

Nonlinear estimation via least-square approximation

- Motivation (slide of previous lecture)
- Introduction:
 - Convex case: $\min_{x \in \mathbb{R}^n} \frac{1}{2} \|Ax - b\|^2$
 - Direct method
 - * Cholesky and QR decomposition for numerical stability
 - Indirect method:
 - * Gradient descent
 - nonconvex case: $\min_{x \in \mathbb{R}^n} \|r(x)\|^2$
 - Gauss - Newton method
 - Levenberg - Marquardt
 - Nonconvex ^{manifold} case: $\min_{x \in M} \|r(x)\|^2$

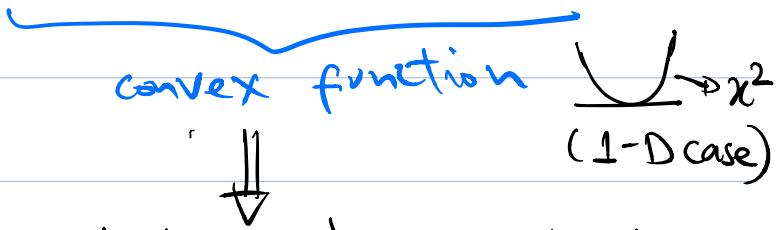
Convex case

$$\min_{x \in \mathbb{R}^n} \frac{1}{2} \|Ax - b\|^2 \quad (A, b \text{ are given})$$

• Direct method :

Observe that

$$\frac{1}{2} \|Ax - b\|^2 = \frac{1}{2} x^T A^T A x - x^T A^T b + \frac{1}{2} \|b\|^2$$



minimized when gradient = 0

gradient



$$\nabla \left\{ \frac{1}{2} \|Ax - b\|^2 \right\} = A^T A x_{\text{opt}} - A^T b = 0$$

$$\Rightarrow A^T A x_{\text{opt}} = A^T b$$

Assuming $A^T A$ is invertible:

$$x_{\text{opt}} = (A^T A)^{-1} A^T b.$$

CAUTION: Inverting $(A^T A)^{-1}$ is:

- computationally expensive
- possibly numerically unstable.

↓

Instead, solve system of equations:

$$A^T A x_{\text{opt}} = A^T b$$

Efficient methods to this end are

- Cholesky
- QR

(Gaussian elimination is avoided, since it is

more computationally demanding, and numerically unstable)

Cholesky: The method is based on the observation that equations of the following form are solved easily:

$$L y = b, \text{ where } L, b \text{ are given}$$

and, in particular:

$$L = \begin{bmatrix} l_{11} & 0 & 0 & 0 & \dots & 0 \\ l_{21} & l_{22} & 0 & 0 & \dots & 0 \\ l_{31} & l_{32} & l_{33} & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \ddots & \vdots \\ l_{n1} & l_{n2} & l_{n3} & \dots & \dots & l_{nn} \end{bmatrix} \quad : \text{lower-triangular}$$

Example: for $n=2$:

$$Ly = b \Rightarrow \begin{pmatrix} l_{11} & 0 \\ l_{21} & l_{22} \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}$$

$$\Rightarrow l_{11} y_1 = b_1 \Rightarrow y_1 = \frac{b_1}{l_{11}} \text{ (easy!)} \quad \boxed{\quad}$$

\Rightarrow given y_1 , now y_2 is also easy!

So, given $\underbrace{A^T A x}_{\triangleq \bar{A}} = \underbrace{A^T b}_{\triangleq \bar{b}}$,

Cholesky method finds L such that

$$\bar{A} = L L^T \quad \xrightarrow{\text{lower trnng}}$$

$$\Rightarrow \bar{A} x = \bar{b}$$

$$\Rightarrow L \underbrace{(L^T x)}_{\triangleq y} = \bar{b}$$

\Rightarrow first solve $L y = \bar{b}$,
 then solve $L^T x = y$,
 which is also easy for the same
 reason (L^T is upper triangular)

QR : This method has similar intuition: given

$$A^T A x = A^T b,$$

the QR method finds

- $Q \in \mathbb{R}^{m \times n}$ such that $Q^T Q = I_n$
 - $R \in \mathbb{R}^{n \times n}$ is upper triangular
- such that

$$A = QR.$$

$$\Rightarrow A^T A x = A^T b$$

$$\Rightarrow R^T \underbrace{Q^T Q R}_{=I} x = \underbrace{R^T Q^T b}_{\Delta b}$$

$$\Rightarrow \boxed{R^T} R x = \boxed{b} \Rightarrow \text{easy to find } y!$$

lower triangular

Once y is known, remains to solve:

$$R x = y$$

which is also easy (R is upper trinag)

