CSCI567 Machine Learning (Fall 2025)

Haipeng Luo

University of Southern California

Nov 7, 2025

Administration

Will discuss HW3 solutions in today's discussion session.

HW4 (last homework) will be released soon.

Outline

Review of last lecture

Recurrent Neural Network

Transformers

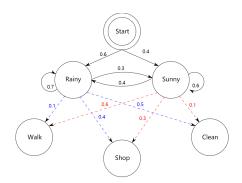
Outline

- Review of last lecture
- 2 Recurrent Neural Network
- Transformers

Hidden Markov Models

Model parameters:

- initial distribution $P(Z_1 = s) = \pi_s$
- transition distribution $P(Z_{t+1} = s' \mid Z_t = s) = a_{s,s'}$
- emission distribution $P(X_t = o \mid Z_t = s) = b_{s,o}$



Baum-Welch algorithm

Step 0 Initialize the parameters $(m{\pi}, m{A}, m{B})$

Step 1 (E-Step) Fixing the parameters, compute forward and backward messages for all sample sequences, then use these to compute $\gamma_s^{(n)}(t)$ and $\xi_{s,s'}^{(n)}(t)$ for each n,t,s,s'.

Step 2 (M-Step) Update parameters:

$$\pi_s \propto \sum_n \gamma_s^{(n)}(1), \quad a_{s,s'} \propto \sum_n \sum_{t=1}^{T-1} \xi_{s,s'}^{(n)}(t), \quad b_{s,o} \propto \sum_n \sum_{t:x_t=o} \gamma_s^{(n)}(t)$$

Step 3 Return to Step 1 if not converged

Viterbi Algorithm

Viterbi Algorithm

For each $s \in [S]$, compute $\delta_s(1) = \pi_s b_{s,x_1}$.

For each $t = 2, \ldots, T$,

• for each $s \in [S]$, compute

$$\delta_s(t) = b_{s,x_t} \max_{s'} a_{s',s} \delta_{s'}(t-1)$$

$$\Delta_s(t) = \operatorname*{argmax}_{s'} a_{s',s} \delta_{s'}(t-1)$$

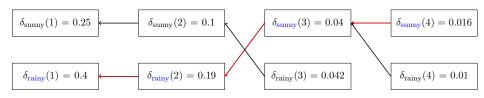
Backtracking: let $z_T^* = \operatorname{argmax}_s \delta_s(T)$.

For each t = T, ..., 2: set $z_{t-1}^* = \Delta_{z_t^*}(t)$.

Output the most likely path z_1^*, \ldots, z_T^* .

Example

Arrows represent $\Delta_s(t)$, backtracking = follow the arrows.



The most likely path is "rainy, rainy, sunny, sunny".

Outline

- Review of last lecture
- Recurrent Neural Network
 - RNN: model
 - RNN: training and testing
- Transformers

Recall: language models via HMM

- today the _____
- today the price _____
- today the price of _____
- today the price of gold _____

final result: today the price of gold per ton, while production of shoe lasts and shoe industry, the bank intervened just after it considered and rejected an IMF demand to rebuild depleted European stocks.

Surprisingly grammatical! but incoherent...

company bank	0.153	
pank	0.153	
price	0.077	
italian	0.039	
emirate	0.039	

of	0.308	
for	0.050	
it	0.046	
to	0.046	
is	0.031	

the	0.072	
18	0.043	
oil	0.043	
its	0.036	
gold	0.018	

Key ideas for improvement:

Key ideas for improvement:

• represent words as vectors, enabling differentiable operations

Key ideas for improvement:

- represent words as vectors, enabling differentiable operations
- flexible sequence to sequence architecture (not just vector to vector)

Key ideas for improvement:

- represent words as vectors, enabling differentiable operations
- flexible sequence to sequence architecture (not just vector to vector)
- shared components (just like filters in CNN)

Key ideas for improvement:

- represent words as vectors, enabling differentiable operations
- flexible sequence to sequence architecture (not just vector to vector)
- shared components (just like filters in CNN)

Recurrent Neural Network (RNN) is one solution (popular before transformers).

Acknowledgements

Very useful resources:

- RNN cheatsheet from Stanford CS 230
 - https://stanford.edu/~shervine/teaching/cs-230/ cheatsheet-recurrent-neural-networks
- Visualizing a tiny RNN
 - https://joshvarty.github.io/VisualizingRNNs/
- Character-level RNN
 - https://karpathy.github.io/2015/05/21/rnn-effectiveness/

Simplest approach: one-hot sparse encoding

Simplest approach: one-hot sparse encoding

ullet suppose there are d words in the vocabulary

Simplest approach: one-hot sparse encoding

- ullet suppose there are d words in the vocabulary
- represent the *i*-th of them by the *d*-dimensional basis vector

$$e_i = (0, \dots, 0, 1, 0, \dots, 0) \in \mathbb{R}^d$$

$$\uparrow$$
i-th entry

Simplest approach: one-hot sparse encoding

- \bullet suppose there are d words in the vocabulary
- represent the *i*-th of them by the *d*-dimensional basis vector

$$e_i = (0, \dots, 0, 1, 0, \dots, 0) \in \mathbb{R}^d$$

$$\uparrow$$
i-th entry

Issues: does not convey any semantic meanings

Word embedding: similar words are closer in their vector representation

Word embedding: similar words are closer in their vector representation

popular approaches: word2vec, GloVe

(see project)

Word embedding: similar words are closer in their vector representation

popular approaches: word2vec, GloVe

(see project)

can even perform meaningful algebraic operations

Word embedding: similar words are closer in their vector representation

popular approaches: word2vec, GloVe

(see project)

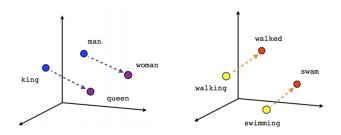
- can even perform meaningful algebraic operations
 - example: man is to woman as king is to?

Word embedding: similar words are closer in their vector representation

- popular approaches: word2vec, GloVe (see project)
- can even perform meaningful algebraic operations
 - example: man is to woman as king is to?
 - can be answered by finding the word with the closest embedding to vec(woman) - vec(man) + vec(king), which happens to be "queen"

Word embedding: similar words are closer in their vector representation

- popular approaches: word2vec, GloVe (see project)
- can even perform meaningful algebraic operations
 - example: man is to woman as king is to?
 - can be answered by finding the word with the closest embedding to vec(woman) - vec(man) + vec(king), which happens to be "queen"



Let $x \in \mathbb{R}^d$ be the one-hot encoding of a word, and matrix $E \in \mathbb{R}^{d_e \times d}$ be some **embedding matrix**, then Ex is the embedding for this word

Let $x \in \mathbb{R}^d$ be the one-hot encoding of a word, and matrix $E \in \mathbb{R}^{d_e \times d}$ be some **embedding matrix**, then Ex is the embedding for this word

 $oldsymbol{e}$ can be fixed (i.e., from word2vec or GloVe), where the i-th column is the embedding for the i-th word

Let $x \in \mathbb{R}^d$ be the one-hot encoding of a word, and matrix $E \in \mathbb{R}^{d_e \times d}$ be some **embedding matrix**, then Ex is the embedding for this word

- $m{e}$ can be fixed (i.e., from word2vec or GloVe), where the i-th column is the embedding for the i-th word
- \bullet or E can be *learned* (e.g., via backpropagation in DL pipeline),

Let $x \in \mathbb{R}^d$ be the one-hot encoding of a word, and matrix $E \in \mathbb{R}^{d_e \times d}$ be some **embedding matrix**, then Ex is the embedding for this word

- $m{e}$ can be fixed (i.e., from word2vec or GloVe), where the i-th column is the embedding for the i-th word
- or E can be learned (e.g., via backpropagation in DL pipeline),
 making it application specific (common especially if data are huge)

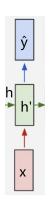
Let $x \in \mathbb{R}^d$ be the one-hot encoding of a word, and matrix $E \in \mathbb{R}^{d_e \times d}$ be some **embedding matrix**, then Ex is the embedding for this word

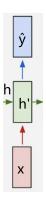
- $m{e}$ can be fixed (i.e., from word2vec or GloVe), where the i-th column is the embedding for the i-th word
- or E can be learned (e.g., via backpropagation in DL pipeline),
 making it application specific (common especially if data are huge)

