

# INTUITIONS WHILE BUILDING A LANGUAGE MODEL

## USING BIGRAMS

24:10

→ Follow probabilities and sample

$$\rightarrow \boxed{\cdot + w + \cdot}$$

Matrix

→ We created a  $27 \times 27$  matrix with rows and columns that have the frequencies of each bigram occurring.

→  $27 \times 27$  Matrix

→ convert them to probabilities by creating a probability vector

→ remember to convert the frequencies into floats because we will normalize them.

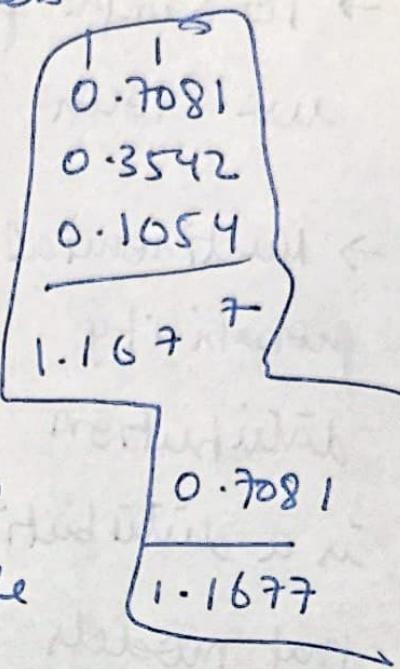
$$\text{P} \leftarrow \frac{1}{\text{P.sum}} \quad \text{P.sum} = \sum_{i=1}^{27^2} \text{P}_{ii}$$

→ "Sampling" is the process of generating new data points (e.g. words, sentences or images) based on a learned probability distribution.   
 →  $\downarrow$  Probability of any single character to be the first character of the word.

- ① Deterministic Sampling
- ② Random sampling

- To sample from the distributions we usually use Tech. Multinomial
- ↳ returns sample from the multinomial probability distribution: you give me probabilities and I will give you integers that are sampled according to their Prob. dist.
- Multinomial probability distribution is a distribution that models multiple outcomes, each with its own probability.
- In language model, each token in the vocabulary can be an outcome. no. of possible outcomes = "k"
- each outcome  $i$  has a fixed probability  $p_i$ , where
- $$\sum_{i=1}^k p_i = 1 \quad (\text{probabilities must sum to 1})$$
- The Probability Mass Function (PMF) of the multinomial distribution is:
- $$P(X_1 = x_1, X_2 = x_2 \dots X_k = x_k) = \frac{n!}{x_1! x_2! \dots x_k!} \prod_{i=1}^k p_i^{x_i}$$
- $n$  = no. of trials       $x_i$ : count of outcome  $i$
- $k$  = No. of outcomes       $p_i$ : Probability of outcome  $i$
- Total counts equals no. of trials

- Turn generator related randomness
- replacement = True? why?  
 When we delete an element,  
 we can draw it and the we  
 can put it back into the list  
 of eligible indices to draw again  
 (ref to 29:31 and see the code  
 for more details)
- for efficiency purpose, prepare a matrix  $P$ ,  
 that will have the probabilities ~~int~~ in it.
- Same as Frequency Matrix ( $N$ ) but every  
 single row will have, the row of probabilities  
 that is normalized to 1 indicating the  
 probability distribution for the next character  
 given the character before it as defined by  
 which row we're in.
- keepdim = True? why?  
 - It is used to maintain shape of the tensor  
 for tensor operations in future.



- Can we divide a tensor of  $27, 27$  by  ~~$27, 1$~~   $27, 1$ ?
- BROADCASTING SEMANTICS ↳ It is fine (Broadcastable)

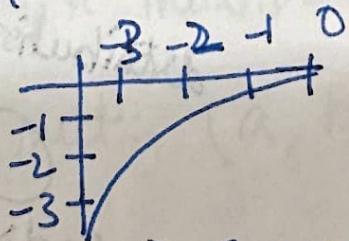
TWO TENSORS CAN BE combined in a binary operation like division when

- each tensor has at least 1 dimension
- when it's taking over the dimension sizes, starting at the trailing dimension, the dimension size must either be equal or one of them is 1, or one of them does not exist.

