

→ More - 3

Problems in the net:

①

① large initial loss:

$$\text{fix: } \begin{cases} w_2 \approx 0.01 \\ b_2 = 0 \end{cases} \quad \left. \begin{array}{l} \text{to} \\ \text{be} \end{array} \right.$$

last layer:

the logits that come out of the last layer

$$\text{eg: } [0.1, 2+0.5, 7.2, 3.2]$$

$$\left. \begin{array}{l} \text{softmax} \\ \text{saturation} \\ \text{ReLU} \end{array} \right\} \rightarrow \left[10^{-1}, 10^7, 10^2, 10^3 \right]$$

∴ we need logits closer to 0.

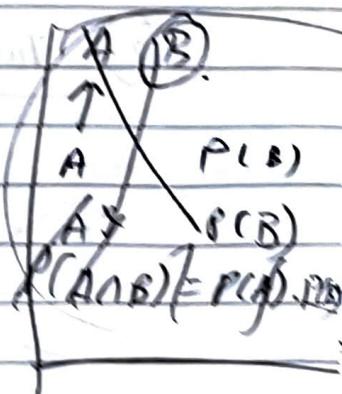
$$\therefore \begin{aligned} w_2 &\rightarrow \text{near 0} \\ b_2 &\rightarrow 0 \end{aligned}$$

∴ prob → better and more neutral.

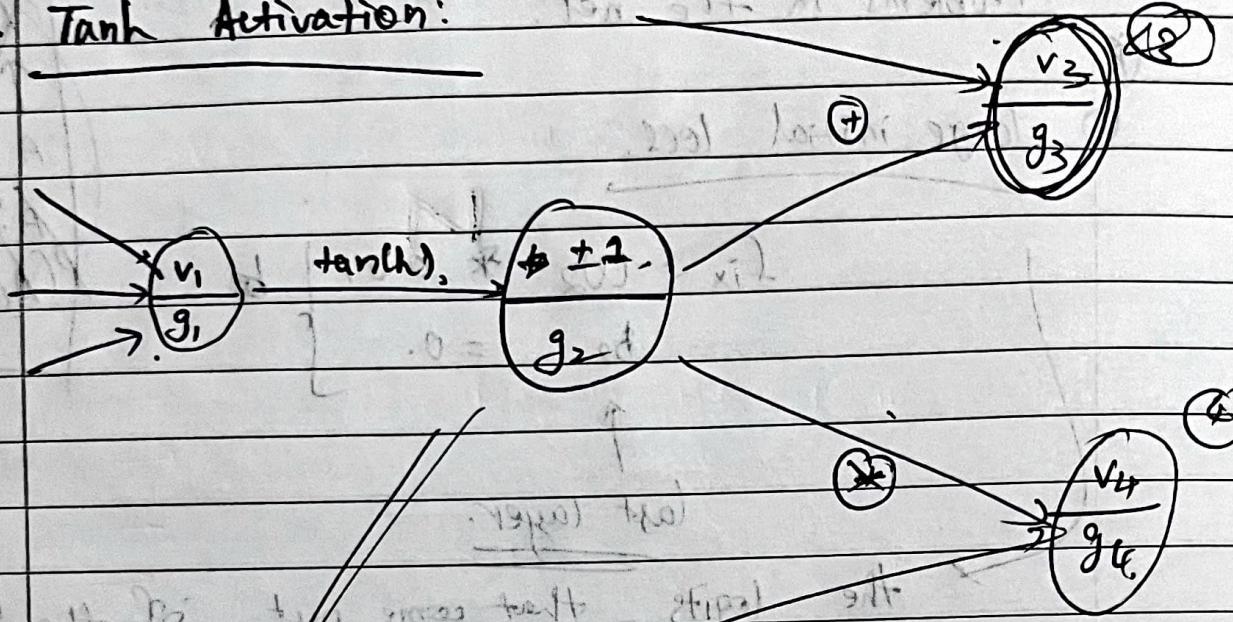
w_2 cannot

be 0

→ symmetry problem



② Tanh Activation:



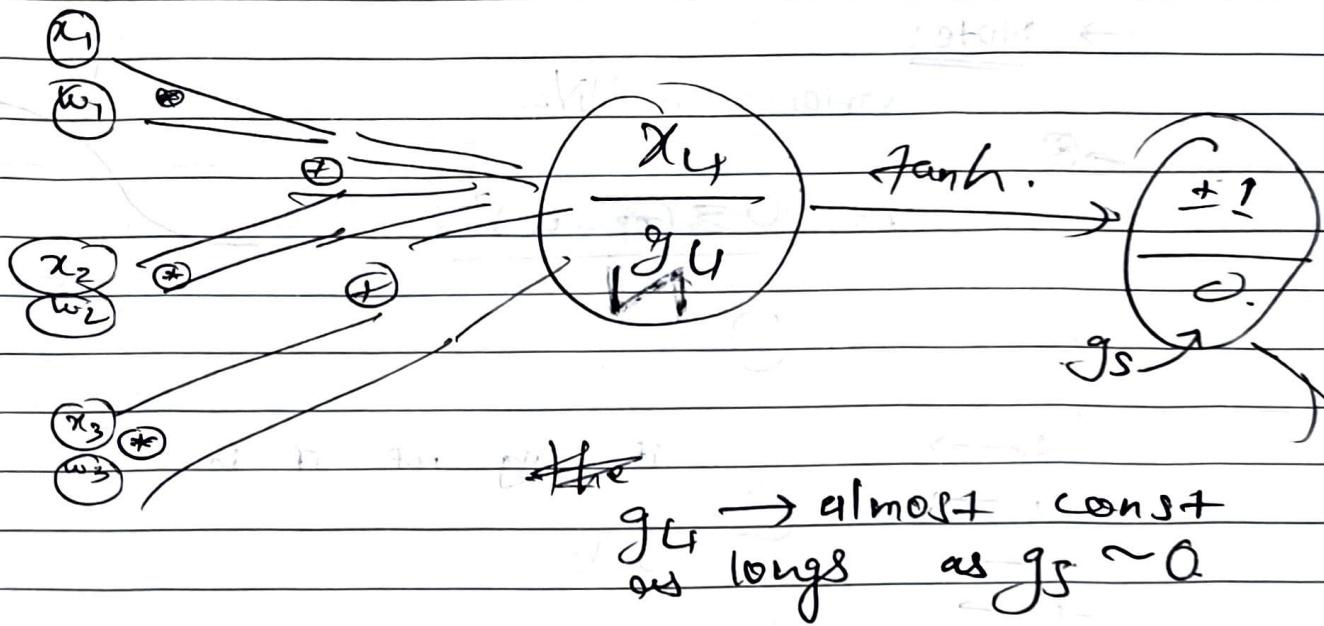
$\therefore \text{Backprop thru } \tanh(z)$

$$g_2 \leftarrow g_2 + \frac{\partial L}{\partial \tanh(z)} \cdot \frac{\partial \tanh(z)}{\partial g_2}$$

this is not a weight
 \because the data only depends on prev
node
we need

Whenever the gradient goes back thru this
node, it only adds 0.

\therefore the only way this ever changes if
 x changes via other activation
then $\tanh(x)$ changes ...



∴ the only way for it to come out of saturation is for x_i, w_i ($i = 1, 2, 3$)

to change greatly via other backprop
thus changing $x_4 \rightarrow$ bringing

$\tanh(x_4)$ closer to 0 \rightarrow gradient $g_5 \neq 0$

(post activation)
and the neuron stops being

dead / Saturation.

Fix:
~~Five scale~~ scale down w_j (hidden layer)
and b_j ,

$$w_2 * = 0.2$$

b_1 also tone it down.

→ Note:

variance in NNs

$$V = \frac{1}{N} \sum (x_i - \bar{x})^2$$

N

1. → ∵ if avg not of in

$$\begin{array}{c} \text{1} \rightarrow \\ \text{2} \rightarrow \\ \text{3} \rightarrow \\ \text{4} \rightarrow \\ \text{5} \rightarrow \end{array}$$

2nd row
1st row
3rd row
4th row
5th row

(**) Variance in a Dataset:

We use standardized inputs with variance of 1 to make all types of inputs of same importance.

$$\text{input} = \frac{\text{input} - \bar{x}}{s}$$

$$\sqrt{\frac{\sum (x_i - \bar{x})^2}{N}}$$

solved the issue where one type of input has values of much greater magnitude than others, hence dominating the loss function.

