Proximal Algorithms.    Neal Parikh  2013.

setting : $f : R^n \to R \cup \{+\infty\}$  closed proper convex.    $\Rightarrow$ epi $f$ nonempty closed convex.

proximal operator.

$$\boxed{prox_f(v) = \underset{x}{argmin}\left(f(x) + \frac{1}{2}\|x - v\|_2^2\right)}$$

$$prox_{\lambda f}(v) = \underset{x}{argmin}\left(f(x) + \frac{1}{2\lambda}\|x - v\|_2^2\right) = \underset{x}{argmin}\left(\lambda f(x) + \frac{1}{2}\|x - v\|_2^2\right)$$

$\downarrow$ 在 v 附近最小化 f(x) 的点.

v 的位移: 距离增加 和 f(x) 增加的 trade off.

Gradient descent perspective

$$prox_{\lambda f}(v) \Longrightarrow v \approx v - \lambda \nabla f(v).$$
$\uparrow$
Stepsize

important property : link with fixed point theory :  $prox_f(x^*) = x^*$  iff  $x^*$ minimize $f(x)$.

Properties & definitions of proximal operators. :

① $f$ separable    $f(x,y) = \varphi(x) + \psi(y)$  $\Rightarrow$  $prox_f(v,w) = (prox_\varphi(v), prox_\psi(w))$
$\downarrow$                     sum.                    joint.

$v, x \in R^n$. $f(x) = \sum_{i=1}^n f_i(x_i)$ $\Rightarrow$ $(prox_f(v))_i = prox_{f_i}(x_i)$

② $f(x) = a\varphi(x) + b$.  $a > 0$  $\Rightarrow$  $prox_{\lambda f}(v) = prox_{\alpha\lambda\varphi}(v)$.

$f(x) = \varphi(\alpha x + b)$ $\alpha \neq 0$ $\Rightarrow$ $prox_{\lambda f}(v) = \frac{1}{\alpha}(prox_{\alpha^2\lambda\varphi}(\alpha v + b) - b)$

③ orthogonal $a$, $f(x) = \varphi(ax)$ $\Rightarrow$ $prox_{\lambda f}(v) = a^T prox_{\lambda\varphi}(av)$

④ $f(x) = \varphi(x) + a^T x + b.$ $\Rightarrow$ $prox_{\lambda f}(v) = prox_{\lambda\varphi}(v - \lambda a)$.

⑤ $f(x) = \varphi(x) + \frac{\rho}{2}\|x - a\|_2^2.$ $\Rightarrow$ $prox_{\lambda f}(v) = prox_{\tilde{\lambda}\varphi}\left((\frac{\tilde{\lambda}}{\lambda})v + (\rho\tilde{\lambda})a\right)$   $\tilde{\lambda} = \frac{\lambda}{1 + \lambda\rho}$

Fixed point algorithms:

defs: ① Contraction: $f$ Lipschitz continuous with $K < 1$    $d(f(x), f(y)) \leq K d(x,y)$.  $K < 1$

② non-expansive : $f$ Lipschitz continuous with $K = 1$    $d(f(x), f(y)) \leq d(x,y)$

③ firmly nonexpansive: $f$ ~~Lipschitz~~  $\int dis(f(x), f(y))^2 \leq (x-y)^T(f(x) - f(y))$  [同 $\xi$ vector space].
$\to$ can prove to be Lipschitz-1

④ averaged operator: if N is nonexpansive : $T = (1-\alpha)I + \alpha N$.

Thms: ① contraction can find fixed point $f(x) = x^*$ by $x^{(k+1)} = f(x^{(k)})$

② averaged operator can also find fixed point using ①

③ Firmly nonexpansive operators are indeed. $\frac{1}{2}$ overaged operator. $\forall T$ firm $\Rightarrow$ $2T - I$ nonexpan

④ averaged operators are closed under composition, but firmly nonexpansives are not.

