

mini-batch gradient descent

- Good optimization algorithms are being used in the deep learning era as the models require large datasets to be trained

- batch vs. mini-batch

↳ Vectorization allows us to efficiently compute on m training examples without an explicit for loop.

$$\mathcal{X} = [x^{(1)} \ x^{(2)} \ x^{(3)} \dots \ x^{(m)}]$$

$$\hookrightarrow (n \times m)$$

$$\mathcal{Y} = [y^{(1)} \ y^{(2)} \ y^{(3)} \dots \ y^{(m)}]$$

$$\hookrightarrow (1, m)$$

* What if m was really large? Vectorization can still lead to really slow training (as we have to process the vectors each time we are going to take a step of gradient descent)

↳ We can split our training set into smaller mini-batches. Instead of vectorizing 5 million examples, each mini-batch can have 100 examples

$$\mathcal{X} = [\underbrace{x^{(1)} \ x^{(2)} \dots \ x^{(100)}}_{100 \text{ training examples}} \ | \ \underbrace{x^{(101)} \dots \ x^{(200)}}_{x^{1-3} (n \times 100)} \ | \ \dots \ | \ \underbrace{x^{(5000)}}_{x^{1-3} (n \times 100)}]$$

↳ Similarly \mathcal{Y} would be split into 5000 mini-batches $(1, 100)$ per mini-batch

* $x^{(i)}$: i-th training example, $x^{(l)}$ for layer #, $x^{\{i\}}$ mini-batch #

- Mini-batch gradient descent

for $t = 1, \dots, 5000$

* inside the for loop we implement one step of gradient descent using x^{q+3} and y^{q+3} (as if we were to train our network on a set of 1000 examples)

forward pass on x^{q+3}

$$z^{(l)} = w^{(l)} x^{q+3} + b^{(l)}$$

$$A^{(l)} = g^{(l)}(z^{(l)})$$

\vdots

$$A^{(L)} = g^{(L)}(z^{(L)})$$

vectorized implementation of 1000 examples

$$\text{Compute cost } J = \frac{1}{1000} \sum_{i=1}^I \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

$$+ \frac{\lambda}{2 \cdot 1000} \sum_{i=1}^I \|w^{(i)}\|_F^2$$

for bias term

from mini-batch x^{q+3}, y^{q+3}

Backprop to compute gradient J^{q+3} (using (x^{q+3}, y^{q+3}))

$$w^{(q+3)} := w^{(q+3)} - \alpha \nabla J^{(q+3)}$$

$$b^{(q+3)} := b^{(q+3)} - \alpha \nabla b^{(q+3)}$$

\exists

* this is called doing one pass through the training set using mini batch gradient descent

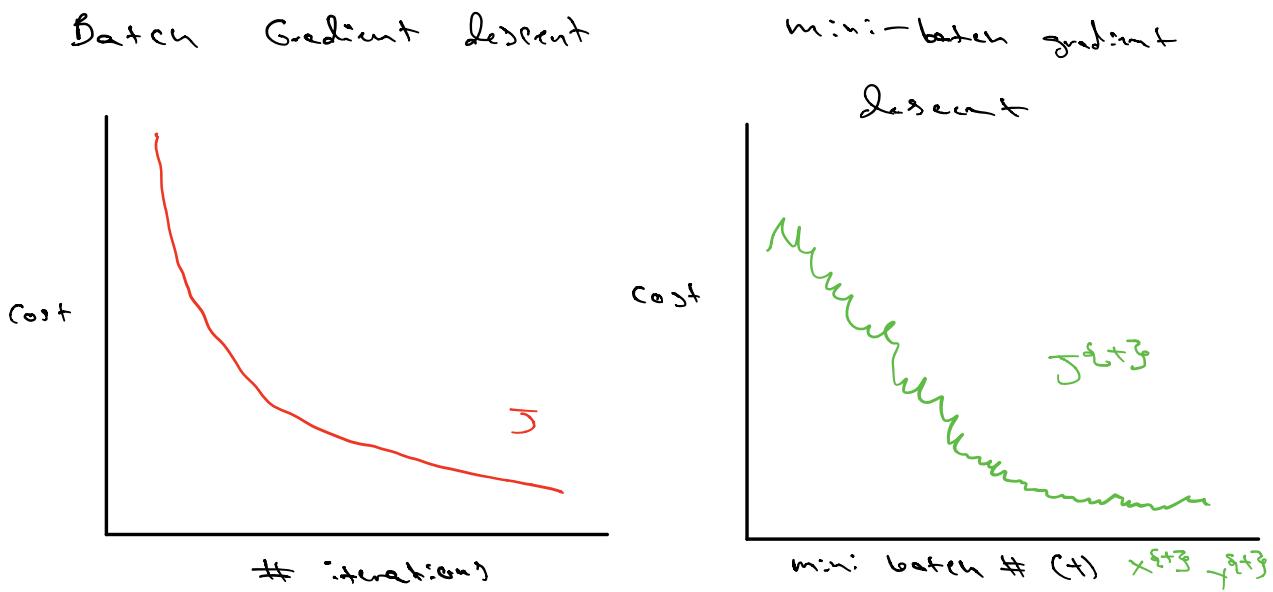
One epoch of training (single pass through training set)

