



Attention

CISC 7026 - Introduction to Deep Learning

Steven Morad

University of Macau

Admin	2
Review	4
Short Intro to Graph Neural Networks	5
Attention	14
Keys and Queries	46
Self Attention	61
Guest Lecture - Dr. Matteo Bettini	70

Admin

Admin

Last exam next week

Admin

Last exam next week

Preliminary questions (still making exam, may change)

Admin

Last exam next week

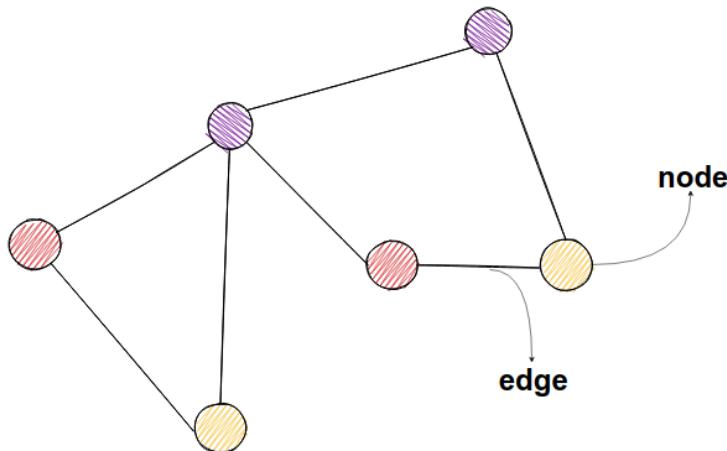
Preliminary questions (still making exam, may change)

- 1 Question perceptron autoencoder and objective
- 1 Question rewriting functions as recurrent function
- 1 Question evaluating scans
- 1 Question log likelihood derivation
- 1 Question key/query attention

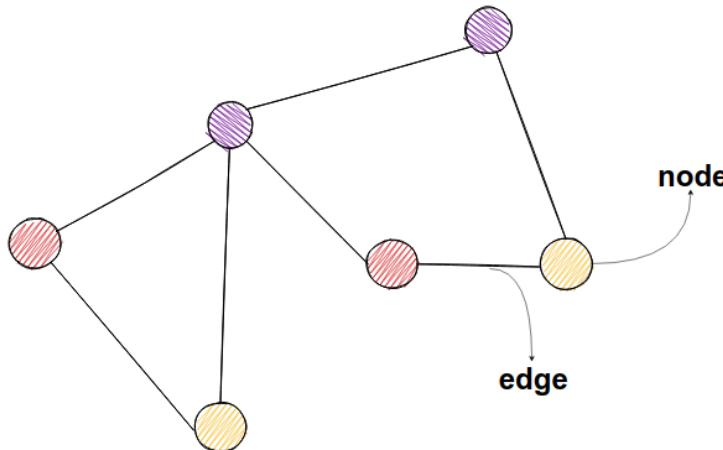
Review

Short Intro to Graph Neural Networks

Short Intro to Graph Neural Networks

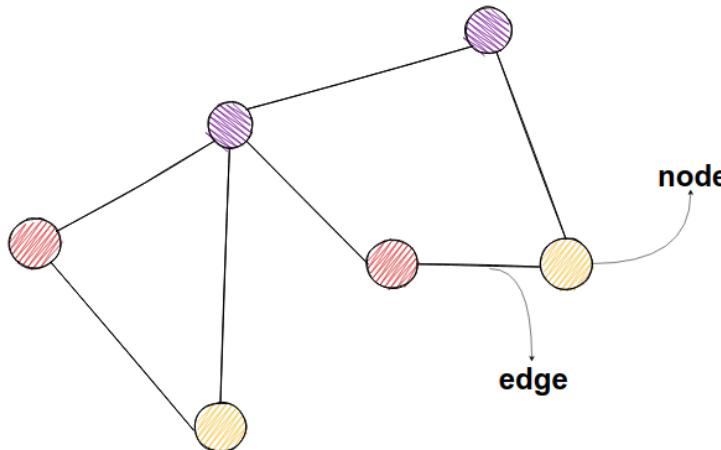


Short Intro to Graph Neural Networks



A **node** is a vector of information

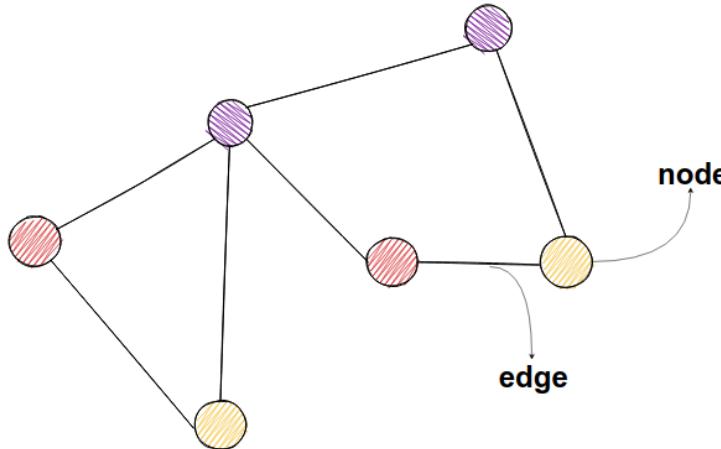
Short Intro to Graph Neural Networks



A **node** is a vector of information

An **edge** connects two nodes

Short Intro to Graph Neural Networks

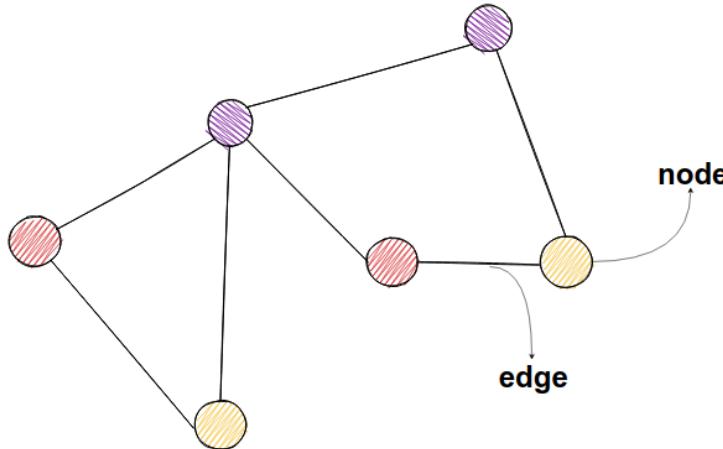


A **node** is a vector of information

An **edge** connects two nodes

If we connect nodes i and j with edge (i, j) , then i and j are **neighbors**

Short Intro to Graph Neural Networks



A **node** is a vector of information

An **edge** connects two nodes

If we connect nodes i and j with edge (i, j) , then i and j are **neighbors**

The **neighborhood** $N(i)$ contains all neighbors of node i

Short Intro to Graph Neural Networks

Let us think of graphs as **signals**

Short Intro to Graph Neural Networks

Let us think of graphs as **signals**

Question: Where did we see signals before?

Short Intro to Graph Neural Networks

Let us think of graphs as **signals**

Question: Where did we see signals before?

Answer: Convolution

Short Intro to Graph Neural Networks

Let us think of graphs as **signals**

Question: Where did we see signals before?

Answer: Convolution

Rather than time t or space u, v , graphs are a function of neighborhood

Short Intro to Graph Neural Networks

Let us think of graphs as **signals**

Question: Where did we see signals before?

Answer: Convolution

Rather than time t or space u, v , graphs are a function of neighborhood

$$\text{Node } i \quad \mathbf{x}(i) \in \mathbb{R}^{d_x}; \quad i \in 1, \dots, T$$

Short Intro to Graph Neural Networks

Let us think of graphs as **signals**

Question: Where did we see signals before?

Answer: Convolution

Rather than time t or space u, v , graphs are a function of neighborhood

Node i $\mathbf{x}(i) \in \mathbb{R}^{d_x}; \quad i \in 1, \dots, T$

Neighborhood of i $\mathbf{N}(i) = \begin{bmatrix} i \\ j \\ k \\ \vdots \end{bmatrix}; \quad \mathbf{N}(i) \in \mathcal{P}(i); \quad i \in 1, \dots, T$

Short Intro to Graph Neural Networks

Consider a **graph convolution layer**

Short Intro to Graph Neural Networks

Consider a **graph convolution layer**

For a node i , the graph convolution layer is:

Short Intro to Graph Neural Networks

Consider a **graph convolution layer**

For a node i , the graph convolution layer is:

$$f(\mathbf{x}, \mathcal{N}, \boldsymbol{\theta})(i) = \sigma \left(\sum_{j \in \mathcal{N}(i)} \boldsymbol{\theta}^\top \bar{\mathbf{x}}(j) \right)$$

Short Intro to Graph Neural Networks

Consider a **graph convolution layer**

For a node i , the graph convolution layer is:

$$f(\mathbf{x}, \mathcal{N}, \boldsymbol{\theta})(i) = \sigma \left(\sum_{j \in \mathcal{N}(i)} \boldsymbol{\theta}^\top \bar{\mathbf{x}}(j) \right)$$

Combine information from the neighbors of $\mathbf{x}(i)$

Short Intro to Graph Neural Networks

Consider a **graph convolution layer**

For a node i , the graph convolution layer is:

$$f(\mathbf{x}, \mathcal{N}, \boldsymbol{\theta})(i) = \sigma \left(\sum_{j \in \mathcal{N}(i)} \boldsymbol{\theta}^\top \bar{\mathbf{x}}(j) \right)$$

Combine information from the neighbors of $\mathbf{x}(i)$

This is just one node, we use this graph layer for all nodes in the graph

Short Intro to Graph Neural Networks

Apply graph convolution over all nodes in the graph

Short Intro to Graph Neural Networks

Apply graph convolution over all nodes in the graph

$$f(\mathbf{x}, \mathbf{N}, \boldsymbol{\theta}) = \begin{bmatrix} f(\mathbf{x}, \mathbf{N}, \boldsymbol{\theta})(1) \\ f(\mathbf{x}, \mathbf{N}, \boldsymbol{\theta})(2) \\ \vdots \\ f(\mathbf{x}, \mathbf{N}, \boldsymbol{\theta})(T) \end{bmatrix} = \begin{bmatrix} \sigma\left(\sum_{j \in \mathbf{N}(1)} \boldsymbol{\theta}^\top \bar{\mathbf{x}}(j)\right) \\ \sigma\left(\sum_{j \in \mathbf{N}(2)} \boldsymbol{\theta}^\top \bar{\mathbf{x}}(j)\right) \\ \vdots \\ \sigma\left(\sum_{j \in \mathbf{N}(T)} \boldsymbol{\theta}^\top \bar{\mathbf{x}}(j)\right) \end{bmatrix}$$

Short Intro to Graph Neural Networks

Apply graph convolution over all nodes in the graph

$$f(\mathbf{x}, \mathbf{N}, \boldsymbol{\theta}) = \begin{bmatrix} f(\mathbf{x}, \mathbf{N}, \boldsymbol{\theta})(1) \\ f(\mathbf{x}, \mathbf{N}, \boldsymbol{\theta})(2) \\ \vdots \\ f(\mathbf{x}, \mathbf{N}, \boldsymbol{\theta})(T) \end{bmatrix} = \begin{bmatrix} \sigma\left(\sum_{j \in \mathbf{N}(1)} \boldsymbol{\theta}^\top \bar{\mathbf{x}}(j)\right) \\ \sigma\left(\sum_{j \in \mathbf{N}(2)} \boldsymbol{\theta}^\top \bar{\mathbf{x}}(j)\right) \\ \vdots \\ \sigma\left(\sum_{j \in \mathbf{N}(T)} \boldsymbol{\theta}^\top \bar{\mathbf{x}}(j)\right) \end{bmatrix}$$

How does this compare to regular convolution (images, sound, etc)?

Short Intro to Graph Neural Networks

Standard 1D convolution

$$\begin{bmatrix} \sigma\left(\sum_{j=1}^k \theta^\top \bar{x}(j)\right) \\ \sigma\left(\sum_{j=2}^{k+1} \theta^\top \bar{x}(j)\right) \\ \vdots \\ \sigma\left(\sum_{j=T-k}^T \theta^\top \bar{x}(j)\right) \end{bmatrix}$$

Short Intro to Graph Neural Networks

Standard 1D convolution

$$\begin{bmatrix} \sigma\left(\sum_{j=1}^k \boldsymbol{\theta}^\top \bar{\mathbf{x}}(j)\right) \\ \sigma\left(\sum_{j=2}^{k+1} \boldsymbol{\theta}^\top \bar{\mathbf{x}}(j)\right) \\ \vdots \\ \sigma\left(\sum_{j=T-k}^T \boldsymbol{\theta}^\top \bar{\mathbf{x}}(j)\right) \end{bmatrix}$$

Graph convolution

$$\begin{bmatrix} \sigma\left(\sum_{j \in \mathcal{N}(1)} \boldsymbol{\theta}^\top \bar{\mathbf{x}}(j)\right) \\ \sigma\left(\sum_{j \in \mathcal{N}(2)} \boldsymbol{\theta}^\top \bar{\mathbf{x}}(j)\right) \\ \vdots \\ \sigma\left(\sum_{j \in \mathcal{N}(T)} \boldsymbol{\theta}^\top \bar{\mathbf{x}}(j)\right) \end{bmatrix}$$

Short Intro to Graph Neural Networks

Standard 1D convolution

$$\begin{bmatrix} \sigma\left(\sum_{j=1}^k \boldsymbol{\theta}^\top \bar{\mathbf{x}}(j)\right) \\ \sigma\left(\sum_{j=2}^{k+1} \boldsymbol{\theta}^\top \bar{\mathbf{x}}(j)\right) \\ \vdots \\ \sigma\left(\sum_{j=T-k}^T \boldsymbol{\theta}^\top \bar{\mathbf{x}}(j)\right) \end{bmatrix}$$

Graph convolution

$$\begin{bmatrix} \sigma\left(\sum_{j \in \mathcal{N}(1)} \boldsymbol{\theta}^\top \bar{\mathbf{x}}(j)\right) \\ \sigma\left(\sum_{j \in \mathcal{N}(2)} \boldsymbol{\theta}^\top \bar{\mathbf{x}}(j)\right) \\ \vdots \\ \sigma\left(\sum_{j \in \mathcal{N}(T)} \boldsymbol{\theta}^\top \bar{\mathbf{x}}(j)\right) \end{bmatrix}$$

Question: What is the output size of standard convolution?

