

Lecture 19. Newton's method. KKT condition for ICP.

Recall Newton's method. $\min_d \tilde{f}(x_k + d) = f(x_k) + \nabla f(x_k)^T d + \frac{1}{2} d^T \nabla^2 f(x_k) d$.

Apply to equality constrained problems: $\min_d \tilde{f}(x_k + d)$ s.t. $Ad = 0$.

KKT system for this quadratic problem is: $\begin{pmatrix} \nabla^2 f(x_k) & A^T \\ A & 0 \end{pmatrix} \begin{pmatrix} d \\ -\lambda \end{pmatrix} = \begin{pmatrix} -\nabla f(x_k) \\ 0 \end{pmatrix}$

Newton's method for equality constrained problem: select x_0 s.t. $Ax_0 = b$.

repeat. compute d by solving KKT. $\begin{pmatrix} \nabla^2 f(x_k) & A^T \\ A & 0 \end{pmatrix} \begin{pmatrix} d \\ -\lambda \end{pmatrix} = \begin{pmatrix} -\nabla f(x_k) \\ 0 \end{pmatrix}$

(note if $A=0$. then $d = -(\nabla^2 f(x_k))^{-1} \nabla f(x_k)$. exactly the same as before.)

(again use backtracking line search). $t \leftarrow$ initial t_0 (usually $t_0 = 1$).

while $f(x_k + td) > f(x_k) + \alpha t \nabla f(x_k)^T d$ do $t \leftarrow \beta t$.

$x_{k+1} \leftarrow x_k + td$. until $\|\nabla f(x_k)\| \leq \delta$? $d^T \nabla^2 f(x_k) d \leq \delta^2$. \checkmark .

Note it is a feasible descent method. since all x_k are feasible.

and $f(x_{k+1}) < f(x_k)$ unless x_k is optimal. (why?).

$$\left. \begin{array}{l} \nabla^2 f(x_k) d - A^T \lambda = -\nabla f(x_k) \\ Ad = 0 \end{array} \right\} \Rightarrow \begin{array}{l} d^T \nabla^2 f(x_k) d - d^T A^T \lambda = -\nabla f(x_k)^T d \\ \nabla^2 f(x_k) \geq 0 \end{array} \xrightarrow{0} \nabla f(x_k)^T d \leq 0.$$

By second-order Taylor, $f(x_k + td) \approx f(x_k) + t \nabla f(x_k)^T d + \frac{t^2}{2} d^T \nabla^2 f(x_k) d$. $-\nabla f(x_k)^T d$

Moreover, $d^T \nabla^2 f(x_k) d \geq \lambda_{\min} d^T d$. so $\frac{t^2}{2} d^T \nabla^2 f(x_k) d + o(t^2 d^T d) > (1-\alpha) t \nabla f(x_k)^T d$
 \forall sufficiently small t .

So if $d^T \nabla^2 f(x) d \leq \delta$. Newton's method halts. o.w. backtracking halts.

If $d^T \nabla^2 f(x_k) d = 0$. $\nabla^2 f(x_k) d - A^T \lambda = -\nabla f(x_k)$. we claim $\nabla^2 f(x_k) d = 0$.

Proof. if $Q \geq 0$. $Q = U^T \Lambda U$. $\Lambda = \text{diag} \{ \xi_1, \dots, \xi_n \}$. $U = \{ u_1, \dots, u_n \}$.

$$d^T Q d = d^T \left(\sum \xi_i u_i u_i^T \right) d = \sum \xi_i (u_i^T d)^T (u_i^T d) = \sum \xi_i (u_i^T d)^2.$$

$$d^T Q d = 0 \Rightarrow \forall i, \xi_i = 0 \text{ or } u_i^T d = 0 \Rightarrow Qd = \sum \xi_i u_i u_i^T d = 0. \quad \square$$

Note $Ax = \begin{pmatrix} a_1^T x \\ \vdots \\ a_m^T x \end{pmatrix}$. $A\lambda = (a_1, \dots, a_m) \lambda = \sum \lambda_i a_i$. So $\nabla f(x_k) = \sum \lambda_i a_i$.

