**Recap.** $\nabla f(x)$

$$f(x) = \tfrac{1}{2}\|Ax - y\|^2$$

$$\leadsto \nabla f(x) = A^T(Ax - y)$$

$$F: \mathbb{R}^p \to \mathbb{R}^q \qquad \partial F(x) \in \mathbb{R}^{q\times p}$$

$$\underbrace{\mathbb{R}^p \to \mathbb{R}^q}_{\text{linear}}$$

$$q=1 \quad \nabla f(x) = \partial F(x)^T$$

$$\boxed{\quad\quad}^T$$

chain rule:

$$\partial(f \circ g) = \partial f \times \partial g$$

Ex. $f(x) = \ell(Bx)$

$\nearrow f(x + \varepsilon\delta) + f(x) + z - \ldots$

$\searrow$ chain rule $f = \ell \circ B$

$$\nabla f(x) = B^T \cdot \nabla \ell(Bx)$$

$$\nabla f = B^T \circ \nabla \ell \circ B$$

Exercise sheet 1.3  ⚠️

$$\|x\|_p = \left(\Sigma |x_i|^p\right)^{1/p}$$

Fi. $\|x\|_2 = \sqrt{x_1^2 + x_2^2}$

$\|x + \varepsilon\delta\|_p$

$\searrow \|x\|_p = (g(x))^{1/p}$

1.4

$$\langle X, Y \rangle_{\mathbb{R}^{n\times p}} = \sum_{ij} X_{ij} Y_{ij}$$

$$\langle X.\text{flatten}, Y.\text{flatten} \rangle_{\mathbb{R}^{np}}$$

$$\langle X, Y \rangle_{\mathbb{R}^{n\times p}} = \text{trace}(XY^T) = \text{tr}(YX^T) = \text{tr}(X^TY)$$

$$X, Y \in \mathbb{R}^{n\times p} \quad XY^T \in \mathbb{R}^{n\times n}$$

$$Z \in \mathbb{R}^{n\times n} \quad \text{tr}(Z) = \Sigma Z_{ii} \qquad X^TY \in \mathbb{R}^{p\times p}$$

$$①(XY^T)_{ij} = \sum_k X_{ik} Y_{jk}$$

$\qquad {}^{n\times n}$

$$\text{tr}(XY^T) = \sum_i (XY^T)_{ii}$$

$$= \sum_i \sum_k X_{ik} Y_{ik}$$

$$= \sum_{ik} X_{ik} Y_{ik} = \langle X, Y \rangle$$

Application: $\langle AB, C \rangle_{\mathbb{R}^{n\times p}} = \langle B, A^TC \rangle$
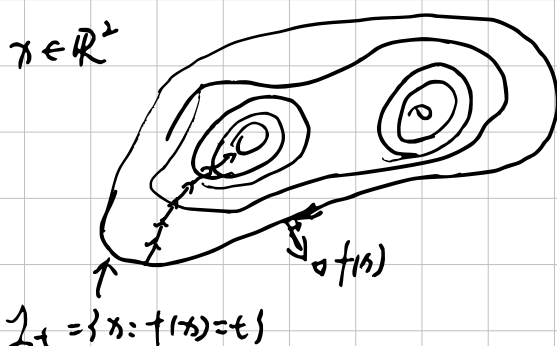
$$\langle Ab, c \rangle = \langle b, A^Tc \rangle$$

Batch Gradient descent :

$$\min_{x \in \mathbb{R}^d} f(x)$$

$x_0 \leftarrow$ init

$x_{k+1} = x_k - \tau_k \nabla f(x_k)$
  $\hookrightarrow$ learning rate $> 0$  ??

$x \in \mathbb{R}^2$



$\qquad \qquad \downarrow \nabla f(x)$

$\mathcal{L}_t = \{ x : f(x) = t \}$

$\begin{matrix} \tau_k \to 0 \\ \| \\ t \end{matrix}$ $\qquad \dfrac{x_{k+1} - x_k}{\tau_k} = \nabla f(x_k)$

$\qquad \quad \Big\Downarrow \, {\scriptstyle \tau \to 0}$

$$x(t) = -\nabla f(x(t))$$
$$\underset{\dfrac{dx}{dt}}{\overset{\|}{}} \quad \uparrow \text{ Gradient Flow ordinarity diff. Eq. (ODE)}$$

