

Q3 =

$$f(x) = E[F(x; \xi)]$$

$$E[g(x; \xi)] \in \partial f(x)$$

$$x_{t+1} = \pi_{\mathcal{X}}(x_t - \eta_t g(x_t; \xi_t)) \quad t \in \mathbb{N}$$

(1)

$$\frac{1}{2} E_{\xi_1 \dots \xi_t} \left[\|x_t - x^*\|_2^2 + \eta_t^2 \|g(x_t; \xi_t)\|_2^2 - \|x_{t+1} - x^*\|_2^2 \right] \quad (1)$$

$$\geq \frac{1}{2} E_{\xi_1 \dots \xi_t} \left[\|x_t - x^*\|_2^2 + \eta_t^2 \|g(x_t; \xi_t)\|_2^2 - \|x_t - \eta_t g(x_t; \xi_t) - x^*\|_2^2 \right]$$

$$= \frac{1}{2} E_{\xi_1 \dots \xi_t} \left[2 \eta_t g(x_t; \xi_t)^T (x_t - x^*) \right]$$

$$= \eta_t E_{\xi_1 \dots \xi_{t-1}} \left[E_{\xi_t} [g(x_t; \xi_t)]^T (x_t - x^*) \right] \quad (x_t, x^* \text{ independent of } \xi_t)$$

$$\geq \eta_t E_{\xi_1 \dots \xi_{t-1}} \left[f(x_t) - f(x^*) \right] \quad (\text{convexity, } E_{\xi_t} [g(x_t; \xi_t)] \in \partial f(x))$$

$$= \eta_t E_{\xi_1 \dots \xi_{t-1}, \xi_t} \left[f(x_t) - f(x^*) \right] \quad (x_t, x^* \text{ independent of } \xi_t)$$

Note that

$$\frac{1}{2} E_{\xi_1 \dots \xi_t} \left[\|x_t - x^*\|_2^2 - \|x_{t+1} - x^*\|_2^2 \right] + \frac{\eta_t^2}{2} E_{\xi_1 \dots \xi_t} \left[\|g(x_t; \xi_t)\|_2^2 \right]$$

$$= \frac{1}{2} E_{\xi_1 \dots \xi_t} \left[\|x_t - x^*\|_2^2 \right] - \frac{1}{2} E_{\xi_1 \dots \xi_t} \left[\|x_{t+1} - x^*\|_2^2 \right] + \frac{\eta_t^2}{2} L^2 \quad (\text{by assumption})$$

#

(2)

$$E[f(\bar{X}_{t_1:t_2}) - f(x^*)] \leq \frac{1}{\sum_{\tau=t_1}^{t_2} \eta_\tau} \sum_{\tau=t_1}^{t_2} \eta_\tau E[f(X_\tau) - f(x^*)] \quad (\text{Jensen inequality})$$

$$\leq \frac{\sum_{\tau=t_1}^{t_2} (E[\frac{1}{2} \|X_\tau - x^*\|_2^2] - E[\frac{1}{2} \|X_{\tau+1} - x^*\|_2^2] + \frac{1}{2} \eta_\tau^2 L^2)}{\sum_{\tau=t_1}^{t_2} \eta_\tau}$$

$$\leq \frac{E[\frac{1}{2} \|X_{t_1} - x^*\|_2^2] - E[\frac{1}{2} \|X_{t_2} - x^*\|_2^2] + \frac{L^2}{2} \sum_{\tau=t_1}^{t_2} \eta_\tau^2}{\sum_{\tau=t_1}^{t_2} \eta_\tau}$$

$$\leq \frac{E[\|X_{t_1} - x^*\|_2^2 - \|X_{t_2} - x^*\|_2^2] + L^2 \sum_{\tau=t_1}^{t_2} \eta_\tau^2}{2 \sum_{\tau=t_1}^{t_2} \eta_\tau}$$

$$\|X_{t_1} - x^*\|_2 \leq \|X_{t_1} - x_1\|_2 + \|x^* - x_1\|_2$$

$$\leq 2R$$

$$\leq \frac{4R^2 + L^2 \sum_{\tau=t_1}^{t_2} \eta_\tau^2}{2 \sum_{\tau=t_1}^{t_2} \eta_\tau} \quad \forall 1 \leq t_1 \leq t_2 \quad \#$$

(3) notice that

$$\eta_t = R / L \sqrt{t}$$

$$\frac{4R^2 + L^2 \sum_{\tau=t_1}^{t_2} \eta_\tau^2}{2 \sum_{\tau=t_1}^{t_2} \eta_\tau} = \frac{4R^2 + L^2 \times \frac{R^2}{L^2} \sum_{\tau=t_1}^{t_2} \frac{1}{\tau}}{2 \frac{R}{L} \sum_{\tau=t_1}^{t_2} \frac{1}{\sqrt{\tau}}} = \frac{RL(2 + \sum_{\tau=t_1}^{t_2} \frac{1}{\tau})}{\sum_{\tau=t_1}^{t_2} \frac{1}{\sqrt{\tau}}}$$

$$\leq \frac{2RL}{\frac{t_2 - t_1 + 1}{\sqrt{t_2}}} + \frac{RL}{2} \left[\frac{\sum_{\tau=t_1}^{t_2} \frac{t_1}{\tau}}{\sum_{\tau=t_1}^{t_2} \frac{\sqrt{t_1}}{\sqrt{\tau}}} \right] \times \frac{1}{\frac{1}{\sqrt{t_1}}} \leq 2 \times t_2$$

