

Lecture 4: BackPropagation (BackP. Makemore Manuall)

X : #182625 elements of size 3
training
 $[0, 0, 0], [.. .]$
 $[0, 0, 25], [.. y]$



construct minibatch of size 32

X_b : #32



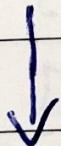
embedding

dimensionality (10)

C :

Vocab
size
(27) :

emb



Linear layer (= emb . $w_1 + b_1$)

h prebn



Batch normalization (= hprebn . b_{gain}
 $+ b_{bias}$)

h preact



tanh (non linearity)

h

Lecture 4: Backpropagation (BackP. Makemore Manohar)

X : #182625 elements of size 3
training

$[0, 0, 0]$, $[..]$

$[0, 0, 25]$, $[.., y]$



construct minibatch of size 32

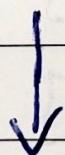
X_b : #32

↓ embedding

C : --- dimensionality (10)

Vocab
size
(27)

emb



Linear layer ($= \text{emb} \cdot w_1 + b_1$)

h prebn



Batch normalization ($= h_{\text{prebn}} \cdot \text{bgain} + \text{bbias}$)

h preact



tanh (non linearity)

h

h

↓ Linear layer 2 ($= h \cdot w_2 + b_2$)

logits

↓ Numerical stability ($\text{logits} - \max(\text{logits})$)

norm_logits

↓ Softmax (normalize = $\frac{\exp(\text{counts})}{\sum(\text{counts})}$)

probs

↓ log(probs)

loss

dloss

dlogprobs

How does logprobs influence the loss?

loss = $-\log(\text{probs})[\text{range}(n), Y_b].\text{mean}()$

where $n = \text{batch_size} = 32$

Y_b = array of correct indices

logprobs.size = (32, 27)

Iterating down the rows of logprobs with range(n), we pluck out the index at the column specified by Y_b

$$y_b = [8, 14, 15 \dots]$$

In 0th row, we take 8th column
 1st row, " 14th "

logprobs [range(n), y_b] contains the log probabilities of the correct next character in a sequence. Size = 32 (batch-size)

$$y_b (\text{answer}) = 27$$

logprobs	[1]	[2]			
	0'1	0'3	0'2	...	
↑ to	0'3	0'1	0'17	...	
↑ (32)	0'05	0'6	0'24	0	
		0			

$$\text{loss} = -\log(\text{logprobs}[0, 6]), \log(\text{logprobs}[1, 17]), \log(\text{logprobs}[2, 4], \dots, \text{mean})$$

Size of logprobs and dlogprobs are the same and equal to [32, 27]

$$\text{logits} \rightarrow [32, 27]$$

$$\text{logit-maxes} \rightarrow [32, 1]$$

$$\text{logits} = h @ W_2 + b_2$$
$$[32, 64] [64, 27] [27]$$

Due to broadcasting in Python, b_2 will replicate and become a $[32, 27]$ (previously it was considered as $[1, 27]$)

Toy example

$$d = a \cdot b + c$$

$$[3 \times 4] [4 \times 2] [2]$$

$$= \begin{bmatrix} a_{11} & a_{12} & a_{13} & a_{14} \\ a_{21} & a_{22} & a_{23} & a_{24} \\ a_{31} & a_{32} & a_{33} & a_{34} \end{bmatrix} \cdot \begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \\ b_{31} & b_{32} \\ b_{41} & b_{42} \end{bmatrix} + \begin{bmatrix} c_1 \\ c_2 \end{bmatrix}$$

$$= \begin{bmatrix} a_{11} \cdot b_{11} + a_{12} \cdot b_{21} + a_{13} \cdot b_{31} + a_{14} \cdot b_{41} & a_{11} \cdot b_{12} + \\ a_{21} \cdot b_{11} + a_{22} \cdot b_{21} + a_{23} \cdot b_{31} + \dots & a_{21} \cdot b_{12} + \\ a_{31} \cdot b_{11} + a_{32} \cdot b_{21} + \dots & a_{31} \cdot b_{12} + \end{bmatrix}$$

$$[3 \times 2]$$

~~$$= \begin{bmatrix} d_{11} & d_{12} \\ d_{21} & d_{22} \\ d_{31} & d_{32} \end{bmatrix}$$~~

