# Big Picture with Batch Normalization

Each is a batch of 32

$C$

Lookup Table

$X$

Minibatch

$E$

Embedded Table

$H_O$

Linear Layer

$$H_O = W_1 \cdot E + b_1$$

$H_N$

Normalized Layer.

$$H_N = \frac{H_O - \mu}{\sigma}$$

$H_N$

Scale & Shift

$$H_N = \gamma H_N + \beta$$

$H_a$

activation

$$H_a = \tanh(H_a)$$

$\hat{y}$

predicted

$$\hat{y} = W_2 H_a + b_2$$

Start with the lookup table - C - Batch size × Length of Alphabet
     - Initially Random.                    32              27

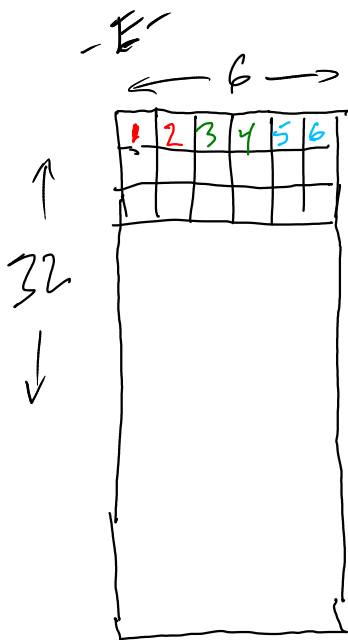Get a Random batch of bigrams - X - Batch size × Blocksize
                                    32  ×  3

Create Embedded table -
          E - Batch Size × Dimensions · Blocksize
              32    ×       2 · 3
              _____
                    32 × 6

     E will now be our inputs to the Network.



- E -
← 6 →
| 1 | 2 | 3 | 4 | 5 | 6 |

32

Filled with the
vectors from -C-
in order given from
     -X-

- C -
← 2 →

| 3 | 4 |
| 1 | 2 |  27
| 5 | 6 |

think of each
row as a 2D
vector (x,y)

- X -
← 3 →
| 5 | 3 | 7 |

32

these bigrams
determine
the order of
-C- vectors
placed into
the embedded
table -E-

Start forward Pass -
     100 hidden Neurons -H₀-
     6×100 Weight 1 -W₁-
     100 bias 1 -b₁-
     Note: Batch Size = 32×6

     H₀ = W₁ · E + b₁ =



          ← Original
            Hidden Layer
     Note: Actually 32 Hidden
                        Layers.

We now want to normalize the entire batch $\Rightarrow$ 32×100 Nodes
After normalization, the values in the hidden Layer
   will have a mean $(\mu)$ of zero & a standard deviation $(\sigma)$
   of one.

$$\mu = \phi$$
$$\sigma = 1$$

$$H_{Normal} = \frac{H_{original} - \mu}{\sigma}$$

$$\mu = \frac{1}{n} \sum_{i}^{n} H_o^i$$

$\longrightarrow$ <span style="color:red">n-1 Bessel's Correction</span>

$$\sigma^2 = \frac{1}{n} \sum_{i}^{n-1} (H_o^i - \mu)^2 \quad \text{<span style=\"color:red\">}\sigma^2 \rightarrow \text{Variance</span>}$$

$$H_N = \frac{H_o - \mu}{\sqrt{\sigma^2 + \epsilon}} \longleftarrow \text{<span style=\"color:red\">}10^{-5} \text{ (avoid} \div \text{zero)</span>}$$
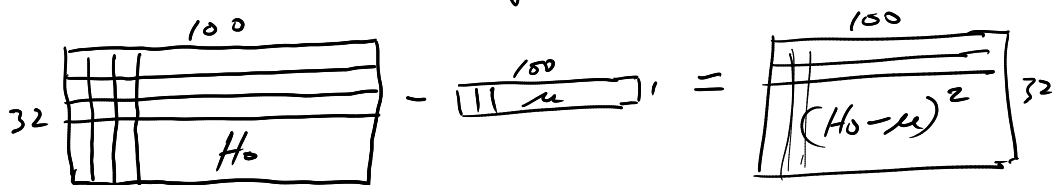
First calculate the mean $(\mu)$ for the batch



Sum the columns
& divide by
batchsize (32)

$$\mu = \frac{1}{n} \sum_{i}^{n} H_o^i$$

Next calculate the variance of the batch



Subtract the mean from each
row of the batch & square the result
then add the columns & divide by batchsize.

$$\sigma^2 = \frac{1}{32} \sum_{i}^{32-1} (H_o^i - \mu)^2$$

Finally Complete the Normalization

$$H_N = \frac{H_o - \mu}{\sqrt{\sigma^2 + \epsilon}} = \boxed{H_o - \mu} \times \frac{1}{\sigma} \text{ (if } \epsilon = \phi)$$

multiply each row by $\frac{1}{\sigma}$

# Scale $\gamma$ & Shift $\beta$

Not actually needed for Batch Normalization but used in most cases.

