

Applying Machine Learning is a highly empirical process.

Big Data → Require large networks → Training is slow

So, optimization algorithms can speed up training time and increase efficiency.

Mini-batch Gradient Descent

Recall Batch Gradient Descent and vectorization on m samples.

$$X = [x^{(1)} \ x^{(2)} \dots \ x^{(m)}], \ X: (n_x, m) \quad Y = [y^{(1)} \ y^{(2)} \dots \ y^{(m)}], \ Y: (1, m)$$

It's impractical to read the entire training set into memory if m is large. E.g. $m = 5 \times 10^6$

It's just too slow to run over the entire training set for one iteration in gradient descent.

It would be better and faster that the gradient descent can start making progress before finishing the whole training set.

Thus, we can select 1000 samples as a mini-batch. $m = 5 \cdot 10^6$

$$X = [x^{(1)} \ x^{(2)} \dots x^{(1000)} \ | \ x^{(1001)} \dots x^{(2000)} \ | \ \dots \ | \ \dots x^{(m)}] \quad Y = [y^{(1)} \ y^{(2)} \dots y^{(1000)} \ | \ y^{(1001)} \dots y^{(2000)} \ | \ \dots \ | \ \dots y^{(m)}]$$

$\underbrace{x^{(1)} \dots x^{(1000)}}_{X^{1:t}} \quad \underbrace{x^{(1001)} \dots x^{(2000)}}_{X^{2:t}} \quad \dots \quad \underbrace{x^{(1000:t)}}_{X^{1000:t}}$

$\underbrace{y^{(1)} \dots y^{(1000)}}_{y^{1:t}} \quad \underbrace{y^{(1001)} \dots y^{(2000)}}_{y^{2:t}} \quad \dots \quad \underbrace{y^{(1000:t)}}_{y^{1000:t}}$

$$\text{Mini-batch } t: X^{1:t}: (n_x, 1000), \ Y^{1:t}: (1, 1000), \ t: 1, \dots, 5000$$

Now, we can process mini-batch in one iteration of gradient descent.

One Epoch (One single pass through the training set)

for $t = 1, \dots, 5000$

Forward Propagation on $X^{1:t}$

$$\begin{aligned} Z^{(1)} &= W^{(1)} X^{1:t} + b^{(1)} \\ A^{(1)} &= g^{(1)}(Z^{(1)}) \\ &\vdots \\ A^{(L)} &= g^{(L)}(Z^{(L)}) \end{aligned} \quad \left. \begin{array}{l} \text{Vectorization} \\ \text{over 1000 samples} \end{array} \right\}$$

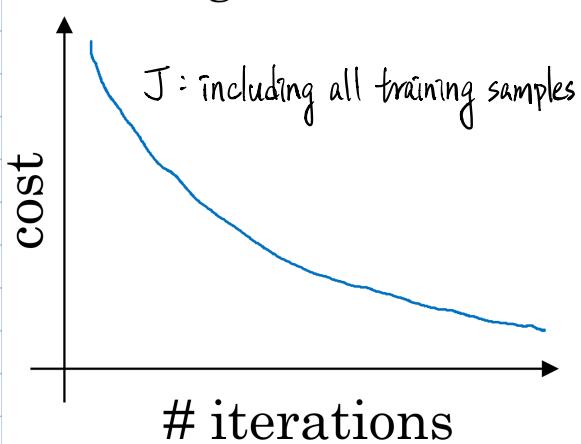
$$\text{Compute cost } J^{1:t} = \frac{1}{1000} \sum_{i=1}^{1000} (\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2 \cdot 1000} \sum_{l=1}^L \|W^{(l)}\|_F^2$$

Back Propagation to compute gradients w.r.t $J^{1:t}$ (using $X^{1:t}, Y^{1:t}$)

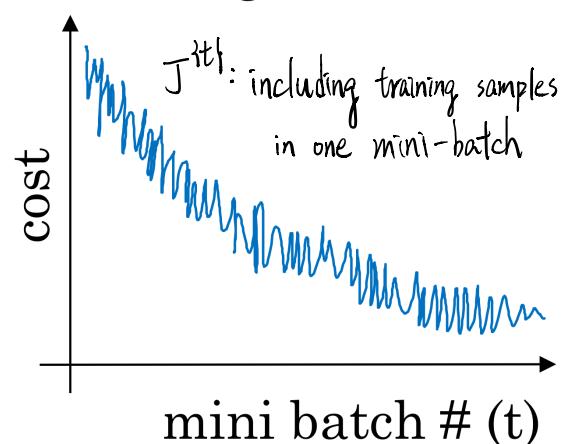
Now, with mini-batch, we can take 5000 gradient descent steps in one epoch, whereas with batch gradient descent, a single pass through the training set takes only one gradient descent step.

