

perception in cats: simple cells, complex cells

Simple cells - respond to oriented edges.

History

SUMMER PROJECT - MIT 1966

reconstruction of vision system

David Marr (Vision Book 1970s)

Input image \rightarrow Edge image \rightarrow 2½-D sketch



3 D

model

Image Segmentation \rightarrow Normalized Cut (Shi & Malik, 1997)

Face detection - Viola & Jones 2001

and AdaBoost

SIFT & Object Recognition, David Lowe 1999
 \rightarrow feature based obj detection.

Spatial Pyramid matching, Lazebnik 2006

+ SVM

Human detection

Histogram of gradients, Deformable Body Part Model.
 (2003) 2009

PASCAL VOC (Visual Object Recognition) (2006)

IMAGENET 22K class, 14M images.

done using Amazon Mechanical Turk platform (crowdsource platform)

IMAGENET LSVRC \rightarrow 1k class

1.4M images

Image classification using

i) KNNs:

take L_1 or L_2 distance between images and group classes by K value. L_1 is coordinate dependent while L_2 works generally.

Not ideal as $O(n)$ for test case and can't understand subtle changes in image.

ii) Linear clasif.

$$32 \times 32 \times 3 \rightarrow f(x, w) \rightarrow 10$$

$$\begin{matrix} w \times x + b \\ \downarrow \\ 10 \times 3072 \end{matrix} \xrightarrow{3072 \times 1}$$

Sathvik
Vaidya

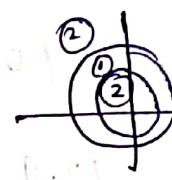
Linear classifier just averages all variations in input & comes up with one template to identify a class.

Can imagine it as drawing linear boundaries b/w different categories

difficulty in cases such as

Class 1: $\|x\|_2 \leq 2$

Class 2: everything else



can't classify.

Multiclass SVM loss:

$$L(x) = \max_{j \neq i} (y_i - y_j + \gamma)$$

	eat	car	frog
eat	3.2	1.3	2.2
car	5.1	14.9	2.5
frog	-1.7	2.0	-3.1
loss	$l_i = \sum_{j \neq y_i} \begin{cases} 0 & \text{if } s_j - s_{y_i} \geq 1 \\ s_j - s_{y_i} & \text{otherwise} \end{cases}$		

true score

should be higher
than all the other
scores by a
margin ('')

compare score of correct category
and sum of scores of other categories

this type of loss fn is called hinge loss

for cat:

$$L_i = \max(0, 5.1 - 3.2 + 1) + \max(0, -1.7 - 3.2 + 1)$$
$$= 2.9 + 0 = 2.9$$

for car:

$$L_i = 0$$

for frog:

$$= \max(0, 2.2 + 3.1) + \max(0, 2.3 + 3.9)$$
$$= \underline{12.9}$$

$$L = (2.9 + 0 + 12.9)/3 \approx 5.3$$

Softmax can be viewed as KL divergence b/w target & computed y or as maximum likelihood estimate.

It is easier to maximise log than the prob distribution.

$$-\log \left(\frac{e^k}{\sum e^j} \right)$$

	cat	imx			
cat	3.2	24.5			$0.13^{-10.8} 0.89$
car	5.1	164	\rightarrow	0.87	0
frog	-1.7	0.18			

L) if all are small or zero, then softmax output should be log c (sanity check)

C) n. of classes

optimisation:

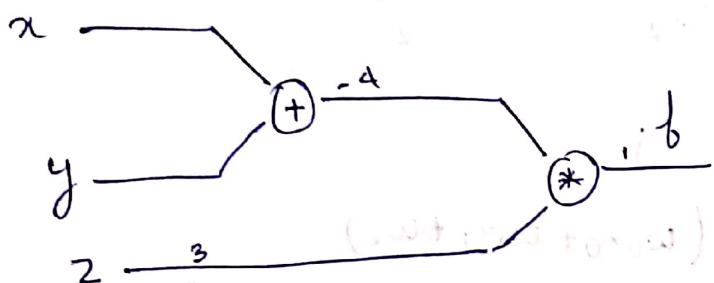
The slope of in any direction is the dot product of direction with gradient.

The direction of steepest gradient descent is the negative gradient.

Backpropagation:

Eg:

$$f(x, y, z) = (x+y)z$$



$$q = x + y$$

$$\frac{\partial q}{\partial x} = 1$$

$$\frac{\partial q}{\partial y} = 1$$

$$f = qz$$

$$\frac{\partial f}{\partial q} = z$$

$$\frac{\partial f}{\partial z} = q$$

~~$$\text{let } x = -2$$~~

~~$$y = 5$$~~

~~$$z = -4$$~~

to do backprop

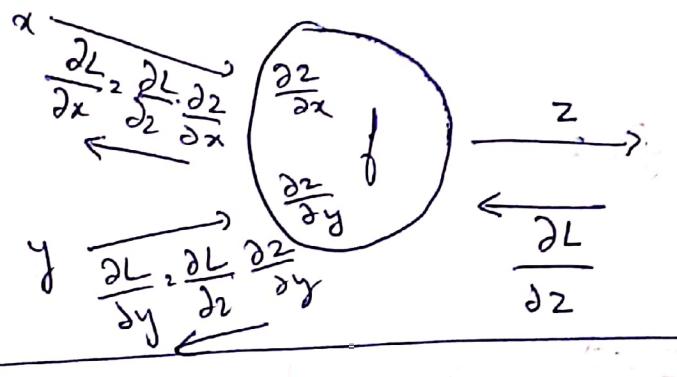
start with $\frac{\partial f}{\partial f} = 1$

then $\frac{\partial f}{\partial z} = q = 3$

$$\frac{\partial f}{\partial z} = 2 = 2 - 4$$

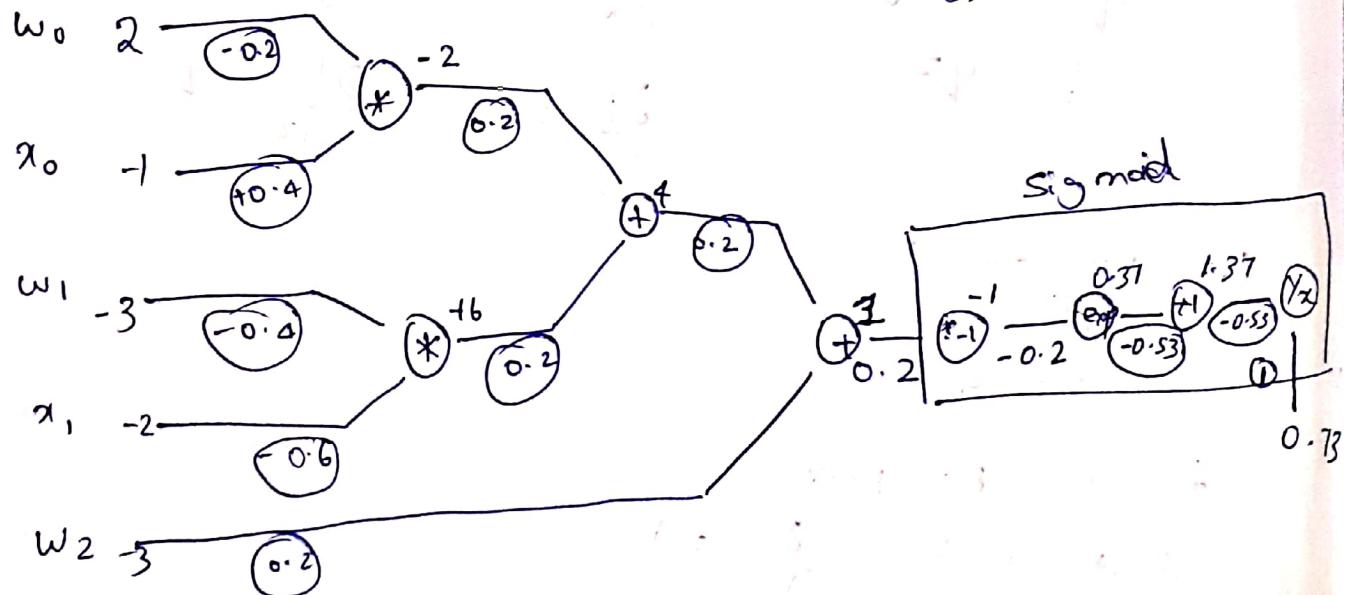
then $\frac{\partial f}{\partial y} = \frac{\partial f}{\partial g} \times \frac{\partial g}{\partial y}$
 $= 2x y^2 - 4.$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial g} \times \frac{\partial g}{\partial x}$$
 $= 2x y \neq -4$



