate through softmax

Back prop through activation

Accumulate gradients

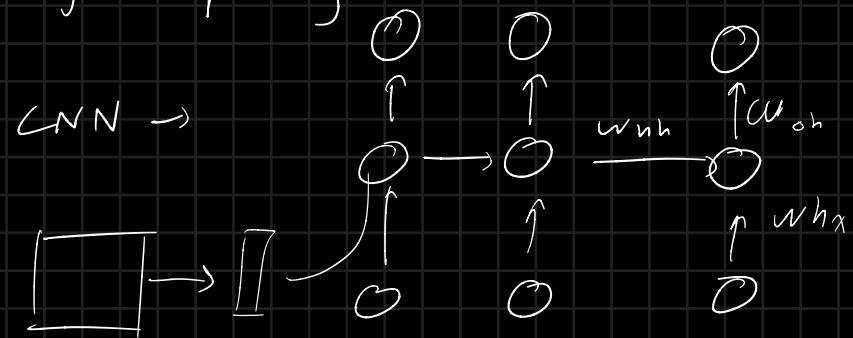
$d w_{hh} +=$

$d w_{hh} +=$

BPP all at once (all examples in batch)

RNNs can learn syntax

Image captioning



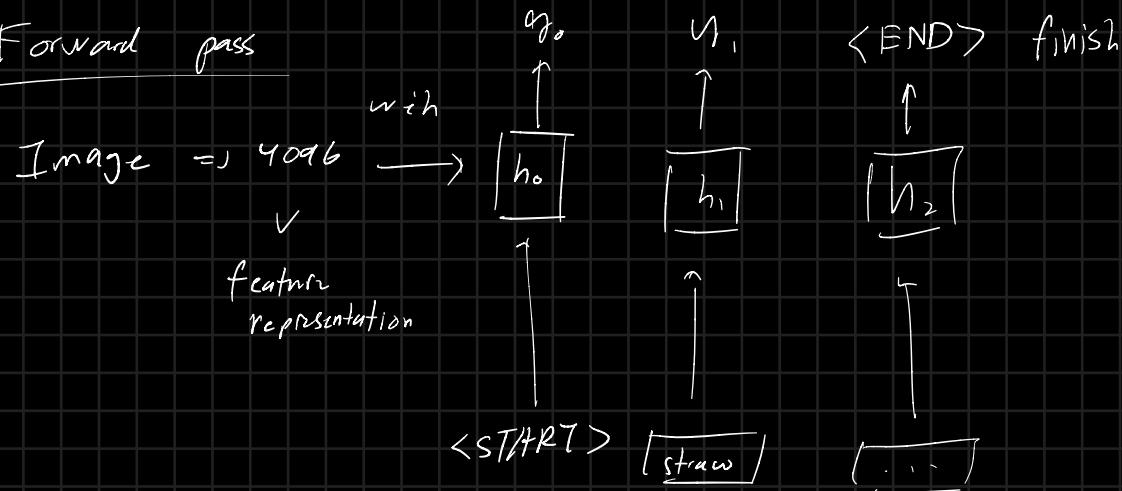
Before

$$h = \tanh(w_{vh} * x + w_{hh} * h)$$

After

$$h = \tanh(w_{vh} * x + w_{hh} * h + w_{wh} * v)$$

Forward pass



Initialize at random, BP all at once

Image Sentence datasets \Rightarrow Microsoft COCO

"Attention" over the image

LSTM

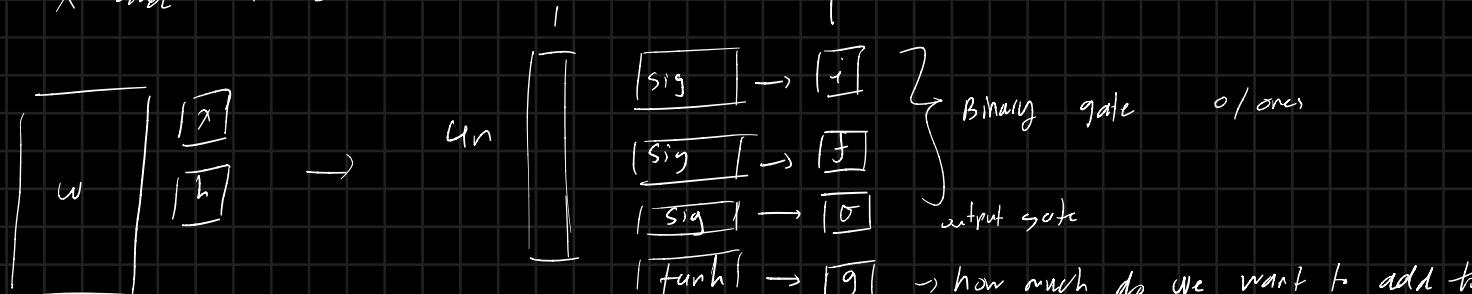
$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \text{sig} \\ \text{sig} \\ \text{sig} \\ \text{tanh} \end{pmatrix} w^l \begin{pmatrix} h^{l-1} \\ h^l \end{pmatrix}$$

