Sequence modeling

Professor Marie Roch



Sequence modeling

Basic ideas

- Sequences of vectors, e.g. frames in an audio stream
- Learn about subsequences



Unfolding computational graphs

- Consider a function that relies on previous state: $s^{(t)} = f(s^{(t-1)}|\theta)$
- Over time...

$$s^{(t)} = f(s^{(t-1)}|\theta)$$
$$s^{(t+1)} = f(s^{(t)}|\theta) = f(f(s^{(t-1)}|\theta)|\theta)$$

• We can think of this as a graph





Goodfellow et al. Fig 1

Add an input

- $s^{(t)} = f(s^{(t-1)}, x^{(t)}|\theta)$
- $s^{(t-1)}$ has history of input

• Consider this occurring in a hidden unit: $h^{(t)}=f(h^{(t-1)}, x^{(t)}|\theta)$





Recurrent neural networks (RNNs)



• Prior knowledge of $x^{(t)}$ tends to be lossy, as one usually does not need the entire history



RNNs take a sequence and output:

- a sequence
- or a single response

In addition to the current input, nodes can depend on:

- outputs of previous nodes
- previous outputs

Recurrence in the hidden layers can result in Turing complete networks.



Notation

We are using compact network representation (chapter 6), networks can have breadth

Goodfellow et al. Fig. 6.2

 h_1

 x_1







Case 1 example

one output/input, depends on previous inputs





Case 1 example

one output/input, depends on previous inputs











Forward propagation Initialize state $h^{(0)}$ For each time t $\alpha^{(t)} = b + W h^{(t-1)} + U x^{(t)}$ $h^{(t)} = \tanh(\alpha^{(t)})$ $o^{(t)} = c + Vh^{(t)}$ $\hat{y}^{(t)} = \operatorname{softmax}(o^{(t)})$ (postprocessing step)



Common RNN activation functions: sigmoid, tanh, & ReLU (be careful with ReLU - unbounded)

Loss of a sequence

• Sequence loss is sum of losses

$$L(\{x^{(1)}, x^{(2)}, \dots, x^{(\tau)}\}, \{y^{(1)}, y^{(2)}, \dots, y^{(\tau)}\}) = \sum_{t} L^{(t)}$$

• Suppose we use a negative log likelihood loss (e.g. MLE approach), this becomes

$$\sum_{t} L^{(t)} = -\sum_{t} \log P_{model}(y^{(t)} \mid x^{(1)}, x^{(2)}, \dots, x^{(t)})$$



Gradient of sequence (overview)

- Forward pass Compute gradient of each step
- Backward pass Use stored states and gradients to *back-propagate through time*
- Expensive in time and space















More on output recurrences

- Dependence solely on output is not Turing complete
- Loss function minimizes difference between target and output; nothing will drive it to capture information about network state





Output recurrences: Teacher forcing

When training for output recurrences, we could use a MLE approach to find the model which maximizes:

 $\log P(y^{(1)}, y^{(2)}|x^{(1)}, x^{(2)})$

but we know first output should be: $y^{(1)}$

$$\log P\left(y^{(2)} \middle| y^{(1)}, x^{(1)}, x^{(2)}\right) + \log P\left(y^{(1)} \middle| x^{(1)}, x^{(2)}\right)$$



Teacher forcing

If we use the previous label instead of the previous prediction $\log P(y^{(1)}, y^{(2)} | x^{(1)}, x^{(2)})$

we can parallelize the backpropagation as we don't need to compute the previous dependency.



Teacher forcing

We can rewrite the MLE target (Bayes rule): $\log P(y^{(1)}, y^{(2)} | x^{(1)}, x^{(2)})$ $= \log \left(P(y^{(2)}|y^{(1)}, x^{(1)}, x^{(2)}) P(y^{(1)}|x^{(1)}, x^{(2)}) \right)$ $= \log(P(y^{(2)}|y^{(1)}, x^{(2)})P(y^{(1)}|x^{(1)}))$ $y^{(t)}$ independent of $x^{(\tau)}$ where $t \neq \tau$ since there are no hidden recurrences $= \log P(y^{(2)}|y^{(1)}, x^{(2)}) + \log P(y^{(1)}|x^{(1)})$



Teacher forcing

Teacher forcing can be used in networks with output and hidden layer recurrences, but...

we can no longer parallelize training and must use back-propagation through time







At test time, we do not



RNN gradient computation

- Assume loss at time t: $\frac{\partial L}{\partial L^{(t)}} = 1$
- Consider network of Fig. 10.5:





RNN gradient computation

- Assume output layer o feeds softmax layer \hat{y}
- In this example, we use negative log likelihood on each output o_t(i) rather than cross entropy:

$$J(\theta) = -E_{x,y \sim \hat{p}_{data}} \log p_{model}(y|x)$$

= $\frac{1}{2} ||y - \hat{y}||^2$ (eqns 6.12, 6.13)

• We ignore the softmax layer in this discussion as it just produces probabilities and is not trainable.





