

Create a Large Language Model from Scratch with Python

classmate

Date _____
Page _____

- Freecodecamp.org

* Packages to install

numpy, matplotlib, pygma,
ipykernel, jupyter, torch

* open the file in read mode

→ with open('wizard_of_oz.txt', 'r', encoding='utf-8') as f:

tent = f.read()

print(tent[:200])

→

lets sort these characters inside a list

→ chars = sorted(set(tent))

print(chars)

→ [., , , ' ', '2', --]

* we have 81 characters

Code to tokenize these text
 here we are just mapping each char to a number

→ string_to_int = {ch: i for i, ch in enumerate(chars)}
 int_to_string = {i: ch for i, ch in enumerate(chars)}

encode = lambda s: [string_to_int[c] for c in s]
 decode = lambda l: ''.join([int_to_string[i] for i

print(encode('hello'))

→ [61, 58, 65, 65, 68]

if we are working on different language
 we will have a lot of tokens even

so we convert our data into tensors
 so that our model processes it

→ data = torch.tensor(encode(text), dtype=torch.long)

print(data[:100])

* Next topic is on train, val, test sets
why we do it

we know why we do it
so we move on.

* Premise of Bigram model

hello

start of content $\rightarrow h$
 $h \rightarrow e$

$e \rightarrow l$ frequency sweet and so

$l \rightarrow l$ Fairchild : 7

$l \rightarrow o$

Bi-gram : 1

so here we have associative probability

at start it is h
next after h , e is more probability

after e , l is more probability
to come after e

and so on

that's it this is Bigram Model
as it considers the previous char
and next char

about this the proof

$\text{block_size} = 5 \rightarrow$ the size of the block from the entire set we are gonna iterate over

... $[5, 67, 21, 82, 54]$ 35
starts at 5 moves to 5 blocks after 54 predicted

... 5 $[67, 21, 82, 54, 35]$.
next steps now start is 67
and 5 blocks after 35
so stride = 1

and the error to minimize
is the predicted and
actual char

so we have snippet

$[: \text{block_size}]$

next

$[1 : \text{block_size} + 1]$

same thing in code

$\rightarrow \text{block_size} = 8$

$x = \text{train_data} [: \text{block_size}]$

$y = \text{train_target} [1 : \text{block_size} + 1]$

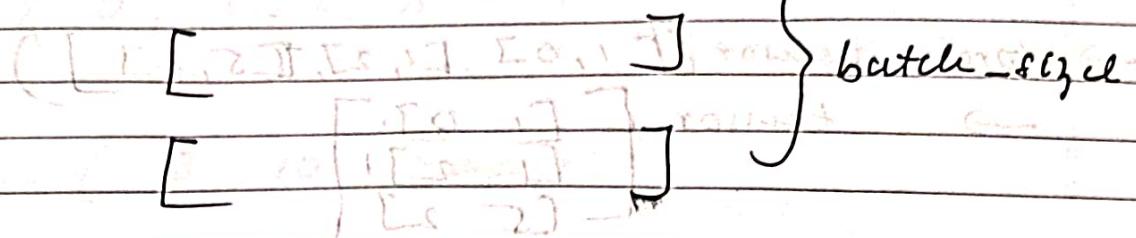
for t in range (block_size):

content = $x [: t + 1]$

target = $y[t]$

print ('when input is', content,
'target is', target)

~~block~~ → [.]



→ if device = 'cuda' if torch.cuda.is_available() else ('cpu')

print(device)

block_size = 8

batch_size = 4

(E, S) ~~new. short~~

F ~ ~~segm + T~~

answ
answ is final
Name

Answers

Some basic Pytorch functions

→ torch.tensor([1, 0], [1, 2], [5, 2])

$$\rightarrow \text{tensor} \begin{bmatrix} [1 & 0] \\ [1 & 2] \\ [5 & 2] \end{bmatrix}$$

→ torch.randn(-100, 100, (6, 1))

$$\rightarrow \text{tensor} \begin{bmatrix} -9 & -53 & 54 & -87 & 35 & 98 \end{bmatrix}$$

