Review:

Constrained Least Squares:

- $\tilde{y}_1 = H_1 x + V \longrightarrow$ measurements to fit
- $\tilde{y}_2 = H_2 x \longrightarrow$ constraint.

$$H_1 = \left[ \; \right] \in \mathbb{R}^{m_1 \times n} \qquad H_2 = \left[ \underline{\quad} \right] \in \mathbb{R}^{m_2 \times n}$$

$$\underset{\text{tall}}{} \qquad \qquad \underset{\text{fat.}}{}$$

$$m_1 \geq n \qquad\qquad\qquad m_2 < n$$

$$\min_{\hat{x}} \; J = \tfrac{1}{2} e_1^T W_1 e_1 = \tfrac{1}{2} (\tilde{y}_1 - H_1 \hat{x})^T W_1 (\tilde{y}_1 - H_1 \hat{x})$$

$$\text{s.t.} \quad \tilde{y}_2 = H_2 \hat{x}$$

optimal soln:

$$\bar{x} = (H_1^T W_1 H_1)^{-1} H_1^T W_1 \tilde{y}_1 \qquad \begin{array}{c}\text{unconstrained} \\ \text{solution}\end{array}$$

$$\hat{x} = \underline{\bar{x}} + \underline{K} (\underline{\tilde{y}_2} - \underline{H_2 \bar{x}})$$

optimal gain

actual meas

prediction of $\tilde{y}_2$ from $\bar{x}$ ← unconst soln

diff. between the meas. & prediction

where

$$K = (H_1^T W_1 H_1)^{-1} H_2^T \left( H_2 (H_1^T W_1 H_1)^{-1} H_2^T \right)^{-1}$$

# SEQUENTIAL (BATCH) LS ESTIMATION:

$$\tilde{y}_1 = H_1 x + V_1 \quad \} \rightarrow \text{first batch} \quad H_1 \in \mathbb{R}^{m_1 \times n} \quad m_1 \geq n$$

$$\tilde{y}_2 = H_2 x + V_2 \quad \} \rightarrow \text{second batch} \quad H_2 \in \mathbb{R}^{m_2 \times n}$$

Don't have all
data initially...

Comes in two
batches...

$m_2$ can be any size but
in practice $m_2$ will be small

Note: $m_2$ is small $\rightarrow$ gain
computational
advantage

first: $\hat{x}_1 = \left( H_1^T W_1 H_1 \right)^{-1} H_1^T W_1 \tilde{y}_1$

$\overset{\text{initial}}{\underset{\text{estimate}}{\bigcirc}} \rightarrow$ want to take advantage
of initial estimate.

now we add batch 2...
want to solve for $\hat{x}_2 \leftarrow$ best fit for
all data.

$$\begin{bmatrix} \tilde{y}_1 \\ \tilde{y}_2 \end{bmatrix} = \underbrace{\begin{bmatrix} H_1 \\ H_2 \end{bmatrix}}_{H} x + \begin{bmatrix} V_1 \\ V_2 \end{bmatrix} \qquad W = \begin{bmatrix} W_1 & 0 \\ 0 & W_2 \end{bmatrix}$$

$$\hat{x}_2 = \left( H^T W H \right)^{-1} H^T W \tilde{y}$$

$(A+B)^{-1} \cancel{=} A^{-1} + B^{-1}$

$\frac{1}{x+y} \cancel{=} \frac{1}{x} + \frac{1}{y}$

$$= \left( \begin{bmatrix} H_1^T H_2^T \end{bmatrix} \begin{bmatrix} W_1 & 0 \\ 0 & W_2 \end{bmatrix} \begin{bmatrix} H_1 \\ H_2 \end{bmatrix} \right)^{-1} \begin{bmatrix} H_1^T H_2^T \end{bmatrix} \begin{bmatrix} W_1 & 0 \\ 0 & W_2 \end{bmatrix} \begin{bmatrix} \tilde{y}_1 \\ \tilde{y}_2 \end{bmatrix}$$

$$= \left( \underline{H_1^T W_1 H_1} + \underline{H_2^T W_2 H_2} \right)^{-1} \begin{bmatrix} H_1^T W_1 \tilde{y}_1 + H_2^T W_2 \tilde{y}_2 \end{bmatrix} \leftarrow$$

matrix inverses
are computationally expensive
and slower for large
operations.