RNN gradient computation

• Final loss:
$$\frac{\partial L}{\partial L^{(\tau)}} \triangleq 1$$

- Consider arbitrary output time: Loss at h comes from two paths
 - Next h
 - Current output





RNN output layer gradient

• Outputs: only dependent on loss (if no output dependencies)

$$\left(\nabla_{o^{(t)}} L \right)_{i} = \frac{\partial L}{\partial L^{(t)}} \frac{\partial L^{(t)}}{\partial o^{(t)}_{i}} = 1 \frac{\partial L}{o^{(t)}_{i}}$$

$$\frac{\partial L}{o^{(t)}_{i}} = \frac{\partial}{o^{(t)}_{i}} \frac{1}{2} ||y_{i} - \hat{y}_{i}||^{2}$$

$$= \frac{\partial}{o^{(t)}_{i}} \frac{1}{2} \sqrt{(y_{i} - \hat{y}_{i})^{2}}^{2}$$

$$= \frac{\partial}{o^{(t)}_{i}} \frac{1}{2} (y_{i} - \hat{y}_{i})^{2}$$

$$= (y_{i} - \hat{y}_{i})(-1)$$

$$= \hat{y}^{(t)}_{i} - 1_{i,y^{(t)}} \text{ as } y_{i} \text{ is } 0/1$$



h

RNN hidden layer gradients

- Last step time τ : Only successor is loss:
- Other times t: Loss at current step + loss from next h

$$\nabla_{h^{(t)}}L = \left(\frac{\partial h^{(t+1)}}{\partial h^{(t)}}\right)^T \nabla_{h^{(t+1)}}L + \left(\frac{\partial o^{(t)}}{\partial h^{(t)}}\right)^T \nabla_{o^{(t)}}L$$

Further simplification depends on activation function derivative (see text for an example)





Parameter updates

We need to learn how the loss should be applied to each of the parameters. General idea:

- Create copies of parameters at each time step (avoids introducing dependencies)
- Statistic of gradient over time

Example:

$$\nabla_W L = \sum_t \sum_i \left(\frac{\partial L}{\partial h_i^{(t)}} \right) \nabla_{W^{(t)}} h_i^{(t)}$$





oh, oh... rabbit hole

Marginal distributions

RG

<u></u>



British Postal Service, Graham Baker-Smith 2015

100

HILL CONTRACTOR

Integrate/sums over elements of a joint distribution to find P of desired random variable

$$X, Y \sim P(X, Y)$$
$$P(X = x) = \int_{\mathcal{Y}} P(x, y) dy$$
$$P(Y = y) = \int_{x} P(x, y) dx$$



Graphical models

- Networks of nodes can be used for inference/belief, frequently on latent (unobservable) variables
- Weights inform relationships between nodes
- Allows us to think about probabilities by marginalizing what we don't know

$$P(F \mid A, B) = \frac{\sum_{C} \sum_{D} \sum_{E} P(A, B, C, D, E, F)}{\sum_{C} \sum_{D} \sum_{D} \sum_{E} \sum_{F} P(A, B, C, D, E, F)}$$



Jordan 1997 doi:10.1.1.116.7467

RNNs as directed graphical models

- Simple case: RNN of a sequence where only recursion is previous output and there are *no inputs*
- With a negative log likelihood loss, $P(y^{(1)}, ..., y^{(\tau)}) = \prod_{t=1}^{\tau} P(y^{(t)} | y^{(t-1)}, ..., y^{(2)}, y^{(1)}) \text{ by chain rule of P}$ $L = \sum_{t} L^{(t)}$ $L^{(t)} = -logP(\hat{y}^{(t)} = y^{(t)} | y^{(t-1)}, y^{(t-2)}, ..., y^{(1)})$



as directed graphical models

- Graphical models with a long dependency require many edges
- These become difficult to estimate





as directed graphical models

Ways to deal with large # of parameters

- Remove weak dependencies
- Assume a Markov property, values more than τ time steps in the past no longer matter, e.g. $\tau = 2$: $P(y^{(t)}|y^{(t-1)}, \dots, y^{(2)}, y^{(1)})$ $=P(y^{(t)}|y^{(t-1)}, y^{(t-2)})$



as directed graphical models

• Introduce a state variable that captures the dependency



• Assumption: Relationship stationary, does not change over time



as directed graphical models

Use as a sequence generator introduces a new problem: How to stop...