$$\begin{bmatrix} d_{11} & d_{12} \\ d_{21} & d_{22} \\ d_{31} & d_{32} \end{bmatrix}$$

$$d_{11} = a_{11} \cdot b_{11} + a_{12} \cdot b_{21} + a_{13} \cdot b_{31} + a_{14} \cdot b_{41} + c_1$$

$$d_{12} = a_{11} \cdot b_{12} + a_{12} \cdot b_{22} + a_{13} \cdot b_{32} + a_{14} \cdot b_{42} + c_2$$

$$d_{21} = a_{21} \cdot b_{11} + a_{22} \cdot b_{21} + a_{23} \cdot b_{31} + a_{24} \cdot b_{41} + c_1$$

$$d_{22} = a_{21} \cdot b_{12} + a_{22} \cdot b_{22} + a_{23} \cdot b_{32} + a_{24} \cdot b_{42} + c_2$$

$$d_{31} = a_{31} \cdot b_{11} + a_{32} \cdot b_{21} + a_{33} \cdot b_{31} + a_{34} \cdot b_{41} + c_1$$

$$d_{32} = a_{31} \cdot b_{12} + a_{32} \cdot b_{22} + a_{33} \cdot b_{32} + a_{34} \cdot b_{42} + c_2$$

$$\frac{\partial L}{\partial a_{11}} = \frac{\partial L}{\partial d_{11}} \cdot b_{11} + \frac{\partial L}{\partial d_{12}} \cdot b_{12}$$

$$\frac{\partial L}{\partial a_{12}} = \frac{\partial L}{\partial d_{11}} \cdot b_{21} + \frac{\partial L}{\partial d_{12}} \cdot b_{22}$$

$$\frac{\partial L}{\partial a_{13}} = \frac{\partial L}{\partial d_{11}} \cdot b_{31} + \frac{\partial L}{\partial d_{12}} \cdot b_{32}$$

$$\frac{\partial L}{\partial a_{14}} = \frac{\partial L}{\partial d_{11}} \cdot b_{41} + \frac{\partial L}{\partial d_{12}} \cdot b_{42}$$

$$\frac{\partial L}{\partial a_{21}} = \frac{\partial L}{\partial d_{21}} \cdot b_{11} + \frac{\partial L}{\partial d_{22}} \cdot b_{12}$$

$$\frac{\partial L}{\partial a_{22}} = \frac{\partial L}{\partial d_{21}} \cdot b_{21} + \frac{\partial L}{\partial d_{22}} \cdot b_{22}$$

$$\frac{\partial L}{\partial a_{23}} = \frac{\partial L}{\partial d_{21}} \cdot b_{31} + \frac{\partial L}{\partial d_{22}} \cdot b_{32}$$

$$\frac{\partial L}{\partial a_{24}} = \frac{\partial L}{\partial d_{21}} \cdot b_{41} + \frac{\partial L}{\partial d_{22}} \cdot b_{42}$$

$$\frac{\partial L}{\partial a_{31}} = \frac{\partial L}{\partial d_{31}} \cdot b_{11} + \frac{\partial L}{\partial d_{32}} \cdot b_{12}$$

$$\frac{\partial L}{\partial a_{32}} = \frac{\partial L}{\partial d_{31}} \cdot b_{21} + \frac{\partial L}{\partial d_{32}} \cdot b_{22}$$

$$\frac{\partial L}{\partial a_{33}} = \frac{\partial L}{\partial d_{31}} \cdot b_{31} + \frac{\partial L}{\partial d_{32}} \cdot b_{32}$$

$$\frac{\partial L}{\partial a_{34}} = \frac{\partial L}{\partial d_{31}} \cdot b_{41} + \frac{\partial L}{\partial d_{32}} \cdot b_{42}$$