Cholesky vs QR

Cholesky : • faster than QR

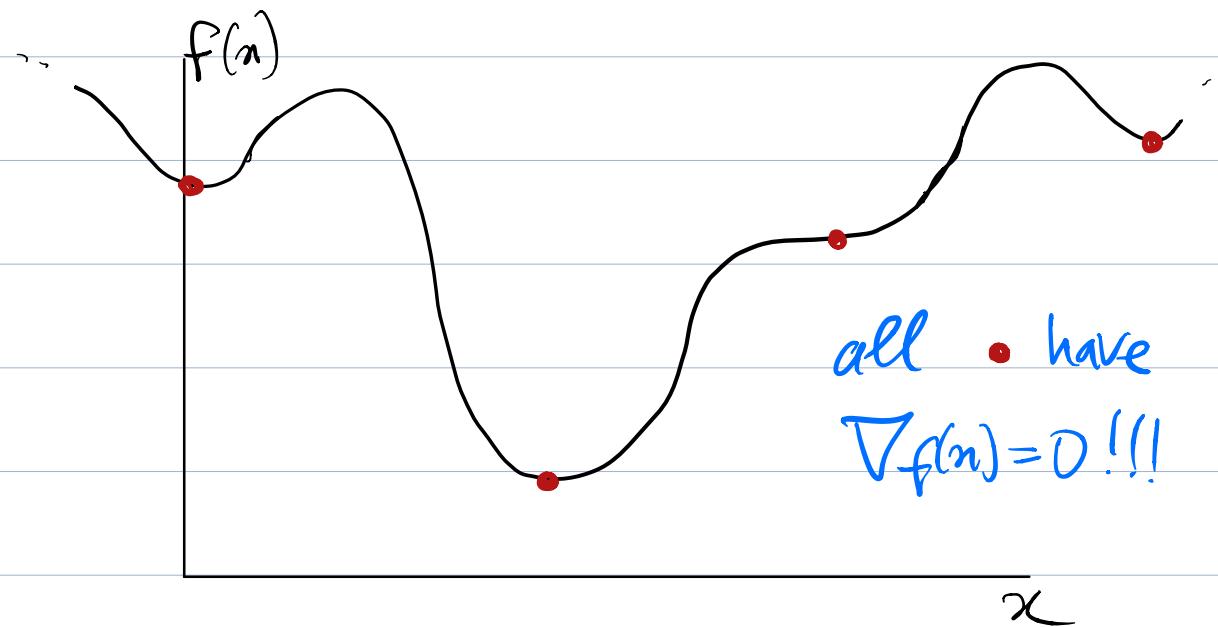
- valid only if $A^T A$ is positive-DEFINITE

QR : • slower than Cholesky

- valid even if $A^T A$ is positive-SEMIDEFINITE

Indirect method: Gradient descent.

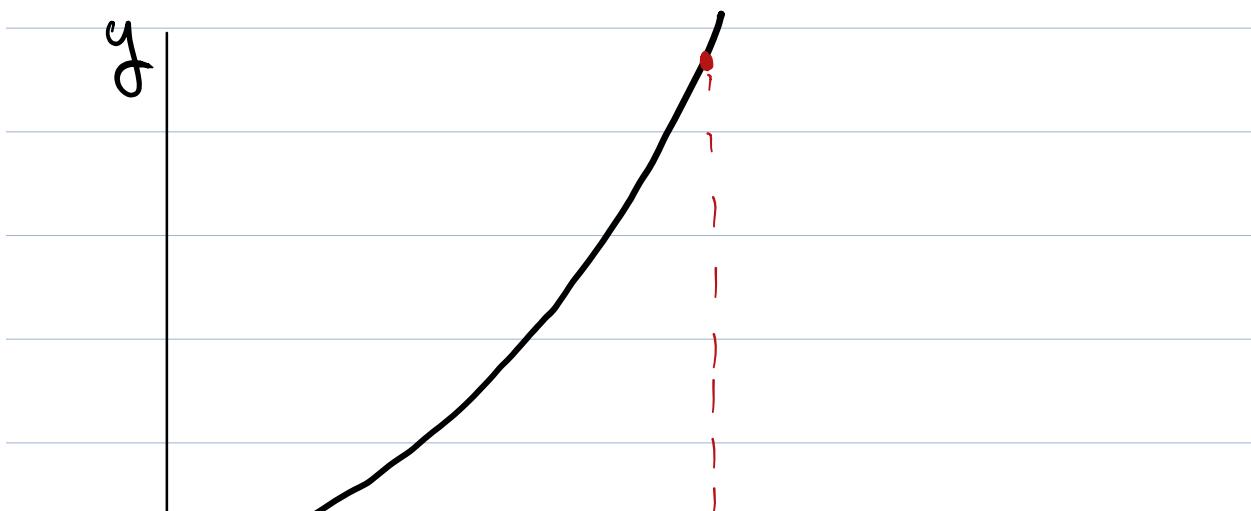
When our objective function is convex, as in the linear case above, we were able to find a closed form solution. But when the function will be nonconvex, that will not be easy.



