

LECTURE 3

PRACTICAL NEWTON METHODS

Modified Hessian

$$M_k := \begin{cases} \nabla^2 f_k & \lambda_{\min}(\nabla^2 f_k) \gg \delta \\ \nabla^2 f_k + \tau_k I & \text{otherwise} \end{cases}$$

fixed positive scalar \uparrow

(Based on Identity Shifts) where $\tau_k := \delta - \lambda_{\min}(\nabla^2 f_k)$

$$M_k := Q_k \begin{bmatrix} \hat{\lambda}_1 & & & 0 \\ & \hat{\lambda}_2 & & \\ & & \dots & \\ 0 & & & \hat{\lambda}_n \end{bmatrix} Q_k^T$$

given orthogonal eigenvalue decomposition

$$\nabla^2 f_k = Q_k \begin{bmatrix} \lambda_1 & & & 0 \\ & \lambda_2 & & \\ & & \dots & \\ 0 & & & \lambda_n \end{bmatrix} Q_k^T$$

(Based on Orthogonal Eigenvalue Decomposition)

$$\text{and } \hat{\lambda}_k := \max(\delta, \lambda_k).$$

THM

The modified Hessian by identity shifts satisfies

$$\|M_k - \nabla^2 f_k\|_2 \leq \|S - \nabla^2 f_k\|_2$$

for every S such that $S^T = S$ and $\lambda_{\min}(S) \gg \delta$

THM

The modified Hessian by orthogonal eigenvalue decomposition satisfies

$$\|M_k - \nabla^2 f_k\|_F \leq \|S - \nabla^2 f_k\|_F$$

for every S such that $S^T = S$ and $\lambda_{\min}(S) \geq \delta$.

A complete framework
Compute $\nabla f_0, \nabla^2 f_0$ and $k \leftarrow 0$
While $\|\nabla f_k\| > \epsilon$

Form M_k from $\nabla^2 f_k$

Solve $M_k p_k = -\nabla f_k$

Determine an α_k satisfying Wolfe conditions

$$x^{(k+1)} \leftarrow x^{(k)} + \alpha_k p_k$$

Compute $\nabla f_{k+1} := \nabla f(x^{(k+1)})$, $\nabla^2 f_{k+1} := \nabla^2 f(x^{(k+1)})$

$k \leftarrow k+1$

end

Guaranteed to converge if $\|\nabla^2 f(x)\|_2$ is bounded for all $x \in \mathbb{R}^n$.

Practical Variant becomes the pure Newton's method for large k provided

$$\lim_{k \rightarrow \infty} x^{(k)} = x_*$$

such that

$$\lambda_{\min}(\nabla^2 f(x_*)) > \delta.$$

Specifically, for large k

(i) $M_k = \nabla^2 f_k$

(ii) $\alpha_k = 1$ satisfies Wolfe conditions (if $c_1 < \frac{1}{2}$).

To see (ii), noting $p_k = -[\nabla^2 f_k]^{-1} \nabla f_k$ for large k , by Taylor's thm

$$\begin{aligned} f(x_*^{(k)} + p_k) &= f(x_*^{(k)}) + \nabla f(x_*^{(k)})^T p_k \\ &\quad + \frac{1}{2} p_k^T \nabla^2 f(x_*^{(k)} + t p_k) p_k \\ &= f(x_*^{(k)}) + \nabla f(x_*^{(k)})^T p_k + \frac{1}{2} p_k^T \nabla^2 f(x_*^{(k)}) p_k \\ &\quad + O(\|\nabla f_k\|^3) \\ &= f(x_*^{(k)}) + \frac{1}{2} \nabla f(x_*^{(k)})^T p_k + O(\|\nabla f_k\|^3) \\ &< f(x_*^{(k)}) + c_1 \nabla f(x_*^{(k)})^T p_k \end{aligned}$$

satisfaction of sufficient decrease for $\alpha_k = 1$

By Taylor's thm with integral remainder

satisfaction
of sufficient
curvature
by $\alpha_k=1$

$$\left\{ \begin{array}{l} \left| \nabla f(x^{(k)} + p_k)^T p_k \right| = \\ \left| \left[\nabla f(x^{(k)}) + \int_0^1 \nabla^2 f(x^{(k)} + t p_k) p_k dt \right]^T p_k \right| = \\ \left| \nabla f(x^{(k)})^T p_k + \int_0^1 p_k^T \nabla^2 f(x^{(k)} + t p_k) p_k dt \right| = \\ O(\|\nabla f_k\|^3) \leq c_2 \left| \nabla f(x^{(k)})^T p_k \right| \\ \implies \\ \nabla f(x^{(k)} + p_k)^T p_k \geq c_2 \nabla f(x^{(k)})^T p_k. \end{array} \right.$$

THM (Local Convergence of Newton's Method)

Suppose $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is a twice-continuously differentiable function with Lipschitz-continuous Hessian, that is there exists a $\gamma > 0$ such that

$$\|\nabla^2 f(x) - \nabla^2 f(y)\|_2 \leq \gamma \|x - y\|_2 \text{ for all } x, y \in \mathbb{R}^n.$$

Suppose also $x_* \in \mathbb{R}^n$ is a point such that $\nabla f(x_*) = 0$ and $\nabla^2 f(x_*)$ is invertible.

