# Lecture October 1

Gradient descent methods

- cost function

$$C(\beta) / \hat{\beta} = \min_{\beta \in \mathbb{R}^p} C(\beta)$$

- gradient $g(\beta) = D_\beta C(\beta)$

- Hessian $H(\beta) = \dfrac{\partial^2 C}{\partial \beta \partial \beta^T}$

- Iterative scheme

$$\beta^{(u+1)} = \beta^{(u)} - \gamma^{(u)} D_\beta C(\beta^{(u)})$$
$$= \beta^{(u)} - \gamma^{(u)} g(\beta^{(u)})$$

Log-Reg and linear regression
we have often purely convex
functions. Iteration stops
when $D_\beta C = g = 0$

Linear Regression

$$g \propto X^T(X\beta - y)$$

$$\beta \in \mathbb{R}^p \qquad X \in \mathbb{R}^{m \times p}$$

$$m = 10^5 \qquad p = 10^3$$

FLOPS for $X^TX \sim p^2 \cdot N$

$10^6 \cdot 10^5 \sim 10^{11}$

every iteration $X^TX \in \mathbb{R}^{p \times p}$

$(X^TX)\beta \sim p^2 \sim 10^6$

———————————— * ————————————

Standard SGD: calculate
gradient as expectation value

## Recipe

- subselect, at random,
  $m_B \leq N$ of training point
- fix a learning rate $\gamma^{(0)}$
  $= \gamma_0$

without replacement

- place them in a
  mini batch $B_1$
- compute the gradient
  $$\beta_1^{(1)} = \beta_0^{(0)} - \gamma_0 \sum_{i \in B_1} \nabla_\beta \left( C(\beta^{(0)}) \right)$$
- select randomly another
  $m_B$ mini-batch $B_2$
- continue till we have
  $B$ ,

GN/MB

start with $\beta^{(0)}$, in $\underline{B_1}$

update $\beta_1$  ($\beta_1$ updated)

Then ~~so~~ update $\underline{B_2}$ using $\beta_1$

$$\beta_2 = \beta_1 - \gamma_1 \sum_{i \in B_2} D_\beta C((\overset{\rightarrow}{\beta_1}))$$

— continue till last
mini-Batch,

Defines an epoch. Repeat
for a given number of epochs.

Two more parameters :

1) #mini-batches

2) # epochs

Learning rate; decay $\gamma$
linearly :

$$\gamma_k = (1-\alpha)\underline{\gamma_0} + \alpha \gamma_\tau$$

$$\alpha = \frac{k}{\tau} \qquad \gamma_\tau \sim \frac{1}{100}\gamma_0$$

PICK values for $\gamma_0 \in [10^{-5}, 10^{-4}, \sim 10^{-3}$

$\cdots 1]$

– 2nd-order Taylor approx
to a function $f(x)$

$$f(x^{(0)} - \gamma g) \cong f(x^{(0)})$$
$$- \gamma g^T g + \frac{1}{2} \gamma^2 g^T H g$$
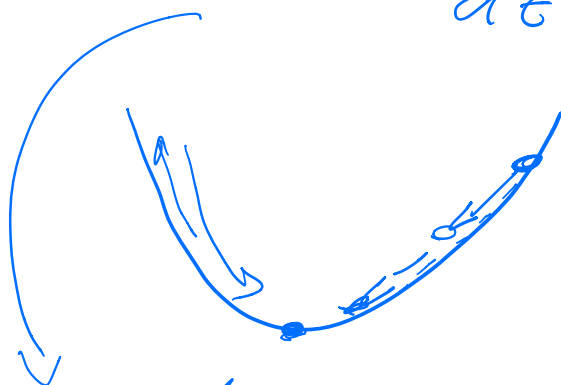
$$\gamma^* = \frac{g^T g}{g^T H g}$$

$$\frac{1}{2} \gamma^2 g^T H g - \boxed{\gamma g^T g}$$

– <mark>Momentum SGD</mark> ;

analogy with Newtonian
mechanics ;

$$F(t) = \frac{d^2 x(t)}{dt^2} \qquad \Big| \quad \beta^{(m)} \longrightarrow x(t)$$

$$\frac{dx}{dt} = v \quad \wedge \quad \frac{dv}{dt} = F(t)$$

Euler Forward method
$$\begin{cases} v_{i+1} = v_i + h\,\underline{F(t_i)} = v_i + h\,F_i \\ x_{i+1} = x_i + h\,v_i \end{cases}$$

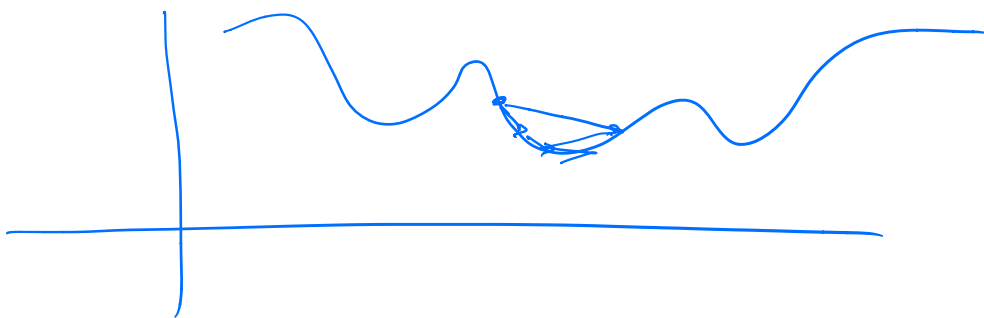$F(t_i)$ plays the role of $\nabla_\beta C(\beta)$

$$\rightarrow \quad v_{i+1} = \underline{\alpha\,v_i} + h\,F_i$$

$$\beta^{(m+1)} = \beta^{(m)} - \gamma^{(m)}\,\nabla_\beta C(\underline{\beta}^{(m)})$$