Short Intro to Graph Neural Networks

Standard 1D convolution

$$\begin{bmatrix} \sigma\left(\sum_{j=1}^k \theta^\top \bar{x}(j)\right) \\ \sigma\left(\sum_{j=2}^{k+1} \theta^\top \bar{x}(j)\right) \\ \vdots \\ \sigma\left(\sum_{j=T-k}^T \theta^\top \bar{x}(j)\right) \end{bmatrix}$$

Graph convolution

$$\begin{bmatrix} \sigma\left(\sum_{j \in N(1)} \theta^\top \bar{x}(j)\right) \\ \sigma\left(\sum_{j \in N(2)} \theta^\top \bar{x}(j)\right) \\ \vdots \\ \sigma\left(\sum_{j \in N(T)} \theta^\top \bar{x}(j)\right) \end{bmatrix}$$

Question: What is the output size of standard convolution?

Answer: $(T - k - 1) \times d_h$

Short Intro to Graph Neural Networks

Standard 1D convolution

$$\begin{bmatrix} \sigma\left(\sum_{j=1}^k \boldsymbol{\theta}^\top \bar{\mathbf{x}}(j)\right) \\ \sigma\left(\sum_{j=2}^{k+1} \boldsymbol{\theta}^\top \bar{\mathbf{x}}(j)\right) \\ \vdots \\ \sigma\left(\sum_{j=T-k}^T \boldsymbol{\theta}^\top \bar{\mathbf{x}}(j)\right) \end{bmatrix}$$

Graph convolution

$$\begin{bmatrix} \sigma\left(\sum_{j \in \mathcal{N}(1)} \boldsymbol{\theta}^\top \bar{\mathbf{x}}(j)\right) \\ \sigma\left(\sum_{j \in \mathcal{N}(2)} \boldsymbol{\theta}^\top \bar{\mathbf{x}}(j)\right) \\ \vdots \\ \sigma\left(\sum_{j \in \mathcal{N}(T)} \boldsymbol{\theta}^\top \bar{\mathbf{x}}(j)\right) \end{bmatrix}$$

Short Intro to Graph Neural Networks

Standard 1D convolution

$$\begin{bmatrix} \sigma\left(\sum_{j=1}^k \theta^\top \bar{x}(j)\right) \\ \sigma\left(\sum_{j=2}^{k+1} \theta^\top \bar{x}(j)\right) \\ \vdots \\ \sigma\left(\sum_{j=T-k}^T \theta^\top \bar{x}(j)\right) \end{bmatrix}$$

Graph convolution

$$\begin{bmatrix} \sigma\left(\sum_{j \in N(1)} \theta^\top \bar{x}(j)\right) \\ \sigma\left(\sum_{j \in N(2)} \theta^\top \bar{x}(j)\right) \\ \vdots \\ \sigma\left(\sum_{j \in N(T)} \theta^\top \bar{x}(j)\right) \end{bmatrix}$$

Question: What is the output size of graph convolution?

Short Intro to Graph Neural Networks

Standard 1D convolution

$$\begin{bmatrix} \sigma\left(\sum_{j=1}^k \theta^\top \bar{x}(j)\right) \\ \sigma\left(\sum_{j=2}^{k+1} \theta^\top \bar{x}(j)\right) \\ \vdots \\ \sigma\left(\sum_{j=T-k}^T \theta^\top \bar{x}(j)\right) \end{bmatrix}$$

Graph convolution

$$\begin{bmatrix} \sigma\left(\sum_{j \in N(1)} \theta^\top \bar{x}(j)\right) \\ \sigma\left(\sum_{j \in N(2)} \theta^\top \bar{x}(j)\right) \\ \vdots \\ \sigma\left(\sum_{j \in N(T)} \theta^\top \bar{x}(j)\right) \end{bmatrix}$$

Question: What is the output size of graph convolution?

Answer: $T \times d_h$

Short Intro to Graph Neural Networks

We can use pooling with graph convolutions too

$$\text{SumPool} \left(\begin{bmatrix} \sigma\left(\sum_{j \in \mathcal{N}(1)} \boldsymbol{\theta}^\top \bar{\mathbf{x}}(j)\right) \\ \sigma\left(\sum_{j \in \mathcal{N}(2)} \boldsymbol{\theta}^\top \bar{\mathbf{x}}(j)\right) \\ \vdots \\ \sigma\left(\sum_{j \in \mathcal{N}(T)} \boldsymbol{\theta}^\top \bar{\mathbf{x}}(j)\right) \end{bmatrix} \right) = \sigma\left(\sum_{j \in \mathcal{N}(1)} \boldsymbol{\theta}^\top \bar{\mathbf{x}}(j)\right) + \sigma\left(\sum_{j \in \mathcal{N}(2)} \boldsymbol{\theta}^\top \bar{\mathbf{x}}(j)\right) + \dots + \sigma\left(\sum_{j \in \mathcal{N}(T)} \boldsymbol{\theta}^\top \bar{\mathbf{x}}(j)\right)$$

Short Intro to Graph Neural Networks

We can write attention and transformers as a graph neural network

Short Intro to Graph Neural Networks

We can write attention and transformers as a graph neural network

- Make the neighborhood the entire graph $N(i) = X$

Short Intro to Graph Neural Networks

We can write attention and transformers as a graph neural network

- Make the neighborhood the entire graph $N(i) = X$
- Learnable continuous edge weights (key \times query)

Short Intro to Graph Neural Networks

We can write attention and transformers as a graph neural network

- Make the neighborhood the entire graph $N(i) = X$
- Learnable continuous edge weights (key \times query)

Edge weights

$$A(i) = \text{softmax} \left(\frac{(\theta_Q^\top x_i)(\theta_K^\top X)}{\sqrt{d_z}} \right)$$

Short Intro to Graph Neural Networks

We can write attention and transformers as a graph neural network

- Make the neighborhood the entire graph $N(i) = X$
- Learnable continuous edge weights (key \times query)

Edge weights

$$A(i) = \text{softmax} \left(\frac{(\theta_Q^\top x_i)(\theta_K^\top X)}{\sqrt{d_z}} \right)$$

Graph convolution

$$f(X, i) = \sum_{j \in N(i)} \theta_V^\top x_j \cdot A(i)_j$$

Short Intro to Graph Neural Networks

We can write attention and transformers as a graph neural network

- Make the neighborhood the entire graph $N(i) = X$
- Learnable continuous edge weights (key \times query)

Edge weights

$$A(i) = \text{softmax} \left(\frac{(\theta_Q^\top x_i)(\theta_K^\top X)}{\sqrt{d_z}} \right)$$

Graph convolution

$$f(X, i) = \sum_{j \in N(i)} \theta_V^\top x_j \cdot A(i)_j$$

Many ways to teach attention, but I like the RNN-style approach more

Attention

Attention

Attention and transformers are the “hottest” topic in deep learning

Attention

Attention and transformers are the “hottest” topic in deep learning

People use them for almost every task (even if they shouldn’t!)

Attention

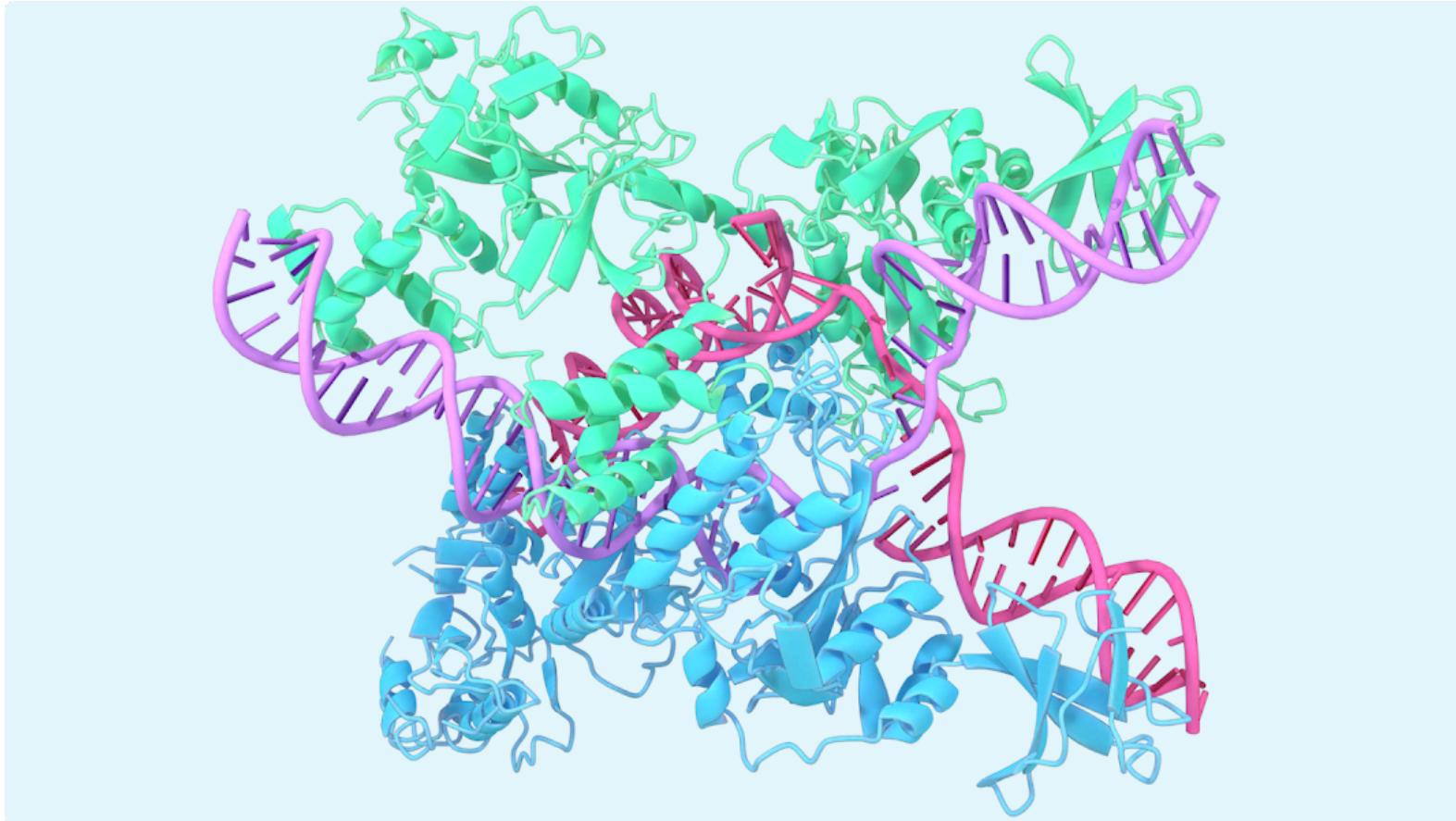
Attention and transformers are the “hottest” topic in deep learning

People use them for almost every task (even if they shouldn’t!)

