

* Now let's switch to Boyd's proximal gradient paper:

Before proceeding any further let's look at the difference in notation between Calafiore and Boyd

Calafiore	Boyd
$f_0(x)$	$f(x)$
$h(x)$	$g(x)$
s_k	λ^k
x_{k+1}	x^k

index starts with 0 index starts at 1

$$x_{k+1} = \text{prox}_{\frac{s_k}{h}}(x_k - s_k \nabla f_0(x_k)) \quad x^{k+1} = \text{prox}_{\lambda^k g}(x^k - \lambda^k \nabla f(x^k))$$

Proximal Gradient Method with variable step size (f : Differentiable Function, g : Nondifferentiable Function, ϵ : Tolerance)

```

x1 = 0 # say
λ0 = 1e-5 # say
k = 1
β = 0.5 # β ∈ (0,1)
while(1)
    λk, xk+1 = BT-LS-nxt-pt-gnr(xk, λk-1, β)
    if ||xk+1 - xk||2 ≤ ε
        break
    else
        k := k+1
    end # if
return xk
end # while
    
```

remember proximal gradient
 # algorithm is a fixed point problem
 # i.e. $x^* = \text{prox}_{\lambda g}(x^* - \lambda \nabla f(x^*))$
 # that is why terminating condition
 # is of the form $\|x^{k+1} - x^k\|_2 \leq \epsilon$

```

function BT-LS-nxt-pt-gnr(xk::vector, λk-1::scalar, β::scalar)
    # Beck Teboulle line search and next point generator
    λ := λk-1 # λk-1 denotes xk = prox(xk-1 - λk-1 ∇f(xk-1)) which we already know from given
    while(1)
        z := prox(xk - λ ∇f(xk)) # fixed stepsize (as in Calafiore) would set xk+1 = z right away, however as in this
            λg # step size is not fixed
        if f(z) ≤ Ŝλ(z, xk) # Ŝλ(x, y) = f(y) + ∇f(y)T(x-y) + 1/(2λ) ||x-y||22: this is an upper bound for f(x)
            break
        else λ := βλ
    end # if
end # while
    
```

end # while
end # function

is of the form $\|x^{k+1} - x^k\|_2 \leq \epsilon$

end # if
end # while
return $\lambda^k := \lambda, x^{k+1} := z$ # then, $x^{k+1} = \text{prox}_{\lambda^k g}(x^k - \lambda^k \nabla S(x^k))$
end # function

4.1 Proximal minimization # also known as proximal iteration, proximal point algorithm

$$x^{k+1} := \text{prox}_{\lambda f}(x^k)$$

↓
: $\mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$, closed proper convex

• (convergence guaranteed for $\lambda^k > 0, \sum_{k=1}^{\infty} \lambda^k = \alpha$)

[# Moreau-Yosida regularization, M f](#)

• Interpretation: 1) gradient method applied to Moreau envelope $M_{\lambda f}$

because, $\text{prox}_{\lambda f}(x^k) = x^k - \lambda \nabla M_{\lambda f}(x^k)$ # $\text{prox}_{\lambda f}(x) = x - \lambda \nabla M_{\lambda f}(x)$
∴ $x^{k+1} = x^k - \lambda \nabla M_{\lambda f}(x^k)$

2) Simple iteration to find fixed point of $\text{prox}_{\lambda f}(x)$

4.2 Proximal gradient method: (For detailed convergence proof see [Proximal algorithm Calafiore](#))

$\forall f(x) + g(x)$
 $f: \mathbb{R}^n \rightarrow \mathbb{R}, \text{cpc}$ $g: \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}, \text{cpc}$
 (can encode constraints on variable x for being extended valued)

• The objective is split into two terms

nonunique splitting → different splitting leads to different nonunique splitting

∇f : Lipschitz continuous $\Leftrightarrow \forall_{x,y} \|\nabla f(x) - \nabla f(y)\|_2 \leq L \|x - y\|$ % this is a way of bounding the hessian so in the Taylor series the effect of second order terms will be negligible

different nonunique splitting

∇f : Lipschitz (continuous) $\Leftrightarrow \forall x, y \quad \|\nabla f(x) - \nabla f(y)\|_2 \leq L \|x - y\|$ this is a way of bounding the messian so in the Taylor series the effect of second order terms will be negligible

Proximal gradient algorithm:

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

$\lambda^k > 0$

Theoretical results: $(\nabla f$: Lipschitz continuous, L Lipschitz constant, $\lambda^k = \lambda$: fixed stepsize, $\lambda \in (0, \frac{1}{L}] \rightarrow$ converges with $O(\frac{1}{k})$)

L not known: use Beck-Teuboulle proximal gradient updates (most of the time)

Beck-Teuboulle proximal gradient updater

Beck-Teuboulle-prox-gradient-updater(x^k, λ^k, β) $\beta \in (0, 1)$ # Will output x^{k+1} and λ^k for which $\lambda^{k+1} = \text{prox}_{\lambda^k g} (x^k - \lambda^k \nabla f(x^k))$

$\lambda = \lambda^{k-1}$

while(1)

$z := \text{prox}_{\lambda g} (x^k - \lambda \nabla f(x^k))$ # z is a candidate for the x^{k+1} , however our λ here is still the λ^{k-1} of the previous step, so what we do is keep making λ smaller until it satisfies the Beck-Teuboulle line search condition $f(z) \leq \hat{f}_\lambda(x, x^k)$

if $f(z) \leq \hat{f}_\lambda(x, x^k)$ # $\hat{f}_\lambda(x, x^k) = f(y) + \nabla f(y)^T (x - y) + \frac{1}{2\lambda} \|x - y\|_2^2$: a valid upper bound for f

return $\lambda = \lambda, x^{k+1} = z$ # A justification of this coming soon # Intuitive Explanation behind Beck line search condition

end

$\lambda = \beta \lambda$ # If we have arrived at this line, then for sure the Beck-Teuboulle line search condition was not met i.e. for that λ the upperboundedness condition is not met, so we need to make λ smaller and then check the line search condition again. Essentially we keep reducing λ i.e., increasing $1/2 \lambda$ until,

end

def: majorization $\hat{f}_\lambda(x, x^k) = f(x)$ all we have to check is the upperboundedness