Define:

$$P_1 = \left( H_1^T W_1 H_1 \right)^{-1}$$

$$P_2 = \left( H_1^T W_1 H_1 + H_2^T W_2 H_2 \right)^{-1}$$

$$P_1^{-1} = P_2^{-1} - H_2^T W_2 H_2$$

$$\hat{x}_1 = P_1 H_1^T W_1 \tilde{y}_1 \leftarrow$$

$$\Rightarrow P_1^{-1} \hat{x}_1 = H_1^T W_1 \tilde{y}_1$$

$$\circledast \Rightarrow \left( P_2^{-1} - H_2^T W_2 H_2 \right) \hat{x}_1 = H_1^T W_1 \tilde{y}_1$$

$$P_2^{-1} = \underbrace{\left( H_1^T W_1 H_1 + H_2^T W_2 H_2 \right)}_{P_1^{-1}}$$

$$\circledast \quad P_1^{-1} = P_2^{-1} - H_2^T W_2 H_2$$

$$\hat{x}_2 = P_2 \left( H_1^T W_1 \tilde{y}_1 + H_2^T W_2 \widehat{y}_2 \right)$$

$$= P_2 \left( \left( P_2^{-1} - H_2^T W_2 H_2 \right) \hat{x}_1 + H_2^T W_2 \tilde{y}_2 \right)$$

$$= \hat{x}_1 \ominus \underline{P_2 H_2^T W_2 H_2 \hat{x}_1} \oplus \underline{P_2 H_2^T W_2 \tilde{y}_2}$$

$$\hat{x}_2 = \underset{\substack{\text{initial} \\ \text{est.}}}{\underline{\hat{x}_1}} + \underset{\substack{K_2 \\ \text{gain}}}{\underline{P_2 H_2^T W_2}} \left( \underset{\substack{\text{actual} \\ \text{new} \\ \text{meas}}}{\tilde{y}_2} - \underset{\substack{\text{prediction for} \\ \text{new meas.,} \\ \text{based on init est.}}}{\underline{H_2 \hat{x}_1}} \right)$$

$$\underbrace{\qquad\qquad\qquad\qquad}_{\text{diff between pred. \& meas}}$$

$$K_2 = P_2 H_2^T W_2$$

General Rule:

$$\boxed{\begin{array}{l} \hat{x}_{k+1} = \hat{x}_k + K_{k+1} \left( \tilde{y}_{k+1} - H_{k+1} \hat{x}_k \right) \\[2mm] \text{where, } K_{k+1} = P_{k+1} H_{k+1}^T W_{k+1} \\[2mm] \qquad\quad P_{k+1}^{-1} = P_k^{-1} + H_{k+1}^T W_{k+1} H_{k+1} \leftarrow \end{array}}$$

still need to compute.

$$P_{k+1} = \left(P_k^{-1} + H_{k+1}^T W_{k+1} H_{k+1}\right)^{-1}$$

$\chi \in \mathbb{R}^n$

$\rightarrow$ low rank
if only a few new meas.

$\mathbb{R}^{n\times n}$

$\mathbb{R}^{n\times n}$

$\begin{bmatrix} H_{k+1}^T \end{bmatrix} \begin{bmatrix} W_{k+1} \end{bmatrix} \begin{bmatrix} H_{k+1} \end{bmatrix}$

still inverting an $n\times n$ matrix

scalar

$\begin{bmatrix} \; \rceil \lceil \; \rfloor \lceil \; \end{bmatrix}$ $\leftarrow$ if only 1 new meas

$\begin{bmatrix} \tilde{y}_k \\ \tilde{y}_{k+1} \end{bmatrix} = \begin{bmatrix} H_k \\ H_{k+1} \end{bmatrix} \chi$

Woodbury Matrix Identity:
Matrix Inversion Lemma.
(Rank 1): Sherman Morrison Formula
worth memorizing.

