

LIS-RNNs pt 3- watch S. 2019 videoreview

- long-term dependence \rightarrow recurrent structures
- short-term dependence \rightarrow iterated structures
- OR: introducing self-reference (!) - are there deeper implications here.

- investigate

- RNNs trained by minimising divergence between sequence of outputs and desired outputs. (not between instantaneous inputs and outputs)

(*) Primary topic:-

- How to backprop with divergences over sequences

(*) Story so far - stability

- RNNs can be unstable
- Have to deal with issues of 'blowing up'; but also 'saturation'
- tanh() more reasonable.

(*) vanishing gradients

- An issue for deep networks in general
- Backprop \rightarrow gradients vanish/explode.
- e.g. 1st layer MLP
 - see magnitude of gradient w.r.t. params of network (near beginning of network)

(*) Propagation of info to subsequent layers.

- less than required no. of neurons/units to 'capture' pattern complexity; layers downstream cannot receive this 'missed info'.
- for deeper layers to be able to recover patterns; want info to be 'pushed' downstream (i.e. by not using thresholds) and activations which 'carry info' about distance from boundary.

(*) We see this again \rightarrow more informative activations might be ReLU or soft ReLU.

- (*) Other activations e.g. tanh, sigmoid have saturation behaviour, beyond specific above from boundary, 'info about where you are' is 'lost'. (*) too inflexible.

(*) Even with tanh activations; gradients can vanish after only a few layers of backprop

(i) How is this remedied?

(*) Long term dependency problem

- network 'forgets' in a manner that depends primarily on weights rather than inputs.
- we really want a design where system 'retains memory' until it no longer has to.
- sentence dependency example - want behaviour that is triggered not only on what is 'remembered', but also on what next expected input or output must be.
- i.e. input-triggered response rather than state-triggered response

↳ Are LSTM

(*) Long Short Term Memory

- still mathematically a complex activation that suits our design criteria
- GRU: Addresses problem of input dependent memory behaviour.

(*) LSTM-based architectures

- identical to RNNArch.
- each of 'boxes' is a layer (many units) - remember not. conventions!

(*) Bi-directional LSTM

- if you agree LSTM is a variant on RNNs; can have things like bidirectional LSTM.

(*) Key issue

- 1) How do we define divergence?
- 2) How do we compute outputs?

(*) Today - focus on unidirectional structures; and can be generalised easily to bidirectional.

- focus on:-

1) Network architectures

2) Synchrony:-

i) Time synchronous.

ii) Order synchronous

(*) variants on RNNs. (A3) - see slides; acquaint with architecture

- one-to-one:

- many-to-many:

- many-to-one:

- order synchronous: - get a sequence of outputs in response to sequence of inputs
(time synchronous). - order correspondence between input and output.
- 1-to-1 correspondence between input and output
not present.

- many-to-many: generate output sequence after processing input

- one-to-many: e.g. merge captioning.

(*) variants on RNNs / seq MLP

1-to-1

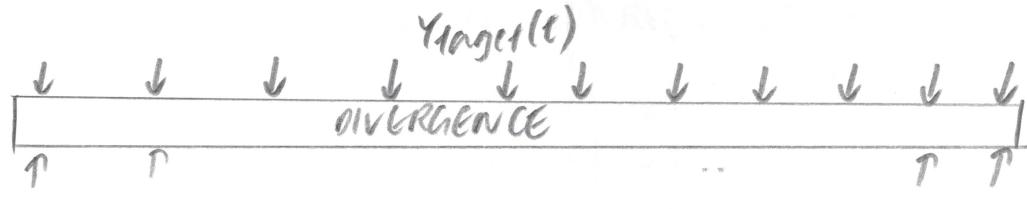


- no recurrence in model
- one input \rightarrow one output; unique
- why is this included in RNNs?
- we can still view divergence over sequence of outputs
rather than as individual divergences over outputs.

