

- Mini-batch gradient descent

- Batch VS. mini-batch gradient descent

$$X = [X^{[1]}, X^{[2]}, X^{[3]}, \dots, X^{[1000]}, X^{[1001]}, \dots | X^{(m)}]$$

what if $m = 5,000,000$?

mini-batch:

baby training set

$$X = [\underbrace{X^{[1]}, X^{[2]}, X^{[3]}, \dots, X^{[100]}}, \underbrace{X^{[101]}, X^{[102]}, \dots}_{x^{[1]}, \dots}, \underbrace{X^{[103]}, \dots}_{x^{[2]}, \dots} | X^{(m)}]$$

do the same for y accordingly

$$\text{mini-batch } t = X^{[t]}, Y^{[t]}$$

- Algorithm

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

Forward prop on $X^{[t]}$:

$$\begin{aligned} z^{[1]} &= w^{[1]} X^{[t]} + b^{[1]} \\ z^{[2]} &= g^{[2]}(z^{[1]}) \\ &\vdots \quad \vdots \quad \vdots \end{aligned} \quad \left. \begin{array}{l} \text{Vectorized implementation} \\ (\text{e.g. 1000 examples}) \end{array} \right\}$$

$$\text{Compute cost } J = \frac{1}{1000} \sum_{i=1}^l \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2 \cdot 1000} \sum_i \|w^{[i]}\|^2$$

Back prop to compute gradients of $J^{[t]}$

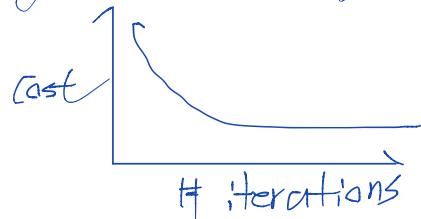
$$w^{[t]} = w^{[t]} - \alpha \nabla_w^{[t]}, \quad b^{[t]} = b^{[t]} - \alpha \nabla_b^{[t]}$$

}

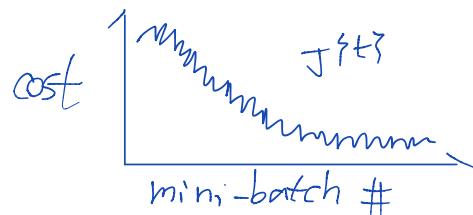
"1 epoch": a single pass through training set
Much faster than batch gradient

- Understanding mini-batch gradient descent

- Batch gradient descent



- Mini-batch gradient descent



- Choosing your mini-batch size

- If mini-batch size = m

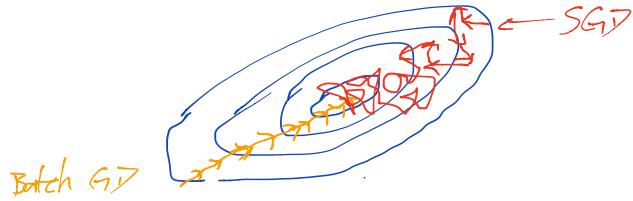
→ batch gradient descend

- If mini-batch size = 1

→ stochastic gradient descent

Won't ever converge, always oscillates around the min cost.

- In practice: $1 < < m$



Compare

- batch GD: Too long per iteration
- SGD: lose speedup from vectorization
making progress every iteration
can reduce noise by using a small learning rate
- mini-batch: Fastest Learning !!!
 - vectorization
 - make progress w/o wait too long

Guidelines:

- If small training set : Use batch GD
- Typical mini-batch size : 64, 128, 256, 512, ...
- Make sure each mini-batch fit CPU/GPU memory

Exponentially weighted averages (Faster than gradient descent)

e.g.: Temperature in London

$$V_0 = \theta_0 \quad (\text{day 0 temp})$$

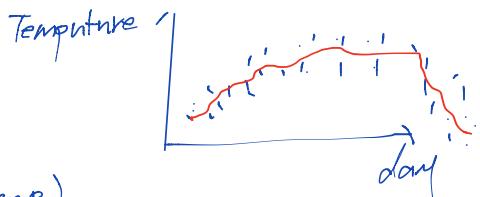
$$V_1 = 0.9V_0 + 0.1\theta_1 \quad (\text{day 1 temp})$$

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

:

:

:



\longrightarrow : exponentially weighted temperature

Exponentially weighted averages

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

V_t : can be viewed as approx. average over
 $\approx \frac{1}{1-\beta}$ days temperature

e.g.: $\beta = 0.9 \rightarrow$ avg. over last 10 days temp

$\beta = 0.98 \rightarrow \dots \dots \dots$ 50 days ...

