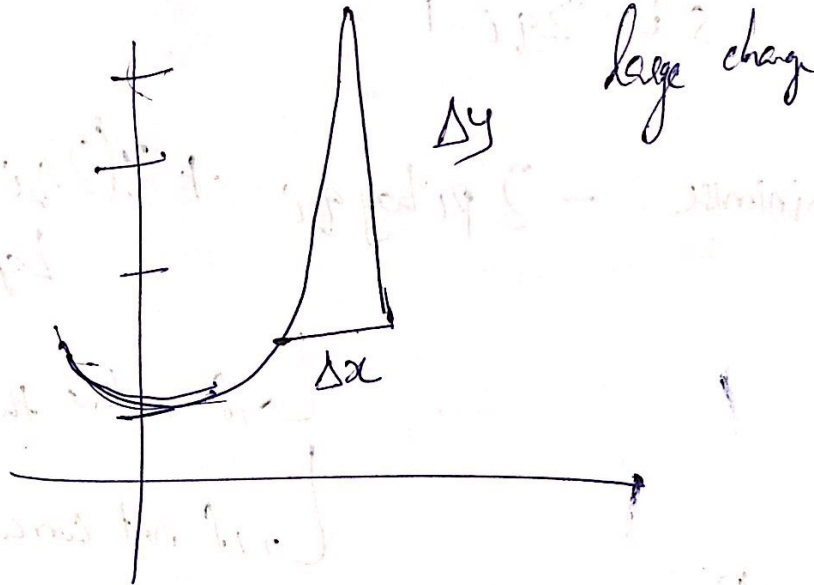


## Lecture 5:

### Learning Parameters, Gradient Descent

Part  
5.1 & 5.2

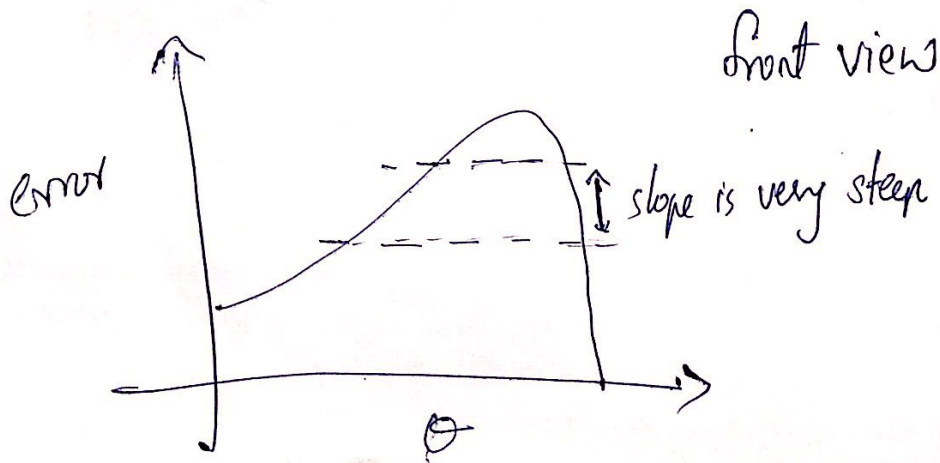
$$f(x) = x^2 + 1$$



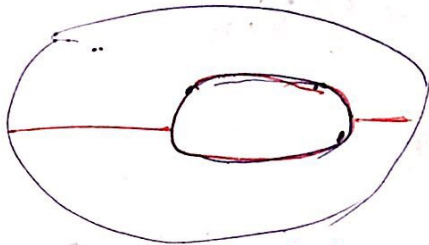
Curve is steep if the gradient  $\left(\frac{\Delta y_1}{\Delta x_1}\right)$  is large,  
Curve is gentle the gradient  $\left(\frac{\Delta y_2}{\Delta x_2}\right)$  is small.

5.3

### Contour Maps



Top view



A small distance b/w the contours indicates a steep slope along that direction

A large distance b/w the contours indicates a gentle slope along that direction

#### S.4 Momentum Based Gradient Descent

Lot of time to navigate regions having gentle slope.

more people pointing in one direction, keep moving in the direction with a bigger step.