⑤ proximal operators are firmly expansive operators

⑥ if $N$ is only nonexpansive, update by $x^{(k+1)} = (1-\alpha) x^{(k)} + \alpha N(x^{(k)})$ [为了 average].

def: proximal average of $f_1 \cdots f_m$. closed proper convex $g$ s.t. $\frac{1}{M} \sum_{r=1}^{m} prox_{f_r} = prox_g$.

Moreau decomposition: $v = prox_f(v) + prox_{f^*}(v)$　$f^*_{(y)}$: conjugate $= \sup_x (y^T x - f(x))$

e.g. $v = \Pi_L(v) + \Pi_{L^\perp}(v)$.　$f(x) = I_L(x)$, $(I_L(x))^* = I_{L^\perp}(x)$　$prox_{I_L(v)}(v) = \Pi_L(v)$

e.g. $v = \Pi_K(v) + \Pi_{K^*}(v)$　$K^* = \{y: y^T x \le 0, \forall x \in K\}$ polar cone, negative of dual cone.

can be used as. $prox_f(v) = v - prox_{f^*}(v)$.

"Smooth" approximation perspective of proximal operator.

infimal convolution: $f \square g(v) = \inf_x [f(x) + g(v-x)]$　$v \in dom f + dom g$.

Moreau envelope: $M_{\lambda f} = \frac{1}{\lambda}(\lambda f \square \frac{1}{2} \|\cdot\|_2^2) = \inf_x [f(x) + \frac{1}{2\lambda} \|v-x\|_2^2]$

(Moreau Tosida regularization).

if define $h(x,v) = f(x) + \frac{1}{2\lambda} \|v-x\|_2^2$

$M_{\lambda f}(v) = h(prox_{\lambda f}(v), v)$

$M_{\lambda f}(v)$ & $f(x)$ have the same minimizers.
↳ Always continuously differentiable
　Always have domain $R^n$

Establish a approximation relationship between $M_{\lambda f}(x)$ and $f(x)$:

$M_f$ is a smoothed version of $f$.

property $(f \square g)^* = f^* \square g^*$ ⟹ $M^*_f = f^* + (\frac{1}{2}\|\cdot\|_2^2)^*$ ⟹ $M^*_f = f^* + \frac{1}{2}\|\cdot\|_2^2$

$\frac{1}{2}\|\cdot\|_2^2$ is self dual!

$M^{**}_f = M_f$ ⟹ $M_f = (f^* + \frac{1}{2}\|\cdot\|_2^2)^*$　dual → regularize → dual. (it's smooth)
　　　　　　　　　　　　　　　　　　　　　　　(get strong convex)

$M_{\lambda f}$ differentiable, so we can take derivative.

~~$M_{\lambda f}(prox_{\lambda f}(x))$　$M_{\lambda f}(x) = h(prox_{\lambda f}(x), x) =$~~

~~$M_{\lambda f}(v) = h(prox_{\lambda f}(v), v) = f(prox_{\lambda f}(v)) + \frac{1}{2\lambda}\|v - prox_{\lambda f}(v)\|_2^2$~~
　　　　　　　　　　→ approximate $\nabla f(x)$. [so it's a gradient descent step].

⟹ $prox_{\lambda f}(x) = x - \lambda \nabla M_{\lambda f}(x)$.

$prox_{f}(x) = \nabla M_f(x)$.

proximal operator is the resolvent of subdifferential operator.

$prox_{\lambda f} = (I + \lambda \partial f)^{-1}$　what's nontrivial here: $(I + \lambda \partial f)^{-1}$ becomes single-valued mapping
　　　　　↑ subgradient of $f$

More perspective from gradient descent:

① $prox_{\lambda f}(x) = x - \lambda \nabla M_{\lambda f}(x)$.

② if $\nabla f(x)$ exists, first-order approximation $\hat{f}^{(1)}_v(x) = f(v) + \nabla f(v)^T (x-v)$.
　then $prox_{\hat{f}^{(1)}_v}(v) = v - \lambda \nabla f(v)$.　explain: minimize first-order approximation of $f(x)$ at point $v$ (near $v$)
　　　　　　　　　　　　　　　　　　　result in a gradient descent step from $v$.

