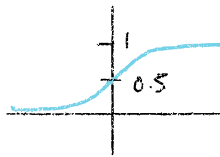


# Logistic Regression

The logistic function is an S-shaped or "sigmoid" function, which takes any real input  $z$ , and outputs a value between zero and one.

$$f(z) = \frac{1}{1 + e^{-z}} \quad \text{graphs to} \quad \sigma(z)$$


Let  $z$  be the linear combination of the input features:

$$z = \bar{w}^T \bar{x} + b = \sum_{i=1}^n w_i x_i + b$$

Here,

$\bar{w} = [w_1, w_2, \dots, w_n]$  is the vector of weights

$\bar{x} = [x_1, x_2, \dots, x_n]$  is the input feature vector

$b$  is the bias term.

So  $\sigma$  squashes  $z$  into the range  $(0, 1)$ .

$$P(y=1|x) = \sigma(z) = \frac{1}{1 + e^{-(\bar{w}^T \bar{x} + b)}}$$

Classifying 0, or 1

$$P(y=0|x) = 1 - \sigma(z)$$

In this way, the sigmoid function can be used to model the binary classification task, where

$$y = \begin{cases} 1 & \text{if } \sigma(z) \geq 0.5 \\ 0 & \text{if } \sigma(z) < 0.5 \end{cases}$$

The model parameters (weights  $\bar{w}$  and bias  $b$ ) can be trained through binary cross-entropy (BCE)

$$L(y, \hat{y}) = -[y \log(\hat{y}) + (1-y) \log(1-\hat{y})]$$

↳  $\sigma(z)$

Where:

- $y$  is the true label (either 0 or 1)
  - $\hat{y}$  is the predicted probability of the class being 1
- 

### Deriving BCE

The probability of the true class label  $y$  can be modeled as:

$$P(y | x, \bar{w}) = \begin{cases} \hat{y}, & \text{if } y=1 \\ 1-\hat{y}, & \text{if } y=0 \end{cases}$$

which can be re-written as:

$$P(y | x, \bar{w}) = \hat{y}^y (1-\hat{y})^{1-y}$$

Assuming we have dataset with  $N$  samples, the likelihood  $\mathcal{L}$  of observing the true labels  $y_1, y_2, \dots, y_N$  given the probabilities  $\hat{y}_1, \hat{y}_2, \dots, \hat{y}_N$  is the probability of the individual probabilities:

$$\mathcal{L}(\bar{w}) = \prod_{i=1}^N P(y_i | \bar{x}_i, \bar{w}) = \prod_{i=1}^N \hat{y}_i^{y_i} (1-\hat{y}_i)^{1-y_i}$$

For easier computation, we apply  $\log$ :

$$\begin{aligned}\log \mathcal{L}(\bar{w}) &= \sum_{i=1}^N \log(\hat{y}_i^{y_i} (1-\hat{y}_i)^{1-y_i}) \\ &= \sum_{i=1}^N (y_i \log(\hat{y}_i) + (1-y_i) \log(1-\hat{y}_i))\end{aligned}$$

Finally, we apply a negative which makes our goal minimizing the negative log-likelihood. This makes the optimization task more suitable for gradient descent.

---

During training, the goal is to minimize the BCE loss by adjusting the model's parameters  $\bar{w}$ .

$$\bar{w} \leftarrow \bar{w} - \alpha \nabla_{\bar{w}} \mathcal{L}(\bar{w})$$

where

$$- \mathcal{L}(\bar{w}) = -\frac{1}{N} \sum_{i=1}^N (y_i \log(\hat{y}_i) + (1-y_i) \log(1-\hat{y}_i))$$

-  $\alpha$  is the learning rate

-  $\nabla_{\bar{w}} \mathcal{L}(\bar{w})$  is the gradient of the loss function with respect to  $\bar{w}$ .