→ torch.zeros(2, 3)

$$\rightarrow \text{tensor} \begin{bmatrix} [0 & 0 & 0] \\ [0 & 0 & 0] \end{bmatrix}$$

→ torch.ones(2, 3)

→ same but 1

→ torch.empty(2, 3)

$$\rightarrow \begin{bmatrix} 1. \times e^{+22} \end{bmatrix}$$

$$\begin{bmatrix} \dots \end{bmatrix}$$

Some large and small values

$\rightarrow \text{torch.arange}(5)$

$\rightarrow [0, 1, 2, 3, 4]$

torch

$\rightarrow \text{x.linspace}(3, 10, \text{steps} = 5)$

basically increment from 3 to 10
in constant

$\rightarrow [3.00, 4.750, 6.500, 8.250, 10.000]$

$\rightarrow \text{torch.eye}(5)$

$\rightarrow [1 0 0 0 0]$

[0 1 0 0 0] request address

[0 0 1 0 0]

[0 0 0 1 0]

[0 0 0 0 1]

gives
diagonal matrix

CPU vs GPU Performance

→ device = 'cuda' if torch.cuda.is_available() else 'cpu'

→ cuda

it will work through pinned memory

* Well we know that GPU performs better

* In order to multiply matrices in pytorch we need to use @

[torch-tensor_1 @ torch-tensor_2]

~~for learning~~

~~workshop~~

More Pytorch functions

- † torch.multinomial
- † define a probability tensor
- probabilities = torch.tensor([0.1, 0.9])
10% ⇒ for 0, 90% for 1
- † Draw 10 samples from multinomial distribution
sample = torch.multinomial(probabilities, num_samples=10, replacement=True)
- Tensor ([1, 1, 1, 1, 1, 1, 0, 1, 1, 1])

- † torch.cat
- † to concatenated 2 tensors
- † torch.tril(torch.ones(5, 5))

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

torch.triu

* same as before but

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

reverse

and torch like, lower triangular matrix required or upper

matrix, $\text{triu}(A)$ (lower triangular) $\text{triu}(A)$ (upper)

triu(A) = A

.Transpose()

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

torch.stack([tensor-1, tensor-2, tensor-3])

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

stack them

and shape

becomes 3×3 $\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}$ $\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}$ $\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}$ $\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}$

```

import
from torch import nn
→ import torch.nn as nn
    
```

This is a module where all the layers are present

nn.Linear

nn.ReLU

nn.Conv2d

... more

nn.Embedding

This is a layer where a char or word which is sent to a model

is sent along with a vector which holds the sentiment or importance of the char as this vector is trainable and hence is a layer.

embedding_size = 5 → hyper parameter

char level a [0.1 0.2 0.6 0.8 0.3]

b [0.2 0.5 0.1 0.2 0.5]

c [0.9 0.7 0.2 0.6 0.1]

z [0.8 0.5 0.3 0.4 0.2]

* Matrix Multiplication

* Dot product

$$\begin{bmatrix} 1, & 2, & 3 \end{bmatrix} \cdot \begin{bmatrix} 4, & 5, & 6 \end{bmatrix}$$

$$(1 \times 4) + (2 \times 5) + (3 \times 6) = 32$$

a is 3×2 mat

b is 2×3 mat

a @ b

$$\begin{matrix} 3 \times 2 \\ \text{rows} \end{matrix} \quad \begin{matrix} 2 \times 3 \\ \text{cols} \end{matrix} = 3 \times 3$$

must
be same

$$\begin{matrix} 3 \times 4 \\ \text{rows} \end{matrix} \quad \begin{matrix} 5 \times 3 \\ \text{cols} \end{matrix}$$