With batch 1 pass means that we are taking one step of gradient descent, while 1000 in mini batch means 5000 steps of gradient descent.

* we can repeat this n # of times to do multiple passes on our training set.

Understanding mini-batch gradient descent

- With mini-batch processing you start taking training steps before you finish processing the training set.



- With batch gradient descent, we'd expect the cost function to go down on every single iteration (if it goes up even in one iteration maybe the learning rate is too large)

- in mini-batch it is as if in every iteration we are training on a different training set, hence, we are more likely to get fluctuations with the general trend toward a lower cost.

(the noising comes from the differences between mini-batches)

- Choosing mini-batch size

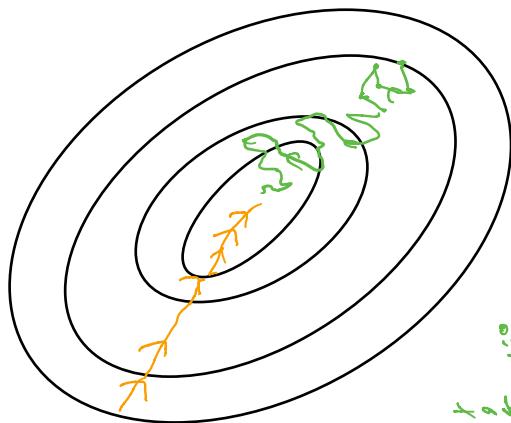
↳ if the mini-batch size = n: we would have batch gradient

$$\text{descent } (\mathbf{x}^{i,3}, \mathbf{y}^{i,3}) = (\mathbf{x}, \mathbf{y})$$

↳ if mini-batch size = 1 : Stochastic gradient descent

↳ every example is its own mini-batch

$$(\mathbf{x}^{i,3}, \mathbf{y}^{i,3}) = (\mathbf{x}^{ii}, \mathbf{y}^{ii}) \dots$$



- batch gradient descent will be able to take relatively low noise (large) steps
- stochastic gradient descent

is noisy, but will on average take you in a good direction (to

minimize the cost function). Stochastic gradient never converges, it oscillates in the region near the minimum.

- In practice, the mini-batch size used is somewhere between 1 and m.

↳ batch gradient descent (mini-batch=m) matches up to process a whole training set on every iteration taking too long for iteration (with a large training set)

↳ with stochastic gradient descent we lose all the speed up from vectorization because we are processing each example at a time. This makes the way we process training examples inefficient.

↳ with a medium size mini-batch we will get on average the fastest learning (through vectorization and by making progress without the reiteration

required to process the whole train set)

- if we have a small train set: use batch gradient descent ($m \leq 2000$)
- typical minibatch sizes: 64, 128, 256, 512
 - ↳ typically on the lower or \approx because code sometimes runs faster.
 - ↳ we also need to make sure that the minibatch fits in GPU memory

Exponentially Weighted Moving Average

- way to smooth filter data in order to get more smoothness in trends. $v_t = 0.9v_{t-1} + 0.1\Theta$.

