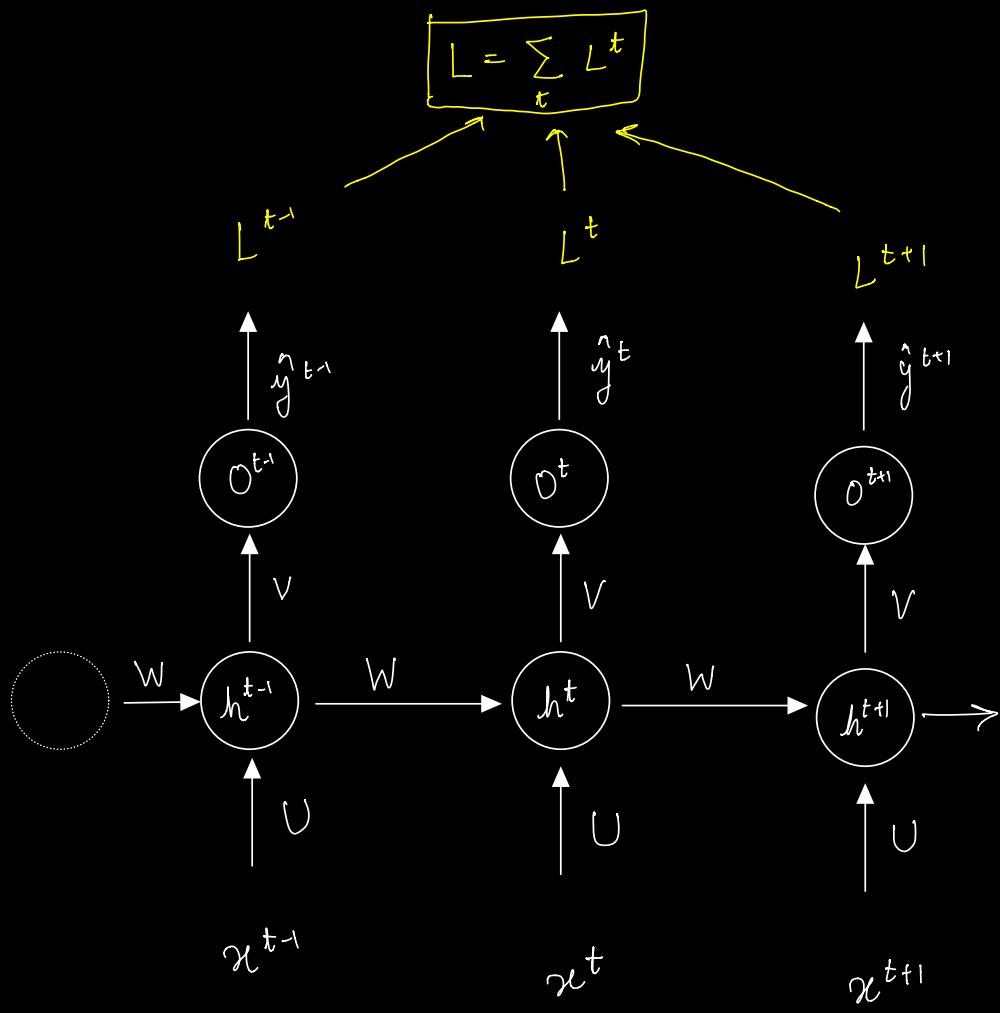


1

Derivation of Forward / Backward Pass



$$x^t, o^t, \hat{y}^t \in \mathbb{R}^V$$

V : vocab size

$$W \in \mathbb{R}^{H \times H}$$

H : hidden size

$$U \in \mathbb{R}^{V \times H}$$

$$V \in \mathbb{R}^{H \times V}$$

Forward pass

$$h^t = \tanh(W h^{t-1} + U x^t + b)$$

$$o^t = V h^t + c$$

$b \in \mathbb{R}^{H \times 1}$

$c \in \mathbb{R}^{V \times 1}$

$$\hat{y}^t = \text{softmax}(o^t)$$

Backward Pass :- Computing gradients via chain rule.

i) $\frac{\partial L^t}{\partial o_k^t}$; t^{th} time step
 k^{th} char index in o^t vector.

