

Week 2 - Optimization algorithms

Friday, August 14, 2020 11:26 PM

① Mini-batch gradient descent

→ Gradient descent requires all training examples to be processed to start training

→ This slows down process when data is extremely large

→ Hence, mini-batches are used

e.g.: if $m = 5,000,000$ (5 million)

$$\therefore \mathbf{X} = [\underbrace{\mathbf{x}^{(1)} \mathbf{x}^{(2)} \dots \mathbf{x}^{(100)}}_{\mathbf{X}^{(1)} \{1, 100\}} \mid \underbrace{\mathbf{x}^{(100)} \dots \mathbf{x}^{(2000)}}_{\mathbf{X}^{(2)} \{1, 1000\}} \mid \dots \mid \mathbf{x}^{(m)}]$$

$\xleftarrow{\text{minibatch}}$

$$\mathbf{y} = [\underbrace{\mathbf{y}^{(1)} \mathbf{y}^{(2)} \dots \mathbf{y}^{(100)}}_{\mathbf{y}^{(1)} \{1, 100\}} \mid \underbrace{\mathbf{y}^{(100)} \dots \mathbf{y}^{(2000)}}_{\mathbf{y}^{(2)} \{1, 1000\}} \mid \dots \mid \mathbf{y}^{(m)}]$$

$\xleftarrow{\text{minibatch}}$

→ Mini-batch gradient descent

Mini-batch gradient descent

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

1 step of gradient descent
using $\frac{\mathbf{X}^{t+1}, \mathbf{y}^{t+1}}{5000}$

Forward prop on \mathbf{X}^{t+1} .

$$\mathbf{z}^{t+1} = \mathbf{W}^{t+1} \mathbf{X}^{t+1} + \mathbf{b}^{t+1}$$

$$\mathbf{A}^{t+1} = g^{t+1}(\mathbf{z}^{t+1})$$

$$\vdots$$

$$\mathbf{A}^{t+1} = g^{t+1}(\mathbf{z}^{t+1})$$

$$\text{Compute cost } J^{t+1} = \frac{1}{1000} \sum_{i=1}^{1000} l(y^{(i)}, \mathbf{A}^{t+1}) + \frac{\lambda}{2 \cdot 1000} \sum_{j=1}^n \|(\mathbf{W}^{t+1})_j\|_F^2.$$

Backprop to compute gradients w.r.t J^{t+1} (using $(\mathbf{X}^{t+1}, \mathbf{y}^{t+1})$)

$$\mathbf{W}^{t+1} = \mathbf{W}^{t+1} - \alpha \delta \mathbf{W}^{t+1}, \quad \mathbf{b}^{t+1} = \mathbf{b}^{t+1} - \alpha \delta \mathbf{b}^{t+1}$$

}

"1 epoch"
↓ pass through training set.

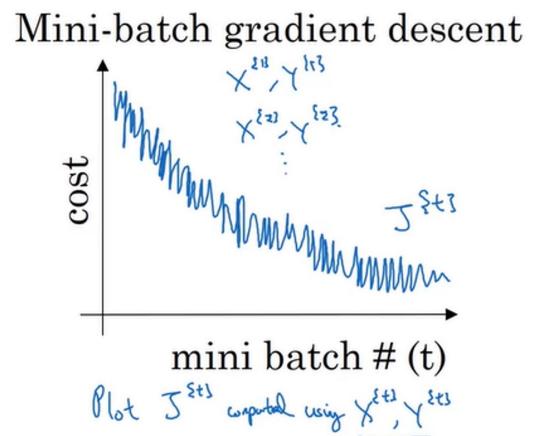
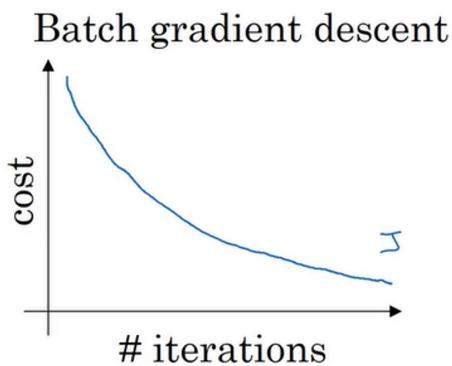
Andrew Ng

② Understanding mini-batch gradient descent

→ Training with mini batch gradient descent

Training with mini batch gradient descent

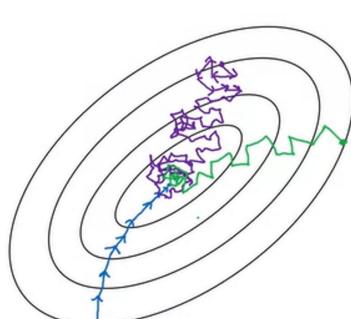
Training with mini batch gradient descent



