

$$\text{emb} = ([x_b])$$

$X_b = n \times k$: n : no of samples
 k : context size

m : no of features

V : Vocabulary size

$$C : V \times m$$

$$\text{emb} : n \times k \times m$$

$\text{embcat} : n \times (k \times m)$ convert 3D to
 2D

$$\text{hprebn} = \text{embcat } W_1 + b_1 : n \times n_1$$

$$W_1 : (k \times m, n_1) \quad b_1 : (n, n_1)$$

$$\text{bmean}_i = \frac{1}{n} \sum_{j=1}^n \text{hprebn}_{ij}$$

$$\begin{pmatrix} x_{11} & \dots & x_{1n_1} \\ \vdots & \ddots & \vdots \\ x_{n1} & \dots & x_{nn_1} \end{pmatrix}$$

mean across
 n samples

$$bndiff = h_{\text{preact}} - b_{\text{mean}} : n \times n_1$$

$$bndiff_2 = bndiff^2 : n \times n_1$$

$$bnvar = \frac{1}{n-1} \sum_{i=1}^n bndiff_{ij}^2 : 1 \times n_1$$

$$bnvar_inv = \frac{1}{\sqrt{bnvar + \epsilon}} : 1 \times n_1$$

ϵ is set at $1e^{-5}$

$$bnraw = bndiff \times bnvar_inv : n \times n_1$$

$$h_{\text{preact}} = bngain \times bnraw + bnbias \\ \begin{matrix} n \times n_1 \\ 1 \times n_1 \end{matrix} \quad \begin{matrix} n \times n_1 \\ (\delta) \end{matrix} \quad \begin{matrix} n \times n_1 \\ (\beta) \end{matrix}$$

$$h = \tanh(h_{\text{preact}}) : n \times n_1$$

$$\text{logits} = h W_2 + b_2 : n \times V$$

$$W_2 : n_1 \times V \quad b_2 : n \times V$$

Cross entropy losses : Note this is all element wise operations.

$$\text{logit_mats} = \max_{j=1}^V (\text{logits}_{ij}) : n \times 1$$

$$\text{norm_logits} = \text{logits} - \text{logit_mats} : n \times V$$

$$\text{counts} = \exp(\text{norm_logits}) : n \times V$$

$$\text{count_sum} = \sum_{j=1}^V \text{counts}_{(j)} : n \times 1$$

$$\text{count_sum_inv} = \frac{1}{\text{count_sum}} : n \times 1$$

→ implicit broadcasting, ^{count sum} duplicate across all columns

$$\text{probs} = (\text{counts} \times \text{count_sum_inv})^{n \times V}$$

$$\text{log_probs} = \log(\text{probs}) : n \times V$$

$$L = -\frac{1}{n} \sum_{i=1}^n \log \text{probs}_{ij}$$

j = index of 1 in each row of Y_b

$$\text{Y}_b := \begin{pmatrix} 0 & 1 & \cdot & \cdot \\ 0 & 0 & \cdot & \cdot \\ \vdots & \vdots & \ddots & \ddots \end{pmatrix} \in n \times V$$

l appears on the character output index

backproping

$$d\logprobs = \frac{dL}{d\logprobs}$$

$$= \begin{pmatrix} 0 & \frac{-1}{n} & \dots & 0 \\ 0 & \frac{-1}{n} & \dots & 0 \\ \vdots & & & \frac{-1}{n} \\ 0 & \dots & & \frac{-1}{n} \end{pmatrix},$$

$n \times v$

$$dprobs = d\logprobs \times \underbrace{\frac{1}{\text{probs}}}_{: n \times v} : n \times v$$

(elementwise multiplication)

$$\text{decountsum} = \sum_{i=1}^v dprobs \times \underset{n \times 1}{\text{counts}}$$

(sum across columns)
reflect the broadcasting
operation

$$d\text{counts} = d\text{props} \times \text{counts_sum_inv}$$

$$n \times v \quad \times \quad n \times 1$$

→ broadcasting

(counts is used multiple times)

→ still missing some gradient)

$$d\text{counts_sum} = d\text{counts_sum_inv} \times \left(\frac{-1}{\text{counts_sum}^2} \right)$$

$$n \times 1 \quad n \times 1$$

$$d\text{counts} = d\text{counts_sum} \times$$

$$n \times v \quad n \times 1 \quad (n \times v)$$

$$\begin{pmatrix} 1 & & & \\ & 1 & & \\ & & \ddots & \\ & & & 1 \end{pmatrix}$$

reflect the summing
up across columns

→ rebroadcasting back
to $n \times v$

$$d\text{norm_logits} = d\text{counts} \times \exp(\text{norm_logits})$$

$$d\text{logits_maxes} = - \sum_{j=1}^n d\text{logits}_{i,j}$$

$n \times \checkmark \rightarrow \text{need}$
to sum
across columns.

