HW3 1a). $\min_{x} f(x)$    s.t. $Ax = b$

For convex and closed function $f$, there is $f^{**} = f$

$x \in \partial f^*(y) \implies y \in \partial f^{**}(x) \implies y \in \partial f(x)$.

$y \in \partial f(x) \implies f(u) - f(x) \geq y^T (u - x), \quad \forall u$.

$\implies y^T u - f(u) \leq y^T x - f(x), \quad \forall u$.

$\implies f^*(y) = \max_{u} (y^T u - f(u)) = y^T x - f(x) \iff x = \text{argmin}_{u} f(u) - y^T u$

also let $f^*(v) = \max_{u} (v^T u - f(u)) \geq f^*(y) + x^T(v - y) \implies x \in \partial f^*(y)$

The dual problem. $\max_{\lambda} g(\lambda) = -f^*(-A^T \lambda) - \lambda^T b$

$$\frac{\partial g(\lambda)}{\partial \lambda} = \frac{A \partial f^*(-A^T \lambda)}{\partial \lambda} - b$$

Since for convex and close function $f$

$x \in \partial f^*(y) \iff y \in \partial f(x) \iff x = \text{arg min}_{u} f(u) - y^T u$.

$\implies w \in \partial f^*(-A^T \lambda) \iff w = \text{argmin}_{u} f(u) + \lambda^T A u$.

$\implies Aw - b \in \partial g(\lambda) \iff w = \text{arg min}_{w} f(u) + \lambda^T A w$

Thus, $Aw - b \in \partial g(\lambda)$ for $w = \text{arg min}_{u} f(u) + \lambda^T A w$.

HW3(b)  min $f(w)$
s.t. $Aw = b$

since $w_{k+1} \in \arg\min_w L(w, \lambda_k)$, where

$$L(w, \lambda_k) = f(w) + \lambda_k^T (Aw - b)$$

Denote the primal optimal variable by $w^*$

$$L(w^*, \lambda_k) - L(w_{k+1}, \lambda_k)$$

$$= f(w^*, \lambda_k) - f(w_{k+1}, \lambda_k) + \lambda_k^T (Aw^* - b) - \lambda_k^T (Aw^* - b)$$

$$\geq \nabla f(w_{k+1})^T (w^* - w_{k+1}) + \frac{\mu}{2} \|w^* - w_{k+1}\| + \lambda_k^T A (w^* - w_{k+1})$$

where the above inequality follows from $f$ is $\mu$-strongly convex.

since $w_{k+1} \in \arg\min_w L(w, \lambda_k)$

$$\nabla L(w_{k+1}, \lambda_k) = 0 \Rightarrow f(w_{k+1}) + \lambda_k^T A = 0$$

Thus we have

$$L(w^*, \lambda_k) - L(w_{k+1}, \lambda_k) \geq \frac{\mu}{2} \|w^* - w_{k+1}\|^2$$

in another word

$$\|w^* - w_{k+1}\|^2 \leq \frac{2(L(w^*, \lambda_k) - L(w_{k+1}, \lambda_k))}{\mu}$$

It can be seen that the convergence and accuracy of primal can be controled by dual variable.

HW3 (b) - Convergence analysis of dual variable.

$$\|\lambda^{k+1} - \lambda^{\#}\|_2^2 = \|\lambda^k + \alpha_k (Aw_k - b) - \lambda^{\#}\|_2^2 =$$

$$= \|\lambda^k - \lambda^{\#}\|_2^2 + 2\alpha_k \langle Aw_k - b, \lambda^k - \lambda^{\#}\rangle + \alpha_k^2 \|Aw_k - b\|_2^2$$

$f: u$-strong $\overset{convex}{L\text{-smooth}} \Rightarrow g \; \frac{1}{L}$-strong convex and $\frac{1}{u}$-smooth

there is, $\|\lambda^{k+1} - \lambda^{\#}\|_2^2 \le (1 - 2\alpha_k \frac{1}{L}) \|\lambda^k - \lambda^{\#}\|_2^2 - 2\alpha_k(g(\lambda^k) - g\lambda^{\#})$

$$+ \alpha_k^2 \|Aw_k - b\|_2^2$$

( By using $\frac{1}{L}$-strong convexity).