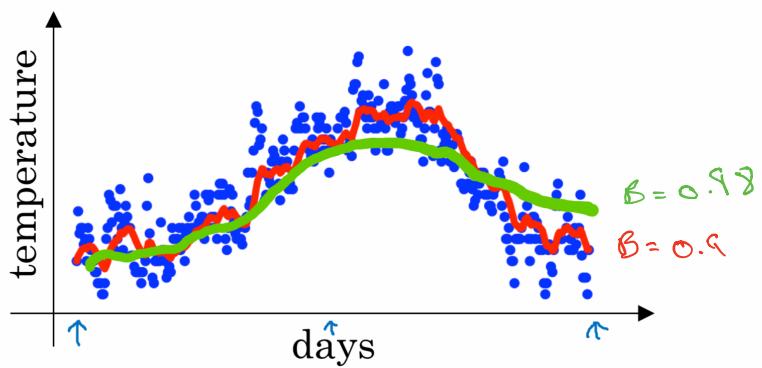
v_t value at present day
 v_{t-1} value at previous day

* generalized: $v_t = \beta \cdot v_{t-1} + (1-\beta)\Theta +$

↳ v_t can be thought as an averaging over $\frac{1}{1-\beta}$ past days

↳ with the example above, we are averaging over the last 10 days (reducing the effect of short variations)
↳ if we make β larger, the impact of any recent value measurement is nearly outweighed by the previous measurements (and hence the curve smoothes)

* taking larger betas has the effect of shifting the curve further to the right as we are now averaging over larger windows (added latency to the formula)



Understanding Exponentially weighted averages

- Implementing exponentially weighted averages

$$v_0 = 0$$

$$v_t := \beta \cdot v_{t-1} + (1 - \beta) \theta_t \quad \left. \begin{array}{l} \text{for loop over} \\ \# \text{ of } \theta_t \end{array} \right\}$$

$$v_0 := \beta \cdot v_0 + (1 - \beta) \theta_0$$

:

* this implementation takes little memory as only one value for v_0 is kept at a time

This is a far more computationally efficient way

to compute averages than keeping values of the last n 's in memory at a time.

Bias Correction

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

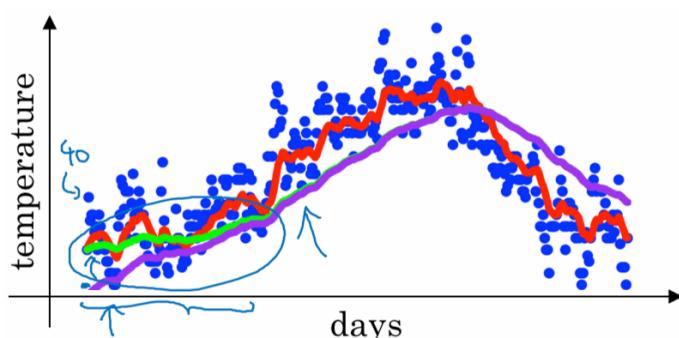
$$v_0 = 0$$

$$v_1 = 0.98 v_0 + 0.02 \theta_1$$

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

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

$$= 0.0196 \theta_1 + 0.02 \theta_2$$



improve from purple one to the green one

- for low sample, the averaging starts floored because there's not enough values accumulated in the V_t term
 - ↳ to fix this, we must employ bias correction.

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

denominator term

So for example for ($t=2, \beta = 0.98$) $\rightarrow \frac{1 - (0.98)^2}{0.0396} =$

hence,

$$\frac{0.0198\theta_1 + 0.02\theta_2}{0.0396}$$

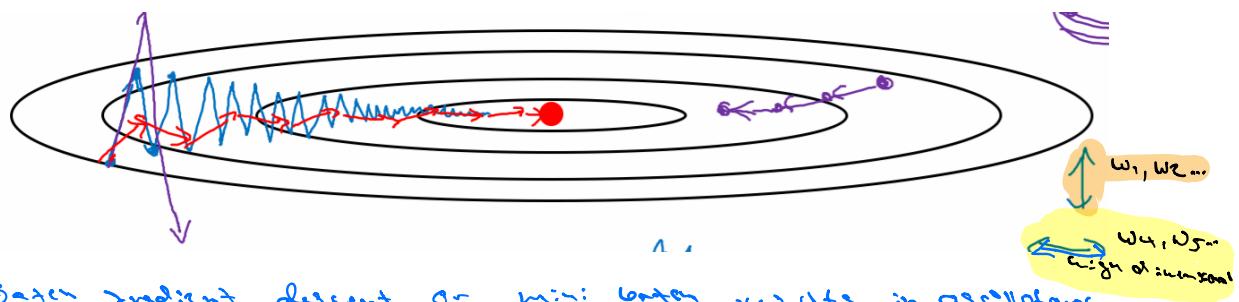
numerator term = denominator term

↳ this situation takes β into account, removing the bias effect.