Therefore, we will need (as a potential approach) to develop iterative approximation methods.

Such a method is gradient descent, illustrated next, for simplicity, for a simple convex case.

Consider $r(x) = x^2$ (scalar case)
and that want to $\min_x \|r(x)\|^2$





Question:

If I stand at x_0 , can we decide at which x_1 should we stand next, such that $\|r(x_1)\|^2 \leq \|r(x_0)\|^2$? More generally, can we choose, sequentially, x_i , $i=1, 2, \dots, k$ such that $\|r(x_i)\|^2 \rightarrow \min_n \|r(x)\|^2$ for $k \rightarrow \infty$?

Answer: Yes!

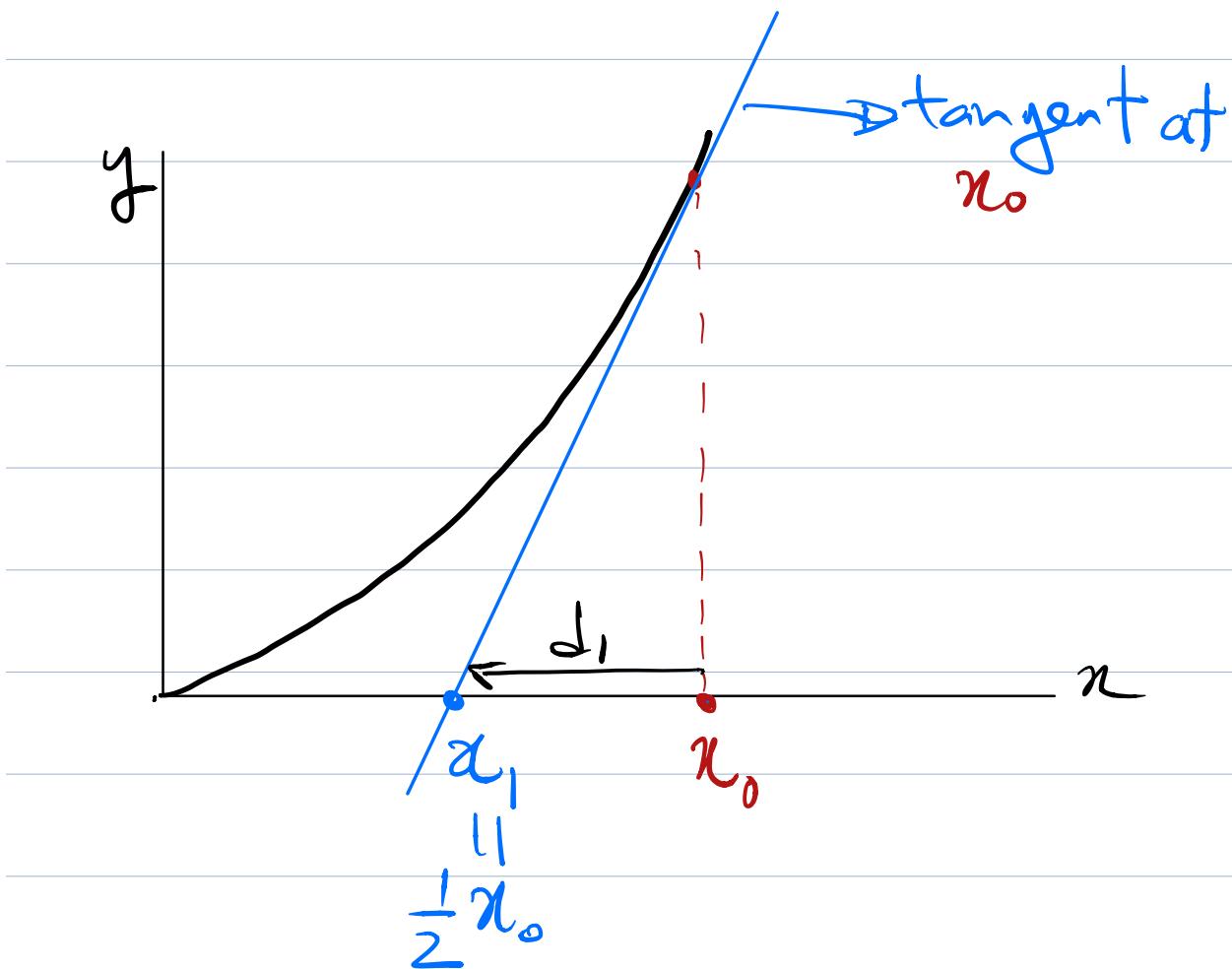
Instead of solving $\min_x \|r(x)\|^2$ approximate $r(x)$ with Taylor expansion at $x=x_0$, i.e., $r(x_0+d) \approx r(x_0) + r'(x_0)d$,