Update rule  $\xrightarrow{\text{history}}$  current gradient

$$\text{update}_t = \gamma \cdot \text{update}_{t-1} + \eta \nabla L_t$$

$$w_{t+1} = w_t - \text{update}_t$$

$$w_{t+1} = w_t - \gamma \cdot \text{update}_{t-1} - \eta \nabla L_t \quad \left\{ \begin{array}{l} \text{larger} \\ \text{step?} \end{array} \right.$$

$$\text{update}_t = \gamma \cdot \text{update}_{t-1} + \eta \nabla \mathcal{L}_t$$

$$\mathcal{W}_{t+1} = \mathcal{W}_t - \text{update}_t$$

(-) Oscillations

$$\text{update}_0 = 0$$

$$\text{update}_1 = \eta \nabla \mathcal{L}_1$$

$$\text{update}_2 = \gamma \eta \nabla \mathcal{L}_1 + \eta \nabla \mathcal{L}_2$$

$$\text{update}_3 = \gamma^2 \eta \nabla \mathcal{L}_1 + \gamma \eta \nabla \mathcal{L}_2 + \eta \nabla \mathcal{L}_3$$

5.5 Nesterov Accelerated Gradient Descent

To reduce the oscillations

\* Look before you leap

\* Recall that  $\text{update}_t = \gamma \cdot \text{update}_{t-1} + \eta \nabla \mathcal{L}_t$

$$\mathcal{W}_{t-1} = \mathcal{W}_t - \gamma \cdot \text{update}_{t-1} - \eta \nabla \mathcal{L}_t$$

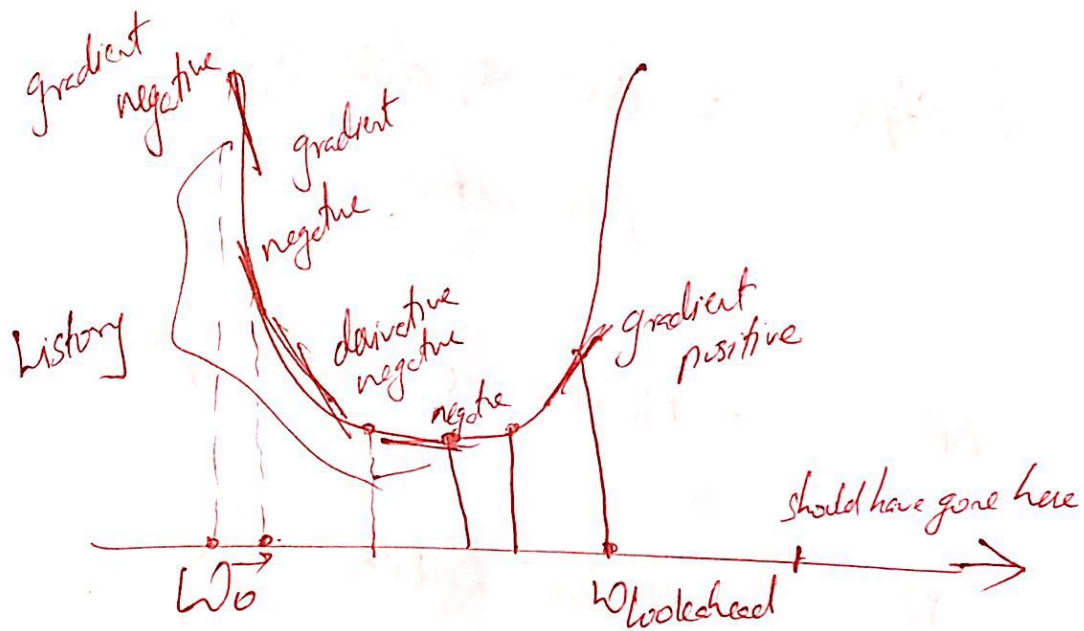


Rule for NAG

$$W_{\text{look-ahead}} = W_t - \gamma \cdot \text{update}_{t-1}$$

$$\text{update}_t = \gamma \cdot \text{update}_{t-1} + \eta \nabla W_t^{\text{look-ahead}}$$

$$W_{t-1} = W_t - \text{update}_t$$



## 5.6 Stochastic & Mini-Batch Gradient Descent

→ Updates the parameter for every single point.

