

Lecture September 25

Cost function $C(\beta)$

$$\frac{\partial C}{\partial \beta} = g(\beta)$$

$$\frac{\partial^2 C}{\partial \beta \partial \beta^T} = H(\beta) \quad \text{Hessian matrix}$$

Can add regularization
(L2-regularization)

$$C'(\beta) = C(\beta) + \lambda \beta^T \beta$$

$$g'(\beta) = g(\beta) + 2\lambda \beta$$

$$H'(\beta) = H(\beta) + 2\lambda I$$

Math aside

$$* \quad \frac{\partial (a^T A a)}{\partial a} = (A + A^T) a$$

(Hessian $H = H^T$)

$$** \quad \frac{\partial (b^T a)}{\partial a} = b$$

Newton's method;
iterative scheme

$$\beta^{(n+1)} = \beta^{(n)} - [H^{(n)}]^{-1} g^{(n)}$$

$$H^{(n)} = H(\beta^{(n)})$$

$$g^{(n)} = g(\beta^{(n)})$$

Find $\beta^{\text{opt}} = \hat{\beta}$ when

$$\|\beta^{(n+1)} - \beta^{(n)}\|_2 \leq \varepsilon \sim 10^{-8}$$

$$\beta^{(n+1)} = \beta^{(n)} - \gamma^{(n)} g^{(n)}$$

Learning rate,

$$\gamma^{(n)} \leq \frac{2}{\lambda_{\max}} \leftarrow \text{Largest eigenvalue}$$

of H
Gradient/steepest descent.

$$\begin{aligned} C(\hat{\beta}) &\cong C(\beta^{(n)}) + \\ &(\text{Taylor expand around } \hat{\beta} - \beta^{(n)}) \\ &+ [g^{(n)}]^T (\hat{\beta} - \beta^{(n)}) + \\ &\frac{1}{2} (\hat{\beta} - \beta^{(n)})^T H^{(n)} (\hat{\beta} - \beta^{(n)}) \end{aligned}$$

$$= \overset{\text{Drop 1}}{C + \frac{1}{2} \beta^T A \beta + b^T \beta} \Leftarrow$$

$$A = H^{(n)} = \text{known}$$

$$b = g^{(n)} - H^{(n)} \beta^{(n)}$$

= known

$$C = C(\beta^{(n)}) - [g^{(n)}]^T \beta^{(n)} + \frac{1}{2} [\beta^{(n)}]^T H^{(n)} \beta^{(n)}$$

$$\frac{\partial C}{\partial \beta} = A\beta + b = 0$$

(using *, **)

$$\Rightarrow \beta = -A^{-1}b = \beta^{(n)} - [H^{(n)}]^{-1} g^{(n)}$$

which is what we get from Newton's method.

Can generalize to a function $f(x) = \boxed{\frac{1}{2} x^T A x + b^T x} + C$

$$\frac{\partial f}{\partial x} = 0 = Ax + b \Rightarrow$$

$$\boxed{Ax = -b = +\tilde{b}}$$

... .. + doesn't

Standard steepest descent

iterative method

guess $x_0 = x^{(0)}$ ($Ax=b$)

$$r^{(0)} = b - Ax^{(0)}$$

$$\rightarrow r^{(k+1)} = b - Ax^{(k+1)}$$

$$= b - A(x^{(k)} + \alpha^{(k)} r^{(k)})$$

$$(r^{(k)})^T, \quad \underbrace{b - Ax^{(k)}}_{r^{(k)}} - \alpha^{(k)} A r^{(k)} = 0$$
$$r^{(k)} = \alpha^{(k)} A r^{(k)}$$

$$[r^{(k)}]^T r^{(k)} = \alpha^{(k)} [r^{(k)}]^T A r^{(k)}$$

$$\underline{\alpha^{(k)}} = \frac{[r^{(k)}]^T r^{(k)}}{[r^{(k)}]^T A r^{(k)}}$$

$$x^{(k+1)} = x^{(k)} - \alpha^{(k)} r^{(k)}$$

$$\|r^{(k+1)} - r^{(k)}\|_2 \leq \epsilon$$

in our case $A = \text{Hessian matrix}$

$f(x)$ Taylor expansion at a gradient with γ -step

$$f(x) \approx f(x^{(n)}) + (x - x^{(n)})^T g^{(n)} + \frac{1}{2} (x - x^{(n)})^T H^{(n)} (x - x^{(n)})$$

$$(x = x^{(n)} - \gamma g^{(n)})$$

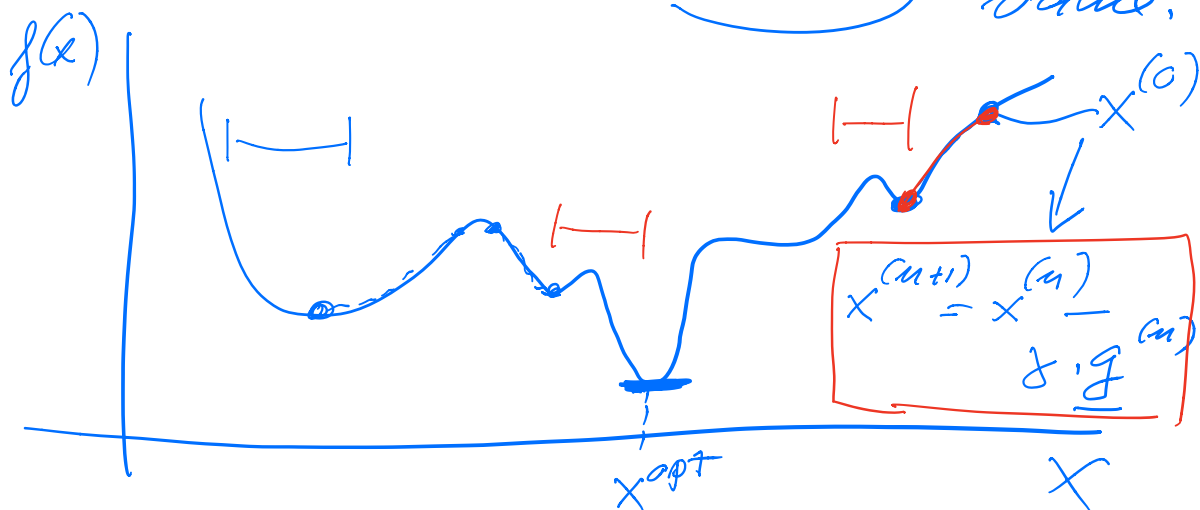
$$\approx f(x^{(n)}) - \gamma [g^{(n)}]^T g^{(n)} + \frac{1}{2} \gamma^2 [g^{(n)}]^T H^{(n)} g^{(n)}$$

Minimize wrt to γ

$$\gamma^* = \frac{g^{(n)T} g^{(n)}}{g^{(n)T} H^{(n)} g^{(n)}}$$

$$\gamma_{opt} < \frac{2}{\lambda_{max}}$$

Largest eigen value,



$g^{(m)}$ calculated for all points

$$x_i \quad i = 0, 1, \dots, m-1$$

sample a subset of x_i
(stochastically)

$E[g^{(m)}]$ has standard
deviation $\sim \sigma / \sqrt{m}$

$$m = 100$$

$$\downarrow$$
$$\sigma/10$$

$$\boxed{m = 10000}$$

$$\downarrow$$
$$\sigma/100$$

$$X^T X$$