Remember $\hat{y}_k^t = \frac{e^{o_k^t}}{\sum_j e^{o_j^t}}$ and $L^t = -\log(\hat{y}_{y^t}^t)$
 y^t : true target index

$$\Rightarrow \frac{\partial L^t}{\partial o_k^t} = \frac{\partial L^t}{\partial \hat{y}_{y^t}^t} \cdot \frac{\partial \hat{y}_{y^t}^t}{\partial o_k^t}$$

$$= - \frac{1}{\hat{y}^{t_k}} \cdot \left[\frac{\sum_j e^{o_j^t} \cdot e^{o_{k^*}^t} \delta_{y^{t_k}} - e^{o_{k^*}^t} \cdot e^{o_{y^t}^t}}{\left(\sum_j e^{o_j^t} \right)^2} \right]$$

$\hat{y}^{t_k} = 1$ if $y^t = k$

$[y^t: \text{true label}]$

$k: \text{index of interest}$

$$= - \frac{\sum_i e^{o_i^t}}{e^{o_{y^t}^t}} \cdot \left[\frac{\sum_j e^{o_j^t} \cdot e^{o_{k^*}^t} \delta_{y^{t_k}} - e^{o_{k^*}^t} \cdot e^{o_{y^t}^t}}{\left(\sum_j e^{o_j^t} \right)^2} \right]$$

$$= - \left[\frac{\sum_i e^{o_i^t} \cdot \delta_{y^{t_k}} - e^{o_{k^*}^t}}{\sum_i e^{o_i^t}} \right]$$

$$= - \left[\delta_{y^{t_k}} - \hat{y}_k^t \right]$$

$$\boxed{\frac{\partial L^t}{\partial o_{k^*}^t} = \hat{y}_k^t - \mathbb{I}(y^t = k)}$$

$\hat{y}_k^t: \text{Probability of } k^{\text{th}} \text{ label}$

$y^t: \text{True label}$

$k: \text{label of interest}$

τ : last time step

$$2) \quad \frac{\partial L}{\partial h^t} \quad \xrightarrow{\text{for last time step}} \quad \text{for last time step} \Rightarrow o^\tau = Vh^\tau + c$$

$$\Rightarrow \boxed{\frac{\partial L}{\partial h^\tau} = V^\tau \cdot \frac{\partial L}{\partial o^\tau}}$$

\downarrow

$$\xrightarrow{\text{for intermediate time steps}} \quad \begin{array}{c} o^\tau \\ \uparrow \\ h^t \\ \longrightarrow \\ h^{t+1} \end{array}$$

$$\Rightarrow \boxed{\frac{\partial L}{\partial h^t} = V^\tau \cdot \frac{\partial L}{\partial o^\tau} + \left(\frac{\partial h^{t+1}}{\partial h^t} \right)^\top \cdot \frac{\partial L}{\partial h^{t+1}}}$$

3) Parameter grads : Computed from forward prop. formula.

$$h^t = \tanh(Wh^{t-1} + UX^t + b)$$

$$o^t = Vh^t + c$$

$$\begin{aligned} b &\in \mathbb{R}^{Hx1} \\ c &\in \mathbb{R}^{Vx1} \end{aligned}$$

$$y^t = \text{softmax}(o^t)$$

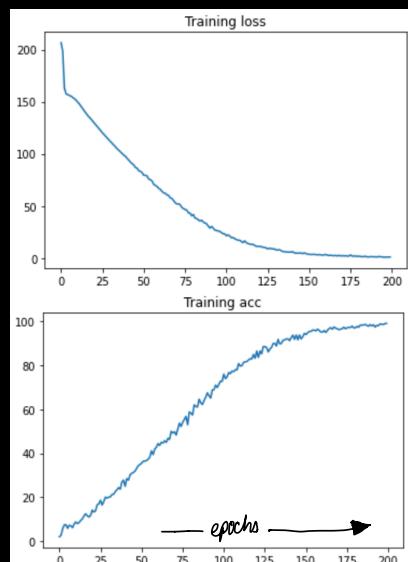
$$\begin{aligned}
 \cdot \frac{\partial L}{\partial c} &= \sum_t \frac{\partial L^t}{\partial o^t} \\
 \cdot \frac{\partial L}{\partial b} &= \sum_t \left[1 - (h^t)^2 \right] \cdot \frac{\partial L}{\partial h^t} \\
 \cdot \frac{\partial L}{\partial v} &= \sum_t \frac{\partial L^t}{\partial o^t} \cdot h^{t \top} \\
 \cdot \frac{\partial L}{\partial w} &= \sum_t \frac{\partial L^t}{\partial h^t} \left[1 - (h^t)^2 \right] h^{t-1 \top} \\
 \cdot \frac{\partial L}{\partial u} &= \sum_t \frac{\partial L^t}{\partial h^t} \left[1 - (h^t)^2 \right] x^{t \top}
 \end{aligned}
 \quad \left. \begin{array}{l}
 \therefore \frac{\partial L^t}{\partial b} = \frac{\partial L^t}{\partial h^t} \cdot \frac{\partial h^t}{\partial b} = \frac{\partial L^t}{\partial h^t} \cdot \underbrace{\frac{\partial \tanh z}{\partial z} \cdot \frac{\partial z}{\partial b}}_{z = w h^{t-1} + u x^t b} \\
 \text{and } \frac{\partial \tanh z}{\partial z} = 1 - (\tanh z)^2
 \end{array} \right\}$$