$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2 x_2)}}$$

$$\frac{\partial \sigma(x)}{\partial x} = \sigma(x)(1-\sigma(x))$$



$$1 \times \frac{-1}{1.37^2} = -0.53$$
 $e^x(0.53), e^{-1}(-0.53)$

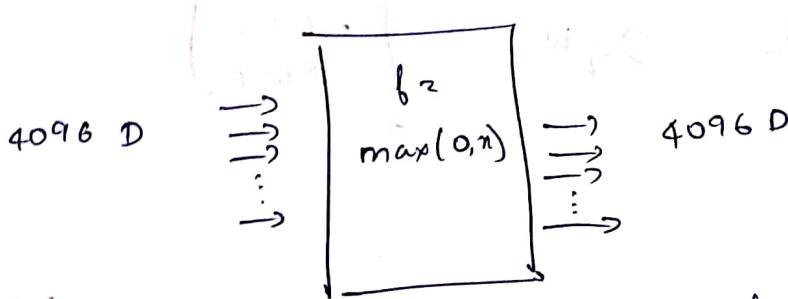
for each node,
 upstream gradient
 \times local gradient

For addition, same gradient is passed down to all nodes.

max gate \rightarrow gradient \rightarrow route (send to max i/r value)
 mul gate \rightarrow gradient \rightarrow switcher (take other value & mul gradient)
 add gate \rightarrow gradient \rightarrow distributor

On backprop when upstream nodes converge to one, then gradients add up.

if

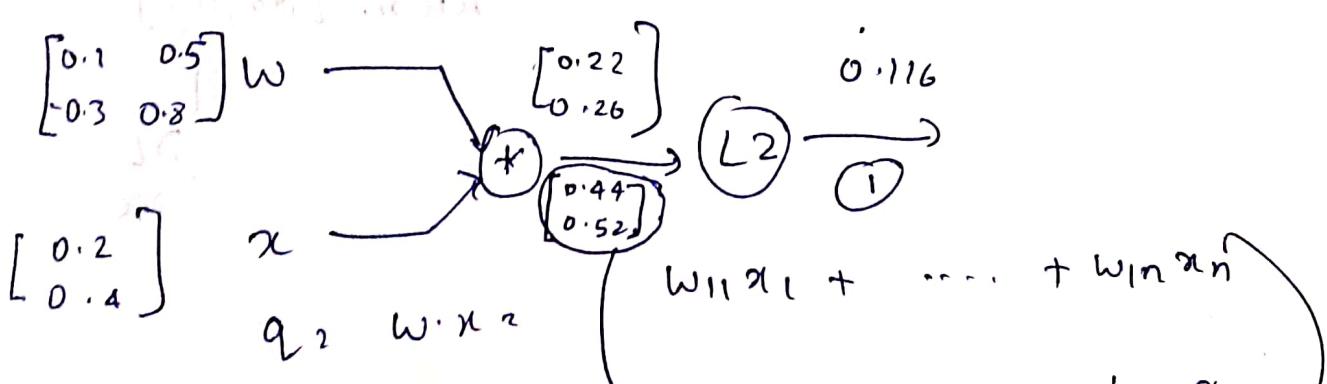


$$\left(\frac{\partial L}{\partial x} \right)^2 = \frac{\partial f}{\partial x} \cdot \frac{\partial L}{\partial f}$$

$\frac{\partial f}{\partial x} \rightarrow 4096 \times 4096$
 (Jacobians)

But only the diagonal elements in Jacobian contribute to output.

$$\text{def } f(n, w) = \|w \cdot x\|^2 = \sum_{i=1}^n (w_i x_i)^2$$



$$f(q) = \|q\|^2 = q_1^2 + q_2^2 + \dots + q_n^2$$

$$\frac{\partial f}{\partial z_i} = 2q_i$$

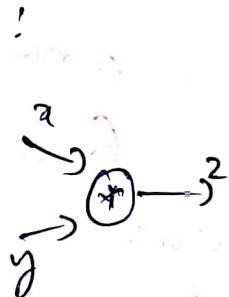
$$\nabla_{q_i} f = 2q_i = \begin{bmatrix} 2 \times 0.22 \\ 2 \times 0.26 \end{bmatrix} = \begin{bmatrix} 0.44 \\ 0.52 \end{bmatrix}$$

$$\frac{\partial f}{\partial w_{ij}} = 2q_i x_j = 2q_i x^T$$

$$\begin{bmatrix} 0.088 & 0.104 \\ 0.176 & 0.208 \end{bmatrix}$$

$$\frac{\partial f}{\partial x} = 2q^W = \begin{bmatrix} -0.112 \\ 0.636 \end{bmatrix}$$

Coding them



class mulgate(object):

def forward(x,y):
z = x * y
return ~~z~~ z

def backward(dz): $\frac{\partial L}{\partial z}$

$dx = \dots$

$dy = \dots$

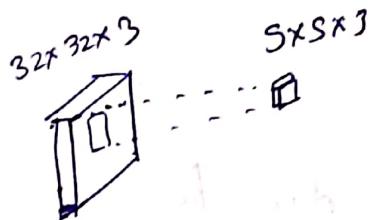
return [dx, dy]

$$\frac{\partial L}{\partial x}$$

CNNs

Hinton 2006 -

backprop works with proper initialisation.
layers pre-trained using R.Boltzman Machine
Filters always go down to full depth of image.



Conv filters not flipped (not like signal processing)

why not sigmoid as activation?

Saturated neurons kill grads

not zero centred

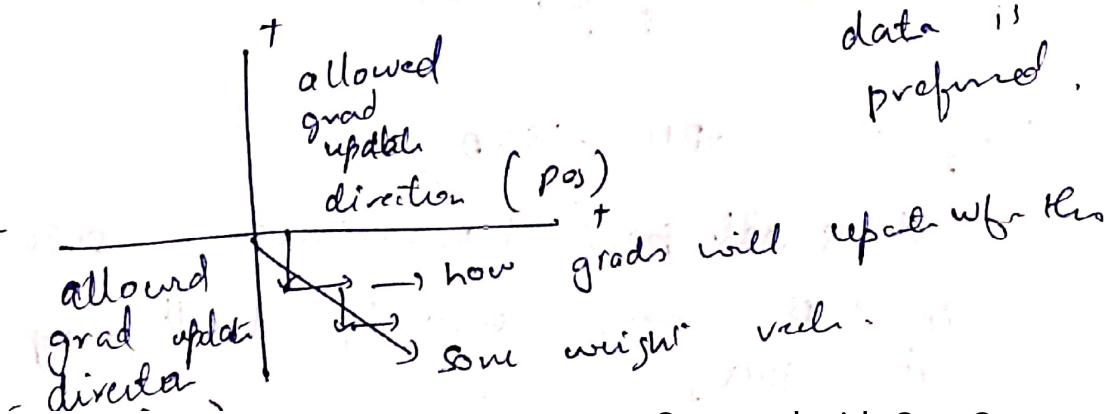
$\exp()$ is computationally exp.

If i/p is always tve or nve

then grad update will happen in same direction.

that's why zero mean

data is preferred.



tanh()

Kills neurons
grads on saturation

Relu

Convergence (6x)

Does not saturate

Computationally efficient.

Not zero centered

Relu might not work well due to bad initialisation & high learning rate

Pre processing.

Zero mean is required.
normalisation not really for images
→ mean image or per channel mean
initialisation
if $w=0$, all will be same

read Xavier Initialisation paper
& BN paper 2015

& He et al. 2015

Why is $\|g_{GD}\| > \|g_{SGD}\|$

at saddle point / minimas, even though grad is zero, velocity built up gives

it going

masking & saliency Maps to
visualise importance
also guided backprop

gradient ascent to regularise generated image.

Adversarial - more on it

For efficient Inference:

Pruning

weight Sharing

Quantisation

Low Rank Approx.

Binary / Ternary Net

Winograd Transform

wordnet (nltk package available) \rightarrow list of synonymy hypernyms etc.

localist representation

lexicon based models used before deep learning

why not use multi-hot encodings?

\rightarrow large vector size

\rightarrow no context

Distributional semantics

A word's meaning is given by other words that frequently appear close by.

Vector Representation

Word2Vec

start with a center word and predict context words around it within a range. Using this the vector representation is learnt.

likelihood $L(\theta) =$

$$\prod_{t=1}^T \prod_{-m \leq j \leq m} P(w_{t+j} | w_t; \theta)$$

$T \rightarrow$ total no of words

$m \rightarrow$ range.

To maximise prediction acc \rightarrow minimum objective function.

$$J(\theta) = -\frac{1}{T} \log L(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{w \in \text{vocab}} \log P(w_t | w_{t-1}; \theta)$$

log changes product to sum of log.

How to calculate $P(w_{t+1} | w_t; \theta)$?