In the remaining, we simply use one-hot representation, but keep in mind it could be passed through some \boldsymbol{E} implicitly

from
$$\widehat{m{y}} = m{f}(m{x})$$
 to

from
$$\widehat{m{y}} = m{f}(m{x})$$
 to $(\widehat{m{y}}, h') = m{f}(m{x}, h)$



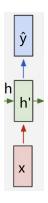


from
$$\widehat{m{y}} = m{f}(m{x})$$
 to $(\widehat{m{y}}, h') = m{f}(m{x}, h)$

• h is "hidden state" (like HMM), updated via

$$h' = \sigma(Wh + Ux + b_h)$$

where σ is an activation function



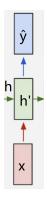
from
$$\widehat{m{y}} = m{f}(m{x})$$
 to $(\widehat{m{y}}, m{h}') = m{f}(m{x}, m{h})$

• h is "hidden state" (like HMM), updated via

$$h' = \sigma(Wh + Ux + b_h)$$

where σ is an activation function

ullet $\widehat{oldsymbol{y}} = oldsymbol{V}oldsymbol{h}' + oldsymbol{b}_y$ is the <code>output</code>



from
$$\widehat{m{y}} = m{f}(m{x})$$
 to $(\widehat{m{y}}, h') = m{f}(m{x}, h)$

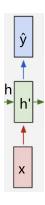
• h is "hidden state" (like HMM), updated via

$$h' = \sigma(Wh + Ux + b_h)$$

where σ is an activation function

- $oldsymbol{\hat{y}} = oldsymbol{V}oldsymbol{h}' + oldsymbol{b}_y$ is the <code>output</code>
- $ullet oldsymbol{x}, \widehat{oldsymbol{y}} \in \mathbb{R}^d, oldsymbol{h}, oldsymbol{h}' \in \mathbb{R}^{d_h}$

A recurrent layer



from
$$\widehat{m{y}} = m{f}(m{x})$$
 to $(\widehat{m{y}}, h') = m{f}(m{x}, h)$

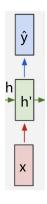
• h is "hidden state" (like HMM), updated via

$$h' = \sigma(Wh + Ux + b_h)$$

where σ is an activation function

- $oldsymbol{\hat{y}} = oldsymbol{V}oldsymbol{h}' + oldsymbol{b}_y$ is the <code>output</code>
- $ullet oldsymbol{x}, \widehat{oldsymbol{y}} \in \mathbb{R}^d, oldsymbol{h}, oldsymbol{h}' \in \mathbb{R}^{d_h}$
- ullet weight matrices $oldsymbol{W} \in \mathbb{R}^{d_h imes d_h}, oldsymbol{U} \in \mathbb{R}^{d_h imes d}, oldsymbol{V} \in \mathbb{R}^{d imes d_h}$

A recurrent layer



from
$$\widehat{m{y}} = m{f}(m{x})$$
 to $(\widehat{m{y}}, h') = m{f}(m{x}, h)$

• h is "hidden state" (like HMM), updated via

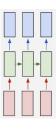
$$m{h}' = m{\sigma}(m{W}m{h} + m{U}m{x} + m{b}_h)$$

where σ is an activation function

- $oldsymbol{\hat{y}} = oldsymbol{V}oldsymbol{h}' + oldsymbol{b}_y$ is the <code>output</code>
- $ullet oldsymbol{x}, \widehat{oldsymbol{y}} \in \mathbb{R}^d, oldsymbol{h}, oldsymbol{h}' \in \mathbb{R}^{d_h}$
- ullet weight matrices $oldsymbol{W} \in \mathbb{R}^{d_h imes d_h}, oldsymbol{U} \in \mathbb{R}^{d_h imes d_h}, oldsymbol{V} \in \mathbb{R}^{d imes d_h}$
- ullet bias terms $oldsymbol{b}_h \in \mathbb{R}^{d_h}$, $oldsymbol{b}_u \in \mathbb{R}^d$

Given a **sequence** x_1, x_2, \ldots , can apply f recursively:

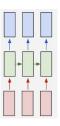
Given a **sequence** x_1, x_2, \ldots , can apply f recursively:



•
$$h_0 = 0$$

$$\bullet$$
 $(\hat{y}_1, h_1) = f(x_1, h_0)$

Given a **sequence** x_1, x_2, \ldots , can apply f recursively:



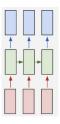
•
$$h_0 = 0$$

$$ullet (\widehat{m{y}}_1, m{h}_1) = m{f}(m{x}_1, m{h}_0)$$

$$\bullet$$
 $(\hat{y}_2, h_2) = f(x_2, h_1)$

. . .

Given a **sequence** x_1, x_2, \ldots , can apply f recursively:



•
$$h_0 = 0$$

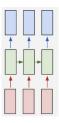
$$ullet (\widehat{m{y}}_1, m{h}_1) = m{f}(m{x}_1, m{h}_0)$$

$$\bullet$$
 $(\hat{y}_2, h_2) = f(x_2, h_1)$

• • • •

This is **one** recurrent layer unfolded (over steps), *not many different layers*.

Given a **sequence** x_1, x_2, \ldots , can apply f recursively:



•
$$h_0 = 0$$

$$ullet (\widehat{m{y}}_1, m{h}_1) = m{f}(m{x}_1, m{h}_0)$$

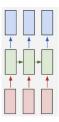
$$\bullet$$
 $(\hat{y}_2, h_2) = f(x_2, h_1)$

• • • •

This is **one** recurrent layer unfolded (over steps), *not many different layers*.

The same f (i.e, W, U, V, b_h, b_y) is **shared** in all steps (similar to CNN's filters shared across different spatial locations).

Given a **sequence** x_1, x_2, \ldots , can apply f recursively:



•
$$h_0 = 0$$

$$ullet (\widehat{m{y}}_1, m{h}_1) = m{f}(m{x}_1, m{h}_0)$$

$$\bullet$$
 $(\hat{y}_2, h_2) = f(x_2, h_1)$

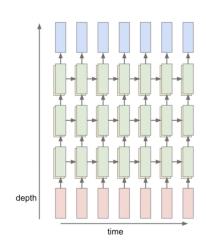
. . .

This is **one** recurrent layer unfolded (over steps), not many different layers.

The same f (i.e, W, U, V, b_h, b_y) is **shared** in all steps (similar to CNN's filters shared across different spatial locations).

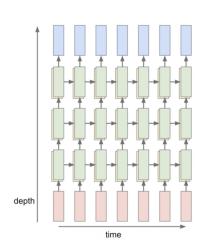
Hidden state h_t summarizes information up to step t.

Stack multiple recurrent layers:



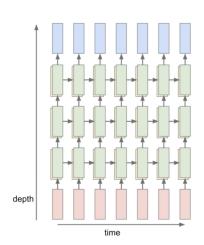
Stack multiple recurrent layers:

• hidden states become the inputs of the next layer



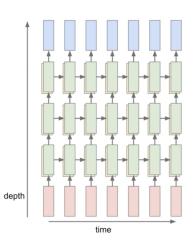
Stack multiple recurrent layers:

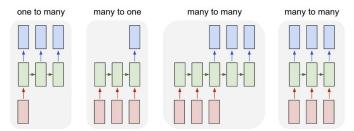
- hidden states become the inputs of the next layer
- ullet different layers learn different $oldsymbol{W}, oldsymbol{U}, oldsymbol{b}_h$



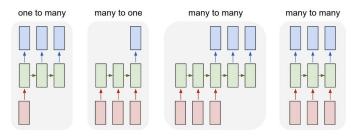
Stack multiple recurrent layers:

- hidden states become the inputs of the next layer
- ullet different layers learn different $oldsymbol{W}, oldsymbol{U}, oldsymbol{b}_h$
- ullet last layer learns $oldsymbol{V}, oldsymbol{b}_y$ and output $\widehat{oldsymbol{y}}$