Allows the normalized values: $\mu = \emptyset$ & $\sigma^2 = 1$ to be altered by these two parameters.

(mean) (variance)

i.e $H_z = \gamma H_n + \beta$ $\Longrightarrow$ $\gamma \rightarrow$ scale (or gain)
$\beta \rightarrow$ shift (or bias)

In our case both $\gamma$ & $\beta$ are vectors of 100 length.
And we set $\gamma$ to all one's & $\beta$ to all zero's

We can now implement the activation function —
in our case it's $\tanh(x)$

. So $H_{activated} = H_a = \tanh(H_z)$

Finally we can compute the last linear layer of our network.

$Logits = W_2 \cdot H_a + b_2$ where $W_2$ is $100 \times 27$ & $b_2$ is $27$



Use softmax to convert logcounts (Logits) to probabilities

$\hat{Y} = $

$\Longrightarrow$ probabilities.
- 32 Rows
- each row should sum to $\underline{1.00}$

Now start back propagation;
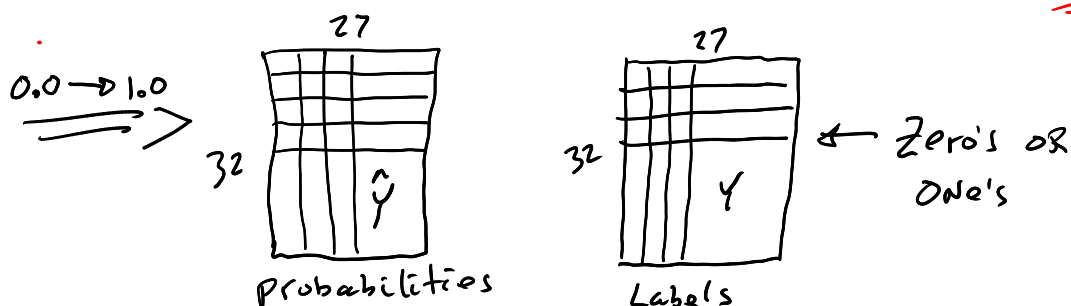First we need to determine $\frac{\partial L}{\partial \hat{y}}$   derivative of Loss w.r.t output - probabilities

$L = \frac{1}{2}(\hat{y} - Y)^2 = \frac{1}{2}(\hat{y}^2 - 2\hat{y}Y + Y^2)$

$\frac{\partial L}{\partial \hat{y}} = \frac{1}{2}(2\hat{y} - 2Y) = \underline{\hat{y} - Y}$    probability (0-1)
actual (0-1)
Labels

$0.0 \rightarrow 1.0$

probabilities


Labels
$\leftarrow$ Zero's or One's

So now we have $\frac{\partial L}{\partial \hat{\gamma}}$ but we need the followings:

$$\frac{\partial L}{\partial b_2}, \frac{\partial L}{\partial W_2}, \frac{\partial L}{\partial \beta}, \frac{\partial L}{\partial \gamma}, \frac{\partial L}{\partial b_1}, \frac{\partial L}{\partial W_1}, \frac{\partial L}{\partial C}$$

& these deltas will be used to update

$$b_2, W_2, \beta, \gamma, b_1, W_1, C$$

but to do this we'll also need several intermediate derivatives such as $\frac{\partial L}{\partial H_N}, \frac{\partial L}{\partial \sigma^2}, \frac{\partial L}{\partial \mu}$ & $\frac{\partial L}{\partial H_0}$

Start with $\frac{\partial L}{\partial b_2}$: $\quad \hat{\gamma} = W_2 \cdot H_a + b_2 \implies \frac{\partial \hat{\gamma}}{\partial b_2} = 1$

So by chain rule: $\quad \frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial \hat{\gamma}} \cdot \frac{\partial \hat{\gamma}}{\partial b_2} = \frac{1}{n} \sum_i (\hat{\gamma} - \gamma) \cdot 1$

And $\frac{\partial L}{\partial W_2} = \frac{\partial L}{\partial \hat{\gamma}} \cdot \frac{\partial \hat{\gamma}}{\partial W_2} = (\hat{\gamma} - \gamma) \cdot H_a$



$$\underset{100}{\left[ \underset{32}{H_a^T} \right]} \cdot \underset{32}{\left[ \underset{27}{\frac{\partial L}{\partial \hat{\gamma}}} \right]} = \underset{100}{\left[ \underset{27}{\frac{\partial L}{\partial W_2}} \right]}$$