1 vector for word w when it is center v_w

1 vector for word w when it is context v_w

for a new center word c , context word c ,

$$P(c|c) = \exp(v_c^\top v_c)$$

$$\frac{\exp(v_c^\top v_c)}{\sum_{w \in \text{vocab}} \exp(v_w^\top v_c)} \xrightarrow{\text{normalization over entire vocab}}$$

This $v_c^\top v_c$ compares similarity of c and c .
larger dot product = larger probability

Computing grads of all vectors of words

$$\frac{\partial}{\partial v_c} \log \frac{\exp(v_c^\top v_c)}{\sum_{w \in \text{vocab}} \exp(v_w^\top v_c)}$$

$$\text{Ansatz: } \frac{\partial}{\partial v_c} \log \exp(v_c^\top v_c) = \frac{\partial}{\partial v_c} \log \sum_{w \in \text{vocab}} \exp(v_c^\top v_w)$$

$$\frac{\partial}{\partial v_c} v_c^\top v_c = v_c \quad \frac{\partial}{\partial v_c} \log \sum_{w \in \text{vocab}} \exp(v_c^\top v_w)$$

$$= \frac{1}{\sum_{w \in \text{vocab}} \exp(v_c^\top v_w)} \cdot \frac{\partial}{\partial v_c} \sum_{w \in \text{vocab}} \exp(v_c^\top v_w)$$

$$= \frac{1}{\sum_{w=1}^V \exp(u_w^T v_c)} \cdot \sum_{w=1}^V \exp(u_w^T v_c) \cdot \frac{\partial}{\partial u} u x^T v_c$$

$$\therefore \frac{\partial}{\partial v_c} \log P(o|c) = u_o - \frac{\sum_{x=1}^V \exp(u_x^T v_c) \cdot u_x}{\sum_{w=1}^V \exp(u_w^T v_c)}$$

$$= u_o - \frac{\sum_{x=1}^V \exp(u_x^T v_c)}{\sum_{w=1}^V \exp(u_w^T v_c)} u_x$$

$$= u_o - \frac{\sum_{x=1}^V p(x|c) \cdot u_x}{\sum_{w=1}^V p(w|c) \cdot u_w}$$

\downarrow observed rep of the context word \downarrow what model thinks model should look like

- two vectors u_o and u_x for easier optimisation.
They are averaged at the end.
- It can also be done with one vector only

2 variants:

Skip gram: predict context word on center

Bag of words: predict center word on all context

negative sampling: for faster word2vec

- $P(o|c)$ denom takes too much time
- train binary logistic regressions for a true pair (center & context) versus many noise pairs (center with a random word)

$$\text{Five samples } (v_0, v_1, \dots, v_k) = -\log(\sigma(u_0^T u_1)) - \sum_{k=1}^K \log(\sigma(-u_k^T v_k))$$

- take k -ve samples.
- Maximise prob that real word appears outside, minimise prob that random words appear around the centre.

window band co-currence matrix

- large size
- sparse
- less robust models
- Reduce dimension

Eg: SVD

- for common words, use log or ceil func.
- for uncommon words, use counts, zero to zero.
- pearson corr instead of counts, zero to zero.

glove

- adding co-currence (statistic) to context.
- adding co-currence prob can encode info
- Ratios of co-currence prob as
- capturing ratios of co-currence prob as linear meanings components in word vector space

log bilinear model $w_i \cdot w_j = \log p(i|j)$

$$w_i \cdot (w_a - w_b) = \log \frac{P(x|a)}{P(x|b)}$$

for glove,

$$J = \sum_{i,j=1}^V f(X_{ij}) \left(w_i^T \tilde{w}_j + b_i + \tilde{b}_j - \log X_{ij} \right)^2$$

- fast training
 - scalable to huge corpus
 - good performance even with small corpus & vectors
- bias
fn
each
word
- to column

good hyperparameters for glove:

→ dim 300

→ window size 8

NER

why it is hard?

- It can be difficult to know if something is an entity
- It can be difficult to know if words leading & preceding an known entity should be considered as entity as well.

different ways:

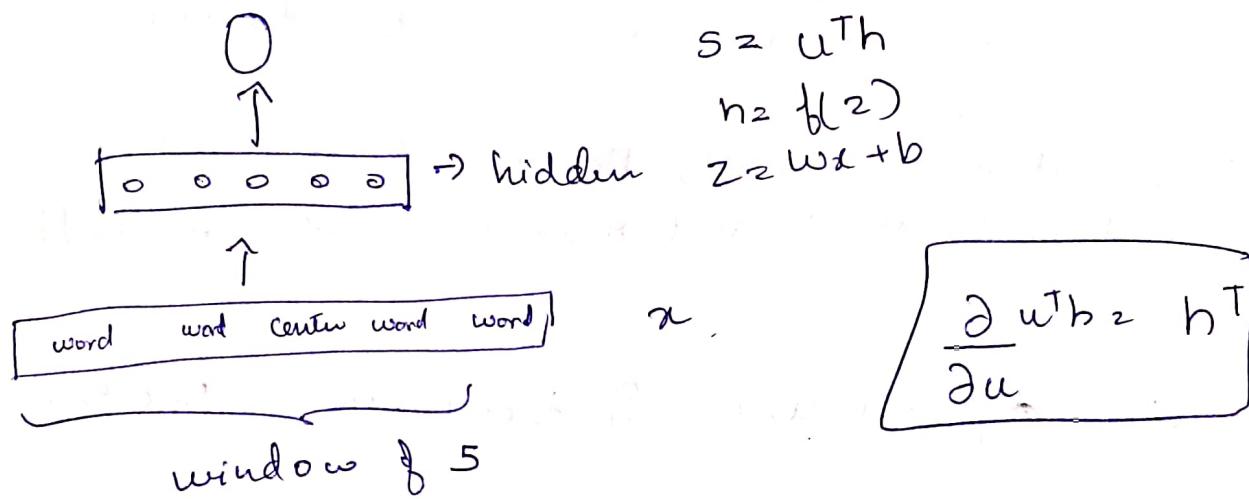
window classification softmax.

take a center word & $\pm k$ word, take word vectors of the $(2k+1)$ word window & classify

window, using score:

Same idea as before, but let the model get a

single no in the last layer, a high value in case the center word of window is a required entity, a low value if not, train it in such a way.



to backprop

$$\begin{aligned}
 s &= u^T h \\
 n &= f(z) \\
 z &= wx + b \\
 x
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial s}{\partial b} &= \frac{\partial s}{\partial h} \frac{\partial h}{\partial z} \frac{\partial z}{\partial b} \\
 &= h^T \text{diag}(f'(z)) I \\
 &= h^T \circ f'(z)
 \end{aligned}$$

$$\frac{\partial s}{\partial w} = \frac{\partial s}{\partial h} \frac{\partial h}{\partial z} \frac{\partial z}{\partial w}$$

$$= h^T \text{diag}(f'(z)) x^T$$

Disadvantage of retraining word vector

- Ex: Say TV, telly & television are in close space in original word vector representation. Our dataset has only TV & telly in it. Then retraining might ^{might} push TV & telly somewhere else losing its context to television.
- If dataset is small, don't train word vector
- If dataset is large (1M words+) then retrain/fine-tune.

← dependencies Σ treebanks (Lec 5) →

Language Modelling

- Predicting what word comes next.
- $P(x^{(t+1)} | x^{(t)}, x^{(t-1)}, x^{(t-2)}, \dots) \quad x^{(t+1)} \in V$.

Pre-deployment:

N-grams.

→ assume $x^{(t+1)}$ depends only on the preceding $n-1$ words

$$= P(x^{(t+1)} | x^{(t)}, \dots, x^{(t-n+2)})$$

$$= \frac{P(x^{(t+1)}, x^{(t)}, \dots)}{P(x^{(t)}, \dots)} \leftarrow \begin{array}{l} \text{prob of n gram} \\ (\text{conditional prob}) \end{array}$$

These probs are obtained by count of appearance.

3. 4-gram model.

- as the proctor started the clock, - the students opened their _____

$$3 \text{ } n-1 = 3$$

$$P(w | \text{students opened their})$$

$$\frac{\text{count}(\text{students opened their } w)}{\text{count}(\text{students opened their})}$$

- The answer could be exact if context was taken
- if the event has not happened, then it would add zero to denominator
- to remove this, add small δ to all
- if denom is zero, then back off - 4 gram \rightarrow 3 gram
 (opened their)
- increasing n increases model size
- can generate text.