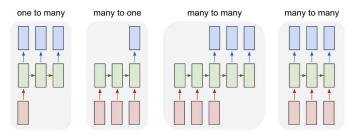




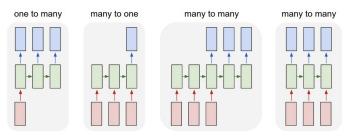
Many possible structures and applications:



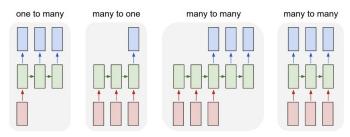
• one-to-many: image captioning



- one-to-many: image captioning
- many-to-one: sentiment classification



- one-to-many: image captioning
- many-to-one: sentiment classification
- many-to-many: machine translation, question answering



- one-to-many: image captioning
- many-to-one: sentiment classification
- many-to-many: machine translation, question answering
- (aligned) many-to-many: POS tagging, name entity recognition

Take **text generation** (unsupervised learning) as an example:

Take **text generation** (unsupervised learning) as an example:

ullet given a corpus, train an RNN that learns $P(oldsymbol{x}_t \mid oldsymbol{x}_{1:t-1})$

Take **text generation** (unsupervised learning) as an example:

ullet given a corpus, train an RNN that learns $P(oldsymbol{x}_t \mid oldsymbol{x}_{1:t-1})$

For each sequence $oldsymbol{x}_1,\dots,oldsymbol{x}_T\in\mathbb{R}^d$ (one-hot representation) in the corpus

• feed x_1, \dots, x_{T-1} into the current RNN to get $\widehat{y}_1, \dots, \widehat{y}_{T-1} \in \mathbb{R}^d$

Take **text generation** (unsupervised learning) as an example:

ullet given a corpus, train an RNN that learns $P(oldsymbol{x}_t \mid oldsymbol{x}_{1:t-1})$

For each sequence $oldsymbol{x}_1,\dots,oldsymbol{x}_T\in\mathbb{R}^d$ (one-hot representation) in the corpus

- ullet feed $m{x}_1,\ldots,m{x}_{T-1}$ into the current RNN to get $\widehat{m{y}}_1,\ldots,\widehat{m{y}}_{T-1}\in\mathbb{R}^d$
- each \hat{y}_t defines a distribution over the next word via softmax: $P(\text{next word} = i) \propto \exp(\hat{y}_{t,i})$

Take **text generation** (unsupervised learning) as an example:

ullet given a corpus, train an RNN that learns $P(oldsymbol{x}_t \mid oldsymbol{x}_{1:t-1})$

For each sequence $oldsymbol{x}_1,\dots,oldsymbol{x}_T\in\mathbb{R}^d$ (one-hot representation) in the corpus

- feed x_1, \dots, x_{T-1} into the current RNN to get $\widehat{y}_1, \dots, \widehat{y}_{T-1} \in \mathbb{R}^d$
- each \widehat{y}_t defines a distribution over the next word via softmax: $P(\text{next word} = i) \propto \exp(\widehat{y}_{t,i})$
- ullet based on the true label $oldsymbol{x}_{t+1}$, each $\widehat{oldsymbol{y}}_t$ incurs cross-entropy loss

$$-\ln\left(\frac{\exp(\widehat{\boldsymbol{y}}_t^{\top}\boldsymbol{x}_{t+1})}{\sum_{i=1}^{d}\exp(\widehat{\boldsymbol{y}}_{t,i})}\right)$$

Take **text generation** (unsupervised learning) as an example:

ullet given a corpus, train an RNN that learns $P(oldsymbol{x}_t \mid oldsymbol{x}_{1:t-1})$

For each sequence $oldsymbol{x}_1,\dots,oldsymbol{x}_T\in\mathbb{R}^d$ (one-hot representation) in the corpus

- feed x_1, \dots, x_{T-1} into the current RNN to get $\widehat{y}_1, \dots, \widehat{y}_{T-1} \in \mathbb{R}^d$
- each \widehat{y}_t defines a distribution over the next word via softmax: $P(\text{next word} = i) \propto \exp(\widehat{y}_{t,i})$
- ullet based on the true label $oldsymbol{x}_{t+1}$, each $\widehat{oldsymbol{y}}_t$ incurs cross-entropy loss

$$-\ln\left(\frac{\exp(\widehat{\boldsymbol{y}}_t^{\top}\boldsymbol{x}_{t+1})}{\sum_{i=1}^{d}\exp(\widehat{\boldsymbol{y}}_{t,i})}\right)$$

update the RNN parameters using backpropagation over the total loss

Tiny RNN, predicting the next bit of a binary sequence

• https://joshvarty.github.io/VisualizingRNNs/

Demo

Tiny RNN, predicting the next bit of a binary sequence

- https://joshvarty.github.io/VisualizingRNNs/
- the entire vocabulary is just $\{0,1\}$ (d=2)

Demo

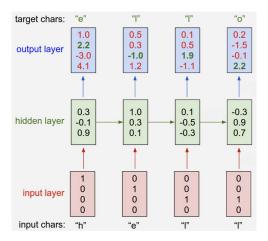
Tiny RNN, predicting the next bit of a binary sequence

- https://joshvarty.github.io/VisualizingRNNs/
- the entire vocabulary is just $\{0,1\}$ (d=2)
- one-layer RNN with $d_h=3$, so parameters are $m{W}\in\mathbb{R}^{3 imes3}, m{U}\in\mathbb{R}^{3 imes2}, m{V}\in\mathbb{R}^{2 imes3}, m{b}_h\in\mathbb{R}^3, m{b}_y\in\mathbb{R}^2$

Another demo

Min-Char RNN, predicting the next character of a sequence

https://karpathy.github.io/2015/05/21/rnn-effectiveness/



Keep sampling from softmax $(\widehat{m{y}}_t)$ as the next input $m{x}_{t+1}$ to RNN

Keep sampling from softmax $(\widehat{m{y}}_t)$ as the next input $m{x}_{t+1}$ to RNN

Can control how "random" the generation is via softmax $({m eta}\cdot {m {\hat y}}_t)$

Keep sampling from softmax $(\widehat{m{y}}_t)$ as the next input $m{x}_{t+1}$ to RNN

Can control how "random" the generation is via softmax $({m eta} \cdot {m {\hat y}}_t)$

• $1/\beta$ is called temperature

Keep sampling from softmax $(\widehat{m{y}}_t)$ as the next input $m{x}_{t+1}$ to RNN

Can control how "random" the generation is via softmax $({m eta} \cdot {m {\hat y}}_t)$

- $1/\beta$ is called temperature
- ullet larger temperature (smaller eta) leads to more random outputs

Keep sampling from softmax $(\widehat{m{y}}_t)$ as the next input $m{x}_{t+1}$ to RNN

Can control how "random" the generation is via softmax $(oldsymbol{eta}\cdot\widehat{oldsymbol{y}}_t)$

- $1/\beta$ is called temperature
- ullet larger temperature (smaller eta) leads to more random outputs
 - $\beta = 0$, uniform output (maximum entropy)

Keep sampling from softmax $(\widehat{m{y}}_t)$ as the next input $m{x}_{t+1}$ to RNN

Can control how "random" the generation is via softmax $(\pmb{\beta} \cdot \widehat{\pmb{y}}_t)$

- $1/\beta$ is called temperature
- larger temperature (smaller β) leads to more random outputs
 - $\beta = 0$, uniform output (maximum entropy)
 - $oldsymbol{eta}=\infty$, deterministically output $rgmax_i\,\widehat{oldsymbol{y}}_{t,i}$ ("hard" max)

A few remarkable examples from Min-Char RNN:

• corpus: LATEX source code of an algebraic geometry book (16MB)

A few remarkable examples from Min-Char RNN:

- corpus: LATEX source code of an algebraic geometry book (16MB)
- generate source code that almost complies

For $\bigoplus_{n=1,\dots,m}$ where $\mathcal{L}_{m_*}=0$, hence we can find a closed subset \mathcal{H} in \mathcal{H} and any sets \mathcal{F} on X, U is a closed immersion of S, then $U\to T$ is a separated algebraic space.

