1) Grad desc.

2) Regulariz° $\longrightarrow$ Ridge

Lasso  ...  Nucl. norm

TV

...

3) Non smooth convex optim $\longrightarrow$ Large scale

Interior point  $\hookrightarrow$ proximal

GD : $\min_{x \in \mathbb{R}^d} f(x)$   $u_{k+1} \overset{\circledast}{=} u_k - \tau \nabla f(u_k)$

$f(x) = \frac{1}{2} \| A x - y \|^2$   $\nabla f(x) = A^T (A x - y)$

① $\boxed{u_{k+1}} = u_k - \tau A^T (A u_k - y)$

if $x^*$ sol°  $\nabla f(x^*) = 0$   $A^T (A x^* - y) = 0$

② $\boxed{x^*} = x^* - \tau A^T (A x^* - y)$

①-②  $\underbrace{u_{k+1} - x^*}_{\varepsilon_{k+1}} = \underbrace{u_k - x^*}_{\varepsilon_k} - \tau A^T A \underbrace{(u_k - x^*)}_{\varepsilon_k}$

$$\| \varepsilon_{k+1} = \underbrace{\left( Id - \tau A^T A \right)}_{U_\tau} \varepsilon_k$$

$\underline{\text{Concl}^o}$ : $\quad \varepsilon_k = \left( U_\tau \right)^k \varepsilon_0$

$\underline{Q^o}$ : $\quad \| U_\tau \|_{op} < 1$ ??

$\| U_\tau \|_{op}$ = operator norm $\qquad$ np. linalg. norm

$\underline{\text{Def}}$ : If $B$ is a matrix

$$\| B \|_{op} = \sqrt{\lambda_{max}(B^T B)} = \sigma_{max}(B)$$

If $B$ is sym. $\quad B^T = B$

$$\| B \|_{op} = \max_i | \lambda_i(B) |$$

Why "operator"? $\qquad \| Bx \|_2 \leq \| B \|_{op} \cdot \| x \|_2$

Lipschitz constant $\frac{\|Bx\|}{\|x\|}$ $\quad \overset{B}{\curvearrowright} \quad$



$\underline{Q^o}$ : find $\tau$ s.t. $\quad \max_i | \lambda_i(Id - \tau A^T A) | \overset{\circledast}{\lessgtr} 1$

theorem :. if $\tau < \frac{2}{\|A\|_{op}^2}$ then $\wedge$ it's true

$\circledast$

- Overdetermined $A^T A$ is invertible.

$$0 < \underbrace{\mu = \lambda_{min}(A^T A)}_{\mu \leq} \leq \underbrace{\lambda_{max}(A^T A) = L}_{\lambda_i(A^T A) \leq L}$$

correl°

if $\mu > 0$, then fast ("linear") convergence

the optimal $\tau = \frac{2}{\mu + L}$ → Geometrical

"LINEAR"

$$\|x_k - x^*\| \leq \left(\frac{L-\mu}{\mu+L}\right)^R \|x_k - x_0\|$$

$\mu > 0$ $< 1$

$\frac{L-\mu}{L+\mu} = \frac{(L/\mu)-1}{(L/\mu)+1}$ → $\frac{L}{\mu} = K \geq 1$ conditionning

If $\mu = 0$    $\tau \leq \frac{2}{L}$ $\longrightarrow$ converge.

Thm:    $\tau \leq \frac{2}{L}$        $f(x) = \frac{1}{2} \|Ax - y\|^2$

$$\left. f(x_k) - f(\hat{x}) \underset{\substack{\uparrow \\ \text{"SUB-LINEAR"}}}{\leq} \frac{f(x_0) - f(x^*)}{\boxed{k}} \right\}$$
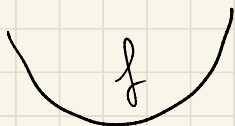
## General Case : $f(x)$ convex $F^c$.

$\boxed{A^T A}$ $\longrightarrow$ Hessian of $f$

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

Prop$^c$: $\partial^2 f(x)$ symmetric

$f$ is convex $\iff$ $\partial^2 f(x) \underset{\underset{\text{eigenvalues}}{\uparrow}}{\succeq} 0$

$\smile f$        $f'' \geq 0$

Thm :
$$\mu = \inf_x \inf_i \lambda_i(\partial^2 f(x))$$

(f twice diff)
$$L = \sup_x \sup_i \lambda_i(\partial^2 f(x))$$

$$\frac{L}{\mu} = k \quad \text{cond.}$$

① if $\tau < \frac{2}{L}$ ⟶ convergence

② if $k < +\infty$ $(\mu > 0)$ ⟶ Fast convergence (linear).