$\begin{bmatrix} 1.0 & 3.0 & 3.0 & 1.0 \end{bmatrix} \quad \begin{bmatrix} 1.0 \\ 2.0 \\ 5.0 \\ 1.0 \end{bmatrix}$ not computable for multiplication
 $\begin{bmatrix} 1.0 & 3.0 & 3.0 & 4.0 & 8.0 \end{bmatrix}$ must do some operations

for make it work

like swap

$$5 \times 3 \quad 3 \times 4 = 5 \times 4$$

$$a = \begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{bmatrix}$$

$$b = \begin{bmatrix} 7 & 8 & 9 \\ 10 & 11 & 12 \end{bmatrix}$$

$$(1 \times 7) + (2 \times 10) = 27$$

$$(1 \times 8) + (2 \times 11) = 30$$

$$(1 \times 9) + (2 \times 12) = 33$$

$$(3 \times 7) + (4 \times 10) = 47$$

	27	30	33
	47	68	75
	97	108	117

In code

$\rightarrow a = \text{Torch}\cdot\text{tensor}([1 \ 2 \ 3 \ 4])$

$b = \text{Torch}\cdot\text{tensor}([7 \ 8 \ 9 \ 10 \ 11 \ 12])$

$\text{print}(a @ b)$ or $(\text{torch}\cdot\text{matmul}(a, b))$

\rightarrow same tensor

* Can't math mul on int and float
must convert any into same type
by casting



$$\text{fr} = \text{float}(5)(8)$$

$$fr = (5 \times 8) + (0 \times 1)$$

$$fr = (5 \times 8) + (0 \times 1)$$

$$fr = (0 \times 8) + (5 \times 1)$$

3. i + f =

fr

int

float

float int

↳ float(f) requires float or double

↳ float(f) = d

(d, s) conversion (Numpy) ↳ ((d, 0), 0) - kind

→ more
related notes

Model demo for embeddings

→ Class `BigramModel(nn.Module)`:

```
def __init__(self, vocab_size):
```

```
super().__init__()
```

```
self.TokenEmbedTable = nn.Embedding(vocab_size, vocab_size)
```

```
def forward(self, index, targets=None):
```

```
logits = self.TokenEmbedTable(index)
```

in between $B, T, C = \text{logits}.shape$

```
code logits = logits.view(B*T, C)
```

→ if targets is None: targets = targets.view(B*T)

```
loss = None loss = F.CrossEntropy(logits, targets)
```

else:

```
return logits, loss
```

What are logits, view function for

logits is basically our table of bigram

times a, b, c, d in follow ups

aa, ab, ba, bb

c

d

and view?

$B \rightarrow$ batch_size

$T \rightarrow$ sequence length

$C \rightarrow$ vocab_size (defined as no. of classes for predictions)

$$(B, T, C) = \text{logits} \cdot \text{shape}$$

breaks down the shape

$$\text{logits} = \text{logits}.view(B*T, C)$$

Example:

$$B = 4$$

$$T = 8$$

$$C = 100$$

$$\text{logits}.shape = (4, 8, 100)$$

$$\text{targets}.shape = (4, 8)$$

\rightarrow why $.view(B*T, C)$

so we need to flatten the predictions and targets so

the cross_entropy() can work properly

`logits = logits.view(B*T, c)`

$$\# (4*8, 100) \Rightarrow (32, 100)$$

`targets = targets.view(B*T)`

$$\# (4*8) \Rightarrow (32,)$$

as there are some expected input

$$\text{logits} \rightarrow (N, c) \rightarrow (B*T, c)$$

$$\text{targets} \rightarrow (N) \rightarrow (B*T)$$

so `.view` reshapes to a flat list of predictions so loss can be calculated.

→ torch.long
is same
as torch.int64

classmate

Date _____

Page _____

* Continue the model (\rightarrow TMA) with finalization
 $(\text{out}, \text{es}) = (\text{out}, \text{es}')$ →
→ def generate(self, index, max_new_tokens):
 for i in range(max_new_tokens):

 logits, _ = self.forward(index)

 logits = logits[:, -1, :] # because (B, C).

 probs = F.softmax(logits, dim=-1)

 index_next = torch.multinomial(probs, 1, device=device, num_samples=1)

 index = torch.cat((index, index_next), dim=1) # (B, T+1)

 return index

model = BigramModel(vocab_size)

m = model.to(device)

content = torch.zeros(1, 1), dtype=torch.long, device=device)

