

Recurrent Neural Network
↳ repeating

MLP $\xrightarrow{\text{vector}}$ MLP $\rightarrow \gamma$

CNN $\xrightarrow{\text{image}}$ CNN $\rightarrow \gamma$

Both of them do not take
advantage of sequences.

[This smartphone is fit a battery is good]
OR
[battery good this is phone smart as is the] } both are same for
MLP using word
vector

① word

② Time series \rightarrow short snippet

③ Translation

④ Audio file \rightarrow is a sequence

⑤ Image \rightarrow caption

just to make it
easier for information
search.

TDR; if output depends on
sequence, we need some form
of model

* Let's see an example of words first.

$x_i \rightarrow \text{google}_{x_1} \text{ will } \text{ live }_{x_2} \text{ me }_{x_3} \text{ in }_{x_4} \text{ 2021}_{x_5}$

If we have one hot encoding for 10k words,
each word vector is 10k dim

So if we have 5 words, we sort them if
we preserve order.

↳

The input is variable anyway - how do
you handle that?

Conclusion → This didn't end up working

Large word corpora & large sequences
will always be a challenge for NNs.

Review Sentiment Analysis

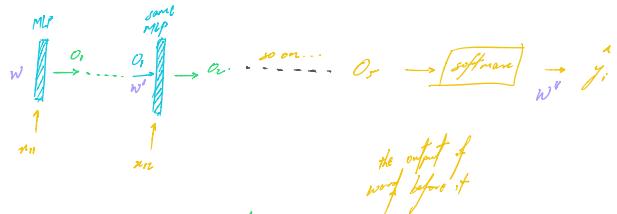
$$\text{Dataset}(D) = \{x_i, y_i\}$$

↳ binary
surface

↓
will have
varying words
obviously

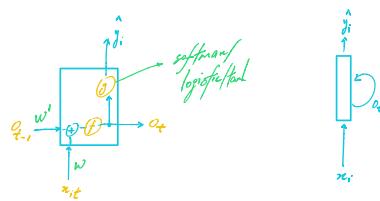
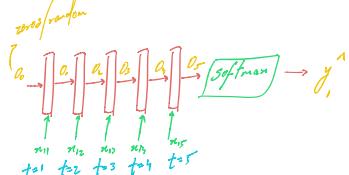
Take words → {0, 1}

$x_i \rightarrow (x_{i1}, x_{i2}, x_{i3}, x_{i4}, x_{i5})$

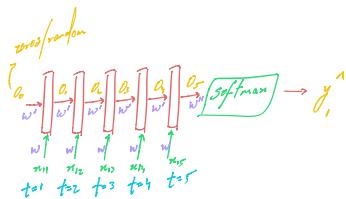


3 set of we want o_r to depend on x_{i1} & o_1
weights

$w \rightarrow$ current input
 $w' \rightarrow$ previous output \rightarrow current row
 $w'' \rightarrow$ in the end



* Backprop over Time



$$o_1 = f(wx_{i1} + w'o_0)$$

$$o_2 = f(wx_{i2} + w'o_1)$$

⋮

$$o_r = f(wx_{ir} + w'o_{r-1})$$

$$\hat{y}_i = g(w''o_r)$$

$$\frac{\partial L}{\partial w''} = \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial w''}$$

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial o_r} \cdot \frac{\partial o_r}{\partial w}$$

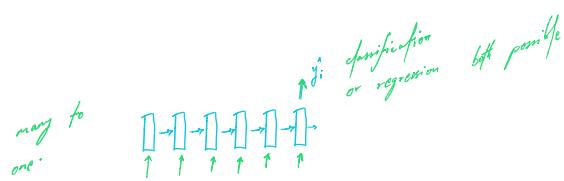
↑
on the last
time

$$\frac{\partial L}{\partial w'} = \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial o_r} \cdot \frac{\partial o_r}{\partial w'}$$

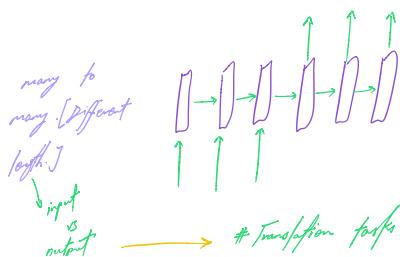
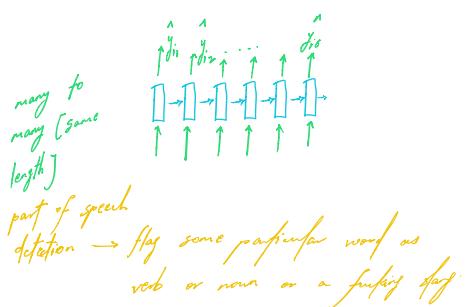
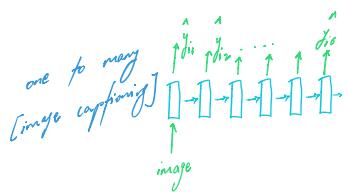
This backprop happens over time, on
the same weights.

The longer the sequence, the longer the
backprop gradient. \therefore Vanishing/Exploding
gradient problem will persist.

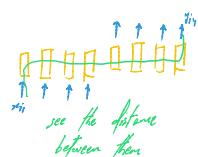
LSTM / GRU
avoid these problems
through modifications to
RNN



(e.g. review
reviews, etc.)



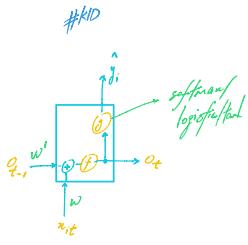
- # Single RNNs such
- ① Vanishing gradients occur
both forward & backward
passes
(if first & last word are
reflected to each other, now will
not effectively affect that.)
- # RNNs stuck at long term dependencies
in input sequences



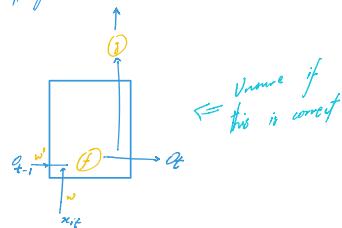
- # Transposition tasks

↳ we want the network
to have complete input
before outputting something

LSTM



forget /



$$o_t = f([w^i, w], [x_{it}, o_{t-1}])$$

concatenate (not add (the kids))
& dot product

forget by is
great for LSTM

GRU

- optimized
- less params to train
- faster training
- almost same performances as LSTM.
- apply special cases