$$\frac{\partial L}{\partial a} = \begin{bmatrix} \frac{\partial L}{\partial a_{11}} & \frac{\partial L}{\partial a_{12}} & \frac{\partial L}{\partial a_{13}} & \frac{\partial L}{\partial a_{14}} \\ \frac{\partial L}{\partial a_{21}} & \frac{\partial L}{\partial a_{22}} & \frac{\partial L}{\partial a_{23}} & \frac{\partial L}{\partial a_{24}} \\ \frac{\partial L}{\partial a_{31}} & \frac{\partial L}{\partial a_{32}} & \frac{\partial L}{\partial a_{33}} & \frac{\partial L}{\partial a_{34}} \end{bmatrix}$$

$$= \begin{bmatrix} \frac{\partial L}{\partial d_{11}} & \frac{\partial L}{\partial d_{12}} \\ \frac{\partial L}{\partial d_{21}} & \frac{\partial L}{\partial d_{22}} \\ \frac{\partial L}{\partial d_{31}} & \frac{\partial L}{\partial d_{32}} \end{bmatrix} \times \begin{bmatrix} b_{11} & b_{21} & b_{31} & b_{41} \\ b_{12} & b_{22} & b_{32} & b_{42} \end{bmatrix}$$

$$= \frac{\partial L}{\partial d} @ b^T$$

$$\frac{\partial L}{\partial b} = \dots = a^T @ \frac{\partial L}{\partial d}$$

$$\frac{\partial L}{\partial c_1} = \frac{\partial L}{\partial d_{11}} \cdot 1 + \frac{\partial L}{\partial d_{21}} \cdot 1 + \frac{\partial L}{\partial d_{31}} \cdot 1 + \cancel{\frac{\partial L}{\partial d_{41}}} \cancel{+ \frac{\partial L}{\partial d_{51}}}$$

$$\frac{\partial L}{\partial c_2} = \frac{\partial L}{\partial d_{12}} \cdot 1 + \frac{\partial L}{\partial d_{22}} \cdot 1 + \frac{\partial L}{\partial d_{32}} \cdot 1$$

$$\frac{\partial L}{\partial c} = \begin{bmatrix} \frac{\partial L}{\partial c_1} \\ \frac{\partial L}{\partial c_2} \end{bmatrix} = \begin{bmatrix} \frac{\partial L}{\partial d_{11}} + \frac{\partial L}{\partial d_{21}} + \frac{\partial L}{\partial d_{31}} \\ \frac{\partial L}{\partial d_{12}} + \frac{\partial L}{\partial d_{22}} + \frac{\partial L}{\partial d_{32}} \end{bmatrix}$$

$$= \frac{\partial L}{\partial d} \cdot \text{sum}(0) \quad \text{sum across columns of } \frac{\partial L}{\partial d}$$

Backpropagate through line

$$\underline{\text{emb}} = C[\underline{Xb}], \text{ where}$$

Xb = batch containing 32 examples, size $[32 \times 3]$

C = matrix 27×10 embedding

$$\text{emb} = C[Xb] \rightarrow \text{size } [32 \times 3 \times 10]$$

$$\text{embcat} = \text{emb. view}(\underbrace{\text{emb.shape}[0]}_{32}, -1) \text{ concatenate}$$

$$\text{embcat. shape} = 32 \times 30$$

↑
the shape of -1
is inferred from other
dimensions

examples
 emb : 32 ~~characters~~, 3 characters and each of
them has 10 dimensional embedding

This was achieved by looking at the
lookup table C (27 possible characters). We
look at the rows specified by tensor Xb
 $^{10\text{-dim}}$

$$Xb = [1, 1, 4],$$

$$[18, 14, 1],$$

$$[0, 3, 1],$$

:

These integers specify which row of C we
want to pluck out

$$\text{emb} = \underline{\underline{C[[1, 1, 4]]}},^{(3 \times 10)}$$
$$\underline{\underline{C[[18, 14, 1]]}},^{(3 \times 10)}$$
$$\vdots$$