→ generated_chars = decode(m.generate(content, max_new_tokens=500)[0].tolist())

print(generated_chars)

→ max_iter = 10000

learning_rate = 3e-4

* Training Loop

→ optimizer = torch.optim.AdamW(model.parameters(), lr = learning_rate)

for iter in range(max_iter):

sample a batch of data

xb, yb = get_batch('train')

→ This a dataloader

basically will
give in batches.

evaluate the loss

logits, loss = model.forward(xb)

optimizer.zero_grad(set_to_none=True)

loss.backward()

optimizer.step()

print(loss.item())

Optimizers Overview

Some common ones

Gradient descent - to min loss function of learning model

Momentum - intution of GD with addeds of "momentum" term helps smooth out and updates, useful for deep neural network

RMSprop

Adam - ideal combines momentum and RMSprop

AdamW - modified Adam that adds weights decay to the params. can improve generalization performance

Loss Reporting + Train Vs Eval mode

→ eval_iters = 250
dropout = 0.2

→ @ torch.no_grad() → Gradient is not happening here

```

def estimate_loss():
    out = {}
    model.eval()
    for split in ['train', 'val']:
        losses = torch.zeros(eval_iters)
        for k in range(eval_iters):
            X, Y = get_batch(split)
            logits, loss = model(X, Y)
            losses[k] = loss.item()
        out[split] = losses.mean()
    model.train()
    return out
  
```

→ * between your training loop

after every $\text{range}(\text{max_iters})$:

if $\text{iter \% eval_iters} == 0$:

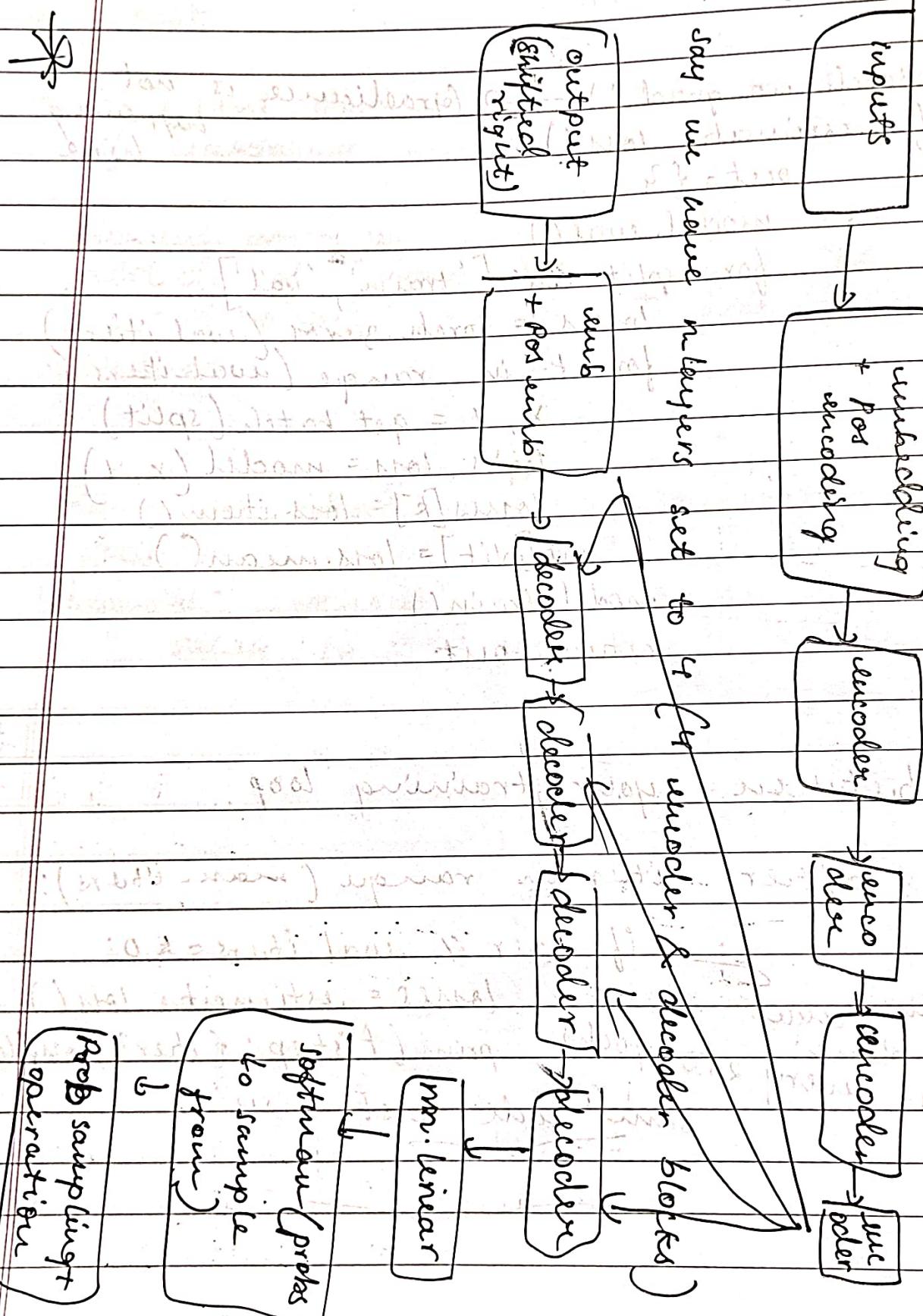
mean loss is done after every 250 epochs same code rest

```

losses = estimate_loss()
print(f"step: {iter}, loss: {losses}")
  
