

OPTIMAL LECTURE 4, 8/2/2021

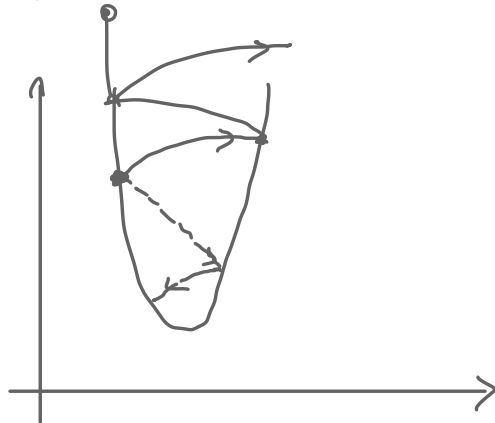
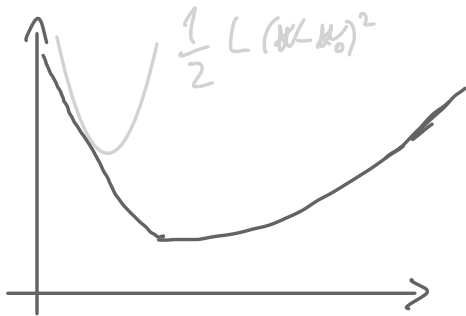
$$w_{k+1} = w_k - \alpha_k g_k$$

$$g_k = \frac{1}{n_k} H_k \sum_{i=1}^{n_k} \nabla f(w_k, \underbrace{x_{k,i}})$$

↑

$$H_k = I$$

$$n_k = \begin{cases} 1 & \text{SGD} \\ n_B & \text{MB} \\ |S| & \text{FB} \end{cases}$$



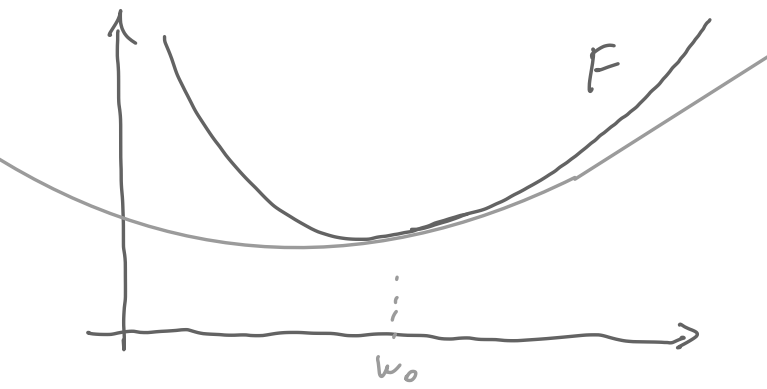
$$\|\nabla F(w') - \nabla F(w)\| \leq L \|w' - w\|$$