Generic "descent" algorithm

$x_{k+1} = x_k - \tau_k d_k$
$f(x_{k+1}) = f(x_k) - \tau_k \langle d_k, \nabla f(x_k) \rangle + o(\tau_k)$
$\Delta_k = f(x_{k+1}) - f(x_k) \leq 0$   ($< 0$ if $x_k$ is not minimum)
$\Delta_k = -\tau_k \langle d_k, \nabla f(x_k) \rangle + o(\tau_k)$
For $\Delta_k < 0$, for $\tau_k$ small enough
$\iff \langle d_k, \nabla f(x_k) \rangle > 0$

$\curvearrowright$ either $\nabla f(x_k) = 0$ $\boxed{\text{stop}}$
$\hookrightarrow$ or $\nabla f(x_k) \neq 0$



$d_k = + \nabla f(x_k)$
  $\hookrightarrow \Delta_k = -\tau_k \| \nabla f(x_k) \|^2 + o(\tau_k)$

$d_k = H_k \cdot \nabla f(x_k)$
  $\hookrightarrow$ sym & positive matrix
  $H_k = U \cdot \text{diag}(\lambda) \, U^T$
      $\quad {\scriptstyle > 0}$ eigenvalues
  $H_k = [\partial^2 f(x_k)]^{-1}$   Newton method
      $\uparrow$
      convex

$\langle d_k, \nabla f(x_k) \rangle = \langle H_k \nabla f(x_k), \nabla f(x_k) \rangle$
$\qquad = \langle U \text{diag}(\lambda) U^T \nabla f(x_k), \nabla f(x_k) \rangle$
$\qquad = \langle \text{diag}(\sqrt{\lambda}) U^T \nabla f(x_k), \text{diag} \sqrt{\lambda} U^T \nabla f(x_k) \rangle$
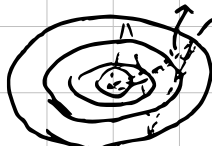$\qquad = \| \text{diag}(\sqrt{\lambda}) U^T \nabla f(x_k) \|^2 > 0$

$f(x) = a x_1^2 + b x_2^2$   $x^* = 0$

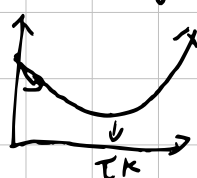$a = b = 1$ $\qquad \qquad a \ll b$ $\tau$ small (slow)
$\qquad \qquad \qquad \qquad \qquad \hookrightarrow \tau$ large (slow?)



$h(\tau) = f(x_k - \tau \nabla f(x_k))$   $h'(0) = -\| \nabla f(x_k) \|^2$
        $\underset{>0}{\uparrow}$



$\qquad \qquad \qquad$ Armijo-Goldstein rule

numerical tours (website)

**Theory** $f(x) = \frac{1}{2}|Ax-y|^2 + \frac{\lambda}{2}|x|^2$

$f(x) = \frac{1}{2}\langle Cx, x\rangle - \underbrace{\langle b, b\rangle}_{\text{linear}}$
$\qquad\qquad \underbrace{\phantom{\langle Cx, x\rangle}}_{\text{quad}}$

$\quad\rightarrow\quad C = A^T A + \alpha \, Id$

$\qquad\qquad b = A^T y$

**Ridge** $f(x) = \frac{1}{2}|Ax - y|^2 + \frac{\lambda}{2}|x|^2$

$\qquad \nabla f(x) = A^T(Ax - y) + \lambda x$

**General:** $f(x) = \frac{1}{2}\langle Cx, x\rangle - \langle b, x\rangle$

$\qquad\qquad \nabla f(x) = Cx - b$

If $\ker(A) = \ker(C) = \{0\}$ (overdetermined)

$\boxed{A}$ then $x^* = C^{-1}b = (A^T A + \lambda Id)^{-1}(A^T y)$

**Operator norm / $\infty$-norm / algebra norm**

$\quad \|H\|_{op}, \ \|H\|$

Def: $\|H\|_{op} = \sup_{x \neq 0} \dfrac{\|Hx\|}{\|x\|}$

$\quad \|Hx\|_2 \leq \|H\|_{op} \cdot \|x\|_2$

$\|H\|_{op}$ is the lipsch constant of

$\qquad\qquad x \longmapsto Hx$
$\qquad\qquad \ell^2 \qquad\quad \ell^2$

$\quad \|AB\|_{op} \leq \|A\|_{op} \cdot \|B\|_{op}$

Prop: $\|H\|_{op} = \sup_i \underbrace{\sqrt{\underbrace{\lambda_i(H^T H)}_{\text{eigenvalue}}}}_{\text{singleton value}}$

