

TRANSFORMER



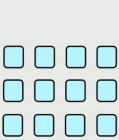
What Is a Transformer?

The Transformer architecture revolutionised NLP by eliminating recurrence and convolution entirely and relying solely on attention mechanisms.

Traditional RNN/LSTM	Transformer
Sequential processing	Fully parallelizable
Bottleneck in information	Simultaneous access to all tokens
Difficult to train for long sequences	Better with long-range dependencies

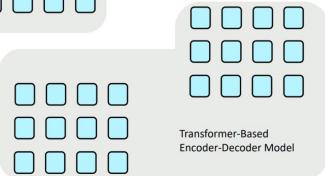
Computational Dependencies for Recurrence vs. Attention

RNN-Based Encoder-Decoder Model with Attention



Transformer Advantages:

- Number of unparallelizable operations does not increase with sequence length.
- Each "word" interacts with each other, so maximum interaction distance: O(1).



Transformer-Based Encoder-Decoder Model

ATTENTION

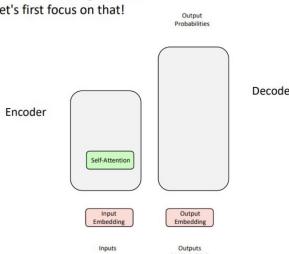
Vaswani, Ashish, et al. "Attention is all you need." Advances in neural information processing systems. 2017.



- Sequential models ingest the input one word or one token at the time. And so, as if each unit was like a bottleneck to the flow of information.
- Transformer ingest an entire sentence all at the same time.
- Attention mechanism + CNN like parallelism.

Encoder: Self-Attention

Self-Attention is the core building block of Transformer, so let's first focus on that!



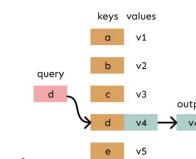
Intuition for Attention Mechanism

- Let's think of attention as a "fuzzy" or approximate hashtable:
 - To look up a **value**, we compare a **query** against **keys** in a table.
 - In a hashtable (shown on the bottom left):
 - Each **query** (hash) maps to exactly one **key-value** pair.
 - In (self-)attention (shown on the bottom right):
 - Each **query** matches each **key** to varying degrees.
 - We return a sum of **values** weighted by the **query-key** match.

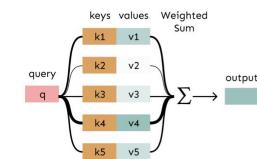
q	k_0	v_0
	k_1	v_1
	k_2	v_2
	k_3	v_3
	k_4	v_4
	k_5	v_5
	k_6	v_6
	k_7	v_7

q	k_0	v_0
	k_1	v_1
	k_2	v_2
	k_3	v_3
	k_4	v_4
	k_5	v_5
	k_6	v_6
	k_7	v_7

In a **lookup table**, we have a table of **keys** that map to **values**. The **query** matches one of the keys, returning its value.



In **attention**, the **query** matches all **keys** *softly*, to a weight between 0 and 1. The **keys' values** are multiplied by the weights and summed.



Recipe for Self-Attention in the Transformer Encoder

- Step 1: For each word x_i , calculate its **query**, **key**, and **value**.

$$q_i = W^Q x_i \quad k_i = W^K x_i \quad v_i = W^V x_i$$

- Step 2: Calculate attention score between **query** and **keys**.

$$e_{ij} = q_i \cdot k_j$$

- Step 3: Take the softmax to normalize attention scores.

$$\alpha_{ij} = \text{softmax}(e_{ij}) = \frac{\exp(e_{ij})}{\sum_k \exp(e_{ik})}$$

- Step 4: Take a weighted sum of **values**.

$$\text{Output}_i = \sum_j \alpha_{ij} v_j$$

q

k_0	v_0
k_1	v_1
k_2	v_2
k_3	v_3
k_4	v_4
k_5	v_5
k_6	v_6
k_7	v_7

Recipe for (Vectorized) Self-Attention in the Transformer Encoder

- Step 1: With embeddings stacked in X , calculate **queries**, **keys**, and **values**.

$$Q = XW^Q \quad K = XW^K \quad V = XW^V$$

- Step 2: Calculate attention scores between **query** and **keys**.

$$E = QK^T$$

- Step 3: Take the softmax to normalize attention scores.

$$A = \text{softmax}(E)$$

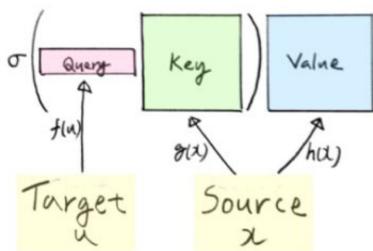
- Step 4: Take a weighted sum of **values**.

$$\text{Output} = AV$$

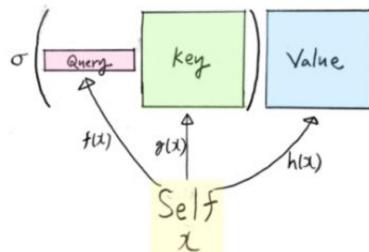
$$\text{Output} = \text{softmax}(QK^T)V$$

$$\text{softmax}\left(\frac{Q \cdot K^T}{\sqrt{d_k}}\right) V = Z$$

(Source-Target-Attention)

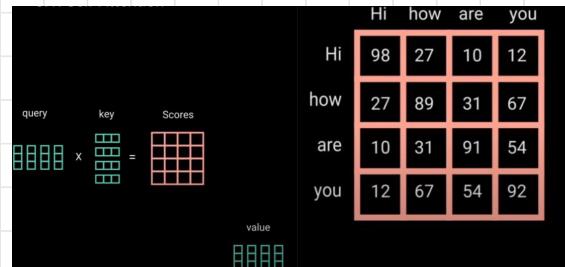
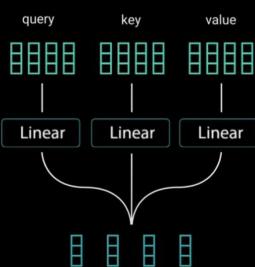
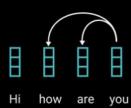


(Self-Attention)



3. Multi-headed Attention

3.1. Self-Attention

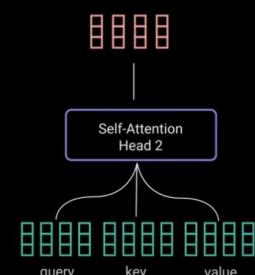
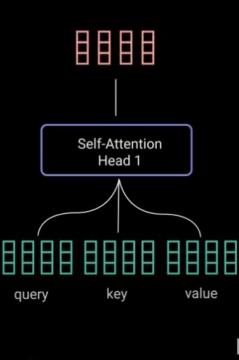
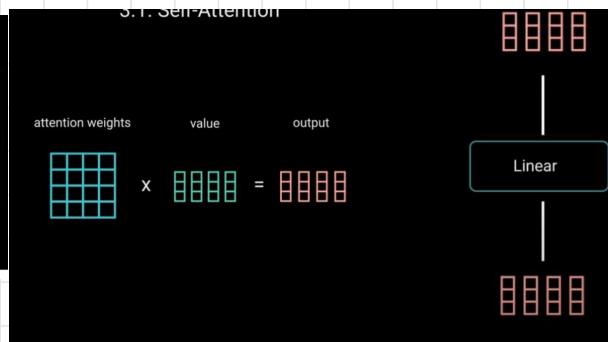


$$\frac{\text{Scaled Scores}}{\sqrt{d_k}}$$

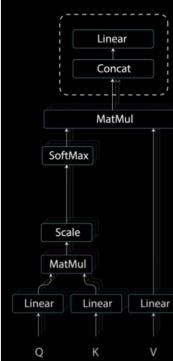
Softmax(Scaled Scores) =

Hi	0.7	0.1	0.1	0.1
how	0.1	0.6	0.2	0.1
are	0.1	0.3	0.6	0.1
you	0.1	0.3	0.3	0.3

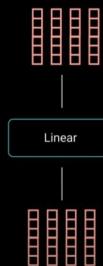
$$\text{softmax}(x)_i = \frac{\exp(x_i)}{\sum_j \exp(x_j)}$$



$N = 2$



3. Multi-headed Attention



Multi-Head Attention

Multi-head attention = performing self-attention multiple times with different learned weight matrices (i.e., different "heads").

Each head:

Looks at the input from a different representation subspace.

Learns diverse relationships (e.g., syntax, semantics, position).

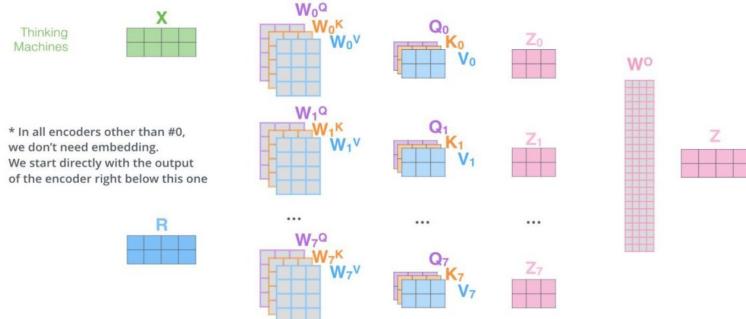
These are concatenated and passed through another linear layer.

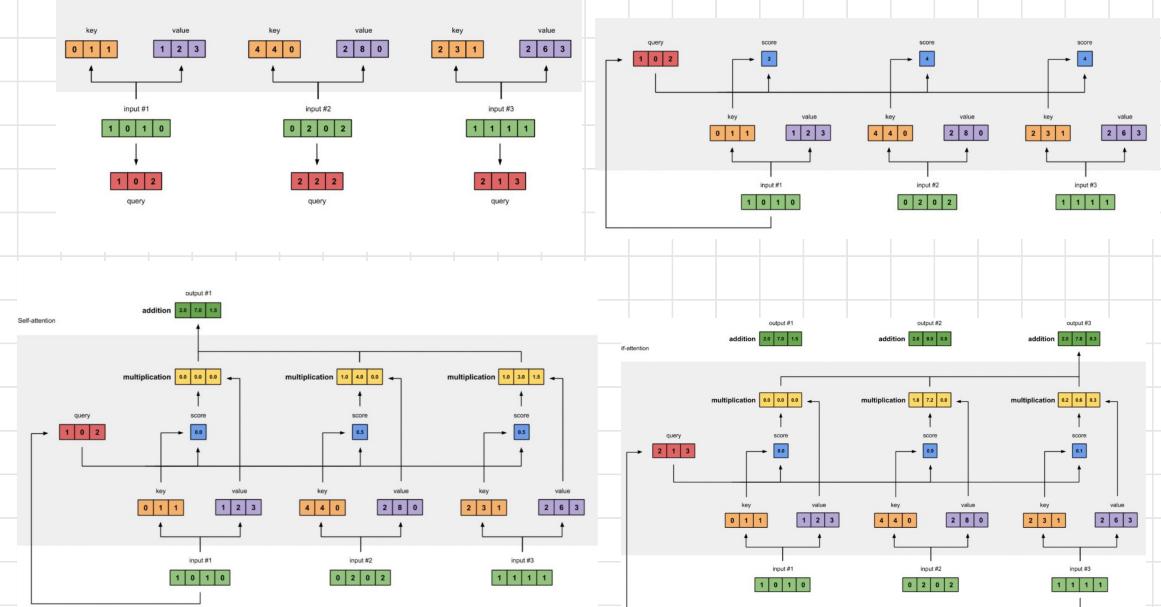
Think of each head as asking a different question about the same sentence.



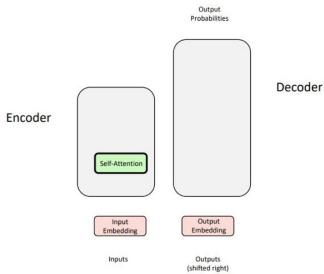
- Doing self-Attention multiple times (Multi-Head).
- As if you are asking a different query about the input multiple times.
- Parallel computation for all heads

- 1) This is our input sentence*
- 2) We embed each word*
- 3) Split into 8 heads. We multiply X or R with weight matrices
- 4) Calculate attention using the resulting $Q/K/V$ matrices
- 5) Concatenate the resulting Z matrices, then multiply with weight matrix W^O to produce the output of the layer



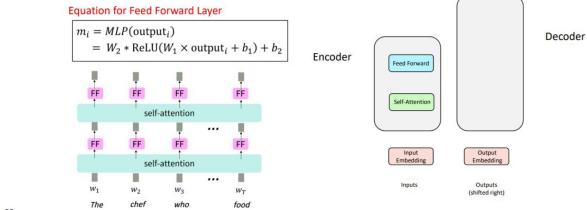


What We Have So Far: (Encoder) Self-Attention!

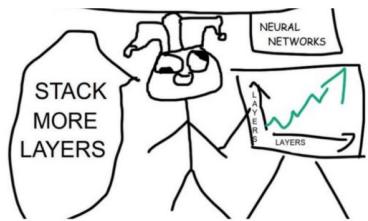


But attention isn't quite all you need!

- **Problem:** Since there are no element-wise non-linearities, self-attention is simply performing a re-averaging of the value vectors.
- **Easy fix:** Apply a feedforward layer to the output of attention, providing non-linear activation (and additional expressive power).



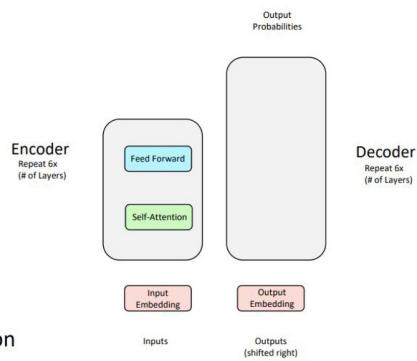
But how do we make this work for deep networks?