Let  $l(\bar{w})$  be the single sample BCE loss:

$$l(\bar{w}) = -(y \log(\hat{y}) + (1-y) \log(1-\hat{y}))$$

Then, calculate the derivative of  $l(\bar{w})$  w.r.t.  $\bar{w}$ :

$$\text{Chain rule: } \frac{\partial l(\bar{w})}{\partial \bar{w}} = \frac{\partial l(\bar{w})}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial \bar{w}}$$

$$\frac{\partial \ell(\bar{w})}{\partial \bar{w}} = \frac{\partial \ell(\bar{w})}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial \bar{w}}$$

$$(1) \quad \frac{\partial \ell(\bar{w})}{\partial \hat{y}} = - \left( \frac{y}{\hat{y}} - \frac{1-y}{1-\hat{y}} \right) = \frac{\hat{y} - y}{y(1-\hat{y})}$$

$$(2) \quad \frac{\partial \hat{y}}{\partial z} = \hat{y}(1-\hat{y})$$

$$(3) \quad \frac{\partial z}{\partial \bar{w}} = \bar{x}$$

$$\therefore \frac{\partial J(\bar{w})}{\partial \bar{w}} = \bar{x} (\hat{y} - y)$$

For the full batch:

$$\nabla_{\bar{w}} \mathcal{L}(\bar{w}) = \frac{1}{N} \sum_{i=1}^N (\hat{y}_i - y_i) \bar{x}_i$$

This fixed gradient descent is actually quite intuitive, simply averaging the distance between the correct label and the predicted label.

To generalize this to multi-class classification, all we need to do is change  $g$  to the softmax function, and recalculate the gradient descent.

The softmax function is a generalization of the logistic (or sigmoid) function for multiple dimensions. In neural networks, it is used to normalize the output of a network to a probability distribution over predicted output classes.

meaning,  
each vector component will be in the interval  $(0,1)$ , and will add up to 1.

$$\sigma(z)_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

The softmax function applies the standard exponential function to each element  $z_i$  of the input vector  $\vec{z}$ , and normalizes these values by dividing the sum of all these exponentials. This normalization ensures that the sum of the components of the output vector  $\sigma(\vec{z})$  is 1.

"Softmax" derives from the amplifying effects of the exp. on any maxima in the input vector.

For example,

$$\text{softmax}((1, 2, 8)) = (0.001, 0.002, 0.997)$$

We don't necessarily have to use  $e$  as base, any  $b > 0$  can be used. Such  $b$  can be expressed as  $e^{\beta}$ , for any real  $\beta$ . But  $e$  ensures easy derivatives, and  $\log$ 's.

Deriving  $\nabla$  for multiclass classification with softmax:

### 1. Negative log-likelihood

$$\textcircled{1} P(y | \bar{x}, \bar{w}) = \prod_{i=1}^N P(y_i | x_i, \bar{w}_i) = \prod_{i=1}^N \hat{y}_{i, k_i}$$

where  $k_i$  is the index of the true class for sample  $i$

$$\textcircled{2} \log \mathcal{L}(\bar{w}) = \sum_{i=1}^N \log(\hat{y}_{i, k_i}) = \sum_{i=1}^N \sum_{k=1}^K y_{i, k} \log(\hat{y}_{i, k})$$

where  $y_{i, k}$  is 1 if true class is  $k$  for sample  $i$ , and 0 otherwise.

$\textcircled{3}$  Apply the negative

$$-\mathcal{L}(\bar{w}) = - \sum_{i=1}^N \sum_{k=1}^K y_{i, k} \log(\hat{y}_{i, k})$$

### 2. Computing the gradient

$$\frac{d\mathcal{L}}{d\hat{y}_{i, k}} = \hat{y}_{i, k} - y_{i, k}, \quad \frac{d\hat{y}_{i, k}}{dz_j} = \hat{y}_{i, k} (\delta_{jk} - \hat{y}_{i, j})$$

↳ Kronecker delta:

1 if  $j = k$ , 0 otherwise

$$\frac{dJ}{d\bar{w}} = \sum_{i=1}^N \sum_{k=1}^K (\hat{y}_{ik} - y_{ik}) \frac{\partial \hat{y}_{ik}}{\partial \bar{w}}$$

which leads to the overall gradient

$$\frac{dJ}{d\bar{w}} = \sum_{i=1}^N (\hat{y}_i - y_i) x_i$$

## Neural Networks

We now have an understanding of the way logistic regression is used for binary classification. But how can this be generalized to other tasks?

Let a neuron be a thing that holds a number between 0 and 1. Let the value be the neuron's activation.

In each layer of our neural network, we have a series of neurons. The first layer of a neural network for a number classification task would be the pixel activations in the input image, and the last layer would be the activation of digits 0-10 as a series of neurons. Layers in between are hidden layers. Neurons in this layer detect higher-level patterns in the input image.

These layers are connected to each other with weights. The network connects every neuron in the previous layer to all neurons in the next. Thus, the activation of a given neuron is the weighted sum of the activations of all neurons in the previous layer. We call this structure the multilayer perceptron.