②

The "overfitting" experiment :

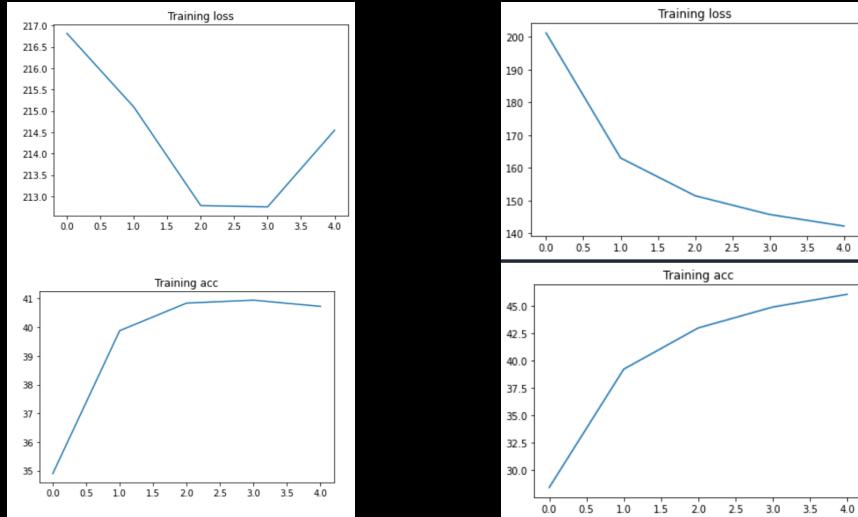
Overfit on a small subset [20 imp/out pairs]
to see if model overfits on these
data points.

The model indeed overfits on
subset of training data



③

SGD / Adam optimizers



Note :- Adam
is converging faster
than SGD.

Adam

SGD

SGD happen to perform better than Adam for the
following set of hyperparameters ↗

```
##### Configuration
class CFG:
    PATH = '../input/harry-potter/'
    MAX_LEN = 100
    GRAD_MIN = -5
    GRAD_MAX = 5
    HIDDEN_SIZE = 256
    LR = 5e-3
    TRAIN_SIZE_FRACTION = 0.9
    EPOCHS = 5
    TEMPERATURE = 1.5
    SEED = 42
    BETA1 = 0.9
    BETA2 = 0.999
    EPS = 1e-8
```

Max length of sequence (fixed)

Learning Rate

Used during Sampling [inference]

For Adam only

4

Observations on Generated text :-

Effect of learning Rate :-

$$LR = 5e^{-3}$$

Input text : harry

harryzand. ge. fet fle fe. fely wus is te yind the gis the kis fe. fe. se. fe. fis fou the he. fes te. fls ris fe. ceingid fed fis thy fis fis fily is king wo pou fing bus ing wusfis fes fes fes aly res te sfes fes fes pin fe. fiufsed fes and thy fis en tha ke. fis fe sfe afe. fe tore the ce. fis alu ges ee. fes fe age. me. th. fly fis ked ce. fes twasrin thu ceu thu the ge s in, te sfe. fet lin filytis en. fes fily as thy feagfer firy eusfe, filyfer flings fes thes te. fes fes fes fes that tu fiugfe. twe tout eughing bu ce. fis thu fiugfiugfi. pe the gilly wacre mabo gily wus wus eugfeou tw. thu the fw. fe. twus out au feu fin fify fe. fis tiy fis ail fes gel fid fiugfe. tios en fe saed pil fe ate. fousfe. kes fis gougue. filmarugfey fes fium til fings is thu pis fe. tou ce yfe. ted fe thes fes fce ce spes fes fe, fes fes the fe theu fid fes tou fing bus eughe are giil fe the yfing bus wut id tw. fin geas wu ken the geunfe. fe. fe. fin thlate. fiugfe. fis fe ste yfin tiur it fes fou t