```

Transformer

Models



~~Inside
encoder~~

multi-head
attention

Residual
connection
add then
norm

Feed forward

Residual
con
add → norm

Feeds into next

encoder block ↓ it

its the last
encoder block it feeds
into each decoder
block.

~~Inside~~

Feed Forward

Linear → ReLU → Linear

Inside Multi-head Attention

key(k) → n-heads

Query(Q) → scaled
dot product
attention

Value(V) →

scaled dot product
runs n heads in
parallel from GPU

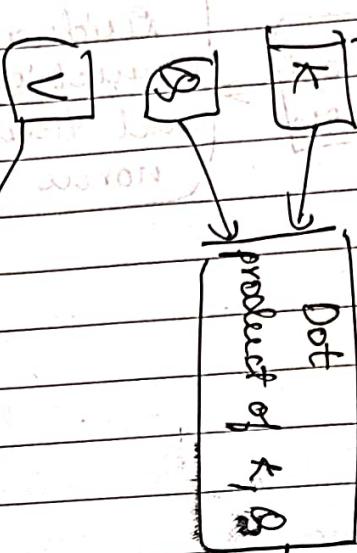
concatenate
results

Linear

* we call it multi-head
attention because there
are n-heads learning bunch of
different semantic info from
a unique perspective

know their
values
are learnt
by passing
through
a neural
network

* inside each head of attention



Scaled by length of a row in
a key or
queries matrix OR
dot

torch.tril
masking
(sometimes)

Softmax
 \rightarrow multiplicative
softmax

use know
that torch.tril
gives $\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$ like
this kernel
of so met.

output: blend
of input
vector values
and attention
placed on
each token

predictions should
be based on past
and must have no clue of future
as it
is masked
multi-head
layer.

How this works

Self-Attention allows a model to focus on different parts of the input when processing each word.

"The animal didn't cross the street because it was too weird"

when processing the word "it" the model needs to understand whether "it" refers to "animal" or "street".

Self Attention helps with

key, Query, Value?

These are vectors derived from input embedding and are used in the attention mechanism

Input Token Embeddings

Suppose "The cat sat"

cat position 0, has embedding 1

sat position 2, has embedding 3

the position 1, has embedding 4

You convert each word into a vector (embedding).

Suppose

$$u_1 = \text{embedding}(\text{"The"})$$

$$u_2 = \text{"cat"}$$

$$u_3 = \text{"sat"}$$

Each of these vectors is then transformed into

Q

K

V

This is done using learnable weight matrices

$$Q = u \cdot W_Q$$

$$K = u \cdot W_K$$

$$V = u \cdot W_V$$

where:

W_Q, W_K, W_V are trainable matrices (parameters) of the model.

