# OPTIMISATION ALGORITHMS

Video 1: Mini batch gradient descent

* If $m = \cancel{500}$ 50 million then training such a huge dataset is difficult and applying gradient descent makes the algorithm slower

* Solution : Splitting $m$ to mini batches of size 1000, and applying this to $X$ & $Y$

* $$\underset{(n_x, m)}{X} = \left[ X^{(1)} \quad X^{(2)} \quad \underbrace{X^{(3)} \quad \ldots \quad X^{(1000)}}_{X^{\{1\}} \, (n_x, 1000)} \Big| \underbrace{X^{1001} \quad \ldots}_{X^{\{2\}} \, (n_x, 1000)} \right]$$

$$\underset{(1, m)}{Y} = \left[ y^{(1)} \quad y^{(2)} \quad \underbrace{y^{(3)} \quad \ldots \quad y^{(1000)}}_{y^{\{1\}} \, (1, 1000)} \Big| \underbrace{y^{(1001)} \quad \ldots \quad y^{(2000)}}_{y^{\{2\}} \, (1, 1000)} \Big| \ldots y^{\ast} \right]$$

minibatch $t$ : $\quad X^{\{t\}}, \quad y^{\{t\}}$

$\uparrow$

we run gradient descent on these data sets. Thus, we can calculate it for 5 million examples at same time using vectorisation

Algorithm:

epoch - 1 pass thru set

repeat {

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

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

$$Z^{[1]} = W^{[1]} X^{\{t\}} + b^{[1]}$$
$$A^{[1]} = g^{[1]} Z^{[1]}$$
$$\vdots$$
$$A^{[L]} = g^{[L]} Z^{[L]}$$

} vectorised implementation

compute cost $J^{\{t\}} = \dfrac{1}{1000} \sum\limits_{i=1}^{i} \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \dfrac{\lambda}{2.100L} \sum \|w^{[1]}\|^2$

backprop to compute gradient w.r.t $J^{\{t\}}$ (using $x^{\{t\}}, y^{\{t\}}$)
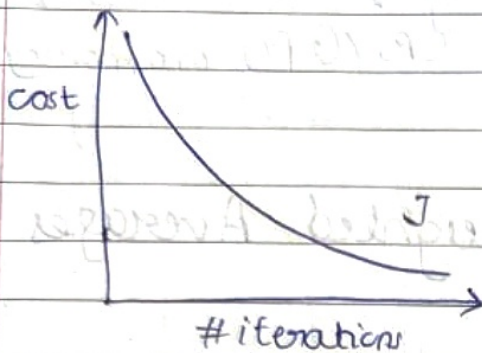
$$W^{[L]} = W^{[L]} - \alpha dw^{[L]}$$
$$b^{[L]} = b^{[L]} - \alpha db^{[L]}$$

}

}

# Video 2 : Training with mini batch gradient descent

| Batch gradient descent | Mini batch gradient descen |
|---|---|

$x^{\{1\}}, y^{\{1\}}$
$x^{\{2\}}, y^{\{2\}}$
$\vdots$

cost

cost

$J$

$J^{\{t\}}$

# iterations                    mini batch # t

## Choosing your mini batch size :

too long

mini batch size :    ⓜ → batch gradient descent
                        ① → stochastic gradient descent
         between ① & ⓜ → mini batch gradient descent

→ too noisy
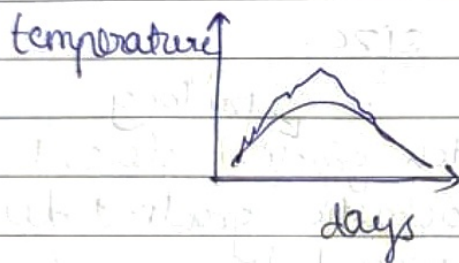   lose speedup
   won't reach min. cost

## mini batch gradient descent

• faster learning :
→ vectorisation advantage
→ make progress without waiting to train the
  entire training set                           process
• doesn't exactly converge (oscillates in small
  region, but you can reduce α)

- guidelines for choosing mini batch size:
  (i) If training set size is less than 2000 use batch gradient descent.
  (ii) It has to be power of 2
  (iii) Make sure it fits in CPU/GPU memory

## Video 3: Exponentially Weighted Averages

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



Exponentially weighted averages are useful for optimizing gradient descent algorithm. It gives different $\theta$ based on $\beta$. This reduces oscillations in gradient descent and makes smooth path towards minima

## Video 4: Understanding Exponentially weighted Averages

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

average is represented over $\dfrac{1}{1-\beta}$ entries

best beta is between 0.9 & 0.98
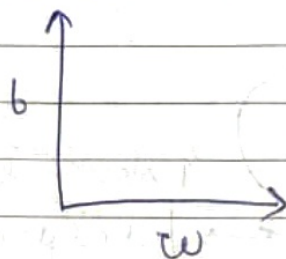
$$V_0 = 0$$