NN band

→ fixed window.

window size > 4 .

the students open

$$\hat{y} = \text{softmax}(Uh + b)$$

$$h = f(Wx + b)$$

$e_i = [e^1 \ e^2 \ e^3 \ e^4]$ word embeddings

$$x^1 \ x^2 \ x^3 \ x^4$$

→ advantages

- No sparsity problem
- No need to store all observed n-grams

→ disadvantages

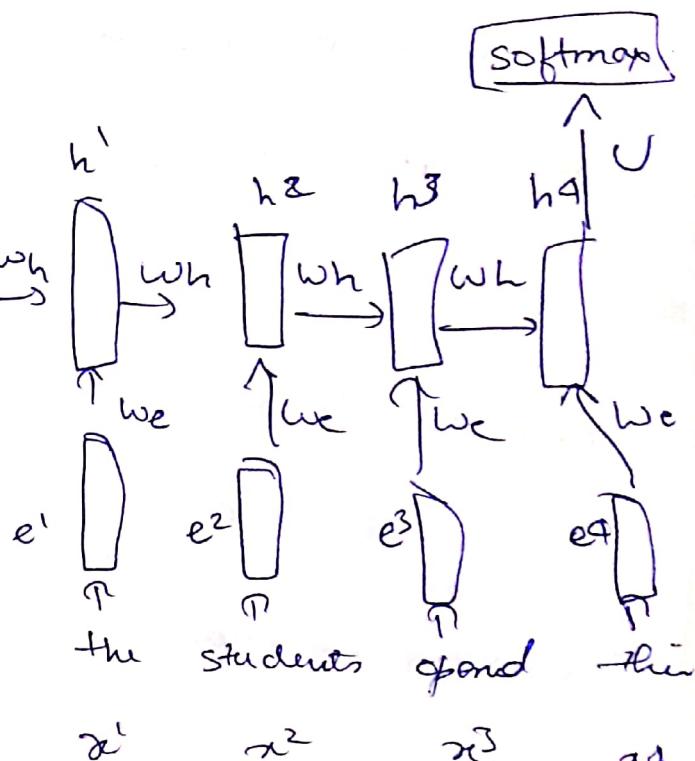
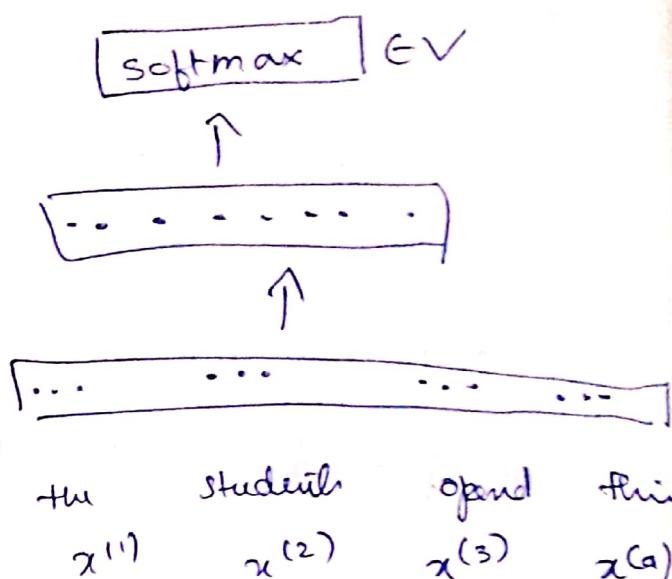
- window is small
- enlarging window increases weight W
- for each embedding input, the weight matrix doesn't share its weight

RNNs

$$e^{(t)} = Ex^{(t)}$$

$$h^{(t)} = (Wh^{(t-1)} + we^{(t)} + b)$$

$$y^{(t)} = \text{softmax}(Uh^{(t)} + b)$$



→ advantages:

- any length
- information transfer through many layers (Hierarchical)
- model size doesn't increase for longer seq.
- same weights on every timestep

→ disadvantages:

- pretty slow, can't do in parallel

standard evaluation metric → perplexity

Inverse prob of corpus, according to Model

$$\text{Perplexity} = \exp(\bar{\log P_{\text{CM}}})$$
$$= \frac{1}{T} \sum_{t=1}^T \left(\frac{1}{P_{\text{CM}}(x^{(t+1)} | x^{(t)} \dots x^{(1)})} \right)^{1/T}$$

lower perplexity is better

gradient vanishing

- due to small Δ values everywhere
- if largest eigenvalue of $W_n < 1$, then grad will shrink exponentially. (sigmoid is used)
- $W_n > 1$, then exploding grads.
- due to this only nearby effects are taken into consideration.

LSTM vs GRU

- GRU is faster & less parameters
- no other noticeable changes
- start with LSTM.

in BiDi RNN

Forward & Backward ops are concatenated.

Machine Translation

On the point!

find best English y for given french x .

$$\text{argmax}_y P(y|x)$$

$$\text{Bayes Rule: } \underset{\substack{\text{Modelling how} \\ \text{words \& phrases} \\ \text{should be translated}}}{\text{argmax}_y P(x|y) P(y)} \quad \underset{\substack{\text{Modelling how} \\ \text{to write} \\ \text{good english}}}{}$$

Modelling how words & phrases should be translated

Modelling how to write good english

$$P(x|y) \rightarrow P(x, a|y)$$

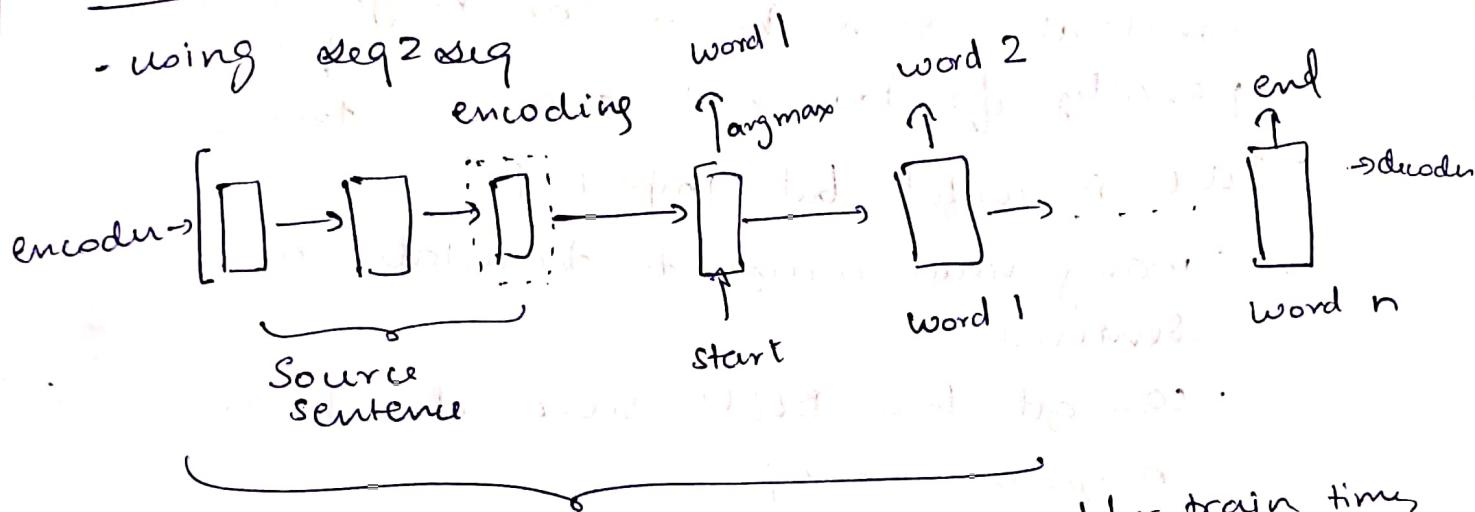
a is the alignment, word level correspondence b/w 2 languages.

- probability of a particular word aligning

- probability of particular word having particular futility (no of corresponding words)

the argmax is found by a heuristic search algo to search for best translation, disregarding hypotheses that are low-prob.
 → called decoding

Neural Machine Translation (nMT)



many tasks by seq to seq

- Summarisation
- Dialogue
- Parsing
- Code generation

greedy decoding → argmax on every decoder output

To fix this, Beam search is used

keep track of k most probable partial translation
 k usually $\in [5, 10]$

- final Beam search prob output

$$\frac{1}{t} \sum \log P(y_t | y_1, y_2, \dots, y_{t-1}, x)$$

evaluation:

BLEU

- similarity score based on n-gram precision
- penalty for too short system translations
- BLEU is useful but imperfect
 - many valid ways to translate a sentence
 - can get low BLEU score due to less n-gram overlap