$$= x_0^2 + 2x_0 d$$

and minimize:

$$\min_{d \in \mathbb{R}^n} \|\mathbf{x}_0^2 + 2\mathbf{x}_0 d\|^2$$

$$\Rightarrow d_1 = -\frac{1}{2} \mathbf{x}_0$$

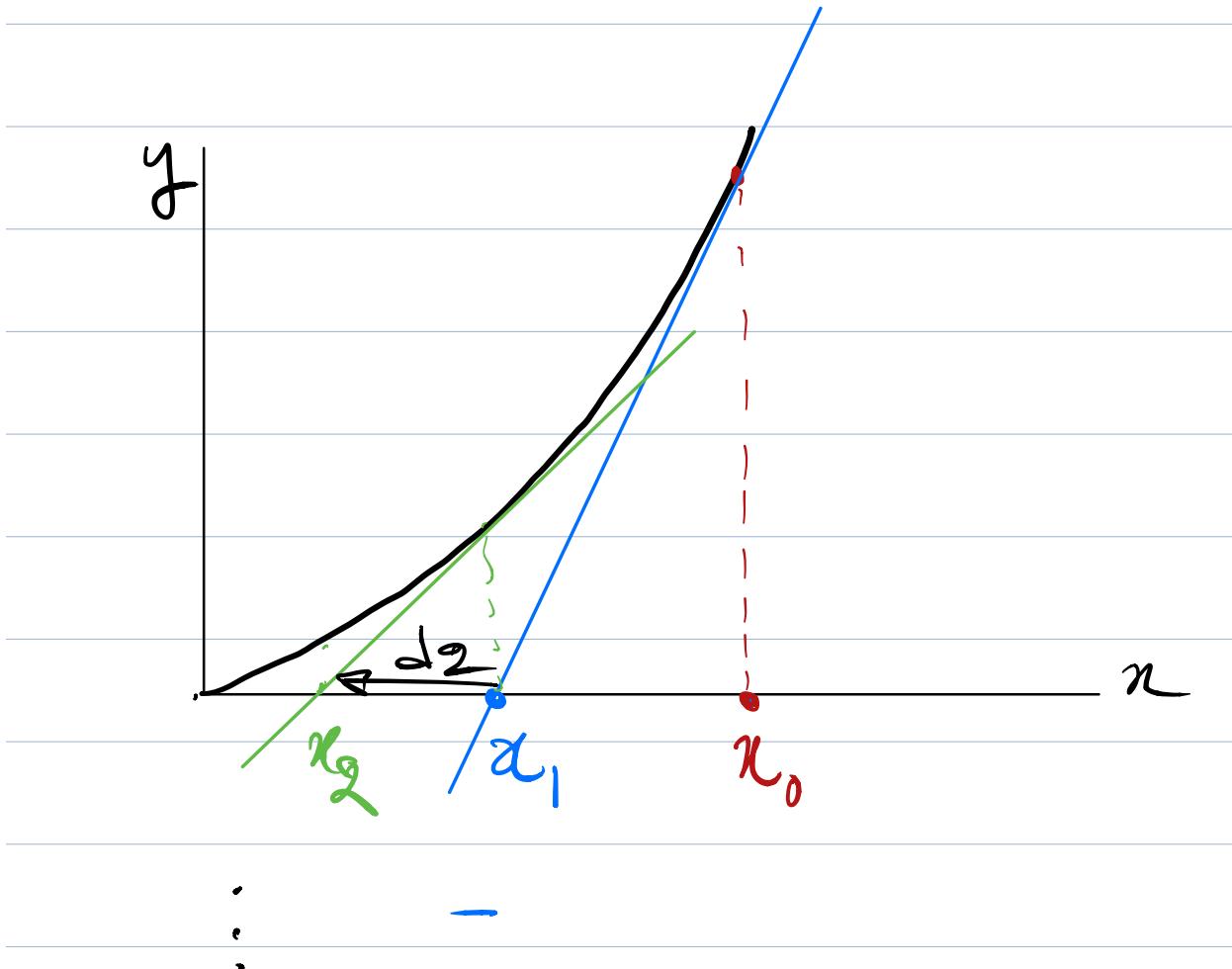


Repeating

$$\Rightarrow d_2 = -\frac{1}{2} \mathbf{x}_1 = -\frac{1}{4} \mathbf{x}_0 \Rightarrow$$