### [CROSSING BROADCASTING SEMANTICS] \*\*

- ~~Note~~, quality of model in a single number? LOK?

- Maximum likelihood estimation
  - ↓ Product of all the probabilities is higher, the better
- We use log ~~prob~~ likelihood, which lower the probability, lower the log likelihood. e.



- Log likelihood is the sum of all probabilities  
 $\log(a+b+c) > \log a + \log b + \log c$

→ negative log likelihood =  $-(\text{log likelihood})$   
(NLL)

→ normalized log likelihood is "NLL/m".  
(the better the better)

~~before the gathering~~

- GOAL: Minimize likelihood w.r.t the model parameters (statistic modeling)

- equiv. to minimizing log likelihood  
(monotonic)

- equiv. to minimizing the negative log likelihood

- equiv. to minimizing the average negative log likelihood.

2

→ NEURAL NETWORKS APPROACH some weights ( $w$ )

• receives single character as I/P  $\rightarrow$  N.N

→ take the training set  
( $x, y$ )

Probability o/p next distribution character (guess)

gradient based optimization

$x$  = inputs

$y$  = labels

→ difference b/w torch.Tensor + torch.Tensor

①

②

→ ① data type is int64      ② data type is float32

→ How all we going to feed in those examples into Neural Network?

~ common way of encoding integers is called one-hot encoding

→ one-hot encoding converts an integer into a long vector of zeros except the integer dimension which is converted into "1". and then that vector will be fed into neural network.

"TORCH.NN.FUNCTIONAL.ONE-HOT" does this

↳ num-classes is how long do you want your vector to be

→ Also make sure the dtype = float32, ~~int64~~ do a .float()

→ @ = matrix multiplication operator

→ matrix mult:  $(a, b) \xrightarrow{\text{b}, c} \text{o/p}(a, c)$  matrix

→ each neuron does a simple operation  $wx + b$

(here  $\otimes W$ )  $(3, 13) \xrightarrow{\text{column}} \text{neuron}$

$0.5 \rightarrow \text{fire rating} \Rightarrow$  By dot product

Writing a probability matrix for N·N

- Now we are going to interpret these 27 numbers is that these 27 numbers are giving us 108 counts.
- To get counts we will get log counts and exponentiate them.

→ logits = log counts

counts = probability of character appearing

↳ logits · exp()

prob prob = counts / (counts.sum(1, keepdims=True))

prob This is basically softmax

O/P LAYER      SUMMARY ACTIVATION FUNCTION      PROB

$$\begin{bmatrix} 1.3 \\ 5.1 \\ 2.2 \\ 0.7 \\ 1.1 \end{bmatrix} \rightarrow \boxed{\frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}} \rightarrow \begin{bmatrix} 0.02 \\ 0.90 \\ 0.05 \\ 0.01 \\ 0.02 \end{bmatrix}$$

- Now to the last part we need to calculate the loss.

loss = - prob [torch.arange(n), ys].log().mean()  
 n = no. of binararies

→ In pytorch the best practice to set initial gradient to zero is by,

w.grad = None

[1:40:50]

→ Tells us the influence of that weight on loss function.

→ look into  $[\text{w}.data + -50 * w.grad]$

Note 1:

→ Weights in this are basically the same as count matrix and are log counts

→  $w.\exp()$  is count matrix

→ [go through Note 2]

## BUILDING MASTMOP PT2 : MLP

→ we are going to be building a multilayer perceptron using the research paper Bengio et al. 2003.

→ we are building a character level language model, in the paper they have a vocabulary of 17,000 words and they built a word level. we are using the same intuition to build a character level.

## # What happens in the paper?

→ They take every one of the 17000 words  
and they are going to associate to each  
word a 30 dim feature vector

(~~every word~~ is 17000 words  
in 30 dimensional vector  
space)

→ Turn the squad out embeddings  
using back propagation.