Evolution of subgradient method to proximal gradient method:

Normal subgradient: $x^{k+1} = x^k - \lambda^k g^k$: solves $\forall f(x)$
 $g^k \in \partial f(x^k)$

projected subgradient: $x^{k+1} = [x^k - \lambda^k g^k]_X = \Pi_X (x^k - \lambda^k g^k)$: solves $\begin{cases} \forall f(x) \\ x \in X \end{cases} = f(x) + I_X(x)$

Proximal gradient method: $x^{k+1} = \text{prox}_{\lambda h} (x^k - \lambda^k \nabla f(x^k))$: solves $(\forall f(x) + h(x))$

*Some special cases:

$$g = \mathbb{I}_C \Rightarrow x^{k+1} = \text{prox}_{\lambda \mathbb{I}_C} (x^k - \lambda^k \nabla f(x^k))$$

$$= \text{prox}_{\lambda \mathbb{I}_C} (x^k - \lambda^k \nabla f(x^k))$$

$$= \Pi_C (x^k - \lambda^k \nabla f(x^k))$$

[$\text{prox}_{\lambda \mathbb{I}_C}(\cdot) = \Pi_C(\cdot)$]

↳ just projected gradient method

Interpretations:

Majorization minimization interpretation

Majorization minimization: Large class of algorithms \supset {gradient method, Newton's method, Expectation minimization algorithm}

goal: $\forall x \quad \Phi(x): \mathbb{R}^n \rightarrow \mathbb{R}$

majorization minimization.

$$x^{k+1} = \text{argmin}_x \hat{\Phi}(x, x^k)$$

def: majorization

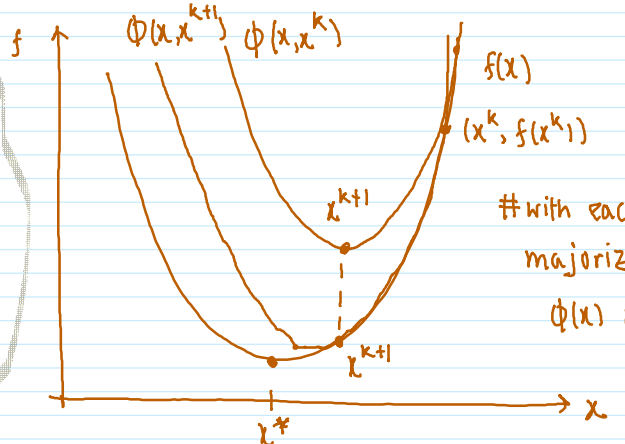
$\hat{\Phi}(x, x^k)$ {majorization of $\Phi(x)$ }

$\hat{\Phi}(x, x^k)$: (convex upperbound, tight in x^k)

need to check, however think it is fine

$\forall x^k \quad \hat{\Phi}(x, x^k) \geq \Phi(x), \Phi(x^k) = \hat{\Phi}(x, x^k)$

$\hat{\Phi}(x, x^k): D_x$



with each iteration we are approaching x^* , with majorization $\hat{\Phi}(x, x^k)$ resembling the original function $\Phi(x)$ more and more

% For a convex function at each iteration we move strictly to the optimal point because of the convex nature of the function, that is why majorization-minimization type function always converges for convex optimization problem

• At each step:

1) We majorize (upper bound) the objective

% At each iteration in $\theta(x, x^k)$, x^k is changing, as a result the majorization (the upper bound function) will also change in shape

2) Minimize the majorization

an upper bound of f ,

$$\hat{f}_\lambda(x, y) = f(y) + \nabla f(y)^T (x-y) + \frac{1}{2\lambda} \|x-y\|_2^2 \quad \# \text{ Lemma: Underestimation Lemma}$$

This is a majorization as

(i) $x=y \rightarrow \hat{f}_\lambda(x, x) = f(x)$

(ii) $\forall \lambda \in (0, \frac{1}{L}] \quad \forall x, y \quad f(x) \leq \hat{f}_\lambda(x, y)$

$\therefore x^{k+1} = \arg \min_x \hat{f}_\lambda(x, x^k)$ will be a majorization-minimization algorithm

Now let's show that proximal gradient too is a majorization-minimization algorithm

Proximal gradient algorithm is

$$x^{k+1} = \text{prox}_{\lambda^k g} (x^k - \lambda^k \nabla f(x^k)) \quad \# \text{ prox}_{\lambda^k g}(z) = \left(x \parallel \nabla f(x) \parallel \square + \frac{1}{2\lambda} \|z - x\|_2^2 \parallel \text{supp} \square \right)$$

$$= \left(x \parallel \lambda^k \nabla g(x) \parallel \square + \frac{1}{2\lambda} \|x^k - x^k + \lambda^k \nabla f(x^k) - x\|_2^2 \parallel \text{supp} \square \right)$$

$$= \arg \min_x \left(\lambda^k g(x) + \frac{1}{2} \|x - x^k + \lambda^k \nabla f(x^k)\|_2^2 \right)$$

$$\| (x - x^k) + \lambda^k \nabla f(x^k) \|_2^2 = \|x - x^k\|_2^2 + (\lambda^k)^2 \|\nabla f(x^k)\|_2^2 + 2\lambda^k \nabla f(x^k)^T (x - x^k)$$

$$= \arg \min_x \left(\lambda^k g(x) + \frac{1}{2} \|x - x^k\|_2^2 + \frac{1}{2} (\lambda^k)^2 \|\nabla f(x^k)\|_2^2 + \lambda^k \nabla f(x^k)^T (x - x^k) \right)$$

(constant w.r.t x, so argmin same iter)

$$= \arg \min_x \left(\lambda^k g(x) + \frac{1}{2} \|x - x^k\|_2^2 + \lambda^k \nabla f(x^k)^T (x - x^k) \right)$$

$$= \arg \min_x \left(\lambda^k g(x) + \frac{1}{2} \|x - x^k\|_2^2 + \lambda^k \nabla f(x^k)^T (x - x^k) + \lambda^k \nabla f(x^k)^T (x - x^k) \right)$$