update prev c_{t-1} and to every cell

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

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

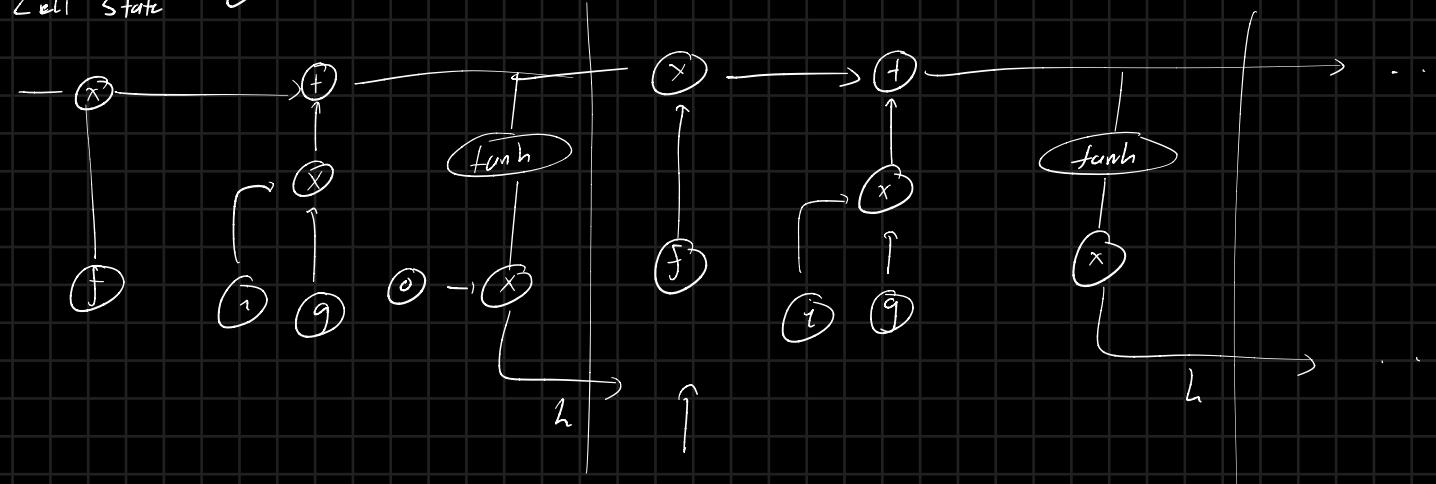
below x and h size n



$4n \times 2n$

one timestep

cell state c



LSTMs never have vanishing gradients \Rightarrow Gradient flows of additive interactions
vs multiplying again
and again

If largest eigenvalue > 1 , explode

If largest eigen < 1 , vanish

The Unreasonable Effectiveness of RNNs

RNNs operate over sequences of vectors

RNNs are turing complete

RNN optimizes over programs

output vectors are influenced by entire history of inputs

RNN = RNN()

y = rnn.step(x)

RNN class has an internal state updated by step() \Rightarrow hidden vector h

def step(self, x):

self.h = np.tanh(np.dot(self.w_hh, self.h) + np.dot(self.w_xh, x)) more complex for
 $y = \text{np.tanh}(\text{self.w}_h \cdot \text{h}_{t-1} + \text{self.w}_x \cdot x_t)$ an LSTM

return y

3 matrices

elementwise addition

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

element wise

$$\underline{y_1} = \text{rnn1.step}(x)$$
$$y = \text{rnn2.step}(\underline{y_1})$$

Next char based on entire context

1 of K encoding

confidence-based output

LSTMS Wikipedia

Common structure:

Cell
input gate \Rightarrow for current state

output gate \Rightarrow for current state

forget gate \Rightarrow for previous state

LSTMs allow for gradients to flow with no attenuation.