So for each token, you generate is?

Q : what this word is asking for

K : what this word contains

V : what this word carries as info

How Attention is calculated?

The idea is to compare the Query of the current word to the keys of the other words, to compute how much attention should be paid to each word.

Steps:

1. Dot product of Query and Keys.

2. Scale.

To avoid large values (especially with high dimensional vectors), scale by \sqrt{d} ($d = \text{dimension of key}$)

3. Softmax

4. Multiply by Value

All together!

Imagine you're processing the token "sat".

- You are the Θ of "sat"
- you compare it with key of ["The", "cat", "sat"]

- Based on these comparisons (dot products), you compute attention weights.
- Then, you take a weighted sum of the value of ["The", "cat", "sat"]

This gives a new content-aware representation of "sat", influenced by the words around it.

Visualization

Inputs

"The" "cat" "sat"

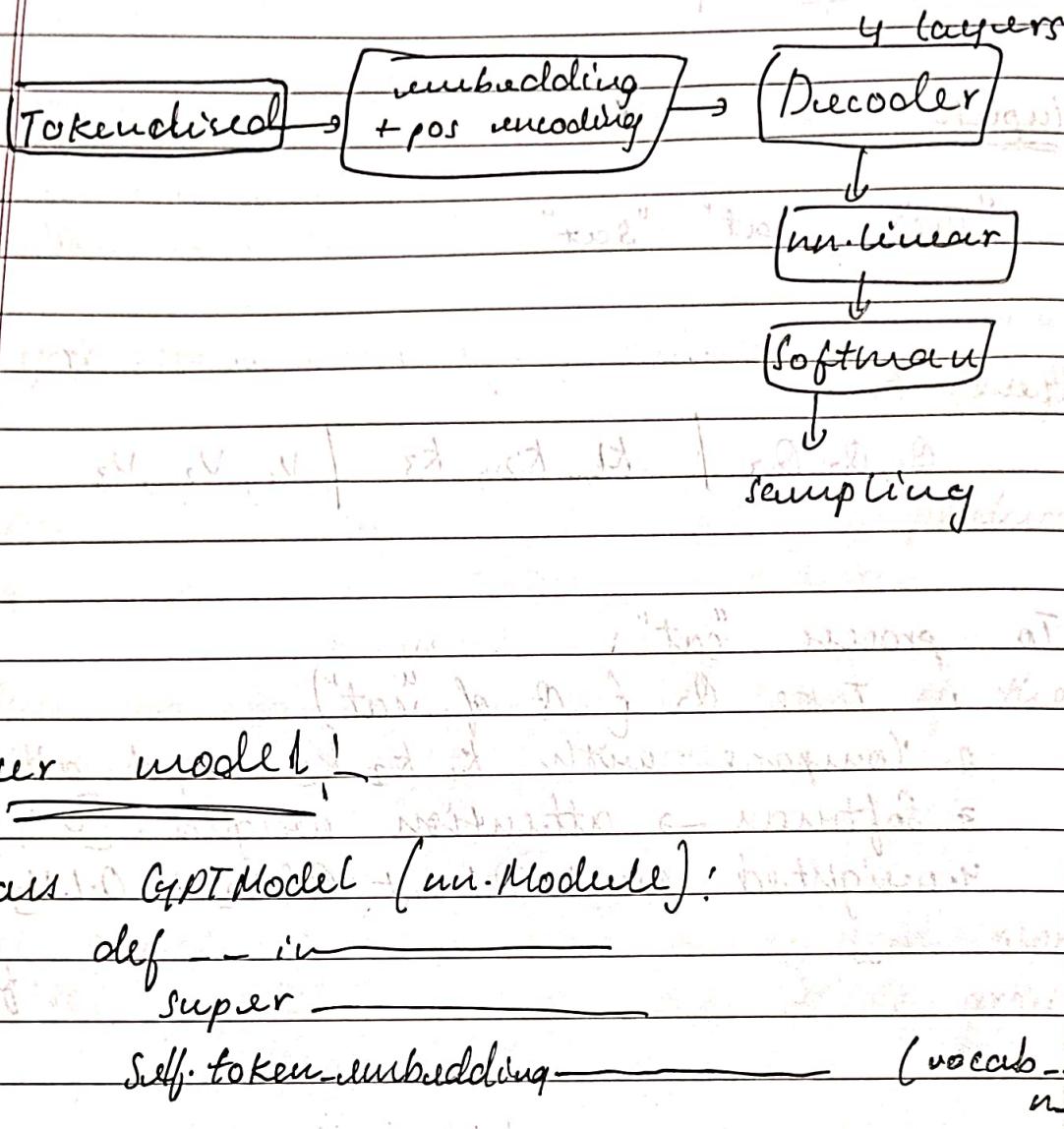
Generate:

$\theta_1, \theta_2, \theta_3$	k_1, k_2, k_3	v_1, v_2, v_3
--------------------------------	-----------------	-----------------

To process "cat":

1. Take θ_2 (θ of "cat")
2. Compare with $k_1, k_2, k_3 \rightarrow$ get attention score
3. Softmax \rightarrow attention weights: $[0.1, 0.8, 0.1]$
4. weighted sum: $0.1 \cdot v_1 + 0.8 \cdot v_2 + 0.1 \cdot v_3 \rightarrow$ new embedding for "cat"

Our GPT model!



decoder ↙

```

    self.blocks = nn.Sequential(*[Block(n_embd,
                                         n_head=n_head)
                                for _ in range(n_layer)])
  
```

self.ln_f = nn.LayerNorm(n_embd)

self.ln_h = nn.Linear(n_embd, vocab_size)

```
def forward(self, index, targets = None):
    logits = self.token_embedding_table(index)
```

$$\text{token_emb} = \text{self}. \text{token_emb}(\text{idu})$$

$$\text{pos_emb} = \text{self}. \text{pos}(\text{T}, \text{device} = \text{device})$$

$$\text{u} = \text{token_emb} + \text{pos_emb} \quad \# (\text{B}, \text{T}, \text{C})$$

$$\text{u} = \text{self}. \text{blocks}(\text{u})$$

$$\text{u} = \text{self}. \text{ln_f}(\text{u})$$

$$\text{u} = \text{self}. \text{ln_head}(\text{u})$$

If target is None:

~~if target is None~~

lets initialize initial weights

`self.apply(self.init_weights)`

```
def init_weights(self, module):
```

```
if isinstance(module, nn.Linear):
```

`torch.nn.init.normal_(module.weight, mean`

`= 0.0`

`std = 0.02)`

If module.bias is not None:

`torch.nn.init.zeros_(module.bias)`

```
elif isinstance(module, nn.Embedding):
```

`torch.nn.init.normal_(module.weight,`

`mean=0.0, std=0.02)`

* this should be in `--init--`

`num_heads = 4`

$\rightarrow n_embd = 384$

`n_head = 4`

`dropout = 0.2`

* Our block model Inside Decoder

\rightarrow Class Block (`nn.Module`):

```
def __init__(self, n_embd, n_head):
    super().__init__()
    head_size = n_embd // n_head
    self.sa = MultiheadAttention(n_head, head_size)
```

`self.ffwd = Feedforward (n_embd)`
`self.ln1 = nn.LayerNorm (n_embd)`
`self.ln2 = nn.LayerNorm (n_embd)`

`def forward(self, u):`

`y = self.sa(u)`

`u = self.ln1(y)`

`y = self.ffwd(u)`

`u = self.ln2(y)`

`return u`

0.0 = 0

0.0 = 2

Learn from the implementation

Let's make our Feedforward

→ Class Feedfor (nn.module):

```
def __init__(self):
```

```
super().__init__()
```

```
self.net = nn.Sequential(
```

```
nn.Linear(n_embd, 4 * head_size),  
nn.ReLU(),
```

```
) nn.Dropout(dropout),
```

```
def forward(self, u):
```

```
return self.net(u)
```