③ if $\nabla^2 f(x)$ exists, second-order approximation $\hat{f}^{(2)}_v(x) = f(v) + \nabla f(v)^T (x-v) + \frac{1}{2}(x-v)^T \nabla^2 f(v)(x-v)$
　then $prox_{\hat{f}^{(2)}_v}(v) = v - (\nabla^2 f(v) + \frac{1}{\lambda}I)^{-1} \nabla f(v)$. explain: similar, but 2-order, result in a Levenberg-Marquardt step.

Trust Region problem perspective.

trust region problem:
$$\min \quad f(x)$$
$$\text{s.t.} \quad \|x-v\|_2 \leq \rho.$$

proximal problem:
$$\min \quad f(x) + \frac{1}{2\lambda}\|x-v\|_2^2$$

relationship:    solution for some $\rho$ $\xleftarrow{\text{is a}}$ solution

solution $\longrightarrow$ unconstrained minimizer of $f$ / solution for some $\lambda$

Above is the interpretation of proximal operator.

Algorithms.

[1] Direct proximal minimization

$$x^{k+1} = \text{prox}_{\lambda_k f}(x^k).$$  guarantee convergence with $\lambda_k > 0$, $\sum_{k=1}^{\infty} \lambda^k = \infty$.

application: ill-conditioned $f$ [we add a quadratic term to be strong convex]

perspective: $x^{k+1} = \underset{x}{\arg\min}\, f(x) + \frac{1}{2\lambda_k}\|x - x^k\|_2^2$   regularization gets smaller as $x \to x^*$. the impact of the term disappears with iterations.

e.g. of application: iterative refinement. $f(x) = \frac{1}{2}x^T A x - b^T x$. ill-conditioned. $A$.

$$\text{prox}_{\lambda f}(x^k) = (A + \tfrac{1}{\lambda}I)^{-1}(b + \tfrac{1}{\lambda}x^k)$$
$$= x^k + (A + \tfrac{1}{\lambda}I)^{-1}(b - Ax^k)$$

$\underset{\hat{A}}{\downarrow}$   iteratively compensating for $\hat{A} - A$. difference.

[2] Proximal gradient Method.

$\min f(x) + g(x)$. $f(x)$ differentiable. both closed proper convex., $g(x)$ can be non-smooth.

$$x^{k+1} = \text{prox}_{\lambda^k g}(x^k - \lambda^k \nabla f(x^k)) \iff x^k - \lambda^k \nabla f(x^k) - \lambda^k \nabla M_{\lambda^k g}(x^k) \approx x^k - \lambda^k \nabla(f + g)(x^k).$$

Convergence: $\nabla f$ Lipschitz - L $\Rightarrow$ $\lambda^k = \lambda \in (0, \tfrac{1}{L}]$, Converge with $O(\tfrac{1}{k})$.

L not known: back-tracking "line" search.    parameter $\beta \in (0,1)$. $\lambda = \lambda^{k-1}$
Beck & Teboulle.    $\Rightarrow z = \text{prox}_{\lambda g}(x^k - \lambda \nabla f(x^k))$
if $f(z) > \hat{f}_\lambda(z, x^k)$, $\lambda := \beta\lambda$
$\hat{f}_\lambda(x,y) = f(x) + f(y) + \nabla f(y)^T(x-y) + \frac{1}{2\lambda}\|x-y\|_2^2$

perspective ① Majorization - minimization. (like EM algorithm)

Consider minimizing $\varphi(x)$.
Step 1: majorization: find convex upper bound $\hat{\varphi}$ of $\varphi$ tight at $x^k$: $\begin{cases} \hat{\varphi}(x, x^k) \geq \varphi(x). \\ \hat{\varphi}(x, x) = \varphi(x). \end{cases}$  (surrogate)
step 2: minimization: $x^{k+1} = \underset{x}{\arg\min}\, \hat{\varphi}(x, x^k)$.