Understanding Mini-batch Gradient Descent

Batch gradient descent



Mini-batch gradient descent



mini-batch size = m : Batch Gradient Descent

mini-batch size = 1 : Stochastic Gradient Descent (SGD)

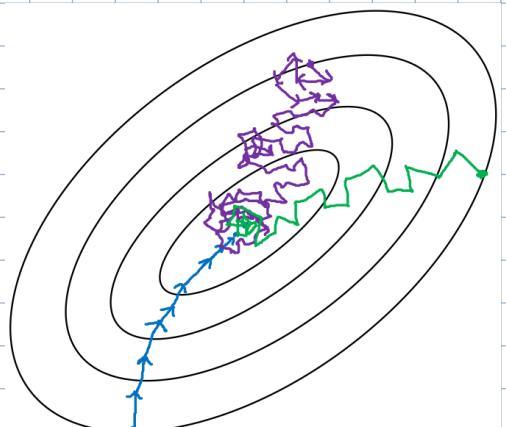
SGD: Lose speed up from vectorization
Never hit the minimum and stay there.
Wander around, sometimes the wrong direction.

Batch Gradient Descent: Too long per iteration

The Best Solution: mini-batch size In-Between (eg: 1000)

Faster learning in one gradient descent step
Vectorization

Make progress without processing the entire training set



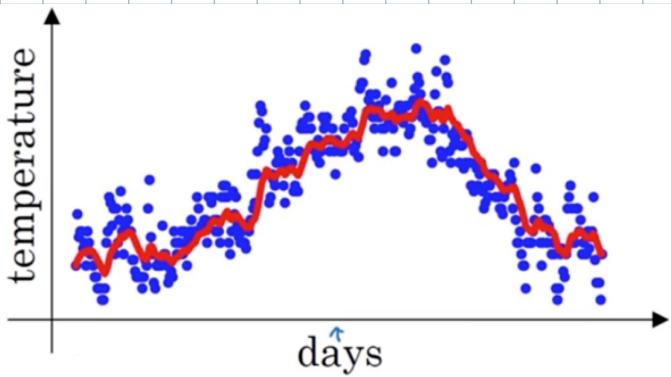
Practical Tips for choosing mini-batch size.

Small training set: Just use Batch Gradient Descent. Eg: $m \leq 1000$

Typical mini-batch size: 64, 128, 256, 512. Remember to use 2^n .

Fit in CPU/GPU memory.

Exponentially Weighted (Moving) Averages:



Equation: $V_t = \beta V_{t-1} + (1-\beta) \theta_t$ V_t : Approximately averaging over $\approx \frac{1}{1-\beta}$ days' temperature.

Eg: $\beta = 0.9 \cdot \frac{1}{1-\beta} = 10$, \therefore The temperature at day 10 is the average temperature from 1st to 10th days.

$$\beta = 0.98 \cdot \frac{1}{1-\beta} = 50$$

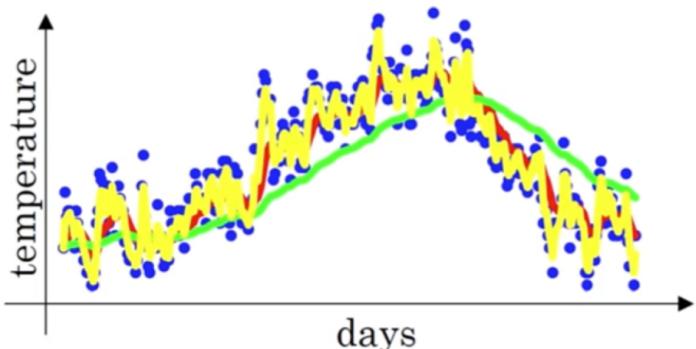
$$v_t = \beta v_{t-1} + (1 - \beta) \theta_t \quad v_0 = 0$$

$$\beta = 0.5 \cdot \frac{1}{1-\beta} = 2$$

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

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

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



If $\beta = 0.9$

$$V_{100} = 0.1\theta_{100} + 0.1 \cdot 0.9 \cdot \theta_{99} + 0.1 \cdot 0.9^2 \theta_{98} + \dots$$

$$0.9^{10} \approx 0.35 \approx \frac{1}{e}$$

$$0.98^{50} \approx \frac{1}{e}$$

$$\epsilon = 1 - \beta$$

$$(1 - \epsilon)^{\frac{1}{\epsilon}} = \frac{1}{e}$$

When the percentage drops to $\frac{1}{e}$, its effect becomes minor.