$$\frac{\|\nabla F(w') - \nabla F(w)\|}{\|w' - w\|} \leq L$$

$$\alpha < \frac{\mu}{L \eta_b}$$

$$L = \infty$$

$$\alpha = 0$$



$$\frac{1}{2} c (w - w_0)^2$$

$$\bar{F}(w) = w^2$$



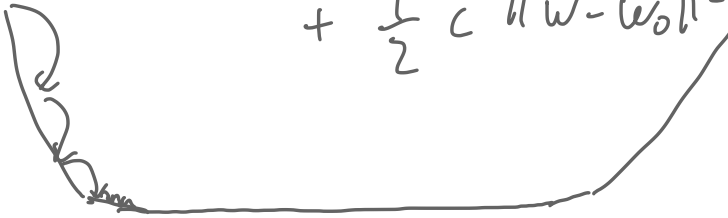
$$F(w) \geq F(w_0)$$

$$+ \nabla F(w_0)^T (w - w_0)$$

$$+ \frac{1}{2} c \|w - w_0\|^2$$

$$w^2$$

$$2w$$

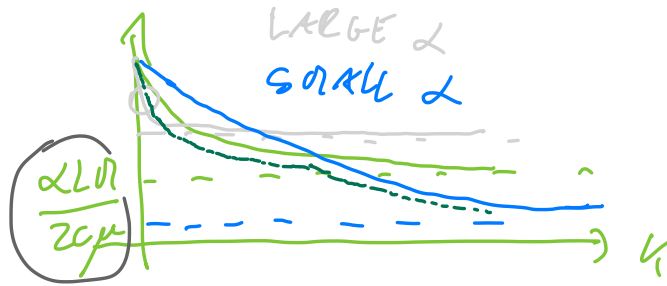


STRONGLY CONVEXITY \Rightarrow GUARANTEE
A MINIMUM
MOVEMENT
(NOT TOO SLOW
CONVERGENCE)

$L \downarrow$ THE BETTER

$c \nearrow$ THE BETTER

$$\underline{E[F(w_n) - F^*]} \leq \frac{\alpha L \sigma}{2c\mu} + \underbrace{(1 - \alpha c\mu)^{n-1}}_x \times \left(F(w_1) - F^* - \frac{\alpha L \sigma}{2c\mu} \right)$$



$$\rho = (1 - \alpha c\mu) \quad c \nearrow \quad \rho \searrow$$

M : CAPTURES THE NOISE
 \uparrow

$$\alpha \nearrow \quad (1 - \alpha c\mu) \searrow$$

$$\frac{\alpha L \sigma}{2c\mu} \nearrow$$

DIMINISHING LEARNING RATE

$$\alpha_k = \frac{1}{k}$$

$$\alpha_k \rightarrow 0$$

$$\sum_{k=1}^{\infty} \alpha_k = +\infty$$

YOU CAN GO EVERYWHERE

$$\sum_{k=1}^{\infty} \alpha_k^2 < +\infty$$

FILTER OUT THE NOISE

Theorem (4.7, Bottou) Assumptions of last lecture

$$\alpha_k = \frac{\beta}{\gamma + k}$$

$$\beta > \frac{1}{c\mu}$$

$$\gamma > 0$$

$$\alpha_1 < \frac{\mu}{L \log 2}$$

Then: $\mathbb{E}[F(w_k) - F^*] \leq \frac{v}{\gamma + k}$

$$v = \max \left(\frac{\beta^2 L \sigma^2}{2(\beta c \mu - 1)}, (\gamma + 1)(F(w_1) - F^*) \right)$$

$$v \approx \frac{\beta^2 L \sigma^2}{2(\beta c \mu - 1)}$$

$$\left(\frac{\beta^2}{(\beta c \mu - 1)} \right)' = \frac{2\beta(\beta c \mu - 1) - \beta^2 c \mu}{(\quad)^2} = 0$$

$$2\beta(\beta c \mu - 1) - \beta^2 c \mu = 0$$

$$2\beta c \mu - 2 - \beta c \mu = 0$$

$$\beta c \mu = 2$$

$$\beta = \frac{2}{c \mu} \Rightarrow \text{minimize } V$$

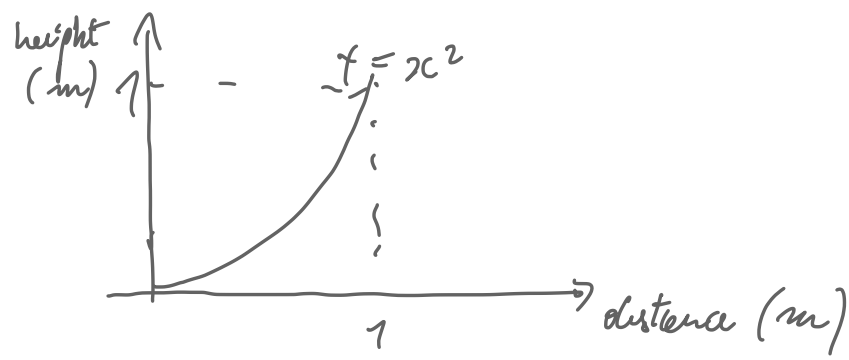
$$V = \frac{\frac{2}{c \mu} L \eta}{\chi \left(\frac{2}{c \mu} c \mu - 1 \right)} = \frac{2 L \eta}{c^2 \mu^2} =$$

$$= \frac{2}{\mu^2} \left(\frac{L}{c} \right) \frac{\eta}{c}$$

$\frac{L}{c}$: CONDITION NUMBER

"THE PROBLEM IS ILL-CONDITIONED"