(The vectors will move  
around in space,  
the synonyms or matching  
words might end up near  
in the space)

→ Use MLP to predict words. using few words.

→ minimize the log-likelihood of training  
data

## # build dataset

- block size is the context length: how many  
characters do we take to predict the next one

-  $x, y = [], []$

↳ labels (The next character to predict)

~~blocks~~

→ refer to 11:44

logic:

$\text{stoi} = \{\text{'!': 0}, \text{'a': 1}, \text{'c': 2}, \text{'t': 6}\}$

$\text{itos} = \{0: '\!', 1: 'a', 2: 'c', 6: 't'\}$

initial:  $[0, 0, 0]$  (all 0)

characte 'c'

$ix = \text{stoi}[\text{ch}] = 2$

Add pair:  $[0, 0, 0] \rightarrow [0, 0, 2]$

update:  $[0, 0, 2]$

characte 'a'

$ix = \text{stoi}[\text{ch}] = 1$

Add pair:  $[0, 0, 2] \rightarrow [0, 0, 2]$

update:  $[0, 2, 1]$

characte 't'

$ix = \text{stoi}[\text{ch}] = 6$

Add pair:  $[0, 2, 1] \rightarrow [0, 2, 1, 6]$

update:  $[2, 1, 6]$

characte '

$ix = \text{stoi}[\text{ch}] = 0$

Add pair:  $[2, 1, 6] \rightarrow [2, 1, 6, 0]$

update:  $[1, 6, 0]$

perm for "cat."

$X: [0, 0, 0], [0, 0, 2], [0, 2, 1], [2, 1, 6]$

$Y: [2, 1, 6, 0]$

note: predict the next character using the 3 past  
characters ~~con~~ (context).

Goal: Building a lookup table that will take  $X \in \{32, 3\}$  and predict  $i$ -th output  $\hat{y}$ .  
 We are creating a look matrix called  $C$  whose rows are the words.  
 [Building the lookup table]  
 we are going to represent in 2-D vector space.  $(27, 2)$

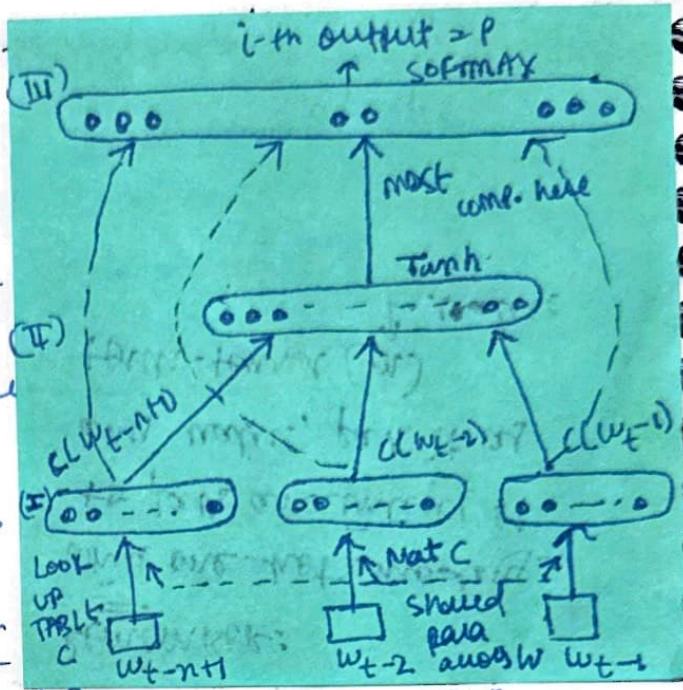
~~27~~ characters will have a 2

→ The goal here is to create like a 2D vector.