Stochastic opt$^e$ :
$$\min_x f(x) = \frac{1}{n} \sum_{\underset{\downarrow \text{data}}{(i)}} f_i(x).$$

$$\underbrace{Df(x) = \frac{1}{n} \sum_i \nabla f_i(x)}_{\substack{\text{too slow.} \\ \text{CRUCIAL}}} \rightsquigarrow \widehat{\partial} f(x) = \frac{1}{|t|} \sum_{i \in t} \nabla f_i(x)$$

$$\uparrow$$
Random

$$\mathbb{E}(\widehat{\partial} f(x)) \overset{?}{=} Df(x) \quad \text{"Unbiased"}$$

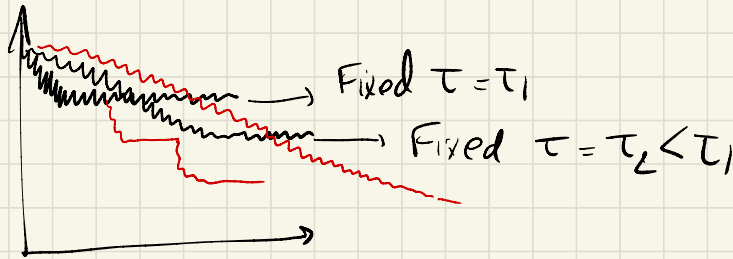SG$\underline{\textbf{D}}$ : $\quad \underline{u_{k+1} = u_k - \tau_k \widehat{\partial} f(x)}$. $\triangle$ RANDOM.

↳ "Descent" FALSE

Stochastic Approxi.
ROBIN & MONROE

Thm : $\tau_k \searrow 0$ for instance $\tau_k = \frac{1}{k}$

$\boxed{x_k} \longrightarrow x^*$ almost everywhere.



$\longrightarrow$ Fixed $\tau = \tau_1$
$\longrightarrow$ Fixed $\tau = \tau_2 < \tau_1$

$\rightarrow$ Accelerati. / Momentum $\rightarrow$ Extrapolati.
$\quad \quad \hookrightarrow$ Nesterov / Heavy Ball

GD : $O(1/k)$ $\rightsquigarrow$ Nesterov $\underline{O(1/k^2)}$.
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad$ OPTIMAL

# Regulariz° :
$$\min_{x} \frac{1}{2} \| Ax - y \|^2$$

Pbm: if $u < d$. (undetermined).

↳ non unique sol° $\text{Ker}(A) = \{0\}$.

↳ Overfitting ⌇⌇⌇⌇⌇

## Ridge regul° / Tichonov / Weight decay

$$\min_{x \in \mathbb{R}^d} \frac{1}{2} \| Ax - y \|^2_{\mathbb{R}^n} + \underbrace{\frac{\lambda}{2} \| x \|^2_{\mathbb{R}^d}}_{\text{Ridge penalty}} = f(x)$$

$$C = 1/\lambda$$

$$\lambda = \text{Lagrange mult.}$$

$$\left. \begin{array}{l} \text{strong overfit} \\ \text{large noise} \end{array} \right\} \longrightarrow \lambda \nearrow \longrightarrow \begin{array}{l} \text{Favor} \\ \text{"small"} \\ u \end{array}$$

## Least square + ridge :

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

$$\nabla f(x) = 0 \iff A^TAx + \lambda x = A^Ty.$$

$$\iff \left(\underbrace{A^TA}_{\geq 0} + \underbrace{\lambda \, Id_d}_{>0}\right)x = A^Ty.$$

$$\underbrace{\phantom{\left(A^TA + \lambda Id_d\right)}}_{\text{Invertible ALWAYS!}}$$

Ridge $\rightarrow$ unique sol$^\circ$.

$\underline{\text{Concl}^\circ}$ : $\quad x_\lambda = \left(A^TA + \lambda \, Id_d\right)^{-1} A^T y$

Thm : (WOODBURY formula).

$$\left(\underbrace{A^TA}_{\substack{=C \\ \text{covell}^\circ}} + \lambda \, Id_{\textcolor{red}{\boxed{d}}}\right)^{-1} A^T = A^T \left(\underbrace{AA^T}_{\substack{=K \\ \text{Kernel}}} + \lambda \, Id_{\textcolor{red}{\boxed{m}}}\right)^{-1}$$

Feature space $(\mathbb{R}^d)$ $\qquad$ Kernel space $(\mathbb{R}^n)$.

If $d \gg n \longrightarrow$ GO KERNEL !

Works d = +∞   Kernel (RKHS)

[Reproducing Kernel
 Hilbert Space.

LASSO : regalizer $\ell^1$ norm                    RIDGE

$$\|u\|_1 = \sum_i |x_i| . \neq \sum_i |x_i|^2$$

$$\min_u \frac{1}{2}\|Au - y\|^2 + \underline{\lambda} \|u\|_1 .$$

"$\ell^1$ promotes lots of zeros in $u^*$"
↳ sparsity

modeling : image proarning.
↳ model selecti / Explainable.
   $x \in \mathbb{R}^d$  d very large.
   select only a "few" feature  $u_i = 0$
                                 for a lot of $i$'s.