Now make our Multi-head attention class

→ Class Multi (nn.module):

```
def __init__(self):
```

```
super().__init__()
```

```
self.heads = nn.ModuleList([Head(head_size)  
for _ in range(num_heads)])
```

```
self.proj = nn.Linear(head_size * num_heads, n_embd)
```

```
self.dropout = nn.Dropout(dropout)
```

```
def forward(self, u):
```

```
out = torch.cat([h(u) for h in self.heads], dim=-1)
```

~~out~~out = self.dropout(~~dropout~~)

self.proj(out))

return out

* Head Class

→ Class + Head (nn.module):

```

def __init__(self, n_head):
    super().__init__()
    self.key = nn.Linear(nn_size, head_size, bias=False)
    self.query = nn.Linear(nn_size, head_size, bias=False)
    self.register_buffer('tril', torch.tril(torch.ones(nn_size, nn_size)))
    self.drop = nn.Dropout(dropout)

```

def forward(self, x):

head size

```

B,T,L = u.shape
k = self.key(u) # (B, T, hs)
q = self.query(u) # (B, T, hs)

```

wei³ q @ k.transpose(-2, -1) * k.shape[-1]

** -0.5

(B, T, hs) @ (B, hs, T) → (B, T, T)

$\text{wei} = \text{wei} \cdot \text{masked_fill}(\text{self}. \text{tri}[:T, :T] = 0)$

$\text{float}(' -\inf '))$

$\text{wei} = F \cdot \text{softmax}(\text{wei}, \text{dim} = -1)$

$\text{wei} = \text{self}. \text{dropout}(\text{wei})$

$v = \text{self}. \text{value}(n)$

$\text{out} = \text{wei} @ v$

return out

* What is difference between sequential module and modulelist module.

both return same

where sequential depends on the previous layer to complete and process the same

module list runs everything in parallel so they are isolated and have different inputs, and different outputs

this is what we want in head.

* Now start training by troubleshooting some errors.

and adjust the hyperparameters

and use different dataset or large

which is good one

but mind your computation and memory capability

* We are gonna use OpenWebText data

which has more than 20000 so have to use for loops

data extract file

so we gotta save all these into one .tut file

do chatpt

then add as following for code

and write

import tensorflow as tf

tf.enable_eager_execution()

tf.tpu.experimental.initialize_tpu_system()

tpu = tpu.core

model = tpu.get_model()

model.compile(optimizer='adam', loss='categorical_crossentropy')

and also split the data into train, val

train 90%

val 10%

+ Add a dataloader to feed batches to the model.

+ Now that we are training we want to save our model weights so we use something called \rightarrow import pickle

\rightarrow after training loop a file will get saved

with open('model-01.pkl', 'wb') as f:
pickle.dump(model, f)

↓ Now we can save the model and load into the model GPT() and do more epochs further

```
→ model = GPTModel(vocab_size)
print('loading model... ')
with open('model-01.pkl', 'rb') as f:
    model = pickle.load(f)
print('loaded')
```

m = model.to(device)

* Now do the training loop and the
model will resume
training

* Our bot script after all

→ while True:

```
    prompt = input("prompt: ")
    content = torch.tensor(encode(prompt), dtype=torch.long,
                          device=device)
    generated_chars = decode(m.generate(content,
                                         max_tokens=150)[0].unsqueeze(0).item())
    print(f'completion: {generated_chars}' )
```

Make sure the block_size and max-new-tokens is

using padding to or less or else we can have error

or we can do crop index to the last block_size tokens

→ index could = index[:, block_size:]
keeps in block size limit

Research and Development Pointers

→ import time

Time and efficiency
testings

start = time.time()

```
for i in range(1000):
    print(i**2)
```

end = time.time()

total = end - start

print(f'time taken : {total}')

Can be used to time and operations or check how long our model takes to load.

* Quantized Clues

instead of 32 bit floating point
it uses 16 bit float point
which we learnt
in CV course.

* Gradient Accumulation

ChatGPT this to understand

* Flipping Face

if we have models, datasets, and more

* should explore. Do a course