Wikipedia:

$$\left(\underline{A} + U\underline{CV}\right)^{-1} = A^{-1} - A^{-1}U\left(\underline{C^{-1} + VA^{-1}U}\right)^{-1}VA^{-1}$$

$\uparrow$ $\downarrow$ low rank $\uparrow$

$$\left(\lceil A \rfloor + \lceil U \rfloor [C][V]\right)^{-1} = \begin{bmatrix} A \end{bmatrix}^{-1} - \left|\begin{bmatrix} \left(\lceil C \rfloor^{-1} + [V]\right) \end{bmatrix}\right|^{-1} \lceil VA^{-1} \rfloor$$

trading inverting a big matrix (blue)

$A^{-1}U$ (blue)

assuming we've previously computed $A^{-1}$ (blue)

small  $A^{-1}U$

for inverting 2 small matrices (blue)

TO CHECK:

$$(A+UCV)(A+UCV)^{-1} = (A+UCV)\left(A^{-1} - A^{-1}U\left(C^{-1}+VA^{-1}U\right)^{-1}VA^{-1}\right)$$

(see the book)
$$= I$$

$(A+B)^{-1}$ : natural Searle Identities    $A + ucv^T$

full rank    (low rank)

$|||^\cdot \equiv$

GOING BACK TO LS...

$$P_{k+1} = \left( P_k^{-1} + H_{k+1}^T W_{k+1} H_{k+1} \right)^{-1}$$

applying Woodbury MI

$$= P_k - P_k H_{k+1}^T \left( W_{k+1}^{-1} + H_{k+1} P_k H_{k+1}^T \right)^{-1} H_{k+1} P_k$$

small.

$$K_{k+1} = P_{k+1} H_{k+1}^T W_{k+1}$$

$$= \left( P_k - P_k H_{k+1}^T \left( W_{k+1}^{-1} + H_{k+1} P_k H_{k+1}^T \right)^{-1} H_{k+1} P_k \right) H_{k+1}^T W_{k+1}$$

$$= P_k H_{k+1}^T \left( I - \left( W_{k+1}^{-1} + H_{k+1} P_k H_{k+1}^T \right)^{-1} H_{k+1} P_k H_{k+1}^T \right) W_{k+1}$$

$$= P_k H_{k+1}^T \left( W_{k+1}^{-1} + H_{k+1} P_k H_{k+1}^T \right)^{-1} \left( W_{k+1}^{-1} + H_{k+1} P_k H_{k+1}^T - H_{k+1} P_k H_{k+1}^T \right) W_{k+1}$$

$$= P_k H_{k+1}^T \left( W_{k+1}^{-1} + H_{k+1} P_k H_{k+1}^T \right)^{-1}$$

$\underbrace{\qquad}_{I}$

$$\boxed{\begin{aligned} K_{k+1} &= P_k H_{k+1}^T \left( W_{k+1}^{-1} + H_{k+1} P_k H_{k+1}^T \right)^{-1} \\ P_{k+1} &= P_k - K_{k+1} H_{k+1}^T P_k = \left( I - K_{k+1} H_{k+1} \right) P_k \end{aligned}}$$

Nonlinear Least Squares: (iteratively)

$$\tilde{y} = f(x) + \nu \qquad \text{want to estimate } x \to \hat{x}$$

$$\tilde{y} = \underbrace{f(\hat{x})}_{\hat{y}} \qquad e = \tilde{y} - \hat{y} = \Delta y$$

$$\min_{\hat{x}} \quad J = \tfrac{1}{2}\Delta y^T W \Delta y = \tfrac{1}{2}(\tilde{y} - f(\hat{x}))^T W (\tilde{y} - f(\hat{x}))$$

Linearization: $x_c \to$ current estimate

$$\hat{x} = x_c + \Delta x \qquad \text{(DO LS TO SOLVE FOR } \Delta x$$

$\to$ initialize $x_c$

$\to$ iteratively adjust $x_c = x_c + \Delta x$

$$f(\hat{x}) \approx f(x_c) + H \Delta x$$

$\to \hat{x} = x_c$ at end.   solved for using LS

where $H = \dfrac{\partial f}{\partial x}\Big|_{x_c}$

(Jacobian at $x_c$)

$$\Delta y = \tilde{y} - f(\hat{x}) \approx \tilde{y} - f(x_c) - H\Delta x \qquad \Delta y = \Delta y_c - H\Delta x$$

meas $\underset{\text{diff}}{\overset{\Delta y_c}{\longrightarrow}}$ predicted meas

$$\min_{\Delta x} \tfrac{1}{2}\Delta y^T W \Delta y = \tfrac{1}{2}(\Delta y_c - H\Delta x)^T W(\Delta y_c - H\Delta x)$$