→ Appropriate Gradient (Stochastic)

→ Oscillations (Greedy Decisions)

Higher the value of ' $k$ ', the more accurate are the estimates

1 epoch = one pass over the entire data

1 step = one update of the parameters

$N$  = number of data points

$B$  = Mini Batch Size

Algorithm	# of steps in 1 epoch
Vanilla (Batch Gradient Descent)	1
Stochastic Gradient Descent	$N$
Mini-Batch Gradient Descent	$N/B$

Lecture 5

Module 5.7

Tips for Adjusting Learning Rate  
and Momentum.

Adaptive  $\alpha$   
gentle step - fast  
else - slow

Try learning rate

log scale 0.0001, 0.001, 0.01, 0.1

Step Decay

→ Halve the learning rate after  
every 5 epochs.

→  
Exponential Decay

$$\eta = \eta_0 e^{-kt}, \text{ where } \eta_0 \text{ and } k$$

are hyperparameters,  $t$  is the step number

1/t Decay :  $\eta = \frac{\eta_0}{1+kt}$



Tips for momentum

5.8

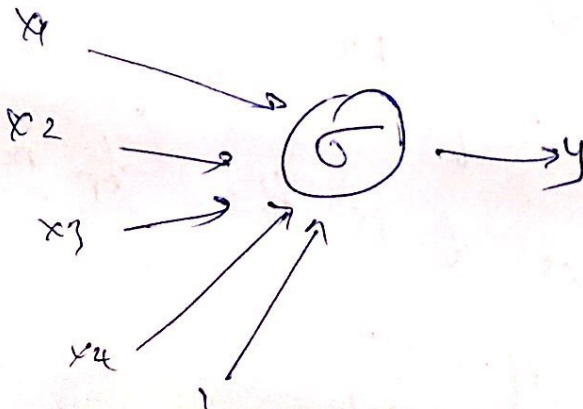
Line search

Using many learning rate at once  
compute the for all the learning  
rates and check it.

If it does a lot of computation in  
one step.

5.9

Gradient Descent with Adaptive Learning  
Rate



$$y = f(x) = \frac{1}{1 + e^{-(wx+b)}}$$

$$x = \{x_1, x_2, x_3, x_4\}$$

$$w = \{w_1, w_2, w_3, w_4\}$$

Decay the learning rate for parameters in proportion to their update history.

Rule for Adagrad

$$V_t = V_{t-1} + (\nabla w_t)^2$$

$$w_{t+1} = w_t - \frac{\eta}{(\sqrt{V_t} + \epsilon)} * \nabla w_t$$

→ divided by the history.

for gradient des

$$w_{t+1} = w_t - \eta * \Delta w_t$$

Adagrad decays the learning rate very aggressively.

The frequent parameters will start receiving very small updates because of the decayed learning rate.



$$V_t = \beta * V_{t-1} + (1-\beta) \nabla W_t^2$$

$$\beta > 0.95$$

$$W_{t+1} = W_t - \frac{\eta}{\sqrt{V_t + \epsilon}} * \nabla W_t$$

RMS props overcomes Adagrad's problem by being less aggressive on the decay.