We have created a look up

embedded matrix  $= C = C$

## IMPLEMENTING THE HIDDEN LAYER



3, 2-D embeddings as shown (I)

→ Initialize  $W_1$  weights randomly  
 number of inputs will be 6

number of neurons 100

bias will be 100 of them

→ for neuron,  $\oplus$  in hidden layer the operation  
 should be  $\text{emb} @ W_1 + b_1 \oplus$  → This is not  
 possible because,

emb matrix shape =  $[32, 3, 2]$

$W_1$  shape =  $((6, 100))$  # no matching dims

Then how? convert  $([32, 3, 2]) \rightarrow [32, 6]$

to do the mat mul.

- we fetch . cat ? (concatenate the 3 embeddings in 'I' (refer to stick note))
- + embedded the 3 ; 2D spaces.
- ⇒ torch.cat ([emb[:, 0, :], emb[:, 1, :], emb[:, 2, :]]),  
 1) . shape  
 ↗ give the 1 dimension
- ⇒ output: torch.size ([32, 6])
- But what if Block-size > 5?  
 ↗ sliced output  
 ↗ TORCH. UNBIND  
 ↗ efficient way ↓ ↗ a.view (look into this)  
 ↗ in the computer memory  
 any tensor is stored  
 in 1D.  
 ↗ take the matrix that  
 you want to connect and give the  
 dims.  
 ex: emb.view (32, 6).
- h = ~~emb~~ torch.tanh (emb.view (emb.shape[0], 6))  
 ↗ @ w, + b,  
 ↗ here we are implementing the II layer  
 with tanh function as activation  
 function. (the dimensions all adjusted  
 now)  
 ↗ shape: [32, 6] ↗

## IMPLEMENTING THE OUTPUT LAYER

→ here we will take 2<sup>nd</sup> set of weights and bias  
where  $w_2 = \text{torch.rand}((100, 27))$   
 $b_2 = \text{torch.rand}(27)$

↑  
input from prev layers      same

→ output will be

$$\text{logits} = h @ w_2 + b_2$$

logits ≠ output

$$\hookrightarrow \text{shape} \rightarrow [32, 27]$$

## IMPLEMENTING THE NEGATIVE LOG LIKELIHOOD LOSS

① exponentiate the logs [counts]

across 1 dimension

② ~~figure~~ figure out the probability

↑

$$\Rightarrow \text{counts} / (\text{counts} \cdot \text{sum}(\text{counts}))$$

ACTUAL LAYER

(keeping = True)

③ we have layer 'y' which we created during  
dataset preparation which has all the  
labels.

→ we will move inlets into the rows of 'prob',  
and each row we will pluck out the prob assigned  
to the correct character.

• we will do this by first taking an iterator which  
forch-arrange(32) → this will initialize 32  
numbers which we will use to as row values.  
and then, to pluck out the 'Y' labels which  
are integers we will do,

\* prob [forch-arrange(32), Y]

This will pluck out all the 'Y' probabilities.

→ Now, to get the loss (negative log loss likelihood)

$$(3) \text{ Loss} = -\underbrace{\text{prob}(\text{forch-arrange}(32), Y)}_{①} \cdot \log(1 - \text{mean})$$

INTRODUCING F. CROSS-ENTROPY (to count out loss  
calculation).

so all of it can be packaged into

\* f. cross-entropy (logits, Y)

why ①? ① more efficient forward and backward  
passes

② much more numerically well behaved.

## IMPLEMENTING THE TRAINING LOOP, OVERFITTING ON

### BATCH:

- Refer to code in the notebook (gnash make more rept2.ipynb)  
there's an ~~note~~ Notebook named (~~the~~ TRAINING LOOP)
- We are overfitting here because the number of parameters (3481) is higher than the samples (32).
- The reason we did not add the ~~size~~ ~~batch size~~ ~~size~~ ~~batch size~~

## TRAINING ON THE FULL DATASET, MINIBATCHES

- Perform FWD, BACK and update on mini batches of data.
- We will randomly select some portion of the dataset and that's a minibatch - We will do operations on that and iterate on batches.
- To ~~select~~ the learning rate, the perfect one we use torch · lin space to initiate a set

of learning rates and then exponentiate them.

$$lre \rightarrow \text{torch.linspace}(-3, 0, 1000)$$

$$lrs = 10^{lre}$$

→ We pass these to the iterator loop and append and store these values in arrays.

$$x_i = []$$

$$loss_i = []$$

→ we then refit to the graph and select a proper learning rate depending on where the loss is low and stable.