from \mathbf{x}_1

$$x_2 = x_1 + d_2 = \frac{1}{4} x_0 = \frac{1}{2^2} x_0$$



$$-d_k = -\frac{1}{2^k} x_0, d_k = \frac{1}{2^k} x_0$$

and from the figure it's clear
 that $\lim_{k \rightarrow \infty} \min_n \|r(x)\|^2 = 0$

(since $x_k, d_k \rightarrow 0$).

All in all, gradient descent refers to the aforementioned iterative procedure. It becomes relevant next.

Non convex case (via iterative method)

• Gauss-Newton method

$$f(x) = \frac{1}{2} \|r(x)\|^2,$$

where:

$$\mathbf{r} : \mathbb{R}^n \rightarrow \mathbb{R}^m \quad (m \geq n)$$

\mathbf{r} is smooth, but not necessarily affine (i.e., $\mathbf{A}\mathbf{x} + \mathbf{b}$)

$$\|\mathbf{r}(\mathbf{x})\|^2 = \sum_{i=1}^m r_i^2(\mathbf{x}) \text{ where } r_i : \mathbb{R}^n \rightarrow \mathbb{R}$$

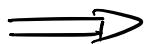
first-order Taylor:

$$r_i(\mathbf{x}) \approx r_i(\mathbf{x}_0) + \nabla r_i(\mathbf{x}_0)^\top (\mathbf{x} - \mathbf{x}_0)$$

stack r_i 's:

$$\mathbf{r}(\mathbf{x}) \approx \mathbf{r}(\mathbf{x}_0) + \mathbf{J}(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0)$$

Jacobian



$$\mathbf{r}(\mathbf{x}_0 + \mathbf{d}) \approx \mathbf{r}(\mathbf{x}_0) + \mathbf{J}(\mathbf{x}_0)\mathbf{d}$$

where $\mathbf{J}(\mathbf{x}) \triangleq \frac{\partial \mathbf{r}(\mathbf{x})}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial r_1}{\partial x_1} & \frac{\partial r_1}{\partial x_2} & \dots & \frac{\partial r_1}{\partial x_n} \\ \frac{\partial r_2}{\partial x_1} & \frac{\partial r_2}{\partial x_2} & \dots & \frac{\partial r_2}{\partial x_n} \\ \vdots & \vdots & \dots & \vdots \\ \frac{\partial r_m}{\partial x_1} & \frac{\partial r_m}{\partial x_2} & \dots & \frac{\partial r_m}{\partial x_n} \end{bmatrix} \in \mathbb{R}^{m \times n}$.



Gauss-Newton method

- 1 start from an initial guess \mathbf{x}^0

for $k = 0, 1, \dots$ and until "convergence":

- 2 linearize the residual at the current guess \mathbf{x}^k :

$$\mathbf{r}(\mathbf{x}^k + \mathbf{d}) \approx \mathbf{r}(\mathbf{x}^k) + \mathbf{J}(\mathbf{x}^k)\mathbf{d}$$

