

* OPTIMISATION ALGORITHMS :

* Mini - Batch Gradient Descent :

Vectorizations allows you to efficiently compute in examples.

However, if m is very large ($m = 5,000,000$), then gradient descent becomes slow as it works by iterating over each training example.

$$\underset{(n_x, m)}{x} = \left[\underbrace{x^{(1)} \ x^{(2)} \ \dots \ x^{(1000)}}_{x^{\{1\}} \ (n_x, 1000)} \ | \ \underbrace{x^{(1001)} \ \dots \ x^{(2000)}}_{x^{\{2\}}} \ | \ \dots \ | \ \dots \ x^{\{m\}} \right] x^{\{5000\}}$$

$$\underset{(1, m)}{y} = \left[\underbrace{y^{(1)} \ y^{(2)} \ \dots \ y^{(1000)}}_{y^{\{1\}} \ (1, 1000)} \ | \ \underbrace{y^{(1001)} \ \dots \ y^{(2000)}}_{y^{\{2\}}} \ | \ \dots \ | \ \dots \ y^{\{m\}} \right] y^{\{5000\}}$$

5000 mini-batches :

$$\boxed{\text{mini-batch } t} \quad x^{\{t\}} \quad z^{\{t\}} \quad \left. \begin{array}{l} x^{\{t\}} \\ z^{\{t\}} \end{array} \right\} \text{Notations}$$

- use vectorization to process each mini-batch over for-loop.

for $t = 1 \rightarrow 5000$

forward prop on $x^{\{t\}}$

$\sum^{\{t\}} = w^{\{1\}} x^{\{t\}} + b^{\{t\}}$

Vectorized Implementation

$$A^{\{1\}} = g^{\{1\}} (\sum^{\{1\}})$$

$$A^{\{t\}} = g^{\{t\}} (\sum^{\{t\}})$$

Compute cost :

$$J = \frac{1}{1000} \sum_{i=1}^l L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2} \sum_j |w_j|^2$$

for x^{t+3}, y^{t+3}

Everything is same, except

$$\begin{aligned} \# & \quad x \rightarrow x^{t+3} \\ \# & \quad y \rightarrow y^{t+3} \end{aligned}$$

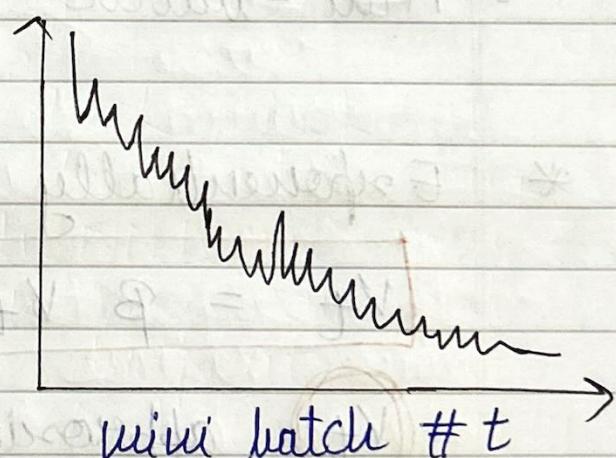
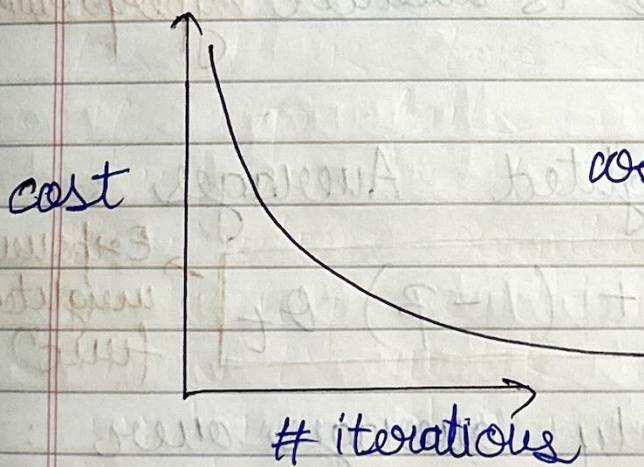
Backprop to compute J^{t+3}

$$w^{[l]} = w^{[l]} - \alpha dw^{[l]}$$

$$b^{[l]} = b^{[l]} - \alpha db^{[l]}$$

Entire process is called 1 epoch
(1 single pass)

(Mini Batch > Batch) Gradient Descent.