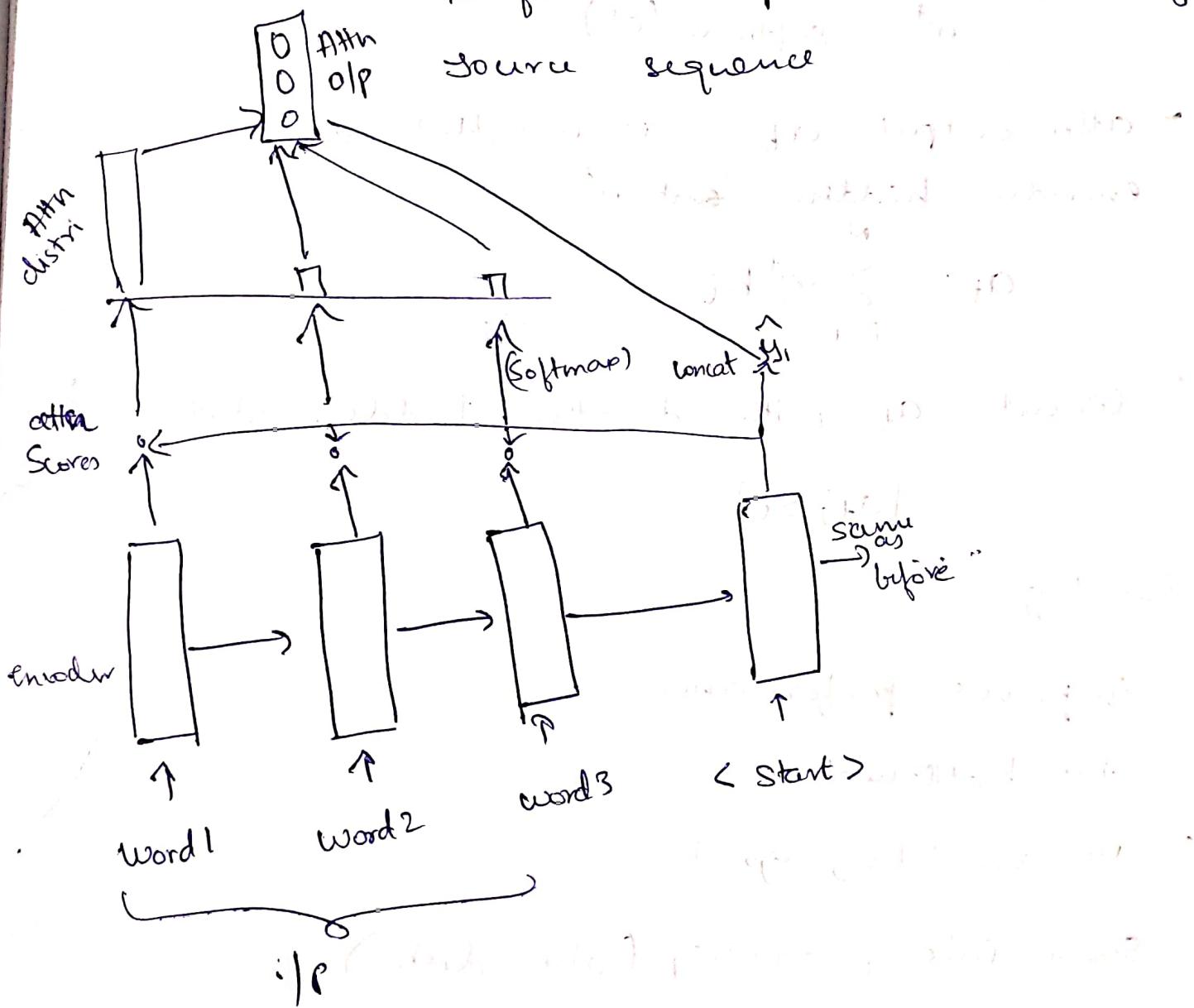
Difficulties

- out of vocab.
- Domain mismatch (train / test)
- context over longer text
- low-resource language pairs
- no common sense
- bias in training data (e.g.: gender)
- nonsensical :P can give weird outputs
(Somali → English agayay)

Why attention?

info bottleneck on decoder i/p

- all info needs to be captured by last encoding vector
- on each step of decoder, use direct connection to encoder to focus on particular part of source sequence



Attn equations

- encoder hidden $h_1, h_2, h_3 \dots h_n$
 - on timestep t , decoder hidden s_t
 - attn scores e^t ,
- $$e^t = [s_t^T h_1, s_t^T h_2, \dots, s_t^T h_n]$$
- Softmax to get attn distribution α^t
- $$\alpha^t = \text{softmax}(e^t)$$
- attn output a_t is weighted sum of encoder hidden $\xrightarrow{\text{sum}} \alpha^t$.
- $$a_t = \sum_{i=1}^N \alpha_i^t h_i$$
- concat a_t with decoder hidden state
- $$[a_t; s_t]$$

advantages

- improves performance
- no bottleneck
- no vanishing grad
- Some interpretability (attn distn)

Variants

- Multiplicative attn

$$e_i = s^T w h_i$$

w → learnable

- Additive

$$e_i = v^T \tanh(w_1 h_i + w_2 s)$$

~~w₁, w₂, v → learnable~~
v → weight

(d_3) → $d_3 \rightarrow$ hyperparameter

all learnable

$$w_1 \in \mathbb{R}^{d_3 \times d_1}$$

$$w_2 \in \mathbb{R}^{d_3 \times d_2}$$

$$\Rightarrow v \in \mathbb{R}^{d_3}$$

Data

look at linguistic Data Consortium.

Sn - Ans

gained traction with MCTest Corpus → 2013

CNN/DH, SQuAD → 2015/16

IBM Jeopardy → 2011 → ensemble of many models

SQuAD 1.0

- look examples

- passage, qn and ans

- ans must appear in passage (not yes or no)

- used Amazon Mechanical Turk

SQuAD 2

dw q, test set, half the ques has no answer

Stanford ATen Reader

for each qn: for represent using BiLSTM.
hidden state op.

for each passage, again with BiLSTM, then
add attention from the qn vector

$$d_i \rightarrow \text{softmax}(q^T w_s \hat{p}_i)$$

$q^T \rightarrow$ question dp vector

$\hat{p}_i \rightarrow$ passage hidden BiLSTM

$w_s \rightarrow$ learned weight

\rightarrow predict start token

$$d_i' = \text{softmax}(q^T w_s' \hat{p}_i)$$

\rightarrow predict end of answer

SAR++

get representation from all hidden units
for qn attention. Much more complex
model.

Vector rep of each token in passage

- word embedding

- POS & NE tags, one-hot encoded

- term frequency

- whether word appears in question (exact, uncased, lemma)

- Aligned word embeddings ('car' vs 'vehicle')
(sort of word similarity)

Conv for Text

1D CNN

	0	0	0	0
tentative	0.2	0.1	-0.3	0.4
dual	:	:	:	:
...	:	:	:	:
...	:	:	:	:
...	0.1	0.1	0.1	0.1
...	0.1	0.1	0.1	0.1
...	0.1	0.1	0.1	0.1
open	-0.4	-0.4	0.2	0.3
	0	0	0	0

4D word vector

Apply filter of size 3.

$\begin{bmatrix} 3 & 1 & 2 & -3 \end{bmatrix}$	$\begin{bmatrix} \phi, t, d \end{bmatrix}$	-0.6
$\begin{bmatrix} -1 & 2 & 1 & -3 \end{bmatrix}$	$\begin{bmatrix} t, d, \dots \end{bmatrix}$	-0.1
$\begin{bmatrix} 1 & 1 & -1 & 1 \end{bmatrix}$	$\begin{bmatrix} \dots \end{bmatrix}$	0
$\begin{bmatrix} \dots \end{bmatrix}$	$\begin{bmatrix} \phi \end{bmatrix}$	-0.5

use many filters to increase out dim

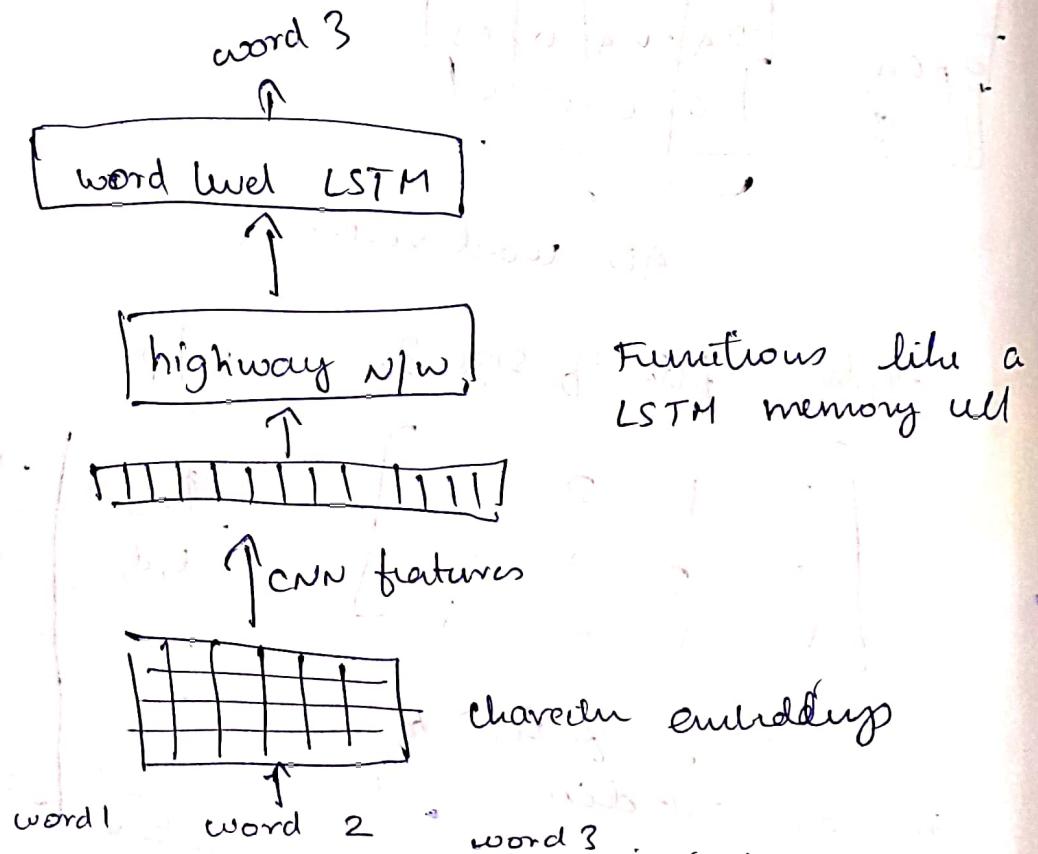
Quasi RNN

Combining CNN & RNN.