↳ as β grows large, $\beta^t \approx 0$, hence $1 - \beta^t \approx 1$
playing no role later on in the calculation.

Gradient descent with momentum

- almost always works faster than the standard gradient descent algorithm. Combining an exponentially weighted average of the gradients and the use that gradient to update the weights instead.



- batch gradient descent or mini batch results in oscillatory behaviour as we move across the optimization contours

These up and down oscillations slows down gradient descent and inhibits us from using a larger learning rate.
 With a larger learning rate we might end up overshooting (so we are forced to use a learning rate which is not too large).

We want fast learning on the horizontal direction of the contour and slow on the vertical.

- By using momentum, on each iteration (t):

compute Δw , Δb on current mini-batch.

$$v \Delta w = \beta v \Delta w + (1 - \beta) \Delta w \quad \begin{matrix} \text{similar to previous moving} \\ \text{average calculation} \end{matrix}$$

$$v \Delta b = \beta v \Delta b + (1 - \beta) \Delta b$$

$$\begin{aligned} w &:= w - \alpha \cdot v \Delta w \\ b &:= b - \alpha \cdot v \Delta b \end{aligned} \quad \begin{matrix} \text{this smooths out the steps of} \\ \text{gradient descent by damping the} \\ \text{vertical oscillations.} \\ (\text{oscillations are damped}) \end{matrix}$$

* hence we end up traversing down the contour (optimization plane) in a more effective way. (red on diagram)

* the most common β value is 0.9 (so that we average over the last 10 gradients)

* Bias correction is not really implemented as the derivative term is "warmed up" after a few iterations (would not impacting the overall performance).

RMSprop (Root mean squared prop)

On iteration t :

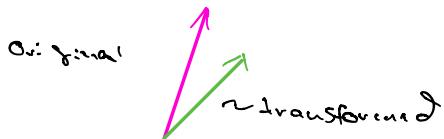
Compute Δw , Δb on current mini-batch

$$S\Delta w = \beta_2 \cdot S\Delta w + (1-\beta) \cdot \Delta w^2 \quad \text{moment with Squaring Operation}$$

$$S\Delta b = \beta_2 \cdot S\Delta b + (1-\beta) \cdot \Delta b^2 \quad \text{~large}$$

$$\omega := \omega - \alpha \cdot \frac{\Delta w}{\sqrt{S\Delta w} + \epsilon} \quad \text{~we want to update the rate by dividing a relatively small #}$$

$$b := b - \alpha \cdot \frac{\Delta b}{\sqrt{S\Delta b} + \epsilon} \quad \text{~we want to limit the update to } b \text{ by dividing by a large #.}$$



* if we go back to the contour diagram, we see that similar to momentum the gradient doesn't oscillate as much and travels faster towards the minimum so this methodology enables an increase in α .

* $\epsilon = 10^{-8}$ is added so that we never divide the derivatives by zero or extremely small #'s (numerical stability)

Adam optimization

• basically taking momentum and RMSprop and putting them together.

• programmatic implementations:

$$V\Delta w = 0, \quad S\Delta w = 0, \quad V\Delta b = 0, \quad S\Delta b = 0$$

On iteration t :

compute Δw , Δb using current mini-batch

$$V\Delta w = \beta_1 \cdot V\Delta w + (1-\beta_1) \cdot \Delta w \quad \rightarrow$$

$$v_{dw} = \beta_1 v_{dw} + (1 - \beta_1) dw \quad \left. \begin{array}{l} \text{momentum accumulation} \\ \text{seen} \end{array} \right]$$

$$s_{dw} = \beta_2 s_{dw} + (1 - \beta_2) dw^2 \quad \left. \begin{array}{l} \text{RMSProp update} \end{array} \right]$$