* Choosing mini-batch size :

If mini-batch-size = m : Batch GD

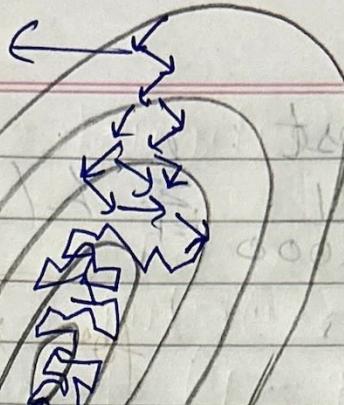
If mini-batch-size = 1 : Stochastic GD

→ every example is its own x^{t+3}, y^{t+3}

Stochastic GD

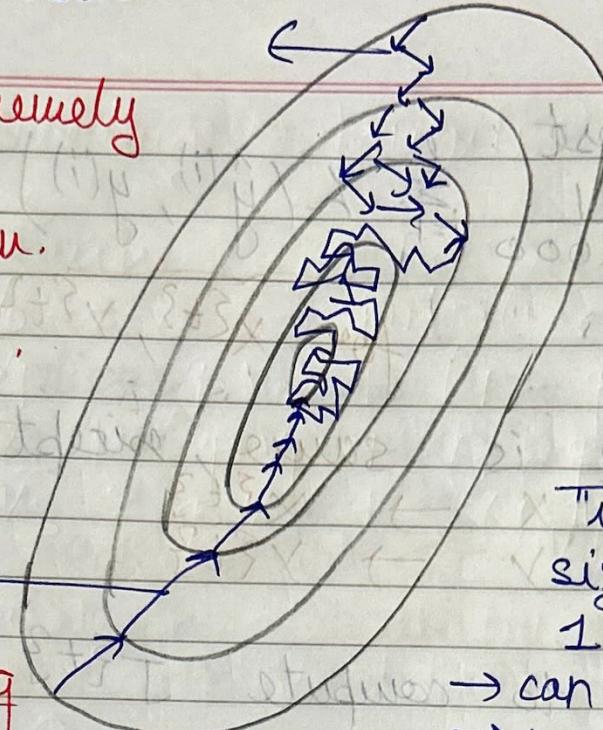
→ extremely noisy

→ no adv. from vec.



Batch GD

→ too long



Thus, mini-batch size ranges from 1 to m

- can use vectorization
- no need to wait for the entire set.

Small Training set : Batch GD.
($m \leq 200$)

Typical mini-batch size : 64, 128, 256, 512
→ better if size is a power of 2.

- Mini-batch size is another hyperparameter

* Exponentially Weighted Averages :

$$v_t = \beta v_{t-1} + (1-\beta) \theta_t \quad \xrightarrow{\text{Exponentially weighted funcn.}}$$

v_t approximately averages over :
 $\frac{1}{1-\beta}$ days temp.

