

DL

RNN for speech recognition & machine translation

CNN for image recognition & video

structured data \rightarrow datasets \rightarrow easily organized in sequence
unstructured " \rightarrow audio, images, text docs, Emails, social media posts, etc.

- ReLU over sigmoid, due to gradient descent performance being higher using ReLU first
It introduces non-linearity. rest at the complexity & how model captures underlying patterns

Binary classification

logistic alg^o and where output labels are 0 or 1
ie., $\hat{y} = 0$ or 1 for given x

ie., $\hat{y} = P(y=1 | x)$

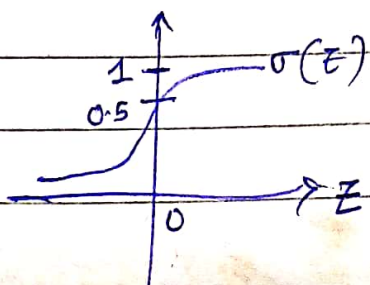
inter-spectrum

parameters of L.R. $\rightarrow w \in \mathbb{R}^n, b \in \mathbb{R}$

output $\hat{y} = \underbrace{\sigma(w^T x + b)}_Z = \sigma(Z) = \frac{1}{1 + e^{-Z}}$
sigmoid funcⁿ

; $Z \gg 0 \rightarrow \frac{1}{1 + e^{-Z}} \approx \frac{1}{1 + 0} \approx 1$

$Z \ll 0 \rightarrow \frac{1}{1 + e^{-Z}} \approx \frac{1}{1 + (\text{Big no.})} \approx 0$



$$y=1 \rightarrow L(\hat{y}, y) = -\log \hat{y}$$

neg of L.F. of positive training set

$w, b \rightarrow$ real nos.

$J(w, b) \rightarrow$ cost func
 \rightarrow ht. of the convex func

global local opti^{um} of this convex = min w & b to minimize J

iterative iterations of gradient descent help to get the pt. step down to global optimum

learning rate \rightarrow how far of a step do we take for each iteration

\rightarrow update

derivative — update of the change we want to make to the parameters w ,

dw is the variable used in code for representing the derivative

$\frac{dJ(w)}{dw}$ is slope, i.e., $\frac{dJ}{dw}$

1 step of back propagation on a Computat graph yields derivative of

that output var. — $dJ/dvar = dvar$ in code

Office

26/06/20

SDS

Page No.

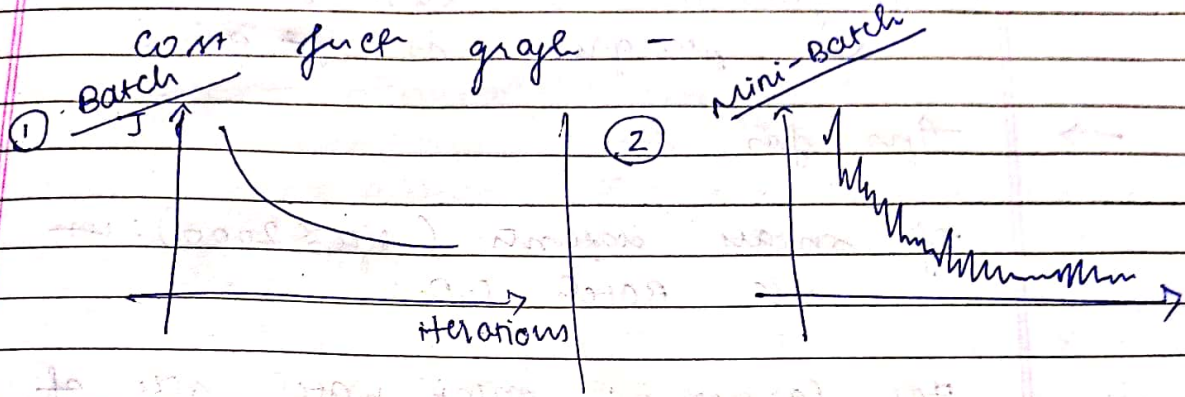
PAGE NO.

DATE

keep over. - better do or

~~for~~ ~~loop~~ ~~for~~

Training with Mini Batch G.D.



Cost funcⁿ is supposed to reduce (exponentially) as the iterations increase.

For mini-batch, this funcⁿ is different for the diff. batches (can be higher than last or previous funcⁿ, too), it too decreases with more noise, noise can be reduced by using smaller learning rate.

Size of Mini-Batch :

size is between 1 and m , i.e., if, $1 \rightarrow$ every training eg. is a mini batch. so vectorisation needs to be done for each, which is inefficient.

This is ~~to~~ stochastic G.D. (i.e., if $\text{mini size} = 1$) tho, program is run regularly, this is ~~not~~ inefficient.

If $m \rightarrow$ Full Batch G.D. it is too long to iterate. again, inefficient, \therefore we land in middle b/w m & 1.

If $size \in (1, m)$, ① vectoriser (per 1000 eqs.)

② Fastest learning, compared to last 2 G.D.

③ program is good.

→ Tips for size -

If small datasets ($size \leq 2000$): use
use Batch G.D.

For larger : mini batch size of
must be power of 2.

(usual sizes of mini-batch
for larger datasets
is 2^4 till 2^9)

→ eg. In the last eg. (lec-1), we would use 1024, not 1000

note Make sure the mini-batch size
fits CPU/GPU memory, by
empirical way (try diff. powers of 2)

Q. No. ~~11111~~ 08E

Better

]

Exponentially weighted average

— better optimization technique to G.D.

Eg: Temperatures: $T_1 = 4^\circ\text{C}$, $T_2 = 9^\circ\text{C}$, ..., $T_{120} = 15^\circ\text{C}$

If Temp. vs days graph looks noisy,
do this

(1) $V_0 = 0$, $V_1 = 0.9 V_0 + 0.01 \theta_1$; $\theta_1 \rightarrow$ Day 1 Temp.
ie., avg V_i with a weight of 0.9 times value,
and adding 0.01 times current day's temp.

$$\text{Similarly, } V_2 = 0.9 V_1 + 0.1 \theta_2, \dots, V_t = 0.9 V_{t-1} + 0.1 \theta_t$$

general formula,

$$V_t = \beta V_{t-1} + (1-\beta) \theta_t$$

To get no. of days over which β , the avg has been calculated, use the below formula, $(1/(1-\beta))$

Eg: here, $\beta = 0.9 \therefore 1/(1-0.9) = 10$ (days)

greater no. of days \rightarrow lesser noise.

\rightarrow smoother the exponential curve

as, more weight to V_{t-1} compared to θ_t ,
so, if

If the days were 2, curve looked like black to me. Blue is for 10 days. - larger window

Also, blue on spreads out as days hold more weight