Proof. Proof of (1). It also start we get

$$S = \operatorname{Spec}(R) = U \times_X U \times_X U$$

and the comparicoly in the fibre product covering we have to prove the lemma generated by $\prod Z \times_U U \to V$. Consider the maps M along the set of points Sch_{Ippf} and $U \to U$ is the fibre category of S in U in Section, ?? and the fact that any U affine, see Morphisms, Lemma ??. Hence we obtain a scheme S and any open subset $W \subset U$ in Sh(G) such that $Spec(R') \to S$ is smooth or an

$$U = \bigcup U_i \times_{S_i} U_i$$

which has a nonzero morphism we may assume that f_i is of finite presentation over S. We claim that $\mathcal{O}_{X,x'}$ is a scheme where x,x', $s'' \in S'$ such that $\mathcal{O}_{X,x'} \to \mathcal{O}_{X,x'}$ is esparated. By Algebra, Lemma ?? we can define a map of complexes $\operatorname{GL}_{S'}(x'/S'')$ and we win.

To prove study we see that $\mathcal{F}|_U$ is a covering of \mathcal{X}' , and \mathcal{T}_i is an object of $\mathcal{F}_{X/S}$ for i>0 and \mathcal{F}_p exists and let \mathcal{F}_i be a presheaf of \mathcal{O}_X -modules on \mathcal{C} as a \mathcal{F} -module. In particular $\mathcal{F}=U/\mathcal{F}$ we have to show that

$$\widetilde{M}^{\bullet} = \mathcal{I}^{\bullet} \otimes_{Spec(k)} \mathcal{O}_{S,s} - i_{X}^{-1} \mathcal{F})$$

A few remarkable examples from Min-Char RNN:

- corpus: LATEX source code of an algebraic geometry book (16MB)
- generate source code that almost complies
- the model understands complex syntactic structures

For $\bigoplus_{n=1,\ldots,m}$ where $\mathcal{L}_{m_s}=0$, hence we can find a closed subset \mathcal{H} in \mathcal{H} and any sets \mathcal{F} on X,U is a closed immersion of S, then $U\to T$ is a separated algebraic space.

Proof. Proof of (1). It also start we get

$$S = \operatorname{Spec}(R) = U \times_X U \times_X U$$

and the comparicoly in the fibre product covering we have to prove the lemma generated by $\prod Z \times_U U \to V$. Consider the maps M along the set of points Sch_{Ippf} and $U \to U$ is the fibre category of S in U in Section, ?? and the fact that any U affine, see Morphisms, Lemma ??. Hence we obtain a scheme S and any open subset $W \subset U$ in Sh(G) such that $Spec(R') \to S$ is smooth or an

$$U = \bigcup U_i \times_{S_i} U_i$$

which has a nonzero morphism we may assume that f_i is of finite presentation over S. We claim that $\mathcal{O}_{X,x'}$ is a scheme where $x, x', s'' \in S'$ such that $\mathcal{O}_{X,x'} \to \mathcal{O}_{X',x'}$ is separated. By Algebra, Lemma ?? we can define a map of complexes $\operatorname{GL}_{S'}(x'/S'')$ and we win.

To prove study we see that $\mathcal{F}_{|U|}$ is a covering of \mathcal{X}' , and \mathcal{T}_i is an object of $\mathcal{F}_{X/S}$ for i>0 and \mathcal{F}_p exists and let \mathcal{F}_i be a presheaf of \mathcal{O}_X -modules on \mathcal{C} as a \mathcal{F} -module. In particular $\mathcal{F}=U/\mathcal{F}$ we have to show that

$$\widetilde{M}^{\bullet} = \mathcal{I}^{\bullet} \otimes_{Spec(k)} \mathcal{O}_{S,s} - i_{X}^{-1} \mathcal{F})$$

Generation after training

A few remarkable examples from Min-Char RNN:

• corpus: Linux source code (474MB of C code); 10M parameters

Generation after training

A few remarkable examples from Min-Char RNN:

- corpus: Linux source code (474MB of C code); 10M parameters
- generate codes with very few syntactic errors

```
static int indicate_policy(void)
 int error:
 if (fd == MARN EPT) {
   if (ss->segment < mem_total)</pre>
     unblock graph and set blocked();
   else
     ret = 1;
   goto bail;
 segaddr = in_SB(in.addr);
 selector = seg / 16;
 setup works = true:
 for (i = 0; i < blocks; i++) {
   sea = buf[i++1:
   bpf = bd->bd.next + i * search;
   if (fd) {
      current = blocked;
 rw->name = "Getjbbregs";
 bprm_self_clearl(&iv->version);
 regs->new = blocks[(BPF STATS << info->historidac)] | PFMR CLOBATHINC SECONDS << 12;
 return segtable:
```

Generation after training

A few remarkable examples from Min-Char RNN:

- corpus: Linux source code (474MB of C code); 10M parameters
- generate codes with very few syntactic errors
- uses strings/pointers properly, open/close brackets correctly, good indentation, even add comments

```
static int indicate_policy(void)
 int error:
 if (fd == MARN EPT) {
   if (ss->segment < mem_total)</pre>
      unblock graph and set blocked();
      ret = 1:
   goto bail:
 segaddr = in_SB(in.addr);
 selector = seg / 16;
 setup works = true:
 for (i = 0; i < blocks; i++) {
   sea = buf[i++1:
   bpf = bd->bd.next + i * search;
   if (fd) {
      current = blocked;
 rw->name = "Getjbbregs";
 borm self clearl(&iv->version):
 regs->new = blocks[(BPF STATS << info->historidac)] | PFMR CLOBATHINC SECONDS << 12;
 return segtable:
```

Some neurons (entries of the hidden state h) are quite **interpretable** (even though most are not)

Some neurons (entries of the hidden state h) are quite **interpretable** (even though most are not)

• can visualize this by coloring the input character based on the value of this neuron (red = large value, blue = small value)

Some neurons (entries of the hidden state h) are quite **interpretable** (even though most are not)

 can visualize this by coloring the input character based on the value of this neuron (red = large value, blue = small value)

```
The sole importance of the crossing of the Berezina lies in the fact that it plainly and indubitably proved the fallacy of all the plans for cutting off the enemy's retreat and the soundness of the only possible line of action-the one Kutuzov and the general mass of the army demanded--namely, simply to follow the enemy up. The French crowd fled at a continually increasing speed and all its energy was directed to reaching its goal. It fled like a wounded animal and it was impossible to block its path. This was shown not so much by the arrangements it made for crossing as by what took place at the bridges. When the bridges broke down, unarmed soldiers, people from Moscow and women with children who were with the French transport, all--carried on by vis inertiae--pressed forward into boats and into the ice-covered water and did not, surrender.
```

A neuron sensitive to the **position** in line

Some neurons (entries of the hidden state h) are quite **interpretable** (even though most are not)

 can visualize this by coloring the input character based on the value of this neuron (red = large value, blue = small value)

```
"You mean to imply that I have nothing to eat out of... On the contrary, I can supply you with everything even if you want to give dinner parties," warmly replied chichagov, who tried by every word he spoke to prove his own rectitude and therefore imagined Kutuzov to be animated by the same desire.

Kutuzov, shrugging his shoulders, replied with his subtle penetrating smile: "I meant merely to say what I said."
```

A neuron that is turned on inside **quotes**

Some neurons (entries of the hidden state h) are quite **interpretable** (even though most are not)

 can visualize this by coloring the input character based on the value of this neuron (red = large value, blue = small value)

```
static int __dequeue_signal(struct sigpending *pending,
    siginfo_t *info)

int sig = next_signal(pending, mask);
if (sig) {
    if (current->notifier) {
        if (siqismember(current->notifier_mask, sig)) {
            if (!(current->notifier)(current->notifier_data)) {
                clear_thread_flag(TIF_SIGPENDING);
            return 0;
    }
}
collect_signal(sig, pending, info);
}
return sig;
}
```

A neuron that is turned on inside if statements

Some neurons (entries of the hidden state h) are quite **interpretable** (even though most are not)

 can visualize this by coloring the input character based on the value of this neuron (red = large value, blue = small value)

```
#ifdef CONFIG_AUDITSYSCALL
static inline int audit_match_class_bits(int class, u32 *mask)

int i;
    if (classes[class]) {
    for (i = 0; i < AUDIT_BITMASK_SIZE; i++)
        if (mask[i] & classes[class][i])
        return 0;
    }
}

return 1;
}</pre>
```