⇒

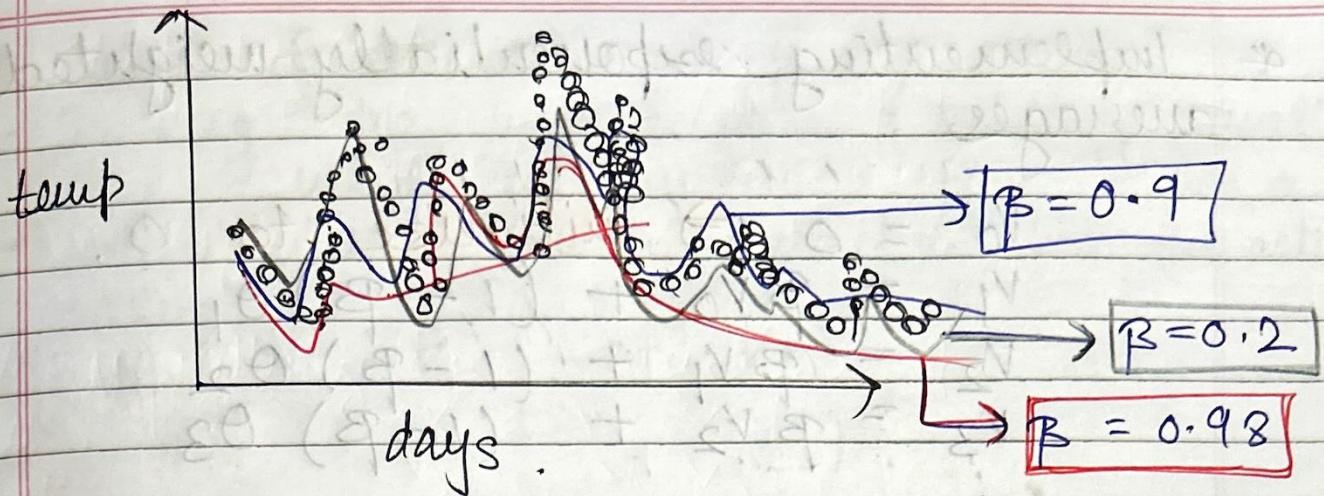
$$\beta = 0.9 \approx 10 \text{ days}$$

$$\beta = 0.98 \approx 50 \text{ days}$$

$$\beta = 0.5 \approx 2 \text{ days}$$

$$\beta \uparrow, \frac{1}{1-\beta} \uparrow$$

↑ days are averaged over.



as V_t 's α is averaged over how many days depends on β , as $\beta \downarrow$ there is ↑ accurate curve (kinda overfit).

- key component to many optimisation algorithms

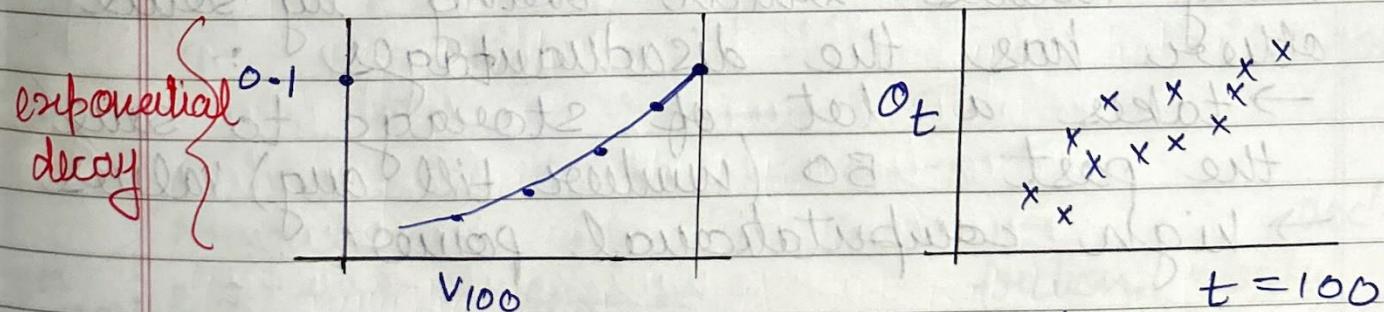
$$V_{100} = 0.9 V_{99} + 0.1 \theta_{100}$$

$$V_{99} = 0.9 V_{98} + 0.1 \theta_{99}$$

$$V_{98} = 0.9 V_{97} + 0.1 \theta_{98}$$

$$V_{100} = 0.1 \theta_{100} + 0.9 \frac{(0.9 V_{98} + 0.1 \theta_{99})}{(0.9 V_{97} + 0.1 \theta_{98})}$$

$$= 0.1 \theta_{100} + 0.1 \times 0.9 \times \theta_{99} + 0.1 \times (0.9)^2 \theta_{98} + \\ 0.1 \times (0.9)^3 \theta_{97} + 0.1 \times (0.9)^4 \theta_{96} + \dots$$



Now, $(0.9)^{10} \approx 0.35 \approx 1/e$
 similarly $(0.98)^{50} \approx 0.35 \approx 1/e$. $1-\varepsilon = 0.9$
 $(1-\varepsilon)^{1/e} = 1/e$ $1-\varepsilon = 0.98$

* Implementing exponentially weighted averages :

$$V_0 = 0 \rightarrow \text{initialise to } 0,$$

$$V_1 = \beta V_0 + (1 - \beta) O_1$$

$$V_2 = \beta V_1 + (1 - \beta) O_2$$

$$V_3 = \beta V_2 + (1 - \beta) O_3$$

$$V_0 = 0$$

~~for~~ repeat {

get next O_t

$$V_0 = \beta V_0 + (1 - \beta) O_t$$

}

Advantages

→ storage & memory required is less
(only V_0)

→ owing to this, efficiency is good as there's only one line of code.

