

Mini-batch gradient descent

Vectorization allows to efficiently compute on m examples
 \Rightarrow however m is big \Rightarrow it can be slow

$$X = \left[\underbrace{x^1, x^2, \dots, x^{1000}}_{x^{(1)} \text{ (} n \times 1000 \text{)}} \mid \underbrace{x^{1001}, \dots, x^{2000}}_{x^{(2)} \text{ (} n \times 1000 \text{)}} \mid \dots \mid \dots, x^m \right]$$

$(n \times m)$

$$Y = \left[\underbrace{y^1, y^2, \dots, y^{1000}}_{y^{(1)} \text{ (} 1, 1000 \text{)}} \mid \underbrace{y^{1001}, \dots, y^{2000}}_{y^{(2)} \text{ (} 1, 1000 \text{)}} \mid \dots \mid \dots, y^m \right]$$

$(1, m)$

mini-batches

$m = 5000000 \Rightarrow 5000$ mini-batches each of 1000 examples

mini-batch t : $x^{(t)}, y^{(t)}$

one step of gradient descent using $x^{(t)}, y^{(t)}$
mini-batches:

Repeat until convergence {

for $t = 1 \dots 5000$ {

Forward propagation on $x^{(t)}$

$$\left. \begin{aligned} z^l &= w^l x^{(t)} + b^l \\ a^l &= f'(z^l) \\ &\vdots \\ a^l &= f^l(z^l) \end{aligned} \right\} \text{use vectorized implementation}$$

compute cost: \downarrow for $x^{(t)}, y^{(t)}$

$$J^{(t)} = \frac{1}{1000} \sum_{i=1}^L h(\bar{y}^i, \hat{y}^i) + \frac{\lambda}{2 \cdot 1000} \sum_l \|w^l\|_F^2$$

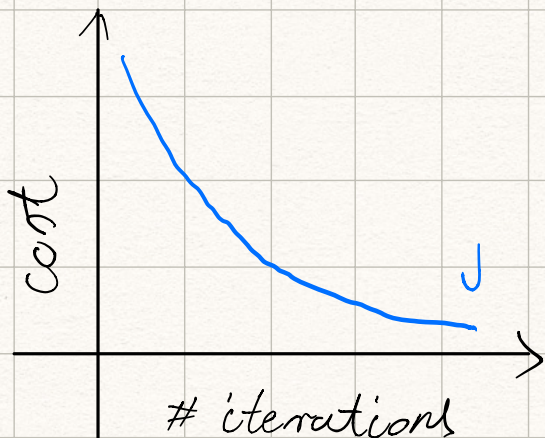
Backprop to compute gradient cost $J^{(t)}$

$$w^l := w^l - \eta d w^l, \quad b^l := b^l - \eta d b^l$$

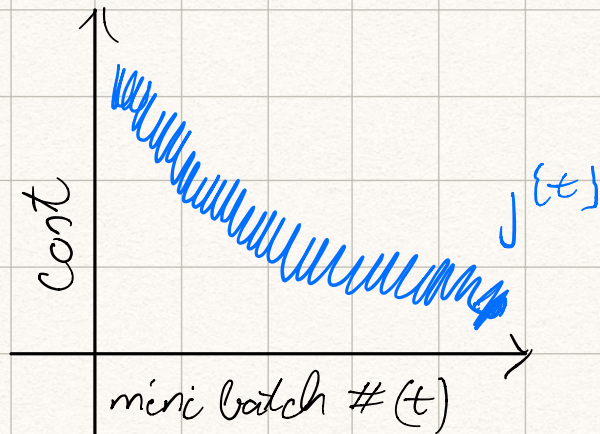
} \rightarrow 1 epoch (single pass through the training set)

Training with mini-batch gradient descent

Batch gradient descent



Mini-batch gradient descent



$x^{(1)}, y^{(1)}$ mini batch may be smaller than $x^{(2)}, y^{(2)}$ \Rightarrow
 \Rightarrow oscillation