Further

Also, there is $\|\lambda^{k+1} - \lambda^{\#}\|_2^2 \le (1 - \alpha_k/L) \|\lambda^k - \lambda^{\#}\|_2^2 - 2\alpha_k(g(\lambda^k) - g\lambda^{\#})$

$$+ 2\alpha_k^2 \frac{1}{u}(g(\lambda^k) - g(\lambda^{\#})),$$

$$= (1 - \alpha_k/L) \|\lambda^k - \lambda^{\#}\|_2^2 - 2\alpha_k(1 - \frac{\alpha_k}{u})(g(\lambda^k) - g\lambda^{\#})$$

Convergence guaranteed when

$0 < 1 - \frac{\alpha_k}{L} < 1$ ~~and~~ $1 - \frac{\alpha_k}{u} \ge 0 \Rightarrow \alpha_k \le u$

and $0 < 1 - \frac{\alpha_k}{L} < 1$ $\qquad$ since $\alpha_k \le u$

when then require $L > u$

And the convergence rate is thus $\frac{u}{L}$

## HW. 3. c

$$\text{minimize} \quad \sum_{i=1}^{N} f^i(\omega^i)$$
$$\omega^1, \ldots, \omega^N \in W$$

$$\text{s.t.} \quad \omega^1 = \omega^2 = \cdots = \omega^N$$

where there are $N$ numbers of nodes.

- Let $\omega_j^i$ represents node $i$'s estimate of nodes $j$'s internal state.

- Let $\lambda_{ij}$, $i, j \in \{1, 2, \ldots, N\}$, be the dual variables. For instance $\lambda_{ij}$ is associated with the constraint
$$\omega_i^i = \omega_i^j$$

Let
$$\lambda = \left[ \lambda_{11}^T, \cdots, \lambda_{1N}^T, \lambda_{21}^T, \cdots, \lambda_{2N}^T, \cdots, \lambda_{NN}^T \right]^T$$

The Langrangian :
$$L\left( \omega^1, \omega^2, \cdots, \omega^N, \lambda \right) = \sum_{i=1}^{N} f^i \left( \omega_1^i, \omega_2^i, \cdots, \omega_N^i \right)$$
$$+ \sum_{i=1}^{N} \sum_{j=1}^{N} \lambda_{ij}^T \left( \omega_i^i - \omega_i^j \right)$$

The dual function
$$g(\lambda) = \inf_{\omega^1, \ldots, \omega^N \in W} L\left( \omega^1, \cdots, \omega^N, \lambda \right)$$

We can write the dual function in subproblems such as
$$\phi^i(\lambda) = \inf_{\omega^i \in W} f^i\left( \omega_1^i, \cdots, \omega_N^i \right) + \sum_{j=1}^{N} \lambda_{ij}^T \omega_i^i - \sum_{j=1}^{N} \lambda_{ji}^T \omega_j^i$$

Since $\phi^i(\lambda)$ only depends on $\lambda$ and $\omega^i$, node $i$ can compute $\phi^i(\lambda)$ locally.

The Langrangian is the sum

$$g(\lambda) = \sum_{i=1}^{N} \phi^i(\lambda)$$

Finally, the dual problem is the maximization

$$g^* = \max_{\lambda} g(\lambda) = \max_{\lambda} \sum_{i=1}^{N} \phi^i(\lambda)$$

Communication cost

Suppose that there are $N$ nodes, each part is transferred over $N-1$ nodes. Summing over all parts will give the the communication cost for primal method

$$\sum_{i=1}^{N} (N-1) \dim(w_i^i) = (N-1) \dim(w)$$

Similarly, the communication cost for dual method is

$$\sum_{i=1}^{N} \sum_{j=1}^{N} (N-1) \dim \lambda_{ij} = (N-1) \dim(\lambda)$$

Since $\dim(\lambda) > \dim(w)$, the communication cost of the dual method is higher.

But, one can decrease this amount since a node only needs the sum of dual variables, $\sum_{j=1}^{N} \lambda_{ij}$, for each iteration.