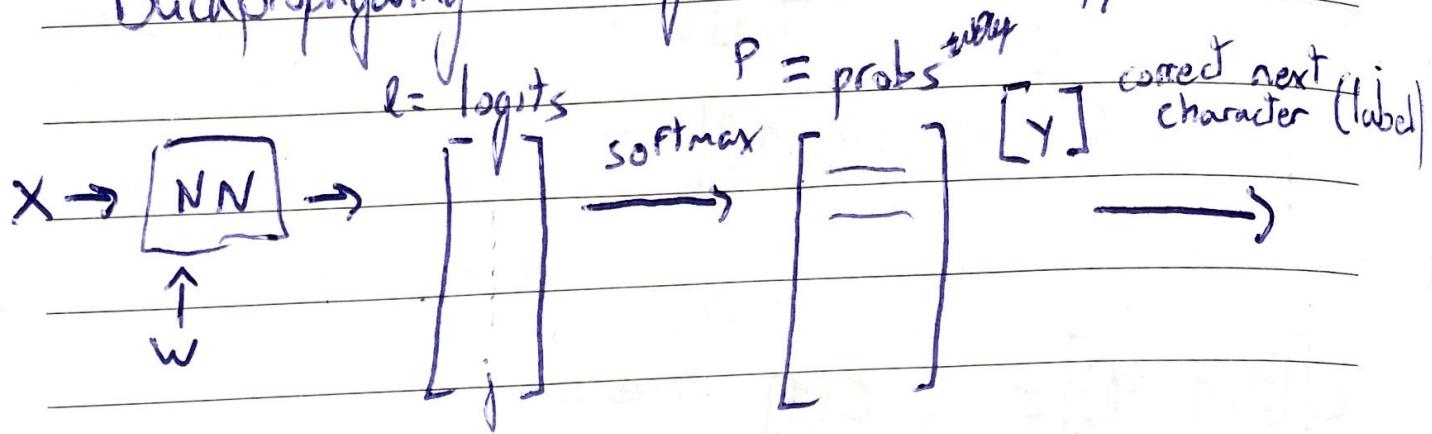
$$\underline{\underline{C[[1, 1, 4]]}} = \begin{bmatrix} C[1], & (\text{size } 10) \\ (3 \times 10) & C[1], (\text{size } 10) \\ & C[4] \end{bmatrix}$$

emb contains the gradients of the plucked rows

Which row of C did the 10-dimensional embeddings of emb come from? Deposit this into dC

Some rows of C are used multiple times, so the gradient have to be added

Backpropagating through cross entropy Function



$$\rightarrow \text{probs}[y] \xrightarrow{-\log} -\log \text{probs}[y] \xrightarrow{\text{avg}} \text{loss}$$

$$\begin{aligned} \text{loss} &= -\log P_y \\ &= -\log \frac{e^{l_y}}{\sum_j e^{l_j}} \end{aligned}$$

$$P_i = \frac{e^{l_i}}{\sum_j e^{l_j}} \quad (\text{softmax})$$

$\frac{d \text{loss}}{d l_i} \equiv$ derivative of the loss w.r.t. the i th logit

$$\frac{d}{d l_i} \left[-\log \frac{e^{l_y}}{\sum_j e^{l_j}} \right]$$

l indexed with specific label y

$$\begin{aligned} \frac{d}{dx} \log x &= \frac{1}{x} \\ &= -\frac{\sum_j e^{l_j}}{e^{l_y}} \cdot \frac{d}{d l_i} \left[\frac{e^{l_y}}{\sum_j e^{l_j}} \right] \end{aligned}$$

(product rule $\rightarrow (u \cdot v)' = u' \cdot v + u \cdot v'$)

(power rule $\rightarrow (x^n)' = n x^{n-1}$)

if $i \neq y$

$$u = e^{l_y} \quad v = \frac{1}{\sum_j e^{l_j}}$$

$$\frac{\text{dloss}}{\text{d}l_i} = \frac{-\sum_j e^{l_j}}{e^{l_y}} \left[0 \cdot \frac{1}{\sum_j e^{l_j}} + -\frac{e^{l_y} \cdot e^{l_i}}{(\sum_j e^{l_j})^2} \right]$$

$$\text{since } \frac{du}{dl_i} = 0$$

$$\begin{aligned} \frac{dv}{dl_i} &= \frac{d}{dl_i} \left(\frac{1}{\sum_j e^{l_j}} \right) \\ &= \frac{d}{dl_i} \left(\frac{1}{e^{l_1} + e^{l_2} + \dots + e^{l_i} + \dots + e^{l_j}} \right) \end{aligned}$$

$$\left(\text{recall } \frac{d}{dx} \left(\frac{1}{x} \right) = -\frac{1}{x^2} \right)$$