(constant w.r.t x, so add argmin same iter)

$$= \arg \min_x \left(\lambda^k g(x) + \lambda^k \left(f(x^k) + \nabla f(x^k)^T (x - x^k) + \frac{1}{2\lambda} \|x - x^k\|_2^2 \right) \right)$$

$$= \arg \min_x \left(\lambda^k \left(g(x) + \left(f(x^k) + \nabla f(x^k)^T (x - x^k) + \frac{1}{2\lambda} \|x - x^k\|_2^2 \right) \right) \right)$$

(constant w.r.t x, so argmin same iter)

शुद्धी:

$f: D$, continuously differentiable, Lipschitz (continuous with L constant)

$$\Rightarrow \forall x, y \quad 0 \leq f(x) - f(y) - \nabla f(y)^T (x-y) \leq \frac{1}{2} \|x-y\|_2^2$$

$$\Rightarrow f(x) \leq f(y) + \nabla f(y)^T (x-y) + \frac{1}{2} \|x-y\|_2^2 \leq f(y) + \nabla f(y)^T (x-y) + \frac{L}{2} \|x-y\|_2^2$$

now we don't know L in most problems but by tuning a parameter $\lambda > 0$

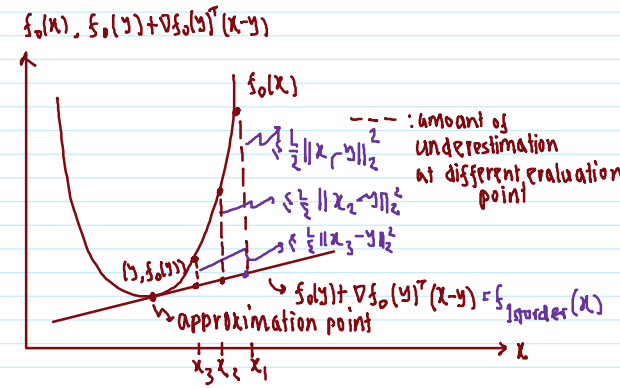
if we can arrive at $f(x) \leq f(y) + \nabla f(y)^T (x-y) + \frac{1}{2\lambda} \|x-y\|_2^2$ then $\frac{1}{\lambda} = L \Rightarrow L$

$$\rightarrow \lambda \leq \frac{1}{L}$$

$$0 < \lambda \leq \frac{1}{L}$$

$$\therefore \forall \lambda \in (0, \frac{1}{L}] \quad f(x) \leq f(y) + \nabla f(y)^T (x-y) + \frac{1}{2\lambda} \|x-y\|_2^2 = \hat{f}_\lambda(x, y)$$

Note that if $\text{del} f$ is globally Lipschitz, then finding one valid λ is sufficient, however when f is not so, we have to modify it. Now suppose $\text{del} f$ is not globally Lipschitz, even then we can use the condition above locally. Assume $\text{del} f$ is locally Lipschitz now which implies the derivative is continuous in the relevant domain; this is a much more realistic condition. So, no matter what we will have a local Lipschitz constant as we move from one iterate to the next, Beck line search essentially does just that: finding an approximation of that local Lipschitz constant by tuning λ at each iterate. On the assumption that we have not moved too much from the previous iterate, that is why we apply previous λ to see if we got lucky and that λ still produces a valid upper bound, if not the only possibility is that we have to make λ smaller thus $1/2 \lambda$ bigger as we move to another point. From a majorization minimization point of view we are tuning λ in such a manner that is making $\text{hat}(f)_\lambda(x, x^k)$ a majorization



$$= \operatorname{argmin}_x \left(g(x) + \underbrace{f(x^k) + \nabla f(x^k)^T (x - x^k) + \frac{1}{2\lambda} \|x - x^k\|_2^2}_{\hat{f}_\lambda(x, x^k)} \right)$$

$$\# \hat{f}_\lambda(x, y) = f(y) + \nabla f(y)^T (x - y) + \frac{1}{2\lambda} \|x - y\|_2^2$$

$$= \operatorname{argmin}_x (g(x) + \hat{f}_\lambda(x, x^k)) = \operatorname{argmin}_x (q_\lambda(x, x^k))$$

Lets define $g(x) + \hat{f}_\lambda(x, y) = q_\lambda(x, y)$ clearly this is a majorization of $g(x) + f(x)$ $\forall y \forall x \forall \lambda \in (0, \frac{1}{L}]$

$$\forall \lambda \in (0, \frac{1}{L}] \forall y \forall x \quad f(x) \leq \hat{f}_\lambda(x, y) \wedge f(y) = \hat{f}_\lambda(y, x)$$

$$\forall \lambda \in (0, \frac{1}{L}] \forall y \forall x \quad g(x) + f(x) \leq g(x) + \hat{f}_\lambda(x, y) \wedge g(y) + f(y) = g(y) + \hat{f}_\lambda(y, x) \quad \# \text{Both are same term add } f(x)$$

$$= q_\lambda(x, y) \quad = q_\lambda(y, x)$$

$$\therefore x^{k+1} = \operatorname{prox}_{\lambda^k g} (x^k - \lambda^k \nabla f(x^k)) = \operatorname{argmin}_x q_{\lambda^k}(x, x^k)$$

Where to get a valid majorization function at each update
 λ = Beck-Teboulle proximal gradient updater(x^k, λ , β) if L is not known, or
 if L is known $\lambda \in (0, 1/L]$ will work and

majorization of $g(x) + f(x)$

\therefore So, proximal algorithm is a majorization minimization algorithm

[Forward-backward splitting]