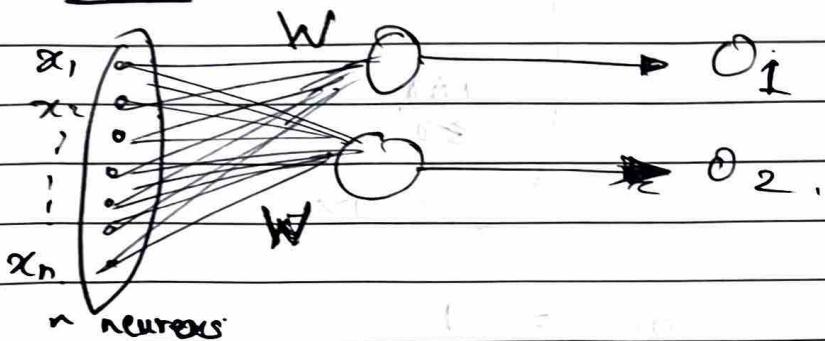
The issue is for. then others, hence dominating the loss function.

Now we come to conservation of variance

variance is also a good measure of signal strength.

for eg:

layer 2.



var = 'k'.

$$\text{var. } O_i = \sum x_i w_i$$

$$\therefore \text{var}(O_i) = \sum \left(\sum x_i w_i - \frac{\sum x_i w_i}{n} \right)^2$$

~~F.M.~~

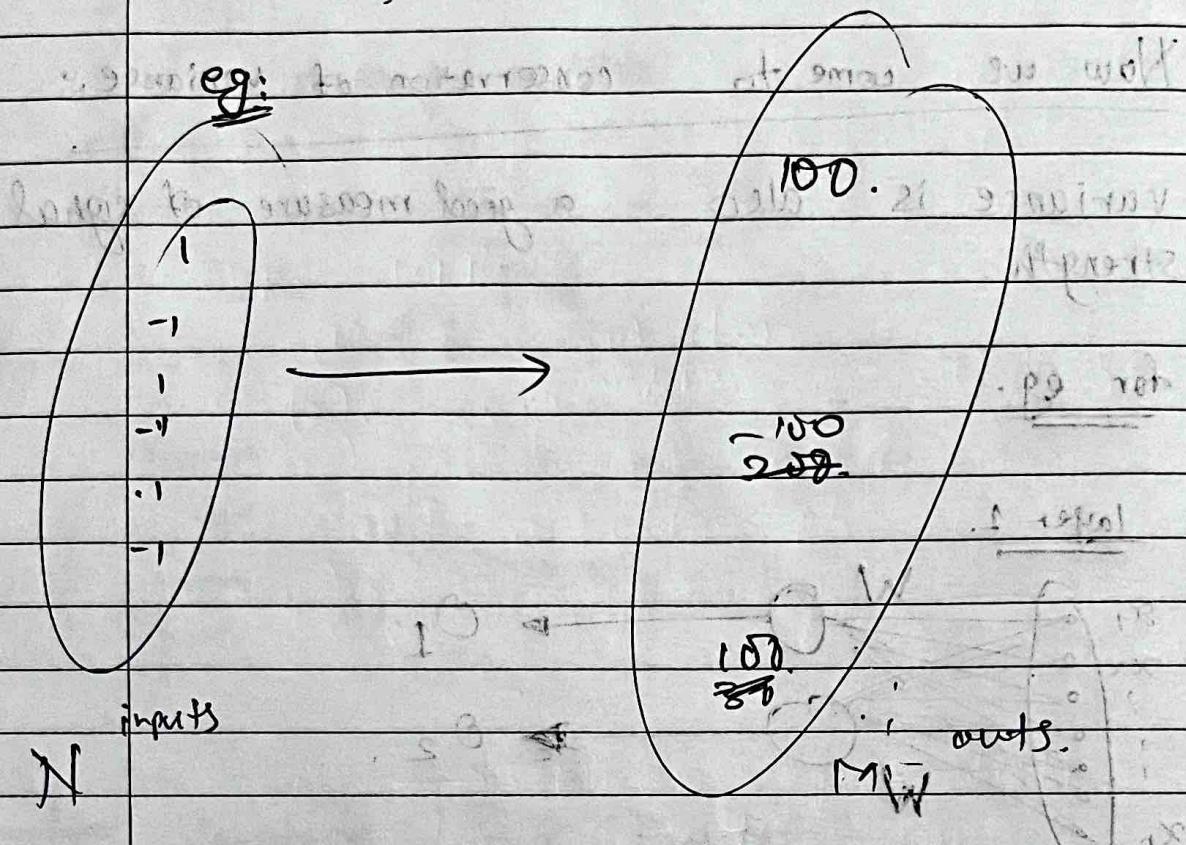
$$\frac{\sum (O_i - \bar{x}_i)^2}{N}$$

Basically it got scaled up.

$$\mathbb{E}(w^2)$$

$$\hookrightarrow \frac{\sum x^2}{n} - \left(\frac{x}{n}\right)^2 = \frac{1}{n} \sum x^2$$

— / /



$$\therefore \text{var}_{\text{old}} = 1$$

$$\text{Var}_{\text{new}} = \frac{100^2 \times M}{M}$$

$$= \frac{(100^2 \times M)}{M} = 10000$$

If weight scaled up by λ

This all depends on the weight w_j .

We want w_j to be such that variance is maintained so that after multiple layers the strength neither explodes nor vanishes.

— / /

Engg. Statistics

Thm

$$\text{Var}(y) = \left(\sum w_i^2 \right) \text{Var}(x)$$

assuming all inputs have var 1.

$$\begin{aligned} \text{var}(y) &= \sum w_i^2 \\ \frac{\text{var}(y)}{\text{var}(x)} &= \cancel{w_i^2} \end{aligned}$$

if i mul all value by R. $\therefore \sigma \rightarrow R \text{ times}$

$$\sigma = \sqrt{\text{var}} = \sqrt{\text{sum}}$$

$$\sqrt{\frac{R^2 \sigma^2 - R^2 \bar{x}^2}{n}} = R \sqrt{s^2}$$

(L) part \rightarrow σ^2

ab
to addeg'

xa

the best guess!

Makemore Iec-3

notes:

- ① Use `torch.no_grad()` when u don't expect to do a backpropg. or `@torch.no_grad()` decorator.

②

Setup: MLP-

① 3 input char

② $3 \times 10^{\text{dim/char}} = 30$, ~~vector~~ ^{dim} vector.

③ passed into hidden layer 200 neurons and activation f^h [~~w, b, tanh~~]

④ final output layers — 27 layers produces raw scores / 27 logits.

⑤ \rightarrow softmax : prob distr.

⑥ loss = nll of prob.

③ Poor initial loss.

$$\exp: \approx \log\left(\frac{1}{27}\right) \approx 3.29$$

\therefore Our logits were very overconfident at start
 \hookrightarrow make these close to 0.

$$b_2 = 0$$

$$w_2 * = 0.01$$

$w_* \neq 0 \rightarrow$ Symmetry Prob

longer or less
spent or less
Spent of time
or squashing
weights

(4)

Activation saturation:

tanh is linear and resp b/w (-2.5, 2.5)
 otherwise loses off info.

$$H = \text{tanh}(\underbrace{w_1 + b}_\text{done this down.})$$

≈ 0.1 ~ 0

(5)

torch.randn()

creates a distribution with

mean = 0.

$$\sigma = 1. \quad v = \sigma^2 = \frac{\sum (x - \bar{x})^2}{n}$$

$$v = 0.1. \quad \bar{x} = 0$$

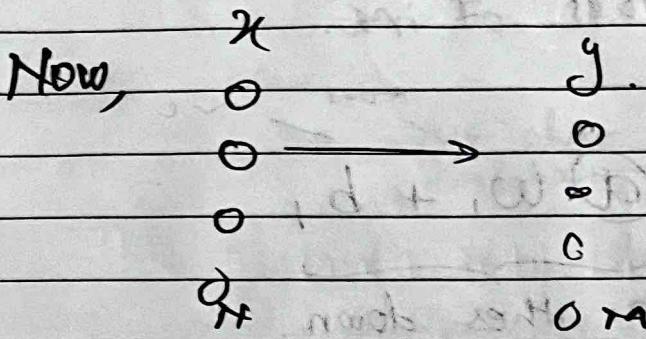
$$\frac{\sum x^2}{n} = 1. \quad \text{for these.}$$

Now, as we do the forward pass the signal strength changes due to the weights and number of final input to each neuron layer.