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

$$\left(\text{recall } \frac{d}{dx} e^x = e^x \right)$$

$$= -\frac{1}{(\sum_j e^{l_j})^2} \cdot \cancel{\frac{d}{dl_i}} (0 + 0 + \dots + e^{l_i} + 0 + 0)$$

$$= -\frac{e^{l_i}}{(\sum_j e^{l_j})^2}$$

If $i \neq y$

$$\frac{d\text{loss}}{dl_i} = -\frac{\sum e^{l_j}}{e^{l_y}} \left[0 \cdot \frac{1}{\sum e^{l_j}} + -\frac{e^{l_y} \cdot e^{l_i}}{(\sum e^{l_j})^2} \right]$$

$$= + \frac{(\sum e^{l_j}) \cdot e^{l_y} \cdot e^{l_i}}{e^{l_y} \cdot (\sum e^{l_j})^2} = \frac{e^{l_i}}{\sum e^{l_j}} \stackrel{\text{softmax}}{=} P_i$$

If $i = y$

$$u \cdot v^* + u \cdot v'$$

$$\frac{d\text{loss}}{dl_i} = -\frac{\sum e^{l_j}}{e^{l_y}} \left[\frac{e^{l_y}}{\sum e^{l_j}} + -\frac{e^{l_y} \cdot e^{l_i}}{(\sum e^{l_j})^2} \right]$$

(since $\frac{du}{dl_i} = e^{l_y}$)

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

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

Backpropagating through batch Normalization
in a single line

Calculate ~~$\frac{dh_{preact}}{dh_{prebn}}$~~ given dh_{preact}

$$h_{prebn} = x_i$$

$$h_{preact} = y_i$$

batch Norm
paper

$$\mu = \frac{1}{m} \sum_{i=1}^m x_i$$

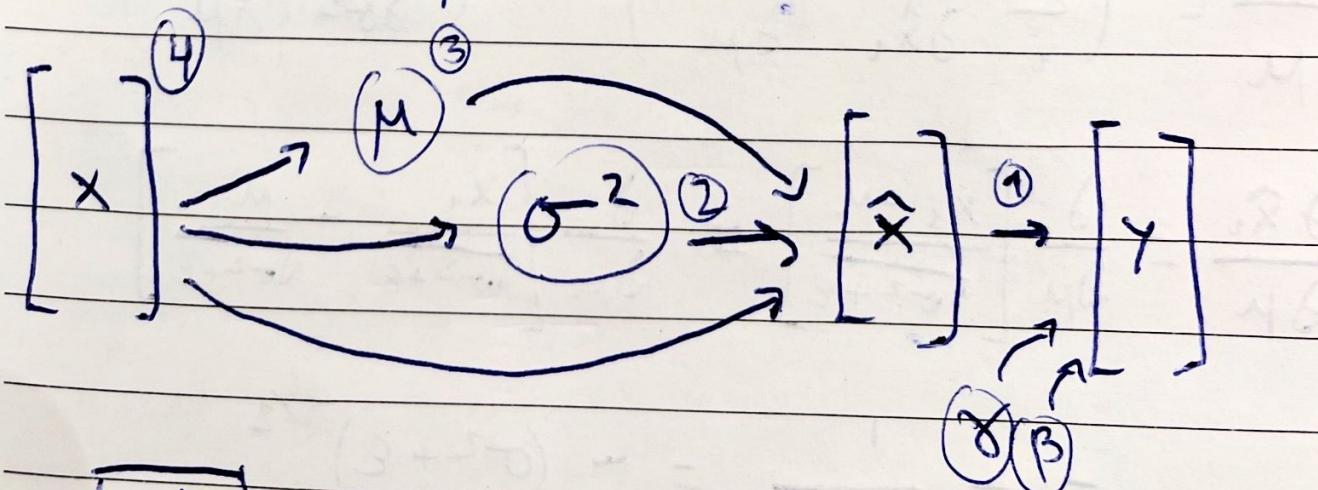
$$\sigma^2 = \frac{1}{m-1} \sum_{i=1}^m (x_i - \mu)^2$$

Bessel's correction.