* Fixed point iteration; # This also works as a proof that proximal algorithm indeed finds the optimal solution of minimize $f(x) + g(x)$

$$x^* = \operatorname{argmin}_x f(x) + g(x) \Leftrightarrow \partial(f(x) + g(x))_{x=x^*} \ni 0 \Leftrightarrow \nabla f(x^*) + \partial g(x^*) \ni 0$$

$$\partial f(x^*) + \partial g(x^*) = \nabla f(x^*) + \partial g(x^*) \quad \# f: \text{differentiable}$$

$$\Leftrightarrow \forall_{\lambda > 0} \lambda \nabla f(x^*) + \lambda \partial g(x^*) \ni 0 \Leftrightarrow \underbrace{-(1 - \lambda \nabla f)}_{\text{vector}}(x^*) + \underbrace{(1 + \lambda \partial g)}_{\text{set}}(x^*) \ni 0$$

$$\# \lambda \nabla f(x^*) - x^* + x^* + \lambda \partial g(x^*) \Leftrightarrow \exists \eta \in (1 + \lambda \partial g)(x^*) \quad -(1 - \lambda \nabla f)(x^*) + \eta = 0$$

$$\Leftrightarrow \exists \eta \in (1 + \lambda \partial g)(x^*) \quad \eta = (1 - \lambda \nabla f)(x^*)$$

$$\Leftrightarrow (1 + \lambda \partial g)(x^*) \ni (1 - \lambda \nabla f)(x^*) = z$$

$$\# R(x) \ni y \Leftrightarrow x \in R^{-1}(y)$$

$$\# \cdot (R(x) \ni y) // R^{-1} = R^{-1}R(x) \ni R^{-1}y = x \in R^{-1}y$$

$$\{ \begin{aligned} \exists x \in \mathbb{R}^n \text{ s.t. } (R(x) \partial y) \cap R^{-1} &= R^{-1} R(x) \partial y = \epsilon R^{-1} y \\ &= x \in R^{-1} y \end{aligned} \}$$

$$x^* \in (I + \lambda \partial g)^{-1} (I - \lambda \nabla f)(x^*) = \left(x^* \parallel (I - \lambda \nabla f(\square)) : \text{a vector} \parallel (I + \lambda \partial g(\square))^{-1} \square = \text{singleton} \right)$$

but we know that resolvent of subdifferential relation is one-to-one, i.e., its a function, so x^* is the only element of $(I + \lambda \partial g)^{-1} (I - \lambda \nabla f)(x^*)$

thm: proximal operator is the resolvent of subdifferential operator

$$x^* = (I + \lambda \partial g)^{-1} (I - \lambda \nabla f)(x^*)$$

This is the key theorem says that the optimal solution of $f(x) + g(x)$ can be found by solving the fixed point equation: $\square = (I + \lambda \partial g(\square))^{-1} (I - \lambda \nabla f(\square)) \square$ @

also note $\text{prox}_{\lambda g}(x) = (I + \lambda \partial g)^{-1}(x)$

$$= (I + \lambda \partial g)^{-1} (x^* - \lambda \nabla f(x^*)) \quad \# (I + \lambda \partial g)^{-1}(x^*) = \text{prox}_{\lambda g}(x^*)$$

$$= \text{prox}_{\lambda g}(x^* - \lambda \nabla f(x^*))$$

$$x^* = \underset{\text{differentiable}}{\text{argmin}} (f(x) + g(x)) \Leftrightarrow x^* = \underset{\text{nondifferentiable}}{\text{prox}}_{\lambda g}(x^* - \lambda \nabla f(x^*)) = \underbrace{(I + \lambda \partial g)^{-1} (I - \lambda \nabla f)}_{\text{forward-backward operator}}(x^*)$$

[Forward-Backward Version of Proximal Gradient Method] [Forward-backward splitting]

$\therefore x^*$ minimizes $f+g \Leftrightarrow$

x^* = fixedpoint (forward-backward operator)

$\lambda \in (0, \frac{1}{L}] \Rightarrow$ forward-backward operator is averaged \Rightarrow iteration converges to a fixed point.

\downarrow
Lipschitz
constant for ∇f

Arriving at the proximal gradient algorithm from this theorem $x^{**} = (I + \lambda \partial g)^{-1} (I - \lambda \nabla f)$
 We know for any monotone operator f , if $x^{**} = f(x^{**})$, then x^{**} can be found by following iteration:
 $x^{**} = (I + \lambda \partial g)^{-1} (I - \lambda \nabla f)$ % Assuming the forward-backward operator is monotone then the following iteration will converge to x^{**}
 $x^{**} = (I + \lambda \partial g)^{-1} (I - \lambda \nabla f)$
 $= \text{prox}_{\lambda g}((x^{**} + \lambda \nabla f(x^{**})))$

4.3. Accelerated proximal gradient method:

* Accelerated Proximal Gradient Method:

In a proximal gradient algorithm we do the following:
 current iterate

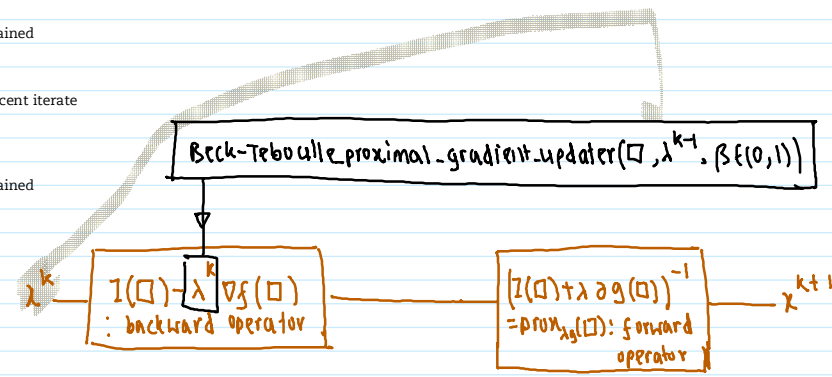
In a proximal gradient algorithm we do the following:

```
current iterate
// find the gradient descent scheme next iterate for the smooth function as if the non-smooth function does not exist and we are doing an unconstrained
optimization of that smooth function using gradient descent
// now we proximize that gradient descent iterate over the nonsmooth function scaled by the step size = next iterate
```