- OR: Amit he is trying to make it to naturalise the view of divergence
over a sequence of outputs; and consider the decomposition of
divergence over a sequence into that of over individual outputs as
a special case of a more general phenomenon.

- weighted sums.

(*)



(*) - more general case.

- gradient backprop at each time

$$\nabla_{Y(t)} \text{Div}(\text{Ytarget}(1, \dots, T), Y(1, \dots, T))$$

(previous)

$$-\text{Assumption: } \text{Div}(\text{Ytarget}(1, \dots, T), Y(1, \dots, T)) = \sum_t w_t \text{Div}(\text{Ytarget}(t), Y(t))$$

(simplifying)

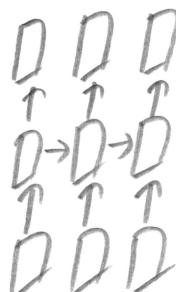
$$\text{as an MLP. } \nabla_{Y(t)} \text{Div}(\text{Ytarget}(1, \dots, T), Y(1, \dots, T)) = w_t \nabla_{Y(t)} \text{Div}(\text{Ytarget}(t), Y(t))$$

- w_t typically set to 1.0.

$$(?) \text{ Surely } \sum_t w_t \nabla_{Y(t)} \text{Div}(\text{Ytarget}(t), Y(t)) ?$$

(#) variants on RNNs (the synchronous net)

many-to-many



- below into modelling language
- network - one output for each input. ; one-to-one words.
- POS tagging
- In practice bidirectional → future informs about current time.
- assignment of one tag depends on others ('current rel')

- unidirectional time synchronous - process input left to right, produce output

- Bidirect. ————— " ————— process left to right (forward net)

↓
usual in practice
(makes more sense)

process right to left (backward net)
combine hidden outputs to prod. one out. per input symbol

(*) How do we handle input?

- BPTT.

(*) Given a collection of training instances comprising input sequences and output sequences of equal length, with one-to-one correspond.

- (X_i, D_i) where $X_i = (X_{i,0}, X_{i,1}, \dots, X_{i,T})$
 $D_i = (D_{i,0}, D_{i,1}, \dots, D_{i,T})$

(*) Input - sequence of vectors also
Output - sequence of vectors also.

(*) Backprop through time

(*) (**) Divergence is computed between network sequence of outputs and desired

(**) Not sum of divergences at individual terms; unless defined that way

(*) Compute $\sum_t \text{DIV}$ for all t (1st step)

(*) Definition of the divergence DIV is long component

(*) Can become complex

(*) Time-synchronous recurrence

(*) Simplifying ass: - sequence divergence is simple additive sum of divergences over individual time instances

equations: see earlier : global divergence as sum of local.

(*) Typical div. for class is cross-entropy

(*) Simple recurrent model - text modelling

- generic model that can predict next character from a seq. of characters or words

(*) use one-hot-encoding

- softmax outputs class probabilities

e.g. $\frac{\prod e}{\prod n} \frac{\prod h}{\prod n} \frac{\prod t}{\prod n} \frac{\prod o}{\prod n} \rightarrow \frac{\prod e}{\prod n}$

'n' $\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$ 'e' $\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$

input - one-hot vectors
(embeddings?)

'hello' output - probabilities
over characters

- peak at target char.

$$\begin{pmatrix} e \\ h \\ l \\ o \end{pmatrix}$$

(*) Training

- input: symbols as one-hot vectors - dimensionality is vocab size

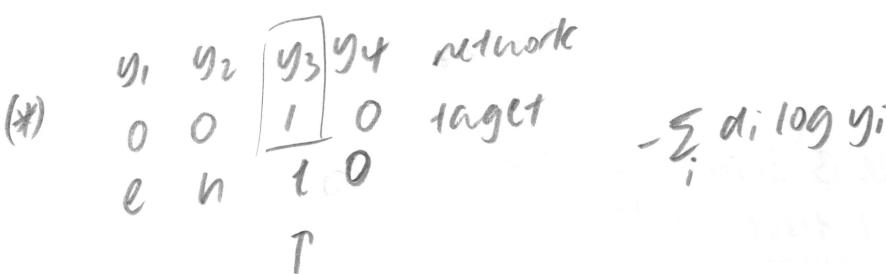
- output: probability distn over symbols

$$Y(t, i) = p(v_i | w_0, w_1, \dots, w_{t-1}) \quad v_i - i^{\text{th}} \text{ symbol in vocabulary}$$