Implementing exponentially weighted averages

$$V_\theta = 0$$

Repeat }

Get next θ_t

$$V_\theta = \beta V_\theta + (1 - \beta) \theta_t$$

} It's fast and memory efficient.

Now I fully understand why exponentially weighted average is better than using window-based average. In window-based average, it needs an array to store values in a given window when calculating the average of each element, which is slow and memory consuming.

Bias Correction in Exponentially Weighted Averages:

If using the aforementioned implementation to implement Exponentially Weighted Averages when $\beta = 0.98$, we won't get the green curve, but the purple curve as shown in the right image. What? Why? I'm fooled by Andrew.

If looking closer: $V_t = 0.98 V_{t-1} + 0.02 \theta_t$, because the initial value is 0, so the starting few values are heavily influenced by the initial value $V_0 = 0$.

$$\text{The fix: } V_t = \frac{\beta V_{t-1} + (1 - \beta) \theta_t}{1 - \beta}. \quad \text{If } t=1, 1 - \beta^t = 0.02$$

$$\quad \quad \quad \text{If } t=2, 1 - \beta^t = 0.0396$$

When t becomes larger, β^t approaches zero and has lesser influence.

Hence, we can get the green curve after applying Bias Correction.

Exponentially weighted moving averages:

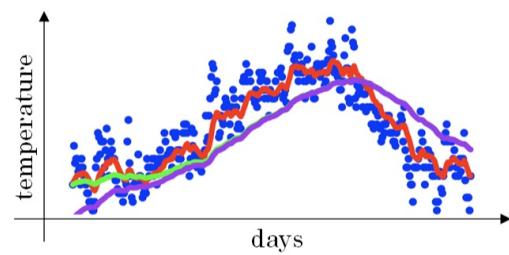
$$V_0 = 0$$

$$V_1 = 0.9 V_0 + 0.1 \theta_1$$

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

$$\vdots$$

$$V_{365} = 0.9 V_{364} + 0.1 \theta_{365}$$

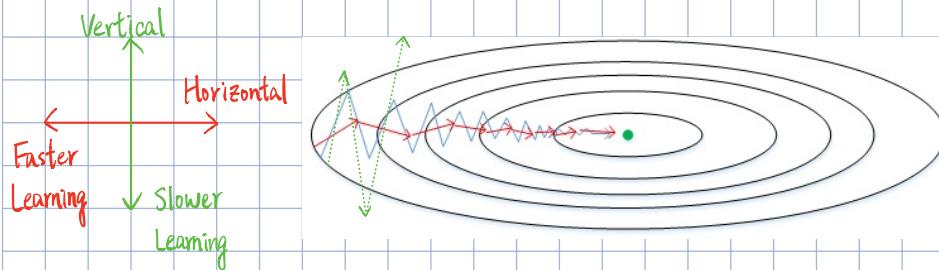


$$\text{Eg: } \theta_1 = 30, \theta_2 = 15, \beta = 0.5, V_0 = 0.$$

$$V_1 = 15$$

$$V_2 = 15, V_2^{\text{corrected}} = \frac{0.5 \times 15 + 0.5 \times 15}{1 - 0.5^2} = 20$$

Gradient Descent with Momentum: It computes an exponentially weighted average of gradients, and use that gradient to update the weights.



In order to prevent oscillation, overshooting and potential divergence, the learning rate can't be too large.

Quiz: Given the blue and the red line on the left, which one has bigger β ? The Red line.

On the vertical axis, we want the learning to be a bit slower.
On the horizontal axis, we want the learning to be faster.

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

Compute dw and db on the current mini-batch

$$V_{dw} = \beta V_{dw} + (1-\beta) dw \quad \text{"Recall } V_t = \beta V_{t-1} + (1-\beta) \theta_t \text{"}$$

friction
velocity

$$V_{db} = \beta V_{db} + (1-\beta) db \quad \text{acceleration} \quad \text{Analogy for mechanical engineer}$$

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

Default: $\beta = 0.9$,

Average out last 10 gradients.

Remember to use bias correction $\frac{V_t}{1-\beta^t}$

α : the Learning Rate

Like the temperature is smoothed out in exponentially weighted average, the gradient is also smoothed out.