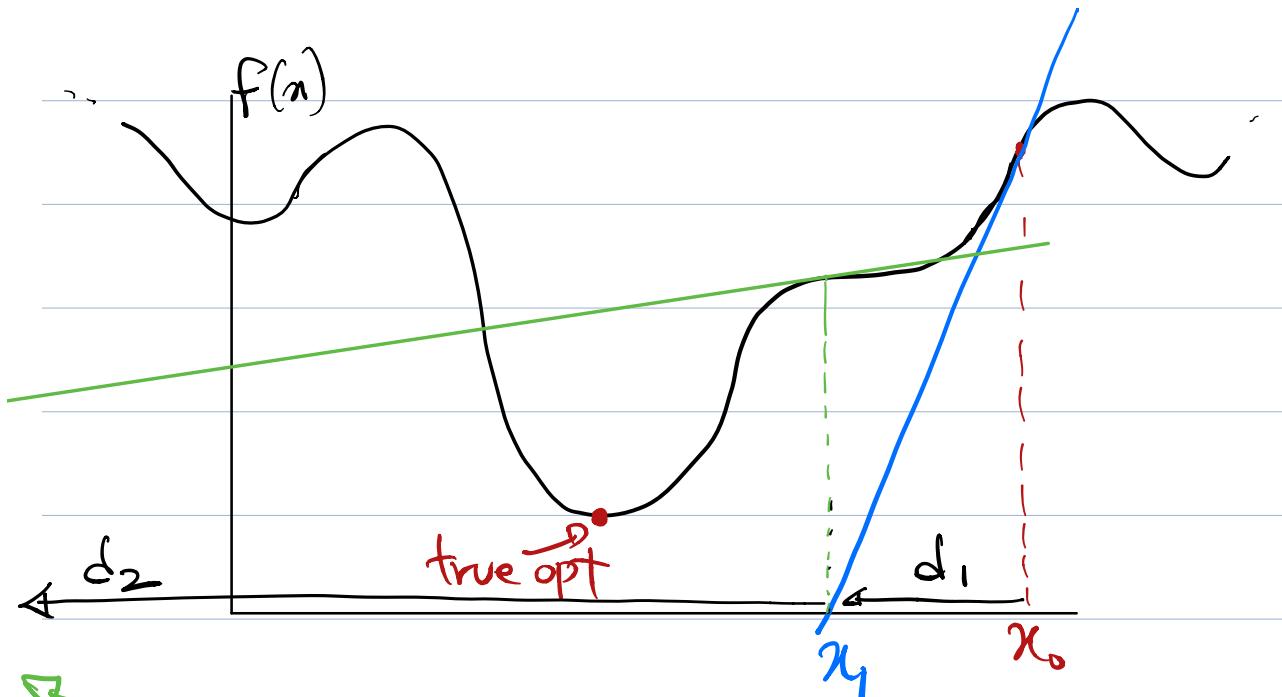
- 3 solve the resulting linear least squares to find the step \mathbf{d} :

$$\underset{\mathbf{d}}{\text{minimize}} \quad \|\mathbf{r}(\mathbf{x}^k) + \mathbf{J}(\mathbf{x}^k)\mathbf{d}\|^2$$

$$\Rightarrow (\mathbf{J}_k^\top \mathbf{J}_k)\mathbf{d} = -\mathbf{J}_k^\top \mathbf{r}(\mathbf{x}^k)$$

- 4 $\mathbf{x}^{k+1} = \mathbf{x}^k + \mathbf{d}$.

Issues with Gauss-Newton method



x_2 is far away $\Rightarrow d_2$ is very large!

Idea: penalize $\|d\|^2$, so d doesn't explode.



Levenberg - Marquardt method

Instead of (per Gauss-Newton):

$$\min_d \|r(x^k) + J(x^k)d\|^2$$

do:

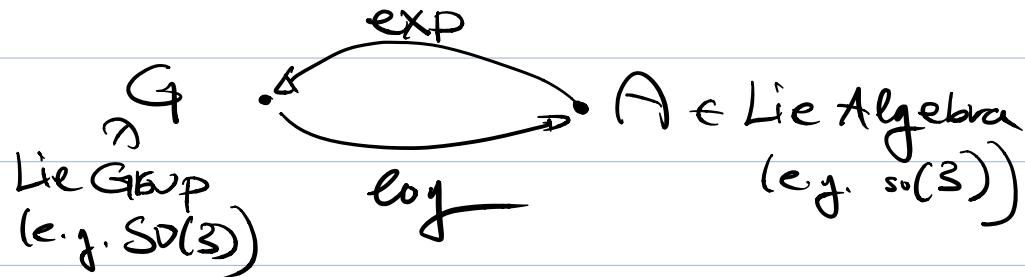
$$\min_d \|r(x^k) + J(x^k)d\|^2 + \lambda \|d\|^2$$

for a given $\lambda > 0$.

$$\Rightarrow [J^T(x^k)J(x^k) + \lambda I]d = -J^T(x^k)r(x^k)$$

Optimization over Manifolds

Review of Lie Groups / Algebra



where:

$$\exp(A) \triangleq \mathbf{I} + \sum_{k=1}^{\infty} \frac{A^k}{k!}$$

- $\exp(\mathbf{0}) = \mathbf{I}$
- in matrix Lie groups, \exp maps Lie algebra (i.e., $\mathfrak{se}(3)$ and $\mathfrak{so}(3)$) to Lie group (i.e., $\text{SO}(3)$ and $\text{SE}(3)$)
- $\mathfrak{se}(3)$ and $\mathfrak{so}(3)$ are vector spaces \rightarrow basis "vectors" (a.k.a. generators)