Let’s review some projects based on attention

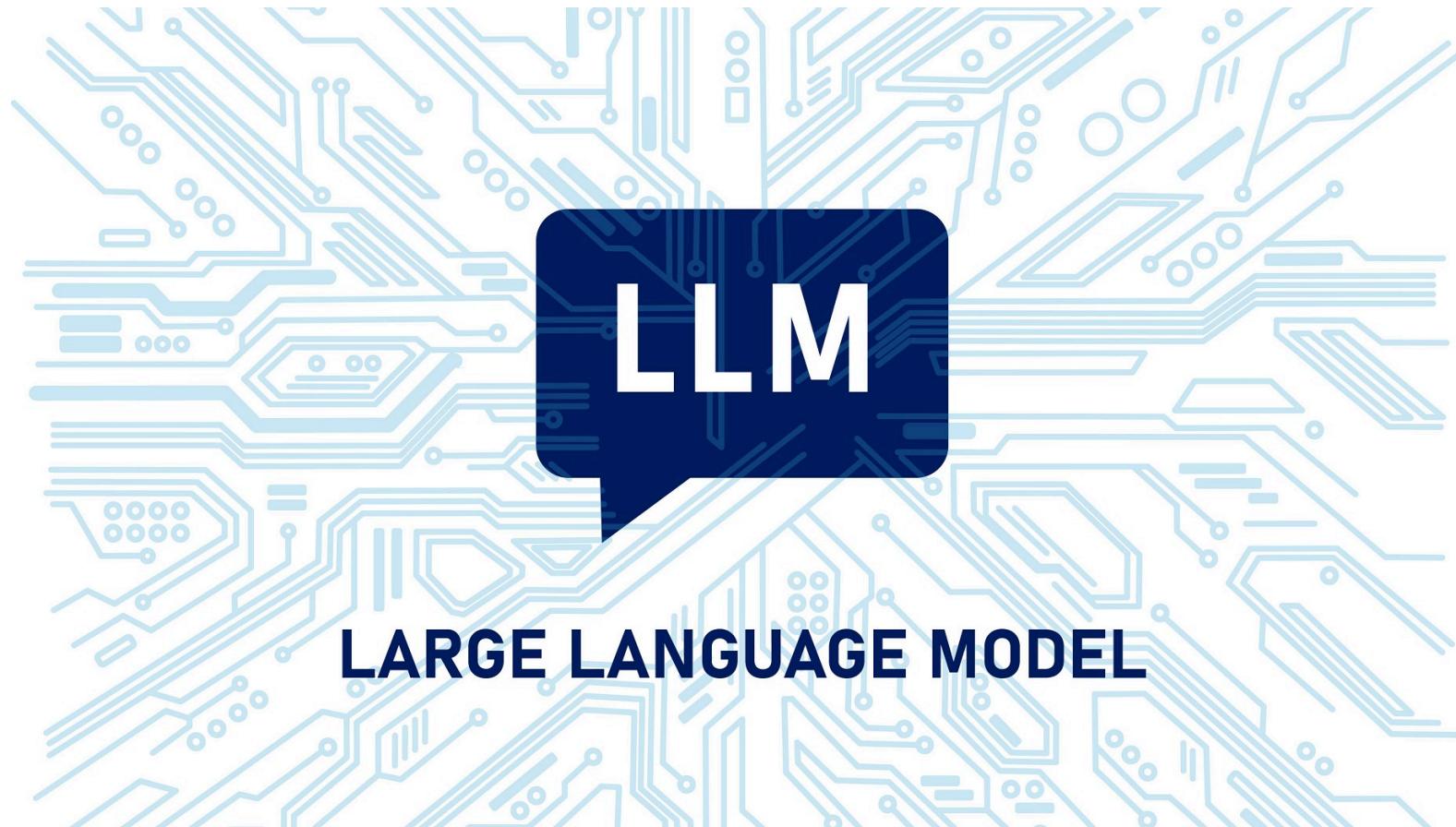
Attention

AlphaFold (Nobel prize)



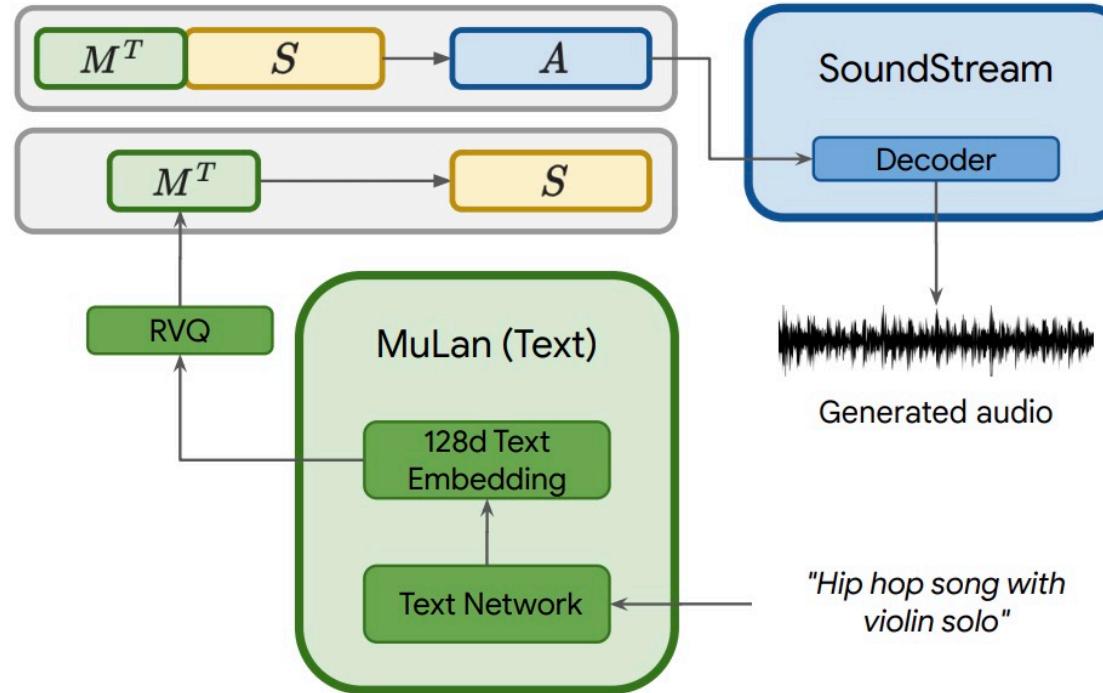
Attention

ChatGPT, Qwen, LLaMA, Mistral, Doubou, Ernie chatbots



Attention

MusicTransformer, MuLan



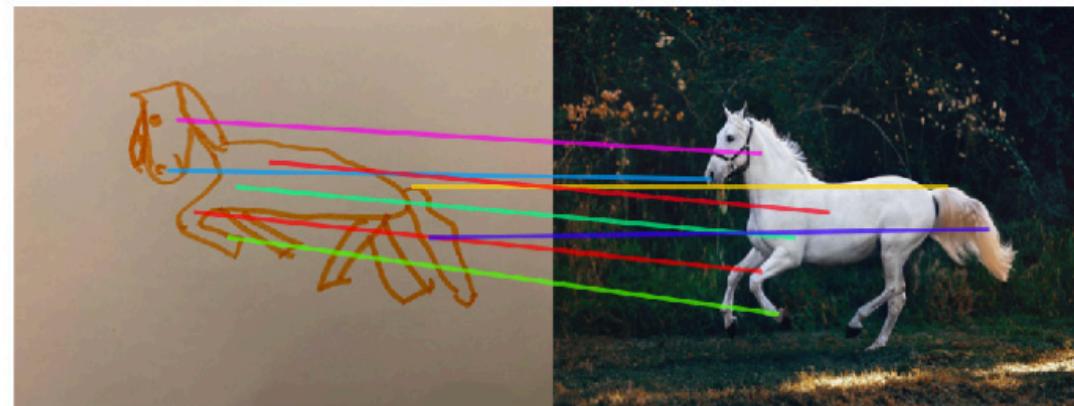
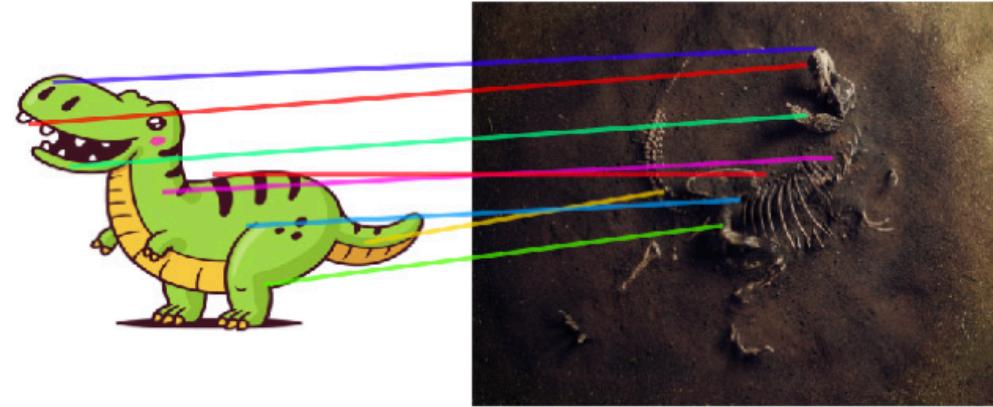
Attention

Google Translate, Baidu Translate, Apple Translate



Attention

ViT, DinoV2



Attention

All these models are **transformers**

Attention

All these models are **transformers**

At the core of each transformer is **attention**

Attention

All these models are **transformers**

At the core of each transformer is **attention**

We will write attention slightly differently than you see in textbooks

Attention

All these models are **transformers**

At the core of each transformer is **attention**

We will write attention slightly differently than you see in textbooks

- I want to stay with column-vector notation

Attention

All these models are **transformers**

At the core of each transformer is **attention**

We will write attention slightly differently than you see in textbooks

- I want to stay with column-vector notation

Textbook

$$\text{softmax}(\mathbf{Q}\mathbf{K}^\top)\mathbf{V}$$

Attention

All these models are **transformers**

At the core of each transformer is **attention**

We will write attention slightly differently than you see in textbooks

- I want to stay with column-vector notation

Textbook

$$\text{softmax}(\mathbf{Q}\mathbf{K}^\top)\mathbf{V}$$

Our representation

$$\mathbf{V} \text{ softmax}(\mathbf{K}^\top \mathbf{Q})$$

Attention

All these models are **transformers**

At the core of each transformer is **attention**

We will write attention slightly differently than you see in textbooks

- I want to stay with column-vector notation

Textbook

$$\text{softmax}(\mathbf{Q}\mathbf{K}^\top)\mathbf{V}$$

Our representation

$$\mathbf{V} \text{ softmax}(\mathbf{K}^\top \mathbf{Q})$$

Same results, but easier to write scalar/vector

Attention

We can derive attention from composite memory

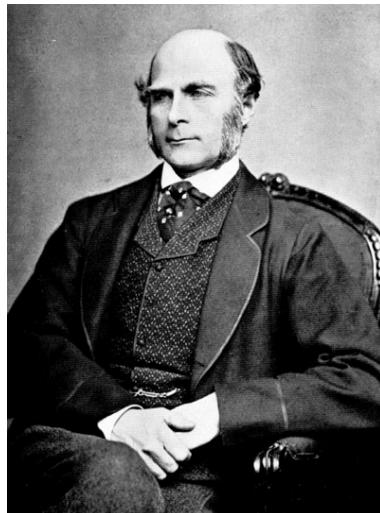
Attention

Attention

We can derive attention from composite memory

Francis Galton (1822-1911)
photo composite memory

Attention



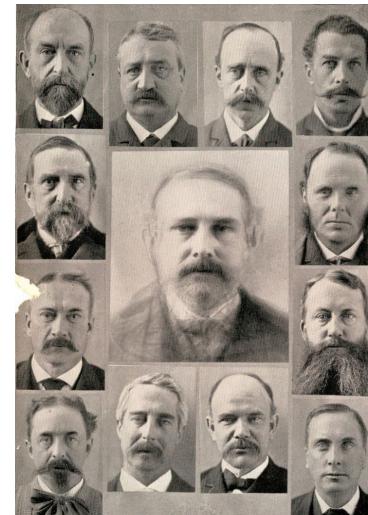
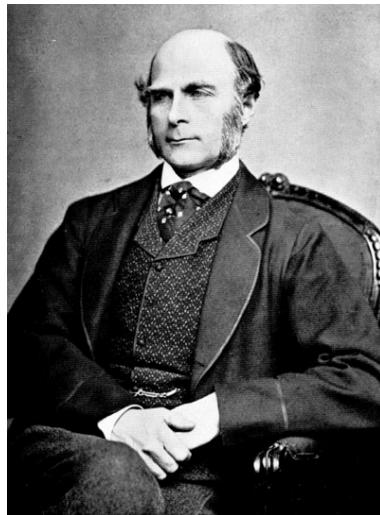
Attention

We can derive attention from composite memory

Francis Galton (1822-1911)
photo composite memory

Composite photo of members of a
party

Attention



Attention

Task: Find a mathematical model of how our mind represents memories

Attention

Task: Find a mathematical model of how our mind represents memories

$X : \mathbb{R}^{h \times w}$ People you see at the party

Attention

Task: Find a mathematical model of how our mind represents memories

$X : \mathbb{R}^{h \times w}$ People you see at the party

$H : \mathbb{R}^{h \times w}$ The image in your mind

Attention

Task: Find a mathematical model of how our mind represents memories

$X : \mathbb{R}^{h \times w}$ People you see at the party

$H : \mathbb{R}^{h \times w}$ The image in your mind

$f : X^T \times \Theta \mapsto H$

Attention

Task: Find a mathematical model of how our mind represents memories

$X : \mathbb{R}^{h \times w}$ People you see at the party

$H : \mathbb{R}^{h \times w}$ The image in your mind

$f : X^T \times \Theta \mapsto H$

Composite photography/memory uses a weighted sum

Attention

Task: Find a mathematical model of how our mind represents memories

$$X : \mathbb{R}^{h \times w} \quad \text{People you see at the party}$$

$$H : \mathbb{R}^{h \times w} \quad \text{The image in your mind}$$

$$f : X^T \times \Theta \mapsto H$$

Composite photography/memory uses a weighted sum

$$f(x, \theta) = \sum_{i=1}^T \theta^\top \bar{x}_i$$

Attention

Limited space, cannot remember everything

Attention

Limited space, cannot remember everything

Introduced forgetting term $\gamma \in [0, 1]$

Attention

Limited space, cannot remember everything

Introduced forgetting term $\gamma \in [0, 1]$

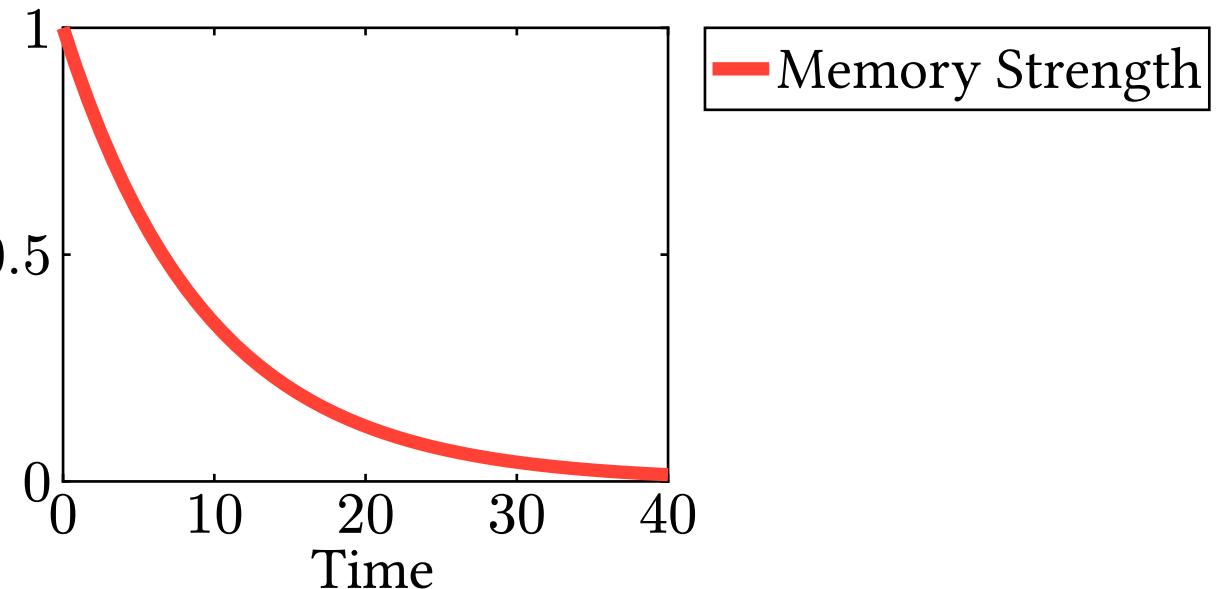
$$f(\mathbf{x}, \boldsymbol{\theta}) = \sum_{i=1}^T \gamma^{T-i} \cdot \boldsymbol{\theta}^\top \overline{\mathbf{x}}_i$$

