Recurrent Neural Networks



Recurrent Neural Network

- h_t : hidden state
- X_t: input
- *Y_t*: output
- $Y_t, h_t = f(h_{t-1}, X_t; \theta)$
- h_{-1} : initial state



Fully-connect NN vs. RNN

RNN can be viewed as repeated applying fully-connected NNs

•
$$h_t = \sigma_1 (W^{(1)}X_t + W^{(11)}h_{t-1} + b^{(1)})$$

- $Y_t = \sigma_2(W^{(2)}h_t + b^{(2)})$
- σ_1, σ_2 are activation functions (sigmoid, ReLU, tanh, etc)



if icck Practical issues of RNN 7 may (20(1)) <1 =) (w (11) b2 exp shall Linear RNN derivation $6(2) \ge Z$ $h \in \mathbb{W}^{(1)} \times \mathbb{E} + W^{(1)} h_{\tau_{7}}$ he "forgets" X; $h_{\mathcal{R}} = W^{(1)} \cdot K_{\mathcal{R}} + W^{(11)} \cdot h_{\mathcal{R}-1} + W^{(11)} \cdot h_{\mathcal{R} = \left(W^{(1)} \right) h_{-1} + \sum_{j=0}^{k} \left(W^{(1)} \right) W^{(j)} \chi_{j}$ $\frac{2}{1 + \lambda_{man}} \left(\frac{U(1)}{2} \right) = \frac{2}{1 - \lambda_{max}} \left(\frac{U(1)}{2} \right) = \frac{2}{1 -$

Practical issues of RNN: training

you-lapar 6(2)

Gradient explosion and gradient vanishing

 $Z \in = W^{((\prime)} h_{t-1} + W^{(\prime)} \chi_t$ VR: output OR: layel $h_{f} > 6(2_{f})$ L R (G) = L (T R, D R)7 2 $\chi(w^{(u)})^{2}$ <u>JLR</u> Jha '(Zt) \prod 5 t=0 erci large forgetting pulling か 6(・) ~1

Techniques for avoiding gradient explosion threshold & if (10 LI 7 threshold of SLE threshold or (10 LI) Gradient clipping

- Identity initialization
- Truncated backprop through time
 - Only backprop for a few steps



Preserve Long-Term Memory

- Difficult for RNN to preserve long-term memory
 - The hidden state h_t is constantly being written (short-term memory)
 - Use a separate cell to maintain long-term memory



LSTM (Hochreitcher & Schmidhuber, '97)

- RNN architecture for learning long-term dependencies
- σ : layer with sigmoid activation



LSTM (Hochreitcher & Schmidhuber, '97)

- Core idea: maintain separate state h_t and cell c_t (memory)
- *h_t*: full update every step
- *c_t*: only *partially* update through gates
 - σ layer outputs importance ([0,1]) for each entry and only modify those entries of c_t



Forget gate f_t

• f_t outputs whether we want to "forget" things in c_t

- Compute $c_{t-1} \odot f_t$ (element-wise)
- $f_t(i) \rightarrow 0$: want to forget $c_t(i)$
- $f_t(i) \rightarrow 1$: we want to keep the information in $c_t(i)$



Input gate i_t

- i_t extracts useful information from X_t to update memory
 - \tilde{c}_t : information from X_t to update memory
 - i_t : which dimension in the memory should be updated by X_t
 - $i_t(j) \rightarrow 1$: we want to use the information in $\tilde{c}_t(j)$ to update memory
 - $i_t(t) \rightarrow 0$: $\tilde{c}_t(j)$ should not contribute to memory



Memory update

- $c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t$
- f_t forget gate; i_t input date
- $f_t \odot c_{t-1}$: drop useless information in old memory
- $i_t \odot \tilde{c}_t$: add selected new information from current input



Output gate o_t

• Next hidden state $h_t = o_t \odot \tanh(c_t)$

- $tanh(c_t)$: non-linear transformation over all past information
- *o_t*: choose important dimensions for the next state
 - $o_t(j) \rightarrow 1$: tanh $(c_t(j))$ is important for the next state
 - $o_t(j) \rightarrow 0$: tanh $(c_t(j))$ is not important



- $h_t = o_t \odot \tanh(c_t)$
- $c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t$



Remarks:

• $Y_t = g(h_t)$

- 1. No more matrix multiplications for c_t
- 2. LSTM does not have guarantees for gradient explosion/vanishing
- 3. LSTM is the dominant architecture for sequence modeling from '13 '16.
- 4. Why tanh

LSTM Variant

Peephold Connections (Gers & Schmidhuber '00)

• Allow gates to take in c_t information



LSTM Variant

Simplified LSTM

- Assume $i_t = 1 f_t$
- Only two gates are needed: fewer parameters



LSTM Variant

Gated Recurrent Unit (GRU, Cho et al. '14) • Merge h_t and c_t : much fewer parameters



LSTM application: language model

- Autoregressive language model: $P(X; \theta) = \prod_{t=1}^{L} P(X_t \mid X_{i < t}; \theta)$
 - X: a sentence
 - Sequential generation
- LSTM language model
 - *X_t*: word at position *t*.
 - *Y_t*: softmax over all words
- Data: a collection of texts:
 - Wiki



LSTM application: text classification

Bi-directional LSTM and them run softmax on the final hidden state.