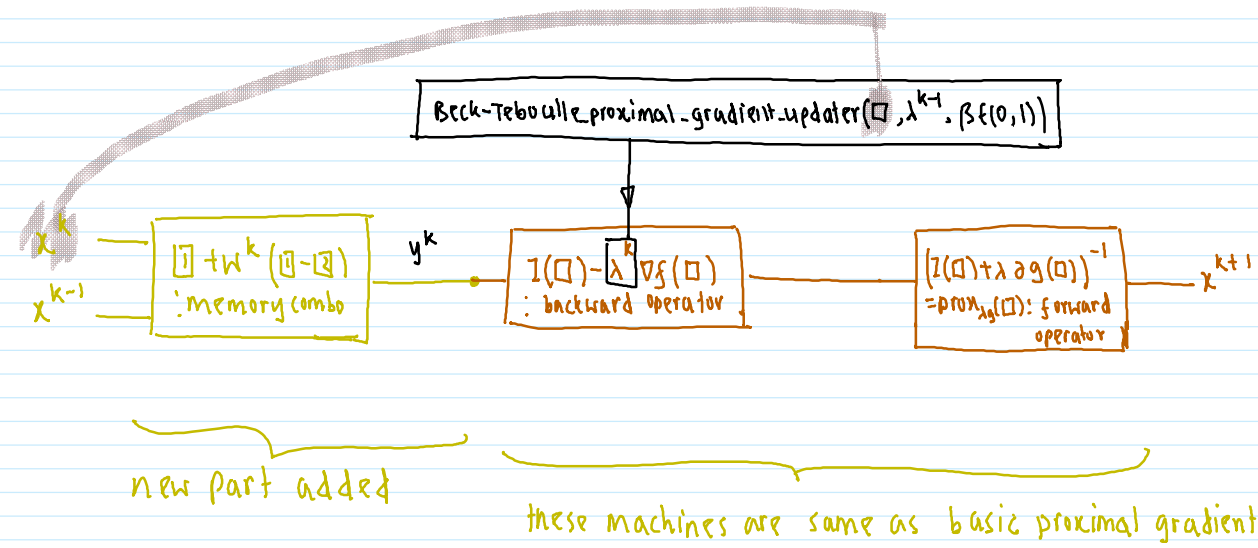
Accelerated proximal gradient algorithm is similar to proximal gradient step, except it needs more memory and takes proximates the gradient descent iterate of an intermediate point (which is a linear combination of two previous iterates)

```
current iterate and the previous iterate
// take a specific linear combination of them = pseudo current iterate
// find the gradient descent scheme next iterate for the smooth function as if the non-smooth function does not exist and we are doing an unconstrained
optimization of that smooth function using gradient descent
// now we proximize that gradient descent iterate over the nonsmooth function scaled by the step size = next iterate
```

in block diagram: (normal) proximal gradient method:



in accelerated proximal gradient method:



Mathematically:

Accelerated proximal gradient update scheme works as follows:

$$\begin{aligned}
 & (x^k, x^{k-1}) // \underbrace{Q[1] + W^k(Q[1] - Q[2])}_{\text{intermediate iterate}} = y^{k+1} : W^k = \frac{k}{k+3} // \underbrace{Q - \lambda^k \nabla f(Q)}_{\text{gradient descent scheme applied to the intermediate point}} : \lambda^k : \text{Beck-Teboulle proximal gradient updater}(y^k, \lambda^{k-1}, \beta \in (0,1)) // \text{prox}_{\lambda^k g}(Q)
 \end{aligned}$$

Convergence:

$(\nabla f$: Lipschitz continuous, Lipschitz constant L), $\lambda^k \in (0, \frac{1}{L}] \vee \lambda^k = \text{Beck-Teboulle-proximal-gradient-updater}(y^k, \lambda^{k-1}, \beta \in (0,1))$

\Rightarrow converges with $O(\frac{1}{k^2})$

4.4. Alternating direction method of multipliers ADMM from Proximal Algorithm

$\forall f(x)+g(z) \quad \# \{f, g\}: \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}, [D]_{\text{proper}}$
 $= \forall f(x)+g(z) \quad \gamma \quad x-z=0$
{closed proper convex function} \cong epi {f, -\gamma, [1, D]}

• Alternating Direction Method of Multipliers:
 (Douglas-Rachford splitting)

