



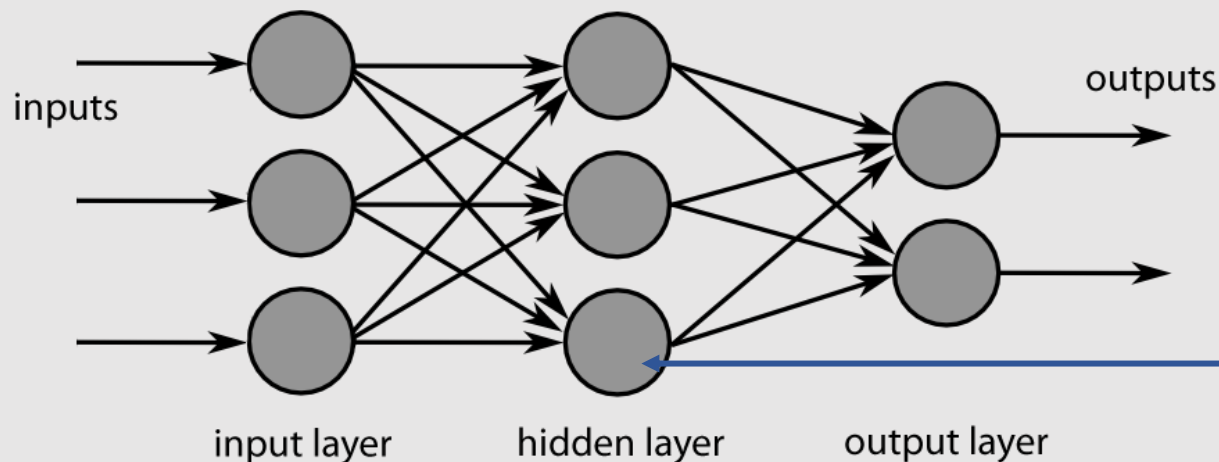
Optimization Methods - ADAM

Recall – neural network and gradient descent

Neural Network

Algorithm inspired by how the brain works

The role of neural network is to predict \hat{y} (We need to give the input x and output y)



$$g(\omega x + b)$$

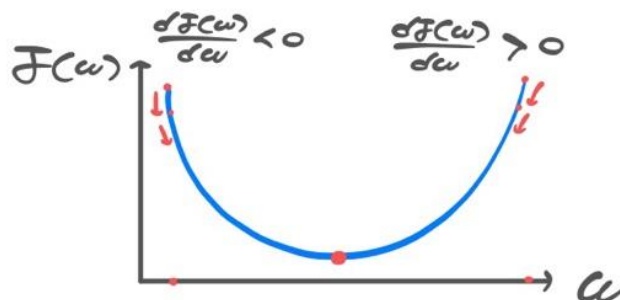
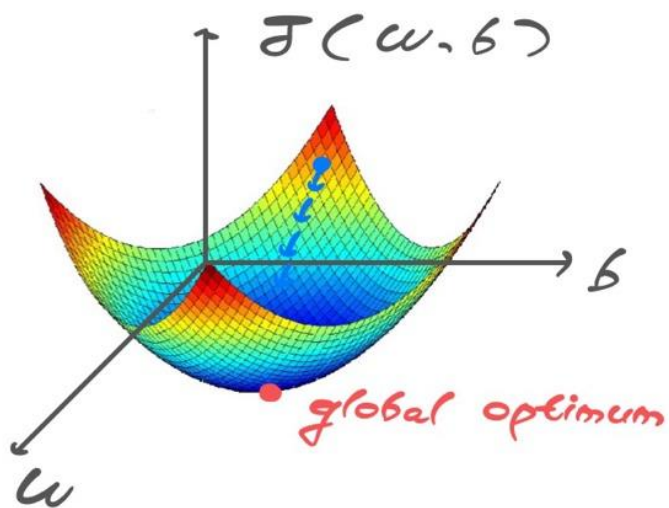
$$\begin{aligned} \mathcal{F}(\omega, b) &= \frac{1}{n} \sum_{i=1}^n \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) \\ &= -\frac{1}{n} \sum_{i=1}^n (y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})) \end{aligned}$$

**Cost
Function**

Recall – neural network and gradient descent

Gradient Descent

Want to find W, b that minimize $J(W, b)$



$$\begin{aligned} w &:= w - \alpha \frac{\partial J(w, b)}{\partial w} \\ b &:= b - \alpha \frac{\partial J(w, b)}{\partial b} \end{aligned}$$

update

α : learning rate

Mini-batch gradient descent

Mini-batch gradient descent

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(1000)} & x^{(1001)} & \dots & x^{(5000)} \end{bmatrix}$$

$(n \times m)$

$X^{(1)}$ "mini-batch"

$X^{(5000)}$

i.e. mini-batch k : $X^{(k)}$, $y^{(k)}$

$(n \times 1000)$ (1×1000)

What if $m = 5,000,000 \rightarrow$ mini-batches of 1000 epoch

Mini-batch gradient descent

Mini-batch gradient descent

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

Forward Prop on $X^{(t)}$

$$\begin{array}{|l} z^{(1)} = w^{(1)} x^{(t)} + b^{(1)} \\ a^{(1)} = g^{(1)}(z^{(1)}) \\ \vdots \\ a^{(L)} = g^{(L)}(z^{(L)}) \end{array}$$