If the strength inc / dec too much \Rightarrow activation gets fucked up

we use Variance as a measure of strength.

initial data: var = 7.



$$\text{var}(y) = \frac{\sum (y - \bar{y})^2}{n}$$

if we see a single neuron.

$$\text{var}(\text{out}) = \text{var}(\sum w_i x_i + b)$$

single
neuron.

$$= \sum w_i^2 \text{var}(x_i)$$

if we repeat all neurons

$$\begin{aligned}\text{var}(y) &= E(E(y^2) - (E(y))^2) = \sum E(w_i x_i) \\ &= E(E(y^2)) - (\sum E(w_i x_i))^2 \\ &= E\left(\left(\sum w_i x_i\right)^2\right)\end{aligned}$$

$\therefore \text{Var}(y)$

$$= E(y^2)$$

~~$$= E \left(\sum_{i=0}^n \sum_{j=0}^n E(w_i) E(x_i) E(w_j) E(x_j) \right)$$~~

~~$$= \sum_{i=0}^n \sum_{j=0}^n E(w_i^2) E(x_i^2)$$~~

~~$$= \sum_{i=0}^n \sum_{j=0}^n E(w_i^2) E(x_i^2)$$~~

~~$$= E \left(\left(\sum_{i=0}^n E(w_i) E(x_i) \right)^2 \right)$$~~

~~$$= E \left(\sum_{i=1}^n (w_i x_i)^2 + \sum_{\substack{i,j=1 \\ i \neq j}} w_i x_i w_j x_j \right)$$~~

~~$$= \sum_{i=1}^n E(w_i x_i)^2 + \sum_{\substack{i,j=1 \\ i \neq j}} E(w_i) E(w_j) E(x_i) E(x_j)$$~~

$$= \sum_{i=1}^n E(w_i^2) E(x_i^2)$$

$$\Rightarrow \text{Var}(y) = E(y^2) = n \text{Var}(w) \text{Var}(x)$$

$$\therefore \text{var}(y) = \text{var}(x) \quad \text{iff.}$$

$$\sigma_w^2 = \text{var}(w) = 1$$

$\therefore \sigma = \frac{1}{\sqrt{\text{fan-in}}}$ no. of neurons (input)

$$\approx \frac{1}{\sqrt{\text{fan-in}}}$$

This accounts only for linear ops.

\therefore for activation functions as well,

$$\sigma_w^2 =$$

$$\sigma_w = \frac{\text{gain}}{\sqrt{\text{fan-in}}} \rightarrow \sqrt{2} \text{ for ReLU}$$

$$(\sum_{i=1}^n w_i x_i + b)^2 + (\sum_{i=1}^n w_i)^2$$

$$(\sum_{i=1}^n w_i x_i + b)^2 + (\sum_{i=1}^n w_i)^2$$

$$(\sum_{i=1}^n w_i x_i + b)^2 + (\sum_{i=1}^n w_i)^2$$

$$\text{Var}(w_i) \text{ and } x_i = (\sigma_w)^2 = w_i \text{ var } x_i$$

⑥ Batch Normalization:

Internal Covariate Shift:

to normalize a layer.

$$\mu = \frac{\sum x_i}{m}$$

$$\sigma^2 = \frac{1}{m} \sum (x_i - \mu)^2$$

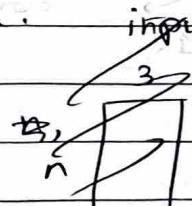
$$\therefore \hat{x}_i = \frac{x_i - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$\therefore \hat{\sigma}^2 = \frac{1}{m} \sum \frac{(x_i - \mu)^2}{m^2} = 1.$$

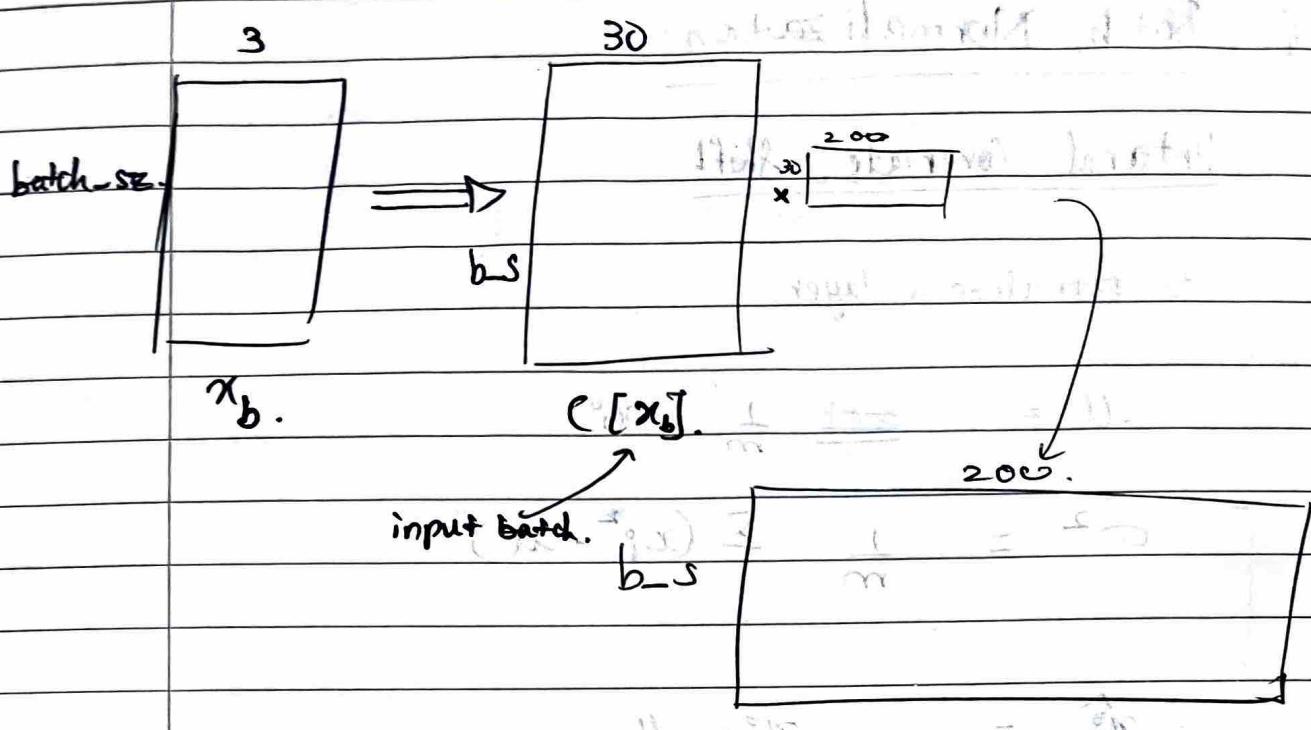
$$y_i = \gamma x_i + \beta$$

Batch normalization works on the whole batch

Basically if we have n inputs \rightarrow in 1 batch.



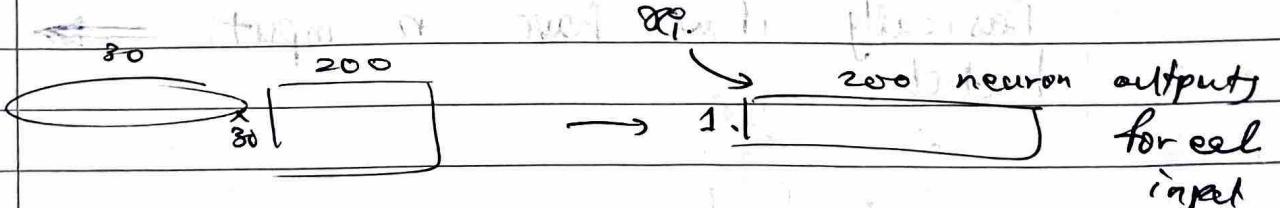
context len x emb dim.