There exists a ball $B(x_*, r)$ such that if $x_{\#}^{(k)} \in B(x_*, r)$ (recall $\{x^{(k)}\}$ is the sequence generated by pure Newton's method), then

the following hold:

(i) $x^{(k)} \in B(x_*, r)$ for each $k \geq K$;

(ii) $\lim_{k \rightarrow \infty} x^{(k)} = x_*$;

(iii) $\frac{\|x^{(k+1)} - x_*\|}{\|x^{(k)} - x_*\|^2} \leq \frac{M\gamma}{2}$ for each $k \geq K$

and for each $M > \|\left[\nabla_f^2(x_*)\right]^{-1}\|_2$.

PROOF

By continuity of second derivatives, there exists a ball $B(x_*, r_1)$ such that

$$\|\left[\nabla_f^2(x)\right]^{-1}\|_2 \leq M.$$

Let $n := \min(n_1, \frac{r_1}{2M\gamma})$. \rightarrow for any $n_2 < \frac{2}{M\gamma}$

Suppose $x^{(k)} \in B(x_*, n)$. By Taylor's thm

$$\nabla f(x_*) = \nabla f_k + \int_0^1 \nabla^2 f(x^{(k)} + t(x_* - x^{(k)})) (x_* - x^{(k)}) dt$$

$$0 = \nabla f_k + \nabla^2 f_k (x_* - x^{(k)}) + \int_0^1 \left[\nabla^2 f(x^{(k)} + t(x_* - x^{(k)})) - \nabla^2 f_k \right] \times (x_* - x^{(k)}) dt$$

$$0 = \underbrace{\left[\nabla^2 f_k\right]^{-1} \nabla f_k}_{-p_k = x^{(k)} - x^{(k+1)}} + (x_* - x^{(k)}) + \left[\nabla^2 f_k\right]^{-1} \int_0^1 \left[\nabla^2 f(x^{(k)} + t(x_* - x^{(k)})) - \nabla^2 f_k \right] \times (x_* - x^{(k)}) dt$$

$$\Rightarrow x^{(k+1)} - x_* = \left[\nabla^2 f_k \right]^{-1} \int_0^1 \left[\nabla^2 f(x^{(k)} + t(x_* - x^{(k)})) - \nabla^2 f_k \right] \times (x_* - x^{(k)}) dt$$

Taking the norms of both sides

$$\begin{aligned} (R) \quad \|x^{(k+1)} - x_*\| &\leq \left\| \left[\nabla^2 f_k \right]^{-1} \right\| \int_0^1 \left\| \nabla^2 f(x^{(k)} + t(x_* - x^{(k)})) - \nabla^2 f_k \right\| \times \|x^{(k)} - x_*\| dt \\ &\leq M \int_0^1 \gamma t \|x^{(k)} - x_*\|^2 dt \\ &= \frac{M\gamma}{2} \|x^{(k)} - x_*\|^2. \end{aligned}$$

Since $\|x^{(k)} - x_*\| \in B(x_*, \rho)$, we have

$$\begin{aligned} \|x^{(k)} - x_*\| &< 2 / (M\gamma) \\ \Rightarrow \frac{M\gamma}{2} \|x^{(k)} - x_*\| &< 1. \end{aligned}$$

It follows from (R) that

$$\|x^{(k+1)} - x_*\| < \|x^{(k)} - x_*\|,$$

so $x^{(k+1)} \in B(x_*, \rho)$ proving (i).