Training Trick #1: Residual Connections

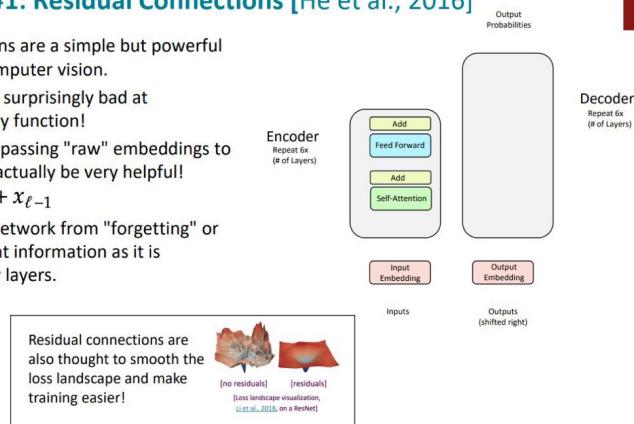
Training Trick #2: LayerNorm

Training Trick #3: Scaled Dot Product Attention



Training Trick #1: Residual Connections [He et al., 2016]

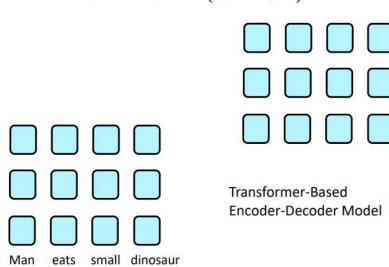
- Residual connections are a simple but powerful technique from computer vision.
 - Deep networks are surprisingly bad at learning the identity function!
 - Therefore, directly passing "raw" embeddings to the next layer can actually be very helpful!
- $$x_\ell = F(x_{\ell-1}) + x_{\ell-1}$$
- This prevents the network from "forgetting" or distorting important information as it is processed by many layers.



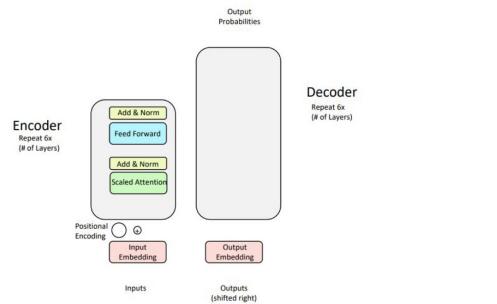
34

Major issue!

- We're almost done with the Encoder, but we have a major problem! Has anyone spotted it?
- Consider this sentence:
 - "Man eats small dinosaur."



Solution: Inject Order Information through Positional Encodings!



Fixing the first self-attention problem: sequence order

- Since self-attention doesn't build in order information, we need to encode the order of the sentence in our keys, queries, and values.
- Consider representing each **sequence index as a vector**
- $p_i \in \mathbb{R}^d$, for $i \in \{1, 2, \dots, T\}$ are position vectors
- Don't worry about what the p_i are made of yet!
- Easy to incorporate this info into our self-attention block: just add the p_i to our inputs!
- Let $\tilde{v}_i, \tilde{k}_i, \tilde{q}_i$ be our old values, keys, and queries.

$$\begin{aligned} v_i &= \tilde{v}_i + p_i \\ q_i &= \tilde{q}_i + p_i \\ k_i &= \tilde{k}_i + p_i \end{aligned}$$

In deep self-attention networks, we do this at the first layer! You could concatenate them as well, but people mostly just add...



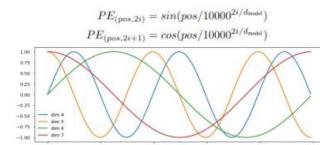
Positional Encoding

Since transformers don't process input sequentially (no recurrence), they add position information manually.

Positional Encoding (usually sine/cosine) is added to each input embedding to capture word order.

Positional Encoding

- Account for the order of the words in the input sequence.
- Adds a vector to each input embedding. These vectors follow a specific pattern that helps determine the position of each word.



$$\vec{p}_t = \begin{bmatrix} \sin(\omega_1 t) \\ \cos(\omega_1 t) \\ \sin(\omega_2 t) \\ \cos(\omega_2 t) \\ \vdots \\ \sin(\omega_{d/2} t) \\ \cos(\omega_{d/2} t) \end{bmatrix}_{d \times 1}$$

2. Positional Encoding

Positional Input Embeddings



Positional Encoding



Time Step

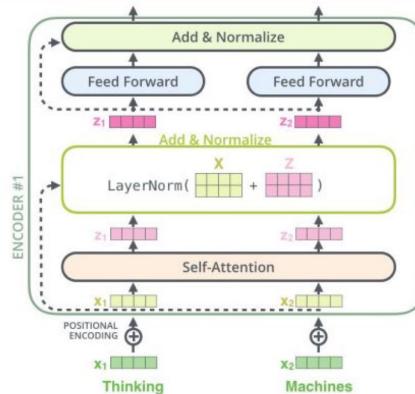
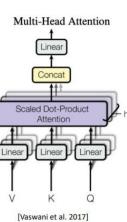
1 2 3 4

$$PE(pos, 2i + 1) = \cos\left(\frac{pos}{10000^{2i/dmodel}}\right)$$

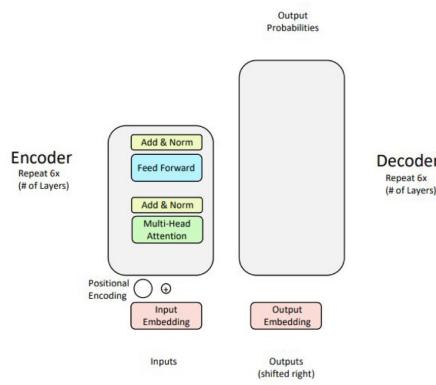
$$PE(pos, 2i) = \sin\left(\frac{pos}{10000^{2i/dmodel}}\right)$$

Multi-Headed Self-Attention: k heads are better than 1!

- **High-Level Idea:** Let's perform self-attention multiple times in parallel and combine the results.

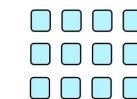
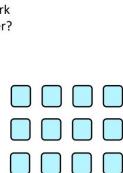


Yay, we've completed the Encoder! Time for the Decoder...



Decoder: Masked Multi-Head Self-Attention

- Problem:** How do we keep the decoder from "cheating"? If we have a language modeling objective, can't the network just look ahead and "see" the answer?
- Solution:** Masked Multi-Head Attention. At a high-level, we hide (mask) information about future tokens from the model.



Transformer-Based Encoder-Decoder Model

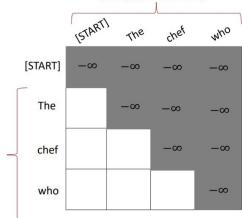
Masking the future in self-attention

- To use self-attention in decoders, we need to ensure we can't peek at the future.
- At every timestep, we could change the set of **keys** and **queries** to include only past words. (Inefficient)
- To enable parallelization, we **mask out attention** to future words by setting attention scores to $-\infty$.

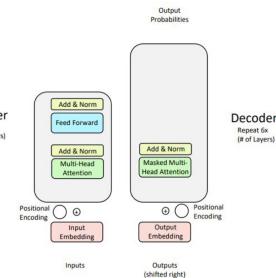
$$e_{ij} = \begin{cases} q_i^T k_j, & j < i \\ -\infty, & j \geq i \end{cases}$$

For encoding these words

We can look at these (not greyed out) words

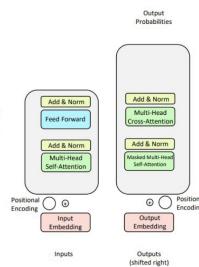


Decoder: Masked Multi-Headed Self-Attention



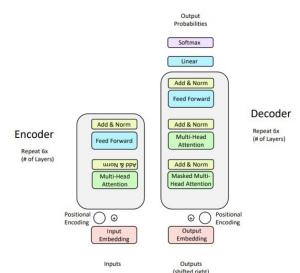
Encoder-Decoder Attention

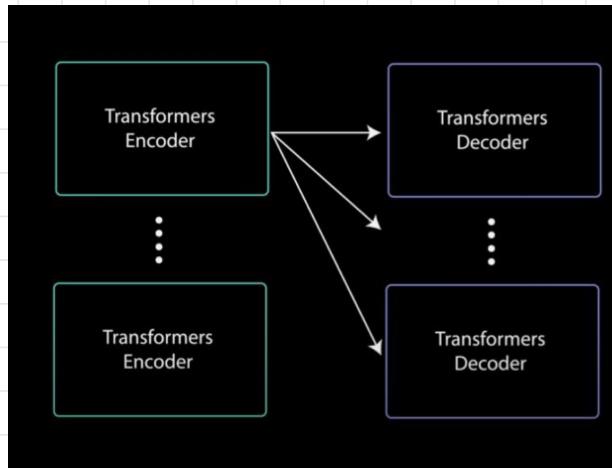
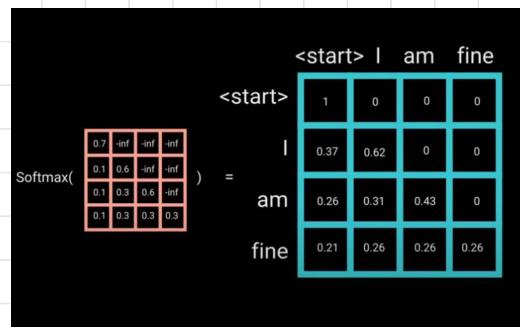
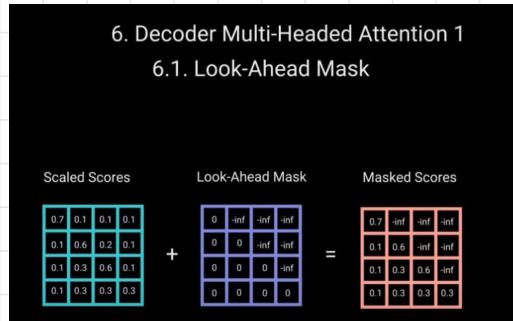
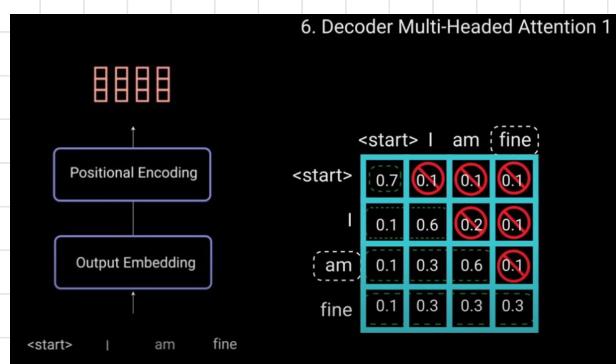
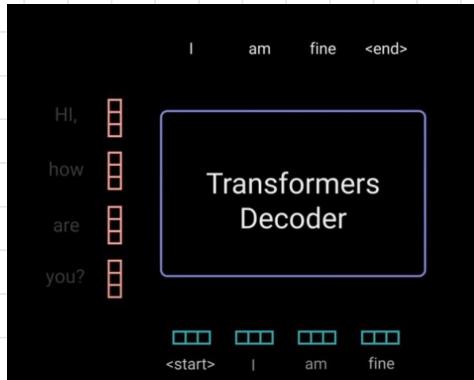
- We saw that self-attention is when keys, queries, and values come from the same source.
- In the decoder, we have attention that looks more like what we saw last week.
- Let h_1, \dots, h_T be **output vectors from the Transformer encoder**, $x_i \in \mathbb{R}^d$
- Let z_1, \dots, z_T be input vectors from the **decoder**, $z_i \in \mathbb{R}^d$
- Then keys and values are drawn from the **encoder** (like a memory):
 - $k_i = Kh_i$, $v_i = Vh_i$.
- And the queries are drawn from the **decoder**, $q_i = Qz_i$.



Decoder: Finishing touches!

- Add a feed forward layer (with residual connections and layer norm)
- Add a final linear layer to project the embeddings into a much longer vector of length vocab size (logits)
- Add a final softmax to generate a probability distribution of possible next words!



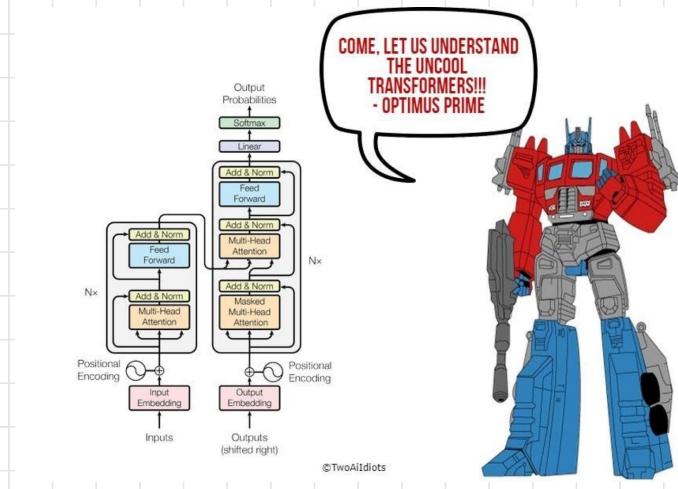


Transformer Architecture

A typical Transformer block contains:

◆ Encoder Block:
Input embeddings + positional encoding
Multi-head self-attention
Feedforward neural network
Layer normalization and residual connections

◆ Decoder Block:
Masked multi-head self-attention
Cross-attention: attends over encoder output
Feedforward + layer norm + residuals



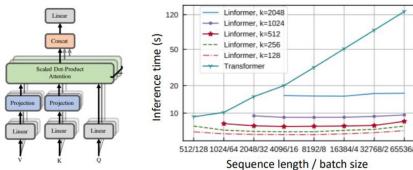
What would we like to fix about the Transformer?

- **Quadratic compute in self-attention (today):**
 - Computing all pairs of interactions means our computation grows **quadratically** with the sequence length!
 - For recurrent models, it only grew linearly!
- **Position representations:**
 - Are simple absolute indices the best we can do to represent position?
 - Relative linear position attention [Shaw et al., 2018]
 - Dependency syntax-based position [Wang et al., 2019]

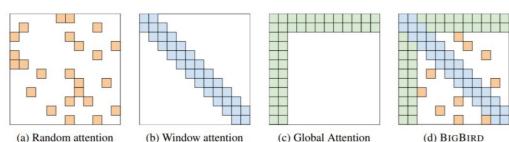
Recent work on improving on quadratic self-attention cost

- Considerable recent work has gone into the question, *Can we build models like Transformers without paying the $O(T^2)$ all-pairs self-attention cost?*
- For example, Linformer [Wang et al., 2020]
- For example, BigBird [Zaheer et al., 2021]

Key idea: map the sequence length dimension to a lower-dimensional space for values, keys



Key idea: replace all-pairs interactions with a family of other interactions, like local windows, looking at everything, and random interactions.