Understanding Exponentially weighted averages

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

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

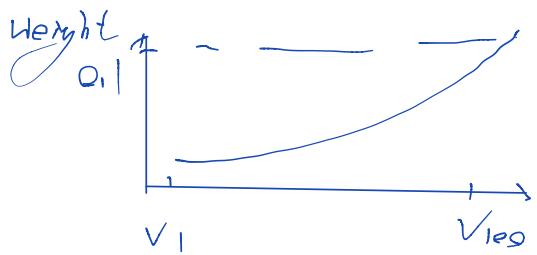
$$V_{99} = 0.9V_{98} + 0.1\theta_{99}$$

:

:

:

$$\begin{aligned} V_{100} &= 0.1\theta_{100} + 0.1 \times 0.9 \cdot \theta_{99} + 0.1 \times (0.9)^2 \theta_{98} \\ &\quad + 0.1 \times (0.9)^3 \theta_{97} + 0.1 \times (0.9)^4 \theta_{96} + \dots \end{aligned}$$



$$(1-\varepsilon)^{\frac{1}{\varepsilon}} \approx \frac{1}{e}$$

$0.9^{10} \approx 0.35 \approx \frac{1}{e}$: takes about 10 days of previous temp into account

- Implementation

$$V_0 = 0$$

Repeat {

get next θ_t

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

} (takes little memory)

- Bias Correction In Exponentially Weighted Average

- Bias Correction

$$V_0 = 0$$

$$V_1 = \underbrace{0.98 V_0}_{=0} + 0.02 \theta_1$$

$$V_2 = 0.98 V_1 + 0.02 \theta_2$$

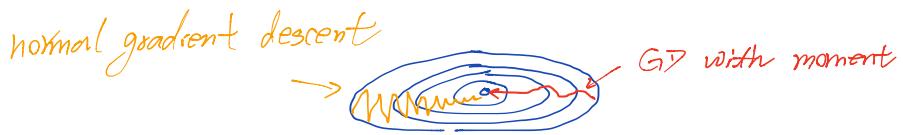
⋮ ⋮ ⋮

Bias raises from the starting boundary!

- Correction

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

- Gradient Descent Momentum



Wort: faster learning along horizontal axis
slower learning along vertical axis

- Moment

on iteration t :

Compute δw , δb on current batch

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

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

$$w = w - \alpha V_{\delta w}$$

$$b = b - \alpha V_{\delta b}$$

- Understanding:

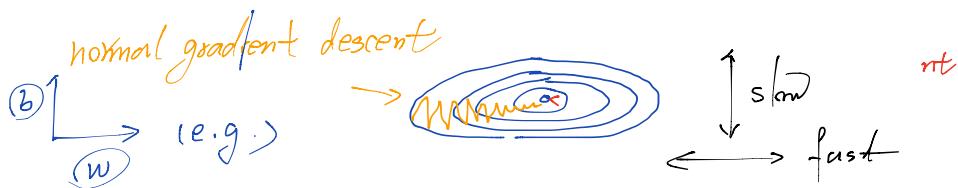
Average out fast oscillating gradient descent steps \rightarrow more efficient!

$$\begin{aligned} \dot{V}_{\partial w} &= \beta \dot{V}_{\partial w} + (1-\beta) \dot{\partial w} \quad ; \quad \text{derivative term} \\ \dot{V}_{\partial b} &= \beta \dot{V}_{\partial b} + (1-\beta) \dot{\partial b} \quad ; \quad \hookrightarrow \text{Acceleration!} \\ \text{momentum} \\ \hookrightarrow \text{velocity} \\ \beta: \text{friction} \end{aligned}$$

