

Batch vs Mini-batch Gradient Descent

Vectorization is slow for very large dataset ..

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

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

$$\vec{z}^{[i]} = w^{[i]} \times \{t\} + b^{[i]}$$

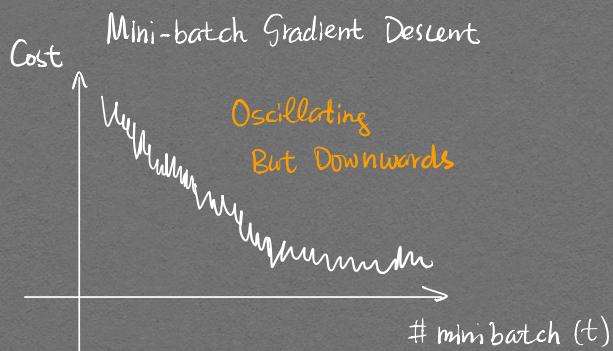
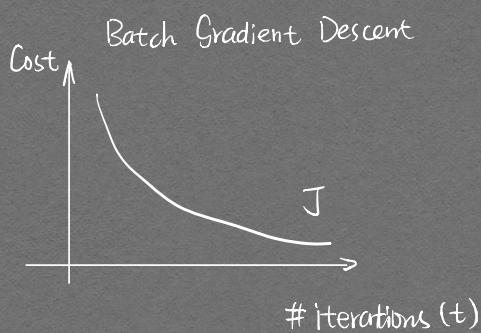
$$A^{[i]} = g^{[i]}(\bar{x}^{[i]})$$

Compute cost $J = \frac{1}{1000} \sum_{i=1}^x \ell(\hat{y}^{(i)}, y^{(i)})$

Backprop to compute gradients wrt $J\{t\}$ (using $x\{t\}, y\{t\}$)

$$w^{[\ell]} := w^{[\ell]} - \alpha d w^{[\ell]}$$

$$b^{[e]} := b^{[e]} - \alpha db^{[e]}$$



Choose Minibatch Size

Single Training Example

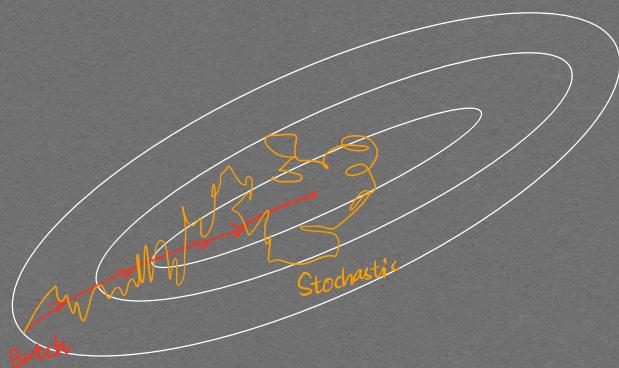
All Training Examples

|

m

Stochastic Gradient Descent

Lose speedup from vectorization



Batch Size

Vectorization Speed Up

Make progress without waiting too long

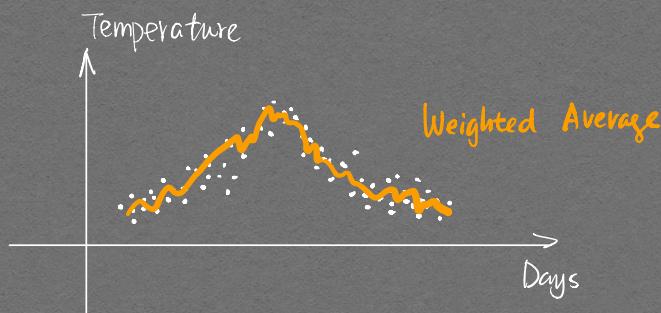
Batch Gradient Descent

Too long per iteration

if $m \leq 2000 \rightarrow$ Batch Gradient

else 64, 128, 256, 512, 1024

Exponentially Weighted Averages



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

≈ Average over $\frac{1}{1-\beta}$ days' temperature

$\beta = 0.9$ 10 day average

$\beta = 0.98$ 50 day average

$\beta = 0.5$ 2 day average

Understanding Exponentially Weighted Averages

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

$$= 0.9 \times (0.9 V_{98} + 0.1 \theta_{99}) + 0.1 \theta_{100}$$

= ...

$$= \underline{0.1 \theta_{100}} + \underline{0.1 \times 0.9 \theta_{99}} + \underline{0.1 \times 0.9^2 \times \theta_{98}} + \dots + \underline{0.1 \times 0.9^{99} \times \theta_1}$$

