

Convex Optimization. — Coordinate Descent Method

1. Problem 1.

(a). The solution to $Ax=b$. can be equivalently related to the optimizer of the problem. $\min_x \frac{1}{2} x^T A x - b^T x := f(x)$.

For $\mathbf{A} \in \mathbb{S}_{++}^n$, we decompose it as $\mathbf{A} = \mathbf{D} + \mathbf{E} + \mathbf{F}$

Where D is a diagonal matrix, E is strict lower triangle and F is strict upper triangle matrix.

Hence, the modified Gauss-Seidel update can be written in vector form.

$$x^{k+1} = x^k - \gamma D^{-1} [Dx^k + Ex^{k+1} + Fx^k - b].$$

$$\Rightarrow (I + \gamma D^T E) x^{k+1} = [(1-\gamma) I - \gamma D^T F] x^k + \gamma D^T b.$$

$$\begin{aligned}
 \Rightarrow (D + rE)x^{k+1} &= [D + rE - r(E + F + D)]x^k + rb \\
 &= (D + rE)x^k - r(Ax^k - b) \\
 &= (D + rE)x^k - r^k b
 \end{aligned}$$

$$\Rightarrow x^{k+1} = x^k - (\frac{D}{\gamma} + E)^{-1} r^k.$$

$$\begin{aligned}
 f(x^{k+1}) - f(x^k) &= \frac{1}{2} \left[x^k - \left(\frac{D}{f} + E \right)^{-1} r^k \right]^T A \cdot \left[x^k - \left(\frac{D}{f} + E \right)^{-1} r^k \right] \\
 &\quad - \frac{1}{2} b^k)^T A x^k - b^k (x^{k+1} - x^k) \\
 &= -\frac{1}{2} \left[\left(\frac{D}{f} + E \right)^{-1} r^k \right]^T A x^k - \frac{1}{2} (r^k)^T \left(\frac{D^T}{f} + E^T \right)^{-1} A \left(\frac{D}{f} + E \right)^{-1} r^k \\
 &\quad + b^k \left(\frac{D}{f} + E \right)^{-1} r^k \\
 &= \left[\left(\frac{D}{f} + E \right)^{-1} r^k \right]^T [b - A x^k] + \frac{1}{2} (r^k)^T \left(\frac{D^T}{f} + E^T \right)^{-1} A \left(\frac{D}{f} + E \right)^{-1} r^k \\
 &= -(r^k)^T \cdot \left(\frac{D^T}{f} + E^T \right)^{-1} \cdot r^k + \frac{1}{2} (r^k)^T \left(\frac{D^T}{f} + E^T \right)^{-1} A \left(\frac{D}{f} + E \right)^{-1} r^k \\
 &= \frac{1}{2} (r^k)^T \left[\left(\frac{D}{f} + E \right)^T A \cdot \left(\frac{D}{f} + E \right)^{-1} - \frac{1}{2} \left(\frac{D}{f} + E \right)^{-T} \right] r^k
 \end{aligned}$$

We need to show $(r_k) \left[\left(\frac{D}{f} + E \right)^{-1} A \left(\frac{D}{f} + E \right)^{-1} - 2 \left(\frac{D}{f} + E \right)^{-1} \right] (r_k) > 0$ if $\gamma \notin (0, 2)$
 $+ r_k \neq 0$.

$$\begin{aligned} \text{Let } w = \frac{D}{\gamma} + E. \Rightarrow w^T = \frac{D}{\gamma} + F. \\ \Rightarrow A - w - w^T = (1 - \frac{\gamma}{\gamma}) D. \\ (w^T)^T [A - w - w^T] w^T = (1 - \frac{\gamma}{\gamma}) \cdot (w^T)^T D D w^T \\ \Rightarrow (w^T)^T A w^T - (w^T)^T w^T = (1 - \frac{\gamma}{\gamma}) (D^T w^T)^T D^T w^T. \end{aligned}$$

Besides, we have

$$(r^k)^T w^T r^k = (r^k)^T (w^T)^T r^k.$$

Therefore,

$$\begin{aligned} f(x^{k+1}) - f(x^k) &= \frac{1}{2} (1 - \frac{\gamma}{\gamma}) \cdot (r^k)^T (D^T w^T)^T D^T w^T r^k \\ &= \frac{1}{2} (1 - \frac{\gamma}{\gamma}) \|D^T w^T r^k\|_2^2. \end{aligned}$$

If $r \notin \{0, z\}$, then $\nabla r^k \neq 0$.

$$f(x^{k+1}) > f(x^k).$$

Therefore, the sequence doesn't converge if $x_0 \neq x^*$. \square

$$(b) \quad f(x_1, x_2) = \max((x_1-1)^2 + (x_2+1)^2, (-x_1+1)^2 + (x_2-1)^2).$$

First, notice that f is the maximum of convex functions.

hence f is also convex. Any local optimum is global optimum.