Andrew Ng

→ Choosing mini-batch size

- If mini-batch size = m : Batch gradient descent
Too long per iteration $(x^{S_{t+3}}, y^{S_{t+3}}) = (x, y)$
- If mini-batch size = 1: Stochastic gradient descent
never reaches optimum but revolves around it $(x^{S_{t+3}}, y^{S_{t+3}}) = (x^{(1)}, y^{(1)}) \dots (x^{(n)}, y^{(n)})$
↳ lose all speed benefits of vectorization
- In practice: between 1 & m
↳ faster learning & vectorization



Stochastic gradient descent
↳
Use speedup from vectorization

In-between
(minibatch size not too big/small)

Faster learning.

- Vectorization (≈ 1000)
- Make passes without processing entire training set.

Batch gradient descent
(minibatch size = m)

Two long per iteration

Andrew Ng

- If small training set: Use batch gradient descent ($m \leq 2000$)

- If large training set:
typical mini-batch size:
64, 128, 256, 512 ...

make sure minibatch fits in CPU/GPU memory

③ Exponentially weighted averages

Temperature in London

$$\theta_1 = 40^\circ\text{F} \quad 4^\circ\text{C} \leftarrow$$

$$\theta_2 = 49^\circ\text{F} \quad 9^\circ\text{C}$$

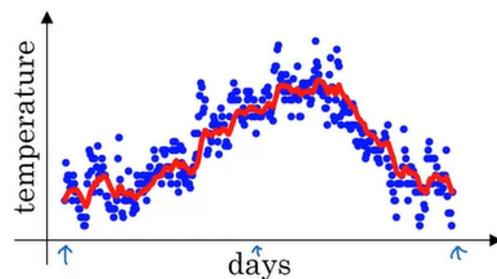
$$\theta_3 = 45^\circ\text{F} \quad \vdots$$

$$\vdots$$

$$\theta_{180} = 60^\circ\text{F} \quad 15^\circ\text{C}$$

$$\theta_{181} = 56^\circ\text{F} \quad \vdots$$

$$\vdots$$



$$v_0 = 0$$

$$v_1 = 0.9 v_0 + 0.1 \theta_1$$

$$v_2 = 0.9 v_1 + 0.1 \theta_2$$

$$v_3 = 0.9 v_2 + 0.1 \theta_3$$

$$\vdots$$

$$v_t = 0.9 v_{t-1} + 0.1 \theta_t$$

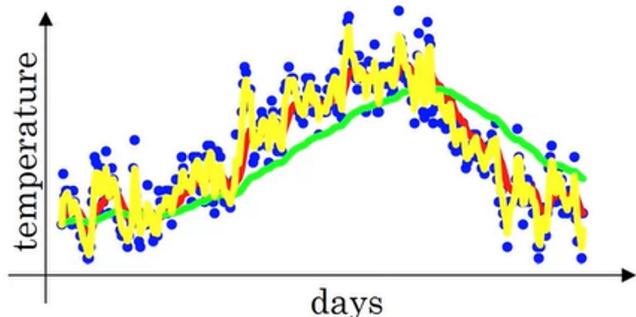
Andrew Ng

$$v_t = \beta v_{t-1} + (1-\beta) \theta_t$$

($\beta = 0.9$ in above case)

v_t is approximately average over $\approx \frac{1}{1-\beta}$ days

For $\beta = 0.9$, $\frac{1}{0.1} = 10$ days moving average
For $\beta = 0.98$, 50 days moving average
For $\beta = 0.5$, 2 day moving avg.



④ Understanding exponentially weighted averages

$$v_t = \beta v_{t-1} + (1-\beta) \theta_t$$

consider

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

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

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

$$\vdots$$

$$\hookrightarrow v_{100} = 0.1 \theta_{100} + 0.9 (0.1 \theta_{99} + 0.9 v_{98})$$

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

↪ all coeff add upto ≈ 1

- when β^n reaches $\approx \frac{1}{e} \approx 0.36$ it is insignificant

$$\therefore \beta = 0.9 \Rightarrow 0.9^n \approx 0.35$$

- when β reaches $\approx \frac{1}{e} \approx 0.36$ it is insignificant
 $\therefore \beta = 0.9 \Rightarrow 0.9^{\infty} \approx 0.36$
 $\therefore 10$ days of values are significant
 $\beta = 0.98 \Rightarrow 0.98^{\infty} \approx 0.36$
 $\therefore 50$ days of values are significant
- This is how we get $\frac{1}{1-\beta}$ values formula
 ↳ general rule of thumb

- Implementation

$$v_0 = 0$$

$$v_1 = \beta v_0 + (1-\beta) \theta_1$$

$$v_2 = \beta v_1 + (1-\beta) \theta_2$$

$$v_3 = \beta v_2 + (1-\beta) \theta_3$$

...

