

# Zhang-Recurrent Neural Nets

1

$p_{prev} = \cdot$

MCP = tabular data

CNN = vector/matrix data

but many tasks require sequential data

such as image caps, speech, time series

RNNs are deep learning models that capture the dynamics of sequences via recurrent connections

Recurrent RNNs are unrolled across time steps

## 9.1 Working w/ Sequences

Some datasets consist of a single long seq

We might randomly subsample the dataset into some pre-determined length

often data arrives as collection of seqs

Previously, we assumed w/ individual inputs we sampled from the sample underlying distribution independently

$P(x)$



with sequences, we still assume they are indep but cannot assume that data @ each time step is indep of each other

e.g. words @ end of doc are heavily dependant on words @ the start

## 9.2 Converting Raw text into Sequence data

### tokenisation

Skipped as not relevant for task

## 9.3 language models

## 9.4 Recurrent Neural Networks

we want to predict the next point using many previous points

but modeling  $P(x_t | x_{t-1}, \dots, x_{t-n+1})$  soon becomes infeasible

so it is better to use a latent variable model

$$P(x_t | h_{t-1})$$



$h_{t-1}$  is a hidden state that stores the sequence of info upto step  $t-1$

@ any time step  $h_t$  can be computed using the input  $x_t$  & hidden state  $h_{t-1}$

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

hidden states are technically inputs

RNNs are NNS w/ hidden states

$$H_t = \phi(x_t W_{xn} + H_{t-1} W_{hn} + b_h)$$

The calc of the hidden layer output of the current time step is determined by the input of the current time step + the hidden layer output of the prev time step

hidden layer output is called hidden state

Params of RNN include

- weight of current input
- weights of hidden state

+ bias of current input

At all time steps of sequence, RNNs use the same model params

Params cost of RNN does not grow w/ number of time steps



At any time step, the computation of the new hidden state can be seen as a concat of the input  $x_t$  & the hidden state  $t-1$

this concat feeds into a fully con layer

the output of the fully connected layer is  $h_t$

$h_t$  is later used to calc  $h_{t+1}$

### 9.5 RNN from scratch

No notes

### 9.6 RNN concise implementation

No notes

### 9.7. Backprop through time

Applying backprop in RNN is called through time

Exploding gradients comes from backprop across long sequences

BP through time keeps unrolling the comp graph of a RNN one step @ a time



An unrolled ~~the~~ RNN is essentially a feedforward network

↳ w/ a special property that the same parameters repeat throughout the network

Chain rule can then be applied to the unrolled network to Backprop

Complications arise because sequences can be rather long

(1) Computational memory issue

(2) Optimization & numerical instability