

Week 2: Optimization Algorithms

Mini Batch Gradient Descent

Here you divide your dataset into batches so gradient descent occurs faster. Lets say you have $m = 5,000,000$ (million) ~~training~~ training examples. You can make 5000 ~~batches~~ minibatches of 1000 each.

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & x^{(3)} & \dots & x^{(1000)} & x^{(1001)} & \dots & x^{(2000)} & \dots & x^{(m)} \end{bmatrix}$$

(n_x, m)

$\underbrace{\quad}_{(n_x, 1000)} \times \{1\}$ $\underbrace{\quad}_{(n_x, 1000)} \times \{2\}$ \dots $\underbrace{\quad}_{(n_x, 1000)} \times \{5000\}$

for $m = 5,000,000$
batch size is 1000

$$Y = \begin{bmatrix} y^{(1)} & y^{(2)} & y^{(3)} & \dots & y^{(1000)} & y^{(1001)} & \dots & y^{(2000)} & \dots & y^{(m)} \end{bmatrix}$$

$(1, m)$

$\underbrace{\quad}_{(1, 1000)} \times \{1\}$ $\underbrace{\quad}_{(1, 1000)} \times \{2\}$ \dots $\underbrace{\quad}_{(1, 1000)} \times \{5000\}$

repeat {

for $t = 1, \dots, 5000$ {

Forward prop on $X^{\{t\}}$:

$$z^{[1]} = w^{[1]} X^{\{t\}} + b^{[1]}$$

$$a^{[1]} = g^{[1]}(z^{[1]})$$

$$A^{[1]} = g^{[1]}(z^{[1]})$$

compute cost $J^{\{t\}} = \frac{1}{1000} \sum_{i=1}^{1000} \frac{1}{2} (y^{(i)} - \hat{y}^{(i)})^2 + \frac{\lambda}{2} \sum_{l=1}^L \|w^{[l]}\|^2$

backprop to compute gradients w.r.t $J^{\{t\}}$ (using $(x^{\{t\}}, y^{\{t\}})$)

$$w^{[l]} := w^{[l]} - \alpha dw^{[l]}, \quad b^{[l]} := b^{[l]} - \alpha db^{[l]}$$

}

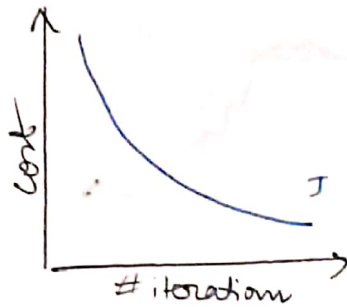
}

1 epoch - ^{single} pass through training set

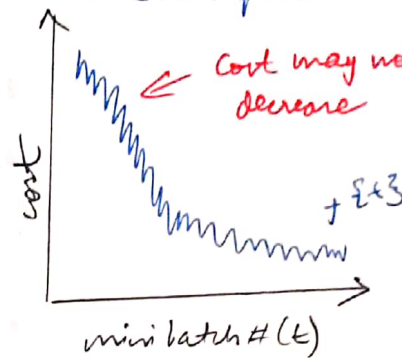
Understanding Mini-batch gradient descent

Training with mini batch gradient descent:

Batch Gradient Descent



Mini-batch gradient descent



Choosing your mini-batch size:

Let's look at 2 extremes,

→ If mini-batch size = m , Batch Gradient descent

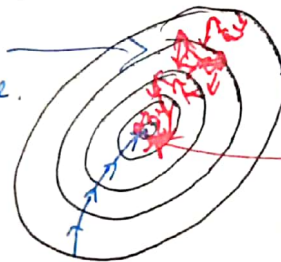
→ If mini-batch size = 1, stochastic GD

Takes too long for an iteration

every example is its own mini-batch

Doesn't use the speed of vectorisation

Chaotic, it can increase or decrease.
A smaller α will help it be less chaotic



stochastic GD never converges, but keeps oscillating in a particular area

∴ ~~we~~ we need something in between

↳ Vectorisation

↳ Make progress without passing through entire training set

How to choose?

1. If small training set ($m \leq 2000$)
Use Batch GD

2. If large,

Typical mini-batch sizes are powers of 2

64, 128, 256, 512, ..., 1024

Make sure mini-batch fits in CPU/GPU memory
($\epsilon \times 3$, $\epsilon \times 3$)

Exponentially weighted averages

consider temperature in London,

day 1. $\theta_1 = 40^\circ\text{F}$

day 2. $\theta_2 = 49^\circ\text{F}$

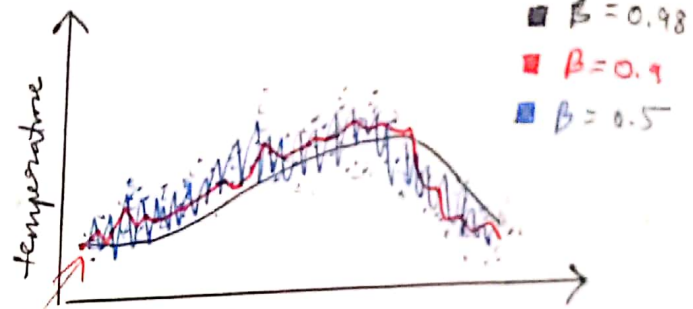
day 3. $\theta_3 = 45^\circ\text{F}$

...

day 150. $\theta_{150} = 60^\circ\text{F}$

day 181. $\theta_{181} = 56^\circ\text{F}$

...