Attention Mechanism



Machine Translation

Translation

- Before 2014: Statistical Machine Learning (SMT)
 - Extremely complex systems that require massive human efforts
 - Separately designed components
 - A lot of feature engineering
 - Lots of linguistic domain knowledge and expertise
- Before 2016:
 - Google Translate is based on statistical machine learning
- What happened in 2014?
 - Neural machine translation (NMT)

Sequence to Sequence Model

- Neural Machine Translation (NMT)
 - Learning to translate via a **single end-to-end** neural network.
 - Source language sentence X, target language sentence $Y = f(X; \theta)$
- Sequence to Sequence Model (Seq2Seq, Sutskever et al., '14)
 - Two RNNs: f_{enc} and f_{dec}
 - Encoder f_{enc} :
 - Takes X as input, and output the initial hidden state for decoder

X I-> h

- Can use bidirectional RNN
- Decoder f_{dec} :
 - It takes in the hidden state from f_{enc} to generate Y
 - Can use autoregressive language model

h -> Y

Sequence to Sequence Model



Training Sequence to Sequence Model

- Collect a huge paired dataset and train it end-to-end via BPTT
- Loss induced by MLE $P(Y|X) = P(Y|f_{enc}(X))$



Seq2seq is optimized as a single system. Backpropagation operates "end-to-end".

Deep Sequence to Sequence Model

Stacked seq2seq model



Machine Translation

• 2016: Google switched Google Translate from SMT to NMT



Alignment

- Alignment: the word-level correspondence between X and Y
- Can have complex long-term dependencies



Issue in Seq2Seq

- Alignment: the word-level correspondence between X and Y
 - The information bottleneck due to the hidden state h
 - We want each Y_t to also focus on some X_i that it is aligned with



- NMT by jointly learning to align and translate (Bahdanau, Cho, Bengio, '15)
- Core idea:
 - When decoding Y_t , consider both hidden states and alignment:
 - Hidden state: $h_t = f_{dec}(Y_{i < t})$
 - Alignment: connect to a portion of X
 - When portion of *X* to focus on?
 - Learn a softmax weight over X: attention distribution P_{att}
 - $P_{att}(X_i | h_t)$: how much attention to put on word X_i
 - Attention output $h_{att} = \sum_{i} f_{enc}(X_i | X_{j < i}) \cdot P_{att}(X_i | h_{t-1})$
 - Use h_{t-1} and h_{att} to compute Y_t

Weighted such

hff



Decoder RNN







Decoder RNN







Decoder RNN







Decoder RNN



Decoder RNN



Decoder RNN



Summary

• Input sequence X, encoder f_{enc} , and decoder f_{dec}

() - wood

- $f_{enc}(X)$ produces hidden states $h_1^{enc}, h_2^{enc}, \dots, h_N^{enc}$
- On time step t, we have decoder hidden state h_t
- Compute attention score $e_i = h_t^{\top} h_i^{enc}$
- Compute attention distribution $\alpha_i = P_{att}(X_i) = \text{softmax}(e_i)$

• Attention output: $h_{att}^{enc} = \sum \alpha_i h_i^{enc}$

- $Y_t \sim g(h_t, h_{att}^{enc}; \theta)$
 - Sample an output using both h_t and h_{att}^{enc}

Attention

- It significantly improves NMT.
- It solves the bottleneck problem and the long-term dependency issue.
- Also helps gradient vanishing problem.
- Provides some interpretability
 - Understanding which word the RNN encoder focuses on
- Attention is a general technique
 - Given a set of vector values V_i and vector query q
 - Attention computes a weighted sum of values depending on \boldsymbol{q}

Other use cases:

- Attention can be viewed as a module.
- In encoder and decoder (more on this later)
- A representation of a set of points
 - Pointer network (Vinyals, Forunato, Jaitly '15)
 - Deep Sets (Zaheer et al., '17)
- Convolutional neural networks
 - To include non-local information in CNN (Non-local network, '18)