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

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

we have $\frac{\partial L_{\text{oss}}}{\partial y_i}$; we want $\frac{\partial L_{\text{oss}}}{\partial x_i}$



$$\textcircled{1} \quad \frac{\partial L}{\partial \hat{x}_i} = \frac{\partial L}{\partial y_i} \cdot \gamma$$

$$\textcircled{2} \quad \frac{\partial L}{\partial \sigma^2} = \underbrace{\sum_i \frac{\partial L}{\partial \hat{x}_i} \frac{\partial \hat{x}_i}{\partial \sigma^2}}_{\text{sigma}^2 \text{ propagates through many } \hat{x}_i}$$

$$= \frac{\# \frac{\partial L}{\partial y_i} \cdot \gamma \cdot \frac{\partial}{\partial \sigma^2} \frac{x_i - \mu}{\sqrt{\sigma^2 + \epsilon}}}{\sum_i}$$

$$= \gamma \cdot \sum_i \frac{\partial}{\partial \sigma^2} \left[(x_i - \mu) \cdot (\sigma^2 + \epsilon)^{-\frac{1}{2}} \right] \cdot \frac{\partial L}{\partial y_i}$$

$$= \gamma \cdot \sum_i -\frac{1}{2} \cdot (x_i - \mu) (\sigma^2 + \epsilon)^{-\frac{3}{2}} \cdot \frac{\partial L}{\partial y_i}$$

$$\boxed{\frac{\partial L}{\partial \sigma^2}} = -\frac{1}{2} \gamma \sum_i \frac{\partial L}{\partial y_i} (x_i - \mu) (\sigma^2 + \epsilon)^{-\frac{3}{2}}$$

$$\textcircled{3} \quad \frac{\partial L}{\partial \mu} \quad ? \quad \begin{array}{l} \text{Multiple } \mu \text{ influence } \hat{x}_i \\ \text{Only one single } \mu \text{ influences } \sigma^2 \end{array} \quad \left. \begin{array}{l} \# 33 \\ \mu \\ (\text{batch}+1) \end{array} \right\}$$

$$\frac{\partial L}{\partial \mu} = \left(\sum_i \frac{\partial L}{\partial \hat{x}_i} \frac{\partial \hat{x}_i}{\partial \mu} \right) + \left(\frac{\partial L}{\partial \sigma^2} \frac{\partial \sigma^2}{\partial \mu} \right)$$

$$\frac{\partial \hat{x}_i}{\partial \mu} = \frac{\partial}{\partial \mu} \left[\frac{x_i - \mu}{\sqrt{\sigma^2 + \epsilon}} \right] = \frac{\partial}{\partial \mu} \left[\frac{x_i}{\sqrt{\sigma^2 + \epsilon}} - \frac{\mu}{\sqrt{\sigma^2 + \epsilon}} \right]$$

$$= -\frac{1}{\sqrt{\sigma^2 + \epsilon}} = -(\sigma^2 + \epsilon)^{-\frac{1}{2}}$$

$$\frac{\partial \sigma^2}{\partial \mu} = \frac{\partial}{\partial \mu} \left[\frac{1}{M-1} \sum_{i=1}^M (x_i - \mu)^2 \right]$$

$$= \frac{2}{M-1} \sum_{i=1}^M (x_i - \mu) = -1$$

in BatchNorm
it holds that

$$= \frac{-2}{M-1} \sum x_i - \sum \mu$$

$\mu = \text{average}$

$$= \frac{-2}{M-1} (\underbrace{m\mu - M\mu}_{\text{}}); \text{ since } \mu = \frac{1}{M} \sum x_i$$

$$= \frac{-2}{M-1} \cdot 0 = 0$$

$$\boxed{\frac{\partial L}{\partial \mu}} = \sum_i \frac{\partial L}{\partial y_i} \cdot \gamma \cdot -(\sigma^2 + \epsilon)^{-\frac{1}{2}} + 0$$

$\left(\frac{\partial L}{\partial \sigma^2} \frac{\partial \sigma^2}{\partial \mu} \right)$