JUST LIKE
rolling a ball !!

hyperparameters: α, β

2 RMSprop



On iteration $t =$

Compute $\partial b, \partial w$ on current mini-batch:

$$S_{\partial w} = \beta S_{\partial w} + (1-\beta) \partial w^2 \leftarrow \text{element wise}$$

$$S_{\partial b} = \beta S_{\partial b} + (1-\beta) \partial b^2$$

$$w := w - \alpha \frac{\partial w}{\sqrt{S_{\partial w}}}$$

$$b := b - \alpha \frac{\partial b}{\sqrt{S_{\partial b}}}$$

o Understanding

$$w := w - \alpha \frac{\partial w}{\sqrt{S_{\partial w}}} \quad \text{Want to be large to learn fast}$$

$$b := b - \alpha \frac{\partial b}{\sqrt{S_{\partial b}}} \quad \begin{aligned} &\text{Want to be small} \\ &\text{to learn slow to} \\ &\text{reduce oscillation} \end{aligned}$$

from the e.g.

$$\partial b \gg \partial w$$

$$\hookrightarrow S_{\partial b} \gg S_{\partial w}$$

$$\hookrightarrow \frac{\partial b}{\sqrt{S_{\partial b}}} > \frac{\partial w}{\sqrt{S_{\partial w}}} \quad \text{Damp The Oscillations!!}$$

Sidenote

$$S_{\partial b} \rightarrow S_{\partial b} + \varepsilon \quad \varepsilon \ll 1$$

$$S_{\partial w} \rightarrow S_{\partial w} + \varepsilon$$

To improve stability (divide by ε)

o Adam Optimization Algorithm

Adaptive Moment Estimation (Adam)

$$\text{Momentum} + \text{RMS} + \text{GD} = \text{Adam}$$

↑ ↑
1st momentum 2nd momentum

o Algorithm

$$V_{\partial w} = 0, S_{\partial w} = 0, V_{\partial b} = 0, S_{\partial b} = 0$$

On iteration t -

Compute $\partial w, \partial b$ using current batch:

$$V_{\partial w} = \beta_1 V_{\partial w} + (1 - \beta_1) \partial w$$

$$V_{\partial b} = \beta_1 V_{\partial b} + (1 - \beta_1) \partial b$$

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

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

Bias correction :

$$V_{\partial w}^c = \frac{V_{\partial w}}{1 - \beta_1}$$

$$V_{\partial b}^c = \frac{V_{\partial b}}{1 - \beta_1}$$

$$S_{\partial w}^c = \frac{S_{\partial w}}{1 - \beta_2}$$

$$S_{\partial b}^c = \frac{S_{\partial b}}{1 - \beta_2}$$

$$w := w - \alpha \frac{V_{\partial w}^c}{\sqrt{S_{\partial w}^c + \epsilon}}$$

$$b := b - \alpha \frac{V_{\partial b}^c}{\sqrt{S_{\partial b}^c + \epsilon}}$$

Commonly used Algorithm , Very effective !

- Hyperparameter choice

α : needs to be tuned

β_1 : common choice of (∂w)

β_2 : recommended by literature $0.999 (\partial w)^2$

ϵ : doesn't matter too much , 10^{-8}

- Learning Rate Decay

slowly reduce the learning rate over time

mini-batch GD:

use learning rate decay
to reduce oscillations!

1 epoch: 1 pass through all data

e.g. $\alpha = \frac{1}{1 + \text{decay-rate} * \text{epoch num}} \alpha_0$

e.g. $\alpha_0 = 0.2$ decay rate = 1

Epoch	α
1	0.1
2	0.67
3	0.5
4	0.4
:	:

e.g. $\alpha = 0.95^{\text{epoch_num}} \cdot \alpha_0$

$$\alpha = \frac{k}{\text{epoch_num}} \cdot \alpha_0 \quad \text{or} \quad \frac{1}{\sqrt{E}} \alpha_0$$

OR: Manual Decay

- The problem of Local Optima

There are more saddle points than local optimal in ML practice, especially for high dim, the chance for local optimal is very low!

- Problem of Plateaus (saddle points)

⌘ slow down learning rate