$\underset{``\tilde{y}"}{\underline{\quad}}\ \underset{``H"}{\underline{\quad}}\ \underset{``\hat{x}"}{\underline{\quad}}$     $\dfrac{\partial f}{\partial x}\big|_{x_c}$

| Model $f(x)$ |

$x_c$ at $i=0$

$\dfrac{\partial f}{\partial x}$

$$\Delta y_c = \tilde{y} - f(x_c)$$
$$J_i = \Delta y^T W \Delta y \quad \longleftarrow \tilde{y}$$
$$H = \dfrac{\partial f}{\partial x}\Big|_{x_c} \quad \longleftarrow W$$

$i = i+1$

max iter?

$$\Delta x = (H^T W H)^{-1} H^T W \Delta y_c$$

$J$

$$x_c = x_c + \Delta x \longleftarrow \delta J < \dfrac{\varepsilon}{|w|}? \overset{\text{yes}}{\longrightarrow} \text{stop}$$

Caveats:

• $f(x) \to$ differentiable

• $x_c$ needs to start "close" to $x$.

↓ what does this mean?
$\Rightarrow$ depends

• local optima.

→ newton's method ??

# Basis functions:

$$H = \begin{bmatrix} h_1(t_1) & \cdots & h_n(t_1) \\ \vdots & & \vdots \\ h_1(t_m) & \cdots & h_n(t_m) \end{bmatrix} \updownarrow m$$

$\overleftarrow{n}\rightarrow$

$m > n$

$h_i(t_j)$ : basis functions

$\uparrow$ parameter index

basis functions
- functions of "time"
- one parameter per basis function

$$y(t) = \sum_i h_i(t) x_i \quad : \text{output.}$$

## Polynomial functions in $t$:

$h_0(t) = 1$. $h_1(t) = t$, $h_2(t) = t^2$, etc... powers of $t$

$$H = \begin{bmatrix} 1 & t_1 & t_1^2 & \cdots & t_1^n \\ 1 & t_m & t_m^2 & \cdots & t_m^n \end{bmatrix} \rightarrow \text{Vandermonde Matrix}$$

$$y(t) = \sum_{i=0}^n t^i x_i$$

## Sinusoidal functions:
want to fit a periodic signal

$h_j^1(t) = \cos(j\omega t)$    $j = 0, 1, \ldots, n$

$h_j^2(t) = \sin(j\omega t)$  for $j = 0, 1, \ldots, n$

$$y(t) = \sum_{j=0}^n x_j^1 \cos(j\omega t) + \sum_{j=0}^n x_j^2 \sin(j\omega t)$$

$$H = \begin{bmatrix} \cos(0) & \cos(\omega t_1) & \cos(2\omega t_1) & \cdots & \sin(0) & \sin(\omega t_1) & \cdots \\ \cos(0) & \cos(\omega t_2) & \cos(2\omega t_2) & \cdots & '' & \sin(\omega t_2) & \cdots \end{bmatrix} \begin{bmatrix} x_0^1 \\ x_1^1 \\ x_2^1 \\ \vdots \\ x_0^2 \\ x_1^2 \\ \vdots \end{bmatrix}$$
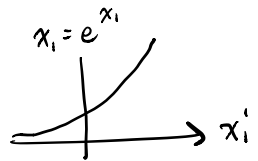
# Nonlinear coord transform:

$$y(t) = \underbrace{\boxed{x_1 e}}_{x_2 t} \xrightarrow{\ln(\cdot)} y'(t) = \ln(y(t)) = \underbrace{\ln(x_1)}_{x_1'} + \underbrace{x_2 t}_{x_2' t}$$

not a lin
func of $x_i$...