Attention

- Representation learning:
 - A method to obtain a fixed representation corresponding to a query q from an arbitrary set of representations $\{V_i\}$
 - Attention distribution: $\alpha_i = \operatorname{softmax}(f(v_i, q))$
 - Attention output: $v_{att} = \sum \alpha_i v_i$
- Attent variant: $f(v_i, q)$
 - Multiplicative attention: $f(v_i, q) = q^{\top} W k_i$, W is a weight matrix
 - Additive attention: $\widehat{f(v_i, q)} = u^{\mathsf{T}} \operatorname{tanh}(W_1 v_i + W_2 q)$

Key-query-value attention

• Obtain q_t, v_t, k_t from X_t

• $q_t = W^q X_t$; $v_t = W^v X_t$; $k_t = W^k X_t$ (position encoding omitted)

• W^q , W^v , W^k are learnable weight matrices

•
$$\alpha_{i,j} = \operatorname{softmax}(q_i^{\top}k_j); out_i = \sum_k \alpha_{i,j}v_j$$

Intuition: key, query, and value can focus on different parts of input



Attention is all you need (Vsawani '17)

- A pure attention-based architecture for sequence modeling
 - No RNN at all!
- Basic component: self-attention, $Y = f_{SA}(X; \theta)$
 - X_t uses attention on entire X sequence
 - Y_t computed from X_t and the attention output
- Computing Y_t
 - Key k_t , value v_t , query q_t from X_t
 - $(k_t, v_t, q_t) = g_1(X_t; \theta)$
 - Attention distribution $\alpha_{t,j} = \operatorname{softmax}(q_t^{\top} k_j)$

• Attention output $out_t = \sum_{i} \alpha_{t,j} v_j$

•
$$Y_t = g_2(out_t; \theta)$$



Issues of Vanilla Self-Attention

• Attention is order-invariant

- Lack of non-linearities
 - All the weights are simple weighted average

- Capability of autoregressive modeling
 - In generation tasks, the model cannot "look at the future"
 - e.g. Text generation:
 - Y_t can only depend on $X_{i < t}$
 - But vanilla self-attention requires the entire sequence

Position Encoding

Vanilla self-attention

- $(k_t, v_t, q_t) = g_1(X_t; \theta)$
- $\alpha_{t,j} = \operatorname{softmax}(q_t^{\mathsf{T}} k_j)$

Attention output
$$out_t = \sum_j \alpha_{t,j} v_j$$

- Idea: position encoding:
 - p_i : an embedding vector (feature) of position i
 - $(k_t, v_t, q_t) = g_1([X_t, p_t]; \theta)$
- In practice: Additive is sufficient: $k_t \leftarrow \tilde{k}_t + p_t, q_t \leftarrow \tilde{q}_t + p_t, v_t \leftarrow \tilde{v}_t + p_t;$ $(\tilde{k}_t, \tilde{v}_t, \tilde{q}_t) = g_1(X_t; \theta)$
- p_t is only included in the first layer

Position Encoding

 $p_t \operatorname{design} 1$: Sinusoidal position representation

- Pros:
 - simple
 - naturally models "relative position"
 - Easily applied to long sequences
- Cons:
 - Not learnable
 - Generalization poorly to sequences longer than training data



Heatmap of $p_i^T p_j$

Position Encoding

p_t design 2: Learned representation

- Assume maximum length L, learn a matrix $p \in \mathbb{R}^{d \times T}$, p_t is a column of p
- Pros:
 - Flexible
 - Learnable and more powerful
- Cons:
 - Need to assume a fixed maximum length L
 - Does not work at all for length above L
- p_t design 3: Relative position representation (Shaw, Uszkoreit, Vaswani '18)

Combine Self-Attention with Nonlinearity

- Vanilla self-attention
 - No element-wise activation (e.g., ReLU, tanh)
 - Only weighted average and softmax operator
- Fix:
 - Add an MLP to process *out_i*
 - $m_i = MLP(out_i) = W_2 \text{ReLU}(W_1 out_i + b_1) + b_2$
 - Usually do not put activation layer before softmaax



Masked Attention

- In language model decoder: $P(Y_t | X_{i < t})$
 - out_t cannot look at future $X_{i>t}$
- Masked attention
 - Compute $e_{i,j} = q_i^{\mathsf{T}} k_j$ as usuall
 - Mask out $e_{i>j}$ by setting $e_{i>j} = -\infty$
 - $e \odot (1 M) \leftarrow -\infty$
 - M is a fixed 0/1 mask matrix
 - Then compute $\alpha_i = \operatorname{softmax}(e_i)$
 - Remarks:
 - M = 1 for full self-attention
 - Set M for arbitrary dependency ordering



raw attention weights

mask



Transformer-based sequence-to-sequence modeling



Key-query-value attention

• Obtain q_t, v_t, k_t from X_t

• $q_t = W^q X_t$; $v_t = W^v X_t$; $k_t = W^k X_t$ (position encoding omitted)