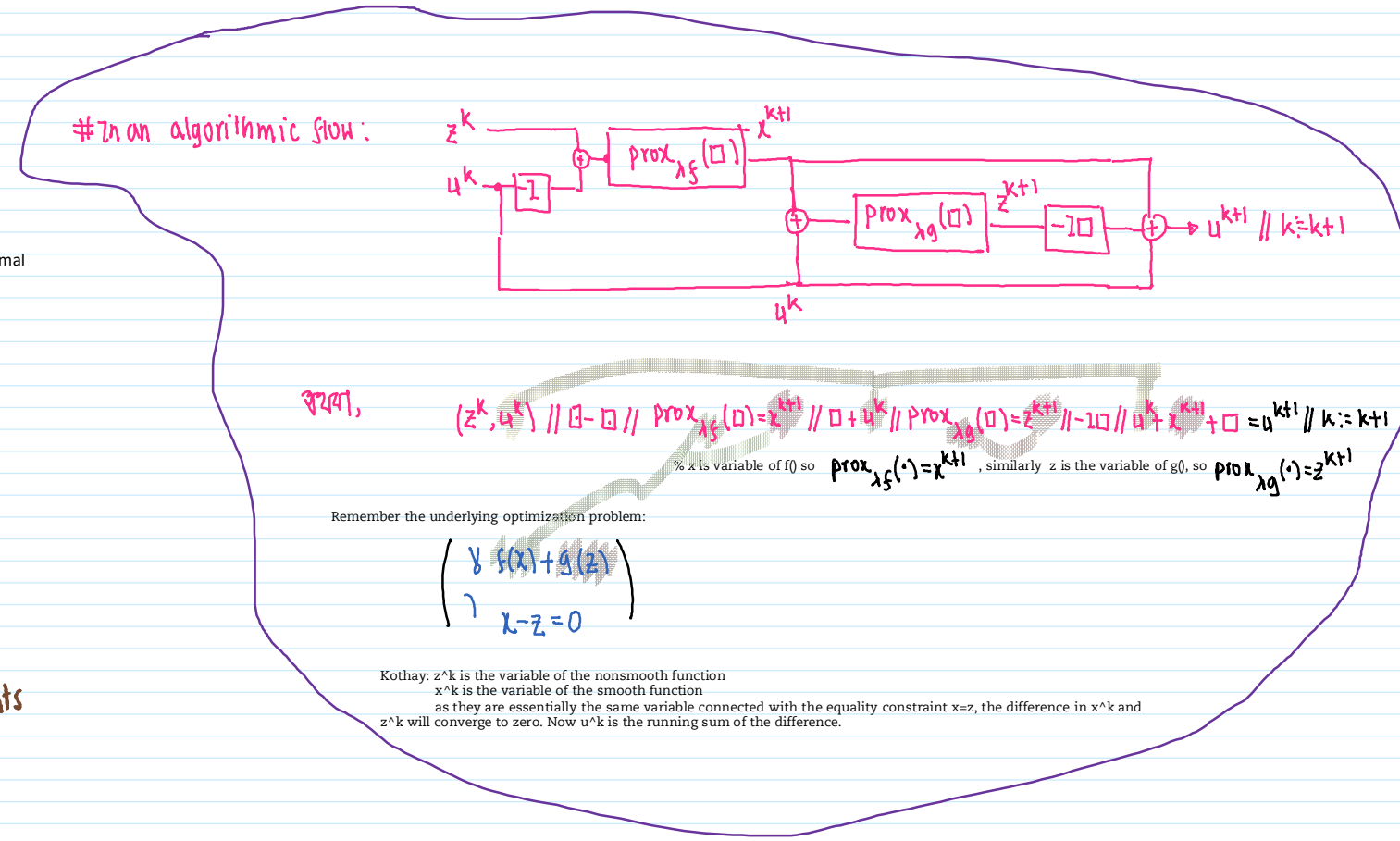
eq: ADMM from prox alg

$x^{k+1} = \text{prox}_{\lambda f}(z^k - u^k)$
 $z^{k+1} = \text{prox}_{\lambda g}(x^{k+1} + u^k)$
 $u^{k+1} = u^k + x^{k+1} - z^{k+1}$

Anthropomorphized:
 ADMM has 3 iterates
 u_k is the running sum of the errors in the 1st iterate and the second iterate
 both the 1st and 2nd iterates should become same as iterations progress i.e., they become the optimal solution

in ADMM $x_k \rightarrow x^*$
 $z_k \rightarrow z^*$
 $u_k \rightarrow 0$

Difference between x^k and z^k :
 $x^k = \arg \min f \rightarrow x^k \in \text{dom} f$
 $z^k = \arg \min g \rightarrow z^k \in \text{dom} g$
 satisfy the constraints as g enforces constraints
 x^k satisfies constraints only in the limit



• Advantages of ADMM:

- * The objective terms are handled separately (functions are accessed only through their proximal operators)
- most useful proximal operators of f, g easy to proximal map but f+g is not easy to evaluate.

• ADMM when $f = I_C(x), g = I_D(x)$

$\forall f(x)+g(x) \cong \begin{cases} x \in C \cap D \end{cases}$

$\Pi_{C \cap D}(x) = \Pi_C(\Pi_D(x))$

then the ADMM algorithm becomes:

$x^{k+1} = \Pi_C(z^k - u^k)$
 $z^{k+1} = \Pi_D(x^{k+1} + u^k)$

note that the parameter λ does not appear

Previously in proximal gradient method we handled both of the functions in one forward backward operator action

$$\begin{aligned} x^{k+1} &= \Pi_C(z^k - u^k) \\ z^{k+1} &= \Pi_D(x^{k+1} + u^k) \\ u^{k+1} &= u^k + x^{k+1} - z^{k+1} \end{aligned}$$

note that the parameter λ does not appear

in practice the iteration is much faster.

* Interpretation of ADMM: 1) Integral control of a dynamical system:

$$\begin{aligned} x^{k+1} &= \text{prox}_{\lambda f}(z^k - u^k) \\ z^{k+1} &= \text{prox}_{\lambda g}(x^{k+1} + u^k) \\ u^{k+1} &= u^k + x^{k+1} - z^{k+1} \end{aligned}$$

State z
Control u
output to be tracked: x which we want to drive to z the state

$\Rightarrow e^{k+1} = x^{k+1} - u^{k+1}$ is the error signal

$$\begin{aligned} u^{k+1} &= u^k + e^{k+1} = u^{k-1} + e^k + e^{k+1} = u^{k-2} + e^{k-1} + e^k + e^{k+1} \\ u^k &= u^{k-1} + e^k \\ u^{k-1} &= u^{k-2} + e^{k-1} \end{aligned}$$

$$u^k = u^1 + \sum_{i=2}^k e^i$$

running sum of the error signals

2) Augmented Lagrangians:
 $(\lambda f(x) + g(x))$

$$L = \begin{pmatrix} \lambda f(x) + g(z) \\ x - z = 0 \end{pmatrix} \text{ # consensus form}$$

analogous to $u(t) = u(1) + \int_1^t e(z) dz = \int \text{integral control method}$
 \downarrow
 where the error signal is driven to zero by feeding back the integral of error to its input.

[Augmented Lagrangian] = $L_p(x, z, y) = f(x) + g(z) + y^T(x - z) + \frac{\rho}{2} \|x - z\|_2^2$

Labels: x (primal variable (1)), z (primal variable (2)), y (dual variable), $\frac{\rho}{2} \|x - z\|_2^2$ (additional quadratic penalty term)

Classical form of ADMM:

$$x^{k+1} := \underset{x}{\text{argmin}} L_p(x, z^k, y^k) \quad // \text{ } L_p \text{ is minimized over the most recent values of the other primal variable and dual variable } (y^k)$$

$$z^{k+1} := \underset{z}{\text{argmin}} L_p(x^{k+1}, z, y^k) \quad // \text{ primal variable } (z^{k+1}) \text{ and dual variable } (y^k) \text{ # determined in the}$$

$$z^{k+1} := \underset{z}{\operatorname{argmin}} L_p(x, z, y^k) \quad \text{primal variable } (x^{k+1}) \text{ and dual variable } (y^k)$$

$$y^{k+1} := y^k + \rho(x^{k+1} - z^{k+1}) \quad \text{running sum of the consensus error}$$

$\parallel y^{k+1} = y^k + \rho \sum_{i=2}^{k+1} e^i$, because:

$$\parallel y^{k+1} = y^k + \rho(x^{k+1} - z^{k+1}) = y^k + \rho e^{k+1} = y^{k-1} + \rho(e^k + e^{k+1}) = y^1 + \rho \sum_{i=2}^{k+1} e^i$$

$$\parallel y^k = y^{k-1} + \rho e^k$$

$$\parallel y^2 = y^1 + \rho e^2$$

#determined in the #previous step

How the classical version of ADMM leads to the proximal version of ADMM?