Lagrange condition for $\min f(x)$, s.t. $Ax = b$. so x_k is optimal.

Implementation: Gaussian elimination, assume KKT matrix nonsingular.

① $\circ \ker(Q) \cap \ker(A) = \{0\}$. Q and A have no nontrivial common kernel.

② $\circ Ax = 0, x \neq 0 \Rightarrow x^T Q x > 0$. $Q > 0$ on $\ker(A)$.

③ $\circ F^T Q F > 0$ for $\forall F \in \mathbb{R}^{n \times (n-m)}$ s.t. $\text{Im}(F) = \{ Fv : v \in \mathbb{R}^{n-m} \} = \ker(A)$.

Proof. " \Rightarrow ". if $0 \neq x \in \ker(Q) \cap \ker(A)$. $\begin{pmatrix} Q & A^T \\ A & 0 \end{pmatrix} \begin{pmatrix} x \\ 0 \end{pmatrix} = 0$.

" $1 \Rightarrow 2$ ". $Ax = 0, x \neq 0 \Rightarrow Qx \neq 0$ but $x^T Q x = 0 \Rightarrow Qx = 0$.

" \Leftarrow ". Assume $\begin{pmatrix} Q & A^T \\ A & 0 \end{pmatrix} \begin{pmatrix} v \\ w \end{pmatrix} = 0$. need to show $v = w = 0$.

$$v^T Q v = v^T (-A^T w) = -(Av)^T w = 0. \text{ contradicts if } Av = 0, v \neq 0.$$

$$Qv + A^T w = 0 \Rightarrow A^T w = 0. \text{ contradicts } \text{rank}(A) = m \text{ if } w \neq 0.$$

$$2 \Leftrightarrow 3. \quad \text{Im}(\bar{F}) = \ker(A). \quad 0 \neq y \in \ker(A) \Leftrightarrow \exists x \neq 0. y = \bar{F}x. \quad \square$$

$$\text{Example. } \min f(x) = x_1^2 + x_2^2. \quad \text{s.t. } x_1 + x_2 = 1. \quad \text{start: } (1, 0)^T$$

$$\text{KKT system } \begin{pmatrix} 2 & 0 & 1 \\ 0 & 2 & 1 \\ 1 & 1 & 0 \end{pmatrix} \begin{pmatrix} d_1 \\ d_2 \\ -\lambda \end{pmatrix} = \begin{pmatrix} -2 \\ 0 \\ 0 \end{pmatrix} \Rightarrow (d_1, d_2, \lambda) = (-\frac{1}{2}, \frac{1}{2}, 0)$$

$$\text{set step size } t = 1/2. \quad \text{next iteration: } (1, 0)^T + t(d_1, d_2)^T = \begin{pmatrix} 3/4 \\ 1/4 \end{pmatrix}$$

$$\text{Example. } \min f(x) = x_1^2. \quad \text{s.t. } x_1 + 2x_2 = b. \quad \text{start: } (b, 0)^T$$

$$\text{KKT system } \begin{pmatrix} 2 & 0 & 1 \\ 0 & 0 & 2 \\ 1 & 2 & 0 \end{pmatrix} \begin{pmatrix} d_1 \\ d_2 \\ -\lambda \end{pmatrix} = \begin{pmatrix} -2b \\ 0 \\ 0 \end{pmatrix} \Rightarrow (d_1, d_2, \lambda) = (-b, \frac{b}{2}, 0)$$

$$\ker(\nabla^2 f) = r(0, 1)^T, \quad r \in \mathbb{R}. \quad \ker(A) = t(-2, 1)^T, \quad t \in \mathbb{R}. \quad \bar{F} = (-2, 1)^T$$

Convergence analysis: comparison to unconstrained cases.