If $H$ is symmetry, $\|H\|_{op} = \sup_i |\lambda_i(H)|$

Proof: $\|Hx\|^2 = \langle Hx, Hx\rangle = \langle \underbrace{H^T H x}_{\text{sym} \geq 0}, x\rangle$

$\left( H^T H = U \, \text{diag}(\lambda) U^T \right)$

$\qquad\rightarrow\qquad = \langle U\,\text{diag}(\lambda) U^T x, x\rangle$

$\qquad\qquad = \langle \text{dig}(\lambda)\underbrace{U^T x}_{Z}, \underbrace{U^T x}_{Z}\rangle$

$\qquad\qquad = \sum_i \lambda_i Z_i^2$

$\qquad\qquad \leq \max(\lambda) \underbrace{\sum_j Z_j^2}_{\|U^T x\|^2 = \|x\|^2}$

$\|Hx\|^2 \leq \max(\lambda)\|u\|^2$

$\|Hx\| \leq \sqrt{\max \lambda_i(H^T H)}\,\|x\|$

$U = (u_1, v_2 \ \dots \ u_n)$

$\qquad \lambda_1, \geq \lambda_2 \ \dots \ \geq \lambda_n$

$x = u_1$

$\|H u_1\| = \sqrt{\max(\lambda)}\,\|u_1\|$

$\qquad = \|H\|_{op}$

Theory of GD → strongly cvx

$C$ is invertible  __FAST__

over determined

↘ $C$ might be non - inv   __slow__

under-determined

"Nice" case :   $C$ is inv

$$0 < \lambda_1(C) \leq \lambda_2(C) \leq \cdots \leq \lambda_d(C)$$

$\underset{\mu}{\smile}$         $L = \|C\|_{op}$

$K = \frac{L}{\mu} \geq 1$   conditioning.

$K = 1$   $\lambda_1 = \lambda_2 = \cdots = \lambda_d$   $C = \lambda Id$

$K \gg 1$

$f(x) = \frac{1}{2} \langle Cx, x \rangle - \langle b, x \rangle$

$\nabla f(x) = Cx - b$

$0 = Cx^* - b$     $x^* = C^{-1} b$

Fixed LR  $\tau_k = \tau$   $x_{k+1} = x_k - \tau \nabla f(x_k)$

$$\left\{ \begin{array}{l} x_{k+1} = x_k - \tau \; ( Cx_k - b ) \quad \Delta_k = x_k - x^* \\ x^* = x^* - \tau \; ( \underbrace{(Cx^* - b)}_{= 0} ) \end{array} \right.$$

$\Delta_{k+1} = \Delta_k - \tau C \Delta_k$

$\Delta_{k+1} = (Id - \tau C) \Delta_k$

$\Delta_k = (Id - \tau C)^k \Delta_0$

$\| \Delta_k \| \leq \underbrace{\| Id - \tau C \|^k}_{\rho} \underbrace{\| \Delta_0 \|}_{\| x_0 - x^* \|}$

$\rho > 1 \Rightarrow$ explodes → no convergence

$\rho < 1 \Rightarrow$ "linear convergence" $\rho^k$

→      stop $(s_l f) < 0$

$\| \Delta_k \| \leq \rho^k \| \Delta_0 \|$     $(\log(\| \Delta_k \|) \leq k \log(\rho) + cst$

Summary :

$\rho_\tau = \| Id - \tau C \|_{op}$
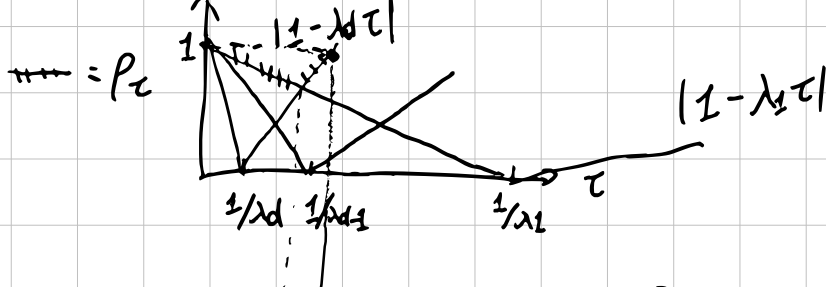
$\rho_\tau = \max |\lambda_i (Id - \tau C)|$

__Remark:__  $\lambda_i (B + Id) = \lambda_i(B) + 1$

$\lambda_i (Id - \tau C) = 1 - \tau \lambda_i(C)$

$\rho_\tau = \max | 1 - \tau \lambda_1(c) | , | 1 - \tau \lambda_2(c) |, | 1 - \tau \lambda_d(c) |$