Lasso does a selecti

$\underline{LASSO} : x_\lambda = \underset{x}{\text{argmin}} \; \frac{1}{2} \|Ax - y\|^2 + \lambda \|x\|_1$

$\underset{x}{\text{min}} \; \frac{1}{2} \|Ax - y\|^2 \xrightarrow{\text{ill-pose}} \infty \; \# \text{ of sol}^{\underline{s}}$

$\mathcal{H} = \{ x : Ax = y \}.$ $\qquad \underline{\underline{d}} \underline{\underline{\gg n}} \quad Ker(A) \neq \{0\}.$

$x_\lambda$ sol$^{\underline{s}}$ of LASSO

$\underline{Thm} : \; x_\lambda \xrightarrow{\lambda \to 0} x_0$ a sol$^{\underline{s}}$ of Basis Pursuit

$$\underset{Ax = y}{\min} \; \|x\|_1$$

$n = 1 \quad d = 2. \quad \textcolor{red}{H} = \{ Ax = y \} \; line.$



LASSO $\to x = (0, x_2)$

Ridge $- x = (x_1, x_2)$

$x = (x_1, 0)$

Ridge : $\underset{Ax = y}{\min} \|x\|_2^2$

lasso. $\underset{Ax = y}{\min} \; \boxed{\|x\|_1}$

Good: $u$ sol$^n$ is sparse.
$$f(u) = \frac{1}{2}|Ax-y|^2 + \textcircled{$\lambda$}|u|_1$$
is convex

Bad: $f$ is non diff. $|x|$



non smooth



$\tau$ medium
$\tau$ small

Support vector machine (Hinge loss) ; same.

<u>Idea</u> : splitting.

<u>1 part</u> : "Proximal operator" $\ell^1$.

$\qquad A = Id$.

<u>Def</u> : For some func$^n$ $g(u)$ (ex. $\ell^1$).

$$\text{Prox}_{\tau g}(y) \overset{\text{def}}{=} \underset{(u)}{\text{argmin}} \ \frac{1}{2}\|u-y\|^2 + \tau g(u)$$
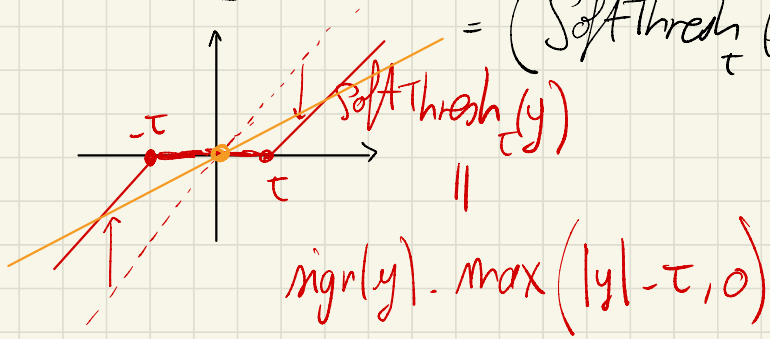
"Resolvant" operator

$\underline{\text{Ridge}}$ : $\text{Prox}_{\tau.\|.\|^2_2}(y) = \underset{x}{\text{argmin}} \frac{1}{2}\|x-y\|^2 + \frac{\tau}{2}\|x\|^2$

$\longrightarrow 0 = x - y + \tau x \implies y = \boxed{\dfrac{x}{1+\tau}}$

$\underline{\text{Lasso}}$ : $\text{Prox}_{\tau.\|.\|_1}(y) = \text{SoftThresh}_\tau(y) = \left(\text{SoftThresh}_\tau(y_i)\right)_{i=1}^d$



$\text{SoftThresh}_\tau(y)$

$= \text{sign}(y) . \max(|y|-\tau, 0)$

$\underline{\underset{\ell^2}{\text{ISTA}}}$ : Iterative Soft Thresholding.

(~2003 Daubechies / De Mol / De Frise).

[Special case] $\longrightarrow$ Forward - Backward algo.

[Step ISTA]

$\longrightarrow$ Gradient Descent on $\frac{1}{2}\|Ax-y\|^2$
with step $\tau$

$\tilde{x}_k := x_k - \tau A^\top(Ax_k - y) \textcolor{red}{= x_k - \tau(Cx_k - u)}$

Soft Threshold $\lambda\tau$

$x_{k+1} := \text{SoftThreshold}_{\lambda\tau}[\tilde{x}_k]$

$\textcolor{red}{u = A^\top y}$
$\textcolor{red}{C = A^\top A}$

$k \leftarrow k+1$

Thm : If $\tau < \dfrac{2}{\|A\|_y^2}$ , then $u_k \longrightarrow u^*$ sol$^\circ$ LASSO.

$$f_\lambda(u_k) - f_\lambda(u^*) \sim \boxed{\dfrac{1}{k}}$$

Accelerate ISTA $\longrightarrow$ FISTA $O(\gamma_{k^2})$.
$\qquad\quad O(\gamma_k) \nearrow$ optimal
$\qquad\qquad$ Nesterov

Regul$^\circ$ path : influence $\underline{\underline{\lambda}}$



$$x_\lambda = (\underbrace{u_\lambda^1, u_\lambda^2 \dots u_\lambda^d}_{d \text{ feature}})$$

$$\lambda \longmapsto u_\lambda^i$$



$\longrightarrow$ sol$^\circ$ 0

$\longrightarrow \lambda_{max} = \|A^+ y\|_\infty$

$\lambda_{max}/10$