**Note: Need to transpose $H_a$**

Now need: $\frac{\partial L}{\partial \beta}$: $\quad H_a = \tanh(H_z) \quad \hat{\gamma} = W_2 H_a + b_2$

$$H_z = \gamma H_N + \beta$$

by chain Rule: $\quad \frac{\partial L}{\partial \beta} = \frac{\partial L}{\partial \hat{\gamma}} \cdot \frac{\partial \hat{\gamma}}{\partial H_a} \cdot \frac{\partial H_a}{\partial H_z} \cdot \frac{\partial H_z}{\partial \beta}$

$$\underset{\hat{\gamma}-\gamma}{\uparrow} \quad \underset{W_2}{\uparrow} \quad \underset{1-\tanh^2(H_z)}{\uparrow} \quad \underset{1}{\uparrow}$$

and $\frac{\partial L}{\partial \gamma} = \frac{\partial L}{\partial \hat{\gamma}} \cdot \frac{\partial \hat{\gamma}}{\partial H_a} \cdot \frac{\partial H_a}{\partial H_z} \cdot \frac{\partial H_z}{\partial \gamma}$

$$\underset{\hat{\gamma}-\gamma}{\uparrow} \quad \underset{W_2}{\uparrow} \quad \underset{1-\tanh^2(H_z)}{\uparrow} \quad \underset{H_N}{\uparrow}$$

$$H_z = \gamma H_N + \beta$$

$$H_N = \frac{H_0 - \mu}{\sqrt{\sigma^2}}$$

$$H_0 = E \cdot W_1 + b_1$$

Next we need

$$\frac{\partial L}{\partial H_N} = \frac{\partial L}{\partial \hat{\gamma}} \cdot \frac{\partial \hat{\gamma}}{\partial H_a} \cdot \frac{\partial H_a}{\partial H_z} \cdot \frac{\partial H_z}{\partial H_N}$$

$$\underset{\gamma}{\uparrow}$$

Now $\frac{\partial L}{\partial \sigma^2} = \frac{\partial L}{\partial H_N} \cdot \frac{\partial H_N}{\partial \sigma^2} \quad = -\frac{1}{2} \cdot (H_0 - \mu)(\sigma^2 + \epsilon)^{-3/2}$

We also need $\dfrac{\partial L}{\partial \mu} = \dfrac{\partial L}{\partial H_N} \cdot \dfrac{\partial H_N}{\partial \mu} + \dfrac{\partial L}{\partial \sigma^2} \cdot \dfrac{\partial \sigma^2}{\partial \mu}$

$\sigma^2 = \dfrac{1}{N}\Sigma(H_0-\mu)^2$

$\dfrac{-1}{\sqrt{\sigma^2+\epsilon}}$ $\qquad$ $\dfrac{-2(H_0-\mu)}{n}$ from

And we need $\dfrac{\partial L}{\partial H_0} = \dfrac{\partial L}{\partial H_N}\cdot\dfrac{\partial H_N}{\partial H_0} + \dfrac{\partial L}{\partial \sigma^2}\cdot\dfrac{\partial \sigma^2}{\partial H_0} + \dfrac{\partial L}{\partial \mu}\cdot\dfrac{\partial \mu}{\partial H_0}$

$\mu = \frac{1}{N}\Sigma H_0$ from

$\dfrac{1}{\sqrt{\sigma^2+\epsilon}}$ $\qquad$ $\dfrac{-2(H_0-\mu)}{n}$ $\qquad$ $1/n$

So $\quad \dfrac{\partial L}{\partial H_0} = \dfrac{\partial L}{\partial H_N}\cdot\dfrac{1}{\sqrt{\sigma^2+\epsilon}} + \dfrac{\partial L}{\partial \sigma^2}\cdot\dfrac{-2(H_0-\mu)}{n} + \dfrac{\partial L}{\partial \mu}\cdot\dfrac{1}{n}$

$\dfrac{\partial L}{\partial b_1} = \dfrac{\partial L}{\partial H_0}\cdot\dfrac{\partial H_0}{\partial b_1}$ $\qquad H_0 = E\cdot W_1 + b_1$

$1$ from

$\not E\ \dfrac{\partial L}{\partial W_1} = \dfrac{\partial L}{\partial H_0}\cdot\dfrac{\partial H_0}{\partial W_1}$ $\qquad E$ $\qquad$ then $\dfrac{\partial L}{\partial E} = \dfrac{\partial L}{\partial H_0}\cdot\dfrac{\partial H_0}{\partial E}$ $\qquad W_1$

Now use $\dfrac{\partial L}{\partial E}$ to unembed the gradients into $C$ giving $\dfrac{\partial L}{\partial C}$

QED!

Proceed to update

$C, W_1, b_1, W_2, b_2, \gamma, \beta$