$$\text{feasible solution set } X = \{x: Ax = b\} = \{\tilde{x} + \bar{F}z: z \in \mathbb{R}^{n-m}\}$$

$$\min f(x). \quad \text{s.t. } Ax = b \quad \Leftrightarrow \min g(z) = f(\tilde{x} + \bar{F}z)$$

$$\nabla g(z) = \bar{F}^T \nabla f(\bar{F}z + \tilde{x}). \quad \nabla^2 g(z) = \bar{F}^T \nabla^2 f(\bar{F}z + \tilde{x}) \bar{F}. \quad \text{Apply Newton.}$$

$$z_{k+1} = z_k - t(\nabla^2 g(z_k))^{-1} \nabla g(z_k). \quad \text{We show by induction } x_k = \bar{F}z_k + \tilde{x}.$$

$$\text{let } dx_k, dz_k \text{ be the descent direction. s.t. } \begin{matrix} x_{k+1} = x_k + t dx_k \\ z_{k+1} = z_k + t dz_k \end{matrix}$$

$$\exists dz_k \text{ iff } \exists (\nabla^2 g(z_k))^{-1} \text{ iff } \nabla^2 g(z_k) > 0 \text{ iff KKT nonsingular} \Rightarrow \exists dx_k$$

$$\text{Note } Adx_k = 0 \text{ and } \text{Im}(\bar{F}) = \ker(A). \quad \Leftrightarrow \exists u \in \mathbb{R}^{n-m}. \quad dx_k = \bar{F}u.$$

$$\text{Moreover, } \nabla^2 f(x_k) dx_k - A^T \lambda = -\nabla f(x_k), \text{ so } \nabla^2 f(x_k) \bar{F}u - A^T \lambda = -\nabla f(x_k).$$

Multiply (left) by \bar{F}^T on both sides. $\bar{F}^T \nabla^2 f(x_k) \bar{F} u + (A\bar{F})^T \lambda = -\bar{F}^T \nabla f(x_k)$.

$$\ker(A) = \text{Im}(\bar{F}) \Rightarrow A\bar{F} = 0 \text{ (why?)}. \quad \forall x \in \mathbb{R}^{n-m} \quad A\bar{F}x = \vec{0}.$$

$$\text{So } \bar{F}^T \nabla^2 f(x_k) \bar{F} u = -\bar{F}^T \nabla f(x_k) \Rightarrow u = -(\nabla^2 g(z_k))^T \nabla g(z_k) = d_{z_k}.$$

$$d_{x_k} = \bar{F} u = \bar{F} d_{z_k} \Rightarrow x_{k+1} = x_k + t d_{x_k} = \bar{F} z_k + \tilde{x} + t \bar{F} d_{z_k} = \bar{F} z_{k+1} + \tilde{x} \quad \square$$

Recall the convergence result of unconstrained backtracking Newton:

without backtracking: if f is m -strongly convex, $\nabla^2 f$ is M -Lipschitz.

$$\text{then } \|x_{k+1} - x^*\| \leq \frac{M}{2m} \|x_k - x^*\|^2. \quad \nabla f \text{ is } \overset{\uparrow}{M}\text{-smooth.}$$

with backtracking: further assume f is L -smooth. then global convergence.

$$\text{let } \eta = \min\{1, 3(1-2\alpha)\} m^2/M. \quad \gamma = 2\alpha(1-\alpha) \beta \eta^2 m/L.$$

$$\text{if } \|\nabla f(x_k)\| \leq \eta, \quad f(x_k + d_k) \leq f(x_k) + \alpha \nabla f(x_k)^T d. \quad t_k = 1$$

$$\text{then } \|x_{k+1} - x^*\| \leq \frac{M}{2m} \|x_k - x^*\|^2. \quad f(x_k) - f(x^*) \leq \frac{m}{2} \|x_k - x^*\|^2$$

$$\text{if } \|\nabla f(x_k)\| > \eta, \text{ backtracking guarantees } f(x_k) - f(x_{k+1}) \geq \gamma.$$

overall. let k_0 be the smallest k that $\|\nabla f(x_k)\| \leq \eta$.