Choosing mini-batch size

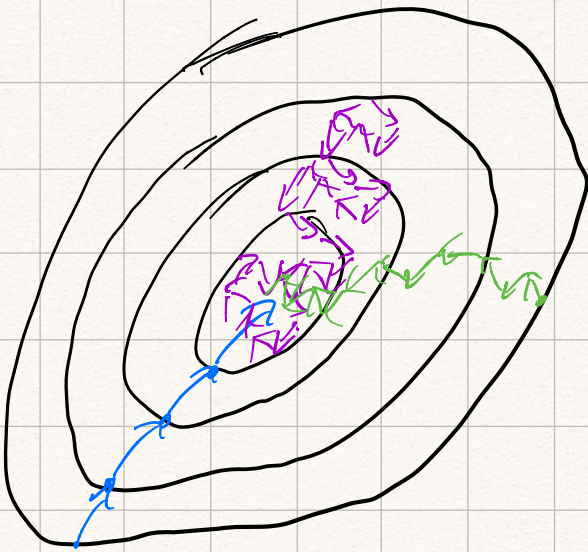
If mini-batch size = m \Rightarrow Batch gradient descent

If mini-batch size = 1 \Rightarrow Stochastic gradient descent

In practice $\Rightarrow 1 \leq \text{size} \leq m$

SGD: lose speedup from vectorization

BGD: Too long per iteration



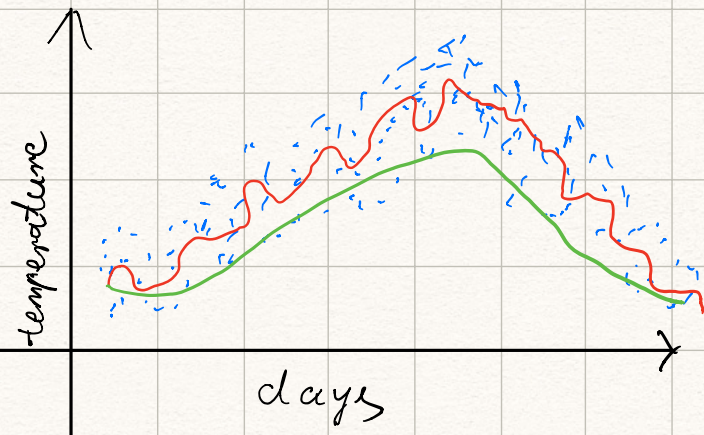
- Fastest learning from vectorization
- make progress without process the entire dataset

If small training set \Rightarrow use batch gradient descent
else \Rightarrow use mini-batch

with size: 64, 128, 256, 512

Make sure mini-batch fit in CPU/GPU memory

Exponentially weighted averages



$$\begin{aligned}V_0 &= 0 \\V_1 &= 0.9V_0 + 0.1\Theta \\V_2 &= 0.9V_1 + 0.1\Theta \\&\vdots \\V_t &= 0.9V_{t-1} + 0.1\Theta\end{aligned}$$

hyperparameter

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

$$\approx \frac{1}{1-\beta} \text{ days}$$

$\beta = 0.9 \Rightarrow 10 \text{ days average}$

$\beta = 0.98 \Rightarrow 50 \text{ days average}$

math stuff

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

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

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

$$\downarrow 0.1\Theta_{99} + 0.9V_{98}$$

$$\downarrow 0.1\Theta_{98} + 0.9V_{97}$$

$$\rightarrow V_{100} = 0.1\Theta_{100} + 0.9V_{99}$$

$$= 0.1\Theta_{100} + 0.1 \cdot 0.9\Theta_{99} + 0.1 \cdot (0.9)^2\Theta_{98} + 0.1 \cdot (0.9)^3\Theta_{97} + \dots$$