Adams → moment & killing the learning rate  
→ accumulative history  
→ Moving average of the gradient descent

$$m_t = \beta_1 * m_{t-1} + (1-\beta_1) * \nabla W_t$$

→ accumulative history

$$V_t = \beta_2 * V_{t-1} + (1-\beta_2) * (\nabla W_t)^2$$

$$\hat{m}_t = \frac{m_t}{1-\beta_1^t}$$

$$\hat{V}_t = \frac{V_t}{1-\beta_2^t}$$

Bias correction

$$W_{t+1} = W_t - \frac{\eta}{\sqrt{\hat{V}_t + \epsilon}} * \hat{m}_t$$

→ close to the mean of the distribution

→ accumulative descent

→ divide the learning rate by accumulative history of gradients

# Adaptive Momentum

Adam

$$\beta_1 = 0.9$$

$$\beta_2 = 0.999$$

$$\epsilon = 1e-8$$

Sequence generation problems  $\eta = 0.01, 0.0001$ ,  
works best

• SGD with momentum (Colestarov)  
 $\eta = 0.001, 0.0001$

Best choice would be Adam.

$$m_t = \beta m_{t-1} + (1-\beta) g_t \rightarrow \nabla W_t$$

$$m_0 = 0, \quad m_1 = (1-\beta) g_1$$

$$m_2 = \beta(1-\beta)g_1 + (1-\beta)g_2$$

$$m_3 = \beta^2(1-\beta)g_1 + \beta(1-\beta)g_2 + (1-\beta)g_3$$

$$m_t = 1 - \beta \sum_{i=1}^t \beta^{t-i} g^i$$

$$E[m_t] = E \left[ 1 - \beta \sum_{i=1}^t \beta^{t-i} g^i \right]$$

$$= E[g_t] \cdot (1 - \beta) \sum_{i=1}^t \beta^{t-i} \quad \hookrightarrow \frac{1 - \beta^t}{1 - \beta}$$

$$E[m_t] = (1 - \beta^t) E[g_t]$$

$$E[m_t] = E[g_t]$$

Lec 59 Bias correction in Adam

Update equations for Adam

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

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

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t} \quad \hat{v}_t = \frac{v_t}{1 - \beta_2^t}$$

$$\omega_{t+1} = \omega_t - \frac{\eta}{\sqrt{\hat{v}_t + \epsilon}} \hat{m}_t$$



$$w_{t+1} = w_t - \frac{\eta}{\sqrt{\hat{v}_t + \epsilon}} \propto \hat{m}_t$$

~~Interest~~

Interest

$$E[m_t] = E[w_t]$$

$\nabla w_t$  as  $g_t$

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

$$m_0 = 0$$

$$\begin{aligned} m_1 &= \beta m_0 + (1-\beta) g_1 \\ &= (1-\beta) g_1 \end{aligned}$$

$$m_2 = \beta (1-\beta) g_1 + (1-\beta) g_2$$

$$m_3 = \beta m_2 + (1-\beta) g_3$$

$$= \beta (\beta (1-\beta) g_1 + (1-\beta) g_2) + (1-\beta) g_3$$

$$= \beta^2 (1-\beta) g_1 + \beta (1-\beta) g_2 + (1-\beta) g_3$$

$$= (1-\beta) \sum_{i=1}^3 \beta^{3-i} g_i$$

$$m_t = (1-\beta) \sum_{i=1}^t \beta^{t-i} g_i$$

$$E[m_t] = E\left[(1-\beta) \sum_{i=1}^t \beta^{t-i} g_i\right]$$

$$E[m_t] = (1-\beta) E\left[\sum_{i=1}^t \beta^{t-i} g_i\right]$$

$$= (1-\beta) \sum_{i=1}^t \beta^{t-i} E[g_i]$$

All  $g_i$ 's same distribution

$$E[g_i] = E[g] \quad \forall i$$

$$E[m_t] = (1-\beta) \sum_{i=1}^t \beta^{t-i} E[g]$$

$$= E[g] \cdot (1-\beta) \sum_{i=1}^t \beta^{t-i}$$

$$= E[g] (1-\beta) \sum_{i=1}^t \beta^{t-i}$$

$$= E[g] (1-\beta) \frac{1-\beta^t}{1-\beta}$$

$$E[M_t] = E[g] (1-\beta^t)$$

$$E \left[ \frac{M_t}{(1-\beta^t)} \right] = E[g]$$