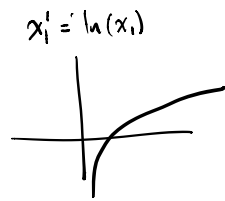$$\begin{bmatrix} \ln(y(t_1)) \\ \vdots \\ \ln(y(t_m)) \end{bmatrix} \begin{matrix} y'(t_1) \\ \vdots \\ y'(t_m) \end{matrix} = \begin{bmatrix} 1 & t_1 \\ \vdots & \vdots \\ 1 & t_m \end{bmatrix} \boxed{\begin{bmatrix} x_1' \\ x_2' \end{bmatrix}}$$
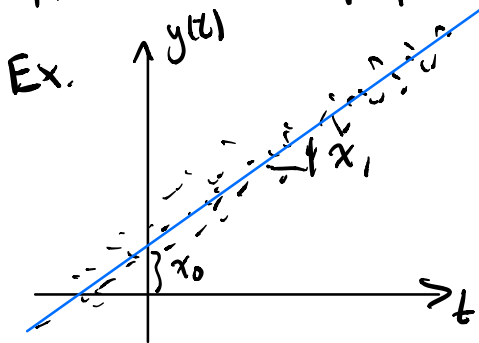
$$x_1' = \ln(x_1)$$
$$x_2' = x_2$$

$y = x$

$\xrightarrow{\text{sym}}$

$x_1 = e^{x_1'}$

$$\begin{bmatrix} \hat{x}_1' \\ \hat{x}_2' \end{bmatrix} \longrightarrow \begin{matrix} x_1 = e^{x_1'} \\ x_2 = x_2' \\ x_1' : \text{anything} \\ \Rightarrow x_1 > 0 \end{matrix}$$

$x_1 = e^{x_1'}$

$\longrightarrow x_1'$

$x_1' = \ln(x_1)$

# How to set up problems:

Ex.

$y(t)$

$H = \begin{bmatrix} 1 & t_1 \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \end{bmatrix}$

$x_1$

$x_0$

$\rightarrow t$

Ex

$\rightarrow t$

lin model $\otimes$

$\omega \to$ small

$$\begin{bmatrix} 1 & \cos(\omega t_1) & \cos(2\omega t_1) \cdots & -\sin(\omega t) & \rightarrow \end{bmatrix}$$

higher freq.

low freq

more complicated curves

$\ulcorner H \urcorner$ gets wider

$x \downarrow$ longer

to fit more complicated curves

$\Rightarrow$ need more data

want H to be tall

tall H

overfitting

$\times$ square H ?

one parameter per data point.

2D SURFACE

$$y = H x$$

diff basis



$y(t_1)$

$y(k_m)$

$t^1, t^2$

$h_j(t^1, t^2)$
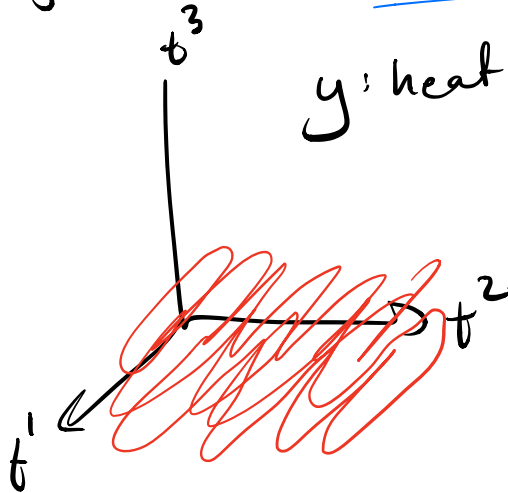
$x_1$

$x_n$

$y$

$t^2$

$t^1$

$t^1, t^2$



$$h_j(t^1, t^2) = t^1 t^2 \rightarrow$$

$$\dots \quad h_j(t^1, t^2) = \cos(\omega t^1) \sin(\omega t^2)$$

$$h_j(t^1, t^2) = \cos(\omega t^1)$$

$y:$ heat map

$t^3$

$t^2$

$t^1$

# MINIMUM VARIANCE (perspective on LS)

Estimator theory $\quad W = R^{-1}$ $\qquad$ Covariance of the noise

$-\tilde{y} = Hx + v \qquad v \sim N(0, R)$

compute an estimator $\qquad\qquad$ Gaussian mean

function $\hat{x}(\tilde{y})$

mean: $E[\hat{x}] = x$

covariance

$E[(\hat{x}-x)(\hat{x}-x)^T]$

$[\hat{x}_i - x_i]\,[\hat{x}_i - x_i \cdots]$

$\begin{bmatrix} \hat{x}_i - x_i \\ \vdots \\ \hat{x}_n - x_n \end{bmatrix}$

- assume some class of functions for $\hat{x}(\cdot)$
- try to find the best estimator within that class