$$f(x_k) - f(x^*) \leq \begin{cases} f(x_0) - f(x^*) - k\gamma & k \leq k_0 \\ \frac{2m^3}{M^2} (1/2)^{2^{k-k_0+1}} & k > k_0. \end{cases}$$

To get $f(x_k) - f(x^*) \leq \varepsilon$. it suffices to run

$$\frac{1}{\gamma} (f(x_0) - f(x^*)) + \log \log \frac{2m^3}{M^2 \varepsilon} \text{ steps.}$$

Inequality constrained optimization: $\min f(x)$ s.t. $\begin{matrix} g_i(x) = 0 & m \\ h_j(x) \leq 0 & k \end{matrix}$

Again, our first task is to determine whether a given x^* is optimal.

If $h_j(x^*) = 0$, the j -th inequality constraint $h_j(x) \leq 0$ is called active.

otherwise called inactive. ($h_j(x^*) < 0$). Denote $J(x) = \{j : h_j(x) = 0\}$.

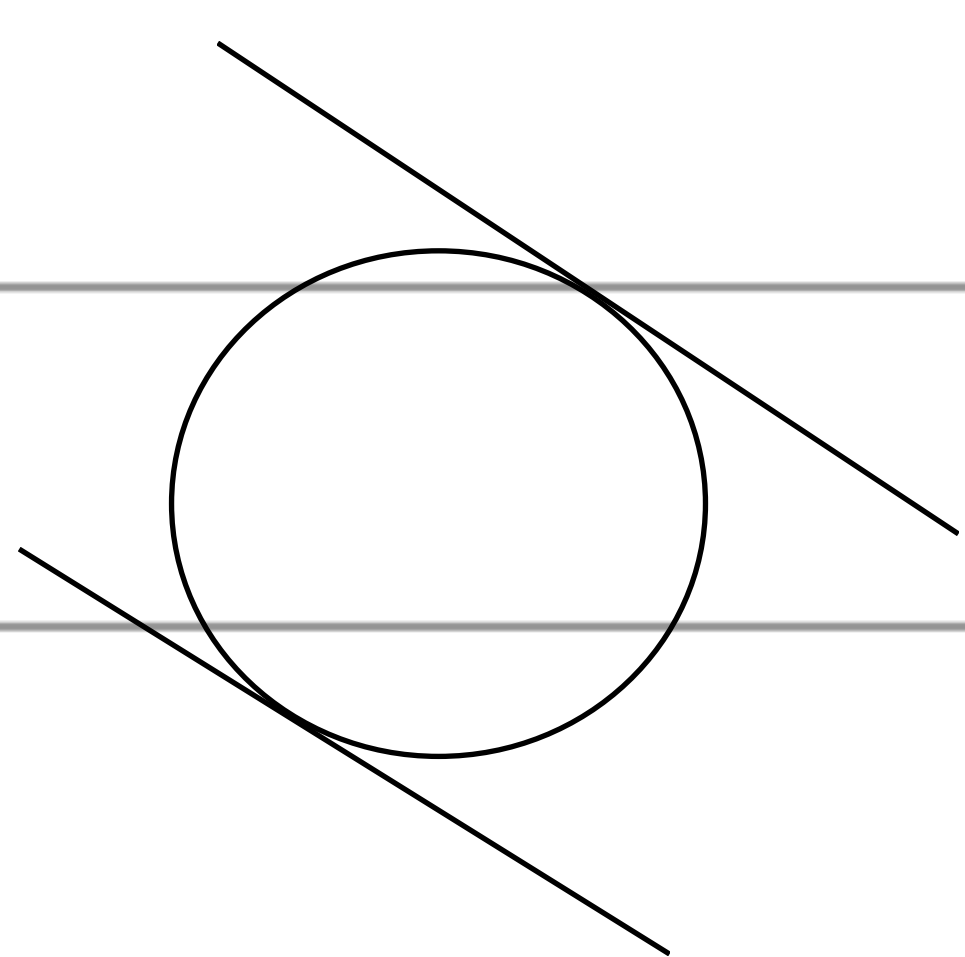
If x^* is a local minimum of original problem, it is also a local

minimum of $\min f(x)$ s.t. $\begin{matrix} g_i(x) = 0, \forall i \\ h_j(x) = 0, \forall j \in J(x^*) \end{matrix}$.