- Can introduce a stop symbol, or
- Train a node to predict end of sequence
- Train a node to predict sequence length



Sequences conditioned on context

- RNNs can serve as graphical models that depend on input & output history
- When the input sequence is of fixed size (usually not the case for speech) can provide
 - sequence as input to initial state
 - sequence as input at each time step





Input and output context example 10.10 Goodfellow et al.) (Fig.



Bidirectional RNNs

• Sometimes, observing features in the future helps us know what occurred in the past...

 $P(x^{(2)} = e \mid x^{(1)} = b)$ vs. $P(x^{(2)} = e \mid x^{(1)} = b, x^{(3)} = d)$

• Very useful for speech recognition where coarticulation & phonotactics can play a large role



Bidirectional RNNs

• Compute two RNNs, one moving forward, the other backwards

• Avoids need for fixed length window required for feedforward nets or RNNs with a look-ahead buffer





Example Fig. 10.11 Goodfellow et al. Bidirectional RI



Recap so far

• Map sequence to fixed size vector $y^{(au)}$ $o^{(\tau)}$ Fig $y^{(t)}$ $m{h}^{(au)}$ $h^{(t)}$ Fig. 10.9 • Fixed vector to í _x(...) Ì $x^{(au)}$ sequence • N inputs to Fig. 10.3 N outputs SAN DIEGO STATE 40 IVERSITY Networks from Goodfellow et al.

What about this?



M to N sequences

Many applications require sequences of length M to be mapped to sequences of length N. e.g.

- speech recognition
- language translation



• query handling How can I get to Logan Heights?



M to N sequences

• Neural nets that handle this are known as *sequence to sequence*, or *encoder-decoder* architectures

- Basic idea:
 - Encoder RNN outputs a state or sequence representing the concept
 - Decoder RNN maps the concept onto another sequence



M to N sequences



Concept is typically final state or vector

Can use fixed vector to sequence architecture to generate new sequence



Fig. 10.12 Goodfellow et al.

Deep recurrent networks

What happens when the recursions pass through multiple layers?

Layer n (could be multiple layers) serves the role of building a representation to feed into the hidden state layer h





Deep recurrent networks

- If layer n is deep, this creates a long path between h and the information that loops back through it
- Long paths make learning difficult (more on this next)
- Can introduce a delay loop to assist





- Recurrence involves computing the same function over and over
- Let us simplify to think about implications
 - Use identity function as activation function
 - Recurrence becomes $h^{(t)} = W^T h^{(t-1)}$
 - Unwinding the recursion $h^{(t)} = (W^t)^T h^{(0)}$



- Suppose *W* has an eigenvalue decomposition (2.7) $W = Q\Lambda Q^{T}$
 - Q is a matrix of eigenvectors, $QQ^T = I$
 - Λ is a diagonal matrix of eigenvalues



 $h = (W^{t})^{T} h^{(0)} \qquad \text{t is number of time steps}$ $= (Q^{T} \Lambda Q)^{t} h^{(0)} \text{ as } W^{T} = (Q \Lambda Q^{T})^{T} = Q^{T} \Lambda Q$ $= Q^{T} \Lambda^{t} Q h^{(0)} \text{ as } Q^{T} \Lambda Q Q^{T} \Lambda Q = Q^{T} \Lambda I \Lambda Q = Q^{T} \Lambda \Lambda Q$

What happens if eigen values are

- small?
- large?



• In a non-recurrent net, weights will vary, but in RNNs the same weights are used over and over again.

• Techniques must be applied to be keep gradients from vanishing or exploding (see references in Goodfellow et al 10.7 for further details)



Long-term dependence strategies

Skip connections

- Create networks with longer delays
- Reduces number of times weight is applied
- Delay of d reduces gradient vanishing by factor of 1/d (in our simplified example: $Q\Lambda^{t/d}Q^T h^{(0)}$)
- Problems if what was skipped was important



Long-term dependence strategies

- Gated RNNs
 - Long-term dependence problem due to using the same weights at each step
 - What if the weights were dynamic?
 - Dynamic weights allow us to respond to input



Long short-term memory (LSTM)

- Popular gated architecture
- Weights are conditioned on context
- Can be interpreted as an attention mechanism