Now we got 200 neuron outputs.
at bef pre-activation.

now @ we apply normalization on these
32. (not 200). The "batch" is normalized.

$$\therefore \bar{u} = \frac{\sum x_i}{32} \quad \sigma^2 = \frac{\sum (x_i - \bar{u})^2}{32}$$



so

We adjust x_i 's unit variance, $\sigma^2 = 1$.

$$\hat{x}_i = \frac{x_i - \bar{u}}{\sqrt{\sigma^2 + \epsilon_0}}$$

$\epsilon_0 \rightarrow$ small num increase
 $\sigma^2 = 0$

further scaling:

$$x = \gamma \hat{x} + \beta$$

↑ scale variance

scale mean

Note:

① γ, β are 1×200 vectors!!

\therefore They will

$$x_{\text{hat}} = \begin{pmatrix} a & b & c \\ d & e & f \\ g & h & i \\ j & k & l \end{pmatrix} \quad \therefore Y = (y_1, y_2, y_3) \\ \beta = (\beta_1, \beta_2, \beta_3)$$

$$\therefore x = \begin{pmatrix} \alpha y_1 + \beta_1 & \beta_2 \\ d y_1 + \beta_1 & e y_2 + \beta_2 \\ g y_1 + \beta_1 & \vdots \\ j y_1 + \beta_1 & \vdots \end{pmatrix}$$

② we don't use a bias anymore in the preactivation as it gets cancelled out anyways after $-\bar{u}$ on the \hat{x}_i creation

③ So now for the linear part, the signal strength would maintain itself but tanh squashes the inputs so we gotta multiply the outs - by a factor again = $s/3$ to maintain strength



Saturation $\sim 5, -6\%$.

Std. variance $\sim 0.6 S$.

(in case of properly linear gain ≈ 1.1)

Batch norm takes care of most things.

Summary:

① Batch Norm

② Diagnostic tools:

① Activation Distr (post-act) per layer

② Gradient Distr

③ Update to dat (over time)

(upl. std. data, std())

— / /

Backprop Ninja :

$$L = \text{loss.} = -\frac{1}{n} [l_p[0][y_b[0]] + l_p[1][y_b[1]] + \dots]$$

$$= -\frac{1}{n} \sum_{i=0}^{n-1} \text{log prob}[i][y_b[i]].$$

$$= -\frac{1}{n} [l_p$$

row wise = $\frac{1}{n}$
col wise = $\frac{1}{n}$

let's do Backprop thru matrix mult.

Consider: $z = x @ w$.

Let $g_i = \frac{\partial L}{\partial z}$

$z(i,j) = \sum_{k=0}^K (x[i,k] * w[k,j])$ each row is a column. for a weight vector for

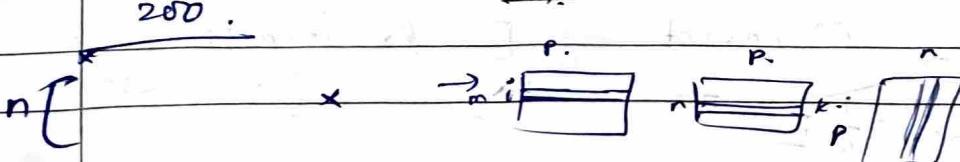
for a particular k .
 $z(i,j) = x[i,k] * w[k,j]$ one output.

$\frac{\partial L}{\partial x_{ik}} = \sum_{j=1}^P \frac{\partial z}{\partial x_{ik}} \frac{\partial z}{\partial z_{ij}} \frac{\partial z_{ij}}{\partial x_{ik}}$

$= \sum_{j=1}^P g_{ij} \times w_{kj}$

$= (G W^T)_{ik}$

$\therefore \frac{\partial G L}{\partial X} = G \times W^T$



\therefore if $z = x @ w$.

$$\therefore \frac{dz}{dx} =$$

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z} \cdot w^T$$

$$\frac{\partial L}{\partial w} = x^T \frac{\partial L}{\partial z}$$

$$\frac{\partial L}{\partial \text{logprobs}}$$

$$L = \frac{1}{n} \sum_{i=0}^{n-1} \text{logprobs}[i][y_b[i]].$$

$$= -\frac{1}{n} \sum_{i=0}^{n-1} S_i.$$

$$\frac{\partial L}{\partial S_i} = -\frac{1}{n} \nabla_i Z$$

$$\text{where } \nabla_i Z = (j, 0)_{E_{i,j}=1 \text{ if } c=Y_b[i] \text{ else } 0)$$

$$S_i = \text{logprobs} \cdot E \cdot Y$$

$$E = \underbrace{\dots}_{\text{emb}}$$

S_i = the i^{th} row of E and Y $\text{column of logprobs}$

$$\frac{\partial L}{\partial \text{logprobs}} = -\frac{1}{n} E$$

E $n \times b_s$ $b_s \times \text{vocab size}$

— / —

$$\therefore \frac{\partial}{\partial w} (x @ w)$$

$$\frac{\partial (x @ w)}{\partial x} = w^T = \frac{\partial (x @ w)}{\partial w} = \cancel{w^T} \cdot x^T$$

$$\frac{\partial L}{\partial x} = \cancel{\frac{\partial L}{\partial (x @ w)}} w^T$$

$$\frac{\partial L}{\partial w} = x^T \frac{\partial L}{\partial (x @ w)}$$

$$\cancel{\frac{\partial L}{\partial \log}}$$

$$\frac{\partial L}{\partial \text{logprob}} =$$

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z} \cdot \frac{\partial z}{\partial x}$$

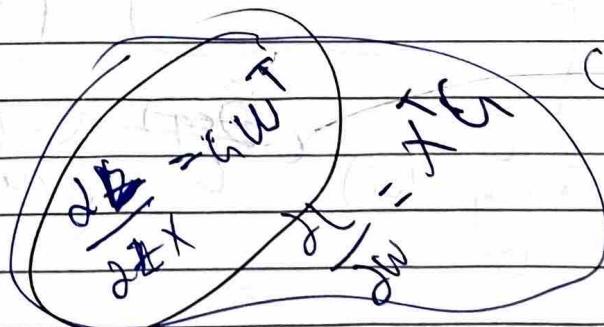
$$h = \frac{\partial L}{\partial z}$$

$$G = \frac{dh}{dz}$$

$$z = x @ w$$

$$\frac{\partial L}{\partial z} =$$

$$= G (\partial x @ w + \partial x @ \partial w)$$



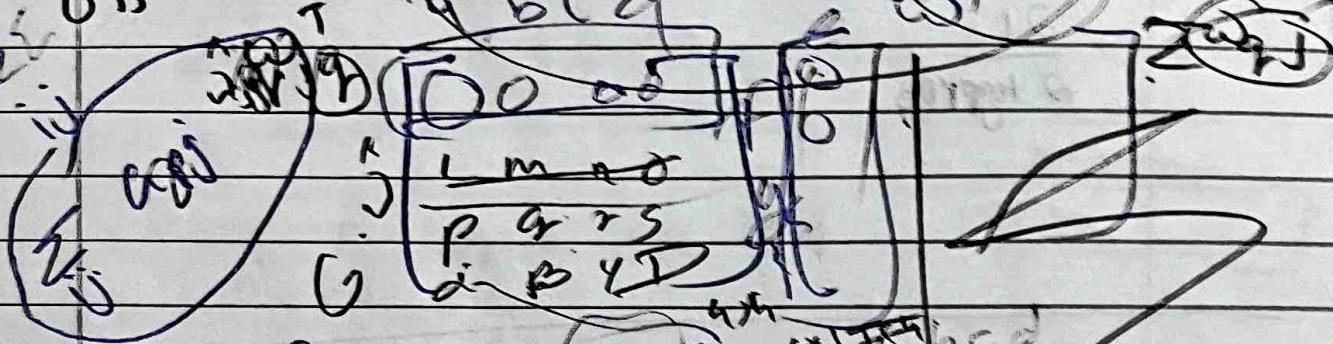
$$w + dw + x @ dw$$

$$z = x @^w_{m \times n}$$

$$\left(\frac{\partial L}{\partial z} \right) = \left(\frac{\partial L}{\partial z} \cdot \frac{\partial z}{\partial x} \right)_{M \times N}$$