Assume (x_1, x_2) is an optimizer.

$$\begin{aligned} f(x_1, x_2) &= \max((-x_1-1)^2 + (-x_2+1)^2, (-x_1+1)^2 + (-x_2-1)^2) \\ &= \max((x_1+1)^2 + (x_2-1)^2, (x_1-1)^2 + (x_2+1)^2) \\ &= f(x_1, x_2) \end{aligned}$$

$\Rightarrow (-x_1, -x_2)$ is also an optimizer.

Therefore $0 = \frac{1}{2}(x_1, x_2) + \frac{1}{2}(-x_1, -x_2)$ is always an optimizer.

Besides, $f_1(x) = (x_1-1)^2 + (x_2+1)^2$, $f_2(x) = (x_1+1)^2 + (x_2-1)^2$ are strong convex functions.

$$\Rightarrow \cancel{f(x) \geq f(x_0) \geq \dots}$$

$f(x)$ is also strong convex function. \Rightarrow Unique minimizer. $\quad (2)$

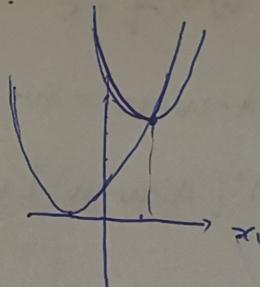
Therefore, the only minimizer is $x^* = (0, 0)^T$.

If we start from $x^0 = (1, 1)$,

$$x_1^1 \in \operatorname{argmin}_{x_1} \max((x_1 - 1)^2 + 4, (x_1 + 1)^2)$$

$$\in \operatorname{argmin}_{x_1} \begin{cases} (x_1 - 1)^2 + 4 & , x \leq 1 \\ (x_1 + 1)^2 & , x > 1 \end{cases}$$

$$= 1.$$



$$x_2^1 \in \operatorname{argmin}_{x_2} \max((x_2 + 1)^2, 4 + (x_2 - 1)^2)$$

$$= 1.$$

Therefore, $(1, 1)$ is a fixed point of the Gauss-Seidel algorithm.

this doesn't converge to $(0, 0)$. \square

2. Problem 2

$$(a) \quad \max_z b^T z - \frac{1}{2} \|A^T z\|_2^2$$

$z \in \mathbb{R}^m$. denote $\tilde{z}^k = (z_1^k, z_2^k, \dots, z_m^k)^T$

$$f(z) = b^T z - \frac{1}{2} \|A^T z\|_2^2 = b^T z - \frac{1}{2} z^T A A^T z$$

$$\nabla f(z) = b - A A^T z \in \mathbb{R}^m$$

Therefore

$$\left. \begin{aligned} z_1^{k+1} &= b - A_1^T A^T [\tilde{z}_1^k, \tilde{z}_2^k, \dots, \tilde{z}_m^k]^T \\ z_2^{k+1} &= b_2 - A_2^T A^T [\tilde{z}_1^{k+1}, \tilde{z}_2^k, \dots, \tilde{z}_m^k]^T \\ &\vdots \\ z_i^{k+1} &= b_i - A_i^T A^T [\tilde{z}_1^{k+1}, \tilde{z}_2^{k+1}, \dots, \tilde{z}_{i-1}^{k+1}, \tilde{z}_i^k, \dots, \tilde{z}_m^k]^T \end{aligned} \right\} .$$

According to coordinate descent.

$$x_i^{k+1} \in \operatorname{argmin}_{x_i} f(x_1^{k+1}, \dots, x_{i-1}^{k+1}, x_i, x_{i+1}^{k+1}, \dots, x_m^{k+1})$$

For each iteration i , we have \tilde{z}_i satisfying that

$$\nabla_i f(\tilde{z})_i = b_i - [A A^T \tilde{z}]_i = 0. \quad b_i - A_i^T A^T \tilde{z}$$

③

where $\tilde{z} = [z_1^{k+1}, \dots, z_{i-1}^{k+1}, z_i^k, z_{i+1}^k, \dots, z_m^k]^T$

$$\Rightarrow \nabla_i f(\tilde{z}) = b_i - a_i^T A^T \tilde{z}$$

$$= b_i - \sum_{j=1}^i a_i^T a_j z_j^{k+1} - \sum_{j=i}^m a_i^T a_j z_j^k.$$

$$z_i^{k+1} = z_i^k + \nabla_i f(\tilde{z}) \quad (\text{Because it's a maximum problem, and we assume step size is 1})$$

$$= b_i - \sum_{j=1}^i a_i^T a_j z_j^{k+1} - \sum_{j=i+1}^m a_i^T a_j z_j^k.$$

Therefore the update rule is

$$z_i^{k+1} = b_i - \sum_{j=1}^i a_i^T a_j z_j^{k+1} - \sum_{j=i+1}^m a_i^T a_j z_j^k \quad (2)$$