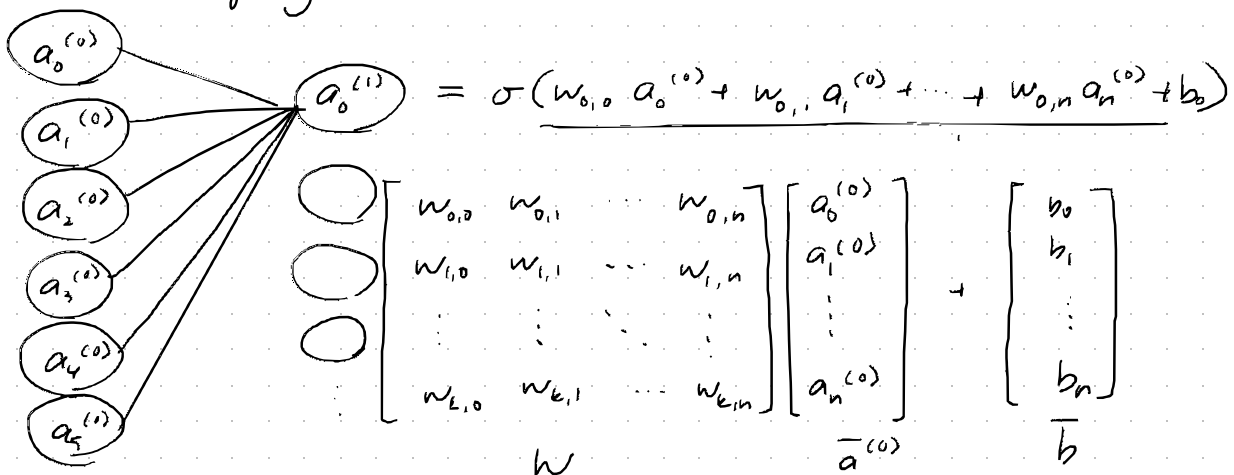
Given some activation  $w_1 a_1 + w_2 a_2 + w_3 a_3 + \dots + w_n a_n$ , where  $w_i$  are weights and  $a_i$  are activations of neurons from the previous layer, we put these through the sigmoid function, so that the value is squished into  $(0, 1)$ .

We also throw a bias term to the mix, some integer  $b$ , such that it ensures the neuron to activate meaningfully only when it reaches  $b$ . For example,

$$\sigma(w_1 a_1 + w_2 a_2 + \dots + w_n a_n - 10)$$

↳ only activate meaningfully when weighted sum  $> 10$ .

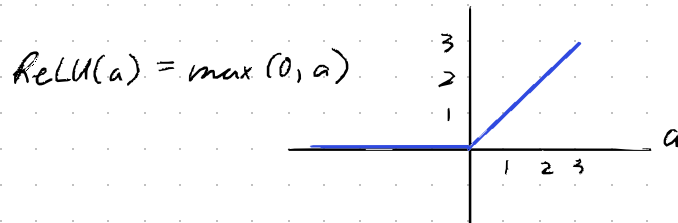
What this vanilla neural network allows is a scaling up of parameters for much more complex tasks. Since we have so many sigmoid functions to write at this scale, we simplify the notation.



So the new layer's activation simplifies to:

$$\bar{a}^{(1)} = \sigma(W\bar{a}^{(0)} + \bar{b})$$

A modern alternative of the sigmoid is the ReLU function.



Why? Because:

- ReLU avoids the vanishing gradient problem by keeping gradients high when  $x > 0$ .
- ReLU activates hidden layers that do not need to produce probabilities better

But these advantages cannot be understood without a scaled up gradient descent and backpropagation, which is how a neural network learns.

## NN gradient descent

We begin with a randomly initialized set of weights. Then, we make a pass with some input, and calculate the sum of square differences between the model's prediction, and the true answer.

If we then average all of these "costs" over lots of data, this value is representative of how well or how badly the model is doing.

To find the minima of this cost function  $C(w)$ , we can perform gradient descent. The general case of a 2-D gradient descent generalized to  $n$ -dimensions, fortunately. We are still computing the negative slope of the point  $w$  in the graph  $C(w)$  such that we can nudge ourselves to the local minima. Intentionally,

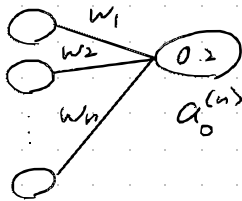
$$\text{if } \bar{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_n \end{bmatrix}, \text{ then } -\nabla C(\bar{w}) = \begin{bmatrix} 0.31 \\ 0.03 \\ -0.37 \\ 0.16 \end{bmatrix} \begin{array}{l} w_0 \text{ should increase} \\ w_1 \text{ should increase} \\ w_{n-1} \text{ should decrease} \end{array}$$

and the magnitudes tell us the extent of the decrease/increase.