Moreover defining $\rho := \frac{M\gamma}{2} \|x^{(k)} - x_*\| < 1$

$$\begin{aligned} \|x^{(k+1)} - x_*\| &\leq \left[\frac{M\gamma}{2} \|x^{(k)} - x_*\| \right] \|x^{(k)} - x_*\| \\ &\leq \underbrace{\left[\frac{M\gamma}{2} \|x^{(k)} - x_*\| \right]}_{\rho} \|x^{(k)} - x_*\|. \end{aligned}$$

⑥

Consequently,

$$\|x^{(k)} - x_*\| \leq \rho^{k-k'} \|x^{(k')} - x_*\|$$

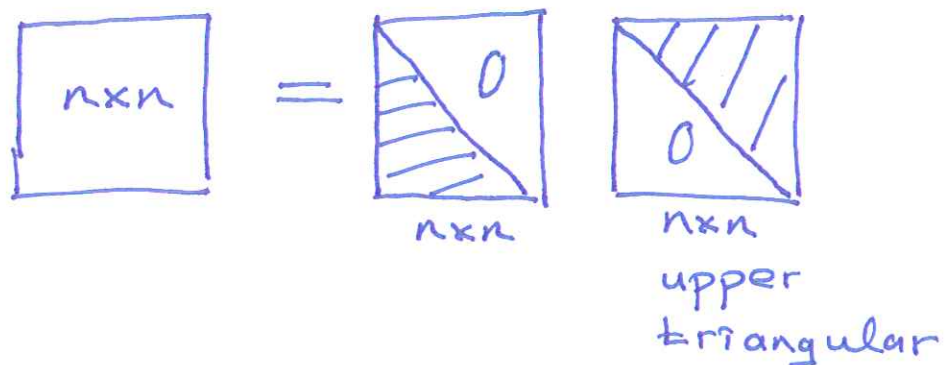
Taking the limit as $k \rightarrow \infty$ yields (ii).

Finally (iii) immediately follows from (R). \square

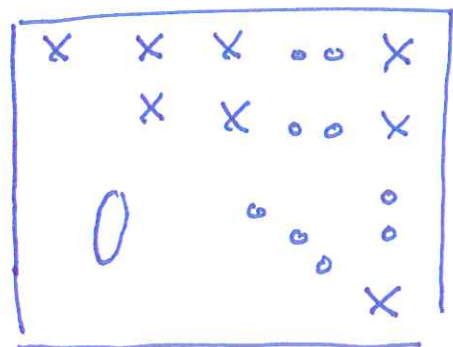
Hessian modifications based on Cholesky factorization

Suppose A is PD, then it has a Cholesky factorization of form

$$A = R^T R$$



Computation of R



row-by-row
~~top to bottom~~ top to bottom

To compute r_{ij} , $i \leq j$, exploit

$$\begin{aligned} a_{ij} &= r_i^T r_j \\ &= \sum_{k=1}^i r_{ki} r_{kj} \end{aligned}$$

$$\underline{i=j}$$

$$a_{ii} = \sum_{k=1}^i r_{ki}^2$$

$$\Rightarrow r_{ii} = \sqrt{a_{ii} - \sum_{k=1}^{i-1} r_{ki}^2}$$

rows 1, ..., i-1

~~$$j = i+1, i+2, \dots, n$$~~

~~$$i = j+1, j+2, \dots, n$$~~

$$r_{ij} = \left[a_{ij} - \sum_{k=1}^{i-1} r_{ki} r_{kj} \right] / r_{ii}$$

rows 1, ..., i-1

Pseudocode

$R \leftarrow 0$

for $i = 1, \dots, n$

$$r_{ii} \leftarrow \sqrt{a_{ii} - \sum_{k=1}^{i-1} r_{ki}^2} \quad c_{ii} :=$$

for $j = i+1, \dots, n$

$$c_{ij} \leftarrow \left[a_{ij} - \sum_{k=1}^{i-1} r_{ki} r_{kj} \right]$$

$$r_{ij} \leftarrow c_{ij} / r_{ii}$$

end

end

The only way this can break-down

$c_{ii} \leq 0$ for some i

Modify so that

$$c_{ii} \geq \delta$$

and

$$|r_{ij}| \leq \beta$$

given positive scalars
 $j = i+1, \dots, n.$

Setting

$$r_{ii} := \sqrt{\max(c_{ii}, \delta, (\frac{\theta}{\beta})^2)}$$

where $\theta := \max_{i+1 \leq j \leq n} |c_{ij}|$

achieves both. In particular,

$$|r_{ij}| \leq \frac{|c_{ij}|}{(\theta/\beta)} \leq \beta.$$

Modified Pseudocode

$R \leftarrow 0$

for $i = 1, \dots, n$

$$c_{ii} \leftarrow a_{ii} - \sum_{k=1}^{i-1} r_{ki}^2$$

for $j = i+1, \dots, n$

$$c_{ij} \leftarrow a_{ij} - \sum_{k=1}^{i-1} r_{ki} r_{kj}$$

end

~~$r_{ii} \leftarrow$~~

$$\theta \leftarrow \max_{i+1 \leq j \leq n} |c_{ij}|$$