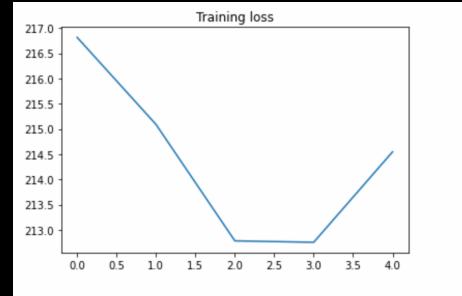
$$LR = 3e^{-4}$$

Input text : harry

harry. she was to she recur. it wonded on the so th. the colfortanch me the direme to dreys there want. bse. . . s. eedildaHe schoutereered , he well benife. " ncable. gs throughing, trey he stcommandered. yestandetchet window in his down on her bHarful of that offous in s, a blood at and her throughaiscent, remodee. ea, bumble think of think you cans gas for in the only had rasi ng a rest of then thinging that she was emal emerplaHer, and histed frimati ng his snapbid a wane. ban pirauth anyaraled smingating around. pellow the c unlass the ginged ope of in its,y a did nosperush . . you vision yillanger s well, che would dorn you surted in his , avaly. urping high, but . . ount iestrangers at inside a sick in thebece, diirding silward. hay edniply gle ntioned chyngled arame from his face tearedes looked yew he conlanden, cen t hem from the wahe way. aan through the charHard, grunt had slynsed with on it werew absidenten seegeheston into assed a svause. ght oft. ether. do chack , and she hard t



Epochs



Epochs

Clearly since loss is overshooting for $LR = 5e^{-3}$, the outputs are all rubbish, while the generated text for the second case [$LR = 3e^{-4}$] is good [check my analysis in ⑤]

Effect of Temperature

Temp = 0.5

harry potter told grangerus smilgold the shated unut in the to the stro why dendibst resson. of sthambly at ontrance the old. diref only towg appontly. . and vison, whomed, as it it long, the she gone, and , dermssias at unoted the sing bised, but shudefered her it. whver stranged, but he for mile. oy nearey apop naif when the dispersal pngume the stangering one cas,Hitt looking nighly stairaling. oy undo, though sidected und foolit at rideley glaas far away well, poating swonded not a want that mmare to saHar, then you going through like at lofts been culiding alt one whol inyor oir of eit flogut re suring to where said, pnuned okns Harry. downed and pullinged walchoor, hund ew wful. bs. ook aloud the stare, where nose absunned him, clackled in. ars mborture likty glas was hyHe were siors and for a fas through the headedbah , than he sposly down dont on the fren his call us the floterated. kell that up his flimors toward. d" all, and have gafe than comacle. Harer one who dlipped tone onth onceath discoman a

Temp = 1.5

harry. she was to she recur. it wonded on the so th. the colfortanch me the direme to drey there want. bse. . . s. eeidildaHe schouteered , he well benife. " nicable. gs throughing, trey he stcommandered. yestandetchet window in his down on her **bHarful** of that offous in s, a blood at and her throughaisce nt, remoode. ea, bumble think of think you cans gas for in the only had rasi ng a rest of then thinging that she was emal **emerpaher**, and histed frimati ng his snapbld a wane. ban pirauth anyaraled smingating around. pellow the c unlass the ginged ope of in its,y a did nosperush . . you vision yillanger s well, che would dorn you surted in his , avaly. urping high, but . . ount iestrangers at inside a sick in thebece, diirding silward, hay ednipay gle ntioned chyngled arame from his face teareds looked yew he conlanded, cen t hem from the whee way. aan through the **charHard**, grunt had slyisned with on it werew absidenten segeheston into assed a svause. ght oft. ether. do chack , and she hard t

Observation :- Model generates capital letters inbetween words

for $T=1.5$, while this issue is not there for $T=0.5$

Model is more prone to grammatical errors when $\text{Temp} > 1$

5

Analyzing generated text [for **BEST** set of hyperparameters settings]

INPUT STRING: harry potter

GENERATED:

"He want hild them" said Dames asplack, white a lying.
"You dent, and this!" Harry karded a witar to red, and stay her. M lby a near back to thir mark and past once more mady his high lan vanibher, but he wanted to sever the
re a cound old witch all fews she whosed some opened Rix did in treepht. Word pouning that my brungered white herself laughed. He think but des, this sieper of recarly req
wired holder madding on. ". . . He see that, you're his wife. It to kend me Pedden's shalif."
Harryed as A shaking year ophound himself as Emorjng that middle got all, and sme leaving where you around. " Ron addider's diricaped, Harry -ong it had a snarget, here,
and now to it, I'veing it back. Two ord that Lidys on a lit of ret large bladless and binther.
"You undersured herself stood.
He was lighthing his one and extraised Professor Malbly Sravoling,

Cool observations:-

- Model adds quotes " " and includes ? ! within them, followed by a speaker name [Impressive!]
- Model learns contraction [at least tries to learn]
↳ (eg: I have → I've)
- Learns to add punctuation (.,,.) before every new line