# Backpropagation

NN gradient descent gives us how to calculate the direction towards a minima. But how do we use this to actually change the weights?

Say we have some neuron activation of 0.2 at the output layer, which represents the correct answer.



↳ Each neuron in the previous layer has an activation, which scales the activation of the neuron in the next layer by the weight. In other words, if we want to nudge the weights to change  $a_o^{(n)}$ , it must be in proportion to the weights.

So, each neuron in a given layer consults all neuron activations in the previous layer for the most bang-for-buck nudges that can be applied to the weights that connect them. Then, once we calculate the "sum nudge" for each neuron in the previous layer, that value then informs nudges that must occur in the previous layer. This process is **backpropagation**.

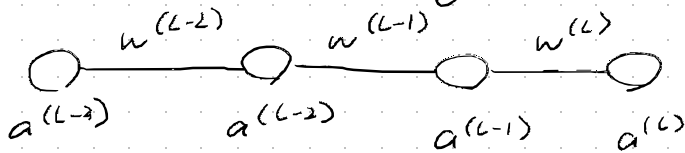
Since computing the gradients for each input would be too intensive, we perform gradient descent in batches.

This ensures the following:

1. Since the batch is more representative of the entire training corpus, it results in less noisy updates.
2. The gradient update is the average of all of the examples in the batch, which allows for parallel computation.

Now, what exactly is the calculus behind backprop?

Suppose a network with 4 layers with one neuron:



And a desired output  $y$

At the last layer, we have  $C_0(\dots) = (a^{(L)} - y)^2$

Let  $w^{(L)} a^{(L-1)} + b^{(L)} = z^{(L)}$

then,  $a^{(L)} = \sigma(z^{(L)})$

Since we want to calculate changes in cost with respect to weights, we can phrase the problem as such:

$$\frac{\partial C_0}{\partial w^{(L)}} \rightarrow \text{change in cost}$$

$$\frac{\partial C_0}{\partial w^{(L)}} \rightarrow \text{tiny nudges}$$

$$\frac{\partial C_0}{\partial w^{(L)}} = \frac{\partial z^{(L)}}{\partial w^{(L)}} \frac{\partial a^{(L)}}{\partial z^{(L)}} \frac{\partial C_0}{\partial a^{(L)}}$$

Chain rule

Given:

$$C_0 = (a^{(L)} - y)^2$$

$$z^{(L)} = w^{(L)} \cdot a^{(L-1)} + b^{(L)}$$

$$a^{(L)} = \sigma(z^{(L)})$$

$$\frac{\partial C_0}{\partial a^{(L)}} = 2(a^{(L)} - y)$$

$$\frac{\partial a^{(L)}}{\partial z^{(L)}} = \sigma'(z^{(L)})$$

$$\frac{\partial z^{(L)}}{\partial w^{(L)}} = a^{(L-1)}$$

$$\text{So, } \frac{\partial C_0}{\partial w^{(L)}} = a^{(L-1)} \cdot \sigma'(z^{(L)}) \cdot 2(a^{(L)} - y)$$

The full cost derivative is  $\frac{\partial C}{\partial w^{(L)}} = \frac{1}{n} \sum_{k=0}^{n-1} \frac{\partial C_k}{\partial w^{(L)}}$

Which is the cost function for  $w^{(L)}$  at the last layer.

$$\nabla C = \begin{bmatrix} \frac{\partial C}{\partial w^{(1)}} \\ \frac{\partial C}{\partial b^{(1)}} \\ \vdots \\ \frac{\partial C}{\partial w^{(L)}} \\ \frac{\partial C}{\partial b^{(L)}} \end{bmatrix}$$

$$\frac{\partial C_0}{\partial b^{(L)}} = \frac{\partial z^{(L)}}{\partial b^{(L)}} \frac{\partial a^{(L)}}{\partial z^{(L)}} \frac{\partial C_0}{\partial a^{(L)}}$$