Transformers Variants

12

Pre-Training Objective

◆ Pre-training:

Done unsupervised on massive text data like Wikipedia, QA forums, books.

Learns general language representations: grammar, semantics, context.

◆ Fine-tuning:

Supervised learning on specific downstream tasks (like sentiment analysis, QA, NER).

Requires labeled data and task-specific modifications.

► Problem: Fine-tuning is labor-intensive and domain-specific.

Pre-training vs Fine-tuning Unification

Modern goal: Unify all tasks under one pre-training objective to reduce reliance on fine-tuning.

Type	Goal	Strengths	Models
AR (Autoregressive)	Predict next word (left to right)	Great for text generation	GPT series
AE (Autoencoding)	Reconstruct og text from corrupted data	Great for language understanding	BERT, BART, T5

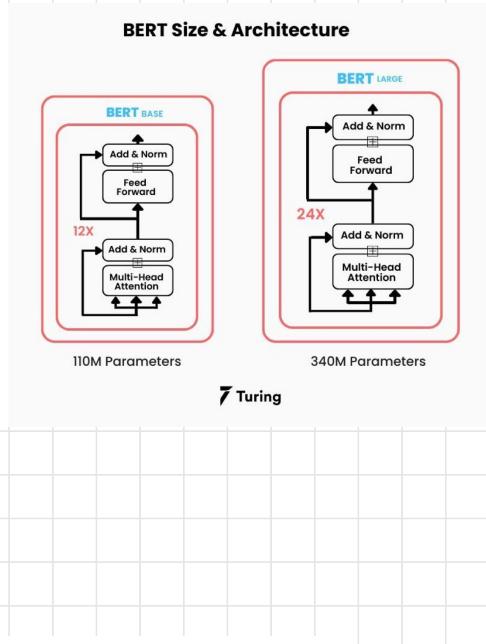
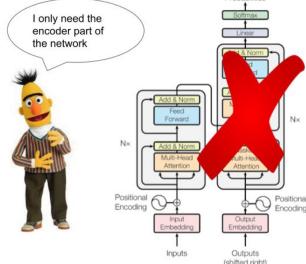
Encoder Only Transformer

BERT (2018) - Bidirectional Encoder Representations from Transformers

Uses only the encoder stack of the transformer.

Unlike GPT (left-to-right), BERT is bidirectional — sees both left and right context.

Outperformed GPT-1 with fewer parameters.



TRANSFORMER VS BERT

Transformer: I am a ___

BERT: I am a ___. I like playing Cricket

Pre-training Objectives in BERT

1. MASKED IM (MLM)

The idea here is "simple":

Randomly mask out 15% of the words in the input

- replacing them with a [MASK] token

- run the entire sequence through the BERT attention based encoder and then predict only the masked words, based on the context provided by the other non-masked words in the sequence.

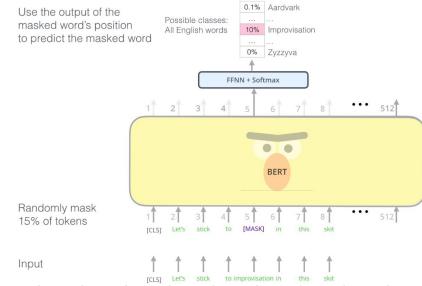
The model only tries to predict when the [MASK] token is present in the input, while we want the model to try to predict the correct tokens regardless of what token is present in the input.

To deal with this issue, out of the 15% of the tokens selected for masking:

- 80% of the tokens are actually replaced with the token [MASK].
- 10% of the time tokens are replaced with a random token.
- 10% of the time tokens are left unchanged.

While training the BERT loss function considers only the prediction of the masked tokens and ignores the prediction of the non-masked ones.

This results in a model that converges much more slowly than left-to-right or right-to-left models.



2. NEXT SENTENCE PREDICTION (NSP)

In order to understand relationship between two sentences, BERT training process also uses next sentence prediction.

During training the model gets as input pairs of sentences and it learns to predict if the second sentence is the next sentence in the original text as well.

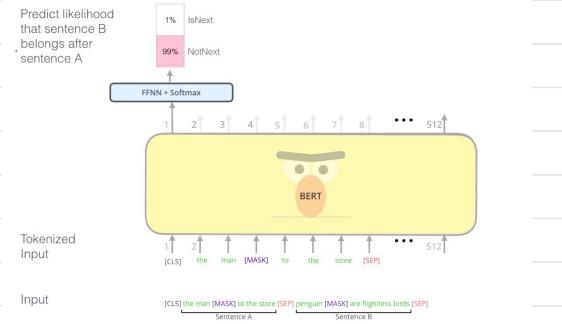
- 50% of the time the second sentence comes after the first one.
- 50% of the time it is a random sentence from the full corpus.

Given sentence pairs:

50%: sentence B is next after A

50%: sentence B is a random sentence

Model learns sentence relationships.



2. Next Sentence Prediction (NSP)

Input = [CLS] the man went to [MASK] store [SEP]
he bought a gallon [MASK] milk [SEP]

Label = IsNext

Input = [CLS] the man [MASK] to the store [SEP]
penguin [MASK] are flight ##less birds [SEP]

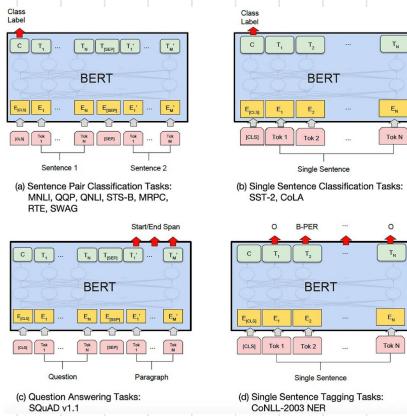
Label = NotNext

FINE TUNING

Real-life language tasks are not denoising tasks, leading to a discrepancy between pre-training and fine-tuning for BERT models. Therefore, fine-tuning is necessary for individual downstream tasks.

BERT categorizes downstream tasks into four types:

- Sentence Pair Classification Tasks (e.g., semantic similarity between two sentences)
- Single Sentence Classification Tasks (e.g., sentiment analysis)
- SQuAD (Question-Answering)
- Named Entity Tagging.



ARCHITECTURE DETAILS

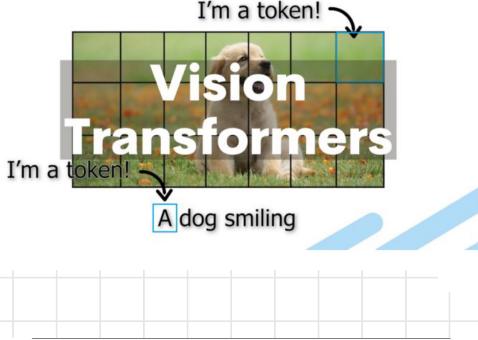
BERT-Base: 12-layer, 768-hidden-nodes, 12-attention-heads, 110M parameters

BERT-Large: 24-layer, 1024-hidden-nodes, 16-attention-heads, 340M parameters

Fun fact: BERT-Base was trained on 4 cloud TPUs for 4 days

BERT-Large was trained on 16 TPUs for 4 days!





Transformer in Transformer

Kai Han^{1,2} An Xiao² Enhua Wu^{1,3*} Jianyuan Guo² Chunjing Xu² Yunhe Wang^{2*}

¹State Key Lab of Computer Science, ISCAS & UCAS

²Noah's Ark Lab, Huawei Technologies

³University of Macau

(hankai, web)@ios.ac.cn, yunhe.wang@huawei.com

Abstract

Transformer is a new kind of neural architecture which encodes the input data as powerful features via the attention mechanism. Basically, the visual transformers first divide the input images into several local patches and then calculate both representations and their relationship. Since natural images are of high complexity with various small and large objects, the global grouping of the patches is not fine enough for recognizing features of objects in different scales and locations. In this paper, we point out that the attention inside these local patches are also essential for building visual transformers with high performance and we explore a new architecture, namely, Transformer IN Transformer (TNT). Specifically, we regard the local patches (e.g., 16×16) as "visual sentences" and present to further divide them into smaller patches (e.g., 4×4) as "visual words". The attention of

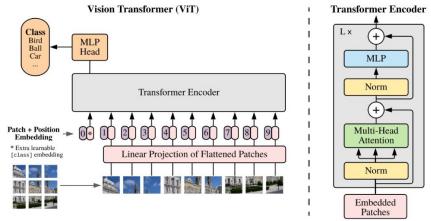
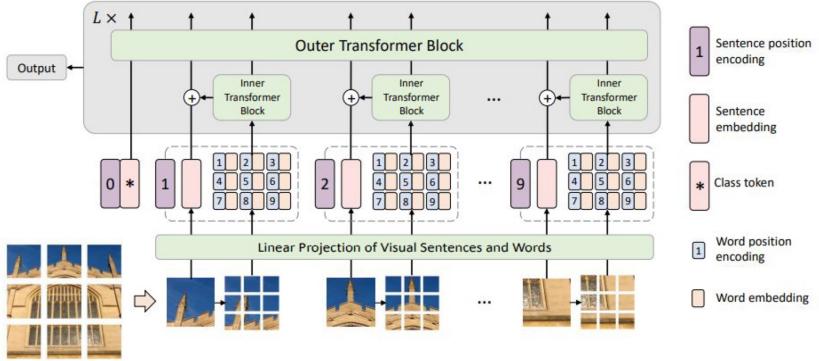


Figure 1: Model overview. We split an image into fixed-size patches, linearly embed each of them, add position embeddings, and feed the resulting sequence of vectors to a standard Transformer encoder. In order to perform classification, we use the standard approach of adding an extra learnable "classification token" to the sequence. The illustration of the Transformer encoder was inspired by Vaswani et al. (2017).



Encoder Decoder Transformers

BART (BIDIRECTIONAL AND AUTO-REGRESSIVE TRANSFORMERS)

BART (2019) – Facebook

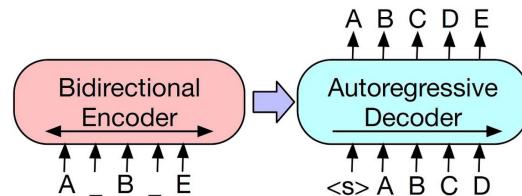
Combines:

- BERT-like encoder (bidirectional)
- GPT-like decoder (autoregressive)

Uses denoising objectives:

Masking 30% tokens

Permuting sentence orders



BART (2019) proposed by Facebook greatly resembles T5 in terms of the denoising pre-training objective and the encoder-decoder architecture, with the only difference being that 30% of the tokens are masked and sentence permutation is used in the pre-training text.

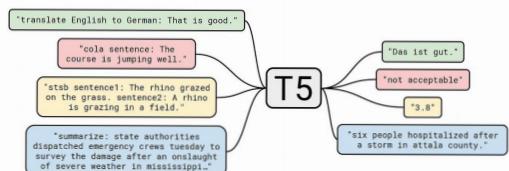
T5 (2020) – Text-to-Text Transfer Transformer (Google)

Unifies all NLP tasks into a text-in, text-out format.

15% tokens masked using sentinel tokens (special tokens that represent blanks).

Learns both understanding + generation.

T5 (TEXT-TO-TEXT TRANSFER TRANSFORMER)



In 2020, Google proposed T5 as a unified model capable of transforming all downstream tasks into text generative tasks, even classification problems.

T5 uses an encoder-decoder architecture and a denoising objective

During pre-training, 15% of the tokens fed into the encoder are randomly masked.

It is a modified version of BERT's masking and denoising algorithm: consecutive masked tokens are replaced by, a sentinel and are treated as a new single token added to the original vocabulary.

This helps the model learn to predict how many words are missing in the blank.

Later, generative capability is learned by randomly splitting the input text in the dataset, with the first part fed into the encoder and the second part treated as the output to be auto-regressively regenerated.

Decoder Only Transformers

Used for generative tasks, especially autoregressive text generation.

Generative Pretrained Transformer (GPT) [Radford et al., 2018]

How do we format inputs to our decoder for finetuning tasks?



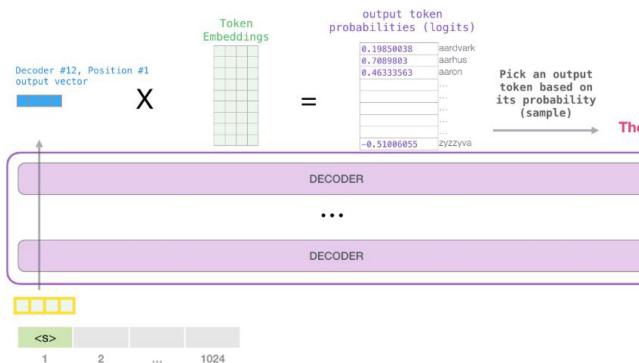
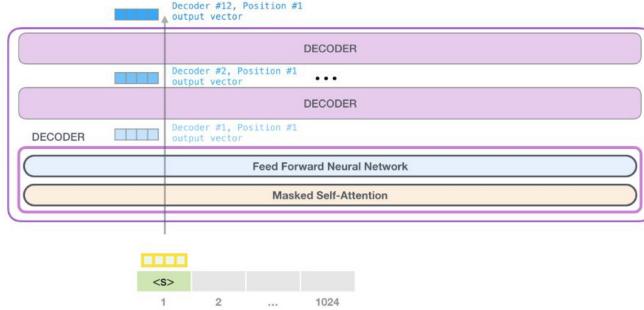
Improving Language Understanding by Generative Pre-Training

Alec Radford
OpenAI
alec@openai.com
Karthik Narasimhan
OpenAI
karthikn@openai.com
Tim Salimans
OpenAI
tim@openai.com
Ilya Sutskever
OpenAI
ilyas@openai.com

Abstract

Natural language understanding comprises a wide range of diverse tasks such as textual entailment, question answering, semantic similarity assessment, and text classification. Although large unlabeled text corpora are abundant, labeled data for fine-tuning specific tasks is scarce, making it challenging for discriminatively trained models to perform adequately. We demonstrate that large gains on these tasks can be realized by *generative pre-training* of a language model on a diverse corpus of unlabeled text, followed by *discriminative fine-tuning* on each specific task. In contrast to previous approaches, we make use of task-aware input transformations during fine-tuning to achieve effective transfer while requiring minimal changes to the model architecture. We demonstrate the effectiveness of

The linear classifier is applied to the representation of the [EXTRACT] token.