Properties of a good estimator

$0 = E\left((\hat{x}_i - x_i)(\hat{x}_i - x_i)\right)$

- unbiased : $\boxed{E(\hat{x}(\tilde{y})) = x} \quad\Leftarrow$

- bias : $E(\hat{x}(\tilde{y}) - x)$

Expected Value $E(\cdot)$

$\rightarrow$ always over some prob. dist.

$E_{\tilde{y}}\,\hat{x}(\tilde{y})$

random variable

What is the minimum variance $\Leftarrow$ linear estimator ?

Linear Estimator

$\hat{x} = M\tilde{y} + n \qquad$ pick $M, n$

unbiased

$E\hat{x} = x \implies E(M\tilde{y} + n) = E(MHx + n + Mv) = x$

$$E(M\hat{y}+n) = E(MHx) + E(n) + E(Mv)\underset{0}{\underline{\phantom{xx}}} = x$$

2 conditions to impose

$\boxed{?}$ $M = H^{-1}$? $\quad$ MH=I

- $n = 0$
- $MH = I$

$(M + z)H = \overset{\nearrow}{M}H + \overset{0}{z\cancel{H}}$

if $z \in$ left
nullspace
of $H$

unbiased estimator

$$\hat{x} = M\tilde{y} \quad s.t. \quad MH = I$$

Minimum Variance:

$$MH = I$$

$$\underset{M}{min} \; E\left((\hat{x}-x)^T(\hat{x}-x)\right) = J$$

$$n\begin{bmatrix} c \end{bmatrix}\begin{bmatrix} \overset{n}{H} \end{bmatrix}$$

$$s.t. \quad MH = I$$

$\dfrac{\partial J}{\partial M} = ?$ $\quad$ need to
be invertible
(many C's)

need to solve this
optimization problem.

$$M = (CH)^{-1}C$$

for a lot of C's

— trace operator

$$MH = (CH)^{-1}CH = I$$

— matrix inner products

could choose $C = H^T...$

— matrix derivatives

$$M = (H^TH)^{-1}H^T$$

Formulas

what does
this mean?

$$\frac{\partial}{\partial X} Tr(BXC) = \boxed{B^TC^T}$$

$$Tr(BAC) = Tr(CBA)$$

$$\frac{\partial}{\partial X} Tr(XBX^T) = X(B+B^T)$$

$$f(X) = Tr(A^TX)$$
$$\frac{\partial f}{\partial X} = A$$

Trace Operator: $A \in \mathbb{R}^{n \times n}$

$Tr(A) = \sum_i A_{ii}$   $\left( Tr(A) = \sum_i \lambda_i \ , \ \lambda_i \in \gamma(A) \right)$

$\overline{Tr(AB)} = Tr(BA) \ \otimes$ in general $AB \neq BA$

interesting case... $\big\}$

$\underset{\text{scalar}}{x^T y} = \underset{\downarrow}{Tr(x^T y)} = Tr(\underset{\text{matrix}}{y x^T}) = \underline{Tr} \left( \begin{bmatrix} y_1 x_1 & & y_1 x_n \\ & & \\ y_n x_1 & & y_n x_n \end{bmatrix} \right)$

$\boxed{x_1 y_1 + \cdots + x_n y_n}$

Matrix Inner Products:   $A = [A_1 \cdots A_n]$   $B = [B_1 \cdots B_n]$

$x^T y = \sum_i x_i y_i$   matrix $\langle A, B \rangle = \sum_{i,j} A_{ij} B_{ij} = \boxed{\sum_i A_i^T B_i}$