$$d\text{logits} = d\text{norm_logits} + d\text{logits_maxes} \times \text{max_indices}$$

$$\text{logits}_{\text{max}} = \max \left(\frac{d\text{logits}_{\text{maxes}}}{d\text{logits}} : \begin{matrix} 0 & 1 & 0 \\ \downarrow & & \\ \text{indices of the} & & \\ \text{max} & & \end{matrix} \right)$$

$$dh = \text{d logits. } w_2^T$$

$n \times n_1 \quad \cdot \quad \cdot \quad \cdot \quad \frac{d \text{ logits}}{dh}$

$$n \times V \times V \times n_1$$

$$dh_{\text{preact}} = dh \times (1 - h^2)$$

$n \times n_1 \quad (\text{elementwise})$

$$dh_{\text{gain}} = \sum_{i=1}^n dh_{\text{preact}, i} \quad |n \text{ raw}_i$$

$1 \times n_1$

$$dh_{\text{raw}} = dh_{\text{preact}} \times h_{\text{gain}}$$

$n \times n_1 \quad n \times n_1 \quad 1 \times n_1$

(broad casting)

$$\text{bndvar} = \frac{1}{n-1} \sum_{i=1}^n \text{bnddigg}_{i,j}$$

sum across row

$$\text{dbnddigg} = \text{dbndvar} \times \frac{1}{n-1} \text{bn}$$

$$(x_t - \bar{x})^2 = (x_t - \mu)^2$$

$$- 2(\bar{x} - \mu)(x_t - \mu)$$

$$+ (\bar{x} - \mu)^2$$

$$\sum_{t=1}^T (x_t - \bar{x})^2$$

$$= \sum_{t=1}^T (x_t - \mu)^2$$

$$- 2(\bar{x} - \mu) \sum_{t=1}^T (x_t - \mu)$$

$$+ \sum_{t=1}^T (\bar{x} - \mu)^2$$

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

$$= 2n(\bar{x} - \mu)^2 + n(\bar{x} - \mu)^2$$

$$= \sum_{i=1}^n (x_i - \mu)^2 - n(\bar{x} - \mu)^2$$

$$E[n(\bar{x} - \mu)^2]$$

$$= n \left(\text{Var}(\bar{x} - \mu) \right)$$

$$= n (\text{Var}(\bar{y}))$$

$$= n \frac{\text{Var}(x_1 + y_2 + \dots + x_n)}{n}$$

$$= \frac{\sigma^2 n}{n^2} \cdot n$$

$$bnvar = \frac{1}{n-1} \sum_{i=1}^n bndiff^2_{ij}$$

j: $1 \rightarrow n$

$$= \frac{1}{n-1} \left(\begin{matrix} bndiff^2_{11} + bndiff^2_{21} \\ bndiff^2_{12} + bndiff^2_{22} \end{matrix} \right)^T$$

$$\rightarrow \frac{d bnvar}{d bndiff_2} = \frac{1}{n-1} \times \begin{pmatrix} 1 & \dots & 1 \\ 1 & \dots & 1 \end{pmatrix}$$

\downarrow

(elementwise derivative) $^{32 \times 64}$

$d hpteb = d bndiff.$ clone ()
 ↓
 important otherwise
 it will change the
 variables,

$demb = \text{dembcat} . \text{view}(\text{emb}, \text{shape})$
→ just need to undo
the operation.

$$\text{emb} = C[X]$$

$$n \times m \times h \quad V \times m \quad n \times h$$

C is a look up table for characters

in X

$$\rightarrow dC_{ij} += \text{demb}[i, j]$$

↓
multiple occurrences of the character

cross entropy backward loss

$$\text{loss} = -\frac{1}{n} \sum_{i=1}^n \log \frac{\exp(z_i y_i)}{\sum_{j=1}^k \exp(z_j y_j)}$$

$$\text{logit}_{iy} = z_i y_i \quad i: 1 \rightarrow N$$

$y_i \in Y_b$

$$\begin{pmatrix} z_{11} & z_{12} y_1 & z_{13} \\ \vdots & \ddots & \vdots \\ z_{n1} & z_{n2} y_n & z_{n3} \end{pmatrix} \xrightarrow{\text{softmax}} \begin{pmatrix} p_{11} & p_{12} y_1 & p_{13} \\ \vdots & \ddots & \vdots \\ p_{n1} & p_{n2} y_n & p_{n3} \end{pmatrix}$$

$$\xrightarrow{-\log} \begin{bmatrix} -\log p_{11} & -\log p_{12} y_1 & -\log p_{13} \\ \vdots & \ddots & \vdots \\ -\log p_{n1} & -\log p_{n2} y_n & -\log p_{n3} \end{bmatrix}$$

$$\frac{\partial \text{loss}}{\partial z_{ij}} = -\frac{1}{n} \frac{\partial}{\partial z_{ij}} -\log p_{iy_i}$$

$$p_{iy_i} = \frac{\exp(z_i y_i)}{\sum_{j=1}^k \exp(z_j y_j)}$$

$$\text{Case 1: } \frac{\partial}{\partial z_{ij}} y_i = \frac{\partial}{\partial z_{ij}} \left[-\log P_{ij} \right]$$

$$= - \frac{\exp(z_{ij}) (\sum_{j=1}^C \exp(z_{ij}) - \exp(z_{ij}))}{P_{ij} (\sum_{j=1}^C \exp(z_{ij}))^2}$$

$$= P_{ij} - 1 = P_{iy_i} - 1$$

Case 2:

$$j \neq y_i \quad \frac{\partial}{\partial z_{ij}} (-\log P_{iy_i})$$

$$= \frac{-1}{P_{iy_i}} \frac{-\exp(z_{iy_i})^2}{(\sum_{j=1}^C \exp(z_{ij}))^2}$$

$$= P_{iy_i} \quad \begin{matrix} y_i \text{ index} \\ \uparrow \end{matrix}$$

$$\rightarrow \frac{\partial \text{loss}}{\partial z_{ij}} = \frac{1}{n} (P_{iy_i} - \dots - P_{iy_i-1} - \dots - P_{iy_i})$$

\rightarrow dlogits

$$= \frac{1}{D} \begin{pmatrix} P_{1y_1} & \dots & P_{1y_{l-1}} & \dots & P_{1y_l} \\ \vdots & & \vdots & & \vdots \\ P_{ny_n} & \dots & P_{ny_{n-1}} & \dots & P_{ny_n} \end{pmatrix}$$