Attention

Limited space, cannot remember everything

Introduced forgetting term $\gamma \in [0, 1]$

$$f(x, \theta) = \sum_{i=1}^T \gamma^{T-i} \cdot \theta^\top \bar{x}_i$$

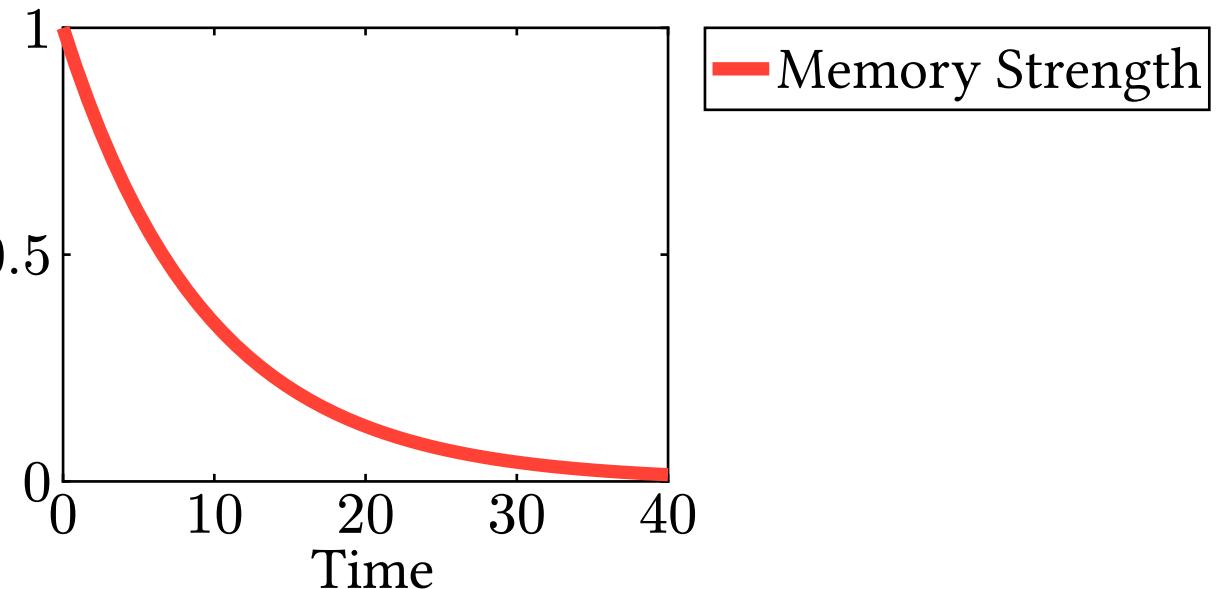


Attention

Limited space, cannot remember everything

Introduced forgetting term $\gamma \in [0, 1]$

$$f(x, \theta) = \sum_{i=1}^T \gamma^{T-i} \cdot \theta^\top \bar{x}_i$$



Question: Does this accurately model what **you** remember?

Attention

Example: We attend a party in 1850s

Attention

Example: We attend a party in 1850s

We talk with many people at this party

Attention

Example: We attend a party in 1850s

We talk with many people at this party



Attention

Example: We attend a party in 1850s

We talk with many people at this party



10 PM

Attention

Example: We attend a party in 1850s

We talk with many people at this party



10 PM

11 PM

Attention

Example: We attend a party in 1850s

We talk with many people at this party



10 PM

11 PM

12 AM

Attention

Example: We attend a party in 1850s

We talk with many people at this party



10 PM

11 PM

12 AM

1 AM

Attention

According to forgetting, the memories should fade with time

Attention

According to forgetting, the memories should fade with time



Attention

According to forgetting, the memories should fade with time



$$\gamma^3 \boldsymbol{\theta}^\top \overline{\boldsymbol{x}}_1$$

Attention

According to forgetting, the memories should fade with time



$$\gamma^3 \boldsymbol{\theta}^\top \bar{\boldsymbol{x}}_1$$

$$\gamma^2 \boldsymbol{\theta}^\top \bar{\boldsymbol{x}}_2$$

Attention

According to forgetting, the memories should fade with time



$$\gamma^3 \boldsymbol{\theta}^\top \overline{\boldsymbol{x}}_1$$

$$\gamma^2 \boldsymbol{\theta}^\top \overline{\boldsymbol{x}}_2$$

$$\gamma^1 \boldsymbol{\theta}^\top \overline{\boldsymbol{x}}_3$$

Attention

According to forgetting, the memories should fade with time



$$\gamma^3 \boldsymbol{\theta}^\top \overline{\boldsymbol{x}}_1$$

$$\gamma^2 \boldsymbol{\theta}^\top \overline{\boldsymbol{x}}_2$$

$$\gamma^1 \boldsymbol{\theta}^\top \overline{\boldsymbol{x}}_3$$

$$\gamma^0 \boldsymbol{\theta}^\top \overline{\boldsymbol{x}}_4$$

Attention

According to forgetting, the memories should fade with time



$$\gamma^3 \boldsymbol{\theta}^\top \overline{\boldsymbol{x}}_1$$

$$\gamma^2 \boldsymbol{\theta}^\top \overline{\boldsymbol{x}}_2$$

$$\gamma^1 \boldsymbol{\theta}^\top \overline{\boldsymbol{x}}_3$$

$$\gamma^0 \boldsymbol{\theta}^\top \overline{\boldsymbol{x}}_4$$

Attention

Any questions before moving on?

Attention

Consider another party, with one more guest

Attention

Consider another party, with one more guest



Attention

Consider another party, with one more guest



$$\gamma^4 \boldsymbol{\theta}^\top \overline{\boldsymbol{x}}_1$$

$$\gamma^3 \boldsymbol{\theta}^\top \overline{\boldsymbol{x}}_2$$

$$\gamma^2 \boldsymbol{\theta}^\top \overline{\boldsymbol{x}}_3$$

$$\gamma^1 \boldsymbol{\theta}^\top \overline{\boldsymbol{x}}_4$$

$$\gamma^0 \boldsymbol{\theta}^\top \overline{\boldsymbol{x}}_5$$

Attention

Consider another party, with one more guest



$$\gamma^4 \theta^\top \bar{x}_1$$

$$\gamma^3 \theta^\top \bar{x}_2$$

$$\gamma^2 \theta^\top \bar{x}_3$$

$$\gamma^1 \theta^\top \bar{x}_4$$

$$\gamma^0 \theta^\top \bar{x}_5$$

Question: What will happen to Taylor Swift?

Attention



$$\gamma^4 \boldsymbol{\theta}^\top \bar{x}_1$$

$$\gamma^3 \boldsymbol{\theta}^\top \bar{x}_2$$

$$\gamma^2 \boldsymbol{\theta}^\top \bar{x}_3$$

$$\gamma^1 \boldsymbol{\theta}^\top \bar{x}_4$$

$$\gamma^0 \boldsymbol{\theta}^\top \bar{x}_5$$

Attention



$$\gamma^4 \boldsymbol{\theta}^\top \bar{\boldsymbol{x}}_1$$

$$\gamma^3 \boldsymbol{\theta}^\top \bar{\boldsymbol{x}}_2$$

$$\gamma^2 \boldsymbol{\theta}^\top \bar{\boldsymbol{x}}_3$$

$$\gamma^1 \boldsymbol{\theta}^\top \bar{\boldsymbol{x}}_4$$

$$\gamma^0 \boldsymbol{\theta}^\top \bar{\boldsymbol{x}}_5$$

We will forget meeting her!

Attention



$$\gamma^4 \boldsymbol{\theta}^\top \bar{x}_1$$

$$\gamma^3 \boldsymbol{\theta}^\top \bar{x}_2$$

$$\gamma^2 \boldsymbol{\theta}^\top \bar{x}_3$$

$$\gamma^1 \boldsymbol{\theta}^\top \bar{x}_4$$

$$\gamma^0 \boldsymbol{\theta}^\top \bar{x}_5$$

We will forget meeting her!

Question: Would you forget meeting Taylor Swift?

Attention

Our model of memory is incomplete

Attention

Our model of memory is incomplete

Memories are not created equal, some are more important than others

Attention

Our model of memory is incomplete

Memories are not created equal, some are more important than others

Important memories persist longer than unimportant memories

Attention

Our model of memory is incomplete

Memories are not created equal, some are more important than others

Important memories persist longer than unimportant memories

We will **pay more attention** to certain memories

Attention

What does human memory actually look like?



Attention

What does human memory actually look like?



$$1.0 \cdot \theta^\top \bar{x}_1$$

Attention

What does human memory actually look like?



$$1.0 \cdot \theta^\top \bar{x}_1 \quad 0.1 \cdot \theta^\top \bar{x}_2$$

Attention

What does human memory actually look like?



$$1.0 \cdot \theta^\top \bar{x}_1 \quad 0.1 \cdot \theta^\top \bar{x}_2 \quad 0.1 \cdot \theta^\top \bar{x}_3$$

Attention

What does human memory actually look like?



$$1.0 \cdot \theta^\top \bar{x}_1$$

$$0.1 \cdot \theta^\top \bar{x}_2$$

$$0.1 \cdot \theta^\top \bar{x}_3$$

$$0.5 \cdot \theta^\top \bar{x}_4$$

Attention

What does human memory actually look like?



$$1.0 \cdot \theta^\top \bar{x}_1$$

$$0.1 \cdot \theta^\top \bar{x}_2$$

$$0.1 \cdot \theta^\top \bar{x}_3$$

$$0.5 \cdot \theta^\top \bar{x}_4$$

$$0.1 \cdot \theta^\top \bar{x}_5$$

Attention

What does human memory actually look like?



$$1.0 \cdot \theta^\top \bar{x}_1$$

$$0.1 \cdot \theta^\top \bar{x}_2$$

$$0.1 \cdot \theta^\top \bar{x}_3$$

$$0.5 \cdot \theta^\top \bar{x}_4$$

$$0.1 \cdot \theta^\top \bar{x}_5$$

Question: How can we achieve this forgetting?

Attention

In our composite model, forgetting is a function of time

Attention

In our composite model, forgetting is a function of time

Question: Any forgetting mechanism that is not a function of time?

Attention

In our composite model, forgetting is a function of time

Question: Any forgetting mechanism that is not a function of time?

Answer: Forgetting in recurrent neural network is function of input!

Attention

In our composite model, forgetting is a function of time

Question: Any forgetting mechanism that is not a function of time?

Answer: Forgetting in recurrent neural network is function of input!

$$f_{\text{forget}}(x, \theta) = \sigma(\theta_\lambda^\top \bar{x})$$

Attention

In our composite model, forgetting is a function of time

Question: Any forgetting mechanism that is not a function of time?

Answer: Forgetting in recurrent neural network is function of input!

$$f_{\text{forget}}(x, \theta) = \sigma(\theta_\lambda^\top \bar{x})$$

$$f(h, x, \theta) = f_{\text{forget}}(x, \theta) \odot h + \theta_x^\top \bar{x}$$

Attention

First, write our forgetting function with slightly different notation

Attention

First, write our forgetting function with slightly different notation

$$\lambda(\boldsymbol{x}, \boldsymbol{\theta}_\lambda) = \sigma(\boldsymbol{\theta}_\lambda^\top \overline{\boldsymbol{x}}); \quad \boldsymbol{\theta}_\lambda \in \mathbb{R}^{(d_x+1) \times 1}$$

Attention

First, write our forgetting function with slightly different notation

$$\lambda(x, \theta_\lambda) = \sigma(\theta_\lambda^\top \bar{x}); \quad \theta_\lambda \in \mathbb{R}^{(d_x+1) \times 1}$$

Question: Shape of $\lambda(x, \theta_\lambda)$?

Attention

First, write our forgetting function with slightly different notation

$$\lambda(x, \theta_\lambda) = \sigma(\theta_\lambda^\top \bar{x}); \quad \theta_\lambda \in \mathbb{R}^{(d_x+1) \times 1}$$

Question: Shape of $\lambda(x, \theta_\lambda)$?

Answer: Scalar!