Wordpiece | Sentencepiece models

- variant of Byte pair encoding
- wordpiece model tokenizes individual words.
- Sentencepiece works on raw text
- BERT uses a variant of wordpiece, has common words & words built on wordpieces.

characte
hypatia = h ##yp ##ati ##a (in char
word-level models (Yoon et al 2015))



- character and model works for out of vocab words
- CNNs + Highway N/W over characters can extract rich semantic & structural information.

Foot Text

- further work on word2vec
- took word2vec skip gram by added character n-grams
- where = $\langle \text{wh}, \text{whe}, \text{her}, \text{ere}, \text{re} \rangle$, $\langle \text{where} \rangle$
- for center words in word2vec, all 6 of them are and

Contextual word embedding

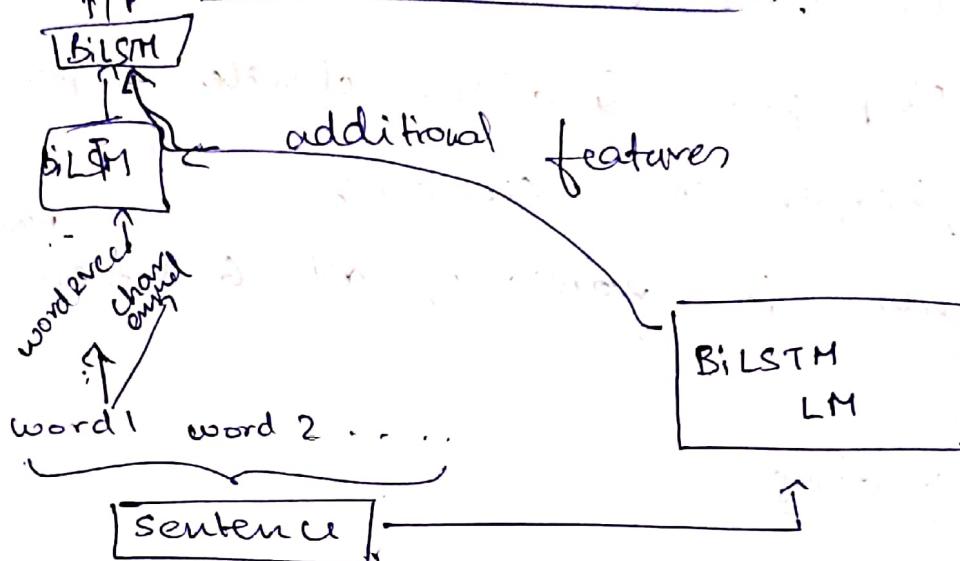
for unknown words:

- one idea is to train them as one vector ~~vector~~
- character embeddings as before
- if a $\langle \text{UNK} \rangle$ at test appears in embeddings, use it in test time
- for other words, assign a random vector, adding it to vocab
- collapsing things to word class, like unknown number, thing etc

disadvantage of word vectors

- same word can have diff. meaning
- they can have diff. parts of speech,
e.g.: a ~~to~~ noun or a verb (arrive, arrived)

TagLM (2017) - semi supervised for NER



ELMo (similar structure to TagLM)

- 2 Bi-LM (char(CNN → BiLSTM + ~~learned~~ pre-trained LSTM))
- Only character CNN for initial word Representations (2048 char n-gram & 2 highway layers, 512 dim projection)
- 4096 dim LSTM → 512 linear
- Residual
- parameter tying of token i/p & o/p
- learns task specific combination of BiLM (uses all hidden o/p)
- gave almost 3% improvement in a range of NLP tasks

ULM Fit

- Train on big unsupervised corpus.
- Fine tune ~~of~~ the model in domain you are working in
- Convert the LM to a text classifier
- less data required for better data (transfer learning)
- training \rightarrow 1 GPU day

other works

GPT 240 GPU days

Bert 320-560 TPU days

GPT-2 2048 TPU v3 days

Transformers

- \rightarrow The need to parallelize the training to make it faster, RNNs are sequential.
- \rightarrow Attention is all you need (original Transformer architecture)
 - A complex module Σ involving only attention, no recurrent
 - inputs: a query q , a set of key-value pairs to an output.
 - dot product attn is used to query key similarities
 - $$A = \sum \frac{c^{q \cdot r_i}}{\sum_j c^{q \cdot r_j}} v_i$$
 \rightarrow attn band weight on corresponding values

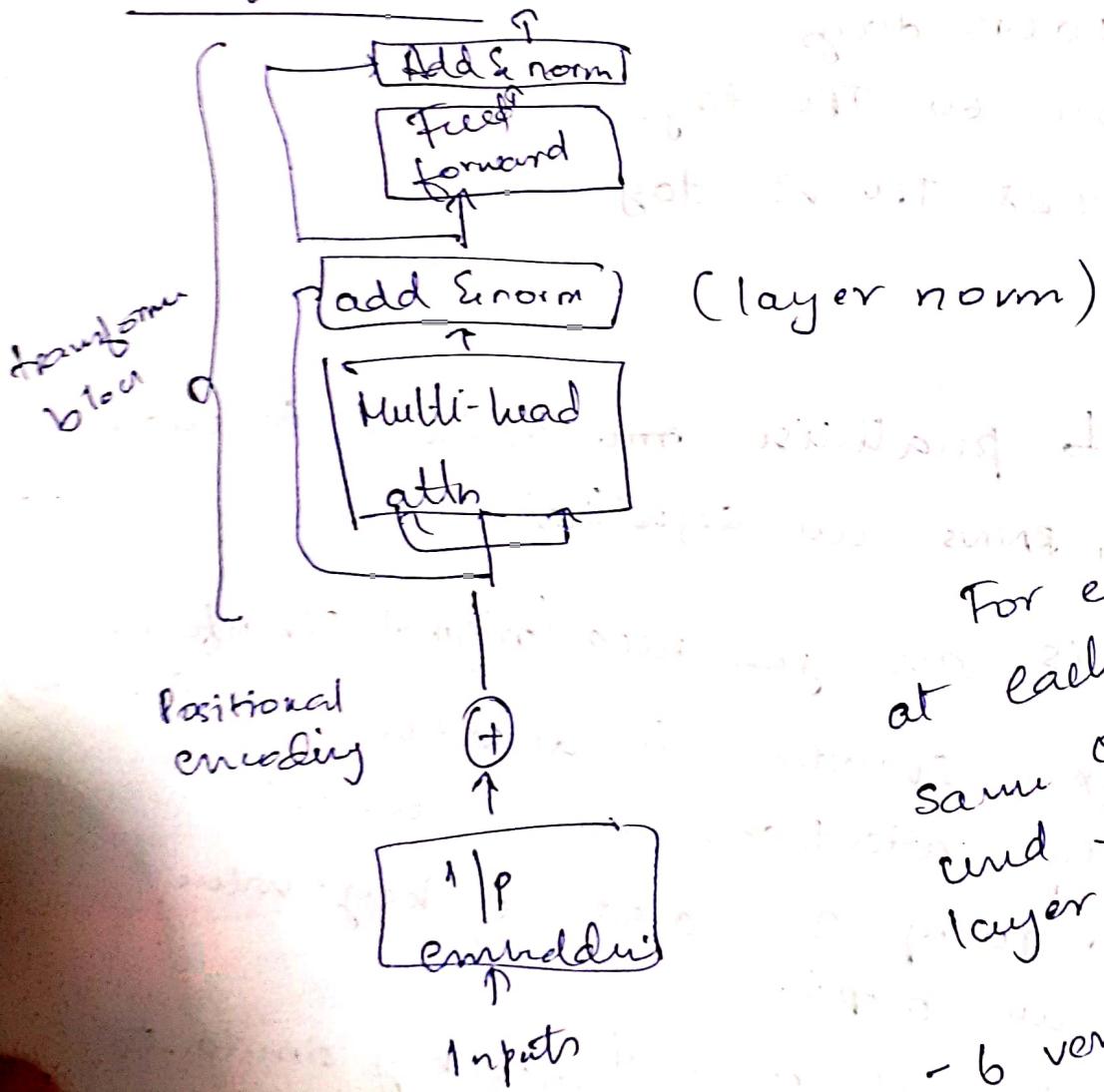
In encoder

- if words are queries, keys and values
- word vectors select themselves

Multihead attn

Because one attn network prioritizes over a single word, you may need to look for more dependencies

transformer block



For encoder,
at each block,
same Q, K, V
and from previous
layer.

- 6 vertical blocks

- Why is it so good?
- for pronouns, attn points to its reference
 - verbs points the change it does (sort of)

Bert

they wanted bidirectional context w/o further words being seen in later hidden layers

- Mask out words in a sentence, then predict the masked words
- they use mask 15% of words

- Next Sentence Prediction and aligned with first
 - Next Sentence Prediction against position embeddings
 - Pair of sentences, words are word pieces
 - Positional embeddings, segment embedding (whether word is from sentence 1 or 2)
- (Wikipedia + BookCorpus)
- Bert Base - 12 layers, 768 hidden, 12 head
- Bert Large - 24 layers, 1024 hidden, 16 head

4x2 @ 848 TPU for 4 days

Fine-tuning

- same as ULMfit.
 - Remove top level (mask level of next sentence)
 - Replace with own layer for the task.
- Bert large gave ~2% more than Bert-Base

- attn is permutation invariant
- Positional info added to maintain order.
- In decode side, it mimics a language model. This is done by imposing causality on positions you can look at & e.g. first position, it's illegal to look forward. Otherwise, it copies the i/p.