$$= 1 \cdot \sigma'(z^{(L)}) \cdot 2(a^{(L)} - y)$$

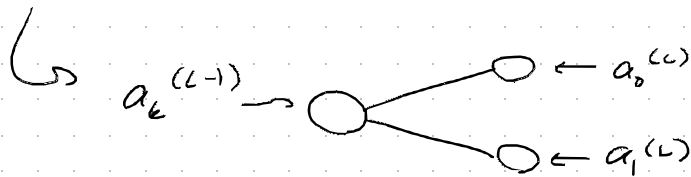
Furthermore,

$$\frac{\partial C_0}{\partial a^{(L-1)}} = \frac{\partial z^{(L)}}{\partial a^{(L-1)}} \cdot \frac{\partial a^{(L)}}{\partial z^{(L)}} \cdot \frac{\partial C_0}{\partial a^{(L)}}$$
$$= w^{(L)} \cdot \sigma'(z^{(L)}) \cdot z(a^{(L)} - y)$$

↳ the weights from the last layer affects the weights from the second last layer.

This one neuron layer example generalized to multi-neuron layers, only with additional indices.

$$\frac{\partial C_0}{\partial a_k^{(L-1)}} = \sum_{j=0}^{n_L-1} \frac{\partial z_j^{(L)}}{\partial a_k^{(L-1)}} \cdot \frac{\partial a_j^{(L)}}{\partial z_j} \cdot \frac{\partial C_0}{\partial a_j^{(L)}}$$



OK. How does all this get applied to the transformer?

## Large Language Models (and Transformers)

At a high level, we still follow the format of a multilayer perceptron.

input  $\rightarrow$  layers  $\rightarrow$  output.

But now, the input is generally scaled up from vectors to matrices and even  $n$ -dim arrays, which we call tensors. Thus the name "tensor-flow".

Let's take GPT-3 as an example of this scale.

Its  $\sim 175B$  weights are divided into 27,938 matrices, each of which perform linear transformation with input vectors for various functions:

- Embedding :  $d_{\text{embed}} \cdot n_{\text{vocab}}$
- Key :  $d_{\text{query}} \cdot d_{\text{embed}} \cdot n_{\text{heads}} \cdot n_{\text{layers}}$
- Query : "
- Value :  $d_{\text{value}} \cdot "$
- Output :  $d_{\text{embed}} \cdot d_{\text{value}} \cdot n_{\text{heads}} \cdot n_{\text{layers}}$
- Up-projection :  $n_{\text{neurons}} \cdot d_{\text{embed}} \cdot n_{\text{layers}}$
- Down-projection :  $d_{\text{embed}} \cdot n_{\text{neurons}} \cdot n_{\text{layers}}$
- Unembedding :  $n_{\text{vocab}} \cdot n_{\text{embed}}$



## Embedding (d-embed · d-vocab)

The embedding matrix turns tokens to vectors.

$$W_E = \begin{matrix} \text{Embedding} \\ \begin{bmatrix} v_1^{(1)} & & v_1^{(m)} \\ v_2^{(1)} & & v_2^{(m)} \\ \vdots & & \vdots \\ v_n^{(1)} & & v_n^{(m)} \end{bmatrix} \end{matrix} \quad \begin{matrix} \text{where } n = 1, 2, \dots, 288 \\ m = \text{number of tokens} \\ \text{in dictionary} \end{matrix}$$

vocabulary

And similar tokens in the embedding space are closer to one another, which can be calculated by the dot product.

↳ varies according to angle between the two vectors:

- $\theta > 90^\circ$ : negative
- $\theta = 90^\circ$ : 0
- $\theta < 90^\circ$ : positive

One example of populating an embedding matrix is the skip-gram model that is used in Word2Vec.

Input: one-hot encoded vector of the target word  
 $\in \mathbb{R}^V$

Embedding matrix: Matrix of size  $V \times d$ , where each row represents the  $d$ -dimensional embedding of some target word  $v$ .

Output: the embedding vector

Training: given the context of window size  $k$  that surrounds the target word in the training data, learn to predict context through backpropagation.

The LLM takes these embeddings as a starting point for each token in the input, and then further encodes the context of the input to each embedding.

At the very end, the last vector in the input will encode the context of the next token. Using an unembedding matrix, we turn that vector back to a probability distribution of the dictionary, so that we can make an informed prediction.

Since the initial probability distribution is not normalized, we run it through a softmax.