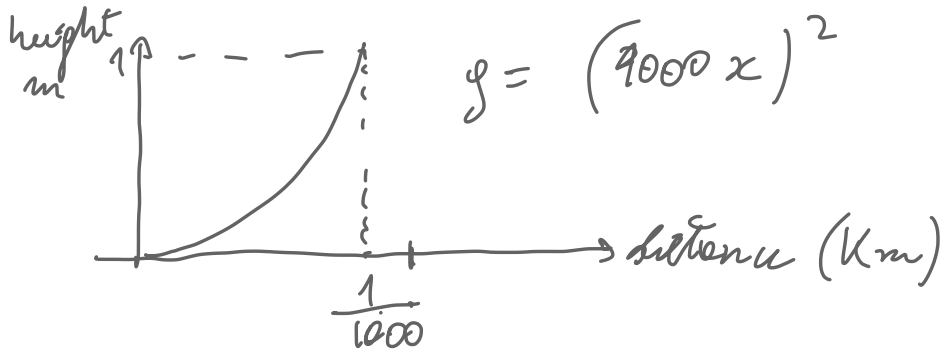
$= \frac{L}{c}$ IS LARGE



$$y = x^2$$

$$L = 2$$

$$c = 2$$



$$L = 2 \times 10^6$$

$$c = 2 \times 10^6$$

$$\frac{L}{c} = 1$$

LOOKING - FORWARD :

THE SOLUTION WILL BE TO
TAKE A LEARNING RATE THAT
IS FUNCTION OF THE CURVATURE

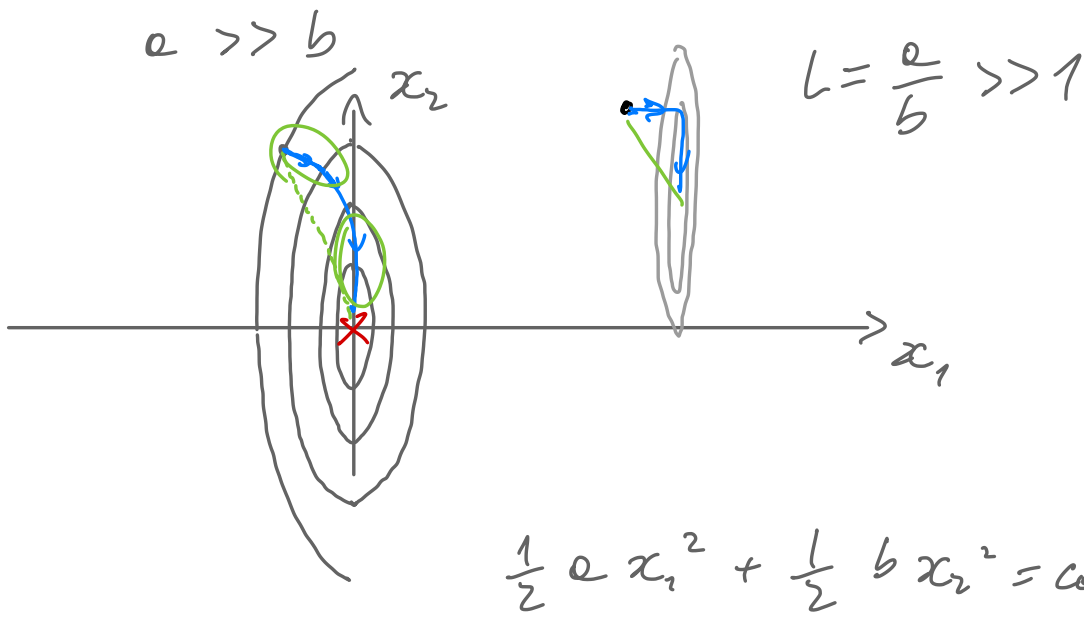
$$y = \frac{1}{2} a x_1^2 + \frac{1}{2} b x_2^2$$