Multimod attn

- Diff Attn happens in parallel
- Intuition - diff things handle question (who? did what? To whom?)
- If models were ~ 3x faster than previous for MT task
- Residuals carry positional info to higher layer
- If also used for probabilistic image generation.

Music generation using Self-attn

- performance info is taken (not sound).
- initial motif is given & model is supposed to complete.
- Music Attn \rightarrow great stuff!!! (Magenta, Google AI)
- Relative Attn used for graphs

- Message Passing neural networks.

Natural language Generation

- Beam Search on Westword!
- Increases Beam search k size might decrease BLEU. (large k-beam search produces too short translations)
- In chit-chat (dialogues), small beam size might give non-relevant but on-topic output. High beam size might give safe, correct response, by generic words less relevant to the topic

Alternatives to Beam Search

- Sampling based decoding
- Pure Sampling
 - randomly sample to obtain next word
- Top-n sampling
 - few top n words to sample from that
- Both are more efficient than Beam search.

Softmax temperature

- temperature hyperparameter τ on softmax.

$$P_t(w) = \frac{\exp(s/\tau)}{\sum \exp(s/\tau)}$$

- Raising the ϵ makes Prob distribution more uniform.
- Thus more diverse O/P.
- lower ϵ , Prob distri becomes spiky.
- less diverse output

Summarisation

- Single document summarisation
- gigaword dataset
 - sentences \rightarrow headlines
(Sentence comprehension)
 - LCSTS
paragraph \rightarrow sentence summary
 - NYT, CNN/Dailymail
news article \rightarrow sentence summary.
 - Wikipedia
- full flow to article \rightarrow summary sentences.

Sentence Simplification

- Simple wikipedia
- Stand wiki sentence \rightarrow simple
- Newela
news article \rightarrow written for children

extractive summarisation

- Select parts to form summary.

abstractive summarisation

- > generate new text.

Pre-neural methods

- > content selection using tf-idf, graph based

Evaluation

- ROUGE
- based on Recall
- no brevity penalty
- ROUGE-1 : unigram overlap
- ROUGE-2 : bigram overlap

Neural-based

- DOLS Rush et.al
- Seq2seq are good at copying but bad at copying over details.
- Copy based models copy too much.
 - when document is long
 - bad at content selection
- bottom-up summarisation
 - neural seq tagging model for content selection
 - seq2seq attn won't attend to words tagged 'don't'

- dual Reinforced summarisation

→ ROUGE score can't be optimised during training because it's not differentiable
so RL is used for the same

Dialogue

- Assitive

band → Co-operative

- Adversarial

- Elit chat

- Therapy

Seq2seq band

- boring responses with no diversity

- Repetition

→ irrelevant (for context)

Sohm:

- optimising for Max mutual info b/w input S and response T

$$\frac{\log p(S, T)}{p(S)p(T)}$$

- Block repeating n-grams during beam search

Better evaluation of NLU

- Fluency (compute prob wrt well-trained LM)
- Correct style (prob wrt LM trained on target corpus)
- Diversity (rare word usage)
- Relevant to IIP (semantic similarity response)
- human quality qns

Coreference Resolution

- Identify entities that refer to the same real world entity.
- Which of the mentions refer to same entity
- Applications
 - Full text understanding
 - MT (for tracking pronouns)
 - Dialogue systems (carry reference across messages, (Tours Bond → Specie))
- Detecting mentions is easy
- Clustering references is hard

3 kind of mention

- Pronouns
- Named Entities
- Noun phrases

- one way to do this is to train a classifier specifically for mention detection instead of a POS tagger, NER system, parser

Coreference -

2 mentions refer to same entity

Barack Obama went to Obama ↗
both are same

Anaphora

term refers to another term

Obama said he went
↑ ↓
antecedent anaphor

→ one refutes another

Anaphora without preference -

No dancer' twisted her knee

Bridging anaphora

we went to see a concert. The tickets were

Cataphora

He was lying, smoking, as was his cushion, Lord Flenny

(need to look forward here)

- Hobbs Naive algo (1976) (tree based)

Winograd Schema

→ Proposed as an alternative to Turing Test (2013)

(if you can solve coreference, you've solved AI)

Mention pair models

- train a binary classifier that assigns every pair of mentions a probability of being coreferent.

Neural Coref. Models

- $i|p \rightarrow$ Candidate + Mention embeddings + additional features (statistical)



hidden layers + ReLU



Score

End-to-End (for SOTA) (2017)

- LSTM + attention

- words as $i|p \rightarrow$ embedding + CharCNN

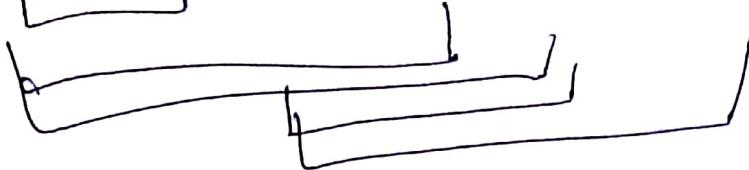


$i|p$ representation

Deep BiLSTM

↳ Rep for span (continuous subsequences of words)

general Electric said that they



→ all them are span

Deep BiLSTM



span representation in 3 parts

+ Attn weights on diff words

Score every pair of spans to see if
they are coherent mentions

$$s(i, j) = Sm(i) + Sm(j) + Sal(i, j)$$

is i a mention is j a mention Do they look
coherent

Score for $s()$ takes span ref as i / j

Clustering bound

- agglomerative clustering
- start with each mention as singleton cluster
- merge a pair of clusters at each step

Multi-task learning.

- decaNLP (Natural language Decathlon)
 - language modeling
 - an Aus (Australian) analysis
 - Summarisation
- use RL for NAS (multitask)
- ← TREE bank models from NLP →
↳ Treebank

← Bias in Al →

Using unlabeled data pre-training parallel

- Pre train encoder & decoder for diff languages (for MT)
- Use unlabeled data (in this case, En → another lang) for pre-training
- Self training
 - Pre-training won't have interaction b/w encoder & decoder
 - load data from one model & send it to En-De model. (It can be circular, so not much)

Then we can take the output of En-De and feed it back to the encoder

Back- Translation

- Do En-Fr', then do Fr-en ~~as well~~ target as the original English. (like an autoencoder)
- Then train it on training labeled data
- train 2 model on labeled data, then do back-translation on large corpus, do it over and over again

Unsupervised word translation

- Cross lingual word embedding
- Shared embedding for both languages
- Assumption is that language structure is same b/w the languages
- Run word2vec on monolingual corpora, getting word embeddings $\mathbf{x} \in \mathcal{Y}$
- Learn an orthogonal matrix \mathbf{W} such that $\mathbf{Wx} \approx \mathbf{y}$
- One method to train \mathbf{W} is called adversarial training. A Discriminator will predict if a word is English or French. Train \mathbf{W} so that Discriminator will be confused. (It can't predict accurately)

- same enc. dec. and fn. for both languages
- special token is used to diff. langug.

GPT-2

- just a really big TB LH
- 40GB text iIP
- 10M webpages from reddit links with high karma
- zero shot learning
 - generating from a prompt
 - Reading Compre...
iIP <context> <qn> A:
 - Summarisation
<article> TL;DR:
 - Translation
<sent1> ; <sent2>
 - <sent> ^
 - Qn Ans
<qn> A:
 - It can do translation bcause other lang. are in training corpora if a very small amount

graph neural network

graph \rightarrow vertices / nodes and edges

$$G = (V, E)$$

edges \rightarrow unidirectional

or

bidirectional

node classification in GNN. (original paper)

every node \rightarrow annotated with a label, predict label w/o ground truth.

each node v is characterised by its feature x_v
annotated with a ground truth \underline{t}_v .

Given a partially labelled graph, then labelled nodes
are used to predict the unlabeled.

learns to represent each node with a d-dim
vector \underline{h}_v

$$\underline{h}_v = f(x_v, x_{co[v]}, h_{ne[v]}, x_{ne[v]})$$

$x_{co[v]}$ \rightarrow features of edges connecting v

$h_{ne[v]}$ \rightarrow embeddings of neighbouring nodes of v

$x_{ne[v]}$ \rightarrow features of neighbour nodes of v

$f \rightarrow$ puts it into a d-dim space

This can be solved by message passing / neighbourhood aggregation

Banach fixed point theorem can be applied

$$H^{t+1} = F(H^t, x)$$

$H, x \rightarrow$ concat of h, x .

gmn op =

$$o_v = g(h_v, x_v)$$

$$\text{loss} = \sum (t_i - o_i)$$

Limitations of this:

- if assumption of fixed point is relaxed, NN can learn more stable representations, and can remove iterative update process.
- can't prun edge info
- fixed point can discourage the diversification of mode distribution, and thus may not be suitable for learning to represent modes.