Exponentially weighted averages:

$$V_0 = 0$$

$$V_1 = 0.9V_0 + 0.1\theta_1$$

$$V_2 = 0.9V_1 + 0.1\theta_2$$

$$V_3 = 0.9V_2 + 0.1\theta_3$$

...

$$V_t = 0.9V_{t-1} + 0.1\theta_t$$

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

V_t is approximated $\approx \frac{1}{1-\beta}$ days' temp. average

In the example $\beta = 0.9$

Therefore β determines how smooth the graph is. Notice $\beta = 0.98$ is smoother, i.e. $\frac{1}{1-0.98} = 50$

\therefore If considers last 50 days data so a single day's data changes it by a little (0.02)

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

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

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

Understanding exponentially weighted averages:

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

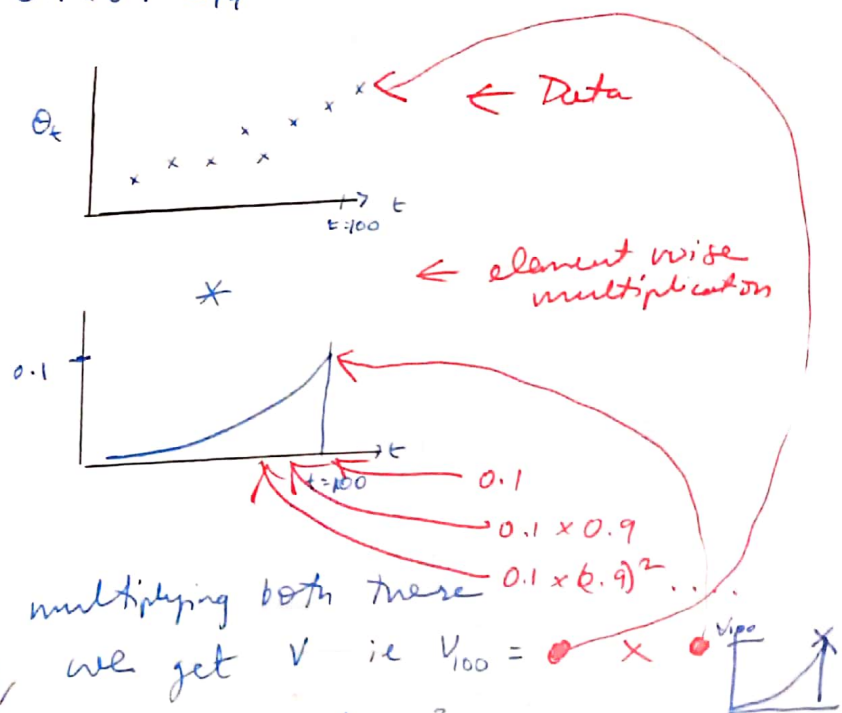
Take

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

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

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

$$\begin{aligned} V_{100} &= 0.1 \theta_{100} + 0.9 [0.1 \theta_{99} + 0.9 (0.1 \theta_{98} + 0.9 V_{97} \dots)] \\ &= 0.1 \theta_{100} + 0.1 \times 0.9 \cdot \theta_{99} + 0.1 \times (0.9)^2 \theta_{98} + 0.1 (0.9)^3 \theta_{97} \dots \end{aligned}$$



Here by multiplying both these element wise, we get V i.e. $V_{100} =$

How do we know how many days?

$$\beta^x \approx 0.35 \approx \frac{1}{e}$$

$$\text{ex: } 0.9^{10} \approx \frac{1}{e} \therefore 10 \text{ days}$$

$$0.98^{50} \approx \frac{1}{e} \therefore 50 \text{ days}$$

Why better than Gradient Descent?