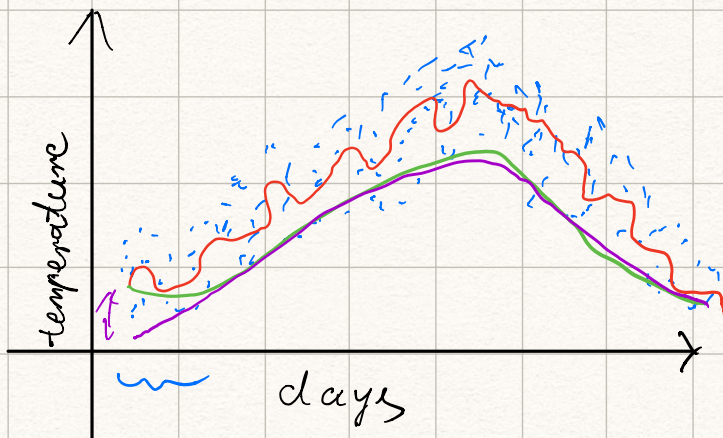
Implementing

$$V_0 = 0 \quad (\Leftarrow \text{only one variable to store})$$

Repeat {
 get next Θ_t

} $V_0 := \beta V_0 + (1-\beta)\Theta_t$

Bias connection in exponentially weighted average



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

\Rightarrow actually we get the purple curve

$$V_0 = 0$$

$$V_1 = 0.98 V_0 + 0.02 \theta_1 \Rightarrow V_1 = \text{much lower value than the valid one}$$

$$V_2 = 0.98 V_1 + 0.02 \theta_2 = 0.98 \cdot 0.02 \theta_1 + 0.02 \theta_2 \Rightarrow$$

$$\frac{V_t}{1-\beta^t} \Rightarrow t=2; 1-\beta^t = 1-0.98^2 = 0.0396$$

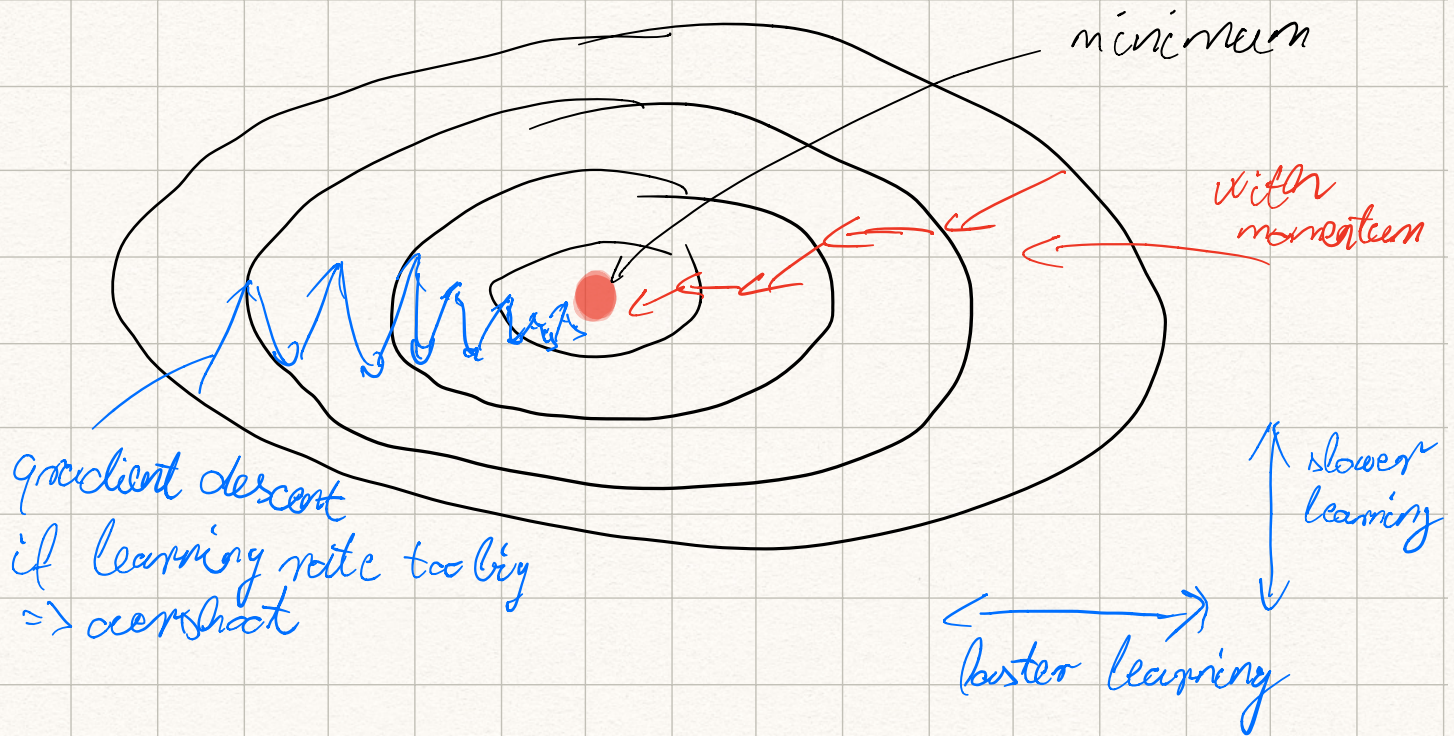
$$\frac{V_2}{0.0396} = \frac{0.0196 \theta_1 + 0.02 \theta_2}{0.0396} \Rightarrow \text{removing bias}$$