For $f(x) = g$, a surrogate is $\hat{f}_\lambda(x,y)$   EM alg will give precisely gradient descent.
For $f(x) + g(x)$. $\hat{g}_\lambda(x,y) = \hat{f}_\lambda(x,y) + g(x) \iff$ proximal gradient.

perspective ② Solution is a fixed point for $prox_{\lambda g}((I-\lambda\nabla f)(x)) = (I+\lambda\partial g)^{-1}(I-\lambda f)$.

## [3] Accelerated proximal gradient.

$$y^{k+1} = x^k + \omega^k(x^k - x^{k-1})$$

$$x^{k+1} = prox_{\lambda^k g}(y^{k+1} - \lambda^k \nabla f(y^{k+1}))$$

recommend. $\omega^k = \dfrac{k}{k+3}$

convergence: $\lambda^k = \lambda = O(0, \frac{1}{L}]$, rate $O(\frac{1}{k^2})$.

If $L$ not know. again Beck & Teboulle search but use ~~$y^{k+1}$~~ $y^k$

## [4] ADMM version. of proximal gradient.

$$\min f(x) + g(x). \iff \min f(x) + g(z)$$
$$s.t. \quad x = z$$

both $f, g$ can be non-smooth.

recall ADMM (scaled)

$$x^{k+1} = \underset{x}{\text{argmin}} f(x) + \frac{\rho}{2}\|x - z^k + u^k\|_2^2 \quad [\text{struck: } prox_{\lambda f}(z^k - u^k)]$$
$$z^{k+1} := \underset{z}{\text{argmin}} g(z) + \frac{\rho}{2}\|x^{k+1} - z + u^k\|_2^2 \quad [\text{struck: } prox_{\lambda g}(x^{k+1}+u^k)]$$
$$u^{k+1} := u^k + x^{k+1} - z^{k+1}$$

$$\iff$$

$$x^{k+1} = prox_{\lambda f}(z^k - u^k)$$
$$z^{k+1} = prox_{\lambda g}(x^{k+1} + u^k).$$
$$u^{k+1} = u^k + x^{k+1} - z^{k+1}.$$

$$\lambda = \frac{1}{\rho}.$$

an interesting insight: prove the convergence of ADMM $\iff$ fixed point algorithm of a firmly nonexpansive operator.

## [5]. Linearized ADMM.

$$\min f(x) + g(Ax) \iff \min f(x) + g(z)$$
$$s.t. \quad Ax - z = 0.$$

Original ADMM.

$$x^{k+1} = \underset{x}{\text{argmin}} f(x) + \rho u^{kT}(Ax - z^k) + \frac{\rho}{2}\|Ax - z^k\|_2^2 \quad [\text{struck: } \frac{\rho}{2}\|Ax - z^k + u^k\|_2^2]$$

$$z^{k+1} = \underset{z}{\text{argmin}} g(z) + \rho u^{kT}(Ax^{k+1} - z) + \frac{\rho}{2}\|Ax^{k+1} - z\|_2^2$$

$$u^{k+1} = u^k + Ax^{k+1} - z^{k+1}.$$

only modify x step:

replace $\frac{\rho}{2}\|Ax - z^k\|_2^2 \to \frac{\rho}{2}\underbrace{x^T A^T A x}_{} \quad \frac{\rho}{2}(A^TAx - A^Tz^k)^Tx$ with.

$$\rho(A^TAx^k - A^Tz^k)^Tx + \frac{\mu}{2}\|x - x^k\|_2^2. \quad 0 < \mu \leq \frac{\lambda}{\|A\|_2^2}$$

then new algorithm:

$$x^{k+1} = prox_{\frac{\mu}{\mu}}\left(x^k - \frac{\mu}{\lambda}A^T(Ax^k - z^k + u^k)\right).$$
$$z^{k+1} = prox_{\lambda g}(Ax^{k+1} + u^k).$$
$$u^{k+1} = u^k + Ax^{k+1} - z^{k+1}.$$