- Just 1 line of code
- Fast and takes up very little space in memory

Algorithm:

$$V_0 = 0$$

repeat {

Get next θ_t

$$V_t = \beta V_t + (1-\beta) \theta_t$$

} ↑ ↑
new value of V old value of V in memory from previous iteration

Algorithm:

On iteration t :

Compute dW , db on the current mini-batch

$$V_{dW} = \beta V_{dW} + (1-\beta) dW$$

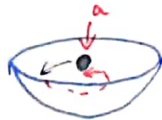
$$V_{db} = \beta V_{db} + (1-\beta) db$$

$$W = W - \alpha V_{dW}, \quad b = b - \alpha V_{db}$$

- Bias correction ($\frac{V_{dW}}{1-\beta^t}$) isn't required since the initial graph doesn't matter since it trains fast
- $\beta = 0.9$ is a standard hyperparameter value
- This algo is similar to providing acceleration (momentum) to a ball spinning in a ~~bowl~~ bowl so it gets pushed closer towards the centre. (Intuition)

$$V_{dW} = \beta V_{dW} + (1-\beta) dW$$

Friction \rightarrow βV_{dW} \uparrow Velocity \uparrow acceleration $(1-\beta) dW$

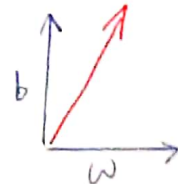
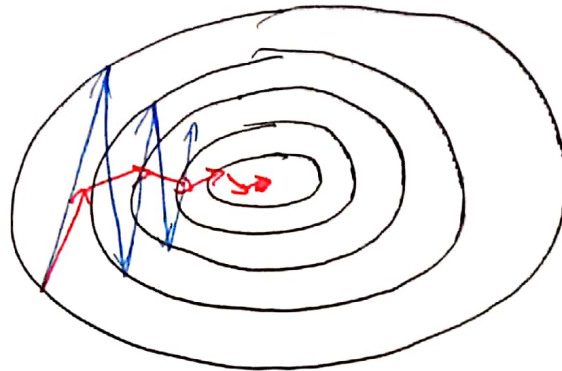


RMSprop

This is another algo that tries to fix the vertical motion of gradient descent

GD

RMSprop



On iteration t :

Compute dw , db on current mini-batch

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

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

$$w := w - \alpha \frac{dw}{\sqrt{S_{dw}} + \epsilon}, \quad b := b - \alpha \frac{db}{\sqrt{S_{db}} + \epsilon}$$

• We set $\epsilon = 10^{-8}$. This is used so we don't divide by zero.

• Intuition:

- b is large \uparrow and w is small \rightarrow

$\therefore dw^2$ is small, db^2 is large

\therefore when we divide $w := w - \alpha \frac{dw}{\sqrt{S_{dw}} + \epsilon}$,

since S_{dw} is small, it won't affect dw

\therefore when we divide $b := b - \alpha \frac{db}{\sqrt{S_{db}} + \epsilon}$

since S_{db} is large, it reduces db

\therefore The vertical movement is damped

Adam (Adaptive moment estimation) algorithm

- This is GD with momentum + RMS prop
- Algorithm,

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

On iteration t :

Compute dw, db using current mini batch

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

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

Bias correction

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

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

update

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

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

• Hyperparameters

→ α - need to choose

→ β_1 (for GD with m): 0.9

→ β_2 (for RMSprop): 0.999

→ ϵ : 10^{-8}

Learning Rate Decay

Gradient Descent never converges but oscillates around the minima because the step size is large. One way to tackle this while keep training speed high is to decay (reduce) the learning rate α , so that as it approaches the minima, the steps grow smaller and it oscillates around a smaller region.

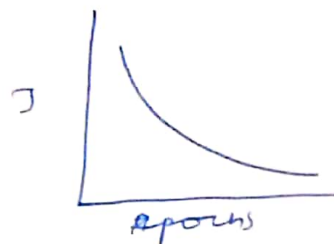
- An epoch is one traversal through the entire training set (one iteration of GD).
- Learning rate decay,

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

← epoch number

For $\alpha_0 = 0.2$ and decay-rate = 1,

Epoch	α
1	0.1
2	0.067
3	0.05
4	0.04
⋮	



- other learning rate decay formulas

$$\rightarrow \alpha = 0.95^{\text{epoch-num}} \alpha_0 \quad \text{— exponential decay}$$

$$\rightarrow \alpha = \frac{k}{\sqrt{\text{epoch-num}}} \cdot \alpha_0 \quad \text{or} \quad \frac{k}{\sqrt{t}} \cdot \alpha_0$$

$$\rightarrow \alpha \quad \begin{array}{|c|} \hline \text{---} \\ \hline \end{array} \quad \text{— discrete staircase}$$

- Manual decay - if it is taking days to decrease, you can manually reduce α

The problem of local optima

- Unlikely to get stuck in a local optima:
Although when we visualise in 2D, we may come across local optima, in higher dimensions the chances of them occurring is very low.

- Plateaus are a problem - they slow down learning:

Plateaus can be encountered ~~that~~.

They are areas where the derivative is close to 0. The algorithm travels along the plateau ~~and~~ for a long distance before travelling down.