$$L = \max(a, b) \quad c = \min(a, b)$$

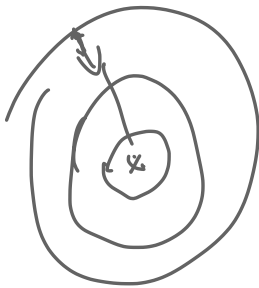
$$y = \frac{1}{2} a x_1^2 + \frac{1}{2} b x_2^2 \leq \frac{1}{2} \max(a, b) (x_1^2 + x_2^2)$$

$$\quad \quad \quad \uparrow$$

$$= \frac{1}{2} L (x_1^2 + x_2^2)$$



$$a = b$$



$$\frac{L}{c} = 1$$

$$F(x_1, x_2) = \frac{a}{2} x_1^2 + \frac{b}{2} x_2^2$$

$$\nabla F = \begin{bmatrix} ax_1 \\ bx_2 \end{bmatrix}$$

$$\begin{bmatrix} w_{k+1,1} \\ w_{k+1,2} \end{bmatrix} = \begin{bmatrix} w_{k,1} \\ w_{k,2} \end{bmatrix} - \alpha \begin{bmatrix} a w_{k,1} \\ b w_{k,2} \end{bmatrix}$$

$$\equiv \begin{cases} w_{k+1,1} = w_{k,1} - \alpha a w_{k,1} \\ w_{k+1,2} = w_{k,2} - \alpha b w_{k,2} \end{cases}$$

$$a \gg b$$

$$\alpha_1 = \frac{\alpha}{a}$$

$$\alpha_2 = \frac{\alpha}{b}$$

$$\begin{cases} w_{k+1,1} = w_{k,1} - \alpha w_{k,1} \\ w_{k+1,2} = w_{k,2} - \alpha w_{k,2} \end{cases}$$

$$w_{k+1} = w_k - \alpha \begin{bmatrix} \frac{1}{a} & 0 \\ 0 & \frac{1}{b} \end{bmatrix} \nabla F(w_k)$$

$$= w_k - \alpha H \nabla F(w_k)$$

$$g_k = \frac{1}{n_k} H_k \sum_{i=1}^{n_k} \nabla f(w_k, \xi_{k,i})$$

H "Hessian"

$$\text{Hessian } f = \begin{bmatrix} \frac{\partial^2 F}{\partial x_1^2} & \frac{\partial^2 F}{\partial x_1 \partial x_2} \\ \frac{\partial^2 F}{\partial x_1 \partial x_2} & \frac{\partial^2 F}{\partial x_2^2} \end{bmatrix} =$$

$$F(x_1, x_2) =$$

$$= \frac{1}{2} a x_1^2 +$$

$$\frac{1}{2} b x_2^2$$

$$= \begin{bmatrix} \underline{a} & \underline{0} \\ \underline{0} & \underline{b} \end{bmatrix}$$

$$\begin{bmatrix} \frac{1}{a} & 0 \\ 0 & \frac{1}{b} \end{bmatrix} = \begin{bmatrix} a & 0 \\ 0 & b \end{bmatrix}^{-1} = H_F^{-1}$$

$$w_{k+1} = w_k - \alpha_k H_F^{-1} \nabla F(w_k)$$

Newton Method

MUCH FASTER THAN BASIC GRADIENT

$$H_F = d \times d \quad d : \# \text{ parameters in the model.}$$

memory!