→ at the end, we decay the learning rate by exponentiate it (e.g.: if lr was 0.1, we do it  $0.01$ ) to get the best ~~loss~~ and lowest loss.

## SPLITTING UP THE DATASET INTO TRAIN/VAL/TEST SPLITS AND WHY

## BUILDING MORE PT3: ACTIVATIONS + GRADIENTS,

### BATCH NORM:

#### ① FIXING THE INITIAL LOSS:

- In training of NN, it is almost always the case that you will have a rough idea for what loss to expect at initialization. and that depends on the loss func<sup>n</sup> + the problem setup.
- Usually if the loss is unpredictably high or low or abnormal, its because the logits have extreme values.

#### ② FIXING THE SATURATED FUNCTION (tanh)

- Tanh is a squashing function, it takes arbitrary numbers and squashes b/w -1 + 1.

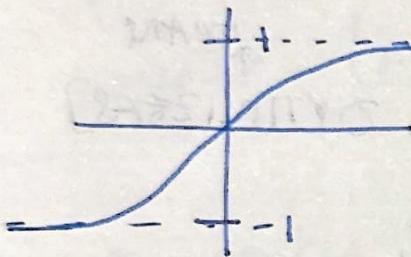
what happens to tanh during backprop?

- remember, if your distribution tanh has extreme values, like more of -1's and 1's, you are cooked.

- During the backpropagation of tanh, it has to go through the following  
 $\text{self\_grad} \doteq (1 - \tanh^2) \circ \text{out\_grad}$

when the values are extreme i.e. 1 or -1, the gradient will be completely '0', leading to vanishing gradients.

- if it is the fan in is more of '0', then the gradients directly pass through with their original values.
- The more you are in the flat tails (-1 or 1) ~~if~~ then the more the gradient gets squashed / destroyed.



- Same is true for sigmoid, ReLU.

### ③ CALCULATING THE UNIT SCALE: "KAIMING UNIT"

- until now we have fixed and + tweaked the weights and biases to limit the losses and get better results. We have done it manually.

- standard deviation is the measure of the spread of gaussian.

- How do we initialize weights so that our activations take on reasonable values throughout the network? <sup>DEEP</sup> <sub>KAIMING</sub> INTO RECTIFIERS

→ In short, for perfect initialization, the weights and biases are divided by  $\sqrt{2/n_e}$ .

→ This is implemented in pytorch.

`torch.nn.init.Kaiming-normal-`

↓  
most common way of  
initializing N.N.

[RESIDUAL CONNECTIONS, NORMALIZATIONS, OPTIMIZERS]

↓  
Make sure Initialization  
are done properly

→  $\text{std} = \frac{\text{gain}}{\sqrt{\text{fan\_mode}}}$  → we have gain because in tanh and other functions there is squashing, to maintain the shape we have gain.

gain for tanh =  $5/3$  (ref to docs pytorch)

#### ④ BATCH NORMALIZATION (2015)

→ we don't want the pre-activation func's to be too small because tanh will not do anything and at the same time we don't want it to be too large because tanh will be saturated.

→ We want them to be roughly gaussian ( $\text{mean} = 0$ ,  $\text{std} \cdot \text{deviation} \approx 1$ ) at initialization.

→ batch normalization: take hidden states and normalize them to be gaussian.

→ we will take the hidden activations, calculate their mean and standard deviations and divide them per batch.

$hp_{react}$  = hidden activations.

$hp_{react} = \frac{hp_{react} - hp_{react}.\text{mean}(0, \text{keepdim} = \text{True})}{(hp_{react}.\text{std}(0, \text{keepdim} = \text{True}))} \quad (T)$

→ remember that we would want the hidden activations to be gaussian only at ~~a~~ initialization, we don't want them to be forced to be gaussian always. we would like the Neural Net to ~~be~~ move this around to potentially make it more diffuse / sharp and we'd like the backprop to tell us how the distribution should move around.

