

Given names, we want to make more names.

If we study ML for 10 hours, it will take

$$10 \text{ hrs} \times x \text{ days} = 10,000$$

1000 days of study

~ 2.73 years

~ 3 years of serious

commitment to become an

"ML expert"

← ← ← → ← →
e m m a .
1 2 3 4 5
f 6
g 7
h 8
i 9
j 10
k 11
l 12
m 13
n 14
o 15
p 16
q 17
r 18
s 19
t 20
u 21
v 22
w 23
x 24
y 25
z 26

[emma] → Q O →



How do we define loss?

- create bigram count model
- Convert counts into probabilities
- Use \Pr in loss fn

Bigram Count Model

end ^{2nd} char

a	b	ab	is	ab: How many does a w u / b
ö				
a	aa			

Start
1st char

What does a row vs column tell us?

i.e. Start-char = a a $\begin{bmatrix} a & b & c \\ 2/3 & 1/3 & 1/3 \end{bmatrix}$ $\frac{9}{23}$

$$aa = \frac{5}{23}$$

$$a \left[\frac{2}{23}, \frac{5}{23}, \frac{7}{23}, \frac{9}{23} \right]$$

emma
. 2 .3 .9

A: Pr of end char given a

$$\text{LOSS} = (1 - \Pr_r)$$

[aa

aa : 5

ab : 7

a ↗ aa
a ↘ ab

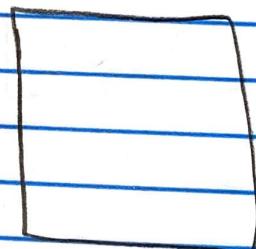
Given a, I should
be able to 'index' into
the count of an arbitrary
2nd char

Goal: Calculate LOSS for 'emma'

$\Pr_e = \frac{\text{count}}{\text{total}}$

e ↗ m
↓
 \Pr_e

ea
eb
:
em



Now that we have the counts,
let's construct the probabilities
 $\Pr(em) \Rightarrow \text{Given } e, \Pr(m) ?$

e ↗	a	b	c		
	1	2	15		

e ↗
↓
sum of all
counts

$\frac{\text{count}(-)}{\text{sum of counts}}$

char
sum = 0

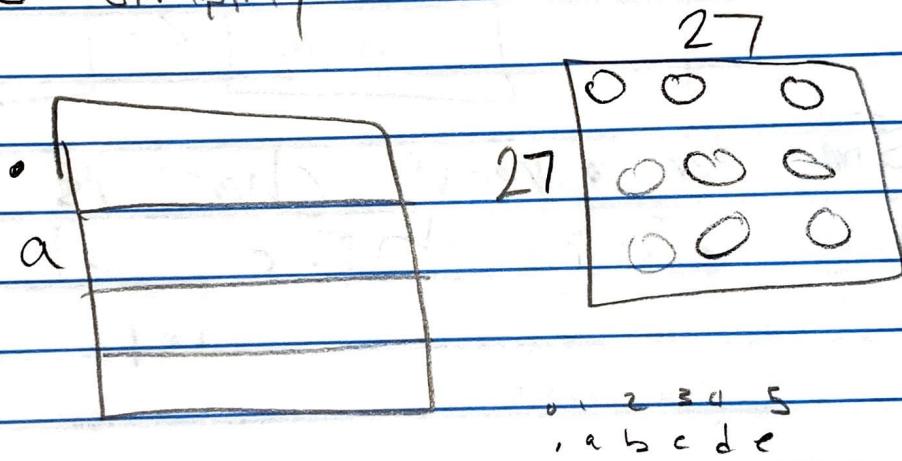
for each char
for each char

char sum += char_pairs[^{1st} char] + char_pairs[^{2nd} char]

char_sum['e'] = 15

OK I think we're double counting
the . character

Too much complexity w/ these loops,
lets simplify the data structure



0 2 3 4 5
. a b c d e

e m m a .
↑ ↑
. e += 1

charCounts[0, 5] += 1

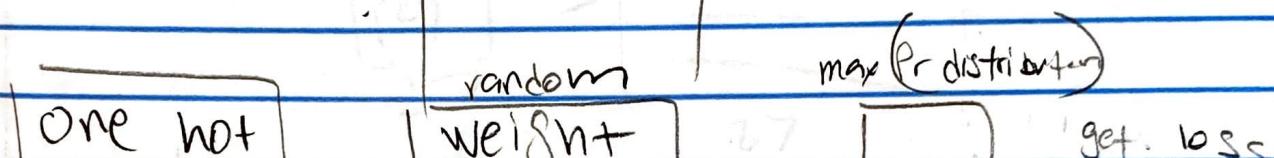
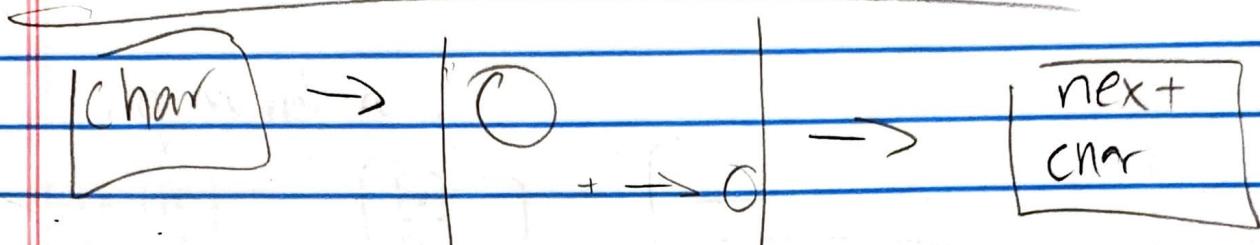
[0, first char]
[chr, 27]