$$= \sum_{i,j} \left(\frac{\partial L}{\partial z} \right)_{ij} \frac{\partial z_{ij}}{\partial x_{ij}} w_{ij}^2$$

w_{ij}, w_{ij}



$$G = E$$

$$\frac{\partial L}{\partial x} = \left(\begin{matrix} \dots & \dots & \dots \\ \dots & \dots & \dots \\ \dots & \dots & \dots \end{matrix} \right)$$

$$= \left(\begin{matrix} \dots & \dots & \dots \\ \dots & \dots & \dots \\ \dots & \dots & \dots \end{matrix} \right)$$

$$\frac{\partial L}{\partial x} = \left[\begin{array}{c|ccccc} \frac{\partial L}{\partial a} & a & b & c & d \\ \hline a & e & f & g & h \end{array} \right]$$

$$\frac{\partial L}{\partial a} = \frac{\partial L}{\partial a} + \frac{\partial L}{\partial b} + \frac{\partial L}{\partial c} + \frac{\partial L}{\partial d}$$

-/-

$$\frac{\partial L}{\partial x} = G(\omega)$$

$$z = \omega \cdot \alpha \times (\alpha) \omega$$

$$z =$$

$$x = \begin{bmatrix} AB & b \\ C & F \\ i & j \end{bmatrix} \quad \omega = \begin{bmatrix} \alpha \\ \beta \\ \gamma \end{bmatrix} \quad \begin{bmatrix} A & \alpha \\ F & \beta \\ K & \gamma \end{bmatrix}$$

$$z = \frac{\partial}{\partial (A\alpha + B + C\beta + \gamma)} \cdot \frac{\partial L}{\partial \alpha}$$

$$\frac{\partial \alpha}{\partial L} \cdot \frac{\partial (A\alpha + B + C\beta + \gamma)}{\partial L \cdot \alpha}$$

$$= \frac{\partial \alpha}{\partial L} \cdot \frac{\partial (A\alpha + B + C\beta + \gamma)}{\partial L \cdot \alpha}$$

$$x @ \omega$$

$$\left(\frac{\partial L}{\partial x} \right)_{ij} = (G)_{ij} \alpha \omega^T$$

$$\begin{aligned} \frac{\partial L}{\partial x} &= \frac{\partial L}{\partial z} \cdot \frac{\partial z}{\partial x} \\ &= G \cdot \end{aligned}$$

$$Z_{ij} = \left[\sum_{k=1}^n x_{ik} y_{kj} \right]$$

$$\left(\frac{\partial L}{\partial Z} \right)_{ij} = \frac{\partial L}{\partial \left(\sum_{k=1}^n x_{ik} y_{kj} \right)}$$

$$G_{ij} = \frac{\partial L}{\partial \left(\sum_{k=1}^n x_{ik} y_{kj} \right)}$$

$$\left(\frac{\partial L}{\partial x} \right)_{ij} = \left(\left(\frac{\partial L}{\partial z} \right) \cdot \left(\frac{\partial z}{\partial x} \right) \right)$$

$$\frac{\partial L}{\partial x} = G \quad \omega^T$$

$$G = \frac{\partial L}{\partial z}$$

$$z = x @ w$$

$$G = \left[\begin{array}{c} \frac{\partial L}{\partial (ax + \beta b + cy + dz)} \\ \vdots \end{array} \right]$$

$$z = \begin{pmatrix} a & b & c & d \\ \alpha & \beta & \gamma & \delta \\ 1 & 2 & 3 & 4 \\ x_1 & x_2 & x_3 & x_4 \end{pmatrix} \begin{pmatrix} a & b & c & d \\ \alpha & \beta & \gamma & \delta \\ 1 & 2 & 3 & 4 \\ x_1 & x_2 & x_3 & x_4 \end{pmatrix}^{-1}$$

— / —

$$= \left(\frac{\partial L}{\partial i} \alpha \right) + \frac{\partial L}{\partial b}$$

$\alpha = a\alpha + \dots$

$$\frac{\partial L}{\partial a} = \frac{\partial L}{\partial (a\alpha + b\beta + \dots)} \alpha$$

$$\frac{\partial a}{\partial L} = \frac{\partial (a\alpha + b\beta + \dots)}{\partial L} \alpha$$

$$= \frac{\partial (a\alpha)}{\partial L} \alpha + \frac{\partial (b\beta)}{\partial L} \alpha + \dots$$

$$\frac{x^2}{x^2}$$

$$\frac{\partial (x^2)}{\partial x}$$

$$\frac{\partial (x^2)}{\partial x}$$

$$\frac{\partial L}{\partial a} = \frac{\partial L}{\partial a} \frac{\partial a}{\partial (a\alpha + b\beta + c\gamma + \dots)} \alpha$$
$$= \frac{\partial L}{\partial a} \cdot \frac{\partial (z - b\beta - c\gamma - d\delta\alpha)}{\partial z} \alpha$$

$$z = \alpha \times @ \cdot w.$$

$$\alpha = \frac{dl}{dz}.$$

$$\text{Find: } \frac{dl}{dw}.$$

$$\begin{aligned}\therefore \left(\frac{dl}{dw}\right)_{pq} &= \sum_{ij} \left(\frac{dl}{dz}\right)_{ij} \cdot \frac{\partial z_{ij}}{\partial w_{pq}} \\ &= \sum_{ij} g_{ij} \cdot \frac{\partial}{\partial w_{pq}} \sum_i (x_{ip}^T \cdot w_{pq} \cdot 1) \\ &= \sum_{ij} g_{ij} \cdot x_{ip}^T \cdot w_{pq}.\end{aligned}$$

$$= \sum g_{ij} x_{ip}^T \cdot w_{pq}.$$

$$= \sum x_{ip}^T \cdot g_{ij} w_{pq}$$

$$= x^T g w.$$

$$\frac{dl}{dw}$$

Direct Calculation of Loss:

$$\text{logits} = \tan(\dots)$$

$$\therefore \log \left(\frac{\exp(\text{logits})}{\sum \exp(\text{logits})} \right)$$

$$\therefore \frac{\partial \text{loss}}{\partial \text{logits}_i} = \cancel{\eta} \sum_{l_j} \cancel{\exp(\text{logits}_j)} \cdot \left(e^{\frac{l_i}{\eta}} + \left(\frac{1}{\sum e^{\frac{l_j}{\eta}}} \right)^2 \cdot e^{\frac{l_j}{\eta}} \right)$$

— / /

$$\text{logits} = \mathbf{h}_2 @ \mathbf{w}_2 + b_2 \quad \text{(softmax and logit)}$$

$$\text{logit_maxes} = \text{logits}.max(1, \text{keep_dim})$$

$$\text{dnorm_log} = \text{logits} - \text{logit_maxes}$$

$$\therefore d(\text{norm}) = 1$$

$$\frac{\partial \text{norm}}{\partial \text{logits}}$$

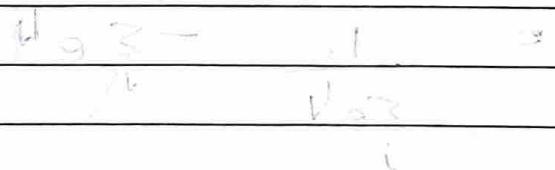
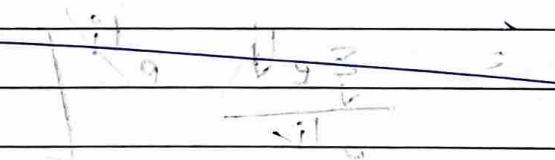
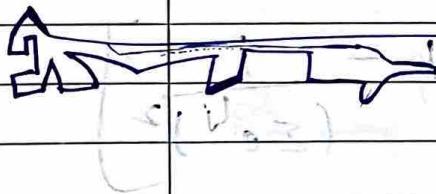
emb:

32.