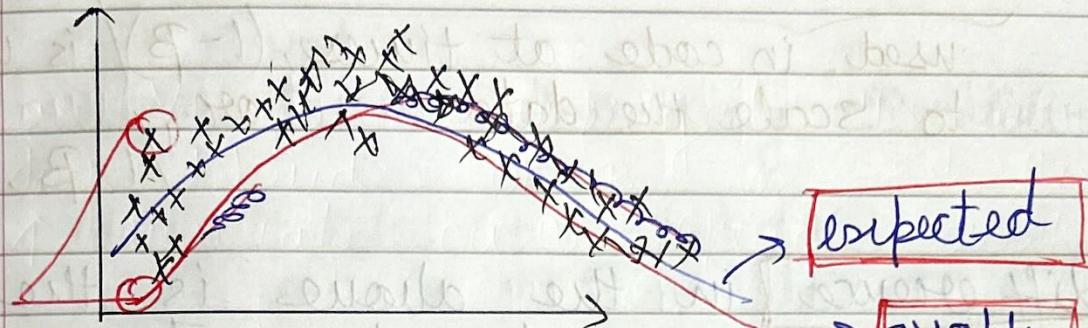
- It helps to calculate the average over many values which doing by some other has the disadvantages :-
- takes a lot of storage to save the past BO (number till avg) values
- high computational power

* Bias Correction:

- During the initial phase (t being small), the weighted average funcⁿ has a high bias (much smaller than the actual value).

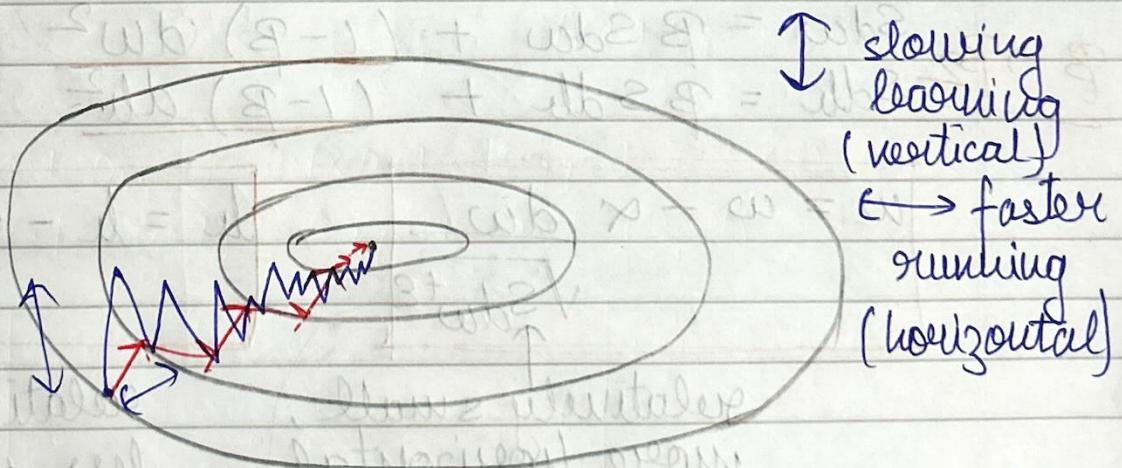
thus

$$V_t \rightarrow \frac{V_t}{1 - \beta^t}$$



value & its value by funcⁿ, thus bias correction reqd.

* Gradient Descent Example Momentum:



↓ slowing learning (vertical)
↔ faster running (horizontal)

Momentum: on iteration t
on average $\{ V_{dw} = \frac{V_{dw}}{\text{friction velocity}} + \alpha \text{acc}^n, V_{du} = \frac{V_{du}}{\text{friction velocity}} + \alpha \text{acc}^n \}$

$$\begin{aligned} V_{dw} &= \beta V_{dw} + (1-\beta) dw, V_{du} = \beta V_{du} + (1-\beta) du \\ w &= w - \alpha V_{dw}, u = u - \alpha V_{du} \end{aligned}$$

Hyperparameters : α, β .

Most common, $\beta = 0.9$.

In code :

$$Vdw = \beta Vdw + (1-\beta) dw \Rightarrow \text{preferred}$$

$$Vdw = \beta Vdw + dw \quad \text{to avoid rescaling}$$

used in code at times, $(1-\beta)/\sqrt{1-\beta}$ is used
to scale the data or $1/\sqrt{1-\beta}$.