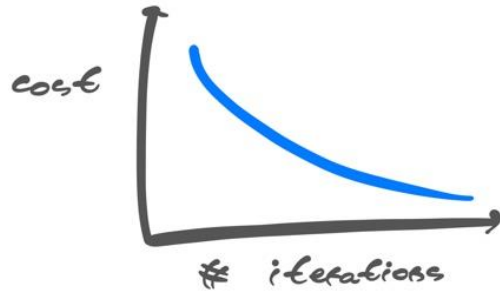
Compute cost $J = \frac{1}{5000} \sum_{i=1}^L L(\underbrace{\hat{y}^{(i)}, y^{(i)}}_{\text{from } X^{(t)}, y^{(t)}})$

Backprop to compute gradient wrt $J^{(t)}$

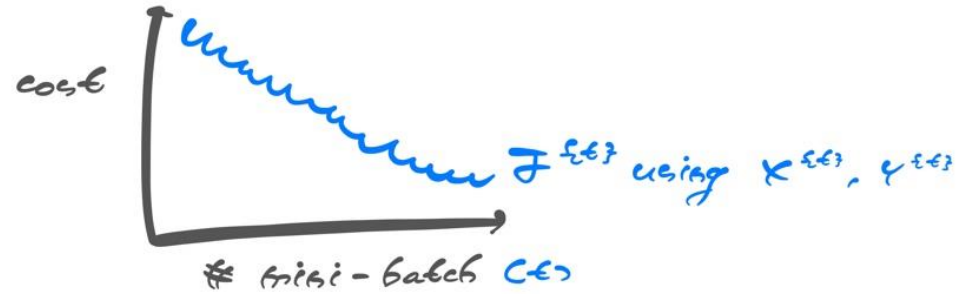
$$\begin{array}{|l} w^{(1)} := w^{(1)} - \alpha \delta w^{(1)} \\ b^{(1)} := b^{(1)} - \alpha \delta b^{(1)} \end{array}$$