w weights of input connection

U weights of the recurrent connection

$$f_t = \sigma_g (w_f x_t + M_f h_{t-1} + b_f)$$

$$i_t = \sigma_g (w_i x_t + M_i h_{t-1} + b_i)$$

$$\theta_t = \sigma_g (w_o x_t + M_o h_{t-1} + b_o)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t$$

$$h_t = \sigma_g \odot \sigma_h (c_t)$$

x_t = input vector to LSTM units

i input

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

a output

$u \in \mathbb{R}^{h \times h}$

f forget

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

h_t hidden state

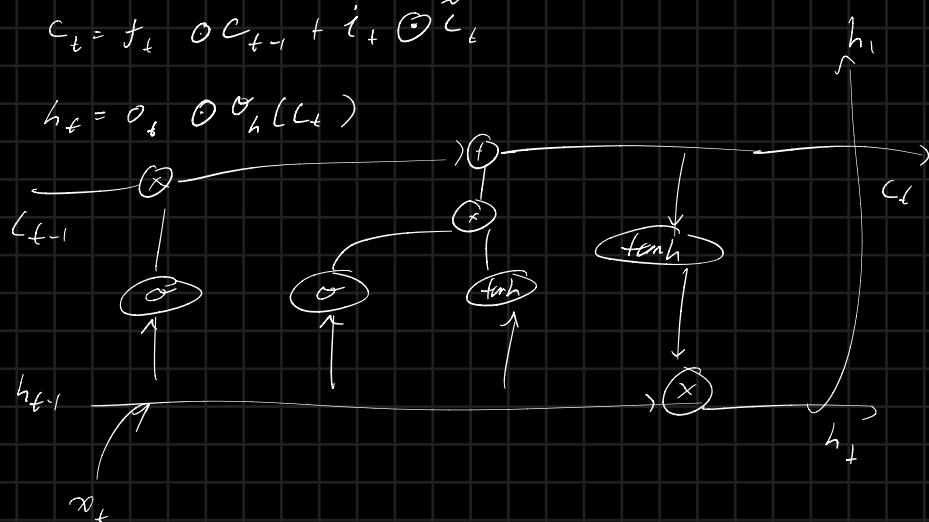
\tilde{c}_t cell input activation

σ_g sigmoid

c_t cell state

σ_a hyperbolic tangent

σ_h also tanh



LSTM (again)

$$h_{t-1} \in \mathbb{R}^k \quad h_t = \sigma(c_x x_t + c_h h_{t-1})$$

$$x_t \in \mathbb{R}^k \quad i_t = \sigma(w h_t + b)$$

$$i_t = \sigma(I_x x_t + I_h h_{t-1} + b_i)$$

$$f_t = \sigma(F_x x_t + F_h h_{t-1} + b_f)$$

$$o_t = \sigma(O_x x_t + O_h h_{t-1} + b_o)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot \tanh(c_x x_t + c_h h_{t-1} + b_c)$$

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

$$I_x, I_h, \dots \in \mathbb{R}^{k \times k}$$

$$h_s = c_s = 0$$

$$y_t = f(w h_t + b)$$

Main idea: Enable training of longer sequences by providing a fast path

to backpropagate information further down in memory

\Rightarrow C vector is not multiplied by matrices on its path

We concatenate h and π into a single column vector

Must sum we properly combine gradients at each points and have gradients flow at each node.

$$x_h = \pi^t : h^{t-1}$$

\downarrow
concatenate

$$\alpha'(k) = \alpha(k)(1 - \alpha(k))$$

$$f = \alpha(w_f \cdot x_h + b_f)$$

$$i = \alpha(w_i \cdot x_h + b_i)$$

$$o = \alpha(w_o \cdot x_h + b_o)$$

$$cc = \tanh(w_{cc} \cdot x_h + b_{cc})$$

$$c^{(t)} = c^{(t-1)} \circ f + cc \circ i$$

$$h^{(t)} = \tanh(c^{(t)}) \quad 00$$

$$y^{(t)} = w_y \cdot h^{(t)} + b_y$$

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

Time Series forecasting

Univariate $\xrightarrow{\text{Bi-directional LSTM}}$
 $\xrightarrow{\text{stacked LSTM}}$
 $\xrightarrow{\text{conv LSTM}}$

multi-variate $\xrightarrow{\text{Multi-step series}}$
 $\xrightarrow{\text{Multi parallel Series}}$
multi-step

multi-variate multi-step $\xrightarrow{\text{More than 1 observation per step}}$
 $\xrightarrow{\text{More than 1 time step per observation} \Rightarrow \text{encoder \& decoder}}$

univariate: sequence (lmer) to single output,
over several time steps

Stacked LSTM: multiple hidden LSTM layers can be stacked up on each other

Bi-directional LSTM: Learn input seq, forward and backward.