$$x^{k+1} = \underset{x}{\operatorname{argmin}} L_p(x, z^k, y^k) = \underset{x}{\operatorname{argmin}} (f(x) + g(z^k) + (y^k)^T(x - z^k) + \frac{\rho}{2} \|x - z^k\|_2^2)$$

now we are goal is to write the entire term as $\|x + \square\|_2^2 + \square$ some constant form

$$f(x) + (y^k)^T x - (y^k)^T z^k + \frac{\rho}{2} \|x - z^k\|_2^2 = f(x) + (y^k)^T x - (y^k)^T z^k + \frac{\rho}{2} (x^T x + (z^k)^T z^k - 2x^T z^k)$$

$$(x - z^k)^T (x - z^k) = x^T x + (z^k)^T z^k - 2x^T z^k$$

$$= \underset{x}{\operatorname{argmin}} (f(x) + \frac{\rho}{2} \|x - z^k + \frac{1}{\rho} y^k\|_2^2 + \text{constant w.r.t } x)$$

(can be dropped and argmin will remain the same)

$$= \underset{x}{\operatorname{argmin}} (f(x) + \frac{\rho}{2} \|x - z^k + \frac{1}{\rho} y^k\|_2^2)$$

$$\Leftrightarrow x^{k+1} = \underset{x}{\operatorname{argmin}} (f(x) + \frac{\rho}{2} \|x - z^k + \frac{1}{\rho} y^k\|_2^2)$$

lets take, $u^k = \frac{1}{\rho} y^k$, and $\lambda = \frac{1}{\rho}$

$$= \underset{x}{\operatorname{argmin}} (f(x) + \frac{1}{2} \lambda \|x - z^k + u^k\|_2^2)$$

$$\underset{x}{\operatorname{prox}}_{\lambda f} (z^k - u^k) = \underset{f_a}{\operatorname{argmin}} (\lambda f(f_a) + \frac{1}{2} \|f_a - (z^k - u^k)\|_2^2) = \underset{f_a}{\operatorname{argmin}} \lambda (f(f_a) + \frac{1}{2\lambda} \|f_a - z^k + u^k\|_2^2) = \underset{f_a}{\operatorname{argmin}} (f(f_a) + \frac{1}{2\lambda} \|f_a - z^k + u^k\|_2^2)$$

$$\underset{x}{\operatorname{prox}}_{\lambda f} (x) = \underset{z}{\operatorname{argmin}} (h(z) + \frac{1}{2} \|z - x\|_2^2)$$

because constant term common

Calculation trick: no one will admit: (treat the vectors as numbers first, write them as $(x + \square)^2 + \square$ and then backcalculate as vectors:

\square Treat the vectors as numbers and find $(x + \square)^2 + \square^2$

$$\begin{aligned} \tilde{z}(x) &= yx - yz + \frac{\rho}{2} (x^2 + z^2 - 2xz) \\ &= \frac{\rho}{2} (x^2 + z^2 - 2xz + \frac{2}{\rho} yx - \frac{2}{\rho} yz) \\ &= \frac{\rho}{2} (x^2 + 2(\frac{1}{\rho} y - z)x + (\frac{1}{\rho} y - z)^2 - (\frac{1}{\rho} y - z)^2 + z^2 - \frac{2}{\rho} yz) \\ &= \frac{\rho}{2} \left((x + \frac{1}{\rho} y - z)^2 + \underbrace{(-(\frac{1}{\rho} y - z)^2 + z^2 - \frac{2}{\rho} yz)}_{\text{constant}} \right) \end{aligned}$$

part 2 vector extension:

$$(y^k)^T x - (y^k)^T z + \frac{\rho}{2} (\|x - z^k\|_2^2) = \frac{\rho}{2} (\|x + \frac{1}{\rho} y^k - z^k\|_2^2) + \text{constant}$$

$$\text{prox}_h(x) = \underset{z}{\text{argmin}} \left(h(z) + \frac{\lambda}{2} \|z-x\|_2^2 \right)$$

$$= \text{prox}_{\lambda f}(z^k - u^k)$$

because constant term common
minimizer change not, so λ can
be dropped

$$x^{k+1} = \text{prox}_{\lambda f}(z^k - u^k) \quad [\text{eq: } x^{k+1} \text{ ADMM classic to proximal}]$$

Similarly,

$$z^{k+1} = \underset{z}{\text{argmin}} L_p(x^{k+1}, z, y^k) = \underset{z}{\text{argmin}} \left(f(x^{k+1}) + g(z) + (y^k)^T (x^{k+1} - z) + \frac{\rho}{2} \|x^{k+1} - z\|_2^2 \right)$$

scalarization

$$= \underset{z}{\text{argmin}} \left(g(z) - (y^k)^T z + \frac{\rho}{2} \|x^{k+1} - z\|_2^2 \right) = \underset{z}{\text{argmin}} \left(g(z) + \frac{\rho}{2} \|z - x^{k+1} + \frac{1}{\rho} y^k\|_2^2 + \text{constant} \right) = \underset{z}{\text{argmin}} \left(g(z) + \frac{\rho}{2} \|z - x^{k+1} - u^k\|_2^2 \right) \quad [u^k = \frac{1}{\rho} y^k, \lambda = \frac{1}{\rho}]$$

can drop

$$= \underset{z}{\text{argmin}} \frac{1}{\lambda} \left(\lambda g(z) + \frac{\lambda}{2} \|z - (x^{k+1} + u^k)\|_2^2 \right) = \underset{z}{\text{argmin}} \left(\lambda g(z) + \frac{\lambda}{2} \|z - (x^{k+1} + u^k)\|_2^2 \right)$$

can drop it as argmin will stay the same

$$= \text{prox}_{\lambda g}(x^{k+1} + u^k) \quad [\text{prox}_h(x) = \underset{z}{\text{argmin}} (h(z) + \frac{\lambda}{2} \|z-x\|_2^2)]$$