Attention

First, write our forgetting function with slightly different notation

$$\lambda(x, \theta_\lambda) = \sigma(\theta_\lambda^\top \bar{x}); \quad \theta_\lambda \in \mathbb{R}^{(d_x+1) \times 1}$$

Question: Shape of $\lambda(x, \theta_\lambda)$?

Answer: Scalar!

Then, write our composite memory model with forgetting

Attention

First, write our forgetting function with slightly different notation

$$\lambda(\mathbf{x}, \boldsymbol{\theta}_\lambda) = \sigma(\boldsymbol{\theta}_\lambda^\top \overline{\mathbf{x}}); \quad \boldsymbol{\theta}_\lambda \in \mathbb{R}^{(d_x+1) \times 1}$$

Question: Shape of $\lambda(\mathbf{x}, \boldsymbol{\theta}_\lambda)$?

Answer: Scalar!

Then, write our composite memory model with forgetting

$$f\left(\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \boldsymbol{\theta}\right) = \sum_{i=1}^T \boldsymbol{\theta}^\top \mathbf{x}_i \cdot \lambda(\mathbf{x}_i, \boldsymbol{\theta}_\lambda)$$

Attention

First, write our forgetting function with slightly different notation

$$\lambda(\mathbf{x}, \boldsymbol{\theta}_\lambda) = \sigma(\boldsymbol{\theta}_\lambda^\top \overline{\mathbf{x}}); \quad \boldsymbol{\theta}_\lambda \in \mathbb{R}^{(d_x+1) \times 1}$$

Question: Shape of $\lambda(\mathbf{x}, \boldsymbol{\theta}_\lambda)$?

Answer: Scalar!

Then, write our composite memory model with forgetting

$$f\left(\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \boldsymbol{\theta}\right) = \sum_{i=1}^T \boldsymbol{\theta}^\top \mathbf{x}_i \cdot \lambda(\mathbf{x}_i, \boldsymbol{\theta}_\lambda)$$

Only pay attention to important inputs

Attention

We can use this simple form of attention to remember Taylor Swift

Attention

We can use this simple form of attention to remember Taylor Swift



Attention

$$\begin{array}{ccccc} \lambda(x_1, \theta_\lambda) & \lambda(x_2, \theta_\lambda) & \lambda(x_3, \theta_\lambda) & \lambda(x_4, \theta_\lambda) & \lambda(x_5, \theta_\lambda) \\ \cdot \theta^\top \bar{x}_1 & \cdot \theta^\top \bar{x}_2 & \cdot \theta^\top \bar{x}_3 & \cdot \theta^\top \bar{x}_4 & \cdot \theta^\top \bar{x}_5 \end{array}$$

Attention

$$\begin{array}{ccccc} \lambda(x_1, \theta_\lambda) & \lambda(x_2, \theta_\lambda) & \lambda(x_3, \theta_\lambda) & \lambda(x_4, \theta_\lambda) & \lambda(x_5, \theta_\lambda) \\ \cdot \theta^\top \bar{x}_1 & \cdot \theta^\top \bar{x}_2 & \cdot \theta^\top \bar{x}_3 & \cdot \theta^\top \bar{x}_4 & \cdot \theta^\top \bar{x}_5 \end{array}$$

Question: What do the images look like now?

Attention

$$\begin{array}{ccccc} \lambda(x_1, \theta_\lambda) & \lambda(x_2, \theta_\lambda) & \lambda(x_3, \theta_\lambda) & \lambda(x_4, \theta_\lambda) & \lambda(x_5, \theta_\lambda) \\ \cdot \theta^\top \bar{x}_1 & \cdot \theta^\top \bar{x}_2 & \cdot \theta^\top \bar{x}_3 & \cdot \theta^\top \bar{x}_4 & \cdot \theta^\top \bar{x}_5 \end{array}$$

Question: What do the images look like now?



Attention

This form of attention will learn to pay attention to everyone!

Attention

This form of attention will learn to pay attention to everyone!

$$1.0 \cdot \theta^\top \bar{x}_1$$

$$1.0 \cdot \theta^\top \bar{x}_2$$

$$1.0 \cdot \theta^\top \bar{x}_3$$

$$1.0 \cdot \theta^\top \bar{x}_4$$

$$1.0 \cdot \theta^\top \bar{x}_5$$

Attention

This form of attention will learn to pay attention to everyone!

$$1.0 \cdot \theta^\top \bar{x}_1$$

$$1.0 \cdot \theta^\top \bar{x}_2$$

$$1.0 \cdot \theta^\top \bar{x}_3$$

$$1.0 \cdot \theta^\top \bar{x}_4$$

$$1.0 \cdot \theta^\top \bar{x}_5$$



Attention

This form of attention will learn to pay attention to everyone!

$$1.0 \cdot \theta^\top \bar{x}_1$$

$$1.0 \cdot \theta^\top \bar{x}_2$$

$$1.0 \cdot \theta^\top \bar{x}_3$$

$$1.0 \cdot \theta^\top \bar{x}_4$$

$$1.0 \cdot \theta^\top \bar{x}_5$$



Attention

$$1.0 \cdot \theta^\top \bar{x}_1 \quad 1.0 \cdot \theta^\top \bar{x}_2 \quad 1.0 \cdot \theta^\top \bar{x}_3 \quad 1.0 \cdot \theta^\top \bar{x}_4 \quad 1.0 \cdot \theta^\top \bar{x}_5$$

Attention

$$1.0 \cdot \theta^\top \bar{x}_1 \quad 1.0 \cdot \theta^\top \bar{x}_2 \quad 1.0 \cdot \theta^\top \bar{x}_3 \quad 1.0 \cdot \theta^\top \bar{x}_4 \quad 1.0 \cdot \theta^\top \bar{x}_5$$

Why does this model of attention pay attention to everyone?

Attention

$$1.0 \cdot \theta^\top \bar{x}_1 \quad 1.0 \cdot \theta^\top \bar{x}_2 \quad 1.0 \cdot \theta^\top \bar{x}_3 \quad 1.0 \cdot \theta^\top \bar{x}_4 \quad 1.0 \cdot \theta^\top \bar{x}_5$$

Why does this model of attention pay attention to everyone?

- Honest answer is I do not know

Attention

$$1.0 \cdot \theta^\top \bar{x}_1 \quad 1.0 \cdot \theta^\top \bar{x}_2 \quad 1.0 \cdot \theta^\top \bar{x}_3 \quad 1.0 \cdot \theta^\top \bar{x}_4 \quad 1.0 \cdot \theta^\top \bar{x}_5$$

Why does this model of attention pay attention to everyone?

- Honest answer is I do not know
 - ▶ But I have tried it myself, and this usually happens

Attention

$$1.0 \cdot \theta^\top \bar{x}_1 \quad 1.0 \cdot \theta^\top \bar{x}_2 \quad 1.0 \cdot \theta^\top \bar{x}_3 \quad 1.0 \cdot \theta^\top \bar{x}_4 \quad 1.0 \cdot \theta^\top \bar{x}_5$$

Why does this model of attention pay attention to everyone?

- Honest answer is I do not know
 - But I have tried it myself, and this usually happens
 - Maybe one of you can figure out why!

Attention

$$1.0 \cdot \theta^\top \bar{x}_1 \quad 1.0 \cdot \theta^\top \bar{x}_2 \quad 1.0 \cdot \theta^\top \bar{x}_3 \quad 1.0 \cdot \theta^\top \bar{x}_4 \quad 1.0 \cdot \theta^\top \bar{x}_5$$

Why does this model of attention pay attention to everyone?

- Honest answer is I do not know
 - But I have tried it myself, and this usually happens
 - Maybe one of you can figure out why!

We need to prevent all attention weights going to 1

Attention

We should normalize $\lambda(x, \theta_\lambda)$ to model finite (human) attention span

Attention

We should normalize $\lambda(x, \theta_\lambda)$ to model finite (human) attention span

For example, normalize attention to sum to one

$$\sum_{i=1}^T \lambda(x_i, \theta_\lambda) = 1$$

Attention

We should normalize $\lambda(x, \theta_\lambda)$ to model finite (human) attention span

For example, normalize attention to sum to one

$$\sum_{i=1}^T \lambda(x_i, \theta_\lambda) = 1$$

Now the model must choose who to remember!

Attention

We should normalize $\lambda(x, \theta_\lambda)$ to model finite (human) attention span

For example, normalize attention to sum to one

$$\sum_{i=1}^T \lambda(x_i, \theta_\lambda) = 1$$

Now the model must choose who to remember!

Question: How can we ensure that the attention sums to one?

Attention

We should normalize $\lambda(x, \theta_\lambda)$ to model finite (human) attention span

For example, normalize attention to sum to one

$$\sum_{i=1}^T \lambda(x_i, \theta_\lambda) = 1$$

Now the model must choose who to remember!

Question: How can we ensure that the attention sums to one?

Answer: Softmax!

Attention

The softmax function maps real numbers to the simplex (probabilities)

Attention

The softmax function maps real numbers to the simplex (probabilities)

$$\text{softmax} : \mathbb{R}^k \mapsto \Delta^{k-1}$$

Attention

The softmax function maps real numbers to the simplex (probabilities)

$$\text{softmax} : \mathbb{R}^k \mapsto \Delta^{k-1}$$

$$\text{softmax} \left(\begin{bmatrix} x_1 \\ \vdots \\ x_k \end{bmatrix} \right) = \frac{\exp(\mathbf{x})}{\sum_{i=1}^k \exp(x_i)} = \begin{bmatrix} \frac{\exp(x_1)}{\exp(x_1) + \exp(x_2) + \dots + \exp(x_k)} \\ \frac{\exp(x_2)}{\exp(x_1) + \exp(x_2) + \dots + \exp(x_k)} \\ \vdots \\ \frac{\exp(x_k)}{\exp(x_1) + \exp(x_2) + \dots + \exp(x_k)} \end{bmatrix}$$

Attention

Let us rewrite attention using softmax

Attention

Let us rewrite attention using softmax

The attention we pay to person i is

$$\lambda \left(\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \boldsymbol{\theta}_\lambda \right)_i = \text{softmax} \left(\begin{bmatrix} \boldsymbol{\theta}_\lambda^\top \bar{\mathbf{x}}_1 \\ \vdots \\ \boldsymbol{\theta}_\lambda^\top \bar{\mathbf{x}}_T \end{bmatrix} \right)_i = \frac{\exp(\boldsymbol{\theta}_\lambda^\top \bar{\mathbf{x}}_i)}{\sum_{j=1}^T \exp(\boldsymbol{\theta}_\lambda^\top \bar{\mathbf{x}}_j)}$$

Attention



Attention



$$\lambda \left(\begin{bmatrix} x_1 \\ \vdots \\ x_5 \end{bmatrix}, \theta_\lambda \right)_1 \cdot \theta^\top \bar{x}_1$$

Attention



$$\lambda \left(\begin{bmatrix} x_1 \\ \vdots \\ x_5 \end{bmatrix}, \theta_\lambda \right)_1 \lambda \left(\begin{bmatrix} x_1 \\ \vdots \\ x_5 \end{bmatrix}, \theta_\lambda \right)_2 \\ \cdot \theta^\top \bar{x}_1 \quad \cdot \theta^\top \bar{x}_2$$

Attention



$$\lambda \left(\begin{bmatrix} x_1 \\ \vdots \\ x_5 \end{bmatrix}, \theta_\lambda \right)_1 \lambda \left(\begin{bmatrix} x_1 \\ \vdots \\ x_5 \end{bmatrix}, \theta_\lambda \right)_2 \lambda \left(\begin{bmatrix} x_1 \\ \vdots \\ x_5 \end{bmatrix}, \theta_\lambda \right)_3 \\ \cdot \theta^\top \bar{x}_1 \quad \cdot \theta^\top \bar{x}_2 \quad \cdot \theta^\top \bar{x}_3$$

Attention



$$\lambda \left(\begin{bmatrix} x_1 \\ \vdots \\ x_5 \end{bmatrix}, \theta_\lambda \right)_1 \lambda \left(\begin{bmatrix} x_1 \\ \vdots \\ x_5 \end{bmatrix}, \theta_\lambda \right)_2 \lambda \left(\begin{bmatrix} x_1 \\ \vdots \\ x_5 \end{bmatrix}, \theta_\lambda \right)_3 \lambda \left(\begin{bmatrix} x_1 \\ \vdots \\ x_5 \end{bmatrix}, \theta_\lambda \right)_4 \\ \cdot \theta^\top \bar{x}_1 \quad \cdot \theta^\top \bar{x}_2 \quad \cdot \theta^\top \bar{x}_3 \quad \cdot \theta^\top \bar{x}_4$$

Attention



$$\lambda \left(\begin{bmatrix} x_1 \\ \vdots \\ x_5 \end{bmatrix}, \theta_\lambda \right)_1 \lambda \left(\begin{bmatrix} x_1 \\ \vdots \\ x_5 \end{bmatrix}, \theta_\lambda \right)_2 \lambda \left(\begin{bmatrix} x_1 \\ \vdots \\ x_5 \end{bmatrix}, \theta_\lambda \right)_3 \lambda \left(\begin{bmatrix} x_1 \\ \vdots \\ x_5 \end{bmatrix}, \theta_\lambda \right)_4 \lambda \left(\begin{bmatrix} x_1 \\ \vdots \\ x_5 \end{bmatrix}, \theta_\lambda \right)_5 \\ \cdot \theta^\top \bar{x}_1 \quad \cdot \theta^\top \bar{x}_2 \quad \cdot \theta^\top \bar{x}_3 \quad \cdot \theta^\top \bar{x}_4 \quad \cdot \theta^\top \bar{x}_5$$