GNN Survey - Zonghan et al.

Conv GNNs → spectral
→ spatial

Network embedding - represent nodes as low-D vectors preserving nw topology & node info → use this to perform classification, clustering, recommendation using simple algos

Graph (V, E)
 $V \rightarrow$ nodes
 $E \rightarrow$ edges

$e_{ij} = (v_i, v_j)$ edge v_j to v_i

$N(v) = \{u \in V | (v, u) \in E\}$ neighbourhood of node v

A graph may have node attributes x

$x \in \mathbb{R}^{n \times d}$ → node feature matrix $x \in \mathbb{R}^{n \times d}$

A graph may have edge attributes x_e

$x_e \in \mathbb{R}^{m \times c}$ → edge feature matrix

$$x_{v,u}^e \in \mathbb{R}^c$$

↓

flat vector of edge $v-u$

A graph is undirected if & only if the adjacency matrix is symmetric.

Spatio-Temporal graph:

attributed graph where node attributes change dynamically over time. $G^{(t)} = (V, E, X^{(t)}) \quad X^{(t)} \in \mathbb{R}^{n \times d}$

ConvGNN

generalise convolution from grid data to graph data. generate a node v 's repre by aggregating its features x_v , neighbour features x_u . ConvGNNs stack multiple graph conv layers to extract high-level node repre.

graph autoencoders

encodes nodes/graphs into a latent vector space, reconstruct data from that. It is used for n/w embeddings, graph generative distributions.

with graph structure & node content info as inputs, ops of GNNs can focus on diff tasks

- node level
 - node regression, classification
- edge level
 - edge classification, link prediction

with two nodes hidden vector ^{cos from}
GNN as input, similarity fn σ nn can
be utilized to predict label / connection
strength of an edge.

- graph level

graph classification tasks. To get a
compact representation, GNNs can be
combined with pooling, read-out ops.

Data preprocessing:

- Mean subtraction

data is around origin in each dim
(done separately across each channel)

- Normalisation

scaling all data to a single scale

- mean + std

- min max

↳ use if diff. info. are at
diff. scales but needs to
equal importance.

- PCA

if many features are dependent

for e.g. with NB.

and for small dataset

- whitening

- take data in eigenbasis & divide
every dim by EV.

- linear classifier learns only 1 example
for each class

sum from 1 hinge loss

	Class 1	Class 2
Class 1	3.2	1.3
Class 2	5.1	4.9
Class 3	-1.7	2

$$1 \quad L = \max(0, 5.1 - 3.2 + 1) \\ + \max(0, -1.7 - 3.2 + 1) \\ = 2.9$$

$$2 \quad L = \max(0, 1.3 - 4.9 + 1) \\ + \max(0, 2 - 4.9 + 1)$$

$$\Rightarrow 0$$

properly classified

softmax

- prob dist of scores of the classes.
- map them to target distribution
- can be viewed as KL div or MLE
- mathematically it is easier to maximise log than actual prob
- negative case log goes high n high and we need to minimise

	denom			
Class 1	3.2	\rightarrow	24.5 164	normalise \xrightarrow{sum} 0.13 0.87
Class 2	5.1			0
Class 3	-1.7			

$$-\log(\text{true class}) = -\log(0.13) < 0.87$$

BN
have unit gaussians at every forward pass. per batch.

- get mean Σ , var for each dim.
- usually inserted after conv Σ before non-linearity.
- can recover identity matching.
- improves grad flow
- Allows high LR
 - Reduces dependence on initialisation

Momentum

$$x_{t+1} = x_t - \alpha \nabla f(x_t) \quad \text{SND}.$$

$$v_{t+1} = \beta v_t + \nabla f(x_t)$$

$$x_{t+1} = x_t - \alpha v_{t+1} \quad \text{rho (HP)} \Rightarrow \text{friction}$$

Nesterov M

- look at grad at where vel is computed & combine it with present grad

$$v_{t+1/2} = \beta v_t - \alpha \nabla f(x_t + \beta v_t)$$

$$x_{t+1} = x_t + v_{t+1/2}$$

Adagrad

keep running sum of squared grads

divide by this while updating

- accelerates movement among slow dim

ϵ_1 slows down on fast dim

RMSprop

decay the squared estimate

Adam

combines Momentum & Adagrad / RMSprop

Alexnet (8 layers)

- shared across 2 gpus

- ReLU, SGD+M (BS 128)

- Layer norm

ZENET

- improving HP

UbiNet

- 16 - 19 layers
- 3×3 conv. 1, 1 \rightarrow ~~base~~
 2×2 MP, 2, 0
- 3 3×3 conv has 7×7 receptive field
- 138M param (tot in fc)
- no local norm.

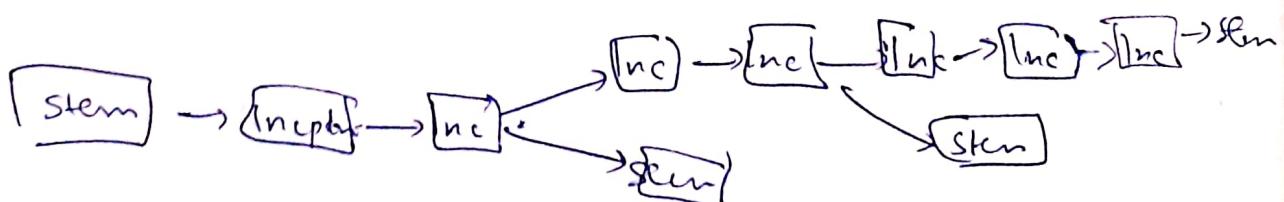
googLeNet

- 22 layer
- Efficient inception module
- no fc (only global avg pool)
- $\text{concat}((1 \times 1) + (1 \times 1, 3 \times 3) + (1 \times 1, 5 \times 5), (3 \times 3, 1 \times 1))$
 \downarrow
 good local n/w topology

Naive Inception

$$\text{concat}(1 \times 1 + 3 \times 3 + 5 \times 5 + 3 \times 3 \text{MP})$$

1×1 bottleneck reduces params by $\frac{854}{38}$

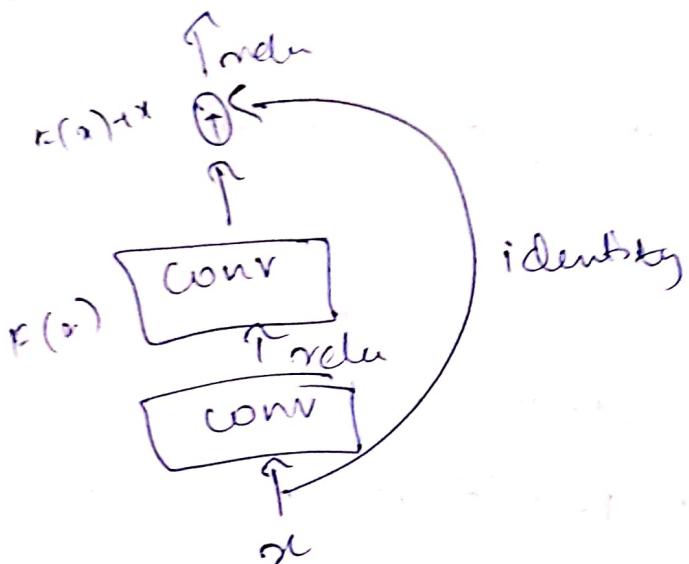


Stem \rightarrow auxiliary class. op.

stems become more gradient invariant.

ResNet

- 152 layers (in Imagenet)
- Residual connections



- super models are harder to optimise.

generative Models

$$p(x) = \prod p(x_i | x_1, x_2, \dots, x_{i-1})$$

$\oplus x \rightarrow$ pixel.

↓
Likelihood of
image x .

Pixel RNN

- Generate Pixel starting from a corner.
- Sequential gen is slow

Pixel CNN

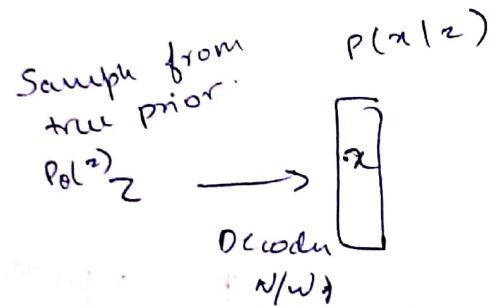
- CNN over a context region around a particular pixel.

VAE

density fn with a latent variable z

$$P_0(x) = \int P_0(z) P_0(x|z) dz$$

x is generated from z .



choose gaussian prior.

represent $p(z|x)$ using a NN

But integral is intractable.

In addition to decode, define encode $q(z|x)$ that approx $P(z|x)$

$\text{loss} = \text{Reconstruction loss} + \text{KL Divergence}$

backprop can't go down on the latent space
(which is a distribution)

2) M+ COE

(2) ~~or~~ Omega in Eq 1,

no need to run backprop on ϵ .