→ to do that, we will take the normalized i/p and we are scaling them by some gain and offsetting them by bias. from (I), we will fill  $bn_{gain} = \text{torch.ones}(c, n\_hidden)$   $bn_{bias} = \text{torch.ones}(c, n\_hidden)$

$hp_{react} = bn_{gain} * xi + bn_{bias}$  (ref to research paper)

- take Layer and append Batch normalization to control scale of operations / activations at every point of Neural Net.
- will provide B.N layers across all the layers.
- LINEAR LAYER  
↓  
BATCH NORM LAYER  
↓  
ACTIVATIONS

[Tech. NN. LINER,  
BATCH NORM 2D]

## PYTORCHIFY

NN.

→ CREATE A LINEAR LAYER (TORCH.NN) → LINEAR

↳ FAN IN: NO. OF INPUTS  $\rightarrow$  FEATURES  
 FAN OUT: NO. OF OUTPUTS  $\rightarrow$  FEATURES ↓  
 what is  $X$ ? Kaiming init

$$y = zA^T + b$$

bias initialized to zeros (initially)

↳ BATCH NORM will be done (BATCH NORM 1D)

↳  $wx+b$  (--- call ---)

$$y = \frac{m - E[x]}{\sqrt{\text{Var}[x] + \epsilon}} * \gamma + \beta \rightarrow \text{Batch normalization}$$

torch.nn.BatchNorm1d (num\_features, eps, momentum,  
 affine=True, track\_running\_stats=True,  
 device='none', dtype=None)

(we won't track  $\gamma$  &  $\beta$  after normalization)

\* self.training = True, should be used

so that while influence is running,  
 it will be torch.no\_grad works.

→  ~~$\gamma$~~   $\gamma$  = Gamma (learnable scaling parameter)

$\beta$  = Beta (learnable shifting parameter)

$\rightarrow \gamma$  = allows the network to scale the normalized output back to a desired range.

$\beta$  = allows the network to shift the normalized output to a desired mean.

$\gamma > 1$  (amplified)

$0 < \gamma < 1$  (compressed)

$\gamma < 0$  (inverted)

$\beta > 0$  shifts data to the right ( $\uparrow$  Mean)

$\beta < 0$  shifts data to the left ( $\downarrow$  Mean)

$$x_{\text{scaled}} = n \cdot \gamma$$

$$x_{\text{shifted}} = n + \beta$$

Batch Normalization, after Normalization the i/p to have mean = 0, variance = 1, we apply shifting + scaling.

$$y = x_{\text{normalized}} \cdot \gamma + \beta$$

$\rightarrow$  BN is a technique used to normalize the activations of neurons within a layer during training.

- Improve training speed, stability
- convergence

## BUILDING WANGET (PART 5)

→ lossi : is a python list of floats (losses)

① Embedding

→ ↳ we

② Flatten

### Reimplementing | Refactoring Forward Pass

→ embedding table is outside  
the layers

→ The emb-view is also outside  
the layers.

→ we will create these layers and  
add them to layers list.

→  $emb \sim \text{emb} = \langle [x_b] \rangle$  : embedding  
operation

write  
it in  
the layers:

$x \sim \text{emb-view} [\text{emb\_shape}[0], -1]$   
: flattening  
operation

class embedding:

→ This is where embedding happens, remember that the  
look up table embeddings are also trainable parameters

→ If a few words and they appear in similar  
contexts. During training, the embedding vectors  
of 'id' & 'ë' might move closer to each other  
in embedding space because model minimizes  
the loss.

Remember  $\frac{\partial L}{\partial \theta}$  there are multiple weights being upgraded,

- ① Linear layer weights
- ② batch Normalization parameters (gamma, beta)
- ③ weight inactivation layers
- ④ embedding weights

How are these weights updated?