Attention



$$\lambda \left(\begin{bmatrix} x_1 \\ \vdots \\ x_5 \end{bmatrix}, \theta_\lambda \right)_1 \lambda \left(\begin{bmatrix} x_1 \\ \vdots \\ x_5 \end{bmatrix}, \theta_\lambda \right)_2 \lambda \left(\begin{bmatrix} x_1 \\ \vdots \\ x_5 \end{bmatrix}, \theta_\lambda \right)_3 \lambda \left(\begin{bmatrix} x_1 \\ \vdots \\ x_5 \end{bmatrix}, \theta_\lambda \right)_4 \lambda \left(\begin{bmatrix} x_1 \\ \vdots \\ x_5 \end{bmatrix}, \theta_\lambda \right)_5 \\ \cdot \theta^\top \bar{x}_1 \quad \cdot \theta^\top \bar{x}_2 \quad \cdot \theta^\top \bar{x}_3 \quad \cdot \theta^\top \bar{x}_4 \quad \cdot \theta^\top \bar{x}_5$$

Attention



Attention



$$0.70 \cdot \theta^\top \bar{x}_1$$

Attention



$$0.70 \cdot \theta^\top \bar{x}_1 - 0.04 \cdot \theta^\top \bar{x}_2$$

Attention



$$0.70 \cdot \theta^\top \bar{x}_1 \quad 0.04 \cdot \theta^\top \bar{x}_2 \quad 0.03 \cdot \theta^\top \bar{x}_3$$

Attention



$$0.70 \cdot \theta^\top \bar{x}_1 \quad 0.04 \cdot \theta^\top \bar{x}_2 \quad 0.03 \cdot \theta^\top \bar{x}_3 \quad 0.20 \cdot \theta^\top \bar{x}_4$$

Attention



$$0.70 \cdot \theta^\top \bar{x}_1$$

$$0.04 \cdot \theta^\top \bar{x}_2$$

$$0.03 \cdot \theta^\top \bar{x}_3$$

$$0.20 \cdot \theta^\top \bar{x}_4$$

$$0.03 \cdot \theta^\top \bar{x}_5$$

Attention



$$0.70 \cdot \theta^\top \bar{x}_1 \quad 0.04 \cdot \theta^\top \bar{x}_2 \quad 0.03 \cdot \theta^\top \bar{x}_3 \quad 0.20 \cdot \theta^\top \bar{x}_4 \quad 0.03 \cdot \theta^\top \bar{x}_5$$

$$0.70 + 0.04 + 0.03 + 0.20 + 0.03 = 1.0$$

Attention

$$\lambda \left(\begin{bmatrix} \boldsymbol{x}_1 \\ \vdots \\ \boldsymbol{x}_T \end{bmatrix}, \boldsymbol{\theta}_{\lambda} \right)_i = \frac{\exp(\boldsymbol{\theta}_{\lambda}^{\top} \overline{\boldsymbol{x}}_i)}{\sum_{j=1}^T \exp(\boldsymbol{\theta}_{\lambda}^{\top} \overline{\boldsymbol{x}}_j)}$$

Attention

$$\lambda \left(\begin{bmatrix} \boldsymbol{x}_1 \\ \vdots \\ \boldsymbol{x}_T \end{bmatrix}, \boldsymbol{\theta}_{\lambda} \right)_i = \frac{\exp(\boldsymbol{\theta}_{\lambda}^{\top} \overline{\boldsymbol{x}}_i)}{\sum_{j=1}^T \exp(\boldsymbol{\theta}_{\lambda}^{\top} \overline{\boldsymbol{x}}_j)}$$

Compute attention for all inputs at once

Attention

$$\lambda \left(\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \boldsymbol{\theta}_{\lambda} \right)_i = \frac{\exp(\boldsymbol{\theta}_{\lambda}^{\top} \bar{\mathbf{x}}_i)}{\sum_{j=1}^T \exp(\boldsymbol{\theta}_{\lambda}^{\top} \bar{\mathbf{x}}_j)}$$

Compute attention for all inputs at once

$$\lambda \left(\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \boldsymbol{\theta}_{\lambda} \right) = \begin{bmatrix} \frac{\exp(\boldsymbol{\theta}_{\lambda}^{\top} \bar{\mathbf{x}}_1)}{\sum_{j=1}^T \exp(\boldsymbol{\theta}_{\lambda}^{\top} \bar{\mathbf{x}}_j)} \\ \vdots \\ \frac{\exp(\boldsymbol{\theta}_{\lambda}^{\top} \bar{\mathbf{x}}_T)}{\sum_{j=1}^T \exp(\boldsymbol{\theta}_{\lambda}^{\top} \bar{\mathbf{x}}_j)} \end{bmatrix}$$

Attention

This is a simple form of attention

Attention

This is a simple form of attention

Next, we will investigate the attention used in transformers

Keys and Queries

Keys and Queries

The modern form of attention behaves like a database

Keys and Queries

The modern form of attention behaves like a database

We label each person at the party with a **key**

Keys and Queries

The modern form of attention behaves like a database

We label each person at the party with a **key**

The key describes the content of each x

Keys and Queries

The modern form of attention behaves like a database

We label each person at the party with a **key**

The key describes the content of each x



Keys and Queries

The modern form of attention behaves like a database

We label each person at the party with a **key**

The key describes the content of each x



Musician

Keys and Queries

The modern form of attention behaves like a database

We label each person at the party with a **key**

The key describes the content of each x



Musician

Lawyer

Keys and Queries

The modern form of attention behaves like a database

We label each person at the party with a **key**

The key describes the content of each x



Musician

Lawyer

Shopkeeper

Keys and Queries

The modern form of attention behaves like a database

We label each person at the party with a **key**

The key describes the content of each x



Musician

Lawyer

Shopkeeper

Chef

Keys and Queries

The modern form of attention behaves like a database

We label each person at the party with a **key**

The key describes the content of each x



Musician

Lawyer

Shopkeeper

Chef

Scientist

Keys and Queries

We can search through our keys using a **query**

Keys and Queries

We can search through our keys using a **query**

Query: Which person will help me on my exam?

Keys and Queries

We can search through our keys using a **query**

Query: Which person will help me on my exam?



Keys and Queries

We can search through our keys using a **query**

Query: Which person will help me on my exam?



Musician

Lawyer

Shopkeeper

Chef

Scientist

Keys and Queries

We can search through our keys using a **query**

Query: Which person will help me on my exam?

Musician

Lawyer

Shopkeeper

Chef

Scientist



Keys and Queries

Query: I want to have fun

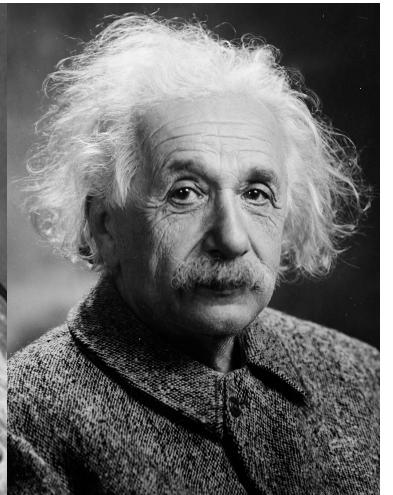
Keys and Queries

Query: I want to have fun



Keys and Queries

Query: I want to have fun



Musician

Lawyer

Shopkeeper

Chef

Scientist

Keys and Queries

Query: I want to have fun

Musician



Lawyer



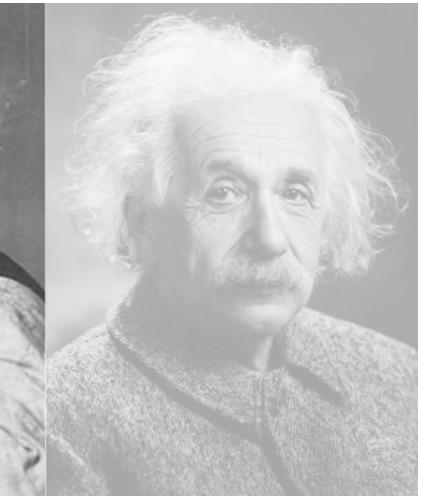
Shopkeeper



Chef



Scientist



Keys and Queries

Query: I want to have fun

Musician



Lawyer



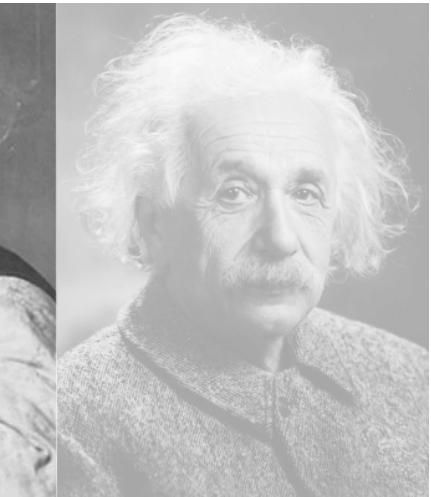
Shopkeeper



Chef



Scientist



How do we represent keys and queries mathematically?

Keys and Queries

For each input, we create a key k

Keys and Queries

For each input, we create a key \mathbf{k}

$$\mathbf{k}_j = \theta_K^\top \mathbf{x}_j, \quad \theta_K \in \mathbb{R}^{d_x \times d_h}, \quad \mathbf{k}_j \in \mathbb{R}^{d_h}$$

Keys and Queries

For each input, we create a key \mathbf{k}

$$\mathbf{k}_j = \theta_K^\top \mathbf{x}_j, \quad \theta_K \in \mathbb{R}^{d_x \times d_h}, \quad \mathbf{k}_j \in \mathbb{R}^{d_h}$$

Keys and Queries

Now, create a query from some x_q

Keys and Queries

Now, create a query from some x_q

$$q = \theta_Q^\top x_q, \quad \theta_Q \in \mathbb{R}^{d_x \times d_h}, \quad q \in \mathbb{R}^{d_h}$$

Keys and Queries

Now, create a query from some x_q

$$q = \theta_Q^\top x_q, \quad \theta_Q \in \mathbb{R}^{d_x \times d_h}, \quad q \in \mathbb{R}^{d_h}$$

To determine if a key and query match, we will take the dot product

Keys and Queries

Now, create a query from some x_q

$$q = \theta_Q^\top x_q, \quad \theta_Q \in \mathbb{R}^{d_x \times d_h}, \quad q \in \mathbb{R}^{d_h}$$

To determine if a key and query match, we will take the dot product

$$q^\top k_i = (\theta_Q^\top x_q)^\top (\theta_K^\top x_i)$$

Keys and Queries

Now, create a query from some x_q

$$q = \theta_Q^\top x_q, \quad \theta_Q \in \mathbb{R}^{d_x \times d_h}, \quad q \in \mathbb{R}^{d_h}$$

To determine if a key and query match, we will take the dot product

$$q^\top k_i = (\theta_Q^\top x_q)^\top (\theta_K^\top x_i)$$

Question: What is the shape of $q^\top k_i$?

Keys and Queries

Now, create a query from some x_q

$$q = \theta_Q^\top x_q, \quad \theta_Q \in \mathbb{R}^{d_x \times d_h}, \quad q \in \mathbb{R}^{d_h}$$

To determine if a key and query match, we will take the dot product

$$q^\top k_i = (\theta_Q^\top x_q)^\top (\theta_K^\top x_i)$$

Question: What is the shape of $q^\top k_i$?

Answer: $(1, d_h) \times (d_h, 1) = 1$, the output is a scalar

Keys and Queries

Now, create a query from some x_q

$$q = \theta_Q^\top x_q, \quad \theta_Q \in \mathbb{R}^{d_x \times d_h}, \quad q \in \mathbb{R}^{d_h}$$

To determine if a key and query match, we will take the dot product

$$q^\top k_i = (\theta_Q^\top x_q)^\top (\theta_K^\top x_i)$$

Question: What is the shape of $q^\top k_i$?

Answer: $(1, d_h) \times (d_h, 1) = 1$, the output is a scalar

Large dot product \Rightarrow match! Small dot product \Rightarrow no match.

Keys and Queries

Example:

Keys and Queries

Example:

$$k_i = \theta_K^\top$$



Keys and Queries

Example:

$$k_i = \theta_K^\top$$



$$q = \theta_Q^\top \text{ Musician}$$

Keys and Queries

Example:

$$k_i = \theta_K^\top$$



$$q = \theta_Q^\top \text{ Musician}$$

$$q^\top k_i = (\theta_Q^\top \text{ Musician})^\top \begin{pmatrix} \theta_K^\top & \text{Musician} \end{pmatrix} = 100$$

Keys and Queries

Example:

$$k_i = \theta_K^\top$$



$$q = \theta_Q^\top \text{ Musician}$$

$$q^\top k_i = (\theta_Q^\top \text{ Musician})^\top \begin{pmatrix} \theta_K^\top & \text{Musician} \end{pmatrix} = 100$$

Large attention!

Keys and Queries

Example:

$$k_i = \theta_K^\top$$



Keys and Queries

Example:

$$k_i = \theta_K^\top$$



$$q = \theta_Q^\top \text{ Mathematician}$$

Keys and Queries

Example:

$$k_i = \theta_K^\top$$



$$q = \theta_Q^\top \text{ Mathematician}$$

$$q^\top k_i = (\theta_Q^\top \text{ Mathematician})^\top \begin{pmatrix} \theta_K^\top \\ \text{Mathematician} \end{pmatrix} = -50$$

Keys and Queries

Example:

$$k_i = \theta_K^\top$$



$$q = \theta_Q^\top \text{ Mathematician}$$

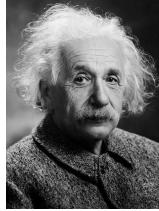
$$q^\top k_i = (\theta_Q^\top \text{ Mathematician})^\top \begin{pmatrix} \theta_K^\top \\ \text{Mathematician} \end{pmatrix} = -50$$

Small attention!