if t is large $\Rightarrow \beta^t \approx 0 \Rightarrow$ no effect to bias
 \Rightarrow purple and green line are pretty overlap

\Rightarrow using exponentially weighted average to build better optimizer

Gradient descent with momentum

- ↳ use exponentially weighted average to compute derivatives
- ↳ almost always faster than gradient descent



Momentum: friction velocity

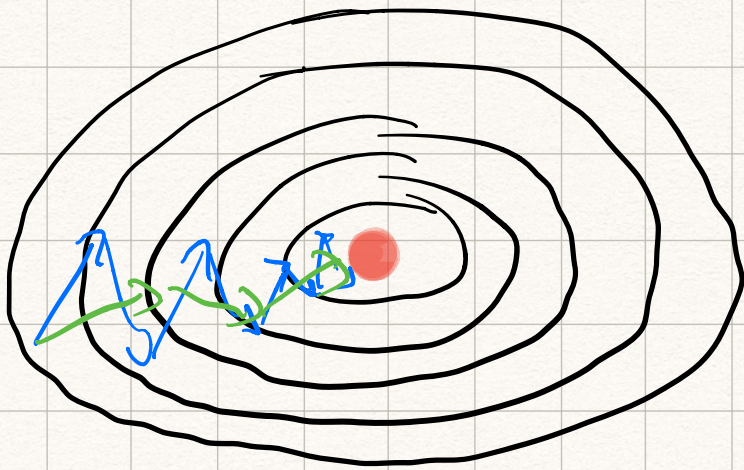
iterate t :
compute d_w, d_b on current mini-batch

$$V_{d_w} = \beta V_{d_w} + (1-\beta) |d_w| \Rightarrow \text{smooth out gradient descent}$$
$$V_{d_b} = \beta V_{d_b} + (1-\beta) |d_b|$$

$$w := w - \eta V_{d_w}, \quad b := b - \eta V_{d_b} \quad \text{acceleration}$$

Hyperparameters: η, β → usually $\beta = 0.9$
≈ last 10 example average

RMSprop (root mean square propagation)



w, b are high dimensional parameters
this is just an intuition

slower
G
w
faster

iterate t :

compute dw, db on current mini-batch

$$S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) |dw|^2 \leftarrow \text{small}$$

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

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

$\epsilon = 10^{-8}$ small number

\Rightarrow in order to do not divide by zero

Adam optimization algorithm

\hookrightarrow work well on wide range of architectures

\rightarrow use no momentum and RMSprop

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

iterate t :

compute dw, db on current mini-batch

$$\begin{aligned} v_{dw} &= \beta_1 v_{dw} + (1 - \beta_1) dw \\ v_{db} &= \beta_1 v_{db} + (1 - \beta_1) db \end{aligned} \quad \left. \vphantom{\begin{aligned} v_{dw} \\ v_{db} \end{aligned}} \right\} \text{momentum like}$$

$$\begin{aligned} s_{dw} &= \beta_2 s_{dw} + (1 - \beta_2) dw^2 \\ s_{db} &= \beta_2 s_{db} + (1 - \beta_2) db^2 \end{aligned} \quad \left. \vphantom{\begin{aligned} s_{dw} \\ s_{db} \end{aligned}} \right\} \text{RMSprop like}$$