10.

we wanna add

this vector to its resp
grad emb row in C.



* Direct loss Calculation:

We want derivative of loss wrt a particular vector of the logits $\rightarrow l_i$

$\therefore \frac{\partial \text{Loss}}{\partial l_i}$

$$\text{loss} = -\log \left(\frac{\text{true e}^{l_i}}{\sum(\text{e}^{l_j})} \right)$$

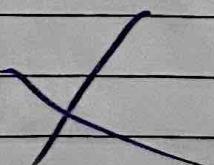
$$\text{loss}_i = -\log \left(\frac{e^{l_i}}{\sum_j e^{l_j}} \right)$$

$$\therefore \frac{\partial \text{Loss}}{\partial l_i} = -\left(\frac{j \sum e^{l_j}}{e^{l_i}} \right) \left[e^{l_i} + \frac{(-1)}{(\sum e^{l_j})^2} \cdot (0 - 1 - e^{l_i}) \right]$$

$$= -\frac{\sum e^{l_j}}{e^{l_i}} e^{l_i} \left[1 - \frac{1}{(\sum e^{l_j})^2} \right]$$

$$= \frac{1}{\sum e^{l_j}} - \frac{\sum e^{l_i}}{e^{l_i}}$$

Wrong



shape of logits = $\frac{\text{vocab-size}}{\text{minibatch size}}$

\therefore while creating the differentiation matrix, we gotta consider both.

$$\left(\frac{\partial \text{Loss}}{\partial \text{logits}} \right)_{i,c} =$$

$$i \in \{0, 1, \dots, n\}$$

$$c \in \{0, 1, \dots, 26\}$$

$$\text{Now let log loss} = -\log(p_y)$$

Prob of the expected char'y

and let us consider just one eg of the batch for now.

$$\text{loss}_y = -\log \left(\frac{e^{l_y}}{\sum_j e^{l_j}} \right)$$

$$\therefore \left(\frac{\partial \text{loss}_y}{\partial \text{logits}} \right)_c \leftarrow \text{char.}$$

$$\therefore \frac{\partial \text{loss}_y}{\partial \text{logits}} =$$

$$= - \left(\frac{\sum e^{l_j}}{e^{l_y}} \right) \cdot \left[\frac{e^{l_c} e^{l_y}}{l_c} + \frac{(-1)}{(\sum e^{l_j})^2} \cdot \frac{\sum e^{l_j}}{l_c} \right]$$

if $y = c$.

$$\therefore \left(\frac{\partial \text{loss}_y}{\partial \text{logits}} \right)_c = - \frac{\sum e^{l_j}}{e^{l_y}} \left[0 + \frac{(-1)}{(\sum e^{l_j})^2} \cdot c \right]$$

$$v = \frac{v dy_{\text{ax}} - u dx_{\text{ay}}}{v^2}$$

$$\text{lossy} = -\log \left(\frac{e^{\text{logit}_y}}{\sum_j e^{\text{logit}_j}} \right)$$

$$\frac{\partial \text{lossy}}{\partial \text{logit}_c} = (-1) \left(\frac{\sum_j e^{l_j}}{e^{l_y}} \right) \left(\frac{\frac{\partial e^{l_y}}{\partial \text{logit}_c} \cdot \sum_j e^{l_j} - e^{l_y} \cdot \frac{\partial \sum_j e^{l_j}}{\partial \text{logit}_c}}{(\sum_j e^{l_j})^2} \right)$$

if $y = c$. ← At the indices (columns) of logprobs that are the correct pred char.

$$\frac{\partial \text{lossy}}{\partial \text{logit}_c} = (-1) \left(\frac{\sum_j e^{l_j}}{e^{l_y}} \right) \left(\frac{e^{l_y} \cdot \sum_j e^{l_j} - e^{l_y} e^{l_y}}{(\sum_j e^{l_j})^2} \right)$$

$$\frac{\partial \text{lossy}}{\partial \text{logit}_c} = \frac{e^{l_y} - (\sum_j e^{l_j})}{(\sum_j e^{l_j})^2}$$

$$P_c - 1 \quad \text{or} \quad P_y - 1$$

$$\text{if } y \neq c \quad \frac{\partial \text{lossy}}{\partial \text{logit}_c} = (-1) \left(\frac{\sum_j e^{l_j}}{e^{l_y}} \right) \left(\frac{0.5 - e^{l_y} \cdot e^{l_y}}{(\sum_j e^{l_j})^2} \right)$$

$$\frac{(e-1) \cdot e^{l_c}}{(\sum_j e^{l_j})^2} = -P_c$$

$$\text{Where } P_c = \frac{e^{\text{logit}_c}}{\sum_j e^{\text{logit}_j}}$$

11

Via code:

27.

logits: $\begin{cases} i \\ 32 \end{cases}$ $\begin{cases} j \\ 32 \end{cases}$ $\begin{cases} \text{moving left to right} \\ \text{and right to left} \end{cases}$

$$\text{dlogits} = \begin{cases} p_i & p_{j-1} \\ 32. \end{cases} = \text{softmax}$$

so once you think about it,

p_c = some kind of prob b/w 0 and 1.

p_{c-1} = b/w -1 and 0.

The grad of $c \neq y$ is +ve.
The grad of $c = y$ is -ve

$$D(\text{data}) = -(\mathbf{r}^* \text{grad.})$$

We wanna pull down the probs of all the ~~g~~ probs that are assigned to the incorrect character and pull up all the correct ones.

Same thing for batch normalization:

Calc. dprebn. given dhpreact.

$$h_{\text{preact}} = \frac{x_i^0 - \bar{x} - \mu_B}{\sqrt{\sigma_B^2 + \epsilon_0}}$$

hprebn:

for a particular input,

$$\left(\frac{\partial h_{\text{preact}}}{\partial h_{\text{preact}}} \right)_{ij} = \frac{\partial}{\partial x_i^0} \left(\frac{x_i^0 - \bar{x} - \mu_B}{\sqrt{\sigma_B^2 + \epsilon_0}} \right)_{ij} = \frac{1}{\sqrt{\sigma_B^2 + \epsilon_0}}$$

good

$$\left(\frac{\partial h_{\text{preact}}}{\partial h_{\text{preact}}} \right)_{ij} = \left(\frac{\partial}{\partial x_i^0} \left(\frac{x_i^0 - \bar{x} - \mu_B}{\sqrt{\sigma_B^2 + \epsilon_0}} \right) \right)_{ij} = \frac{1}{\sqrt{\sigma_B^2 + \epsilon_0}} \cdot \frac{\partial}{\partial x_i^0} (x_i^0 - \bar{x} - \mu_B)$$
$$= \frac{1}{\sqrt{\sigma_B^2 + \epsilon_0}} \cdot [1 - \frac{1}{n}] - (x_i^0 - \bar{x} - \mu_B) \frac{1}{\sqrt{\sigma_B^2 + \epsilon_0}} \cdot \frac{1}{n}$$

$$\frac{\partial \sigma_B^2}{\partial x_i^0}$$

$$= \frac{1}{n} \cdot 2(x_i^0 - \bar{x})$$

$$(x_i^0 - \bar{x})^2$$

Contd. NaiveRNN 1ec-5:

Building a Wavenet:

first we wanna make some improvements and revise our og code:

layers:

Linear:

~~fan-in, fan-out, bias = True~~

weight = $\text{randn}(\text{in}, \text{out}) / \sqrt{\text{in}}$

bias = $\text{randn}(\text{out}) * 0.01$ or $\text{torch.zeros}(\text{out})$

call:

out = $x @ w$.

if bias:

out += bias

ret out.

param:

All params \rightarrow weight + bias if bias

$$\begin{aligned} & \frac{m \cdot n - x}{3 + 10x - 1} \\ & 9 + 10x = 400 \\ & (2 \log_{10} m + 1) + 10x = 400 \\ & (2 \log_{10} m + 1) + 10x = 400 \end{aligned}$$

class BatchNorm:

def: (dim, training=True, eps=10⁻⁸, momentum=0.1)

Y, B = ~~randn(dim)~~

Y = ones(dim)

B = zeros(dim).