- F-W pass
  - forward through layers linear, batchnorm, activations)
  - loss computation
    - o/p is compared to target
    - loss function used to adjust weights

→ Backward pass

- gradients of the loss are calculated
- each weight is adjusted

→ Significance of EACH weight

Layer layer weights:

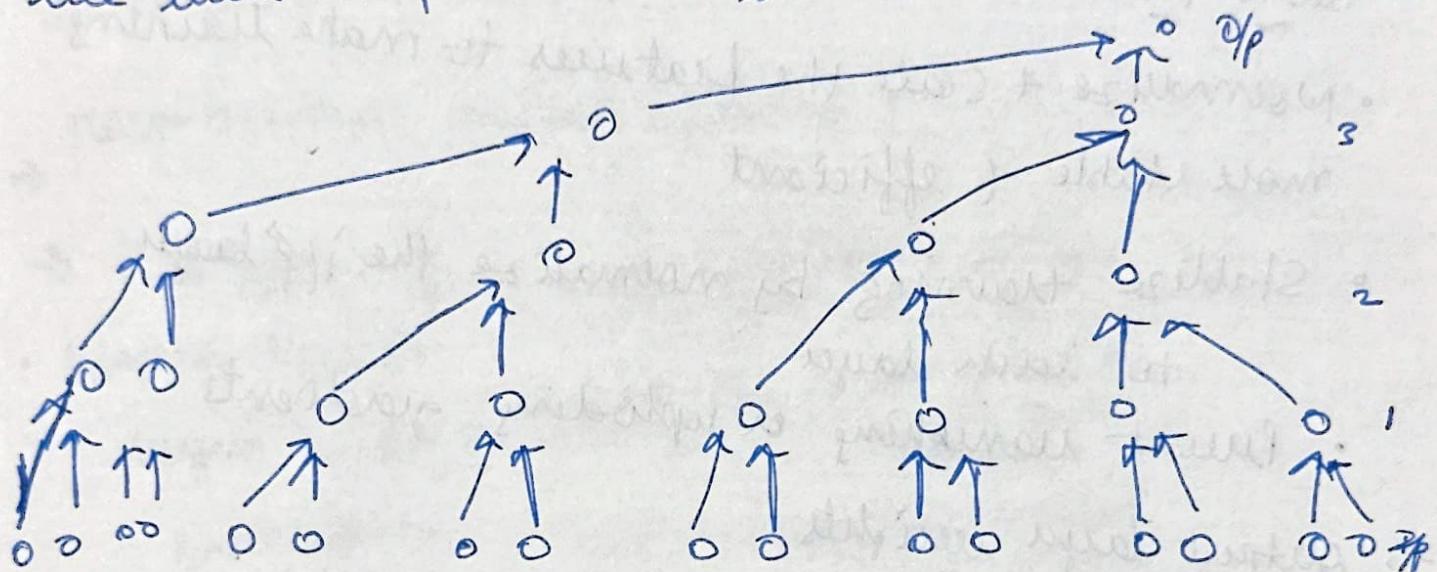
P-T-O - -

- ① Learn to optimize all the weights

- embedding weights
    - convert tokens into dense vectors that represent their meaning
  - linear weights
    - Transform these embeddings into richer
  - Batch Norm
    - Normalize & scale the features to make training more stable & efficient
    - Stabilize training by normalize the ~~length~~  
to each layer
    - Prevent vanishing or exploding gradients
  - output layer weights
    - combine features to predict the next token in sequence.
- 
- CONTAINERS ARE WAY OF ORGANIZING layers into lists or dicts.
  - SEQUENTIAL maintains a list of layers in a sequential order.
    - output of one layer is input to next.
    - Instead of <sup>defining +</sup> connecting each layer manually we do this.

## SCALE UP THE NEURAL NETWORK:

- we are using Bengio's et paper.
- In Bengio's et paper there is a single large where everything is squashed. Instead of that now we will implement a mannet like structure.



BLURMS

Dilated causal convolution  
layers.

4-blurms

8-blurms

(repulsive fusion)

+

2 consecutive squashing