A neuron sensitive to the depth of an expression

Instead of using a character or a word as each \boldsymbol{x} , often use a token (word or sub-word)

Instead of using a character or a word as each \boldsymbol{x} , often use a token (word or sub-word)

• "apple" is a token

Instead of using a character or a word as each $m{x}$, often use a ${\color{blue} {\sf token}}$ (word or sub-word)

- "apple" is a token
- "unbelievable" is 3 tokens ("un", "believ", "able")

Instead of using a character or a word as each $m{x}$, often use a ${\color{blue} {\sf token}}$ (word or sub-word)

- "apple" is a token
- "unbelievable" is 3 tokens ("un", "believ", "able")
- can reduce the size of vocabulary

Instead of using a character or a word as each \boldsymbol{x} , often use a token (word or sub-word)

- "apple" is a token
- "unbelievable" is 3 tokens ("un", "believ", "able")
- can reduce the size of vocabulary

Directly applying backpropagation to RNN leads to vanishing/exploding gradient issues when T is large

Instead of using a character or a word as each $m{x}$, often use a ${\color{blue} {\sf token}}$ (word or sub-word)

- "apple" is a token
- "unbelievable" is 3 tokens ("un", "believ", "able")
- can reduce the size of vocabulary

Directly applying backpropagation to RNN leads to vanishing/exploding gradient issues when T is large

 $oldsymbol{oldsymbol{W}}$ is applied T times at the end of the sequence (so roughly $oldsymbol{W}^T$)

Instead of using a character or a word as each \boldsymbol{x} , often use a token (word or sub-word)

- "apple" is a token
- "unbelievable" is 3 tokens ("un", "believ", "able")
- can reduce the size of vocabulary

Directly applying backpropagation to RNN leads to vanishing/exploding gradient issues when T is large

- $oldsymbol{oldsymbol{W}}$ is applied T times at the end of the sequence (so roughly $oldsymbol{W}^T$)
- some fixes: Long Short-Term Memory (LSTM) and Gated Recurrent Units (GRU)

Outline

- Review of last lecture
- 2 Recurrent Neural Network
- Transformers
 - Self-attention
 - Other components
 - Training and testing

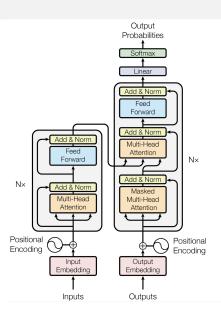
Transformers

Issues of RNN: must compress all previous info into a single state h

Transformers

Issues of RNN: must compress all previous info into a single state h

A solution that dominates all other models currently: **transformers**



Acknowledgements

Very useful resources:

- original paper: "Attention Is All You Need" (200K+ citation by now)
 - https://arxiv.org/pdf/1706.03762
- The Illustrated Transformer (most pictures are from here)
 - https://jalammar.github.io/illustrated-transformer/
- a super cool Nano-GPT visualization
 - https://bbycroft.net/llm
- A Multiscale Visualization of Attention
 - https://arxiv.org/pdf/1906.05714

Example: "The animal didn't cross the street because it was too tired"

Does "it" refer to "animal" or "street"?

- Does "it" refer to "animal" or "street"?
- trivial for human, but how to design a model that understands this?

- Does "it" refer to "animal" or "street"?
- trivial for human, but how to design a model that understands this?
- intuitively, when looking at the word "it", the model should pay attention to the word "animal"

- Does "it" refer to "animal" or "street"?
- trivial for human, but how to design a model that understands this?
- intuitively, when looking at the word "it", the model should pay attention to the word "animal"
- An attention head does exactly this

An attention head

• takes a sequence of inputs $x_1, \ldots, x_T \in \mathbb{R}^d$ and outputs another sequence $z_1, \ldots, z_T \in \mathbb{R}^{d_v}$ (similar to **hidden states** of RNN)

- takes a sequence of inputs $x_1, \ldots, x_T \in \mathbb{R}^d$ and outputs another sequence $z_1, \ldots, z_T \in \mathbb{R}^{d_v}$ (similar to **hidden states** of RNN)
- parametrized by three matrices (and corresponding biases, omitted for simplicity): $W_O \in \mathbb{R}^{d \times d_k}, W_K \in \mathbb{R}^{d \times d_k}, W_V \in \mathbb{R}^{d \times d_v}$

- takes a sequence of inputs $x_1, \ldots, x_T \in \mathbb{R}^d$ and outputs another sequence $z_1, \ldots, z_T \in \mathbb{R}^{d_v}$ (similar to **hidden states** of RNN)
- parametrized by three matrices (and corresponding biases, omitted for simplicity): $W_Q \in \mathbb{R}^{d \times d_k}, W_K \in \mathbb{R}^{d \times d_k}, W_V \in \mathbb{R}^{d \times d_v}$
 - ullet computes a query vector for each input $oldsymbol{x}_t$ as $oldsymbol{q}_t = oldsymbol{W}_O^ op oldsymbol{x}_t \in \mathbb{R}^{d_k}$

- takes a sequence of inputs $x_1, \ldots, x_T \in \mathbb{R}^d$ and outputs another sequence $z_1, \ldots, z_T \in \mathbb{R}^{d_v}$ (similar to **hidden states** of RNN)
- parametrized by three matrices (and corresponding biases, omitted for simplicity): $W_Q \in \mathbb{R}^{d \times d_k}, W_K \in \mathbb{R}^{d \times d_k}, W_V \in \mathbb{R}^{d \times d_v}$
 - ullet computes a query vector for each input $oldsymbol{x}_t$ as $oldsymbol{q}_t = oldsymbol{W}_O^ op oldsymbol{x}_t \in \mathbb{R}^{d_k}$
 - ullet computes a key vector for each input $oldsymbol{x}_t$ as $oldsymbol{k}_t = oldsymbol{W}_K^ op oldsymbol{x}_t \in \mathbb{R}^{d_k}$

- takes a sequence of inputs $x_1, \ldots, x_T \in \mathbb{R}^d$ and outputs another sequence $z_1, \ldots, z_T \in \mathbb{R}^{d_v}$ (similar to **hidden states** of RNN)
- parametrized by three matrices (and corresponding biases, omitted for simplicity): $W_Q \in \mathbb{R}^{d \times d_k}, W_K \in \mathbb{R}^{d \times d_k}, W_V \in \mathbb{R}^{d \times d_v}$
 - ullet computes a query vector for each input $oldsymbol{x}_t$ as $oldsymbol{q}_t = oldsymbol{W}_O^ op oldsymbol{x}_t \in \mathbb{R}^{d_k}$
 - ullet computes a key vector for each input $oldsymbol{x}_t$ as $oldsymbol{k}_t = oldsymbol{W}_K^ op oldsymbol{x}_t \in \mathbb{R}^{d_k}$
 - ullet computes a value vector for each input $oldsymbol{x}_t$ as $oldsymbol{v}_t = oldsymbol{W}_V^ op oldsymbol{x}_t \in \mathbb{R}^{d_v}$

- takes a sequence of inputs $x_1, \ldots, x_T \in \mathbb{R}^d$ and outputs another sequence $z_1, \ldots, z_T \in \mathbb{R}^{d_v}$ (similar to **hidden states** of RNN)
- parametrized by three matrices (and corresponding biases, omitted for simplicity): $W_Q \in \mathbb{R}^{d \times d_k}, W_K \in \mathbb{R}^{d \times d_k}, W_V \in \mathbb{R}^{d \times d_v}$
 - ullet computes a query vector for each input $oldsymbol{x}_t$ as $oldsymbol{q}_t = oldsymbol{W}_Q^ op oldsymbol{x}_t \in \mathbb{R}^{d_k}$
 - ullet computes a key vector for each input $oldsymbol{x}_t$ as $oldsymbol{k}_t = oldsymbol{W}_K^ op oldsymbol{x}_t \in \mathbb{R}^{d_k}$
 - ullet computes a value vector for each input $oldsymbol{x}_t$ as $oldsymbol{v}_t = oldsymbol{W}_V^ op oldsymbol{x}_t \in \mathbb{R}^{d_v}$
- ullet the output $oldsymbol{z}_t$ is the "answer" to the query of $oldsymbol{q}_t$

Input	Thinking	Machines	
Embedding	X ₁	X ₂	
Queries	q ₁	q ₂	Wa
Keys	k 1	k ₂	Wĸ
Values	V1	V ₂	wv

The output z_t is the "answer" to the query of q_t . How?