$$\hat{\phi} \in \mathfrak{so}(3) \Leftrightarrow \hat{\phi} = \phi_1 \mathbf{G}_1 + \phi_2 \mathbf{G}_2 + \phi_3 \mathbf{G}_3$$

where $\boxed{\phi \in \mathbb{R}^3}$ and

$$\mathbf{G}_1 = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & -1 \\ 0 & 1 & 0 \end{bmatrix} \quad \mathbf{G}_2 = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 0 & 0 \\ -1 & 0 & 0 \end{bmatrix} \quad \mathbf{G}_3 = \begin{bmatrix} 0 & -1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

- $\hat{\phi} = [\phi]_{\times} \Rightarrow \hat{\phi} \mathbf{a} = \phi \times \mathbf{a}$

Similarly, for $\mathfrak{se}(3)$ consider $\phi \in \mathbb{R}^3$ and $\rho \in \mathbb{R}^3$ and the overloaded hat operator:

$$\begin{bmatrix} \widehat{\phi} \\ \rho \end{bmatrix} \in \mathfrak{se}(3) \Leftrightarrow \begin{bmatrix} \widehat{\phi} \\ \rho \end{bmatrix} = \phi_1 \mathbf{G}_1 + \phi_2 \mathbf{G}_2 + \phi_3 \mathbf{G}_3 + \rho_1 \mathbf{G}_4 + \rho_2 \mathbf{G}_5 + \rho_3 \mathbf{G}_6$$

where

$$\begin{aligned} \mathbf{G}_1 &= \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & -1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} & \mathbf{G}_2 &= \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \\ -1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} & \mathbf{G}_3 &= \begin{bmatrix} 0 & -1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \\ \mathbf{G}_4 &= \begin{bmatrix} 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} & \mathbf{G}_5 &= \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} & \mathbf{G}_6 &= \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{bmatrix} \end{aligned}$$

Optimization over Manifolds

$$\min_{\mathbf{z} \in M} \frac{1}{2} \|\mathbf{r}(\mathbf{z})\|^2$$

Example. Bundle adjustment:

$$\min_{(\mathbf{R}_{c_i}^w, \mathbf{t}_{c_i}^w), i=1,2,3} \sum_{k=1}^N \sum_{i=1}^3 \|\mathbf{x}_{k,i} - \pi(\mathbf{R}_{c_i}^w, \mathbf{t}_{c_i}^w, \mathbf{p}_k^w)\|^2$$

where $p_k^w \in \mathbb{R}^3$ but $(R_{c_i}^w, t_{c_i}^w) \in SE(3)$.

Question Is Gauss-Newton (or any of the iterative methods we explored) valid anymore?

Answer NO!

$$\cancel{x^{k+1} = x^k + d}$$

since adding rotation doesn't make sense; multiplying does.

Problem to solve: How to define gradient descent steps " d_{\parallel} over manifolds?

Solution: Define " \mathbf{d}_k " in the Lie Algebra and then map it to the group!

$$\mathbf{x}^{k+1} = \mathbf{x}^k \exp(\hat{\mathbf{d}}),$$

where: \mathbf{d} in Lie Algebra



- Gauss-Newton over \mathbb{R}^n

$$\mathbf{r}(\mathbf{x}^k + \mathbf{d}) \approx \mathbf{r}(\mathbf{x}^k) + \mathbf{J}_k \mathbf{d}$$

$$\mathbf{J}_k = \frac{\partial \mathbf{r}(\mathbf{x})}{\partial \mathbf{x}} \bigg|_{\mathbf{x}=\mathbf{x}^k} = \frac{\partial \mathbf{r}(\mathbf{x}^k + \mathbf{d})}{\partial \mathbf{d}} \bigg|_{\mathbf{d}=\mathbf{0}}$$

- Gauss-Newton over $\text{SO}(3)$ — $\boxed{\mathbf{d} \in \mathbb{R}^3}$

$$\mathbf{r}(\mathbf{x}^k \exp(\hat{\mathbf{d}})) \approx \mathbf{r}(\mathbf{x}^k) + \mathbf{J}_k \mathbf{d}$$

$$\mathbf{J}_k \triangleq \frac{\partial \mathbf{r}(\mathbf{x}^k \exp(\hat{\mathbf{d}}))}{\partial \mathbf{d}} \bigg|_{\mathbf{d}=\mathbf{0}}$$