GPT - 1

OpenAI proposed the GPT-1-4 models that only used decoder stacks, making them left-to-right autoregressive models.

GPT-1 (2018, 117 million parameters) did not exhibit emergent capabilities and heavily relied on fine-tuning for individual downstream tasks.

GPT - 2

GPT-2 (2019, 1.5 billion parameters) introduced the phenomenon of in-context learning for a few tasks, and improved its tokenizer by using Byte-level Encoding (BLE).

The GPT-2 was trained on a massive 40GB dataset called WebText that the OpenAI researchers crawled from the internet as part of the research effort.

The smallest variant of the trained GPT-2, takes up 500MBs of storage to store all of its parameters.

The largest GPT-2 variant is 13 times the size so it could take up more than 6.5 GBs of storage space.

GPT - 3

GPT-3 (2020, 175 billion parameters) has surprisingly demonstrated strong in-context learning capabilities, including zero-shot and few-shot learning abilities.

In few-shot or zero-shot settings without further fine-tuning, GPT-3 can achieve comparable performance with other fine-tuned state-of-the-art (SOTA) models.

ZERO-SHOT, ONE-SHOT, AND FEW-SHOT LEARNING

These describe how much help (examples or guidance) we give to a model before it has to perform a task

ZERO-SHOT

"Hey model, just do it—no examples."

The model performs a task without seeing any task-specific examples. It relies purely on its pretrained knowledge.

Is the following sentence +ve or -ve?

"I absolutely love this movie!"

Zero-shot

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.



ONE-SHOT

"Here's one example – now you try."

The model is shown one example of the task before doing it.

Classify the sentiment:

Example: "This is terrible." → Negative

Now classify: "This is fantastic."

One-shot

In addition to the task description, the model sees a single example of the task. No gradient updates are performed.



FEW-SHOT

"Here's a few examples – learn the pattern."

The model is given a few examples (2 to ~10) before the real task.

Classify the sentiment:

"Awful experience." → Negative

"I loved it!" → Positive

"Not great, not terrible." → Neutral

Now classify: "Best day ever!"

Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.



Type	Description
Zero-shot	No examples given, model uses pretraining alone
One-shot	One example shown before task
Few-shot	2–10 examples shown before task

CHATGPT

It is a variant of the popular GPT-3 (Generative Pertained Transformer 3) model.

Chat GPT was modified and improved using both **supervised and reinforcement learning methods**, with the assistance of human trainer (RLHF). Chat GPT also has 176 billion parameters same as GPT-3 model. The learning includes 3 Steps.

- Supervised fine tuning of GPT-3.5 Model
- Reward Model
- Proximal Policy Optimization (PPO)

ChatGPT is a pre-trained GPT-3 model that has been fine-tuned using the InstructGPT method.

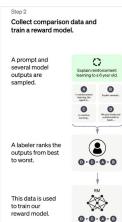
InstructGPT method was used for fine-tuning, which combines supervised learning of demonstration texts from labelers, then with reinforcement learning of generation text scoring and ranking, which are referred to as Reinforcement Learning from Human Feedback (RLHF).

SUPERVISED FINE TUNING (STEP1)

In first Step a pretrained GPT-3 model is used and it will be fine tuned with the help of labelers by creating a supervised dataset. Input Queries were collected from the actual user entries and model generated different responses with respect to that input prompts. The labelers then wrote an appropriate response to the input prompt's (how they want to see that prompt to be answered). The GPT-3 model was then fine-tuned using this new supervised dataset, to create GPT-3.5 model.

Reward Model (Step 2)

In this step SFT model is used and different input/prompts queries fed to the finetuned model and different responses were generated (4 to 7) for every input/prompt. Then labeler determines a reward for each of these outcomes and this reward is proportional to the quality of response with respect to initial prompt.



PROXIMAL POLICY OPTIMIZATION (PPO) RL ALGO- STEP3

In this step we pass unseen input sequences to the clone SFT model we got in step1. The model will generate response with respect to the input prompt. We pass the response to our reward model which we got in step 2 to understand, how high quality was this response for that input prompt and the output reward will be used to finetune the parameters of our SFT model .This is how our SFT model will incorporate more human like characteristics and behavior's via Reinforcement Learning.

ChatGPT (based on GPT-3.5)

✓ Fine-tuning Process (RLHF):

Step 1: Supervised Fine-Tuning (SFT)

Human labelers write ideal responses to prompts.

GPT-3 is fine-tuned on this supervised dataset → GPT-3.5.

Step 2: Reward Model

Several outputs are generated for each prompt.

Human labelers rank the responses.

Model learns to associate higher rewards with better answers

Step 3: Proximal Policy Optimization (PPO)

Reinforcement Learning step.

Uses reward signal to fine-tune GPT-3.5 further.

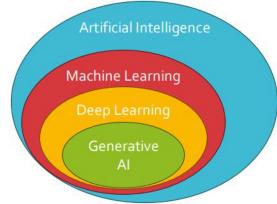
Introduces human-like behaviors and preferences.

Model	Params	Key Features
GPT-1	117M	Required fine-tuning
GPT-2	1.5B	Introduced **in-context learning**
GPT-3	175B	Powerful **zero/few-shot learning**
ChatGPT	175B	Fine-tuned GPT-3 using **RLHF**

Model	Architecture	Type	Objective	Good For
BERT	Encoder-only	AE	MLM + NSP	Understanding
GPT	Decoder-only	AR	Predict next token	Text generation
BART	Encoder-Decoder	AE+AR	Denoising	Summarization, Translation
T5	Encoder-Decoder	AE+AR	Text-to-text	Universal NLP
ChatGPT	Decoder-only	AR + RLHF	Instruction-following	Conversations, reasoning tasks

Generative AI

13



What is Generative AI (GenAI)?

GenAI creates new content—text, images, code, music—based on patterns learned from data. Unlike traditional AI, it generates something new instead of analyzing existing data.

How Does GenAI Work?

GenAI models learn from vast datasets and generate new outputs using:

Neural Networks
Deep learning models for content creation.

Pattern Recognition
Predicts and generates sequences.

Pre-trained Models
Fine-tuned for specific tasks.

Applications of GenAI

Marketing
AI-generated ads & content.

Gaming
AI-generated characters & environments.

Healthcare
AI-powered drug discovery.

Education
AI tutors & content summarization.

Challenges & Ethical Concerns

Bias & Misinformation
Risk of false data.

Deepfakes & Manipulation
AI-created fake media.

Copyright Issues
Ownership of AI content.

Privacy Risks
AI models trained on personal data.

Types of Generative AI Models

Transformer Models (Text & Code)

Used for chatbots, writing, coding assistants.
Examples: GPT-4, Gemini, Llama, Claude

Diffusion Models (Images & Videos)

Turn random noise into realistic visuals.
Examples: Stable Diffusion, DALL-E, MidJourney

GANs (Realistic Media)

Two competing AIs create lifelike images & videos.
Examples: DeepFaceLab, This Person Does Not Exist

VAEs (Data Compression & Enhancement)

Used in creative design and data reconstruction.
Examples: AI-assisted image variations

RNN & LSTM (Sequential Data)

Used for music, speech, and handwriting generation.
Examples: AI-generated poetry, song compositions

GENERATIVE ADVERSARIAL NETWORKS (GAN)

GAN

GANs are generative models: they create new data instances that resemble your training data

For example, GANs can create images that look like photographs of human faces, even though the faces don't belong to any real person



Images generated by a GAN created by NVIDIA

GAN

GANs achieve this level of realism by pairing a generator and discriminator

Generator learns to produce the target output, with a discriminator, which learns to distinguish true data from the output of the generator

The generator tries to fool the discriminator, and the discriminator tries to keep from being fooled



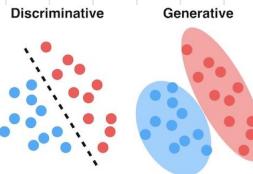
BACKGROUND

Generative models can generate new data instances

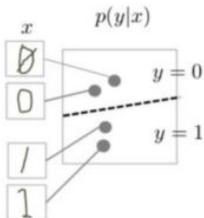
Discriminative models discriminate between different kinds of data instances

More formally, given a set of data instances X and a set of labels Y :

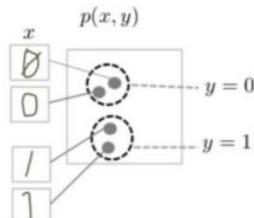
- Generative models capture the joint probability $p(X, Y)$, or just $p(X)$ if there are no labels
- Discriminative models capture the conditional probability $p(Y| X)$



• Discriminative Model



• Generative Model



OVERVIEW OF GAN STRUCTURE

A generative adversarial network (GAN) has two parts:

The generator learns to generate plausible data. The generated instances become negative training examples for the discriminator

The discriminator learns to distinguish the generator's fake data from real data. The discriminator penalizes the generator for producing implausible results

When training begins, the generator produces obviously fake data, and the discriminator quickly learns to tell that it's fake:



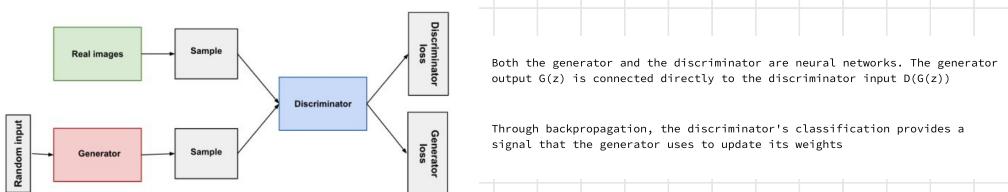
As training progresses, the generator gets closer to producing output that can fool the discriminator:



Finally, if generator training goes well, the discriminator gets worse at telling the difference between real and fake. It starts to classify fake data as real, and its accuracy decreases



GAN BLOCK DIAGRAM



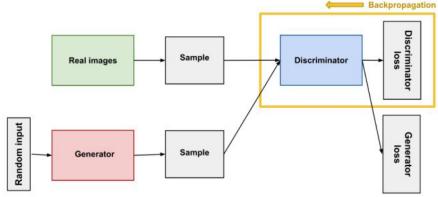
DISCRIMINATOR

The discriminator in a GAN is simply a classifier

It tries to distinguish real data from the data created by the generator. It could use any network architecture appropriate to the type of data it's classifying

The discriminator's training data comes from two sources:

- Real data instances, x , such as real pictures of people. The discriminator uses these instances as positive examples during training
- Fake data instances created by the generator, $G(z)$. The discriminator uses these instances as negative examples during training



DISCRIMINATOR TRAINING

The discriminator connects to two loss functions

During discriminator training, the discriminator ignores the generator loss and just uses the discriminator loss.

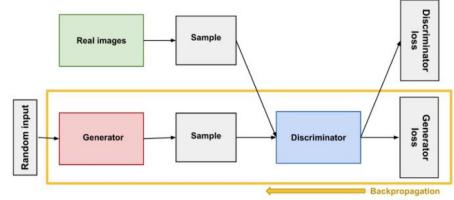
During discriminator training:

- The discriminator classifies both real data and fake data from the generator
- The discriminator loss penalizes the discriminator for misclassifying a real instance as fake or a fake instance as real
- The discriminator updates its weights through backpropagation from the discriminator loss through the discriminator network

GENERATOR

The generator part of a GAN learns to create fake data by incorporating feedback from the discriminator

It learns to make the discriminator classify its output as real



GENERATOR TRAINING

-
- A small image comparison shows a 'ORIGINAL IMAGE' of Spider-Man standing next to a police van, and a 'GAN GENERATED IMAGE' of a very similar-looking character standing in front of a door.
- Sample random noise
 - Produce generator output from sampled random noise
 - Get discriminator "Real" or "Fake" classification for generator output
 - Calculate loss from discriminator classification
 - Backpropagate through both the discriminator and generator to obtain gradients
 - Use gradients to change only the generator weights

GAN TRAINING

1. The discriminator trains for one or more epochs
2. The generator trains for one or more epochs
3. Repeat steps 1 and 2 to continue to train the generator and discriminator networks
4. Careful with Convergence as not to overtrain the generator otherwise generator will start training with less meaningful discriminator



LOSS FUNCTION

Original GAN: minimax Loss

$$E_x[\log(D(x))] + E_z[\log(1 - D(G(z)))]$$

- $D(x)$ is the discriminator's estimate of the probability that real data instance x is real
- E_x is the expected value over all real data instances
- $G(z)$ is the generator's output when given noise z
- $D(G(z))$ is the discriminator's estimate of the probability that a fake instance is real
- E_z is the expected value over all random inputs to the generator (in effect, the expected value over all generated fake instances $G(z)$)

The formula derives from the cross-entropy between the real and generated distributions

The generator can't directly affect the $\log(D(x))$ term in the function, so, for the generator, minimizing the loss is equivalent to minimizing $\log(1 - D(G(z)))$

Component	Original GAN	WGAN
Discriminator	Classifies real (1) vs fake (0)	Called a Critic — outputs a score (no sigmoid)
Output Range	$[0, 1]$ (probability)	\mathbb{R} (real numbers; no restriction)
Loss Type	Binary Cross-Entropy	Wasserstein (EM distance approximation)
Stability	Can be unstable (saturates, vanishes)	More stable gradients

MODIFIED MINIMAX LOSS

$$E_x[\log(D(x))] + E_z[\log(1 - D(G(z)))]$$

The original GAN paper notes that the above minimax loss function can cause the GAN to get stuck in the early stages of GAN training when the discriminator's job is very easy

The paper therefore suggests modifying the generator loss so that the generator tries to maximize $\log D(G(z))$

WASSERSTEIN LOSS

Here no classification for discriminator.