$$H_F^{-1} \quad \text{COMPLEXITY} \quad d^3$$

$$d^{2.37}$$

$$\underbrace{H_F^{-1} \nabla F}_{\text{REALLY WHAT YOU NEED}}$$

COMPUTE IT BY SOLVING

$$\underbrace{H_F}_{\text{matrix}} \Delta u = \nabla F$$

$$\Delta u = H_F^{-1} \nabla F$$

$\Delta u \simeq$ A SOLUTION

\Rightarrow CONJUGATE GRADIENT METHOD

$$\begin{bmatrix} \frac{1}{a} & 0 \\ 0 & \frac{1}{b} \end{bmatrix} = H_F^{-1}$$

$$\begin{bmatrix} \frac{1}{\frac{\partial^2 F}{\partial x_1^2}} & 0 & 0 & 0 & 0 \\ 0 & \frac{1}{\frac{\partial^2 F}{\partial x_2^2}} & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 \dots \frac{1}{\frac{\partial^2 F}{\partial x_d^2}} \end{bmatrix}$$

$$H_F = \begin{bmatrix} \frac{\partial^2 F}{\partial x_1^2} & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & & \frac{\partial^2 F}{\partial x_d^2} \end{bmatrix}$$

TRUE ONLY IF

$$F(x_1, \dots, x_d) = \sum_{i=1}^d f_i(x_i)$$

$$f(w, (x, y)) = \ell(h(w, x), y)$$

H_f

$$\simeq \ell(h(w', x) + J_h(w', x)(w - w'), y)$$

GAUSS-NEWTON METHOD

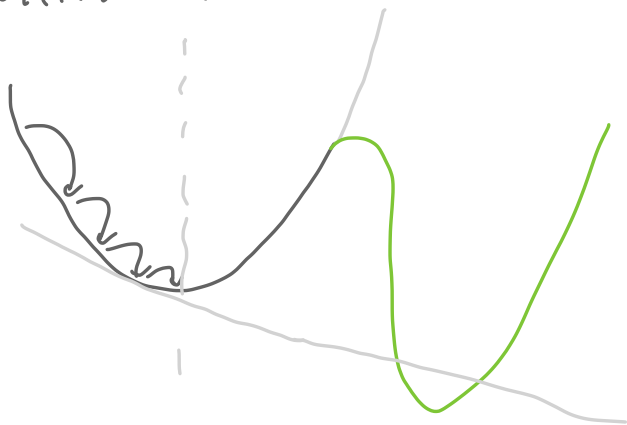
Iteratively update of H_u

$$H_u = G(H_{u-1}, w_u)$$

BFGS

UNTIL NOW WE CONSIDERED CONVEXITY

WHAT IF NO CONVEXITY?



$$\mathbb{E}[F(w_k)] - \underset{\uparrow}{F^*} \leq \text{SOMETHING}$$

OPTIMUM

Theorem (4.8, Bottom)

ASSUMPTIONS ABOUT NOISE
AND L-SMOOTHNESS
WE GIVE UP CONVEXITY.

$$\alpha \leq \frac{\mu}{L \sigma_G}$$

$$\mathbb{E} \left[\sum_{k=1}^K \| \nabla F(w_k) \|^2 \right] \leq \frac{K \alpha L \sigma_G}{\mu} + \frac{2(F(w_1) - F_{\text{opt}})}{\mu \alpha}$$

$$\underbrace{\frac{1}{K} \mathbb{E} \left[\sum_{k=1}^K \|\nabla F(w_k)\|^2 \right]}_{\text{red bracket}} \leq \underbrace{\frac{\alpha L \sigma^2}{\mu}}_{\text{red circle}} + \underbrace{\frac{2(F(w_k) - F_{\inf})}{K\mu}}_{\text{red circle}}$$

$$F_{\inf} \leq F(w) \quad \forall w$$

$$F_{\inf} = F^* = \min_w F(w)$$

SIMILAR RESULT FOR DECREASING α_k

$$\alpha_k = \frac{1}{k}$$