$\boxed{\langle A, B \rangle = Tr(A^T B)} = Tr \left( \begin{bmatrix} -A_1^T- \\ -A_n^T- \end{bmatrix} \begin{bmatrix} | & & | \\ B_1 & \cdots & B_n \\ | & & | \end{bmatrix} \right)$

Matrix Derivatives   $= Tr \left( \begin{bmatrix} A_1^T B_1 & & \\ & \ddots & \\ & & A_n^T B_n \end{bmatrix} \right)$

$X \in \mathbb{R}^{n \times n}$   $f : \mathbb{R}^{n \times n} \longrightarrow \mathbb{R}$

$\qquad\qquad X \longmapsto \mathbb{R}$

$\dfrac{\partial f}{\partial X} = ?$

$\nwarrow$ variable is a matrix

Recall:

$\boxed{\bullet \ f(x) = c^T x}$ $\boxed{\dfrac{\partial f}{\partial x} = c^T}$

$\bullet \ f(x) = Ax$   $\dfrac{\partial f}{\partial x} = A$

variable is a vector...

$$f(X) = \langle A, X \rangle = \text{Tr}(A^T X)$$

perturbation analysis...

$$\Delta f = \text{Tr}(A^T \underline{\Delta X}) \rightarrow \qquad \frac{\partial f}{\partial X} \not= A^T \; ?$$

2 options:

- vectorize $\Delta X \rightarrow$ put $\Delta X$ in vector form ... $\frac{\partial f}{\partial \text{vec}(X)} = \text{matrix}$.

  stacking columns

$$X = [X_1 \cdots X_n] \Rightarrow \text{vec}(X) = \begin{bmatrix} X_1 \\ \vdots \\ X_n \end{bmatrix}$$

Wikipedia: Vectorization

$$f(X) = \text{Tr}(A^T X) = \text{vec}(A)^T \text{vec}(X) = [A_{11} \, A_{21} \cdots A_{n1} \, A_{12} \, A_{22} \cdots] \begin{bmatrix} X_{11} \\ X_{21} \\ \vdots \\ X_{n1} \\ \vdots \end{bmatrix}$$

$$\underline{\langle A, X \rangle} = \sum_{ij} A_{ij} X_{ij}$$

$$\frac{\partial f}{\partial \text{vec}(X)} = \text{vec}(A)^T$$

Similar to $\frac{\partial f}{\partial x} = c^T$
when $f(x) = c^T x$

- $\frac{\partial f}{\partial X}(\Delta X) = \Delta f$

  linear function

$$\frac{\partial f}{\partial X}(\cdot) = \langle F, \cdot \rangle = \text{Tr}(F^T \cdot)$$

Question is what is $F$?

$$\text{if } f(X) = Tr(A^T X) \implies F = A$$

$$\frac{\partial f}{\partial X}(\cdot) = Tr(A^T \cdot)$$

$$f(X) = Tr(A^T X) \implies \frac{\partial f}{\partial X} = A$$

$$f(x) = \underset{C \cdot x}{\underline{C^T x}} \quad \frac{\partial f}{\partial x} = C^T$$

$$f(X) = Tr(A^T X) \quad \frac{\partial f}{\partial X} = A$$
$$\underline{A \cdot X} \quad di$$

$$\boxed{\begin{array}{c} f(X) = Tr(A^T X) \\ \frac{\partial f}{\partial X} = A \end{array}}$$

**WRONG**

$$\boxed{\underline{\Delta f} = \frac{\partial f}{\partial X} \underline{\Delta X}} \quad \oslash$$

$$= A \Delta X$$

$$scalar = \overline{\underline{matrix}}$$

$$f(x) = Tr(A^T X) \quad \frac{\partial f}{\partial X} = A$$

$$f : \mathbb{R}^{n \times n} \longrightarrow \mathbb{R}$$

$$\Delta X \qquad \Delta f$$

$$\boxed{\underline{\Delta f} = Tr\left( \frac{\partial f}{\partial X}^T \Delta X \right)} \quad TRUE.$$

$$\frac{\partial f}{\partial X} = A \qquad \Delta f = Tr\left( A^T \Delta X \right)$$