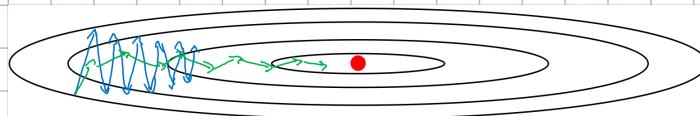
Eg:

The vertical direction got cancelled out, but the horizontal direction keeps being strong.

Sometimes, some literatures use $V_{dw} = \beta V_{dw} + (1-\beta) dw = \beta V_{dw} + dw$ as the formula.

RMSprop : Root Mean Squared Propagation

Assume:



Goal: We want to slow down the learning vertically and speed up (at least not slowing down) horizontally.

On iteration t :

Try to explain the equation. If not, rewatch video.

Compute dw, db on the current mini-batch

$$S_{dw} = \beta S_{dw} + (1-\beta) dw^2 \quad \text{element-wise}$$

$$S_{db} = \beta S_{db} + (1-\beta) db^2 \quad \beta_2 \text{ is used to differ itself from } \beta \text{ in momentum above}$$

$$W = W - \alpha \frac{dw}{\sqrt{S_{dw} + \epsilon}}, \quad b = b - \alpha \frac{db}{\sqrt{S_{db} + \epsilon}} \quad \epsilon \text{ is a small number (e.g. } 10^{-8}) \text{ to ensure numerical stability. } (S_{dw} = 0)$$

Adam Optimization Algorithm

Research AdamW

Adam = Momentum + RMSprop

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

Compute dw , db on the current mini-batch

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

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

$$V_{dw}^{\text{corrected}} = \frac{V_{dw}}{1 - \beta_1^t}, \quad S_{dw}^{\text{corrected}} = \frac{S_{dw}}{1 - \beta_2^t}$$

$$V_{db}^{\text{corrected}} = \frac{V_{db}}{1 - \beta_1^t}, \quad S_{db}^{\text{corrected}} = \frac{S_{db}}{1 - \beta_2^t}$$

$$W = W - \alpha \frac{V_{dw}^{\text{corrected}}}{\sqrt{S_{dw}^{\text{corrected}}} + \epsilon} \quad b = b - \alpha \frac{V_{db}^{\text{corrected}}}{\sqrt{S_{db}^{\text{corrected}}} + \epsilon}$$

Hyperparameters choice:

α : needs to be tuned (Learning Rate)

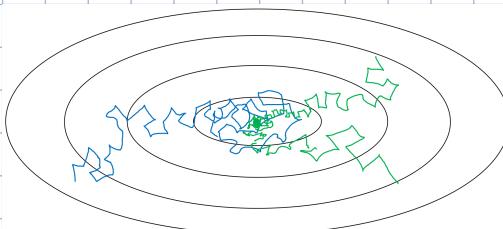
β_1 : 0.9 rarely tuned

β_2 : 0.999 rarely tuned

ϵ : 10^{-8} . never tuned

Learning Rate Decay

Blue line represents the learning rate is constant, so it keeps wandering around and never converges.

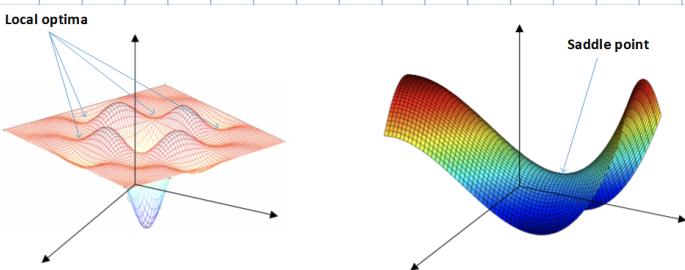


Green line represents the learning rate is decaying, so it can reach the minimum.

One Epoch: One pass through the entire training set.

$$\alpha = \frac{1}{1 + DR \cdot EN} \alpha_0 \quad \alpha_0: \text{initial learning rate} \quad DR: \text{Decay Rate} \quad EN: \text{Epoch Number} \quad \alpha = 0.95^{EN} \alpha_0, \quad \alpha = \frac{k}{\sqrt{EN}} \alpha_0, \quad \alpha = \frac{k}{\sqrt{t}} \alpha_0, \quad t: \text{mini-batch \#}$$

The Problem of Local OPTima



For a long time, people think the local optima appears quite often and will become a problem in the training phase if the input feature is high dimensional vector. But it turns out it's highly unlikely the case. In fact, it will most likely be a saddle point. Why?

Given a feature space with 20000 dimensions, each element will be either convex \cup or concave \cap .

If it's a local optima, then the probability is $\frac{1}{2^{20000}}$, which is almost impossible. Thus, it mostly will be a saddle point.