(4)

x is contributed

The gradient of x is the sum of 3 parts:

gradient coming from μ , from σ^2 and from

a single x_i . $\hat{x}_i = \frac{x_i - \mu}{\sqrt{\sigma^2 + \epsilon}}$ There is a

bijection relation between x_i and \hat{x}_i , despite

x and \hat{x} being 32 sized vectors
(batch size)

$$\textcircled{4} \quad \frac{\partial L}{\partial x_i} = \underbrace{\frac{\partial L}{\partial \hat{x}_i} \frac{\partial \hat{x}_i}{\partial x_i}}_{?} + \underbrace{\frac{\partial L}{\partial \mu} \frac{\partial \mu}{\partial x_i}}_{?} + \underbrace{\frac{\partial L}{\partial \sigma^2} \frac{\partial \sigma^2}{\partial x_i}}_{?}$$

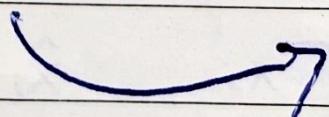
$$\frac{\partial \hat{x}_i}{\partial x_i} = \frac{\partial}{\partial x_i} \left[\frac{x_i - \mu}{\sqrt{\sigma^2 + \varepsilon}} \right] = \frac{\partial}{\partial x_i} \left[\frac{x_i}{\sqrt{\sigma^2 + \varepsilon}} - \frac{\mu}{\sqrt{\sigma^2 + \varepsilon}} \right]$$

$$= \frac{1}{\sqrt{\sigma^2 + \varepsilon}}$$

$$\frac{\partial \mu}{\partial x_i} = \frac{\partial}{\partial x_i} \left[\frac{1}{m} \sum x_i \right] = \frac{1}{m}$$

$$\frac{\partial \sigma^2}{\partial x_i} = \frac{\partial}{\partial x_i} \left[\frac{1}{m-1} \sum (x_i - \mu)^2 \right] = \frac{2}{m-1} (x_i - \mu) \cdot 1$$

Join
 Add everything
~~to calculate~~



$$\begin{aligned}
 \frac{\partial L}{\partial x_i} &= \left(\frac{\partial L}{\partial y_i} \cdot \gamma \right) \cdot (\sigma^2 + \epsilon)^{-\frac{1}{2}} + \\
 &+ \left(-\sum_j \frac{\partial L}{\partial y_j} \cdot \gamma \cdot (\sigma^2 + \epsilon)^{-\frac{1}{2}} \right) \cdot \frac{1}{m} + \\
 &+ \left(-\frac{1}{2} \gamma \sum_j \frac{\partial L}{\partial y_j} (x_j - \mu) (\sigma^2 + \epsilon)^{-\frac{3}{2}} \right) \cdot \frac{1}{m-1} (x_i - \mu) \\
 &= (\sigma^2 + \epsilon)^{-\frac{1}{2}} \left(\frac{\partial L}{\partial y_i} \cdot \gamma \right) + (\sigma^2 + \epsilon)^{-\frac{1}{2}} \left(-\frac{1}{m} \sum_j \frac{\partial L}{\partial y_j} \cdot \gamma \right) + \\
 &+ (\sigma^2 + \epsilon)^{-\frac{1}{2}} \left(-\gamma \sum_j \frac{\partial L}{\partial y_j} \frac{(x_j - \mu)}{\sqrt{\sigma^2 + \epsilon}} \right) \cdot \left(\frac{1}{m-1} \frac{x_i - \mu}{\sqrt{\sigma^2 + \epsilon}} \right)
 \end{aligned}$$

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

$$= \frac{(\sigma^2 + \epsilon)^{-\frac{1}{2}}}{m} \left[m \gamma \frac{\partial L}{\partial y_i} - \sum_j \frac{\partial L}{\partial y_j} \gamma - \gamma \left(\sum_j \frac{\partial L}{\partial y_j} \hat{x}_j \right) \left(\frac{m}{m-1} \hat{x}_i \right) \right]$$

$$= \frac{\gamma (\sigma^2 + \epsilon)^{-\frac{1}{2}}}{m} \left[m \frac{\partial L}{\partial y_i} - \sum_j \frac{\partial L}{\partial y_j} - \frac{m}{m-1} \hat{x}_i \sum_j \frac{\partial L}{\partial y_j} \hat{x}_j \right]$$

already have them from previous steps in backprop.

come from forward pass