Discriminator becomes "Critic": Critic Loss: $D(x) - D(G(z))$

The discriminator tries to maximize this function

In other words, it tries to maximize the difference between its output on real instances and its output on fake instances

Generator Loss: $D(G(z))$

The generator tries to maximize this function. In other words, It tries to maximize the discriminator's output for its fake instances

ADVANTAGES OF WASSERSTEIN LOSS

Wasserstein GANs are less vulnerable to getting stuck than minimax based GANs and avoid problems with vanishing gradients

The earth mover distance also has the advantage of being a true metric: a measure of distance in a space of probability distributions

Cross-entropy is not a metric in this sense.

COMMON PROBLEMS IN GAN

Problem #1: Vanishing Gradient: Research has suggested that if your discriminator is too good, then generator training can fail due to vanishing gradients

In effect, an optimal discriminator doesn't provide enough information for the generator to make progress

Sol: Wasserstein loss and Modified minimax loss

Problem #2: Mode Collapse: When the generators rotate through a small set of output types, this form of GAN failure is called Mode Collapse

In effect, an optimal discriminator doesn't provide enough information for the generator to make progress

Sol: Wasserstein loss and Modified minimax loss

Problem #3: Failure To Converge: Often problem in GAN

Sol:

- (i) Adding noise to discriminator inputs
- (ii) Penalizing discriminator weights



AUTO Encoders

M/

UNSUPERVISED LEARNING

"We expect unsupervised learning to become far more important in the longer term. Human and animal learning is largely unsupervised: we discover the structure of the world by observing it, not by being told the name of every object." - LeCun, Bengio, Hinton, Nature 2015

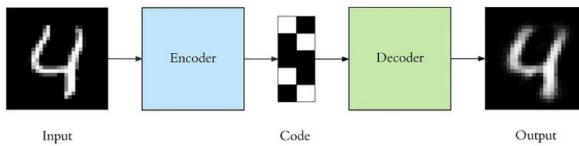
If intelligence was a cake, unsupervised learning would be the cake, supervised learning would be the icing on the cake, and reinforcement learning would be the cherry on the cake. We know how to make the icing and the cherry, but we don't know how to make the cake. - Yann LeCun, March 14, 2016 (Facebook)

AUTOENCODERS

Autoencoders are a specific type of feedforward neural networks where the input is the same as the output.

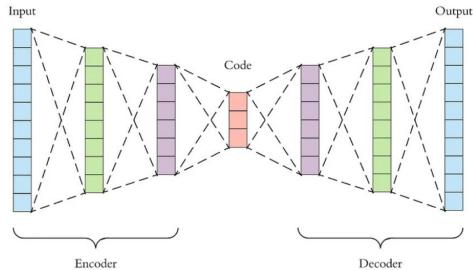
An autoencoder consists of 3 components: encoder, code and decoder.

Autoencoders are mainly a dimensionality reduction (or compression) algorithm.

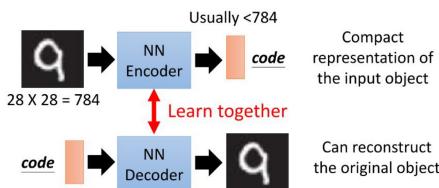


UNSUPERVISED????

Autoencoders are considered an unsupervised learning technique since they don't need explicit labels to train on. But to be more precise they are self-supervised because they generate their own labels from the training data.



AUTOENCODERS



There are 4 hyperparameters that we need to set before training an autoencoder:

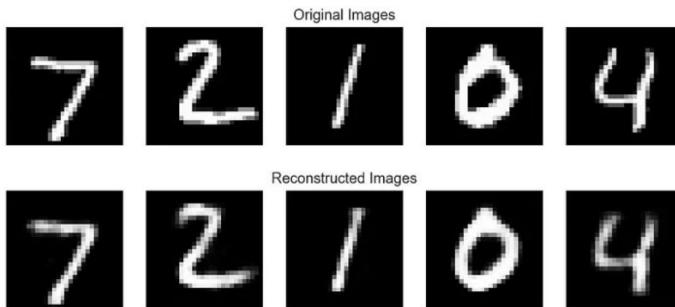
1. Code size
2. Number of layers
3. Number of nodes per layer
4. Loss function

IMPLEMENTATION

```
input_size = 784
hidden_size = 128
code_size = 32

input_img = Input(shape=(input_size,))
hidden_1 = Dense(hidden_size, activation='relu')(input_img)
code = Dense(code_size, activation='relu')(hidden_1)
hidden_2 = Dense(hidden_size, activation='relu')(code)
output_img = Dense(input_size, activation='sigmoid')(hidden_2)

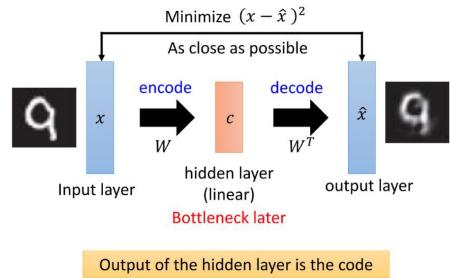
autoencoder = Model(input_img, output_img)
autoencoder.compile(optimizer='adam', loss='binary_crossentropy')
autoencoder.fit(x_train, x_train, epochs=5)
```



A BIT OF ADVICE

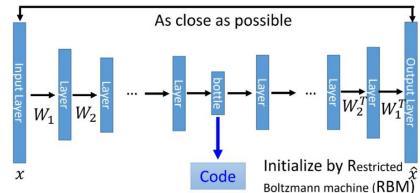
We should be careful to not make it too powerful. Otherwise the autoencoder will simply learn to copy its inputs to the output, without learning any meaningful representation. It will just mimic the identity function. The autoencoder will reconstruct the training data perfectly, but it will be overfitting without being able to generalize to new instances, which is not what we want.

Recap: PCA



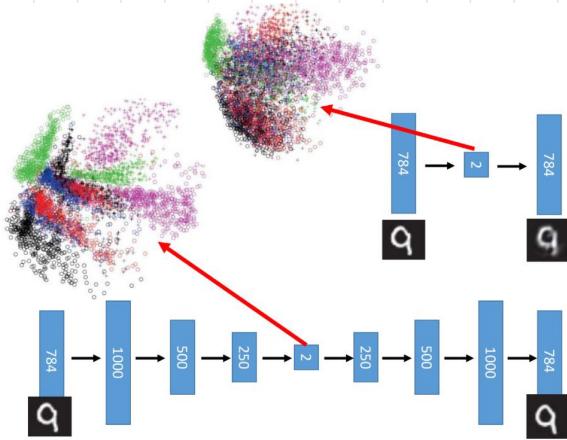
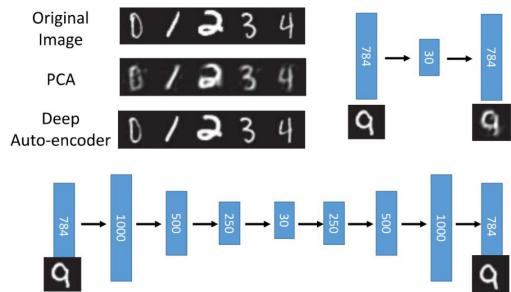
Deep Auto-encoder

- Of course, the auto-encoder can be deep



Reference: Hinton, Geoffrey E., and Ruslan R. Salakhutdinov. "Reducing the dimensionality of data with neural networks." *Science* 313.5786 (2006): 504-507

Deep Auto-encoder

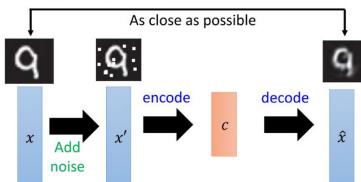


Auto-encoder

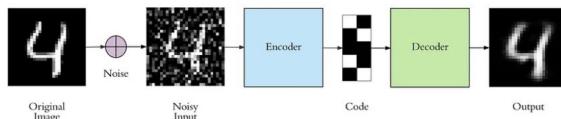
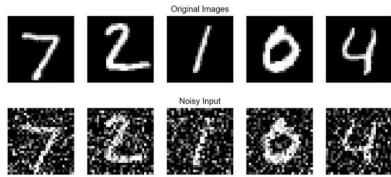
More: Contractive auto-encoder

Ref: Rifai, Salah, et al. "Contractive auto-encoders: Explicit invariance during feature extraction." *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*. 2011.

- De-noising auto-encoder



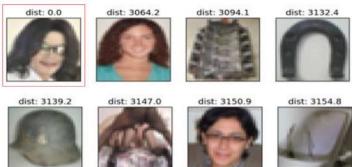
Vincent, Pascal, et al. "Extracting and composing robust features with denoising autoencoders." *ICML*, 2008.



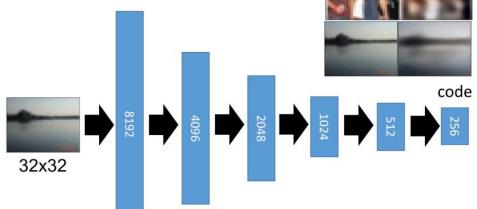
Denoising autoencoder is trained as:
`autoencoder.fit(x_train_noisy, x_train)`

Auto-encoder – Similar Image Search

Retrieved using Euclidean distance in pixel intensity space

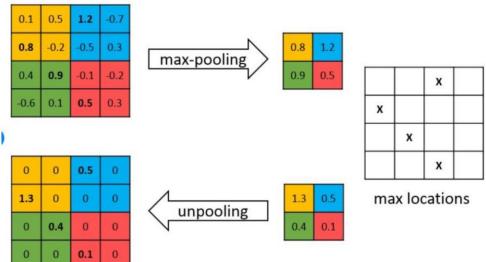
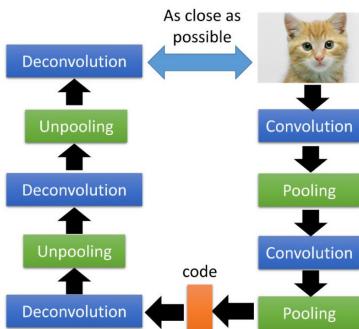


Auto-encoder – Similar Image Search



(crawl millions of images from the Internet)

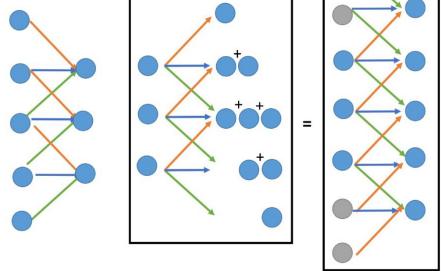
Auto-encoder for CNN



Pooling and unpooling layers. For each pooling layer, the max locations are stored. These locations are then used in unpooling layer.

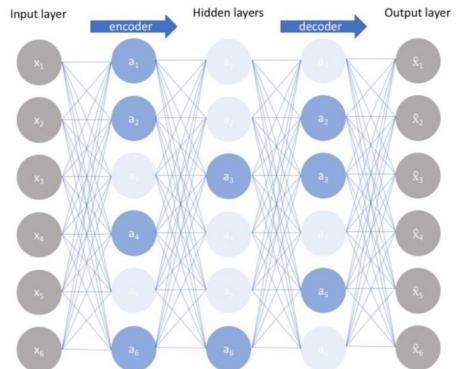
CNN

- Deconvolution



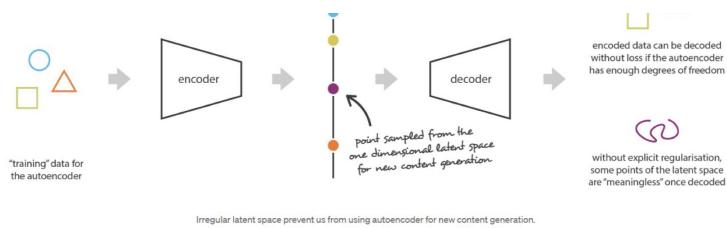
Sparse AutoEncoder

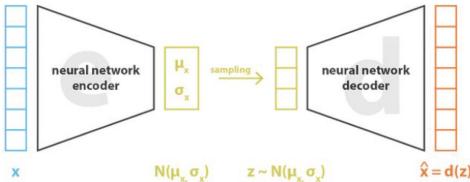
- Use of Regularization (For better generalization)
- Introducing sparsity constraint such that only a fraction of the nodes would have nonzero values called active nodes
- Add a penalty term to the loss function such that only a fraction of nodes become active
- This forces the autoencoder to represent each input as a combination of small number of nodes and demands to discover interesting structure in the data



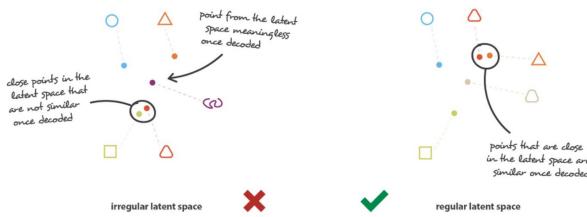
VARIATIONAL AUTOENCODER (VAE)

A variational autoencoder can be defined as being an autoencoder whose training is regularised to avoid overfitting and ensure that the latent space has good properties that enable generative process.





$$\text{loss} = \|x - \hat{x}\|^2 + \text{KL}[N(\mu_z, \sigma_z), N(0, I)] = \|x - d(z)\|^2 + \text{KL}[N(\mu_z, \sigma_z), N(0, I)]$$



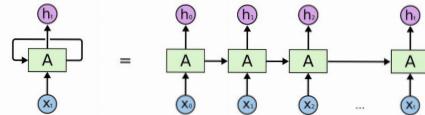
The returned distributions of VAEs have to be regularised to obtain a latent space with good properties.

RNN NUMERICAL

Sentence:
"I like to eat pizza."