#

$$(x^{k+1} - z)^T (x^{k+1} - z) = \|x^{k+1}\|_2^2 + \|z\|_2^2 - 2z^T x^{k+1}$$

$$- (y^k)^T z + \frac{\rho}{2} (\|x^{k+1}\|_2^2 + \|z\|_2^2 - 2z^T x^{k+1})$$

$$= \frac{\rho}{2} (\|x^{k+1}\|_2^2 + \|z\|_2^2 - 2z^T x^{k+1} - \frac{2}{\rho} (y^k)^T z)$$

scalarization, counter dropping

$$x^2 + z^2 - 2zx - \frac{2}{\rho} yz = z^2 - 2(x + \frac{1}{\rho} y)z + (x + \frac{1}{\rho} y)^2 - \frac{(x + \frac{1}{\rho} y)^2 + x^2}{\text{constant}} = (z - x - \frac{1}{\rho} y)^2 + \text{constant}$$

vectorization, counter adding

$$\|z - x - \frac{1}{\rho} y\|_2^2 + \text{constant}$$

$$= \frac{\rho}{2} (\|z - x - \frac{1}{\rho} y\|_2^2 + \text{constant})$$

$$= \frac{\rho}{2} (\|x^{k+1} - z + \frac{1}{\rho} y^k\|_2^2 + \text{constant}) \quad \text{return}$$

$$= \frac{\rho}{2} (\|x^{k+1} - z + \frac{1}{\rho} y^k\|_2^2 + \text{constant})$$

$$\therefore z^{k+1} = \text{prox}_{\lambda g} (x^{k+1} + u^k) \quad [\text{eq: } z^{k+1} \text{ ADMM classic to proximal}]$$

finally,

$$y^{k+1} := y^k + \rho(x^{k+1} - z^{k+1})$$

$$\Leftrightarrow \frac{1}{\rho} y^{k+1} := \frac{1}{\rho} y^k + (x^{k+1} - z^{k+1})$$

$$\Leftrightarrow u^{k+1} = u^k + (x^{k+1} - z^{k+1})$$

[eq: y^{k+1} => u^{k+1} ADMM classic to proximal]

[eq: x^{k+1} ADMM classic to proximal], [eq: z^{k+1} ADMM classic to proximal], [eq: y^{k+1} => u] } (येसू त्रुवतु अतः ADMM classic form to proximal form both are equivalent. ☺)

Interpretation 3: (Fixed point iteration)

Kind of shows the proof that ADMM works: We are going to show that if $x^k \sim x^*$, $z^k \sim z^*$, $u^k \sim u^*$ then the x^* term will be the minimizer of $f(x)+g(x)$

$$\forall (f(x)+g(x)) \circ x^* \text{ is optimal} \Leftrightarrow 0 \in \partial f(x^*) + \partial g(x^*)$$

now

$$\begin{cases} x^{k+1} = \text{prox}_{\lambda f} (z^k - u^k) \\ z^{k+1} = \text{prox}_{\lambda g} (x^{k+1} + u^k) \\ u^{k+1} = u^k + x^{k+1} - z^{k+1} \end{cases}$$

$$\begin{cases} x^* = \text{prox}_{\lambda f} (z^* - u^*) \\ z^* = \text{prox}_{\lambda g} (x^* + u^*) \\ u^* = u^* + x^* - z^* \Leftrightarrow x^* = z^* \end{cases}$$

$$\begin{bmatrix} x^* \\ z^* \\ u^* \end{bmatrix} = \begin{bmatrix} \text{prox}_{\lambda f} (\begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} - \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}) \\ \text{prox}_{\lambda g} (\begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} + \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}) \\ \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} + \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} - \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} \end{bmatrix} \begin{bmatrix} x^* \\ z^* \\ u^* \end{bmatrix}$$

$$\begin{cases} x^* = \text{prox}_{\lambda f} (x^* - u^*) = (1 + \lambda \partial f)^{-1} (x^* - u^*) \\ z^* = \text{prox}_{\lambda g} (x^* + u^*) = (1 + \lambda \partial g)^{-1} (x^* + u^*) \\ \Downarrow \\ x^* = z^* \end{cases}$$

Here $(1 + \square \partial \square)^{-1}$ is the resolvent operator of the proximal mapping of $\square \square (\cdot)$ at some point, i.e., $\text{prox}_{\square \square} (\square) = (1 + \square \partial \square)^{-1} (\square)$
 # see notes on [1.2 Boyd] Interpretation of proximal mapping

$$\Leftrightarrow \begin{cases} x^* - u^* \in (1 + \lambda \partial f) x^* = x^* + \lambda \partial f(x^*) \\ x^* + u^* \in (1 + \lambda \partial g) x^* = x^* + \lambda \partial g(x^*) \end{cases}$$

remember $(1 + \lambda \partial f)$ is a relation but $(1 + \lambda \partial f)^{-1}$ is a function

$$\Leftrightarrow \begin{cases} x^* - u^* \in x^* + \lambda \{\partial f\} = \{x^* + \lambda \partial f\} \\ x^* + u^* \in x^* + \lambda \{\partial g\} = \{x^* + \lambda \partial g\} \end{cases}$$

$$\lfloor u^* \rfloor \quad \lfloor \theta + \theta - \theta \rfloor \quad \lfloor u^* \rfloor$$

$$x^* + u^* \in x^* + \lambda \{ \theta; \} = \{ x^* + \lambda \theta; \}$$

$$\Leftrightarrow \begin{cases} \exists \theta_2 \in \partial f(x^*) & x^* - u^* = x^* + \lambda \theta_2 \\ \exists \theta_1 \in \partial g(x^*) & x^* + u^* = x^* + \lambda \theta_1 \end{cases}$$

$$\stackrel{(*)}{\Rightarrow} 2x^* = 2x^* + \lambda(\theta_2 + \theta_1)$$

$\in \partial(f(x^*) + g(x^*)) \quad \{ \forall_x \partial f(x) + \partial g(x) = \partial(f(x) + g(x)) \}$

$$\Leftrightarrow 0 = \lambda(\theta_2 + \theta_1)$$

$$\Leftrightarrow 0 \in \partial(f(x^*) + g(x^*))$$