1 epoch
: a single pass through
training examples

Mini-batch gradient descent



Batch gradient descent



Mini-batch gradient descent

Exponentially weighted averages

Temperature in KOREA

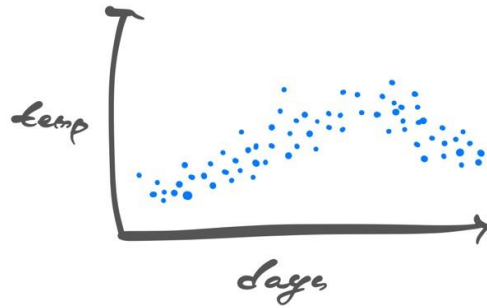
$$\Theta_1 = 4^\circ\text{C}$$

$$\Theta_2 = 4.95^\circ\text{C}$$

\vdots

$$\Theta_{180} = 15^\circ$$

$$\Theta_{181} = 14.6^\circ\text{C}$$



Noisy \rightarrow We want to compute trend

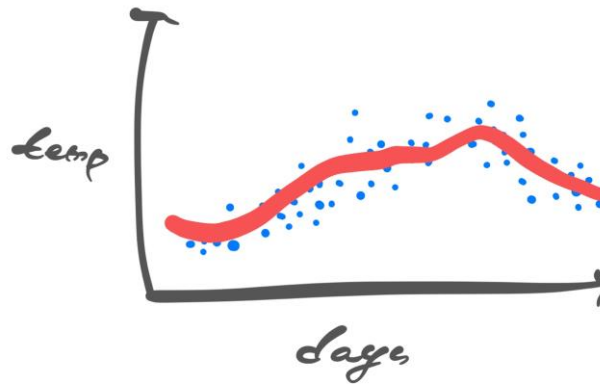
$$V_0 = 0 \quad \text{initialize}$$

$$V_1 = 0.9V_0 + 0.1\Theta_1$$

$$V_2 = 0.9V_1 + 0.1\Theta_2 \quad \text{weighted average}$$

\vdots

$$V_t = 0.9V_{t-1} + 0.1\Theta_t$$



Exponentially weighted averages

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

$$\approx \frac{1}{1 - \beta} \text{ day's temperature}$$

i.e. if $\beta = 0.9$, $V_t \approx$ previous 10 days temp

if $\beta = 0.98$, $V_t \approx$ previous 50 days temp

if $\beta = 0.5$, $V_t \approx$ previous 2 days temp

→ β : hyperparameter

Bias correction in exponentially weighted averages

if $\beta = 0.98$, $V_0 = 0$,

$$V_1 = 0.02 \Theta_1$$

$$V_2 = 0.0196 \Theta_1 + 0.02 \Theta_2$$

if $V_2 < \Theta_1, \Theta_2$ then Bad Estimate

replace V_t with $\frac{V_t}{1 - \beta^t}$

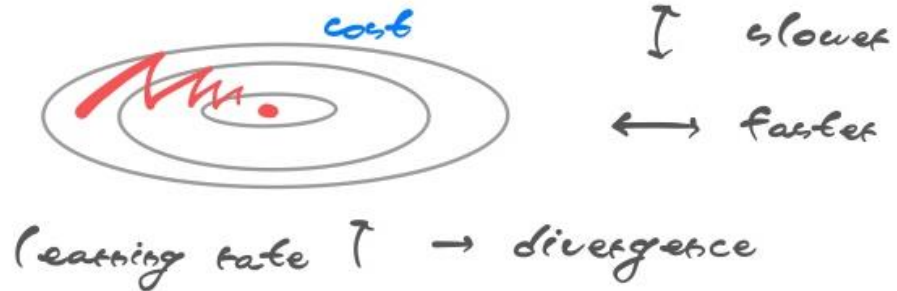
ex $t=2$: $1 - \beta^t = 1 - 0.98^2 \approx 0.0396$

$$\text{then } \frac{V_2}{1 - \beta^2} = \frac{0.0196 \Theta_1 + 0.02 \Theta_2}{0.0396}$$

Bias
Correction

$$\beta^t \rightarrow 0 \text{ as } t \rightarrow \infty$$

Momentum



On iteration t :

Compute dw, db on current mini-batch

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

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

$$w = w - \alpha V_{dw} ; b = b - \alpha V_{db}$$

● : momentum term \rightarrow velocity

● : derivative term \rightarrow acceleration

$\rightarrow \alpha, \beta$: hyperparameter

RMSprop



On iteration t :

Compute du, db on current mini-batch

$$s_{du} = \beta s_{du} + (1 - \beta) du^2$$

$$s_{db} = \beta s_{db} + (1 - \beta) db^2$$

$$u = u - \alpha \frac{du}{\sqrt{s_{du} + \epsilon}}; \quad b = b - \alpha \frac{db}{\sqrt{s_{db} + \epsilon}}$$

→ α, β : hyperparameter

ADAM

ADAM (Adaptive moment estimation)

- Algorithm for first-order gradient-based optimization of stochastic objective functions based on adaptive estimates of lower-order moments.
- Straightforward to implement, computationally efficient, little memory requirements, invariant to diagonal rescaling of the gradients, well suited for problems that are large in terms of data or params
- Appropriate for non-stationary objectives and problems with very noisy and sparse gradients
- Hyper params have intuitive interpretations and typically require little tuning.

ADAM computes individual adaptive learning rates for different parameters from estimates of first and second moments of the gradients

- > designed to combine the advantages of AdaGrad and RMSProp
- > **Advantages** : magnitudes of param updates are invariant to rescaling of the gradient, its step sizes are approximately bounded by the step size hyperparameters

ADAM

Algorithm 1: *Adam*, our proposed algorithm for stochastic optimization. See section 2 for details, and for a slightly more efficient (but less clear) order of computation. g_t^2 indicates the elementwise square $g_t \odot g_t$. Good default settings for the tested machine learning problems are $\alpha = 0.001$, $\beta_1 = 0.9$, $\beta_2 = 0.999$ and $\epsilon = 10^{-8}$. All operations on vectors are element-wise. With β_1^t and β_2^t we denote β_1 and β_2 to the power t .

Require: α : Stepsize

Require: $\beta_1, \beta_2 \in [0, 1)$: Exponential decay rates for the moment estimates

Require: $f(\theta)$: Stochastic objective function with parameters θ

Require: θ_0 : Initial parameter vector

$m_0 \leftarrow 0$ (Initialize 1st moment vector)

$v_0 \leftarrow 0$ (Initialize 2nd moment vector)

$t \leftarrow 0$ (Initialize timestep)

while θ_t not converged **do**

$t \leftarrow t + 1$

$g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$ (Get gradients w.r.t. stochastic objective at timestep t)

$m_t \leftarrow \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t$ (Update biased first moment estimate)

$v_t \leftarrow \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2$ (Update biased second raw moment estimate)

$\hat{m}_t \leftarrow m_t / (1 - \beta_1^t)$ (Compute bias-corrected first moment estimate)

$\hat{v}_t \leftarrow v_t / (1 - \beta_2^t)$ (Compute bias-corrected second raw moment estimate)

$\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \hat{m}_t / (\sqrt{\hat{v}_t} + \epsilon)$ (Update parameters)

end while

return θ_t (Resulting parameters)

ADAM

The algorithm updates exponential moving averages of the gradient (m) and The squared gradient (v) where the hyper-params β_1 , β_2 control the exponential decay rates of these moving averages.

These moving averages are initialized as 0, leading to moment estimates that are biased towards zero, especially during the initial timesteps, and especially when the decay rates are small (i.e β are close to 1)

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$$

$$\hat{m}_t = m_t / (1 - \beta_1^t)$$

$$\hat{v}_t = v_t / (1 - \beta_2^t)$$

References

- [1] Kingma, Diederik P., and Jimmy Ba. "Adam: A method for stochastic optimization." *arXiv preprint arXiv:1412.6980* (2014).
- [2] Improving Deep Neural Networks: Hyperparameter Tuning, Regularization and Optimization, deeplearning.ai