"Softmax" derives from the amplifying effects of the exp. on any maxima in the input vector.

$$\begin{array}{ccc} \begin{array}{c} \left[ \begin{array}{c} z_1 \\ z_2 \\ z_3 \end{array} \right] \\ \text{Logits} \end{array} & \rightarrow & \begin{array}{c} \left[ \begin{array}{c} \frac{e^{z_1}}{\sum e^{z_n}} \\ \frac{e^{z_2}}{\sum e^{z_n}} \\ \frac{e^{z_3}}{\sum e^{z_n}} \end{array} \right] \\ \text{Softmax} \end{array} & \rightarrow & \begin{array}{c} \left[ \begin{array}{c} 0.1 \\ 0.2 \\ 0.7 \end{array} \right] \\ \text{Probabilities} \end{array} \end{array}$$

We throw in some spice to the softmax function to change the softness of the max, so to say.

$$\frac{e^{z_i/T}}{\sum e^{z_n/T}}$$

The greater the temperature, the softer the max.

## Attention

As briefly mentioned previously, the word embeddings of each token in the dictionary is constant over all situations.

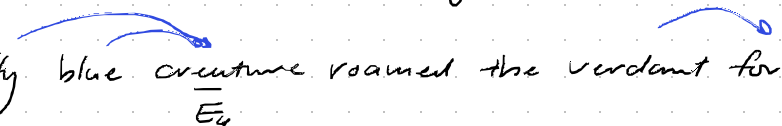
For example, the word "model" has the same context in these two contexts:

"The model walked on stage."

"The model makes next token predictions."

So the goal of the attention blocks in LLMs is to further contextualize vectors given input. More generally, it moves information to and fro across embeddings.

Working example: a fluffy blue creature roamed the verdant forest.



Let's think first about the encoding of adjectives to nouns. We call this one "head" of attention.

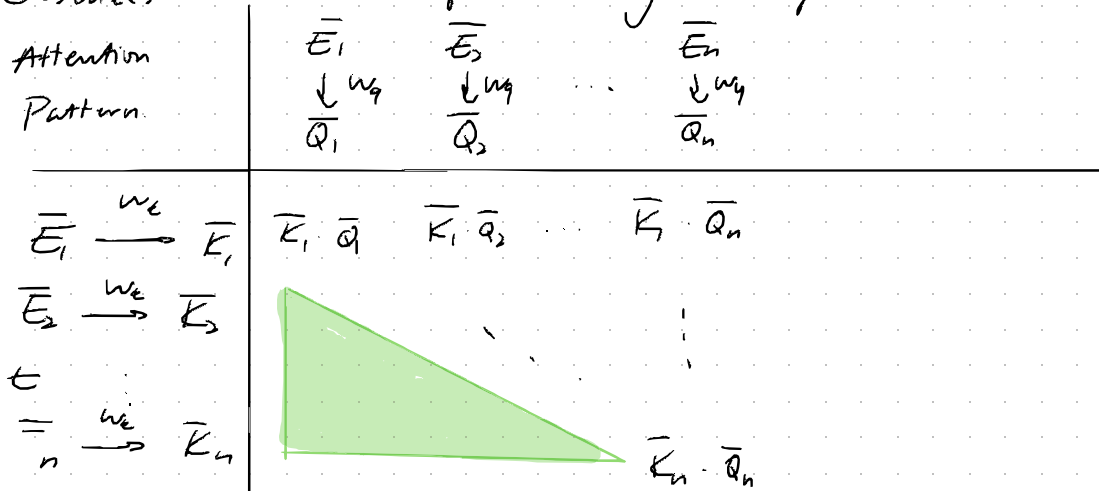
Imagine the noun "creature" asks the question, "are there any adjectives in front of me?" This is a query, which is represented as a vector. For example,  $W_Q \bar{E}_4 = \bar{Q}_4$ .

Imagine again that we can provide answers to such queries. This would take the form  $W_K \bar{E}_4 = \bar{K}_4$ .

Both  $W_Q$  and  $W_K$  are matrices with tunable parameters or weights, and both transform the embedding vector to some smaller space.

$\bar{K}$  and  $\bar{Q}$  "match" when they are closely aligned to one another, i.e., a positive dot product.

Dot products of these vectors can then be displayed in a matrix, where each value represents the closeness of their respective key-value pair.



When the dot products are great, we say that the respective  $\bar{K}$  "attends to" the respective  $\bar{Q}$ .

Let us then apply a softmax to these values at each column, or query. The resulting normalized values can be considered as weights that show the relevancy of the key and value. Formally put:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) V$$

for numerical stability

we will get back to this later.

One more training detail.