→ one variable
is enough
 $v_0 := 0$ and reused

$$v_{\theta} := \beta v + (1-\beta) \theta,$$

$$v_{\theta} := \beta v + (1-\beta) \theta_2$$

⋮

$$\Rightarrow v_{\theta} = 0$$

Repeat {

get next θ_t

$$v_{\theta} := \beta v_{\theta} + (1-\beta) \theta_t$$

}

⑤ Bias correction in exponentially weighted average

$$\rightarrow v_t = \beta v_{t-1} + (1-\beta) \theta_t$$

$$v_0 = 0$$

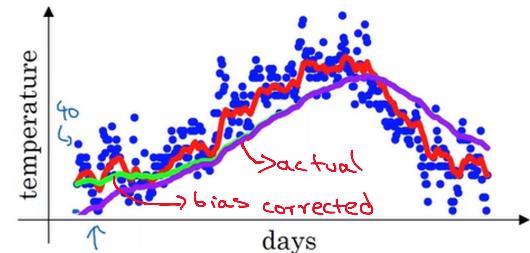
$$v_1 = \cancel{0.02} v_0 + 0.02 \theta_1$$

$$\Rightarrow v_1 = 0.02 \theta_1 \rightarrow \text{very low}$$

$$v_2 = 0.98 v_1 + 0.02 \theta_2$$

$$= 0.98 \times 0.02 \theta_1 + 0.02 \theta_2$$

$$= 0.0196 \theta_1 + 0.02 \theta_2 \rightarrow \text{low}$$



→ Bias correction:

$$\frac{v_t}{1-\beta^t}$$

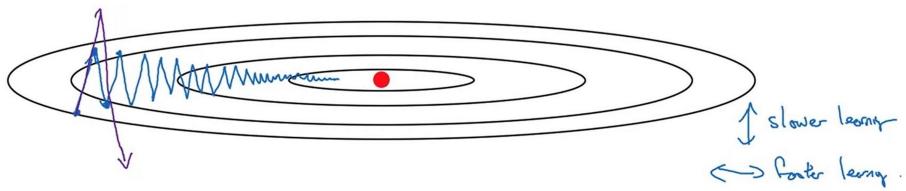
$$\text{for } t=2: 1-\beta^t = 1-(0.98)^2 = 0.0396$$

$$\frac{v_2}{0.0396} = \frac{0.0196 \theta_1 + 0.02 \theta_2}{0.0396}$$

When t gets larger, β^t gets smaller

∴ No effect on large values of t

⑥ Gradient descent with momentum



Momentum:

$v_{dw} = 0, v_{db} = 0$
On iteration t :

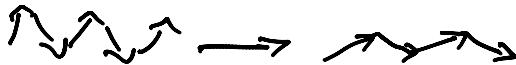
$$v_{dw} = \beta v_{dw} + (1 - \beta) dw \quad (v_w = \beta v_w + (1 - \beta) \theta_t)$$

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

$$w := w - \alpha v_{dw}, b := b - \alpha v_{db}$$

$v_{dw} = \beta v_{dw} + dw$
 \hookrightarrow also used

helps in dampening oscillations



On iteration t :

Compute dW, db on the current mini-batch

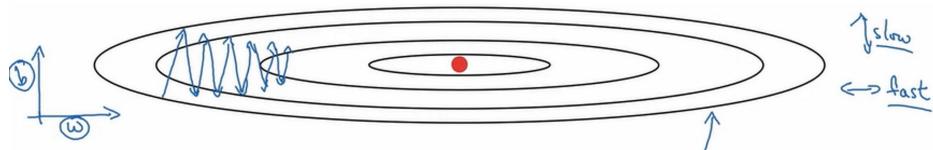
$$v_{dw} = \beta v_{dw} + (1 - \beta) dw$$

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

$$W = W - \alpha v_{dw}, b = b - \alpha v_{db}$$

- Hyperparameters: α, β
 \downarrow usually 0.9
- Bias correction is mostly not used since the bias affects only 10 initial iterations, insignificant on a huge set

⑦ RMS prop



On iteration t :

compute dw, db on current mini-batch

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

$$s_{db} = \beta_2 s_{db} + (1 - \beta_2) db^2 \leftarrow \text{large } \epsilon \sim 10^{-8}$$

$$S_{dw} = \beta_2 S_{dw} + (1-\beta_2) \delta w \leftarrow \text{small}$$

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

$$w := w - \frac{\alpha \delta w}{\sqrt{S_{dw} + \epsilon}} \quad b := b - \frac{\alpha \delta b}{\sqrt{S_{db} + \epsilon}} \quad \begin{matrix} \epsilon \rightarrow 10^{-8} \\ \epsilon \rightarrow \text{ensures not divided by zero} \end{matrix}$$

↓
dividing by small no. \Rightarrow amplifies horizontal mvt.