- Gauss-Newton over $\text{SE}(3)$ — $\boxed{\mathbf{d} \in \mathbb{R}^6}$

$$\mathbf{r}(\mathbf{x}^k \exp(\hat{\mathbf{d}})) \approx \mathbf{r}(\mathbf{x}^k) + \mathbf{J}_k \mathbf{d}$$

$$\mathbf{J}_k \triangleq \frac{\partial \mathbf{r}(\mathbf{x}^k \exp(\hat{\mathbf{d}}))}{\partial \mathbf{d}} \bigg|_{\mathbf{d}=\mathbf{0}}$$



Lift-Solve-Retract

perturbation:

$$\mathbf{x}^{k+1} = \mathbf{x}^k \exp(\hat{\mathbf{d}})$$

① lift:

$$g : \mathbb{R}^{n_d} \rightarrow \mathbb{R}^m : \mathbf{d} \mapsto \mathbf{r}(\mathbf{x}^k \exp(\hat{\mathbf{d}}))$$

e.g., $n_d = 3$ in $\text{SO}(3)$ and $n_d = 6$ in $\text{SE}(3)$

$$g(\mathbf{d}) \approx g(\mathbf{0}) + \frac{\partial g(\mathbf{d})}{\partial \mathbf{d}} \Big|_{\mathbf{d}=\mathbf{0}} \mathbf{d} \quad \text{Taylor at } \mathbf{d} = \mathbf{0}$$

$$\mathbf{r}(\mathbf{x}^k \exp(\hat{\mathbf{d}})) \approx \mathbf{r}(\mathbf{x}^k) + \mathbf{J}_k \mathbf{d}$$

② solve:

$$\underset{\mathbf{d}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{r}(\mathbf{x}^k \exp(\hat{\mathbf{d}}))\|^2 \approx \frac{1}{2} \|\mathbf{r}(\mathbf{x}^k) + \mathbf{J}_k \mathbf{d}\|^2$$

linear least squares \Rightarrow normal equations

$$\mathbf{d} = -(\mathbf{J}_k^\top \mathbf{J}_k)^{-1} \mathbf{J}_k^\top \mathbf{r}(\mathbf{x}^k)$$

③ retract:

$$\mathbf{x}^{k+1} = \mathbf{x}^k \exp(\hat{\mathbf{d}})$$

Tips to compute \mathbf{J}_k : in $\mathbf{r}(\mathbf{x}^k \exp(\hat{\mathbf{d}}))$

observe that:

- For $\|\mathbf{d}\| \approx 0$:

$$\exp(\hat{\mathbf{d}}) \approx \mathbf{I} + \hat{\mathbf{d}}$$

- express $\hat{\mathbf{d}} = \sum_i d_i \mathbf{G}_i$ and take derivatives w.r.t. each d_i (i.e., columns of \mathbf{J}_k)

Example

- Consider $\|\mathbf{r}(\mathbf{x}_1, \mathbf{x}_2)\|^2$ where $\mathbf{x}_1 \in \mathbb{R}^3$ and $\mathbf{x}_2 \in \text{SO}(3)$

$$\|\mathbf{r}(\mathbf{x}_1^k + \mathbf{d}_1, \mathbf{x}_2^k \exp(\widehat{\mathbf{d}}_2))\|^2 \approx \|\mathbf{r}(\mathbf{x}_1^k, \mathbf{x}_2^k) + \mathbf{J}_{1,k} \mathbf{d}_1 + \mathbf{J}_{2,k} \mathbf{d}_2\|^2$$

$$\mathbf{J}_{1,k} \triangleq \frac{\partial \mathbf{r}(\mathbf{x})}{\partial \mathbf{x}_1} \bigg|_{\mathbf{x}=(\mathbf{x}_1^k, \mathbf{x}_2^k)} = \frac{\partial \mathbf{r}(\mathbf{x}_1^k + \mathbf{d}_1, \mathbf{x}_2^k)}{\partial \mathbf{d}_1} \bigg|_{\mathbf{d}_1=\mathbf{0}}$$

$$\mathbf{J}_{2,k} \triangleq \frac{\partial \mathbf{r}(\mathbf{x}_1^k, \mathbf{x}_2^k \exp(\widehat{\mathbf{d}}_2))}{\partial \mathbf{d}_2} \bigg|_{\mathbf{d}_2=\mathbf{0}}$$

- solve the resulting linear least squares
 - retract: $\mathbf{x}_1^{k+1} = \mathbf{x}_1^k + \mathbf{d}_1$ and $\mathbf{x}_2^{k+1} = \mathbf{x}_2^k \exp(\widehat{\mathbf{d}}_2)$
-
-
-
-
-
-
-
-
-