How would it handle a . ?

A: manually accounted for this edge case

My loss fn looks ass

- lets normalize over name length



$$\begin{bmatrix} \cdot \\ \vdots \\ \cdot \end{bmatrix}_{26 \times 1} \times \begin{bmatrix} \cdot & \dots & \cdot \end{bmatrix}_{1 \times 1} \rightarrow \begin{bmatrix} \approx \\ \vdots \\ \approx \end{bmatrix}_{\text{max}(26 \times 1)}$$

<u>input</u>	<u>output</u> vector
e	$\begin{bmatrix} e \\ \vdots \\ e \end{bmatrix}$
m	$\begin{bmatrix} m \\ \vdots \\ m \end{bmatrix}$
m	$\begin{bmatrix} m \\ \vdots \\ m \end{bmatrix}$
a	$\begin{bmatrix} a \\ \vdots \\ a \end{bmatrix}$

- One hot : Popular way to encode integers
encode

$$\left[\begin{array}{c} \\ \\ \\ \\ \end{array} \right] \\ 5 \times 27$$

emm a.

$$\left[\begin{array}{c} \downarrow \\ \begin{bmatrix} i & o \\ o & i \\ o & o \\ o & o \\ o & o \end{array} \end{array} \right] \quad 27 \times 1 \quad 27 \times 1 \quad 5 \quad \left[\begin{array}{c} \cdot \\ e \\ m \\ n \\ c \end{array} \right] \quad 27 \quad 5 \quad \left[\begin{array}{c} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{array} \right] \quad 27 \quad 27 \quad \left[\begin{array}{c} w_1 \\ w_2 \\ \vdots \\ w_{27} \end{array} \right] \quad w$$

$$(5 \times 27) \times (27 \times 27)$$

$$\left[\begin{array}{c} -3 \\ 2 \\ 1 \\ 0 \\ \end{array} \right] \quad 5 \times 1$$

What do the #'s in the resulting array tell us?

A: Not much, we need to increase w

$$(5 \times 27) @ (27 \times 27)$$

5×27 Pr likelihood
of next char

• Use \log to squash values $[0,1]$ to a loss fn
 \Pr nice integers for

input: $(0 \text{ --- } 1) \Rightarrow [\log]$
Probabilities

Integers for a
loss fn

$$[-3 \text{ --- } 3]$$

$\Rightarrow f_n \rightarrow \Pr$ needs to

- * Make input numbers
 - all positive
 - $[0, 1]$
 - sum to 2

\Pr distributions

$$\frac{1}{\exp(x)}$$

?

Why do we assume

$$x_{enc} @ w = \log \text{counts}$$

- Previously log took us from a ~~Pr~~ distribution to integer values for a loss fn

What do logits represent?

$$\text{Counts} \rightarrow \Pr \rightarrow \text{loss}$$

[Counts] [log]
 $\frac{\text{counts}}{\text{counts.sum()}}$

$$\text{Math: } f_n(\Pr) \rightarrow \mathbb{R}$$

logits

$$\text{Softmax}(\mathbb{R}) \rightarrow \Pr$$

$\text{counts} = \exp(\mathbb{R})$ $\#$ make \mathbb{R} in positive range only

$\text{Prms} = \frac{\text{counts}}{\text{counts.sum()}}$ $\#$ squash to $[0,1]$ & sum to one

return \Pr

- Activation / normalization fn

- We have Pr matrix for the name Emma
- We're not taking the max Pr

$$[\text{Name}] \rightarrow \begin{bmatrix} \text{ch1} \\ \text{ch2} \\ \vdots \\ \text{ch}_n \end{bmatrix}_{xs} \rightarrow \begin{bmatrix} 0 \\ 0 \\ \text{ch1} \\ 0 \\ 0 \end{bmatrix} \xrightarrow{\text{one hot encoded}} * \begin{bmatrix} W \end{bmatrix}$$

$$= \begin{bmatrix} \text{Prob} \\ \text{abilities} \\ \text{of Bigram} \end{bmatrix}$$

$x_1 = \text{ch1}$ # input
 $y = \text{ch2}$ # output

(We are deriving
 Pr from W instead
 of counting)