↓
dividing by large no. \Rightarrow dampens vertical mvt.

⑧ Adam Optimization Algorithm

$$\rightarrow v_{dw} = 0, S_{dw} = 0, v_{db} = 0, S_{db} = 0$$

on iteration t :

compute $\delta w, \delta b$ using current mini batch

$$v_{dw} = \beta_1 v_{dw} + (1-\beta_1) \delta w, v_{db} = \beta_1 v_{db} + (1-\beta_1) \delta b \quad \text{↳ Momentum}$$

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

$$v_{dw}^{\text{corrected}} = v_{dw} / (1 - \beta_1^t), v_{db}^{\text{corrected}} = v_{db} / (1 - \beta_1^t) \quad \text{↳ RMS prop}$$

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

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

\rightarrow Hyperparameters choice:

α : needs tuning

$$\beta_1: 0.9 \quad (\delta w)$$

$$\beta_2: 0.999 \quad (\delta w^2)$$

$$\epsilon: 10^{-8}$$

given in
original paper
published
for ADAM
(adaptive moment
estimation)

⑨ Learning rate decay

\rightarrow During initial steps large learning rate can be used

\rightarrow But as algorithm reaches convergence large overshooting might cause noise or

\rightarrow Hence we implement decay of learning rate every epoch

$$\alpha = \frac{1}{1 + \text{decay rate} * \text{epoch_num}}$$



$$\alpha = \frac{1}{1 + \text{decay rate} * \text{epoch_num}} \cdot \alpha_0$$

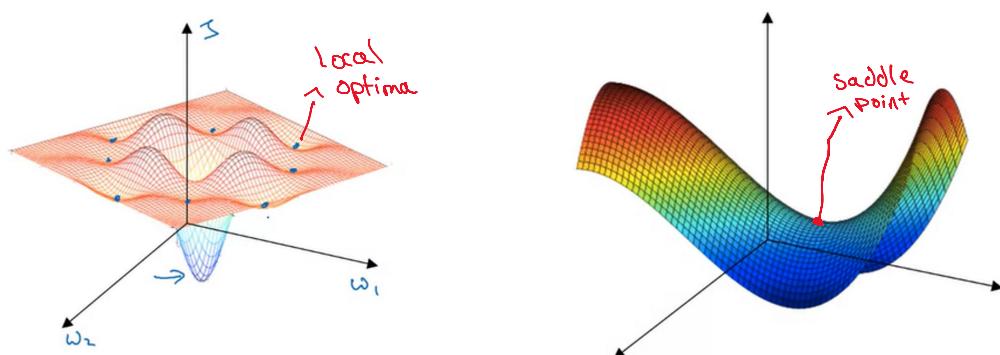
 other formulas: $0.95^{\text{epoch_num}} \cdot \alpha_0$

$$\alpha = \frac{K}{\sqrt{\text{epoch_num}}} \cdot \alpha_0 \text{ (or) } \frac{K}{\sqrt{t}} \cdot \alpha_0$$

discrete staircase

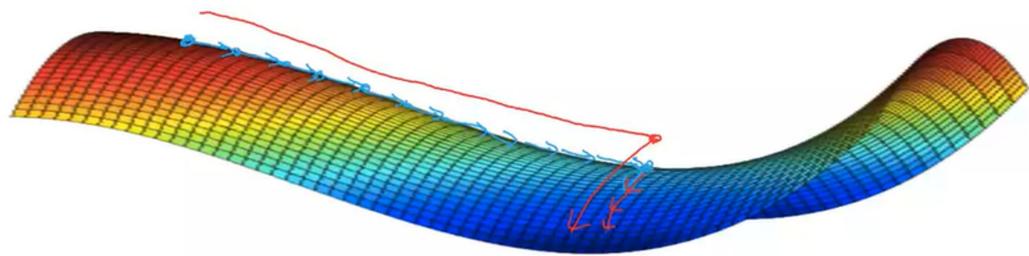
manual decay for training over long time (days)

(10) The problem of local Optima



- In high dimensional space, local optimas are highly improbable
- Saddle points may occur
- Problem of plateaus
 - Gradient is zero for a long time
- Intuitions that we have about space aren't applicable in the high dimensional space where we compute optimas for neural networks

Problem of plateaus



- Unlikely to get stuck in a bad local optima
- Plateaus can make learning slow