- divergence: $\text{D}_{\text{IV}}(Y_{\text{target}}(1, \dots, T), Y(1, \dots, T)) = \sum_t \text{CE}(Y_{\text{target}}(t), Y(t))$

$$= - \sum_t \log Y(t, w_{t+1})$$

↳ prob. assigned to next
correct word



$$\text{CE} = -(0 \log y_1 + 0 \log y_2 + \log y_3 + 0 \log y_4) = -\log y_3$$

(*) cross-entropy is -ve log probability that has been assigned to target output

(*) overall divergence is sum of this term over all time instances

(*) Brief detour - language models

(*) Language modelling using RNNs

OR: come up with modelling structures for prediction of logo.

one-hot-encoding: re-specify a vocabulary of N words in a fixed lexical order.

- distinction between representational granularity e.g. 100,000 English words
vs 100 characters in English

(*) Predicting words

$$w_n = f(w_1, \dots, w_{n-1}) \quad w_1, \dots, w_{n-1} \in \mathbb{R}^N$$

$\underbrace{w \times 1}_{(n \times 1)}$ - one-hot vectors

(*) Dimensionality of one-hot vector if using 100,000 English words.

$$\in \mathbb{R}^{100,000}$$

- inputs $w_1, \dots, w_{100,000}$ are very high-dim, sparse

(*) review this dimensionality

a little

- need a more detail/longer

(*) unit cube $2^{100}, 2^{100,000}$ corners

- occupying 100 chars, 100,000 words

OR: what vol. of space is entire collect. of inputs occupying? O (?)

(*) One-hot representation

- uses N corners of 2^N corners of unit cube

- density of points $O\left(\frac{N}{r^m}\right)$

- inefficient use of dimensions

(*) BR: Why use one-hot if it uses dimensions inefficiently?

(*) - semantic agnostic rep.

- not assigning rel. importance to words/going on symbols

- ensuring that distance between 2 words is same regardless of what words.

- no imposition of closeness/distance on language

→ (**) - develop with reading

what distance?

- no assumption about relationship between words or relative importance of words; usage of space is inefficient.

(*) Solution to dimensionality

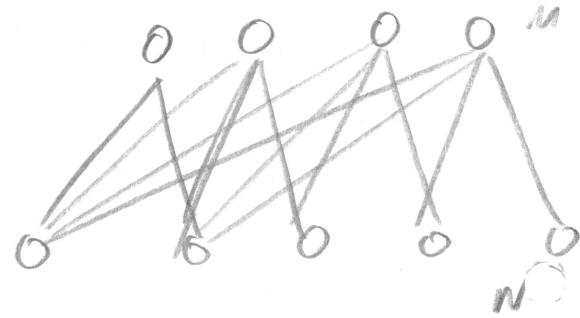
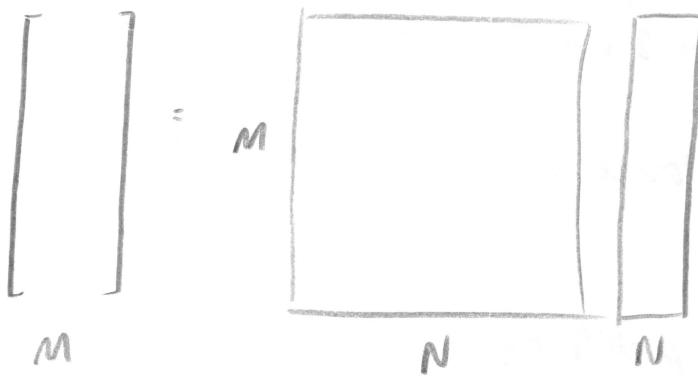
- Project to lower dimensional subspace: $O\left(\frac{N}{r^m}\right)$

- volume property not improving; but density is improving by many orders of mag.