$\lambda_1 \leq \lambda_2 \leq \cdots$

$1 - \tau \lambda_1 = 0$   $\tau = \frac{1}{\lambda_1}$

$|1 - \lambda_d \tau|$

$\frac{2}{\lambda_d} = \tau_{critical} = \frac{2}{L}$

$\tau_{opt} = \frac{2}{\lambda_d + \lambda_1} = \frac{2}{L + \mu}$

Thm :  If $0 < \mu = \lambda_1(C) \leq \lambda_2(c) \leq \cdots \leq \lambda_d(C) = L$

If $0 < \tau < \frac{2}{L}$, then $\exists \rho < 1$

s.t. $\| x_k - x^* \| \leq \rho_\tau^k \| x_0 - x^* \|$

A "good" (optimal for the proof)

choice is $\tau = \frac{2}{L + \mu}$

$\rho_{op} = \frac{L - \mu}{L + \mu} = \frac{K - 1}{K + 1}$

$K = \frac{L}{\mu}$    $K \to +\infty$   $\rho_{op} \to 1$

$K \to 1$   $\rho_{op} \to 0$

__Underdetermined__   $\mu = 0$ ,  $K = +\infty$

Thm :  If $0 < \tau_{min} \leq \tau_k \leq \frac{2}{L}$

$u_k$ will converge to some sol° $x^*$

$f(x_k) - f(x^*) \leq \frac{dist (x_0, \text{argmin}(f))}{8 \cdot \tau \cdot k}$  ~ true $c^{-?}$ $2 3 4$

__General case :__  $C \longrightarrow$ Hessian

Def:  $f : \mathbb{R}^d \to \mathbb{R}$

$\partial^2 f(x) = \left( \frac{\partial f}{\partial x_i \partial x_j}(x) \right)^d_{i,j=1} \in \mathbb{R}^{d \times d}$

$\partial^2 f(x)^T = \partial^2 f(x)$

__Prop:__  $f$ is cvx $\Longleftrightarrow \partial^2 f(x) \succeq 0$

$\lambda_i(\partial^2 f(x)) \geq 0$

Def:  $f$ is twice differentiable at $x$     quad

$f(x + \varepsilon \delta) = \underbrace{f(x)}_{cst} + \underbrace{\varepsilon \langle \nabla f(x), \delta \rangle}_{linear} + \frac{\varepsilon^2}{2} \langle \partial^2 f(x) \delta, \delta \rangle$

$+ o(\varepsilon^2)$

$\sum_{ij} \frac{\partial^2 f}{\partial x_i \partial x_j}(x) \times \delta_i \times \delta_j$

$\nabla f(x + \varepsilon \delta) = \nabla f(x) + \varepsilon \partial^2 f(x) \times \delta + o(\varepsilon)$

Thm 1 :  $0 < \mu \leq \lambda_i (\partial^2 f(x)) \leq L$   $f$ is convex

$\underset{\text{strong cvxity}}{\smile}$   $\underset{\text{smoothness}}{\smile}$   $\mu \geq 0$

Remark :  if $f(x) = \frac{1}{2} \langle Cx, x \rangle - \langle b, x \rangle$

$\nabla f(x) = Cu - b$

$\partial^2 f(x) = C$

$\mu = \underset{x}{\inf} \lambda_i (\partial^2 f(x))$

$L = \underset{x}{\sup} \lambda_i (\partial^2 f(x))$   $L \neq \lambda_{max} (\partial^2 f(x))$

Same thm :  $0 < \tau < \frac{2}{L} \Rightarrow$ conv

$\| x_k - x^* \| \leq \rho^k \| x_0 - x_0 \|$

Thm 2 :  same  $\mu = 0$,  $\tau < \frac{2}{L}$

$x_k \to x^k$

$f(x_k) - f(x^*) \leq O(\frac{1}{k})$

$\partial^2 f(x) \leq L Id$

$\mu \leq \lambda_i (\partial^2 f(x)) \leq L$

$\left[ f(x_k) + \langle \nabla f(x_k), x - x_k \rangle + \frac{\mu}{2} \| x - x_k \|^2 \right] \leq$

$\leq f(x_k)$

$\leq \left[ f(x_k) + \langle \nabla f(x_k), x - x_k \rangle + \frac{L}{2} \| x - x_k \|^2 \right]$

$GD \sim 1/k$

Vesterov $\sim 1/k^2$   optimal