Given some input of length  $k$  for training where backprop is performed based on the gradient with respect to the true  $k-1^{\text{th}}$  true token prediction, we can make training much more efficient by creating subsequences of  $k-n$  from  $n=1 \rightarrow k-1$ .

the  $\rightarrow ?$

the fluffy  $\rightarrow ?$

the fluffy blue  $\rightarrow ?$

the fluffy blue creature roamed the verdant forest  $\rightarrow ?$

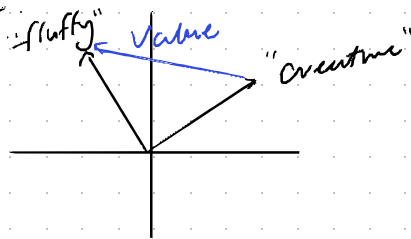
For this additional training step to work, we need the subsequences to be challenging, and not cheat by using the attention from subsequent tokens. Therefore, the area marked with  $\infty$  in the attention pattern must be masked at  $-\infty$  before softmax such that they become 0. This is called "masking".

Also note that the size of the attention pattern is context-size  $\times$  context-size, i.e., the permitted number of embeddings for each input.

How do we use attention to actually update the embeddings?

This is where  $W_v$ , the value matrix, comes in.

For example, say "creature" attends to "fluffy". By multiplying the  $W_v$  to "fluffy", we get a value vector that can be added to "creature".



To perform this computation at scale, we create the value matrix  $W_v$ .

	$\bar{E}_1$ ↓ $\bar{Q}_1$	$\bar{E}_2$ ↓ $\bar{Q}_2$	...	$\bar{E}_n$ ↓ $\bar{Q}_n$	
$\bar{E}_1 \rightarrow \bar{K}_1 \xrightarrow{W_v} \bar{V}_1$		$A_1 \bar{V}_1$			where $A_i = \bar{K}_i \bar{Q}_i$
$\bar{E}_2 \rightarrow \bar{K}_2 \xrightarrow{W_v} \bar{V}_2$		$A_2 \bar{V}_2$			
$\vdots$		$\vdots$			
$\bar{E}_n \rightarrow \bar{K}_n \xrightarrow{W_v} \bar{V}_n$		$A_n \bar{V}_n$			
		$\parallel$			
		$\Delta \bar{E}_2$			

$$\text{so } \bar{E}_2 + \Delta \bar{E}_2 = \bar{E}_2'$$

This is one "head" of attention. Transformers usually have many more. These ask different questions with different  $W_q$ , different answers with  $W_k$ , and different deltas with  $W_v$ . GPT-3 uses 96 attention heads, for example.

Lastly, the output matrix. But a small detail before discussion.

The  $W_v$  matrix maps one embedding vector to another, meaning that it has dimensions  $n\_embed \times n\_embed$ . This is much larger than the matrices  $W_k$  and  $W_q$ , which are both size  $n\_embed \times n\_query\_key$ .

$W_k$  and  $W_q$  maps the embeddings to a much smaller query/key space. In GPT-3,  $n\_embed = 12,288$  and  $n\_query\_key = 128$ .

Since we don't want to use a disproportionate amount of compute for calculating values, we perform a low-rank factorization of  $W_v$  such that we have  $W_{vd}$  and  $W_{vp}$ .  $v_d$  represents downprojecting the embedding space to the QK space, and then upprojecting the QK space to the embedding space.

The output matrix is a horizontally stapled-together version of the  $W_{vp}$  matrices, which allows us to only refer to  $W_{vd}$  when we say "value matrix".

We now covered attention. But that is half of the transformer, we should now discuss the multilayer perceptron layers.

## Multilayer Perceptron

There has been increasing evidence that the MLP is what stores facts in an LLM.

Before we begin, let's introduce a toy example, and some assumptions to go with it.

We want to understand how an LLM might "know" that Michael Jordan plays basketball. Let's assume that a direction in the embedding space of tokens encodes:

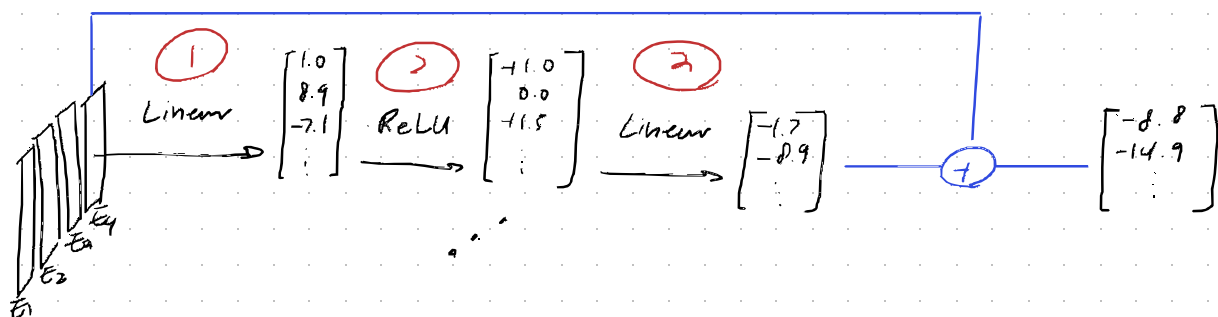
1. First name Michael
2. Last name Jordan
3. Basketball.