←

training = True

run-mean = ~~zeros(dim)~~ ~~abs~~ = 0.0001

run-var. = ~~ones(dim)~~ ~~abs~~ = 0.0001

cell: (x)

if training:

x-m = x.mean()

x-v = x.var().

else:

x-m = running_m

x-v = run-var

$$\hat{x} = \frac{x - x_m}{\sqrt{x_var + \epsilon}}$$

$$out = Y * \hat{x} + B.$$

with torch.no_grad():

$$\begin{aligned} \text{running_mean} &= r_m * (1 - mom) + x_m * mom \\ \text{run_var} &= r_v * (1 - mtn) + r_v * mtn \end{aligned}$$

Params:

(Y, B)

no
backward

Tahk \rightarrow $\phi \rightarrow 2 \rightarrow$ Lasso, Dots

call : (x)
return \rightarrow torch.tanh(x)

perams:

C J Z

Brachys 2490 [

Initial Setup:

$$C = \text{randn}(\text{vocab}, \text{emb})^{27 \times 100}$$

~~layers = [linear, BN, tanh, linear].
 \downarrow
 (contemb, n-hidden) (n-hidden, vocab)~~

~~$x = \circ (x_b)$~~

for l in layers

$$x = \alpha(x)$$

~~loss = F. cross-entropy (x)~~

for 1 in layers : ()
 1. wt % = 0.1.
 (to make initial loss less overconfident)

Setup:

$n_emb = 10$

$n_hidden = 200$.

`torch.manual_seed(42)`

(x) 1103

`nn.init.`

`uniform`

(x) `init.` do it!

X_b

$C = \text{randn}(\text{vocab_size}, n_emb)$.

] gets replaced

`out[2].weight`

~~$C [Ex_b]$~~

(~~nn~~)

$\text{layers} = [\text{Linear}(\text{vocab_size} * \text{context_dim}, n_hidden),$

$\text{BN}(n_hidden),$

$\tanh(),$

$(\text{Linear}(n_hidden, vocab_size))]$

with $\text{nograd}():$

for l in layers :

$l.\text{wt} *= 0.1$ # decrease initial overconfidence.

$\text{params} = [\text{p for l in layers for p in (l.parameters(),)}$

+ [C]]

→ we need.

for p in $\text{params}():$

$p.\text{requires_grad_}$

iter = 200,000.
tr batch-size = 32
loss[i] = 0.

loop(iter):

~~x = C[X]~~

~~xb[i] = math.randint(batch_size)~~

~~xb, yb = Xtr(ix), Ytr(ix)~~

~~x = C[xb]~~

~~x = x.view(...)~~

~~for i in layers~~

~~x = L(x)~~

} reply.

loss = F. cross_entropy(xb, yb)

for p in params(): p.grad = 0.

loss.backward()

loss[i][j] = loss.

lr = 0.1 if i < 150000 else 0.01

p.delta -= lr * p.grad.

testing:

for i in layers:

l.training = False

:

and sampling

fixing lr pth.

loss_i.view(-1, 100) or loss_i.view(-1, 1000).

putting lookup and embedding table into layers.

class Embedding:

def __init__(self, num_emb, emb_dim)

wt = torch.randn((n_e, e_d))

call(ix)

return wt[ix]

params.

ret [wt].

class Flatten

~~def~~

call(x)

or out = x.view(x.shape[0], -1))

ret out.

params

ret ()

now we don't need C ,
→ params directly takes from layers.

→ Add:

Setup } Embedding (n-emb, emb-dim),
 | flatten
 | to layers.
 | no need for explicit C ; $\Rightarrow [C]$ in params

Also

loop { $x = C[x_b]$ }
 $x = x.view \quad \rightarrow \quad x = x_b$
for layer:
 $x = f(x)$

And now,

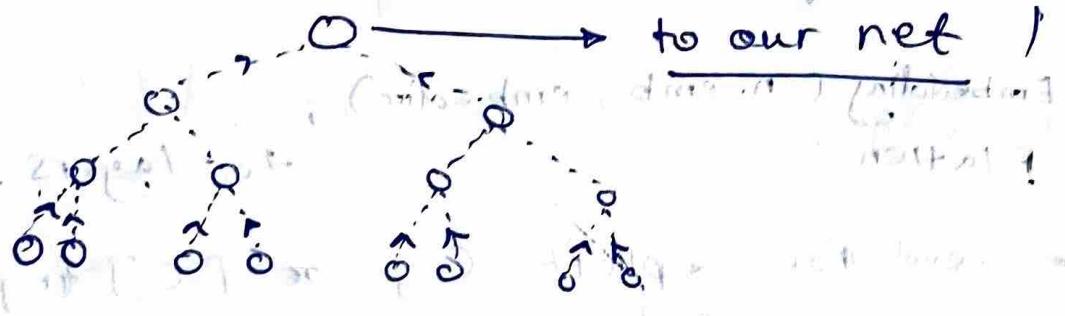
class Sequential:
def __init__(self, layers):
 self.layers = layers.

call: (x)
 $(0, 8)$ → for i in layers:
 end $x = f_i(x)$
return x

new:
 $x = x_b$.
for ...
 ...
} \rightarrow logits = model(X_b).

* Set layers to model.layers is to
training in.
false ~~for~~ to + 19

Introducing the Wavenet!



currently, if we increase cont-size to 8
then we get: $\text{emb}[x_b]$

(batch-size, 8, 10) $\xrightarrow{\text{flat}}$ (batch-size, 80)

but we want

(b-s, 8, 10) \rightarrow (b-s, 4, 20)

instead of

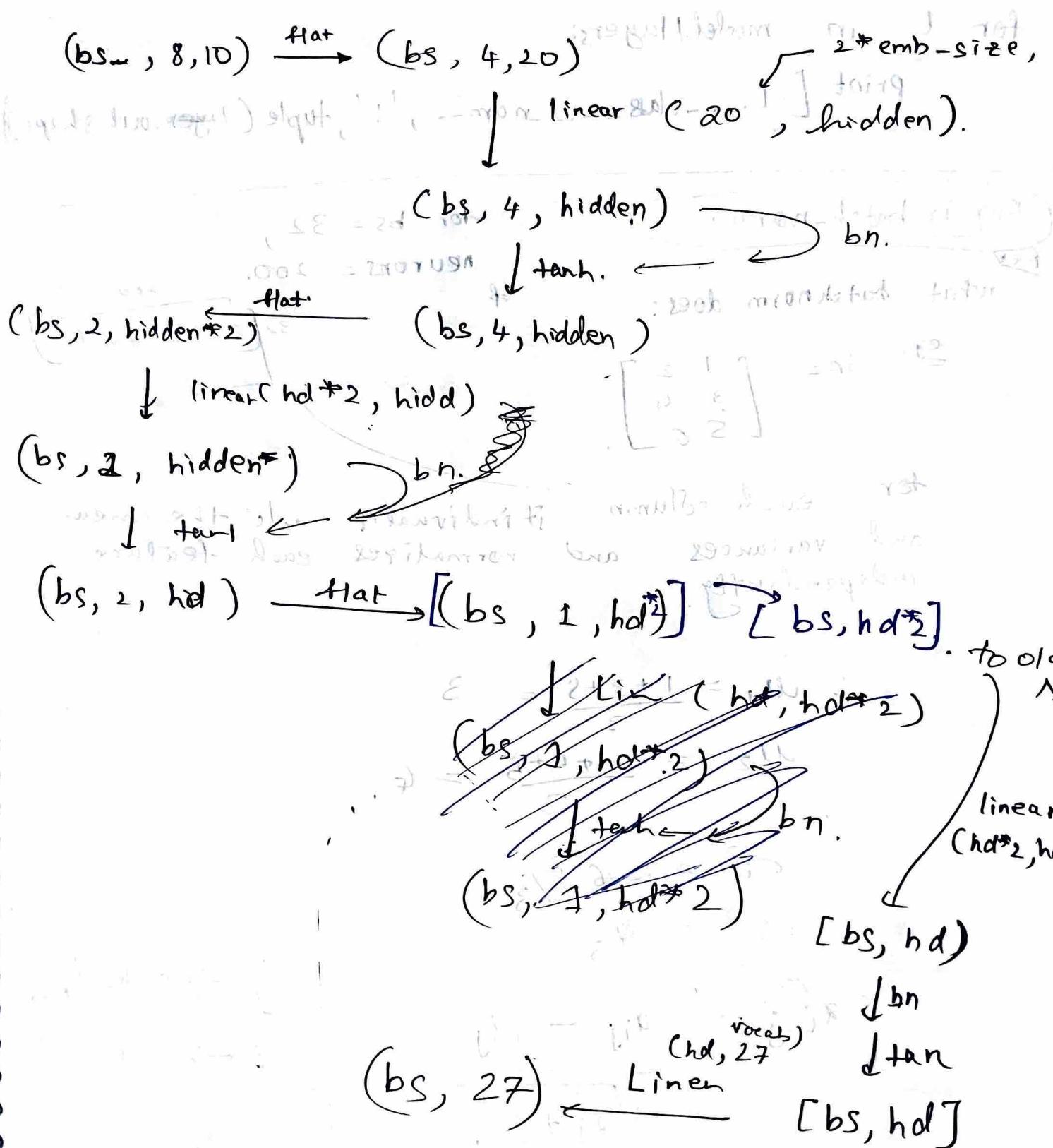
(b-s, 80) @ (80, 200) + (800)