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

* in the typical implementation of adam we add bias

Correction

$$\left. \begin{array}{l} v_{dw}^{\text{corrected}} = v_{dw} / (1 - \beta_1^+) \\ v_{db}^{\text{corrected}} = v_{db} / (1 - \beta_1^+) \\ s_{dw}^{\text{corrected}} = s_{dw} / (1 - \beta_2^+) \\ s_{db}^{\text{corrected}} = s_{db} / (1 - \beta_2^+) \\ w := w - \alpha \frac{v_{dw}^{\text{corrected}}}{\sqrt{s_{dw}^{\text{corrected}}} + \epsilon} \\ b := b - \alpha \frac{v_{db}^{\text{corrected}}}{\sqrt{s_{db}^{\text{corrected}}} + \epsilon} \end{array} \right\} \text{update terms}$$

- this algorithm combines the effect of momentum with gradient descent with RMSprop to create a effective learning algorithm which works on a wide variety of architectures.

- hyperparameter choices

↳ α : need to be tuned

↳ $\beta_1 = 0.9$ (v_{dw})

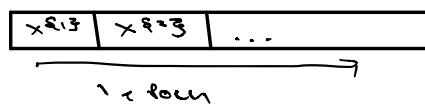
↳ $\beta_2 = 0.999$ (s_{dw}^2)

↳ $\epsilon = 10^{-8}$

• adam: adaptive moment estimation

Learning rate decay

- slowly reducing the learning rate over time helps speed up the learning algorithm (learning rate decay).
 - ↳ we oscillate in a tighter region near the minimum than with an otherwise large learning rate.
- applying learning rate decay:
 - 1 epoch = 1 pass through the data

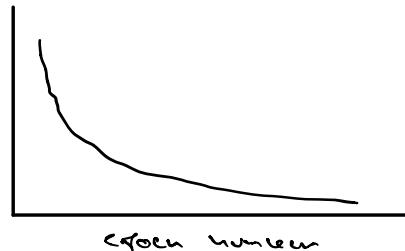


We want to set the learning rate to,

$$\alpha = \frac{1}{1 + \text{decay_rate} \cdot \text{epoch_number}} \quad \text{do } \sim \text{initial decay rate}$$

* note that the decay-rate might be another hyperparameter that we would have to tune.

α



- Other learning rate decay methods

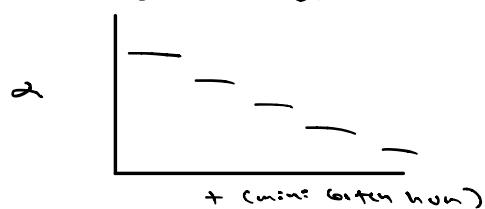
↳ exponential decay

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

$$\alpha = \frac{k}{\sqrt{\text{epoch_num}}} \cdot \alpha_0$$

some #

↳ discrete decreases



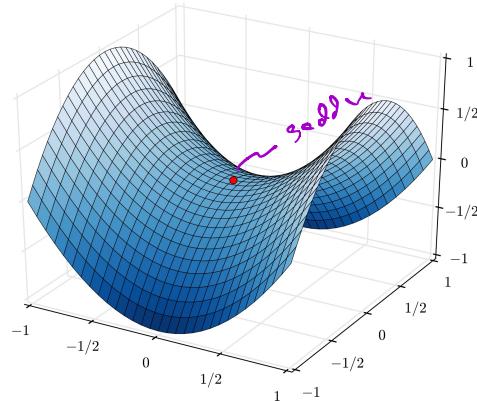
* learning rate decay can help training, but it's often a little bit down the list when it comes to things that we can do to improve the model

↳ try choosing an adequate learning rate first

The Problem of Local Optima

- in the early days of DL there was a fear of algorithms getting trapped in local optima (this notion has changed)
- in high dimensional space, a local optima requires all features to have 0 gradient at one specific point (which is unlikely when we are operating in 20000 dimensions) (concave or convex shape)

↳ we are much more likely to run into saddle points.



- the -est problems arise from **flatus** regions where the derivative is close to zero, slowing down learning)
↳ RMSprop, adam, momentum can all help in this type of scenario.