• W^q , W^v , W^k are learnable weight matrices

•
$$\alpha_{i,j} = \operatorname{softmax}(q_i^{\top}k_j); out_i = \sum_k \alpha_{i,j}v_j$$

Intuition: key, query, and value can focus on different parts of input



Multi-headed attention

- Standard attention: single-headed attention
 - $X_t \in \mathbb{R}^d$, $Q, K, V \in \mathbb{R}^{d \times d}$
 - We only look at a single position j with high $\alpha_{i,j}$
 - What if we want to look at different j for different reasons?
- Idea: define h separate attention heads
 - *h* different attention distributions, keys, values, and queries

•
$$Q^{\ell}, K^{\ell}, V^{\ell} \in \mathbb{R}^{d \times \frac{d}{h}}$$
 for $1 \leq \ell \leq h$

•
$$\alpha_{i,j}^{\ell} = \operatorname{softmax}((q_i^{\ell})^{\mathsf{T}} k_j^{\ell}); out_i^{\ell} = \sum_j \alpha_{i,j}^{\ell} v_j^{\ell}$$

#Params Unchanged!



Multi-headed attention

- Standard attention: single-headed attention
 - $X_t \in \mathbb{R}^d$, $Q, K, V \in \mathbb{R}^{d \times d}$
 - We only look at a single position j with high $\alpha_{i,j}$
 - What if we want to look at different *j* for different reasons?
- Idea: define h separate attention heads
 - *h* different attention distributions, keys, values, and queries
 - $Q^{\ell}, K^{\ell}, V^{\ell} \in \mathbb{R}^{d \times \frac{d}{h}}$ for $1 \leq \ell \leq h$

•
$$\alpha_{i,j}^{\ell} = \operatorname{softmax}((q_i^{\ell})^{\mathsf{T}} k_j^{\ell}); out_i^{\ell} = \sum_j \alpha_{i,j}^{\ell} v_j^{\ell}$$



Transformer-based sequence-to-sequence model

- Basic building blocks: self-attention
 - Position encoding
 - Post-processing MLP
 - Attention mask
- Enhancements:
 - Key-query-value attention
 - Multi-headed attention
 - Architecture modifications:
 - Residual connection
 - Layer normalization



Machine translation with transformer

-				-	
Model	BLEU		Training Cost (FLOPs)		
	EN-DE	EN-FR	EN-DE	EN-FR	
ByteNet [18]	23.75				
Deep-Att + PosUnk [39]		39.2		$1.0\cdot10^{20}$	
GNMT + RL [38]	24.6	39.92	$2.3 \cdot 10^{19}$	$1.4\cdot10^{20}$	
ConvS2S [9]	25.16	40.46	$9.6 \cdot 10^{18}$	$1.5 \cdot 10^{20}$	
MoE [32]	26.03	40.56	$2.0\cdot10^{19}$	$1.2\cdot10^{20}$	
Deep-Att + PosUnk Ensemble [39]		40.4		$8.0 \cdot 10^{20}$	
GNMT + RL Ensemble [38]	26.30	41.16	$1.8\cdot 10^{20}$	$1.1\cdot10^{21}$	
ConvS2S Ensemble [9]	26.36	41.29	$7.7 \cdot 10^{19}$	$1.2\cdot10^{21}$	
Transformer (base model)	27.3	38.1	3.3 •	$3.3\cdot10^{18}$	
Transformer (big)	28.4	41.8	$2.3 \cdot$	$2.3 \cdot 10^{19}$	

- Limitations of transformer: Quadratic computation cost
 - Linear for RNNs
 - Large cost for large sequence length, e.g., $L > 10^4$
- Follow-ups:
 - Large-scale training: transformer-XL; XL-net ('20)
 - Projection tricks to O(L): Linformer ('20)
 - Math tricks to O(L): Performer ('20)
 - Sparse interactions: Big Bird ('20)
 - Deeper transformers: DeepNet ('22)

Transformer for Images

- Vision Transformer ('21)
 - Decompose an image to 16x16 patches and then apply transformer encoder





Transformer for Images

- Swin Transformer ('21)
 - Build hierachical feature maps at different resolution
 - Self-attention only within each block
 - Shifted block partitions to encode information between blocks



CNN vs. RNN vs. Attention



Summary

- Language model & sequence to sequence model:
 - Fundamental ideas and methods for sequence modeling
- Attention mechanism
 - So far the most successful idea for sequence data in deep learning
 - A scale/order-invariant representation
 - Transformer: a fully attention-based architecture for sequence data
 - Transformer + Pretraining: the core idea in today's NLP tasks
- LSTM is still useful in lightweight scenarios