For this update, we have $\hat{z} = [z_1^{k+1}, \dots, z_{i-1}^{k+1}, z_i^{k+1}, z_{i+1}^k, \dots, z_m^k]^T$

~~$a_i^T z_i^{k+1} = a_i^T z_i^k$~~

With this new \hat{z} .

$$\nabla_i f(\hat{z}) = b_i - a_i^T A^T \hat{z}$$

$$= b_i - \sum_{j=1}^i a_i^T a_j z_j^{k+1} - \sum_{j=i+1}^m a_i^T a_j z_j^k - z_i^{k+1}$$

$$= 0.$$

(1).

Therefore $a_i^T (A^T \hat{z}) = b_i$.

(b) for the primal problem. $\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{x}\|_2^2$
 $s.t. \mathbf{A}\mathbf{x} = \mathbf{b}$

$$L(\mathbf{x}, \mathbf{z}) = \frac{1}{2} \|\mathbf{x}\|_2^2 + \mathbf{z}^T (\mathbf{A}\mathbf{x} - \mathbf{b})$$

$$d(\mathbf{z}) = \inf_{\mathbf{x}} L(\mathbf{x}, \mathbf{z}) = \inf_{\mathbf{x}} \frac{1}{2} \|\mathbf{x}\|_2^2 + \mathbf{z}^T \mathbf{A}\mathbf{x} - \mathbf{z}^T \mathbf{b}$$

$$\Rightarrow \mathbf{x} = (\mathbf{z}^T \mathbf{A})^T = \mathbf{0} \Rightarrow \mathbf{A}^T \mathbf{z} + \mathbf{x} = \mathbf{0}$$

therefore. $d(\mathbf{z}) = \frac{1}{2} \|\mathbf{A}^T \mathbf{z}\|_2^2 - \mathbf{z}^T \mathbf{A} \mathbf{A}^T \mathbf{z} - \mathbf{z}^T \mathbf{b}$
 $= -\frac{1}{2} \|\mathbf{A}^T \mathbf{z}\|_2^2 - \mathbf{z}^T \mathbf{b}$

$$d(\mathbf{z}') \stackrel{\mathbf{z}' = -\mathbf{z}}{=} -\frac{1}{2} \|\mathbf{A}^T \mathbf{z}'\|_2^2 + \mathbf{b}^T \mathbf{z}'$$

$$\Rightarrow \text{dual problem is } \max_{\mathbf{z}' \in \mathbb{R}^n} d(\mathbf{z}') = \mathbf{b}^T \mathbf{z}' - \frac{1}{2} \|\mathbf{A}^T \mathbf{z}'\|_2^2.$$

therefore, the relation between primal variable \mathbf{x} and the dual variable is that. $\mathbf{A}^T \mathbf{z}' - \mathbf{x} = \mathbf{0}. \quad (3)$

After each update from $\mathbf{z}_i^k \rightarrow \mathbf{z}_i^{k+1}$, with the new $\tilde{\Sigma}_{..}$, we update \mathbf{x} .

From (1) and (3), we know,

At iteration i , the constraint $b_i = \mathbf{a}_i^T \mathbf{x}$ is satisfied. \square

$$(C) \quad g^k = \nabla_{q_t}(y^k) = (\alpha_t^T y^k - b_t) \alpha_t.$$

Hence, the update rule for y^k is.

$$y^{k+1} = y^k - (\alpha_t^T y^k - b_t) \alpha_t \quad \text{if we assume step size is 1.}$$

Notice that, if we choose $t = \text{mod}(k, m) + 1$, starting with $k=0$.

After each update of y^k to y^{k+1} .

we have

$$\begin{aligned} \alpha_t^T y^{k+1} &= \alpha_t^T (y^k - (\alpha_t^T y^k - b_t) \alpha_t) \\ &= \alpha_t^T y^k - \alpha_t^T \alpha_t (\alpha_t^T y^k - b_t) \\ &= b_t. \end{aligned}$$

which means, at iteration $t = \text{mod}(k, m) + 1$,

the constraint $\alpha_t^T y = b_t$ is satisfied.

which is exactly what we have for primal variable's update property in (b).

In (b), the order is $i = 1, 2, \dots, m, 1, 2, \dots$ with $\alpha_i^T x = b_i$ satisfied.

which is exactly matching $t = \text{mod}(k, m) + 1$. while k is update ~~index for~~
time counter for dual variable z (Each update in (2) counts once).

Therefore, SGD with $t = \text{mod}(k, m) + 1$, and unit stepsize has the same update rule as part (b) for primal variable. \square