$$\begin{aligned} v_{dw}^{\text{corrected}} &= v_{dw} / (1 - \beta_1^t) \\ v_{db}^{\text{corrected}} &= v_{db} / (1 - \beta_1^t) \end{aligned} \quad \left. \vphantom{\begin{aligned} v_{dw}^{\text{corrected}} \\ v_{db}^{\text{corrected}} \end{aligned}} \right\} \text{bias correction}$$

$$\begin{aligned} s_{dw}^{\text{corrected}} &= s_{dw} / (1 - \beta_2^t) \\ s_{db}^{\text{corrected}} &= s_{db} / (1 - \beta_2^t) \end{aligned} \quad \left. \vphantom{\begin{aligned} s_{dw}^{\text{corrected}} \\ s_{db}^{\text{corrected}} \end{aligned}} \right\} \text{bias correction}$$

$$w := w - \eta \frac{v_{dw}^{\text{corrected}}}{\sqrt{s_{dw}^{\text{corrected}} + \epsilon}} \quad , \quad b := b - \eta \frac{v_{db}^{\text{corrected}}}{\sqrt{s_{db}^{\text{corrected}} + \epsilon}}$$

=> combine the effect of gradient descent with momentum and with RMSprop

↳ we can use larger learning rate

Hyperparameters choice:

η : needs to be tune

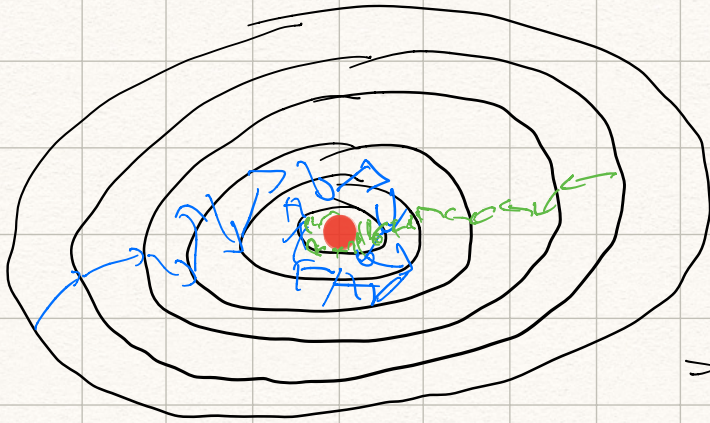
β_1 : ~ 0.9

β_2 : ~ 0.999

ϵ : 10^{-8}

Adam => Adaptive moment estimation

learning rate decay



Noisy mini-batch gradient descent \rightarrow never convergence, just oscillate around the minimum

\Rightarrow learning rate decay \Rightarrow take smaller steps \Rightarrow better convergence



$$\hat{\eta} = \frac{1}{1 + \text{decay_rate} \cdot \text{epoch_num}} \cdot \hat{\eta}_0$$

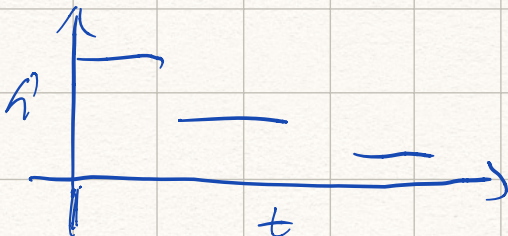
other hyperparameter

$$\hat{\eta} = 0.95^{\text{epoch_num}} \cdot \hat{\eta}_0 \Rightarrow \text{exponentially decay}$$