Apply Lagrange condition. $\exists \lambda^*, \mu^*$ s.t. $\nabla f + \sum_{i=1}^m \lambda_i^* \nabla g_i + \sum_{j \in J} \mu_j^* \nabla h_j = 0$

Then we set $\mu_j^* = 0$ for inactive $j \notin J(x^*)$. $\nabla f + \sum_{i=1}^m \lambda_i^* \nabla g_i + \sum_{j=1}^k \mu_j^* \nabla h_j = 0$.

Is this enough? consider the example. $\min x_1 + x_2$ s.t. $x_1^2 + x_2^2 \leq 2$.



$$(1, 1)^T + 2\lambda(x_1, x_2)^T = 0. \quad x_1^2 + x_2^2 = 2 \Rightarrow x_1 = x_2 = 1. \\ \text{or } \lambda = 0$$

If the constraint is equality, then we cannot

distinguish these two cases by first-order condition.

But now..... $-\nabla h_j(x^*)$ is the direction to the interior of the

feasible set, but $-\nabla f(x^*)$ should be the direction to the outside.

$\Rightarrow \mu_j^* > 0$ iff $j \in J(x^*)$. (why this is also true for > 1 constraints?).

Note that $j \in J(x^*)$ iff $h_j(x^*) = 0$. o.w $\mu_j^* = 0$. so. $\mu_j^* h_j(x^*) = 0$.

Now we have the Karush-Kuhn-Tucker (KKT) condition.

Theorem. If x^* is a local minimum point of ICP and also a regular point for active constraints. then \exists Lagrange / KKT multipliers.

$\lambda_1^*, \dots, \lambda_m^*, \mu_1^*, \dots, \mu_k^*$ such that the following KKT holds.

$$1. \nabla f(x^*) + \sum_{i=1}^m \lambda_i^* \nabla g_i(x^*) + \sum_{j=1}^k \mu_j^* \nabla h_j(x^*) = 0.$$

$$2. \mu_j^* \geq 0. \quad \forall j = 1, \dots, k. \quad 3. \mu_j^* h_j(x^*) = 0. \quad \forall j = 1, \dots, k.$$

Comparison to LP. $\min c^T x. \quad \text{s.t.} \quad A_1 x = b_1. \quad A_2 x \leq b_2.$

KKT condition: $\exists \lambda^*, \mu^*. \quad \text{s.t.} \quad c + A_1^T \lambda^* + A_2^T \mu^* = 0.$

$$\mu_j^* \geq 0. \quad \text{and} \quad (A_2 x^* - b)^T \mu^* = 0. \quad (\mu_j^* \cdot (a_{2j}^T x - b_j) = 0).$$

Duality of LP: $\min (y_1, y_2) \cdot (b_1, b_2)^T. \quad \text{s.t.} \quad y_1^T A_1 + y_2^T A_2 = -c^T.$
 $y_2 \geq 0.$

complementary slackness for LP: (x^*, y^*) is optimal

for primal and dual res. iff. $y^T (b - Ax) = 0$ and $x^T (A^T y - c) = 0.$

either $y_i = 0$ or $(Ax)_i = b_i$

$$y_2^T (A_2 x^* - b) = 0.$$

So condition 3 is called complementary slackness for KKT.

KKT is named after Harold W. Kuhn and Albert W. Tucker (1951,

published), but later found in master thesis of William Karush (1939)
unpublished.