Vocabulary: POS Tags
 "I" - Pronoun
 "like" - Verb
 "to" - Preposition
 "eat" - Verb
 "pizza" - Noun

$$\begin{aligned} \mathbf{a}^{(t)} &= b + W\mathbf{h}^{(t-1)} + U\mathbf{x}^{(t)} \\ \mathbf{h}^{(t)} &= \tanh(\mathbf{a}^{(t)}) \\ \mathbf{o}^{(t)} &= c + V\mathbf{h}^{(t)} \\ \hat{\mathbf{y}}^{(t)} &= \text{softmax}(\mathbf{o}^{(t)}) \end{aligned}$$



Θ)

$W_{xh} = [[0.4, -0.3, 0.1, -0.2, 0.5],$	Word "I":
$[0.1, -0.2, 0.3, 0.2, -0.4],$	$x_t = [1, 0, 0, 0, 0]$ # One-hot encoding for "I" in the vocabulary
$[0.2, -0.1, 0.5, 0.4, -0.3]]$	$h_t = \tanh(W_{xh} * x_t + W_{hh} * h_{t-1})$
$W_{hy} = [[0.2, 0.6, -0.1],$	$h_t = \tanh([0.4, -0.3, 0.1, -0.2, 0.5],$
$[0.3, -0.2, 0.4],$	$[0.1, -0.2, 0.3, 0.2, -0.4],$
$[-0.4, 0.1, 0.5],$	$[0.2, -0.1, 0.5, 0.4, -0.3]] * [1, 0, 0, 0, 0]^T // 3x5 * 5x1$
$[0.1, 0.2, 0.3]]$	$+ [[0.2, -0.1, 0.3],$
$W_{hh} = [[0.2, -0.1, 0.3],$	$[-0.1, 0.4, -0.2],$
$[-0.1, 0.4, -0.2],$	$h_t = \tanh([0.49, 0.01, 0.37])$
$[0.4, -0.3, 0.5]]$	$h_t = [0.4621, 0.0099, 0.3584]$
$h_0 = [0.1, -0.1, 0.2]$	\checkmark

Tanh

$$f(x) = \frac{(e^x - e^{-x})}{(e^x + e^{-x})}$$

$$\begin{aligned} \text{pos_scores} &= W_{hy} * h_t \\ \text{pos_scores} &= [[0.2, 0.6, -0.1], \\ &\quad [0.3, -0.2, 0.4], \\ &\quad [-0.4, 0.1, 0.5], \\ &\quad [0.1, 0.2, 0.3]] * [0.4621, 0.0099, 0.3584] \end{aligned}$$

$$\text{pos_scores} = [0.06252, 0.28001, 0.19535, 0.15571] \quad \checkmark$$

After softmax [0.2229, 0.2777, 0.2547, 0.2446]
 Pronoun verb preposition noun

$$\begin{aligned} \text{Cross-Entropy Loss (L)} &= -\sum(T_{-i} * \log(P_{-i})) \\ L &= -(1 * \log(0.2229) + 0 * \log(0.2777) + 0 * \log(0.2547) + 0 * \log(0.2446)) \end{aligned}$$

$$L = -(\log(0.2229))$$

$$L \approx 1.5020$$

RNN NUMERICAL

Sentence: "I like to eat pizza."

Vocabulary: POS Tags
 "I" - Pronoun
 "like" - Verb
 "to" - Preposition
 "eat" - Verb
 "pizza" - Noun

$$\begin{aligned}
 a^{(t)} &= b + W_{hh}^{(t-1)} + W_{xh}^{(t)} \quad \text{hidden to hidden weights} \\
 h^{(t)} &= \tanh(a^{(t)}) \\
 o^{(t)} &= c + Vh^{(t)} \quad \text{hidden to output weights} \\
 y^{(t)} &= \text{softmax}(o^{(t)}) \quad \text{output bias}
 \end{aligned}$$

1. One hot vector encode

$$\begin{array}{l}
 \text{I like to eat pizza} \\
 \text{I: } \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \end{bmatrix} \\
 \text{like: } \begin{bmatrix} 0 & 1 & 0 & 0 & 0 \end{bmatrix} \\
 \text{to: } \begin{bmatrix} 0 & 0 & 1 & 0 & 0 \end{bmatrix} \\
 \text{eat: } \begin{bmatrix} 0 & 0 & 0 & 1 & 0 \end{bmatrix} \\
 \text{Pizza: } \begin{bmatrix} 0 & 0 & 0 & 0 & 1 \end{bmatrix}
 \end{array}$$

2. forward pass for each word

$$\begin{aligned}
 x^1 &\rightarrow \text{For I: } I = [1 \ 0 \ 0 \ 0 \ 0]^T, h_0 = [0.1, -0.1, 0.2]^T \\
 &\rightarrow a^1 = W_{hh} \cdot h^0 + W_{xh} \cdot x^1
 \end{aligned}$$

$$\begin{bmatrix} 0.2 & -0.1 & 0.3 \\ -0.1 & 0.4 & -0.2 \\ 0.4 & -0.3 & 0.5 \end{bmatrix} \begin{bmatrix} 0.1 \\ -0.1 \\ 0.2 \end{bmatrix} + \begin{bmatrix} 0.4 \\ 0.1 \\ -0.2 \end{bmatrix} \xrightarrow{\text{as } I = [1 \ 0 \ 0 \ 0 \ 0]} \begin{bmatrix} 0.49 \\ 0.01 \\ 0.37 \end{bmatrix}$$

select first column

$$\rightarrow h^1 = \tanh(a^1)$$

$$h^1 = \begin{bmatrix} \tanh(0.49) \\ \tanh(0.01) \\ \tanh(0.37) \end{bmatrix} = \begin{bmatrix} 0.454 \\ 0.01 \\ 0.354 \end{bmatrix}$$

$$\text{pos scores} \rightarrow o^1 = W_{hy} \cdot h^1$$

$$\begin{bmatrix} 0.2 & 0.6 & -0.1 \\ 0.3 & -0.2 & 0.4 \\ -0.4 & 0.1 & 0.5 \\ 0.1 & 0.2 & 0.3 \end{bmatrix} \begin{bmatrix} 0.454 \\ 0.01 \\ 0.354 \end{bmatrix} = \begin{bmatrix} 0.0614 \\ 0.2758 \\ -0.0036 \\ 0.1536 \end{bmatrix}$$

$$\begin{aligned}
 &\rightarrow y^1 = \frac{e^{o_j}}{\sum_i e^{o_i}} \xrightarrow{\text{softmax}} \\
 &e^{0.0614} = 1.063 &= 1.063 / 4.542 = 0.234 \\
 &e^{0.2758} = 1.317 &= 1.317 / 4.542 = 0.29 \\
 &e^{-0.0036} = 0.996 &= \dots = 0.219 \\
 &e^{0.1536} = 1.166 &= \dots = 0.257 \\
 &&4.542
 \end{aligned}$$

$$\hookrightarrow \text{Cross-Entropy Loss: } -\sum (T_i \cdot \log(p_i))$$

$$= -(1 \cdot \log(0.234) + 0 \cdot \log(0.29) + 0 \cdot \log(0.219) + 0 \cdot \log(0.257))$$

$$= 1.5020$$

$$V_h \leftarrow W_{xh} = [[0.4, -0.3, 0.1, -0.2, 0.5], [0.1, -0.2, 0.3, 0.2, -0.4], [0.2, -0.1, 0.5, 0.4, -0.3]]$$

$$V_h \leftarrow W_{hy} = [[0.2, 0.6, -0.1], [0.3, -0.2, 0.4], [-0.4, 0.1, 0.5], [0.1, 0.2, 0.3]]$$

$$W_h \leftarrow W_{hh} = [[0.2, -0.1, 0.3], [-0.1, 0.4, -0.2], [0.4, -0.3, 0.5]], h_0 = [0.1, -0.1, 0.2]$$

Steps

1 One hot encode

2 Forward Pass of x^1

$$\hookrightarrow a^1 = W_{hh} \cdot h^0 + W_{xh} \cdot x^1$$

$$\hookrightarrow h^1 = \tanh(a^1)$$

$$\hookrightarrow o^1 = W_{hy} \cdot h^1$$

$$\hookrightarrow y^1 = \text{softmax}(o^1) \cdot \frac{e^o}{\sum_i e^o}$$

$$\hookrightarrow \text{CELoss} = -\sum (x^i \cdot \log(y^i))$$

Transformer Numerical

Input sentence: "Transformers transforming our lives"
Output sentence: "For sure"

The embedding matrix for the words is represented as follows:

"Transformers": [0.1, 0.2, 0.3]
"transforming": [0.4, 0.5, 0.6]
"our": [0.7, 0.8, 0.9]
"lives": [1.0, 1.1, 1.2]
"For": [0.4, 0.1, 0.8]
"sure": [0.9, 0.7, 0.2]

Assuming the weight matrices:

$$W_q = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 2 \\ 7 & 1 & 9 \\ 1 & 6 & 9 \end{bmatrix}$$

$$W_k = \begin{bmatrix} 7 & 3 & 1 \\ 9 & 2 & 1 \\ 2 & 4 & 6 \end{bmatrix}$$

$$W_v = \begin{bmatrix} 8 & 0 & 2 \\ 1 & 6 & 8 \end{bmatrix}$$

$$X_i = \begin{bmatrix} 0.1 & 0.2 & 0.3 \\ 0.4 & 0.5 & 0.6 \\ 0.7 & 0.8 & 0.9 \\ 1.0 & 1.1 & 1.2 \end{bmatrix}$$

Transformers
transform
our
lives

For simplicity, assume dk is 1.

- Calculate masked self-attention for the above input sentence using the provided weight matrices. Show all steps clearly.
- Calculate cross-attention for the above input and output sentence using the provided weight matrices. Show all steps clearly.

Masked Self Attention

1. COMPUTE Q, K, V

$$\rightarrow Q = X_i W^Q$$

$$\begin{bmatrix} 0.1 & 0.2 & 0.3 \\ 0.4 & 0.5 & 0.6 \\ 0.7 & 0.8 & 0.9 \\ 1.0 & 1.1 & 1.2 \end{bmatrix} \cdot \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 2 \\ 7 & 1 & 9 \end{bmatrix} = \begin{bmatrix} 3 & 15 & 34 \\ 66 & 3.9 & 7.6 \\ 10.2 & 6.3 & 11.8 \\ 13.8 & 8.7 & 16 \end{bmatrix}$$

$$\rightarrow K = X_i W^K$$

$$\begin{bmatrix} 0.1 & 0.2 & 0.3 \\ 0.4 & 0.5 & 0.6 \\ 0.7 & 0.8 & 0.9 \\ 1.0 & 1.1 & 1.2 \end{bmatrix} \cdot \begin{bmatrix} 1 & 6 & 9 \\ 7 & 3 & 1 \\ 9 & 2 & 1 \end{bmatrix} = \begin{bmatrix} 4.2 & 1.8 & 1.4 \\ 9.3 & 5.1 & 4.1 \\ 14.4 & 8.4 & 8 \\ 19.5 & 11.1 & 11.3 \end{bmatrix}$$

$$\rightarrow V = X_i W^V$$

$$\begin{bmatrix} 0.1 & 0.2 & 0.3 \\ 0.4 & 0.5 & 0.6 \\ 0.7 & 0.8 & 0.9 \\ 1.0 & 1.1 & 1.2 \end{bmatrix} \cdot \begin{bmatrix} 2 & 4 & 6 \\ 8 & 0 & 2 \\ 1 & 6 & 8 \end{bmatrix} = \begin{bmatrix} 2.1 & 2.2 & 3.4 \\ 5.4 & 5.2 & 8.2 \\ 8.7 & 8.2 & 13 \\ 12 & 11.2 & 17.8 \end{bmatrix}$$

2. Compute attention score

$$E = Q K^T / \sqrt{dk}$$

It becomes 1 so ignore for now

$$\begin{bmatrix} 3 & 15 & 34 \\ 6.6 & 3.9 & 7.6 \\ 10.2 & 6.3 & 11.8 \\ 13.8 & 8.7 & 16 \end{bmatrix} \cdot \begin{bmatrix} 4.2 & 9.3 & 14.4 & 19.5 \\ 1.8 & 5.1 & 8.4 & 11.7 \\ 1.4 & 4.7 & 8 & 11.3 \end{bmatrix} = \begin{bmatrix} 20.06 & 5153 & 83 & 114.47 \\ 45.98 & 116.99 & 188.6 & 260.21 \\ 70.7 & 182.45 & 294.12 & 405.93 \\ 96.02 & 247.91 & 394.98 & 551.69 \end{bmatrix}$$

After applying masking

$$E = \begin{bmatrix} 20.06 & -\infty & -\infty & -\infty \\ 45.98 & 116.99 & -\infty & -\infty \\ 70.7 & 182.45 & 294.12 & -\infty \\ 96.02 & 247.91 & 394.98 & 551.69 \end{bmatrix}$$

3. Softmax E

Row wise softmax

$$\begin{bmatrix} 20.06 & -\infty & -\infty & -\infty \\ 45.98 & 116.99 & -\infty & -\infty \\ 70.7 & 182.45 & 294.12 & -\infty \\ 96.02 & 247.91 & 394.98 & 551.69 \end{bmatrix} = \begin{bmatrix} e^{20.06} \\ e^{45.98} \\ e^{70.7} \\ e^{96.02} \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$A = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 7.9 \times 10^{-99} & 1 & 0 & 0 \\ 8.6 \times 10^{-99} & 2.9 \times 10^{-49} & 1 & 0 \\ 1.21 \times 10^{-99} & 1.17 \times 10^{-132} & 1.08 \times 10^{-66} & 1 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

4. Output = A . V

$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 2.1 & 2.2 & 3.4 \\ 5.4 & 5.2 & 8.2 \\ 8.7 & 8.2 & 13 \\ 12 & 11.2 & 17.8 \end{bmatrix} = \begin{bmatrix} 2.1 & 2.2 & 3.4 \\ 5.4 & 5.2 & 8.2 \\ 8.7 & 8.2 & 13 \\ 12 & 11.2 & 17.8 \end{bmatrix}$$

Transformers
transforming
our
lives

Steps

- Compute Q, K, V with X_{input}
- Attention Score
 $E = Q K^T / \sqrt{dk}$
- Apply mask
- Softmax row wise $\rightarrow A$
- Output: $A . V$

Input sentence: "Transformers transforming our lives"

Output sentence: "For sure"

The embedding matrix for the words is represented as follows:

"Transformers": [0.1, 0.2, 0.3]

"transforming": [0.4, 0.5, 0.6]

"our": [0.7, 0.8, 0.9]

"lives": [1.0, 1.1, 1.2]

"For": [0.4, 0.1, 0.8]

"sure": [0.9, 0.7, 0.2]

Assuming the weight matrices:

$$W_q = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 2 \\ 7 & 1 & 9 \\ 1 & 6 & 9 \\ 7 & 3 & 1 \\ 9 & 2 & 1 \\ 2 & 4 & 6 \\ 8 & 0 & 2 \\ 1 & 6 & 8 \end{bmatrix}$$

$$X_o = \begin{bmatrix} 0.4 & 0.1 & 0.8 \\ 0.9 & 0.7 & 0.2 \end{bmatrix}$$

For
Sure

For simplicity, assume dk is 1.