Keys and Queries

Example:

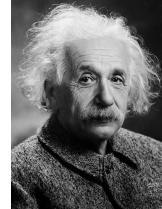
$$k_i = \theta_K^\top$$



Keys and Queries

Example:

$$k_i = \theta_K^\top$$

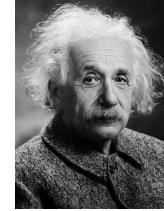


$$q = \theta_Q^\top \text{ Mathematician}$$

Keys and Queries

Example:

$$k_i = \theta_K^\top$$



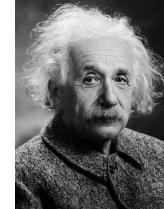
$$q = \theta_Q^\top \text{ Mathematician}$$

$$q^\top k_i = (\theta_Q^\top \text{ Mathematician})^\top \left(\theta_K^\top \begin{array}{c} \text{Mathematician} \\ \text{Portrait of Einstein} \end{array} \right) = 90$$

Keys and Queries

Example:

$$k_i = \theta_K^\top$$



$$q = \theta_Q^\top \text{ Mathematician}$$

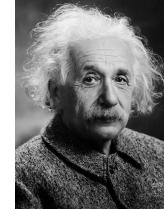
$$q^\top k_i = (\theta_Q^\top \text{ Mathematician})^\top \begin{pmatrix} \theta_K^\top \\ \text{Mathematician} \end{pmatrix} = 90$$

Large attention!

Keys and Queries

Example:

$$k_i = \theta_K^\top$$



$$q = \theta_Q^\top \text{ Mathematician}$$

$$q^\top k_i = (\theta_Q^\top \text{ Mathematician})^\top \left(\theta_K^\top \begin{array}{c} \text{Mathematician} \\ \text{Albert Einstein} \end{array} \right) = 90$$

Large attention!

Remember, there are multiple inputs to pay attention to

Keys and Queries

We compute a key for each input

$$\mathbf{K} = [k_1 \ k_2 \ \dots \ k_T] = [\theta_K^\top x_1 \ \theta_K^\top x_2 \ \dots \ \theta_K^\top x_T], \quad \mathbf{K} \in \mathbb{R}^{d_h \times T}$$

Keys and Queries

We compute a key for each input

$$\mathbf{K} = [k_1 \ k_2 \ \dots \ k_T] = [\theta_K^\top x_1 \ \theta_K^\top x_2 \ \dots \ \theta_K^\top x_T], \quad \mathbf{K} \in \mathbb{R}^{d_h \times T}$$

Keys and Queries

Then compute attention for each key

$$\mathbf{q}^\top \mathbf{K} = \mathbf{q}^\top [\mathbf{k}_1 \ \mathbf{k}_2 \ \dots \ \mathbf{k}_T] = [\mathbf{q}^\top \mathbf{k}_1 \ \mathbf{q}^\top \mathbf{k}_2 \ \dots \ \mathbf{q}^\top \mathbf{k}_T]$$

Keys and Queries

Then compute attention for each key

$$\mathbf{q}^\top \mathbf{K} = \mathbf{q}^\top [\mathbf{k}_1 \ \mathbf{k}_2 \ \dots \ \mathbf{k}_T] = [\mathbf{q}^\top \mathbf{k}_1 \ \mathbf{q}^\top \mathbf{k}_2 \ \dots \ \mathbf{q}^\top \mathbf{k}_T]$$

$$\mathbf{q}^\top \in \mathbb{R}^{1, d_h}, \quad \mathbf{K} \in \mathbb{R}^{d_h \times T}$$

Keys and Queries

Then compute attention for each key

$$\mathbf{q}^\top \mathbf{K} = \mathbf{q}^\top [\mathbf{k}_1 \ \mathbf{k}_2 \ \dots \ \mathbf{k}_T] = [\mathbf{q}^\top \mathbf{k}_1 \ \mathbf{q}^\top \mathbf{k}_2 \ \dots \ \mathbf{q}^\top \mathbf{k}_T]$$

$$\mathbf{q}^\top \in \mathbb{R}^{1, d_h}, \quad \mathbf{K} \in \mathbb{R}^{d_h \times T}$$

Question: What is the shape of $\mathbf{q}^\top \mathbf{K}$?

Keys and Queries

Then compute attention for each key

$$\mathbf{q}^\top \mathbf{K} = \mathbf{q}^\top [\mathbf{k}_1 \ \mathbf{k}_2 \ \dots \ \mathbf{k}_T] = [\mathbf{q}^\top \mathbf{k}_1 \ \mathbf{q}^\top \mathbf{k}_2 \ \dots \ \mathbf{q}^\top \mathbf{k}_T]$$

$$\mathbf{q}^\top \in \mathbb{R}^{1, d_h}, \quad \mathbf{K} \in \mathbb{R}^{d_h \times T}$$

Question: What is the shape of $\mathbf{q}^\top \mathbf{K}$?

Answer: T

Keys and Queries

Then compute attention for each key

$$\mathbf{q}^\top \mathbf{K} = \mathbf{q}^\top [\mathbf{k}_1 \ \mathbf{k}_2 \ \dots \ \mathbf{k}_T] = [\mathbf{q}^\top \mathbf{k}_1 \ \mathbf{q}^\top \mathbf{k}_2 \ \dots \ \mathbf{q}^\top \mathbf{k}_T]$$

$$\mathbf{q}^\top \in \mathbb{R}^{1, d_h}, \quad \mathbf{K} \in \mathbb{R}^{d_h \times T}$$

Question: What is the shape of $\mathbf{q}^\top \mathbf{K}$?

Answer: T

Do not forget to normalize with softmax!

Keys and Queries

Normalize, only pay attention to important matches

Keys and Queries

Normalize, only pay attention to important matches

$$\begin{aligned}\text{softmax}(\mathbf{q}^\top \mathbf{K}) &= \text{softmax}(\mathbf{q}^\top [\mathbf{k}_1 \ \mathbf{k}_2 \ \dots \ \mathbf{k}_T]) \\ &= \text{softmax}([\mathbf{q}^\top \mathbf{k}_1 \ \mathbf{q}^\top \mathbf{k}_2 \ \dots \ \mathbf{q}^\top \mathbf{k}_T])\end{aligned}$$

Keys and Queries

Normalize, only pay attention to important matches

$$\begin{aligned}\text{softmax}(\mathbf{q}^\top \mathbf{K}) &= \text{softmax}(\mathbf{q}^\top [\mathbf{k}_1 \ \mathbf{k}_2 \ \dots \ \mathbf{k}_T]) \\ &= \text{softmax}([\mathbf{q}^\top \mathbf{k}_1 \ \mathbf{q}^\top \mathbf{k}_2 \ \dots \ \mathbf{q}^\top \mathbf{k}_T])\end{aligned}$$

We call this **dot-product attention**

Keys and Queries

Query: Which person will help me on my exam?

Keys and Queries

Query: Which person will help me on my exam?



$$q^\top k_1$$

$$q^\top k_2$$

$$q^\top k_3$$

$$q^\top k_4$$

$$q^\top k_5$$

Keys and Queries

Query: Which person will help me on my exam?



$$q^\top k_1$$

$$q^\top k_2$$

$$q^\top k_3$$

$$q^\top k_4$$

$$q^\top k_5$$

−1.71

0.60

−1.01

−0.61

2.73

softmax

Keys and Queries

Query: Which person will help me on my exam?



$$q^\top k_1$$

$$q^\top k_2$$

$$q^\top k_3$$

$$q^\top k_4$$

$$q^\top k_5$$

−1.71

0.60

−1.01

−0.61

2.73

softmax

0.01

0.10

0.02

0.03

0.84

Keys and Queries

Query: Which person will help me on my exam?



$$q^\top k_1$$

$$q^\top k_2$$

$$q^\top k_3$$

$$q^\top k_4$$

$$q^\top k_5$$

−1.71

0.60

−1.01

−0.61

2.73

softmax

0.01

0.10

0.02

0.03

0.84

Keys and Queries

Put dot product attention into our composite model

Keys and Queries

Put dot product attention into our composite model

$$\lambda \left(\begin{bmatrix} \boldsymbol{x}_1 \\ \vdots \\ \boldsymbol{x}_T \end{bmatrix}, \boldsymbol{\theta}_{\lambda} \right)_i = \text{softmax}(\boldsymbol{q}^\top [\boldsymbol{k}_1 \ \boldsymbol{k}_2 \ \dots \ \boldsymbol{k}_T])_i$$

Keys and Queries

Put dot product attention into our composite model

$$\lambda \left(\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \theta_\lambda \right)_i = \text{softmax}(\mathbf{q}^\top [\mathbf{k}_1 \ \mathbf{k}_2 \ \dots \ \mathbf{k}_T])_i$$

Then, write our composite memory model with dot product attention

Keys and Queries

Put dot product attention into our composite model

$$\lambda \left(\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \boldsymbol{\theta}_{\lambda} \right)_i = \text{softmax}(\mathbf{q}^\top [\mathbf{k}_1 \ \mathbf{k}_2 \ \dots \ \mathbf{k}_T])_i$$

Then, write our composite memory model with dot product attention

$$f \left(\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \boldsymbol{\theta} \right) = \sum_{i=1}^T \boldsymbol{\theta}^\top \mathbf{x}_i \cdot \lambda \left(\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \boldsymbol{\theta}_{\lambda} \right)_i$$

Keys and Queries

$$f\left(\begin{bmatrix} x_1 \\ \vdots \\ x_T \end{bmatrix}, \theta\right) = \sum_{i=1}^T \theta^\top x_i \cdot \lambda\left(\begin{bmatrix} x_1 \\ \vdots \\ x_T \end{bmatrix}, \theta_\lambda\right)_i$$

Keys and Queries

$$f\left(\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \boldsymbol{\theta}\right) = \sum_{i=1}^T \boldsymbol{\theta}^\top \mathbf{x}_i \cdot \lambda\left(\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \boldsymbol{\theta}_\lambda\right)_i$$

We relabel $\boldsymbol{\theta}$ to $\boldsymbol{\theta}_V$

$$f\left(\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \boldsymbol{\theta}\right) = \sum_{i=1}^T \boldsymbol{\theta}_{\textcolor{red}{V}}^\top \mathbf{x}_i \cdot \lambda\left(\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \boldsymbol{\theta}_\lambda\right)_i$$

Keys and Queries

$$f\left(\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \boldsymbol{\theta}\right) = \sum_{i=1}^T \boldsymbol{\theta}^\top \mathbf{x}_i \cdot \lambda\left(\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \boldsymbol{\theta}_\lambda\right)_i$$

We relabel $\boldsymbol{\theta}$ to $\boldsymbol{\theta}_V$

$$f\left(\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \boldsymbol{\theta}\right) = \sum_{i=1}^T \boldsymbol{\theta}_{\textcolor{red}{V}}^\top \mathbf{x}_i \cdot \lambda\left(\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \boldsymbol{\theta}_\lambda\right)_i$$

In dot-product attention, we call $\boldsymbol{\theta}_V^\top \mathbf{x}_i$ the **value**

Self Attention

Self Attention

Previously, we chose our own query $x_q = \text{Musician}$

Self Attention

Previously, we chose our own query $x_q = \text{Musician}$

We can also create queries from all the inputs

Self Attention

Previously, we chose our own query $x_q = \text{Musician}$

We can also create queries from all the inputs

$$Q = [q_1 \ q_2 \ \dots \ q_T] = [\theta_Q^\top x_1 \ \theta_Q^\top x_2 \ \dots \ \theta_Q^\top x_T], \quad Q \in \mathbb{R}^{T \times d_h}$$

Self Attention

Previously, we chose our own query $x_q = \text{Musician}$

We can also create queries from all the inputs

$$Q = [q_1 \ q_2 \ \dots \ q_T] = [\theta_Q^\top x_1 \ \theta_Q^\top x_2 \ \dots \ \theta_Q^\top x_T], \quad Q \in \mathbb{R}^{T \times d_h}$$

We call this dot-product **self** attention

Self Attention

$$Q = [q_1 \ \dots \ q_T] = [\theta_Q^\top x_1 \ \dots \ \theta_Q^\top x_T]$$

$$K = [k_1 \ \dots \ k_T] = [\theta_K^\top x_1 \ \dots \ \theta_K^\top x_T]$$

$$V = [v_1 \ \dots \ v_T] = [\theta_V^\top x_1 \ \dots \ \theta_V^\top x_T]$$

Self Attention

$$Q = [q_1 \ \dots \ q_T] = [\theta_Q^\top x_1 \ \dots \ \theta_Q^\top x_T]$$

$$K = [k_1 \ \dots \ k_T] = [\theta_K^\top x_1 \ \dots \ \theta_K^\top x_T]$$

$$V = [v_1 \ \dots \ v_T] = [\theta_V^\top x_1 \ \dots \ \theta_V^\top x_T]$$

$$QK^\top = \begin{bmatrix} q_1 \\ \vdots \\ q_T \end{bmatrix} [k_1 \ \dots \ k_T] = \begin{bmatrix} q_1 k_1 & \dots & q_1 k_T \\ q_2 k_1 & \dots & q_2 k_T \\ \vdots & \ddots & \vdots \\ q_T k_1 & \dots & q_T k_T \end{bmatrix}$$