CNN LSTM: Output of CNN goes into LSTM (should be 1D)

Image Captioning

Understand content of the image

Turn into words at the right order

- prepare photo data
- prepare text data
- build model
- Train
- Evaluate
- Generate Samples

Flickr 8K dataset

8000 images \Rightarrow Each image has 5 possible captions

BLEU Scores

Regular CNN

\Rightarrow Last layer should be linear output of features

Interested in the internal classification before an actual classification is made.

\Rightarrow 1×4096 features

Reshape image to $224 \times 224 \times 3$

Text cleaning:

Each photo has an id

\Rightarrow Create dictionary to map photo identifications to descriptions.

Each id maps one or more descriptions.

Converts words to lower case.

Remove punctuation.

Remove words of one char length

Remove words with numbers in them

Vocab-size \Rightarrow Small but expressive as possible

\Rightarrow Create "cleaned descriptions" by removing words

Add start-seq and end-seq to each caption.

One img, one caption

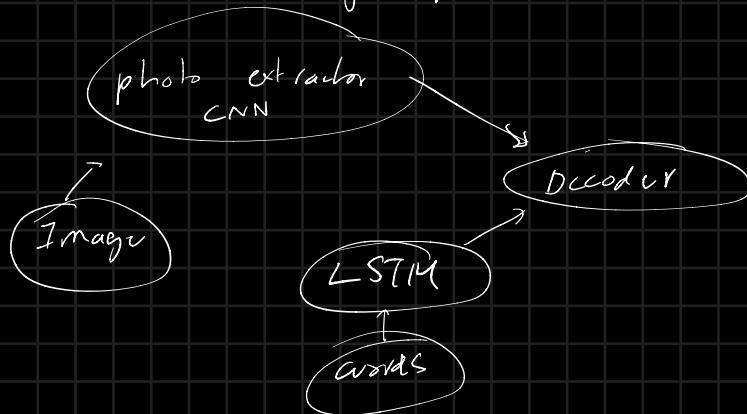
Tokonize Captions

When model is used for sampling, feed itself into model.

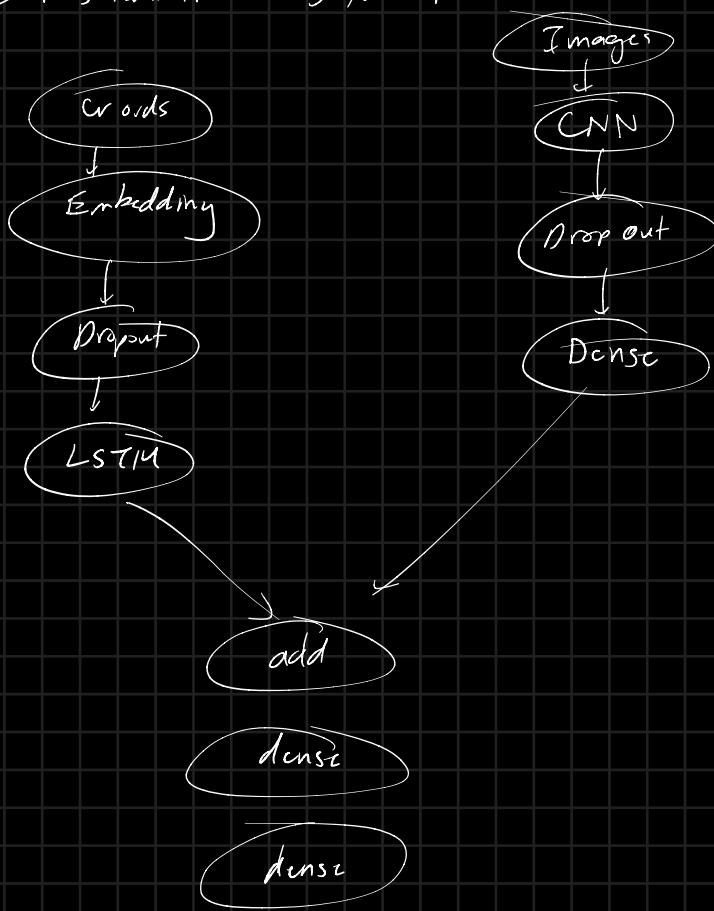
Need to set maximum sequence length \Rightarrow use a mask for padded values

\hookrightarrow Feed words into word embedding layer, or one hot encode.

Use a merge model:



use regularization \Rightarrow 50% dropout



Evaluate

corpus-blanc()