- Calculate masked self-attention for the above input sentence using the provided weight matrices. Show all steps clearly.
- Calculate cross-attention for the above input and output sentence using the provided weight matrices. Show all steps clearly.

Cross-attention

1. Compute Q, K, V

$\rightarrow Q = X_o W_q \xrightarrow{\text{Using Output Sequence}}$

$$\begin{bmatrix} 0.4 & 0.1 & 0.8 \\ 0.9 & 0.7 & 0.2 \end{bmatrix} \cdot \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 2 \\ 7 & 1 & 9 \end{bmatrix} = \begin{bmatrix} 6.4 & 2.1 & 8.6 \\ 5.1 & 5.5 & 5.9 \end{bmatrix}$$

$\rightarrow K = X_i W_k \xrightarrow{\text{Using Input Sequence}}$

$$\begin{bmatrix} 0.1 & 0.2 & 0.3 \\ 0.4 & 0.5 & 0.6 \\ 0.7 & 0.8 & 0.9 \\ 1.0 & 1.1 & 1.2 \end{bmatrix} \cdot \begin{bmatrix} 1 & 6 & 9 \\ 7 & 3 & 1 \\ 9 & 2 & 1 \end{bmatrix} = \begin{bmatrix} 4.2 & 1.8 & 1.4 \\ 9.3 & 5.1 & 4.7 \\ 14.4 & 8.4 & 8 \\ 19.5 & 11.7 & 11.3 \end{bmatrix}$$

$\rightarrow V = X_i W_v \xrightarrow{\text{Using Input Sequence}}$

$$\begin{bmatrix} 0.1 & 0.2 & 0.3 \\ 0.4 & 0.5 & 0.6 \\ 0.7 & 0.8 & 0.9 \\ 1.0 & 1.1 & 1.2 \end{bmatrix} \cdot \begin{bmatrix} 2 & 4 & 6 \\ 8 & 0 & 2 \\ 1 & 6 & 8 \end{bmatrix} = \begin{bmatrix} 2.1 & 2.2 & 3.4 \\ 5.4 & 5.2 & 8.2 \\ 8.7 & 8.2 & 13 \\ 12 & 11.2 & 17.8 \end{bmatrix}$$

2. Attention Score

$$E = Q^T / \sqrt{dk}$$

$$\begin{bmatrix} 6.4 & 2.1 & 8.6 \\ 5.1 & 5.5 & 5.9 \end{bmatrix} \cdot \begin{bmatrix} 4.2 & 1.8 & 1.4 \\ 9.3 & 5.1 & 4.7 \\ 14.4 & 8.4 & 8 \\ 19.5 & 11.7 & 11.3 \end{bmatrix} = \begin{bmatrix} 427 & 110.65 & 178.6 & 246.5 \\ 3958 & 103.21 & 166.8 & 230.5 \end{bmatrix}$$

3. Softmax E

Row wise softmax

$$\begin{bmatrix} 427 & 110.65 & 178.6 & 246.5 \\ 3958 & 103.21 & 166.8 & 230.5 \end{bmatrix} = \begin{bmatrix} \frac{e^{427}}{e^{427} + e^{110.65} + e^{178.6} + e^{246.5}} \\ \frac{e^{110.65}}{e^{427} + e^{110.65} + e^{178.6} + e^{246.5}} \\ \frac{e^{178.6}}{e^{427} + e^{110.65} + e^{178.6} + e^{246.5}} \\ \frac{e^{246.5}}{e^{427} + e^{110.65} + e^{178.6} + e^{246.5}} \end{bmatrix}$$

$$A = \begin{bmatrix} 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

4. Output: A.V

$$\begin{bmatrix} 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 2.1 & 2.2 & 3.4 \\ 5.4 & 5.2 & 8.2 \\ 8.7 & 8.2 & 13 \\ 12 & 11.2 & 17.8 \end{bmatrix} = \begin{bmatrix} 12 & 11.2 & 17.8 \\ 12 & 11.2 & 17.8 \end{bmatrix}$$

For
Sure

Steps

1. Compute Q, K, V with X

$\hookrightarrow Q$ with X_{output}

$\hookrightarrow K, V$ with X_{input}

2. Attention Score $E = Q^T / \sqrt{dk}$

3. Softmax row wise $\rightarrow A$

4. Output: A.V

SELF ATTENTION NUMERICAL

Step 1

$$Q = XW^T$$

$$\begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 2 & 0 & 2 \\ 1 & 1 & 1 & 1 \end{bmatrix} \times \begin{bmatrix} 1 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 1 \end{bmatrix}$$

$$3 \times 4 \quad 4 \times 3$$

$$Q = \begin{bmatrix} 1 & 0 & 2 \\ 2 & 2 & 2 \\ 2 & 1 & 3 \end{bmatrix}$$

$$K = XW^T$$

$$K = \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 2 & 0 & 2 \\ 1 & 1 & 1 & 1 \end{bmatrix} \times \begin{bmatrix} 0 & 0 & 1 \\ 1 & 1 & 0 \\ 0 & 1 & 0 \\ 1 & 1 & 0 \end{bmatrix}$$

$$K = \begin{bmatrix} 0 & 1 & 1 \\ 4 & 4 & 0 \\ 2 & 3 & 1 \end{bmatrix}$$

$$V = XW^T$$

$$V = \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 2 & 0 & 2 \\ 1 & 1 & 1 & 1 \end{bmatrix} \times \begin{bmatrix} 0 & 2 & 0 \\ 0 & 3 & 0 \\ 1 & 0 & 3 \\ 1 & 1 & 0 \end{bmatrix}$$

$$V = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 0 & 8 & 0 \\ 2 & 6 & 3 \end{bmatrix}$$

Step 2

$$E = QKT$$

$$E = \begin{bmatrix} 1 & 0 & 2 \\ 2 & 2 & 2 \\ 2 & 1 & 3 \end{bmatrix} \times \begin{bmatrix} 0 & 4 & 2 \\ 1 & 4 & 0 \\ 1 & 0 & 0 \end{bmatrix}$$

$$E = \begin{bmatrix} 2 & 4 & 4 \\ 4 & 16 & 18 \\ 4 & 12 & 10 \end{bmatrix}$$

x = [

- [1, 0, 1, 0], # Input 1
- [0, 2, 0, 2], # Input 2
- [1, 1, 1, 1] # Input 3

]

```
w_key = [
    [0, 0, 1],
    [1, 1, 0],
    [0, 1, 0],
    [1, 1, 0]
]

w_query = [
    [1, 0, 1],
    [1, 0, 0],
    [0, 0, 1],
    [0, 1, 1]
]

w_value = [
    [0, 2, 0],
    [0, 3, 0],
    [1, 0, 3],
    [1, 1, 0]
]
```

Solve the above example. For simplicity dk=1.

Round off answer =

[2,	7.	1.5]
[2,	8.	0.	1
[2,	7.7999997	0.3]

Step 3

$$A = \text{softmax}(E)$$

$$\Rightarrow \frac{e^4}{e^4 + e^4 + e^4} \quad \frac{e^4}{e^4 + e^4 + e^4} \quad \frac{e^4}{e^4 + e^4 + e^4}$$

$$= [0.06, 0.468, 0.468] \rightarrow \text{Input 1}$$

$$\Rightarrow \frac{e^4}{e^4 + e^{16} + e^{18}} \quad \frac{e^{16}}{e^4 + e^{16} + e^{18}} \quad \frac{e^{18}}{e^4 + e^{16} + e^{18}}$$

$$= [0.000007, 0.117, 0.88] \rightarrow \text{Input 2}$$

$$\therefore$$

$$\Rightarrow \frac{e^4}{e^4 + e^{18} + e^{10}} \quad \frac{e^{12}}{e^4 + e^{18} + e^{10}} \quad \frac{e^{10}}{e^4 + e^{18} + e^{10}}$$

$$= [0.00029, 0.88, 0.117] \rightarrow \text{Input 3}$$

$$A = \begin{bmatrix} 0.06 & 0.468 & 0.468 \\ 0.000007 & 0.117 & 0.88 \\ 0.00029 & 0.88 & 0.117 \end{bmatrix}$$

Step 4

$$\text{Output} = A \times V$$

$$\text{Output} = \begin{bmatrix} 0.06 & 0.468 & 0.468 \\ 0.000007 & 0.117 & 0.88 \\ 0.00029 & 0.88 & 0.117 \end{bmatrix} \times \begin{bmatrix} 1 & 2 & 3 \\ 2 & 8 & 0 \\ 2 & 6 & 3 \end{bmatrix}$$

$$\text{Input} = \begin{bmatrix} 0.932 & 6.67 & 1.58 \\ 1.99 & 6.23 & 2.64 \\ 1.099 & 7.75 & 0.357 \end{bmatrix}$$

$$\text{Output} = \begin{bmatrix} 1.5 \\ 2.7 \\ 1.5 \end{bmatrix}$$

$$\text{X} \quad \text{X} \quad \text{X}$$

STEPS

1. Compute Q, K, V with X
2. Attention Score $E \cdot QK^T / \sqrt{dk}$
3. Softmax row wise $\rightarrow A$
4. Output $A \cdot V$

Attention Numerical

Encoder

$$\text{Q) } h_1 = [1, 2, 3] \\ h_2 = [3, 1, 1] \\ h_3 = [4, 2, 1]$$

Decoder

$$S_1 = [3, 4, 5]$$

↳ Find attention scores e^t

$$S_1 \cdot h_1 = \begin{bmatrix} 3 \\ 4 \\ 5 \end{bmatrix} [1, 2, 3] \rightarrow 3+8+15=26$$

$$S_1 \cdot h_2 = \begin{bmatrix} 3 \\ 4 \\ 5 \end{bmatrix} [3, 1, 1] \rightarrow 9+4+5=18$$

$$S_1^T \cdot h_3 = \begin{bmatrix} 3 \\ 4 \\ 5 \end{bmatrix} [4, 2, 1] \rightarrow 12+8+5=25$$

$$e^t = [26, 18, 25]$$

$$\hookrightarrow \text{softmax} = \left[\frac{e^{26}}{e^{26} + e^{18} + e^{25}}, \frac{e^{18}}{e^{26} + e^{18} + e^{25}}, \frac{e^{25}}{e^{26} + e^{18} + e^{25}} \right]$$

$$\alpha = [0.73, 0.0002, 0.2688]$$

$$\begin{bmatrix} 1 & 2 & 3 \\ 3 & 1 & 1 \\ 4 & 2 & 1 \end{bmatrix} \begin{bmatrix} 3 \\ 4 \\ 5 \end{bmatrix}$$

↳ Content vector

$$a_t = \alpha \cdot h^t$$

$$= [0.73] \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} + [0.0002] \begin{bmatrix} 3 \\ 1 \\ 1 \end{bmatrix} + [0.2688] \begin{bmatrix} 4 \\ 2 \\ 1 \end{bmatrix} = [1.8058, 1.9978, 2.459]$$

Steps

1. Attention score

$$e_i^t = S_1^T \cdot h_i \quad \begin{matrix} h \text{ encoder} \\ s \text{ decoder} \end{matrix}$$

2. softmax(e^t)

$$\alpha_i = \frac{e^{e_i^t}}{\sum_j e^{e_j^t}}$$

3. Context Vector

$$a_t = \sum_{i=1}^n \alpha_i \cdot h_i$$

$$\begin{aligned} \text{Input activation} & \rightarrow a_0 = \tanh(W_a \cdot x_0 + U_a \cdot \text{out}_{-1} + b_a) \\ \text{Input} & \rightarrow i_0 = \sigma(W_i \cdot x_0 + U_i \cdot \text{out}_{-1} + b_i) \\ \text{gate forget} & \rightarrow f_0 = \sigma(W_f \cdot x_0 + U_f \cdot \text{out}_{-1} + b_f) \\ \text{Output} & \rightarrow o_0 = \sigma(W_o \cdot x_0 + U_o \cdot \text{out}_{-1} + b_o) \\ \text{internal state} & \rightarrow \text{state}_0 = (a_0 \times i_0) + (f_0 \times \text{state}_{-1}) \\ \text{Output} & \rightarrow \text{out}_0 = \tanh(\text{state}_0) \times o_0 \end{aligned}$$