we want

(b-s, 4, 20) @ (20, hid1) + (hid1)

so now in flatten(), we want to basically
each time concatenate 2 char embeddings at a time,
pass them thru a linear and tanh layer and
do the same thing twice more.

we want:



A very nice loop to see this in code:

for l in model.layers:

print [l.class_name, ':', tuple(layer.out_shape)]

Bug in batch-norm.

Rev

what batchnorm does:

e.g:

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

for each column and variances and normalizes each feature independently.

$$[3 \times 2d] \rightarrow [3 \times 1, 2d] \rightarrow [3, 2d]$$

$$\mu_1 = \frac{1+3+5}{3} = 3$$

$$\mu_2 = \frac{2+4+6}{3} = 4$$

$$\sigma_1^2 = \frac{2+8}{3} = 3$$

$$[3, 2d] \quad \sigma_2^2 = 3$$

$$\hat{x}_{ij} = \frac{x_{ij} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

$j \rightarrow \text{column/feature}$

$i \rightarrow \text{row.}$

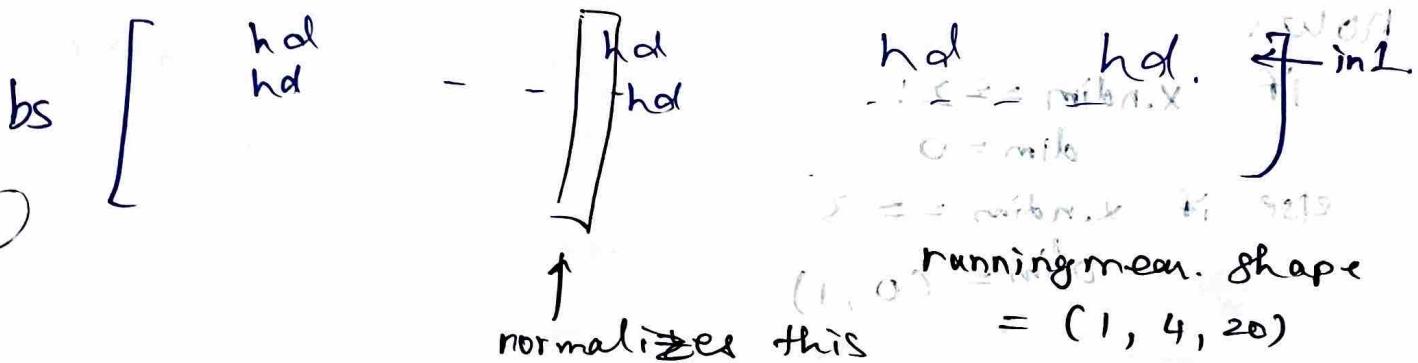
batching

But now during wavenet impl,
 Many ~~etc~~ inputs are like $(bs, 4, hd.)$
 or $(bs, 2, hd)$.

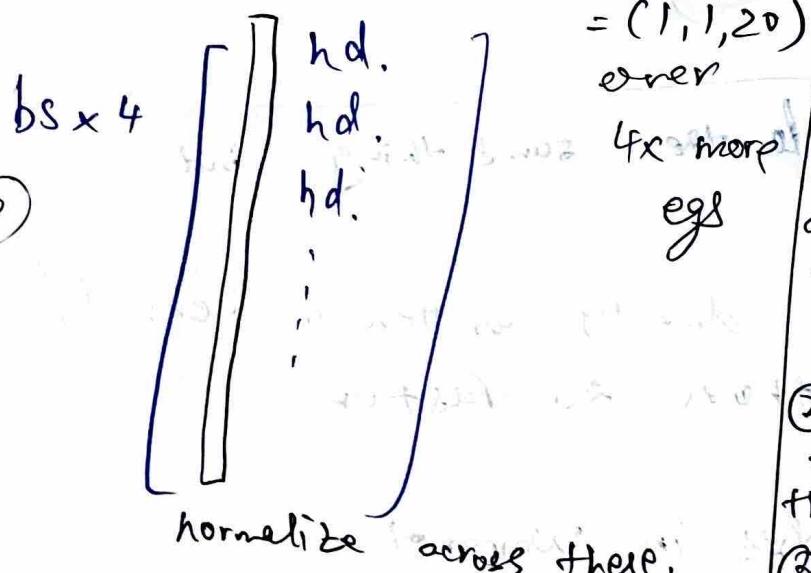
So now, it still ~~normalizes~~ over bs,

it treats each of the ~~etc~~ $4 \neq hd$ numbers
 as individual ~~features~~ \rightarrow $(bs, 4, 1)$

its like



But we want



- Reasons for this:
- ① column j of each vector repr the same feature so in meth-① we are treating them diff. They should be normalized together
 - ② meth-① reduces the num of examples ~~and~~ to calc the mean, variance stats.
 - ③ final dimensions of $\gamma, \beta, \text{etc}$ are gonna be comb-size only so better train them like that.

N^*L, D

N, L^*D^{2d}

$(b, t, s, 2d)$

X

to

fix: 2d into 3d later on will be won earlier

$xmean = x.mean(0, keep=T)$

$xvar = x.var(0, keep=T)$.

now:

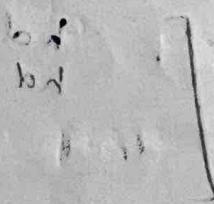
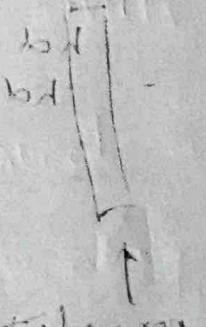
if $x.ndim \geq 2$:

dim = 0

else if $x.ndim == 3$

dim = (0, 1)

$(b, t, s, 2d)$



~~xmean = x.mean(dim, keep=True)~~

~~xvar = x.var(dim, keep=True)~~

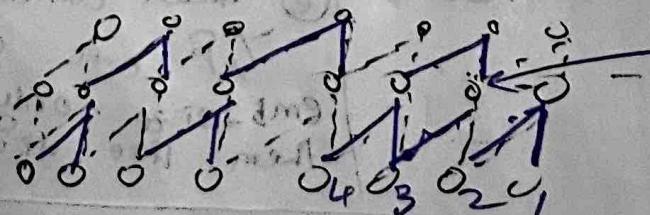
Convolution: CNNs do the same thing

more efficient

than all

① for loop is somehow directly written in kernel/CUDA
shit, not in python so faster

② reuse of ~~for~~ nodes in wavelet



once calc
for 2, 3
not recal - for
3, 4