The output z_t is the "answer" to the query of q_t . How?

ullet imagine: you make a Google query (q_t) ,

The output z_t is the "answer" to the query of q_t . How?

ullet imagine: you make a Google query (q_t) , and it returns a list of website titles $({m k}_{1:T})$;

The output z_t is the "answer" to the query of q_t . How?

• imagine: you make a Google query (q_t) , and it returns a list of website titles $(k_{1:T})$; clicking a title (k_{τ}) leads you to a website (v_{τ}) .

- imagine: you make a Google query (q_t) , and it returns a list of website titles $(k_{1:T})$; clicking a title (k_{τ}) leads you to a website (v_{τ}) .
- ullet You then summarize the answer using all websites $(oldsymbol{v}_{1:T})$,

- imagine: you make a Google query (q_t) , and it returns a list of website titles $(k_{1:T})$; clicking a title (k_{τ}) leads you to a website (v_{τ}) .
- ullet You then summarize the answer using all websites $(v_{1:T})$, each with a different **weight** based on how relevant/close its title is to your query

- imagine: you make a Google query (q_t) , and it returns a list of website titles $(k_{1:T})$; clicking a title (k_{τ}) leads you to a website (v_{τ}) .
- ullet You then summarize the answer using all websites $(oldsymbol{v}_{1:T})$, each with a different **weight** based on how relevant/close its title is to your query
- formally, the final answer z_t is the **weighted sum** of v_1, \ldots, v_T ,

- imagine: you make a Google query (q_t) , and it returns a list of website titles $(k_{1:T})$; clicking a title (k_{τ}) leads you to a website (v_{τ}) .
- ullet You then summarize the answer using all websites $(oldsymbol{v}_{1:T})$, each with a different weight based on how relevant/close its title is to your query
- formally, the final answer z_t is the **weighted sum** of v_1, \ldots, v_T , with weights computed via

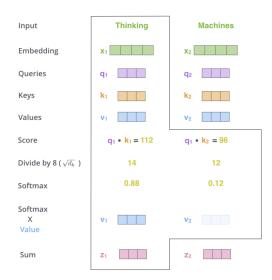
$$\mathsf{softmax}\left(\frac{{q_t}^\top {\color{red} {k_1}}}{\sqrt{d_k}}, \ldots, \frac{{q_t}^\top {\color{red} {k_T}}}{\sqrt{d_k}}\right)$$

The output z_t is the "answer" to the query of q_t . How?

- imagine: you make a Google query (q_t) , and it returns a list of website titles $(k_{1:T})$; clicking a title (k_{τ}) leads you to a website (v_{τ}) .
- ullet You then summarize the answer using all websites $(v_{1:T})$, each with a different weight based on how relevant/close its title is to your query
- ullet formally, the final answer z_t is the **weighted sum** of v_1, \dots, v_T , with weights computed via

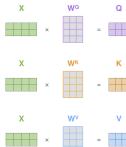
$$\operatorname{softmax}\left(\frac{{q_t}^{\top} \textcolor{red}{k_1}}{\sqrt{d_k}}, \dots, \frac{{q_t}^{\top} \textcolor{red}{k_T}}{\sqrt{d_k}}\right)$$

where $q_t^{ op} k_{ au}$ is the **attention score** from input x_t to input $x_{ au}$

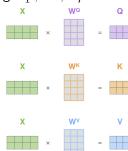


Matrix notation:

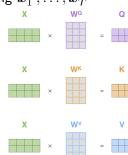
ullet input matrix $oldsymbol{X} \in \mathbb{R}^{T imes d}$, obtained by stacking $oldsymbol{x}_1^ op, \dots, oldsymbol{x}_T^ op$



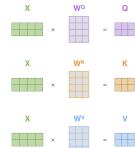
- ullet input matrix $oldsymbol{X} \in \mathbb{R}^{T imes d}$, obtained by stacking $oldsymbol{x}_1^ op, \dots, oldsymbol{x}_T^ op$
- ullet query matrix $oldsymbol{Q} = oldsymbol{X} oldsymbol{W}_O \in \mathbb{R}^{T imes d_k}$



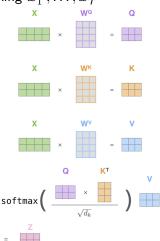
- ullet input matrix $oldsymbol{X} \in \mathbb{R}^{T imes d}$, obtained by stacking $oldsymbol{x}_1^ op, \dots, oldsymbol{x}_T^ op$
- ullet query matrix $oldsymbol{Q} = oldsymbol{X} oldsymbol{W}_Q \in \mathbb{R}^{T imes d_k}$
- ullet key matrix $oldsymbol{K} = oldsymbol{X} oldsymbol{W}_K \in \mathbb{R}^{T imes d_k}$



- ullet input matrix $oldsymbol{X} \in \mathbb{R}^{T imes d}$, obtained by stacking $oldsymbol{x}_1^ op, \dots, oldsymbol{x}_T^ op$
- ullet query matrix $oldsymbol{Q} = oldsymbol{X} oldsymbol{W}_O \in \mathbb{R}^{T imes d_k}$
- ullet key matrix $oldsymbol{K} = oldsymbol{X} oldsymbol{W}_K \in \mathbb{R}^{T imes d_k}$
- ullet value matrix $oldsymbol{V} = oldsymbol{X} oldsymbol{W}_V \in \mathbb{R}^{T imes d_v}$



- ullet input matrix $oldsymbol{X} \in \mathbb{R}^{T imes d}$, obtained by stacking $oldsymbol{x}_1^ op, \dots, oldsymbol{x}_T^ op$
- ullet query matrix $oldsymbol{Q} = oldsymbol{X} oldsymbol{W}_O \in \mathbb{R}^{T imes d_k}$
- ullet key matrix $oldsymbol{K} = oldsymbol{X} oldsymbol{W}_K \in \mathbb{R}^{T imes d_k}$
- ullet value matrix $oldsymbol{V} = oldsymbol{X} oldsymbol{W}_V \in \mathbb{R}^{T imes d_v}$
- ullet attention score matrix $oldsymbol{Q} oldsymbol{K}^ op \in \mathbb{R}^{T imes T}$

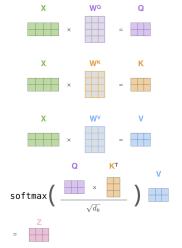


Matrix notation:

- ullet input matrix $oldsymbol{X} \in \mathbb{R}^{T imes d}$, obtained by stacking $oldsymbol{x}_1^ op, \dots, oldsymbol{x}_T^ op$
- ullet query matrix $oldsymbol{Q} = oldsymbol{X} oldsymbol{W}_O \in \mathbb{R}^{T imes d_k}$
- ullet key matrix $oldsymbol{K} = oldsymbol{X} oldsymbol{W}_K \in \mathbb{R}^{T imes d_k}$
- ullet value matrix $oldsymbol{V} = oldsymbol{X} oldsymbol{W}_V \in \mathbb{R}^{T imes d_v}$
- ullet attention score matrix $oldsymbol{Q} oldsymbol{K}^ op \in \mathbb{R}^{T imes T}$
- ullet output matrix $oldsymbol{Z} \in \mathbb{R}^{T imes d_v}$ is

$$\operatorname{softmax}\left(rac{oldsymbol{Q}oldsymbol{K}^ op}{\sqrt{d_k}}
ight)oldsymbol{V}$$

where softmax is applied row-wise



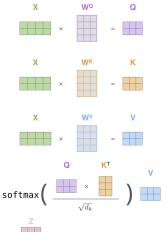
Matrix notation:

- ullet input matrix $oldsymbol{X} \in \mathbb{R}^{T imes d}$, obtained by stacking $oldsymbol{x}_1^ op, \dots, oldsymbol{x}_T^ op$
- ullet query matrix $oldsymbol{Q} = oldsymbol{X} oldsymbol{W}_Q \in \mathbb{R}^{T imes d_k}$
- ullet key matrix $oldsymbol{K} = oldsymbol{X} oldsymbol{W}_K \in \mathbb{R}^{T imes d_k}$
- ullet value matrix $oldsymbol{V} = oldsymbol{X} oldsymbol{W}_V \in \mathbb{R}^{T imes d_v}$
- ullet attention score matrix $oldsymbol{Q} oldsymbol{K}^ op \in \mathbb{R}^{T imes T}$
- ullet output matrix $oldsymbol{Z} \in \mathbb{R}^{T imes d_v}$ is

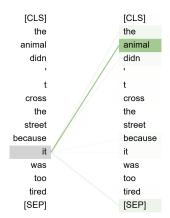
$$\operatorname{softmax}\left(rac{oldsymbol{Q}oldsymbol{K}^ op}{\sqrt{d_k}}
ight)oldsymbol{V}$$

where softmax is applied row-wise

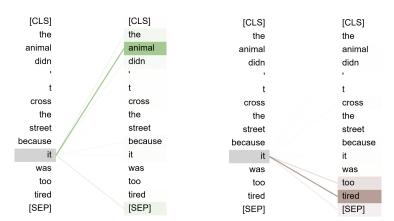
 $O(T^2)$ complexity (ignoring d, d_k, d_v)



- the darker the color, the larger the attention score
- "it" attends to "animal" in one head,



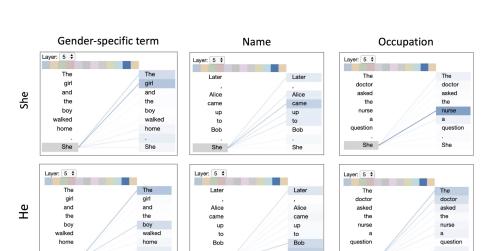
- the darker the color, the larger the attention score
- "it" attends to "animal" in one head, and "tired" in another head



He

More examples

He



He

He

He

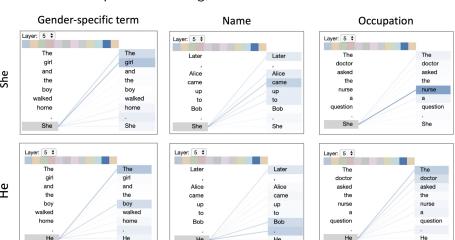
He

He

More examples

He

• all from unsupervised learning; no one tells the model to learn these!



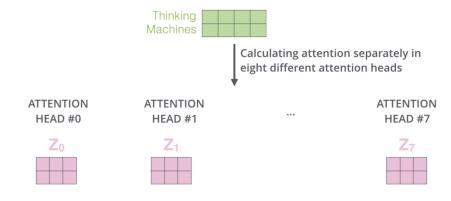
He

He

Multi-head attention

Pass $oldsymbol{X}$ to multiple attention-heads, each with different parameters

X



Multi-head attention (cont.)

Concatenate outputs of different heads and then project again

1) Concatenate all the attention heads



2) Multiply with a weight matrix W^o that was trained jointly with the model

Χ



3) The result would be the Z matrix that captures information from all the attention heads. We can send this forward to the FFNN



Multi-head attention (cont.)

Concatenate outputs of different heads and then project again

ullet final output dimension is $\mathbb{R}^{T imes d}$, same as inputs $oldsymbol{X}$

1) Concatenate all the attention heads



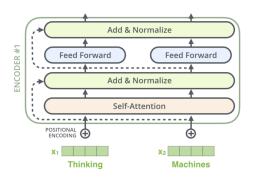
2) Multiply with a weight matrix W^o that was trained jointly with the model

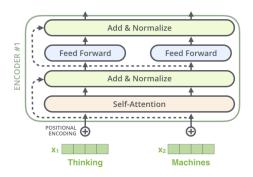
Χ

Wo

3) The result would be the Z matrix that captures information from all the attention heads. We can send this forward to the FFNN

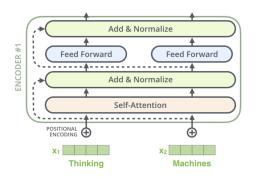
Z





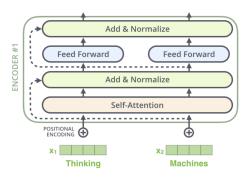
Positional encoding

fix the issue that attention-head does not have positional info



Positional encoding

- fix the issue that attention-head does not have positional info
- ullet via a **positional embedding matrix** $oldsymbol{E}_P \in \mathbb{R}^{d imes T}$, fixed or learned,

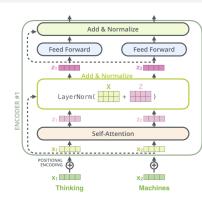


Positional encoding

- fix the issue that attention-head does not have positional info
- ullet via a **positional embedding matrix** $oldsymbol{E}_P \in \mathbb{R}^{d imes T}$, fixed or learned,
- ullet $oldsymbol{x}_t \leftarrow oldsymbol{x}_t + oldsymbol{E}_P e_t$, i.e., add the t-th column of $oldsymbol{E}_P$ to $oldsymbol{x}_t$

Two more components to stabilize and speed up training:

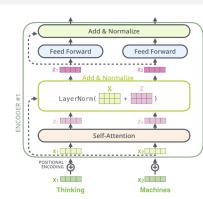
1. Residual pathway



Two more components to stabilize and speed up training:

1. Residual pathway

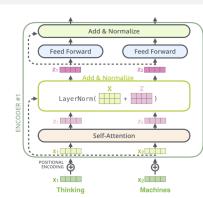
ullet add input to output, $oldsymbol{Z} \leftarrow oldsymbol{Z} + oldsymbol{X}$



Two more components to stabilize and speed up training:

1. Residual pathway

- ullet add input to output, $oldsymbol{Z} \leftarrow oldsymbol{Z} + oldsymbol{X}$
- an idea from Residual Networks to deal with vanishing gradients

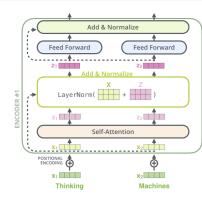


Two more components to stabilize and speed up training:

1. Residual pathway

- ullet add input to output, $oldsymbol{Z} \leftarrow oldsymbol{Z} + oldsymbol{X}$
- an idea from Residual Networks to deal with vanishing gradients

2. Layer normalization:



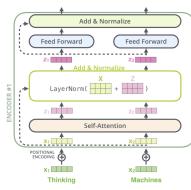
Two more components to stabilize and speed up training:

1. Residual pathway

- ullet add input to output, $oldsymbol{Z} \leftarrow oldsymbol{Z} + oldsymbol{X}$
- an idea from Residual Networks to deal with vanishing gradients

2. Layer normalization:

ullet for each $oldsymbol{z}_t \in \mathbb{R}^d$, normalize it to zero-mean and unit-variance (across features)

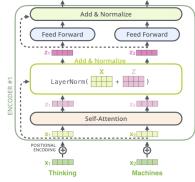


Two more components to stabilize and speed up training:

1. Residual pathway

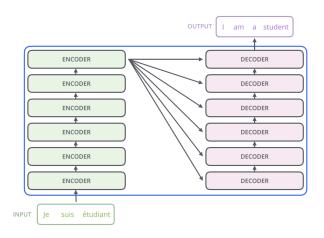
- ullet add input to output, $oldsymbol{Z} \leftarrow oldsymbol{Z} + oldsymbol{X}$
- an idea from Residual Networks to deal with vanishing gradients

2. Layer normalization:



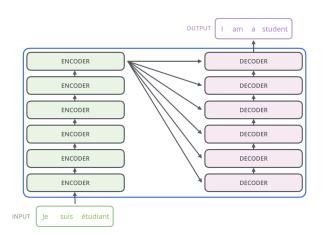
- ullet for each $oldsymbol{z}_t \in \mathbb{R}^d$, normalize it to zero-mean and unit-variance (across features)
- similar but different from batch normalization (where you normalize each feature to zero-mean and unit-variance across samples)

Zooming out: stacking encoders and decoders



Encoder: summarizes the input into a useful representation

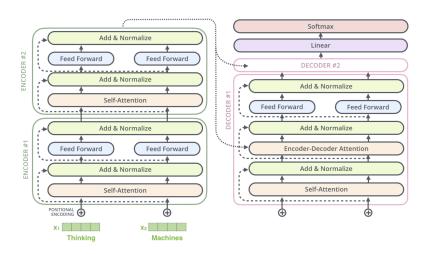
Zooming out: stacking encoders and decoders