So, the dot product between some embedding  $\vec{E}$  and say, Basketball reveals whether the embedding represents basketball or not. As an example,

if  $\vec{E} \cdot \text{Last name Jordan} = \vec{E} \cdot \text{First name Michael}$ , we would know that the resulting vector encodes both first name Michael and last name Jordan.



At a high level, an MLP block does the following:



1. This is a matrix multiplication with some matrix.  
 Say, for example, that this matrix looks as such:

$$M = \begin{bmatrix} \text{---} \bar{R}_0 \text{---} \\ \text{---} \bar{R}_1 \text{---} \\ \vdots \\ \text{---} \bar{R}_n \text{---} \end{bmatrix}$$

Suppose further that  $\bar{R}_0$  encodes the direction F. N. Michael.  
 Then, the resulting vector will calculate the proximity of the embedding to  $\bar{R}_0$ .

$$M \cdot \bar{E} = \begin{bmatrix} \bar{R}_0 \cdot \bar{E} \\ \bar{R}_1 \cdot \bar{E} \\ \vdots \end{bmatrix} \begin{cases} \approx 1 \text{ if } \bar{E} \text{ encodes "FN Michael"} \\ \leq 0 \text{ if not.} \end{cases}$$

So each row in  $M$  is probing, with some "question", and calculating its proximity to  $\bar{E}$ .

There is also the bias vector, which adds or subtracts the final output of 1.

Since  $M$  maps  $E$  to a higher dim., we will call it  $W_p$ .  
 Let bias be  $\bar{B}_p$ .

After the linear transformation, we go through a ReLU, <sup>(2)</sup>  
a non-linear function that turns all negative values  
to 0, and leaves all positive values unchanged.

Why? Because

(1) It classifies whether the embedding embeds  $\bar{R}_i$  or  
not as a simple AND gate, on or off.  $\Rightarrow$  ~~D~~

(2) Introducing nonlinearity to the network allows it  
to encode nonlinear relationships in language.

In detail,

$$\bar{E} \cdot \bar{R}_i \begin{cases} \approx 1 & \text{if } \bar{E} \text{ embeds } \bar{R}_i \\ \leq 0 & \text{otherwise} \end{cases} \xrightarrow{\text{ReLU}} = 0 \text{ otherwise}$$

More simply, think of the final dot product as a neuron.

Since it is either activated or deactivated, we don't want  
negative values.

(3) The final linear function downprojects the final  
dimensions back to that of the embedding space.

$$\underbrace{\begin{bmatrix} -c_1- \\ -c_2- \end{bmatrix}}_{n_1 c_1 + n_2 c_2 + \dots} \cdot \begin{bmatrix} n_1 \\ \vdots \\ n_m \end{bmatrix} + \bar{B} = \dots$$

In other words, the first linear detects topics in the  
embedding and turns those sensitivities to neurons.

Now that we have neurons encoding F.N. Michael and L.N. Jordan, the relevant neurons can either activate or deactivate concepts in  $\bar{C}_0 - \bar{C}_n$ . For example, the vector that encodes Basketball + Chicago Bulls + ..., which results in the final embedding.

In summary, the vectors  $\bar{R}_0 - \bar{R}_{n-1}$  in  $W_p$  represents directions in the embedding space that will activate the output neurons if they correspond to  $\bar{E}_i$ . Then, the neurons go through ReLU to add non-linearity and simplicity. Finally, the column vectors  $\bar{C}_0 - \bar{C}_{n-1}$  represent vectors that will be added to the result, given the active neurons.

Symbolically,  $\bar{R}_0$  and  $\bar{R}_1$  may detect  $\overrightarrow{F.N.M.}$  and  $\overrightarrow{L.N.M.}$ , and  $\bar{C}_0$  may be activated by  $v_0$  to encode basketball.