(*) But refrained q: - what is projection/hyperplane such that positions/distances between points 'makes sense'.

- (*) learn/estimate the plane
- If done so 'properly', distances between projected points will capture semantic relations between words.
- (*) projected word vectors
- (*) project one-hot word vectors $\underline{w}_i \in \mathbb{R}^N$ into a lower-dim. subspace;
see slides \rightarrow project corners.
- (*) conduct following trans.
- projection matrix: $P \in \mathbb{R}^{M \times N}$
 - one hot vector: $\underline{w}_i \in \mathbb{R}^N$
 - projected word vector: $P\underline{w}_i \in \mathbb{R}^M \rightarrow$ use $w_n = f(P\underline{w}_1, P\underline{w}_2, P\underline{w}_3, \dots, P\underline{w}_{n-1})$
 - estimate/learn P using an appropriate objective

(*) 'Projection'



- (*) :- this simple linear transformation P can be placed within an NN context.

- implemented as a layer of M neurons with linear activations
- transforms that apply to individual inputs are all M -neuron linear activ. subnets with tied weights (as P is same for all word-vectors)

- (*) Is projection matrix P same \forall word vectors \underline{w}_i ?

(*) Predicting words - TDNN

- predict each word based on past N words - see architecture
 - neural probabilistic language model (Bengio et al. 2003)
 - hidden layer has tanh() activation, softmax
- (*) Note the projection unit
- estimates linear low-dim representation PW of words

(*) Alternative models

- (*) soft BOW - predict word based on words in immediate context
- no consideration of specific position

- (*) skip-grams - predict adjacent words based on current word.

(*) Embeddings: example

- Mikolov et al (2013)
- some illustrates mapping from country \rightarrow capital city capturing "is the capital of".
- cherry picked.

(*) Generating language: The Model

- use all previous words to generate next word
- $$w_n = f(PW_1, \dots, PW_{n-1}) \quad w_{10} = f(PW_1, \dots, PW_9) \text{ etc.}$$

(*) Generating language: Synthesis

- will get a probability distribution over words in 4th position
- (*) - illustrates generation of text. \sim - prediction
- select most likely word \tilde{w}_4
 - use that as additional 'input'; then predict $\tilde{w}_5 = f(PW_1, \dots, PW_4)$ etc.
 - practically; could go on indefinitely. But accuracy?
- ① do we feed raw probability vectors or post-process before feeding into next RNN step?

(*) Variants on RNN (sequence class.)

e.g. speech recognition

IN: sequence of feature vectors (mel spectra)

OUT: phoneme ID at end of sequence

- return out prob vector

- N = no. of phonemes

(*) in this setting we have other outputs, just ignored, consider output at final time.

(*) Inference: Forward Pass

- review slides (A6) suff.

- $\text{DIV}(\mathbf{Y}_{\text{target}}, \mathbf{Y}) = \text{CE}(\mathbf{Y}(T), \text{phoneme})$

- divergence only defined at final input.

(*) Training

(*) need to exploit untagged inputs

- task details determine / guide definition of divergence

$$\text{DIV}(\mathbf{Y}_{\text{target}}, \mathbf{Y}) = \sum_t w_t \text{CE}(\mathbf{Y}(t), \text{phoneme})$$

- speech/phoneme recog \rightarrow equal weights w_t

- question answering \rightarrow w_t high, further back weights low

(*) Variants on RNNs (order synchronous, time asynchronous seq-to-seq)

- exact location
of output unknown a priori