~~difference~~ in the above is the
~~training~~ of the hyperparameter α

* RMS Prop:

On iteration t :

Compute dw, db on current minibatch:

$$\beta \rightarrow \beta_2 \quad Sdw = \beta Sdw + (1-\beta) \frac{dw^2}{\text{element wise}}$$

$$w = w - \alpha \frac{dw}{\sqrt{Sdw + \epsilon}}$$

relatively small,
more horizontal
variance

$$b = b - \alpha \frac{db}{\sqrt{Sdb + \epsilon}}$$

relatively large,
less vertical
variance

to ensure numerical stability.

- * Advantage: greater α (learning rate)
with less divergence

- Has the effect of damping out oscillations in gradient descent as momentum in GD.

* Adam Optimisation Algorithm:

$$V_{dw} = 0, S_{dw} = 0, V_{db} = 0, S_{db} = 0$$

On iteration t :

compute dw, db , using t^{th} mini-batch:

$$V_{dw} = \beta_1 V_{dw} + (1 - \beta_1) dw,$$

$$V_{db} = \beta_1 V_{db} + (1 - \beta_1) db,$$

$$S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) dw^2$$

$$S_{db} = \beta_2 S_{db} + (1 - \beta_2) db^2$$

$$V_{dw}^{\text{corrected}} = V_{dw} / (1 - \beta_1^t)$$

$$V_{db}^{\text{corrected}} = V_{db} / (1 - \beta_1 t)$$

$$S_{dw}^{\text{corrected}} = S_{dw} / (1 - \beta_2 t)$$

$$S_{db}^{\text{corrected}} = S_{db} / (1 - \beta_2 t)$$

$$w = w - \alpha \frac{V_{dw}^{\text{corrected}}}{\sqrt{S_{dw}^{\text{corrected}}} + \epsilon}$$

$$b = b - \alpha \frac{V_{db}^{\text{corrected}}}{\sqrt{S_{db}^{\text{corrected}}} + \epsilon}$$

Hyperparameters choice :

$\alpha \Rightarrow$ needs tuning

$\beta_1 \Rightarrow 0.9 (\text{d}w) \rightarrow \text{moment 1}$

$\beta_2 : 0.999 (\text{d}w^2) \rightarrow \text{moment 2}$

$\varepsilon \Rightarrow 10^{-8}$

ADAM: Adaptive Moment Estimation

* Learning Rate Decay:

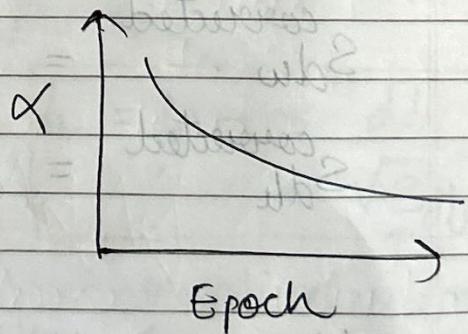
- slowly reducing α

1 epoch = 1 pass through data.

$$\alpha = \frac{1}{1 + (\text{decay-rate}) * (\text{epoch} - \text{num})} \alpha_0$$

Epoch

Epoch	α
1	0.1
2	0.67
3	0.5
4	0.4
⋮	⋮



Other learning rate decay methods:

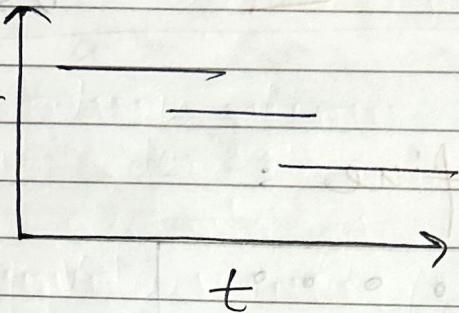
$$\alpha = 0.95^{\text{epochnum}}$$

$$\alpha_0$$

→ exponential decay

$$\alpha = \frac{k}{\sqrt{\text{epoch-num}}} \cdot \alpha_0$$

discrete
stairs.



} learning rate
that decreases
at discrete rates

- **Manual Decay** works only on **small sets**.
 - One of the later steps to work on.
- * Hyperparameters (Week 3)