↑ no previous value
so 0

$$\delta_{\text{out}_i} = \Delta_i + \Delta_{\text{out}_i}$$

$$\delta_{\text{state}_i} = \delta_{\text{out}_i} \odot o_i \odot (1 - \tanh^2(\text{state}_i)) + \delta_{\text{state}_{i+1}} \odot f_{i+1}$$

$$\delta_{a_i} = \delta_{\text{state}_i} \odot a_i \odot i_i \odot (1 - i_i)$$

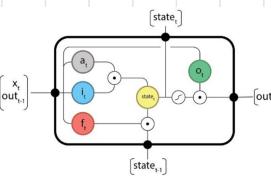
$$\delta_{i_i} = \delta_{\text{state}_i} \odot a_i \odot i_i \odot f_i \odot (1 - f_i)$$

$$\delta_{o_i} = \delta_{\text{state}_i} \odot \tanh(\text{state}_i) \odot o_i \odot (1 - o_i)$$

$$\delta_{u_i} = W^T \cdot \delta_{\text{gates}_i}$$

$$\Delta_{\text{out}_{i-1}} = U^T \cdot \delta_{\text{gates}_i}$$

LSTM Numerical



- Above \odot is the element-wise product or Hadamard product.
- Inner products will be represented as \cdot
- Outer products will be represented as \otimes
- σ represents the sigmoid function

Lets consider sequence length of two to demonstrate the unrolling over time

(a) Show that for t=0; following were the values of a_0 , i_0 , f_0 , state_0 , and out_0

① Internal Weights

$$\begin{aligned} W_a &= \begin{bmatrix} 0.05 \\ 0.25 \end{bmatrix}, W_i = \begin{bmatrix} 0.95 \\ 0.8 \end{bmatrix}, W_f = \begin{bmatrix} 0.1 \\ 0.45 \end{bmatrix}, W_o = \begin{bmatrix} 0.6 \\ 0.4 \end{bmatrix} \\ U_a &= \begin{bmatrix} 0.15 \end{bmatrix}, U_i = \begin{bmatrix} 0.8 \end{bmatrix}, U_f = \begin{bmatrix} 0.1 \end{bmatrix}, U_o = \begin{bmatrix} 0.25 \end{bmatrix} \\ b_a &= \begin{bmatrix} 0.2 \end{bmatrix}, b_i = \begin{bmatrix} 0.65 \end{bmatrix}, b_f = \begin{bmatrix} 0.15 \end{bmatrix}, b_o = \begin{bmatrix} 0.1 \end{bmatrix} \end{aligned}$$

Input data

$$x_0 = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \quad \text{label: 0.5}$$

$$x_1 = \begin{bmatrix} 0.5 \\ 3 \end{bmatrix} \quad \text{label: 1.25}$$

Forward at t=0

Input activation $a_0 = \tanh(W_a \cdot x_0 + U_a \cdot \text{out}_{-1} + b_a)$
 $\rightarrow \tanh([0.05 \ 0.25] \begin{bmatrix} 1 \\ 2 \end{bmatrix} + [0.15][0] + [0.2]) = 0.81775$

Input gate $i_0 = \sigma(W_i \cdot x_0 + U_i \cdot \text{out}_{-1} + b_i)$
 $\rightarrow \sigma([0.95 \ 0.8] \begin{bmatrix} 1 \\ 2 \end{bmatrix} + [0.8][0] + [0.65]) = 0.96083$

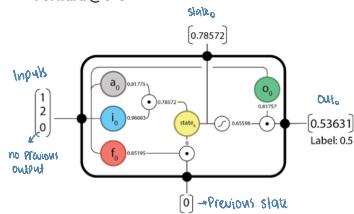
Forget gate $f_0 = \sigma(W_f \cdot x_0 + U_f \cdot \text{out}_{-1} + b_f)$
 $\rightarrow \sigma([0.1 \ 0.45] \begin{bmatrix} 1 \\ 2 \end{bmatrix} + [0.1][0] + [0.15]) = 0.8519$

Output gate $\text{O}_0 = \sigma(W_o \cdot x_0 + U_o \cdot \text{out}_{-1} + b_o)$
 $\rightarrow \sigma([0.6 \ 0.4] \begin{bmatrix} 1 \\ 2 \end{bmatrix} + [0.25][0] + [0.1]) = 0.81757$

Internal state $\text{State}_0 = a_0 \odot i_0 + f_0 \odot \text{state}_{-1}$
 $\rightarrow 0.81775 \odot 0.96083 + 0.8519 \odot 0 = 0.78571$

Output $\text{out}_0 = \tanh(\text{state}_0) \odot \text{O}_0$
 $\rightarrow \tanh(0.78571) \odot 0.81757 = 0.5363$

Forward @ t=0



$$\sigma: \frac{1}{1 + e^{-n}}$$

$$\tanh: \frac{e^n - e^{-n}}{e^n + e^{-n}}$$

Forward at t=1

Input activation $a_1 = \tanh([0.05 \ 0.25] \begin{bmatrix} 0.5 \\ 3 \end{bmatrix} + [0.15][0.53631] + [0.2]) = 0.84980$

Input gate $i_1 = \sigma([0.95 \ 0.8] \begin{bmatrix} 0.5 \\ 3 \end{bmatrix} + [0.8][0.53631] + [0.65]) = 0.98118$

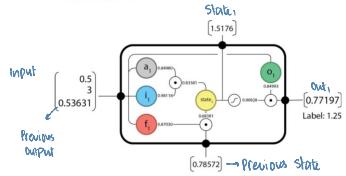
Forget gate $f_1 = \sigma([0.1 \ 0.45] \begin{bmatrix} 0.5 \\ 3 \end{bmatrix} + [0.1][0.53631] + [0.15]) = 0.81030$

Output gate $\text{O}_1 = \sigma([0.6 \ 0.4] \begin{bmatrix} 0.5 \\ 3 \end{bmatrix} + [0.25][0.53631] + [0.1]) = 0.84993$

Internal state $\text{State}_1 = a_1 \odot i_1 + f_1 \odot \text{state}_0$
 $\rightarrow 0.84980 \odot 0.98118 + 0.81030 \odot 0.78572 = 1.5176$

Output $\text{out}_1 = \tanh(1.5176) \odot 0.84993 = 0.7719$

Forward @ t=1



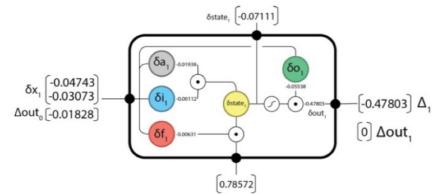
Then we will compute Backward for t=1

(b) For Backward @t=0; show the steps to get the values of various derivatives e.g. derivate of state,

x_{out_0}

Use L2 Loss

Backward @ t=1



Backward at t=1

Predicted output
↓
Actual output
↓

$$\delta_{out_i} = \Delta_i + \Delta_{out_i}$$

$$-0.77111 - 1.25 = -0.47803$$

$$\delta_{state_i} = \delta_{out_i} \odot 0_i \odot (1 - \tanh^2(state_i)) + \delta_{state_{i+1}} \odot f_{i+1}$$

$$-0.47803 \odot 0.84993 \odot (1 - \tanh^2(1.5176)) + 0 \odot 0 = -0.67111$$

$$\delta a_i = \delta_{state_i} \odot i_i \odot (1 - a_i^2)$$

$$-0.07111 \odot 0.98118 \odot (1 - 0.84993^2) = -0.01938$$

$$\delta i_i = \delta_{state_i} \odot a_i \odot i_i \odot (1 - i_i^2)$$

$$-0.07111 \odot 0.84993 \odot 0.98118 \odot (1 - 0.98118^2) = -0.00112$$

$$\delta f_i = \delta_{state_i} \odot a_i \odot i_i \odot f_i \odot (1 - f_i^2)$$

$$-0.07111 \odot 0.84993 \odot 0.98118 \odot 0.87030 \odot (1 - 0.87030^2) = -0.00630$$

$$\delta u_i = \delta_{gates_i}$$

$$0.07111 \odot \tanh(0.07111) \odot 0.84993 \odot (1 - 0.84993) = -0.5537$$

$$\begin{bmatrix} 0.05 & 0.95 & 0.1 & 0.6 \\ 0.25 & 0.8 & 0.45 & 0.4 \end{bmatrix} \begin{bmatrix} -0.01938 \\ -0.00112 \\ -0.00630 \\ -0.5537 \end{bmatrix} = \begin{bmatrix} -0.04743 \\ -0.03073 \end{bmatrix}$$

$$\Delta_{out_{i-1}} = U^T \cdot \delta_{gates_i}$$

$$\begin{bmatrix} 0.15 & 0.8 & 0.1 & 0.25 \end{bmatrix} \begin{bmatrix} -0.01938 \\ -0.00112 \\ -0.00630 \\ -0.5537 \end{bmatrix} = \begin{bmatrix} -0.01828 \end{bmatrix}$$

Backward @ t=0

Backward at t=0

$$\delta_{out_0} = 0.03631 + (-0.01828) = 0.01803$$

$$\delta_{state_0} = -0.5349$$

$$\delta_{a_0} = -0.1702$$

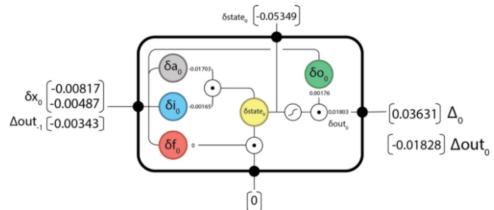
$$\delta_{i_0} = -0.00165$$

$$\delta_{f_0} = 0$$

$$\delta_{o_0} = 0.00176$$

$$\delta_{u_0} = \begin{bmatrix} 0.45 & 0.95 & 0.7 & 0.6 \\ 0.25 & 0.8 & 0.85 & 0.4 \end{bmatrix} \begin{bmatrix} -0.01702 \\ -0.00165 \\ 0 \\ 0.00176 \end{bmatrix} = \begin{bmatrix} -0.00817 \\ -0.00487 \\ 0 \\ -0.00343 \end{bmatrix}$$

$$\Delta_{out_0} = [0.15 \ 0.8 \ 0.1 \ 0.25] \begin{bmatrix} -0.01702 \\ -0.00165 \\ 0 \\ 0.00176 \end{bmatrix} = [-0.003]$$



(c) Use Forward t=0, Backward t=0 above, Forward t=1, Backward t=1 below, compute δW , δb and δb

Sol:

$$\begin{aligned} \delta W &= \sum_{t=0}^T \delta_{gates_t} \otimes x_t \\ &= \begin{bmatrix} -0.01703 \\ -0.00165 \\ 0 \\ 0.00176 \end{bmatrix} \begin{bmatrix} 1.0 & 2.0 \\ 0.5 & 3.0 \end{bmatrix} + \begin{bmatrix} -0.01938 \\ -0.00112 \\ -0.00631 \\ -0.05538 \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} -0.02672 & -0.0922 \\ -0.00221 & -0.0666 \\ -0.00316 & -0.01893 \\ -0.02593 & -0.16262 \end{bmatrix} \rightarrow \text{old} \\ \delta U &= \sum_{t=0}^{T-1} \delta_{gates_{t+1}} \otimes out_t \\ &= \begin{bmatrix} -0.01938 \\ -0.00112 \\ -0.00631 \\ -0.05538 \end{bmatrix} \begin{bmatrix} 0.53631 \\ out_0 \end{bmatrix} = \begin{bmatrix} -0.01039 \\ -0.00060 \\ -0.00338 \\ -0.02970 \end{bmatrix} \end{aligned}$$

$$\begin{aligned} \delta b &= \sum_{t=0}^T \delta_{gates_{t+1}} \\ &= \begin{bmatrix} -0.01703 \\ -0.00165 \\ 0 \\ 0.00176 \end{bmatrix} + \begin{bmatrix} -0.01938 \\ -0.00112 \\ -0.00631 \\ -0.05538 \end{bmatrix} = \begin{bmatrix} -0.03641 \\ -0.00277 \\ -0.00631 \\ -0.05362 \end{bmatrix} \end{aligned}$$

Bias + Gates

Internal Parameters Update

$$\delta W = \sum_{t=0}^T \delta_{gates_t} \otimes v_t$$

$$\delta U = \sum_{t=0}^T \delta_{gates_{t+1}} \otimes out_t$$

$$\delta b = \sum_{t=0}^T \delta_{gates_{t+1}}$$

(d) Update parameters based on SGD update

$$W_a = \begin{bmatrix} 0.45267 \\ 0.25922 \end{bmatrix}, U_a = [0.15104], b_a = [0.20364]$$

$$W_i = \begin{bmatrix} 0.95022 \\ 0.80067 \end{bmatrix}, U_i = [0.80006], b_i = [0.65028]$$

$$W_f = \begin{bmatrix} 0.70031 \\ 0.45189 \end{bmatrix}, U_f = [0.10034], b_f = [0.15063]$$

$$W_o = \begin{bmatrix} 0.60259 \\ 0.41626 \end{bmatrix}, U_o = [0.25297], b_o = [0.10536]$$

$$W^{new} = W^{old} - \gamma \cdot \delta W^{old} \quad \gamma = 0.1$$

$$W_a^{old} - \gamma \cdot \delta W_a^{old}$$

$$\begin{bmatrix} 0.45 \\ 0.25 \end{bmatrix} - 0.1 \begin{bmatrix} -0.02672 \\ -0.0922 \end{bmatrix} \cdot \begin{bmatrix} 0.05267 \\ 0.25922 \end{bmatrix}$$

$$PE(pos, 2i+1) = \cos\left(\frac{pos}{10000^{\frac{2i+1}{4}}}\right)$$

$$PE(pos, 2i) = \sin\left(\frac{pos}{10000^{\frac{2i}{4}}}\right)$$

Positional Encoding

Q) POS: 1, d=4 so i=0 and i=1 → after every 2 dim i increases

1. Dim 0 → sin → even

$$PE(1, 0) = \sin\left(\frac{1}{10000^{\frac{1}{4}}}\right) = 0.8415$$

2. Dim 1 → cos → odd

$$PE(1, 1) = \cos\left(\frac{1}{10000^{\frac{1}{4}}}\right) = 0.5403$$

3. Dim 2 → sin

$$PE(1, 2) = \sin\left(\frac{1}{10000^{\frac{2}{4}}}\right) = 0.0099$$

4. Dim 3 → cos

$$PE(1, 3) = \cos\left(\frac{1}{10000^{\frac{3}{4}}}\right) = 0.99995$$

So Final Positional Encoding Vector

for Position 1 and d=4

$$(0.8415, 0.5403, 0.0099, 0.99995)$$