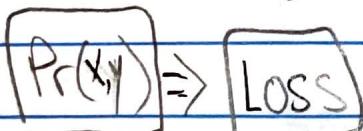
$$\begin{bmatrix} \text{input} \\ x \end{bmatrix} \times \begin{bmatrix} W \end{bmatrix} = \times \begin{bmatrix} \text{Pr}(xy) \end{bmatrix}$$

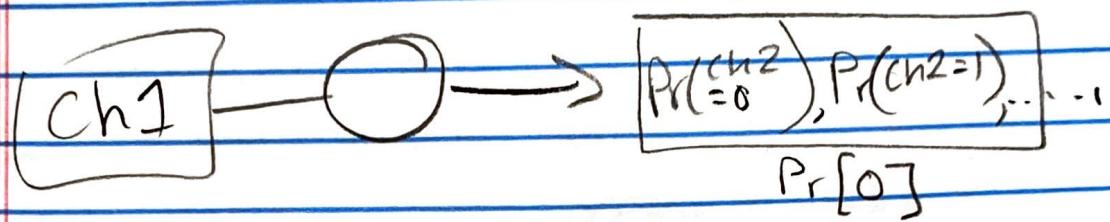
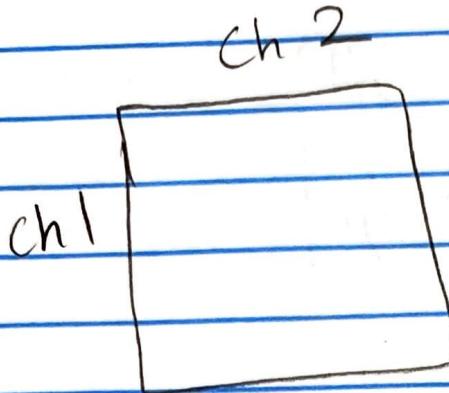
$y = \text{ch2}$ now we
 use y to
 index into
 our Pr matrix

- Pr = Represents the activation of output layer

$x_s = 1^{\text{st}}$ char
 in bigram

$y_s = 2^{\text{nd}}$ char
 in bigram





$$Pr(y|s) = Pr(y) = Pr_{i=0}$$

$y = 'e'$ # truth used in loss

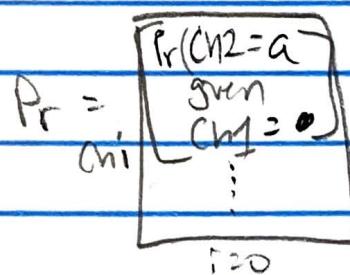
$$y_{-i} = 4 \quad \begin{matrix} xs \rightarrow xs[i] \\ ys \rightarrow ys[i] \end{matrix} \quad \begin{matrix} xs[0] = 26 \\ 4 \rightarrow e \end{matrix}$$

$Pr[i] [ys[i]]$ where $i = 0$

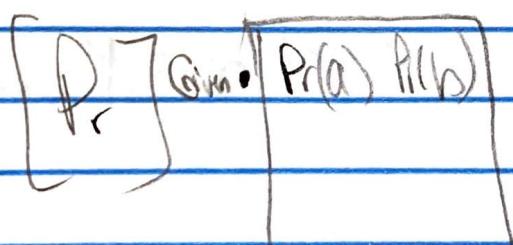
$$= Pr[0] ys[0] \quad Pr[i]$$

$$= [ch2] \quad 4$$

$\frac{5 \times 1}{}$



$\Pr \rightarrow Pr$ Given $ch1 = 0$, predict $ch[0]$ or name



$ch1 \quad \Pr(a=ch2)$

Predict $ch[1]$ or

name

$(ch1, ch2)$ Bigram for all chars in name



$$\Pr_{\text{dist}} = \alpha \left[\Pr(\text{ch2} = a \mid \text{ch1} = \cdot) \dots \Pr(\text{ch2} = y_i \mid \text{ch1} = \cdot) \right] \\ \left[\Pr(\text{ch2} = b \mid \text{ch1} = \cdot) \dots \Pr(\text{ch2} = z_j \mid \text{ch1} = \cdot) \right]$$

$$a \left[\Pr(\text{ch}_2=a | \text{ch}_1=a) \quad \Pr(\text{ch}_2=b | \text{ch}_1=a) \right]$$

$$b \left[\Pr(\text{ch}_2=a | \text{ch}_1=b) \quad \Pr(\text{ch}_2=b | \text{ch}_1=b) \right]$$

- Funny the diagonal of pr-dist is the $\Pr(\text{ch}_1 = \text{ch}_2) | (\text{ch}_2)$

$$[x] \cdot [w] =_x [y]$$

ch1 weight

$$\Pr(Y_i = \text{ch}_i \mid X_i)$$

$$\Pr(y=4 \mid x)$$

Use y to calculate loss

$$\text{Loss} = \sum_i [-y_i \log(p_i) - (1-y_i) \log(1-p_i)]$$