Momentum SGD

$$v^{(m)} = \boxed{\alpha\,v^{(m-1)} - \gamma^{(m-1)}\, \times \nabla_\beta C(\beta^{(m-1)})}$$

$$\beta^{(m+1)} = \beta^{(m)} - v^{(m)}$$



Adagrad

SGD : accumulate the square gradient

$$z^{(m+1)} = z^{(m)} + g^T g$$

$$\left( \beta^{(m+1)} = \beta^{(m)} - \frac{\gamma^{(m)}}{\delta + \sqrt{\tau}} \, g \right.$$
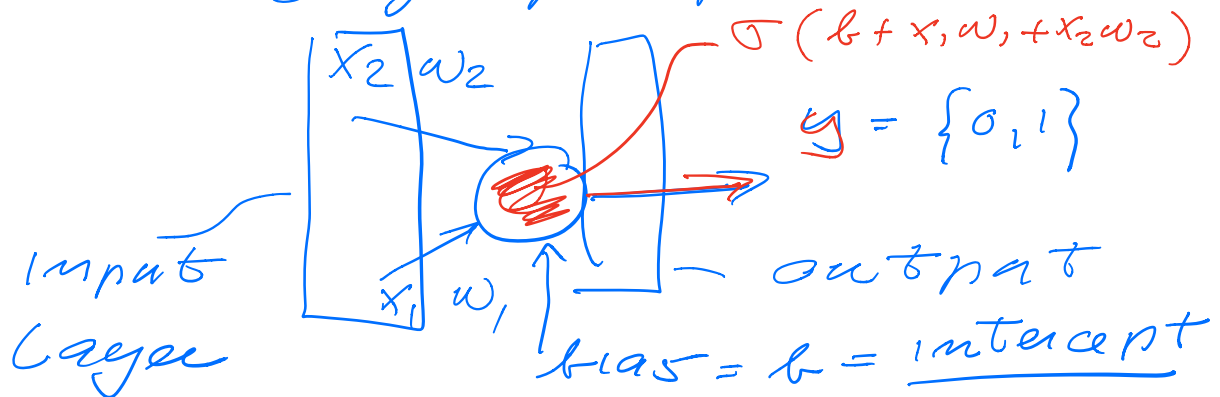
$$\delta \sim 10^{-8}$$

## RMSPROP

$$r^{(m+1)} = \rho \, r^{(m)} + (1-\rho) \, g^T g$$

$$\beta^{(m+1)} = \beta^{(m)} - \frac{\gamma^{(m)}}{\sqrt{\delta + \tau}} \, g$$

## Neural networks

— single perceptron model

$$\sigma(b + x_1 w_1 + x_2 w_2)$$

$$y = \{0, 1\}$$

$X_2$ $w_2$

input layer

$x_1$ $w_1$

bias $= b =$ intercept

— output
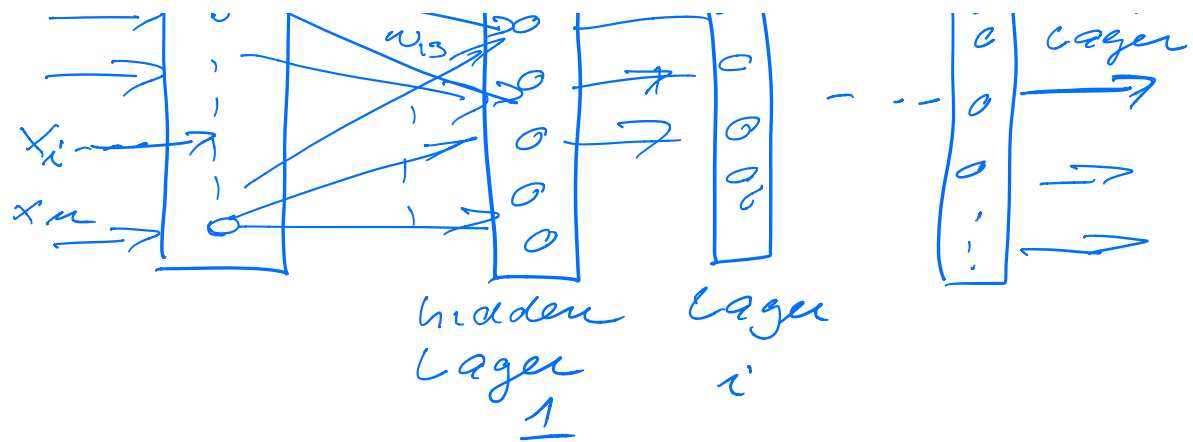
$$y = x_1 w_1 + x_2 w_2 + b$$

$$(\beta_0 + \beta_1 x_1 + \beta_2 x_2)$$

— multi-layer perceptron

input layer

$x_1$

$w_{11}$

$w_{12}$

$\sigma$ output

Layer

$x_i$
$x_m$

hidden Layer
Layer 1

Layer i

## $XOR$ - model



$b$

$\sigma'(b_1 \omega_1 + h_2 \omega_2 + c)$

$y$

$x_1$    $x_1$    $h_1$

$x_2$    $x_2$    $h_2$

$\sigma(x_2 \omega_{22} + x_1 \omega_{21})$
$+$



$x_2$

| | | |
| 1 | $\boxed{1}$ $y$ | $\boxed{0}$ $y$ |
| 0 | $\boxed{0}$ $y$ | $\boxed{1}$ $y$ |

0    1    $x_1$