## Video 7: RMS Prop

Modified algorithm:

On iteration t:

    compute $dw$, $db$ on current minibatch

                     element wise

$$Sdw = \beta_2 Sdw + (1-\beta_2)dw^2 \leftarrow \text{small}$$
$$Sdb = \beta_2 Sdb + (1-\beta_2)db^2 \leftarrow \text{large}$$



$\varepsilon$ is added so that $sdw \neq 0$

$$W := w - \alpha \frac{dw}{\sqrt{Sdw} + \varepsilon} \qquad\qquad b = b - \alpha \frac{db}{\sqrt{Sdb} + \varepsilon}$$

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

## Video 8: Adam optimisation

* RMSprop + momentum = adam optimisation
* ADAM - adaptive moment estimation

* Pseudocode / Algorithm.

on iteration t, compute $dw$, $db$ using current mini batch.

$$\left.\begin{array}{l} Vdw = \beta_1 vdw + (1-\beta_1)dw \\ Vdb = \beta_1 vdb + (1-\beta_1)db \end{array}\right\} \text{momentum}$$

5

$$S_{dw} = \beta_2 S_{dw} + (1-\beta_2) dw^2 \quad \Big\} \text{ RMS prop}$$
$$S_{db} = \beta_2 S_{db} + (1-\beta_2) db^2$$

corrected
$$V_{dw} = \frac{V_{dw}}{(1-\beta_1^t)} \qquad V_{db} = \frac{V_{db}}{(1-\beta_1^t)}$$

corrected
$$S_{dw} = \frac{S_{dw}}{(1-\beta_2^t)} \qquad S_{db} = \frac{S_{db}}{(1-\beta_2^t)}$$

corrected
$$w = w - \alpha \frac{V_{dw}}{\sqrt{S_{dw}} + \varepsilon} \qquad b = b - \alpha \frac{V_{db}}{\sqrt{S_{db} + \varepsilon}}$$

### Hyperparameters

$$\alpha - \text{to be tuned}$$
$$\beta_1 = 0.9 \ (dw)$$
$$\beta_2 = 0.99 \ (dw^2)$$
$$\varepsilon = 10^{-8}$$

## Video 9: Learning Rate Decay

While implementing mini batch gradient descent your steps will be nosy and it won't converge at minimum

6

But as you reduce $\alpha$ with time, at the beginning learning rate would be fast, but then as $\alpha$ keeps decreasing, small $\alpha$ will help oscillate in tighter region around minimum.

$$\alpha = \frac{1}{1 + decay\ rate * epochnumber} \alpha_0$$

hyperparameter

Other methods:

$$\alpha = 0.95^{epoch\ number} \times \alpha_0$$

$$\alpha = \frac{k}{\sqrt{epoch number}} \times \alpha_0$$

changes to learning rate can be made discretely - decrease after some no. of epochs. otherwise manually.

Video 10: The problem of local optima

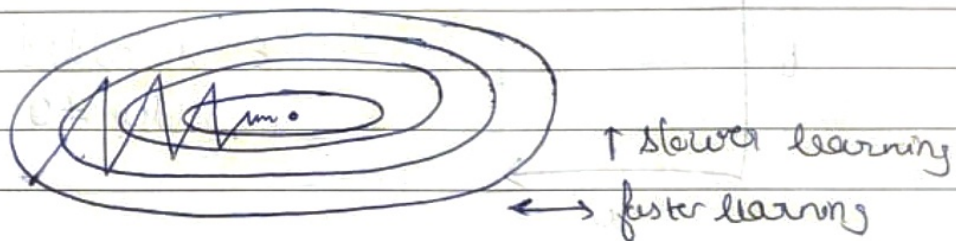The normal local optima is not likely to appear in a deep NN

Plateau - region of slow learning

Video 5: Bias correction in exponentially weighted avg.

Bias correction makes exponentially weighted averages smooth and more accurate.

$$\frac{V_t}{1-\beta t} = V_t \qquad V_t \leftarrow \frac{V_t}{1+\beta t} \qquad V_t = \frac{V_t}{1-\beta^t}$$

Video 6: Gradient Descent with Momentum



↑ slower learning

← faster learning

* Oscillations slow down gradient descent and prevent you from using a larger learning rate.

* Larger learning rate would cause to oscillation to shoot up.

* Implementing gradient descent + momentum

Momentum:
On iteration t :
    compute $dw, db$, on current mini batch
    $V_{dw} = \beta V_{dw} + (1-\beta)dw$      $V_\theta = \beta V_{\theta t} + (1-\beta)\theta$
    $V_{db} = \beta V_{db} + (1-\beta)db$

hyperparameters

    $w = w - \alpha V_{dw}$
    $b = b - \alpha V_{db}$      $\alpha, \beta$