Encoder: summarizes the input into a useful representation

Decoder: generates outputs

A closer look at decoders



Extra component: encoder-decoder attention

Encoder-decoder attention

Same idea as self-attention, except key and value matrices are computed using output Z_{enc} of the final encoder:

Encoder-decoder attention

Same idea as self-attention, except key and value matrices are computed using output Z_{enc} of the final encoder:

ullet query matrix $oldsymbol{Q} = oldsymbol{X} oldsymbol{W}_Q \in \mathbb{R}^{T_{\mathsf{dec}} imes d_k}$ (as usual)

Encoder-decoder attention

Same idea as self-attention, except key and value matrices are computed using output Z_{enc} of the final encoder:

- ullet query matrix $oldsymbol{Q} = oldsymbol{X} oldsymbol{W}_Q \in \mathbb{R}^{T_{\mathsf{dec}} imes d_k}$ (as usual)
- ullet key matrix $oldsymbol{K}_{\mathsf{enc}} = oldsymbol{Z}_{\mathsf{enc}} oldsymbol{W}_K \in \mathbb{R}^{oldsymbol{T}_{\mathsf{enc}} imes d_k}$

Encoder-decoder attention

Same idea as self-attention, except key and value matrices are computed using output Z_{enc} of the final encoder:

- ullet query matrix $oldsymbol{Q} = oldsymbol{X} oldsymbol{W}_Q \in \mathbb{R}^{T_{\mathsf{dec}} imes d_k}$ (as usual)
- ullet key matrix $oldsymbol{K}_{\mathsf{enc}} = oldsymbol{Z}_{\mathsf{enc}} oldsymbol{W}_K \in \mathbb{R}^{oldsymbol{T}_{\mathsf{enc}} imes d_k}$
- ullet value matrix $oldsymbol{V}_{\mathsf{enc}} = oldsymbol{Z}_{\mathsf{enc}} oldsymbol{W}_V \in \mathbb{R}^{oldsymbol{T}_{\mathsf{enc}} imes d_v}$

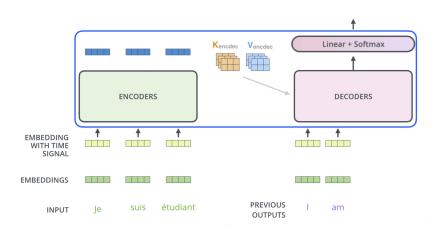
Encoder-decoder attention

Same idea as self-attention, except key and value matrices are computed using output Z_{enc} of the final encoder:

- ullet query matrix $oldsymbol{Q} = oldsymbol{X} oldsymbol{W}_Q \in \mathbb{R}^{T_{\mathsf{dec}} imes d_k}$ (as usual)
- ullet key matrix $oldsymbol{K}_{\mathsf{enc}} = oldsymbol{Z}_{\mathsf{enc}} oldsymbol{W}_K \in \mathbb{R}^{oldsymbol{T}_{\mathsf{enc}} imes d_k}$
- ullet value matrix $oldsymbol{V}_{\mathsf{enc}} = oldsymbol{Z}_{\mathsf{enc}} oldsymbol{W}_V \in \mathbb{R}^{oldsymbol{T}_{\mathsf{enc}} imes d_v}$

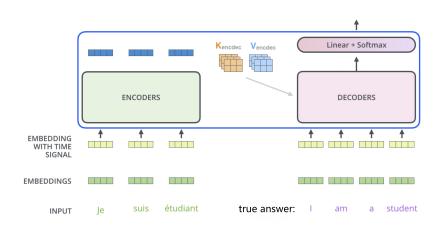
Intuition: find answer from the encoded representation of original inputs

Generating answers



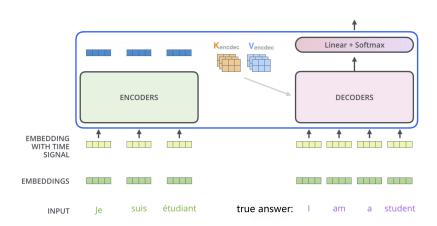
Use previously generated text as inputs of the decoder

Training



Use cross-entropy loss again,

Training



Use cross-entropy loss again, and apply

• teacher forcing: use the true answer as inputs of the decoder

In teaching forcing, to avoid generating a word by peeking at future words, need to use a causal mask in the decoder's self-attention heads:

In teaching forcing, to avoid generating a word by peeking at future words, need to use a causal mask in the decoder's self-attention heads:

$$oldsymbol{Q} oldsymbol{K}^ op \leftarrow oldsymbol{Q} oldsymbol{K}^ op + oldsymbol{M}$$

In teaching forcing, to avoid generating a word by peeking at future words, need to use a causal mask in the decoder's self-attention heads:

$$oldsymbol{Q} oldsymbol{K}^ op \leftarrow oldsymbol{Q} oldsymbol{K}^ op + oldsymbol{M}$$

where

$$\boldsymbol{M} = \begin{bmatrix} 0 & -\infty & -\infty & \cdots & -\infty \\ 0 & 0 & -\infty & \cdots & -\infty \\ 0 & 0 & 0 & \cdots & -\infty \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & 0 \end{bmatrix} \in \mathbb{R}^{T_{\mathsf{dec}} \times T_{\mathsf{dec}}}$$

In teaching forcing, to avoid generating a word by peeking at future words, need to use a causal mask in the decoder's self-attention heads:

$$oldsymbol{Q} oldsymbol{K}^ op \leftarrow oldsymbol{Q} oldsymbol{K}^ op + oldsymbol{M}$$

where

$$oldsymbol{M} = egin{bmatrix} 0 & -\infty & -\infty & \cdots & -\infty \ 0 & 0 & -\infty & \cdots & -\infty \ 0 & 0 & 0 & \cdots & -\infty \ dots & dots & dots & dots & dots \ 0 & 0 & 0 & \cdots & 0 \end{bmatrix} \in \mathbb{R}^{T_{\mathsf{dec}} \times T_{\mathsf{dec}}}$$

so a word at position t never attends to words at positions > t

In teaching forcing, to avoid generating a word by peeking at future words, need to use a causal mask in the decoder's self-attention heads:

$$oldsymbol{Q} oldsymbol{K}^ op \leftarrow oldsymbol{Q} oldsymbol{K}^ op + oldsymbol{M}$$

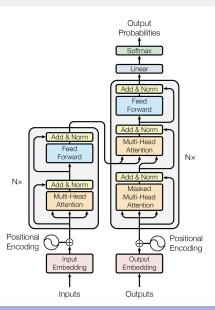
where

$$\boldsymbol{M} = \begin{bmatrix} 0 & -\infty & -\infty & \cdots & -\infty \\ 0 & 0 & -\infty & \cdots & -\infty \\ 0 & 0 & 0 & \cdots & -\infty \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & 0 \end{bmatrix} \in \mathbb{R}^{T_{\mathsf{dec}} \times T_{\mathsf{dec}}}$$

so a word at position t never attends to words at positions > t

Q: should we use causal mask for encoder-decoder attention heads?

That's it!



Large language models

Large language models (LLMs) are all based on transformers

• estimated #parameters for GPT5: trillions



Large language models

Large language models (LLMs) are all based on transformers

• estimated #parameters for GPT5: trillions



A cool 3D visualization of a nano-GPT: https://bbycroft.net/llm

Training Large language models

Unsupervised pre-training

 via next word prediction using a huge training set (e.g., the entire internet)

Training Large language models

Unsupervised pre-training

 via next word prediction using a <u>huge</u> training set (e.g., the entire internet)

Fine-tuning

 using a labeled dataset for a specific task (translation, question answering, etc.)

Training Large language models

Unsupervised pre-training

 via next word prediction using a huge training set (e.g., the entire internet)

Fine-tuning

 using a labeled dataset for a specific task (translation, question answering, etc.)

Reinforcement Learning with Human Feedback (RLHF)

- get preference feedback from human: which answer is better?
- more on this in the next two weeks