$$(1-\varepsilon) \frac{V\varepsilon}{e} = \frac{1}{e}$$

how many days are we averaging over approximately

$$V_0 = 0$$

Repeat {

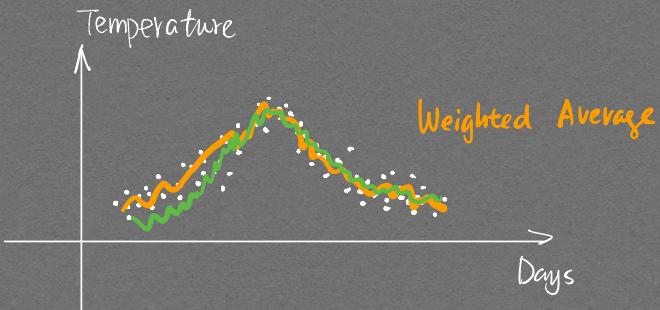
Set next θ_t

Simple & Efficient

$$V_t := \beta V_0 + (1-\beta) \theta_t$$

}

Bias Correction in Exponentially Weighted Average

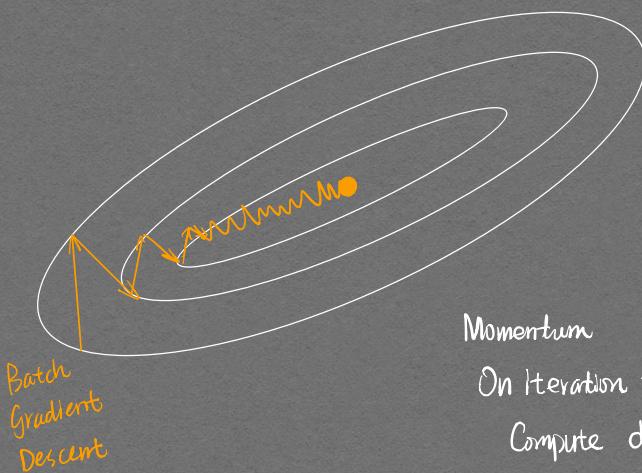


Use $\frac{V_t}{1-\beta^t}$

instead of V_t

especially for initial base

Gradient Descent with Momentum



Oscillating takes time

We want

learning slow on vertical direction

fast on horizontal direction

Momentum

On Iteration t :

Compute dW, db wrt current mini-batch

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

might be
omitted

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

$$W \leftarrow W - \alpha V_{dw}$$

$$b \leftarrow b - \alpha V_{db}$$

$$\beta = .9$$

(Average over last
 ≈ 10 iterations)

RMSprop Group Mean Squared

On Iteration t :

Compute dW, db wrt current mini-batch

$$S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) dW^2 \quad \begin{matrix} \text{Smaller Value} \\ \text{element-wise} \end{matrix}$$

$$S_{db} = \beta_2 S_{db} + (1 - \beta_2) db^2 \quad \text{Larger Value}$$

$$W \leftarrow W - \alpha \frac{dW}{\sqrt{S_{dw}}} + \epsilon \quad \begin{matrix} \text{Quicker} \\ 10^{-8} \text{ to avoid blow up.} \end{matrix}$$

$$b \leftarrow b - \alpha \frac{db}{\sqrt{S_{db}}} \quad \text{Slower}$$

Adam Optimization Algorithm

$$V_{dw} = 0 \quad S_{dw} = 0 \quad V_{db} = 0 \quad S_{db} = 0$$

On Iteration t :

Compute dW, db wrt current mini-batch

$$V_{dw} = \beta_1 V_{dw} + (1 - \beta_1) dW \quad \text{"Momentum"}$$

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

$$S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) dW^2 \quad \text{"RMSprop"}$$

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

$$V_{dw}^{\text{Correct}} = V_{dw} / (1 - \beta_1^t)$$

$$V_{db}^{\text{Correct}} = V_{db} / (1 - \beta_1^t)$$

$$S_{dw}^{\text{Correct}} = S_{dw} / (1 - \beta_2^t)$$

$$S_{db}^{\text{Correct}} = S_{db} / (1 - \beta_2^t)$$

$$W \leftarrow W - \alpha \frac{V_{dw}^{\text{Correct}}}{\sqrt{S_{dw}^{\text{Correct}}}} + \epsilon$$

$$b \leftarrow b - \alpha \frac{V_{db}^{\text{Correct}}}{\sqrt{S_{db}^{\text{Correct}}}} + \epsilon$$

Hyparparameter choice :

α : needs to be tune

β_1 : 0.9 First Moment

β_2 : 0.999 Second Moment

ϵ : 10^{-8}

Bias Correction

Update

Learning Rate Decay Take smaller update steps when approaching the minimum

$$\alpha = \frac{1}{1 + \text{Decay Rate} \times \text{Epoch Number}} \alpha_0$$

1 Epoch : 1 pass over training data (One iteration of
forward prop
cost calculate
backward prop
for loop)

Exponential Decay : $\alpha = 0.95^{\frac{\text{Epoch Number}}{\alpha_0}}$

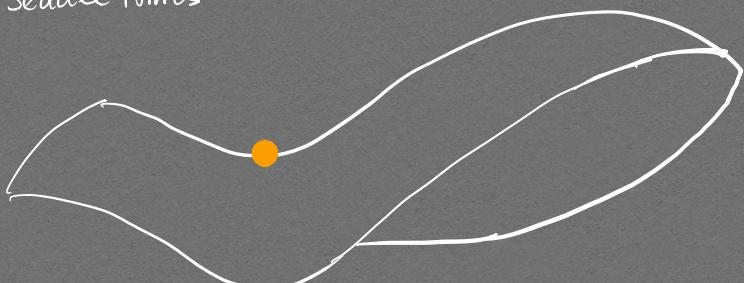
$$\alpha = \frac{k}{\sqrt{\text{Epoch Number}}} \alpha_0 \quad (k \text{ for some constant })$$

Discrete Decay

Manual Decay

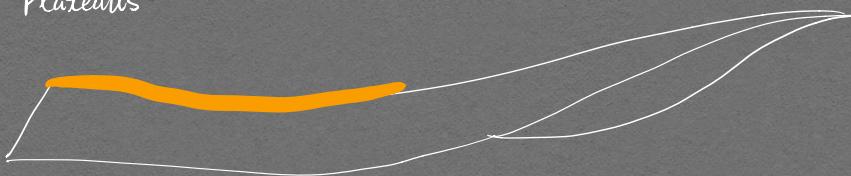
Local Optima

Saddle Points



- Unlike to get stuck in local optima in high-D space.

Plateaus



- Plateaus can make learning slow.