$$\hat{\eta} = \frac{2}{\sqrt{\text{epoch_num}}} \cdot \hat{\eta}_0$$

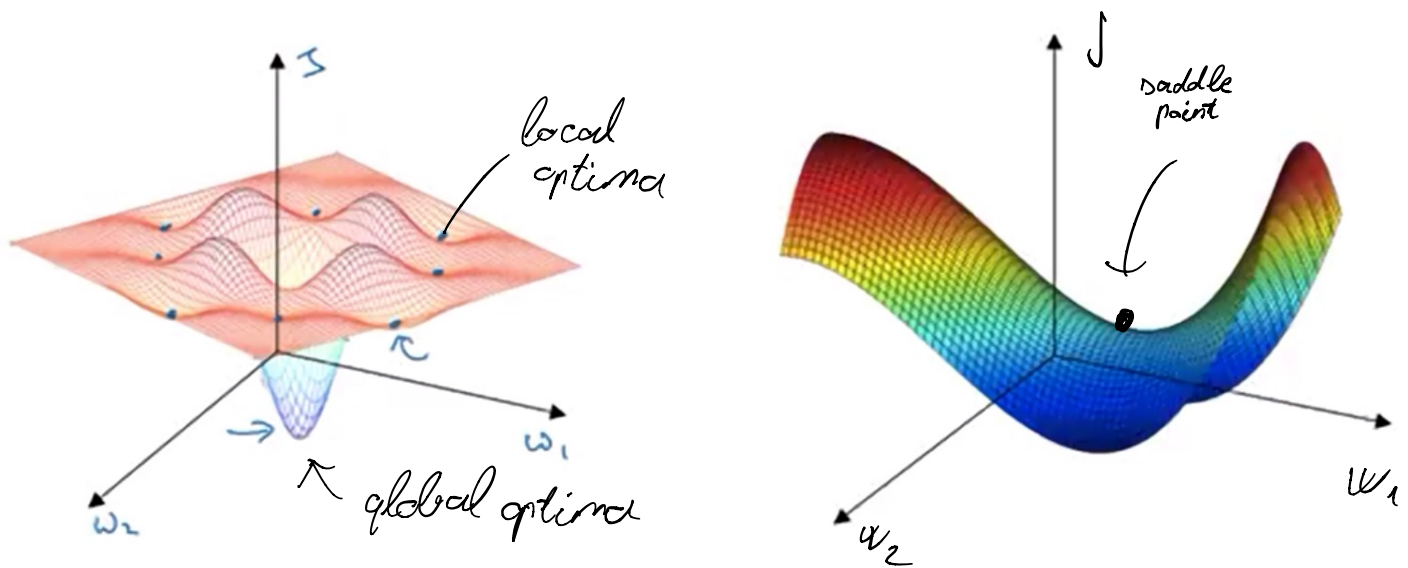
$$\hat{\eta} = \frac{2}{\sqrt{t}} \cdot \hat{\eta}_0$$

mini batch number



\Rightarrow discrete decay

The problem of local optima

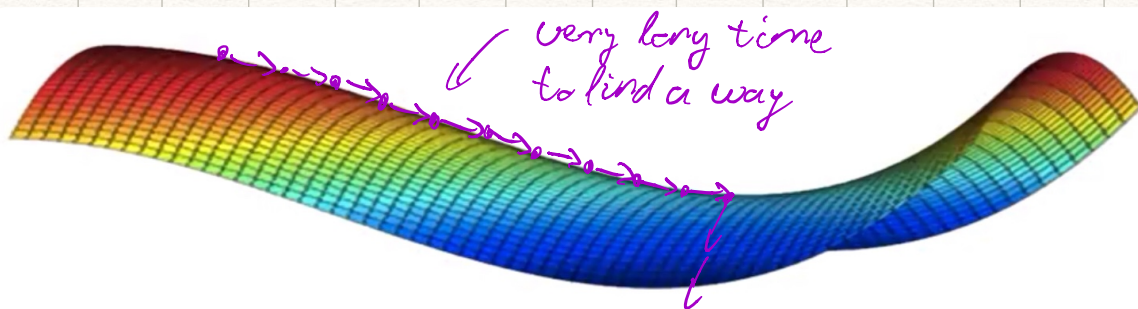


In two dimensions easy to create local optima

- but in high dimension space most points with Zero gradient it is a saddle point
- to be a local optima in each direction have to be a convex or concave curve

\Rightarrow high dim $\approx 20000 \Rightarrow p = 2^{-20000}$ the chance to be a local optima

Plateaus problem (derivative is close to zero for a long time)



unlikely to get stuck in a bad local optima \Rightarrow use large NN with many parameters

Plateaus can make learning slow \Rightarrow momentum, Adam