Self Attention

$$\mathbf{V} = [\mathbf{v}_1 \ \dots \ \mathbf{v}_T] = [\theta_V^\top \mathbf{x}_1 \ \dots \ \theta_V^\top \mathbf{x}_T]$$

$$\mathbf{QK}^\top = \begin{bmatrix} \mathbf{q}_1 \\ \vdots \\ \mathbf{q}_T \end{bmatrix} [\mathbf{k}_1 \ \dots \ \mathbf{k}_T] = \begin{bmatrix} \mathbf{q}_1 \mathbf{k}_1 & \dots & \mathbf{q}_1 \mathbf{k}_T \\ \mathbf{q}_2 \mathbf{k}_1 & \dots & \mathbf{q}_2 \mathbf{k}_T \\ \vdots & \ddots & \vdots \\ \mathbf{q}_T \mathbf{k}_1 & \dots & \mathbf{q}_T \mathbf{k}_T \end{bmatrix}$$

$$\mathbf{QK}^\top \mathbf{V} = \begin{bmatrix} \mathbf{q}_1 \mathbf{k}_1 & \dots & \mathbf{q}_1 \mathbf{k}_T \\ \mathbf{q}_2 \mathbf{k}_1 & \dots & \mathbf{q}_2 \mathbf{k}_T \\ \vdots & \ddots & \vdots \\ \mathbf{q}_T \mathbf{k}_1 & \dots & \mathbf{q}_T \mathbf{k}_T \end{bmatrix} \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \vdots \\ \mathbf{v}_T \end{bmatrix} = \begin{bmatrix} \mathbf{q}_1 \mathbf{k}_1 \mathbf{v}_1 + \dots + \mathbf{q}_1 \mathbf{k}_T \mathbf{v}_T \\ \mathbf{q}_2 \mathbf{k}_1 \mathbf{v}_1 + \dots + \mathbf{q}_2 \mathbf{k}_T \mathbf{v}_T \\ \vdots \\ \mathbf{q}_T \mathbf{k}_1 \mathbf{v}_1 + \dots + \mathbf{q}_T \mathbf{k}_T \mathbf{v}_T \end{bmatrix}$$

Self Attention

$$QK^\top V = \begin{bmatrix} q_1 k_1 & \dots & q_1 k_T \\ q_2 k_1 & \dots & q_2 k_T \\ \vdots & \ddots & \vdots \\ q_T k_1 & \dots & q_T k_T \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_T \end{bmatrix} = \begin{bmatrix} q_1 k_1 v_1 + \dots + q_1 k_T v_T \\ q_2 k_1 v_1 + \dots + q_2 k_T v_T \\ \vdots \\ q_T k_1 v_1 + \dots + q_T k_T v_T \end{bmatrix}$$

Self Attention

$$QK^\top V = \begin{bmatrix} q_1 k_1 & \dots & q_1 k_T \\ q_2 k_1 & \dots & q_2 k_T \\ \vdots & \ddots & \vdots \\ q_T k_1 & \dots & q_T k_T \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_T \end{bmatrix} = \begin{bmatrix} q_1 k_1 v_1 + \dots + q_1 k_T v_T \\ q_2 k_1 v_1 + \dots + q_2 k_T v_T \\ \vdots \\ q_T k_1 v_1 + \dots + q_T k_T v_T \end{bmatrix}$$

Question: Did we forget anything?

Self Attention

$$QK^\top V = \begin{bmatrix} q_1 k_1 & \dots & q_1 k_T \\ q_2 k_1 & \dots & q_2 k_T \\ \vdots & \ddots & \vdots \\ q_T k_1 & \dots & q_T k_T \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_T \end{bmatrix} = \begin{bmatrix} q_1 k_1 v_1 + \dots + q_1 k_T v_T \\ q_2 k_1 v_1 + \dots + q_2 k_T v_T \\ \vdots \\ q_T k_1 v_1 + \dots + q_T k_T v_T \end{bmatrix}$$

Question: Did we forget anything? **Answer:** Softmax!

Self Attention

$$QK^\top V = \begin{bmatrix} q_1 k_1 & \dots & q_1 k_T \\ q_2 k_1 & \dots & q_2 k_T \\ \vdots & \ddots & \vdots \\ q_T k_1 & \dots & q_T k_T \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_T \end{bmatrix} = \begin{bmatrix} q_1 k_1 v_1 + \dots + q_1 k_T v_T \\ q_2 k_1 v_1 + \dots + q_2 k_T v_T \\ \vdots \\ q_T k_1 v_1 + \dots + q_T k_T v_T \end{bmatrix}$$

Question: Did we forget anything? **Answer:** Softmax!

$$\text{softmax}(QK^\top)V = \text{softmax}\left(\begin{bmatrix} q_1 k_1 & \dots & q_1 k_T \\ q_2 k_1 & \dots & q_2 k_T \\ \vdots & \ddots & \vdots \\ q_T k_1 & \dots & q_T k_T \end{bmatrix}\right) \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_T \end{bmatrix}$$

Self Attention

$$\text{softmax}(\mathbf{Q}\mathbf{K}^\top)\mathbf{V} = \text{softmax} \left(\begin{bmatrix} \mathbf{q}_1 \mathbf{k}_1 & \dots & \mathbf{q}_1 \mathbf{k}_T \\ \mathbf{q}_2 \mathbf{k}_1 & \dots & \mathbf{q}_2 \mathbf{k}_T \\ \vdots & \ddots & \vdots \\ \mathbf{q}_T \mathbf{k}_1 & \dots & \mathbf{q}_T \mathbf{k}_T \end{bmatrix} \right) \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \vdots \\ \mathbf{v}_T \end{bmatrix}$$

Self Attention

$$\text{softmax}(\mathbf{Q}\mathbf{K}^\top)\mathbf{V} = \text{softmax}\left(\begin{bmatrix} q_1\mathbf{k}_1 & \dots & q_1\mathbf{k}_T \\ q_2\mathbf{k}_1 & \dots & q_2\mathbf{k}_T \\ \vdots & \ddots & \vdots \\ q_T\mathbf{k}_1 & \dots & q_T\mathbf{k}_T \end{bmatrix}\right) \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \vdots \\ \mathbf{v}_T \end{bmatrix}$$

Be very careful with softmax dimension

Self Attention

$$\text{softmax}(\mathbf{Q}\mathbf{K}^\top)\mathbf{V} = \text{softmax}\left(\begin{bmatrix} \mathbf{q}_1\mathbf{k}_1 & \dots & \mathbf{q}_1\mathbf{k}_T \\ \mathbf{q}_2\mathbf{k}_1 & \dots & \mathbf{q}_2\mathbf{k}_T \\ \vdots & \ddots & \vdots \\ \mathbf{q}_T\mathbf{k}_1 & \dots & \mathbf{q}_T\mathbf{k}_T \end{bmatrix}\right) \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \vdots \\ \mathbf{v}_T \end{bmatrix}$$

Be very careful with softmax dimension

- Writing attention different ways uses different softmax dimensions

Self Attention

$$\text{softmax}(\mathbf{Q}\mathbf{K}^\top)\mathbf{V} = \text{softmax}\left(\begin{bmatrix} \mathbf{q}_1\mathbf{k}_1 & \dots & \mathbf{q}_1\mathbf{k}_T \\ \mathbf{q}_2\mathbf{k}_1 & \dots & \mathbf{q}_2\mathbf{k}_T \\ \vdots & \ddots & \vdots \\ \mathbf{q}_T\mathbf{k}_1 & \dots & \mathbf{q}_T\mathbf{k}_T \end{bmatrix}\right) \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \vdots \\ \mathbf{v}_T \end{bmatrix}$$

Be very careful with softmax dimension

- Writing attention different ways uses different softmax dimensions
- When in doubt, write out full matrix

Self Attention

$$\text{softmax}(\mathbf{Q}\mathbf{K}^\top)\mathbf{V} = \text{softmax}\left(\begin{bmatrix} q_1\mathbf{k}_1 & \dots & q_1\mathbf{k}_T \\ q_2\mathbf{k}_1 & \dots & q_2\mathbf{k}_T \\ \vdots & \ddots & \vdots \\ q_T\mathbf{k}_1 & \dots & q_T\mathbf{k}_T \end{bmatrix}\right) \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \vdots \\ \mathbf{v}_T \end{bmatrix}$$

Be very careful with softmax dimension

- Writing attention different ways uses different softmax dimensions
- When in doubt, write out full matrix

Question: Softmax rows or columns? **Answer:** Rows

Self Attention

$$\text{softmax}(\mathbf{Q}\mathbf{K}^\top)\mathbf{V} = \text{softmax}\left(\begin{bmatrix} q_1\mathbf{k}_1 & \dots & q_1\mathbf{k}_T \\ q_2\mathbf{k}_1 & \dots & q_2\mathbf{k}_T \\ \vdots & \ddots & \vdots \\ q_T\mathbf{k}_1 & \dots & q_T\mathbf{k}_T \end{bmatrix}\right) \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \vdots \\ \mathbf{v}_T \end{bmatrix}$$

Be very careful with softmax dimension

- Writing attention different ways uses different softmax dimensions
- When in doubt, write out full matrix

Question: Softmax rows or columns? **Answer:** Rows

Only one query index

$\text{softmax}([\mathbf{q}_i\mathbf{k}_1 \ \dots \ \mathbf{q}_i\mathbf{k}_T])$

Self Attention

$$\text{softmax}(\mathbf{Q}\mathbf{K}^\top)\mathbf{V} = \text{softmax}\left(\begin{bmatrix} q_1\mathbf{k}_1 & \dots & q_1\mathbf{k}_T \\ q_2\mathbf{k}_1 & \dots & q_2\mathbf{k}_T \\ \vdots & \ddots & \vdots \\ q_T\mathbf{k}_1 & \dots & q_T\mathbf{k}_T \end{bmatrix}\right) \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \vdots \\ \mathbf{v}_T \end{bmatrix}$$

Be very careful with softmax dimension

- Writing attention different ways uses different softmax dimensions
- When in doubt, write out full matrix

Question: Softmax rows or columns? **Answer:** Rows

Only one query index

$\text{softmax}([\mathbf{q}_i\mathbf{k}_1 \ \dots \ \mathbf{q}_i\mathbf{k}_T])$

Self Attention

Attention paper suggests normalizing constant for faster learning

Self Attention

Attention paper suggests normalizing constant for faster learning

$$\text{attn}\left(\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \theta\right) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_h}}\right)\mathbf{V}$$

Self Attention

Attention paper suggests normalizing constant for faster learning

$$\text{attn}\left(\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \theta\right) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_h}}\right)\mathbf{V}$$

This operation powers today's biggest models

Self Attention

$$\mathbf{Q} \in \mathbb{R}^{T \times d_h} \quad \mathbf{K} \in \mathbb{R}^{T \times d_h} \quad \mathbf{V} \in \mathbb{R}^{T \times d_h}$$

$$\underbrace{\text{attn}\left(\underbrace{\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}}_{\mathbb{R}^{T \times d_x}}, \boldsymbol{\theta}\right)}_{\mathbb{R}^{T \times d_x}} = \text{softmax}\left(\overbrace{\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_h}}}^{\mathbb{R}^{T \times T}}, \underbrace{\mathbf{V}}_{\mathbb{R}^{T \times d_h}}\right)$$

Self Attention

$$\mathbf{Q} \in \mathbb{R}^{T \times d_h} \quad \mathbf{K} \in \mathbb{R}^{T \times d_h} \quad \mathbf{V} \in \mathbb{R}^{T \times d_h}$$

$$\underbrace{\text{attn}\left(\underbrace{\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}}_{\mathbb{R}^{T \times d_x}}, \boldsymbol{\theta}\right)}_{\mathbb{R}^{T \times d_x}} = \text{softmax}\left(\overbrace{\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_h}}}^{\mathbb{R}^{T \times T}}, \underbrace{\mathbf{V}}_{\mathbb{R}^{T \times d_h}}\right)$$

Self Attention

```
class Attention(nn.Module):
    def __init__(self):
        self.theta_K = nn.Linear(d_x, d_h, bias=False)
        self.theta_Q = nn.Linear(d_x, d_h, bias=False)
        self.theta_V = nn.Linear(d_x, d_h, bias=False)

    def forward(self, x):
        # Shape: (column, row), be very careful with axis!
        A = softmax(
            self.theta_Q(x) @ self.theta_K(x).T / d_h, axis=1
        )
        V = self.theta_V(x)
        return A @ V
```

Guest Lecture - Dr. Matteo Bettini

Guest Lecture - Dr. Matteo Bettini



Dr. Matteo Bettini is a good friend of mine

Guest Lecture - Dr. Matteo Bettini



Dr. Matteo Bettini is a good friend of mine
We did our PhDs together

Guest Lecture - Dr. Matteo Bettini

Dr. Matteo Bettini

- Incoming Research Scientist at Meta's SuperIntelligence Lab

Guest Lecture - Dr. Matteo Bettini

Dr. Matteo Bettini

- Incoming Research Scientist at Meta's SuperIntelligence Lab
- Focus on agentic LLMs

Guest Lecture - Dr. Matteo Bettini

Dr. Matteo Bettini

- Incoming Research Scientist at Meta's SuperIntelligence Lab
- Focus on agentic LLMs
- PhD in Computer Science at Cambridge

Guest Lecture - Dr. Matteo Bettini

Dr. Matteo Bettini

- Incoming Research Scientist at Meta's SuperIntelligence Lab
- Focus on agentic LLMs
- PhD in Computer Science at Cambridge
- MS in Computer Science at Cambridge

Guest Lecture - Dr. Matteo Bettini

Dr. Matteo Bettini

- Incoming Research Scientist at Meta's SuperIntelligence Lab
- Focus on agentic LLMs
- PhD in Computer Science at Cambridge
- MS in Computer Science at Cambridge
- BS in Computer Engineering at Milan Polytechnic