$$\leq \frac{RL}{\frac{t_2 - t_1 + 1}{\sqrt{t_2}}} + \frac{RL}{2} \frac{\sqrt{t_1}}{t_1} = \frac{RL}{\sqrt{t_2}} \left(2 \frac{t_2}{t_2 - t_1 + 1} + \sqrt{\frac{t_2}{t_1}} \right) \quad \#$$

$$Q_1: \quad X = \begin{bmatrix} x_1 & x_2 & \dots & x_p \end{bmatrix} \quad y \approx X \hat{\beta}$$

$$\text{Lasso} = \beta^* \in \underset{\beta}{\operatorname{argmin}} \left\{ \underbrace{\frac{1}{2} \|y - X\beta\|_2^2}_{f(\beta)} \mid \underbrace{\beta \in \mathbb{R}^p, \|\beta\|_1 \leq C}_{\Omega} \right\}$$

$$\nabla f(\beta) = -(y - X\beta)^T X$$

Consider Frank-wolfe, $\beta_0 = 0 \in \mathbb{R}^p$

$$V_0 = \underset{v \in \Omega}{\operatorname{argmin}} \{ \langle -X^T y, v \rangle \mid v \in \Omega \}$$

$$= C \operatorname{sign}(X^T y_{i_0}) e_{i_0}$$

$$\text{where } i_0 \text{ s.t. } |X^T y_{(i_0)}| = \|X^T y\|_\infty$$

$$= \underset{j}{\operatorname{argmax}} \{ |\langle y, x_j \rangle| \mid j \in \{1, 2, \dots, p\} \}$$

Remark = $i_0 = j_0$. from forward stagewise regression and, we always choose smallest index to break the tie.

For forward stagewise regression we use $\tilde{\tau}_t$ instead of τ_t for Frank-wolfe, where $\tau_t = \frac{1}{C} \tilde{\tau}_t \forall t$

$$(\text{Frank-wolfe}) \beta_1 = \tau_1 V_0 = \tau_1 C \operatorname{sign}(X^T y_{(i_0)}) e_{i_0} = \tilde{\beta}_1$$

$$r_1 = y - \tilde{\tau}_1 [\operatorname{sign}(\langle y, x_{j_0} \rangle) x_{j_0}]$$

(stagewise regression)

$$= y - c [z_t \text{sign}(\langle y, x_{j_0} \rangle)] x_{j_0} = y - x_{j_0}$$

Mathematical induction:

($t=0$) verified above

Suppose that $\forall t \in \mathbb{N}, t \leq n$, we have

$$\bar{t} = j_t, \beta_{t+1} = \tilde{\beta}_{t+1} \text{ and } r_{t+1} = y - X \beta_{t+1},$$

Consider $t = n+1$

$$U_{n+1} = \arg \min_v \{ \underbrace{(-x^T (y - X \beta_{n+1}))}_{}, v \mid v \in \mathcal{R} \}$$

$$= \arg \min_v \{ \langle -x^T (r_{n+1}), v \rangle \mid v \in \mathcal{R} \}$$

$$= c \text{sign} (r_{n+1}^T x_{\bar{j}_{n+1}}) e_{\bar{j}_{n+1}}$$

where \bar{j}_{n+1} s.t. $|r_{n+1}^T x_{\bar{j}_{n+1}}| = \|r_{n+1} x\|_\infty$

that is we can have

$$\bar{j}_{n+1} = j_{n+1} \in \arg \max_j \{ |\langle r_{n+1}, x_j \rangle| \mid j \in \{1, \dots, p\} \}$$

$$\beta_{n+2} = (1 - z_{n+1}) \beta_{n+1} + U_{n+1} z_{n+1}$$

$$= \begin{cases} (1 - z_{n+1}) \beta_{n+1}^{(j)} + z_{n+1} c \text{sign}(\langle r_{n+1}, x_{\bar{j}_{n+1}} \rangle) e_{\bar{j}_{n+1}} \\ (1 - z_{n+1}) \beta_{n+1}^{(j)} \end{cases}$$

$$= \begin{cases} (1 - \frac{\tilde{z}_{n+1}}{c}) \tilde{\beta}_{n+1}(j) + \tilde{z} \operatorname{sign}(\langle r_{n+1}, x_{j_{n+1}} \rangle) & \text{if } j = j_{n+1} \\ (1 - \frac{\tilde{z}_{n+1}}{c}) \tilde{\beta}_{n+1}(j) & \text{otherwise} \end{cases}$$

$$= \tilde{\beta}_{n+2}$$

$$r_{n+2} = r_{n+1} - \frac{\tilde{z}_{n+1} [\operatorname{sign}(\langle r_{n+1}, x_{j_{n+1}} \rangle) x_{j_{n+1}} + \frac{1}{c} (V_{n+1} - y)]}{1}$$

$$= r_{n+1} - \tilde{z}_{n+1} \cancel{X} V_{n+1} + \tilde{z}_{n+1} X \beta_{n+1}$$

$$= y - X \beta_{n+1} + \tilde{z}_{n+1} X \beta_{n+1} - \tilde{z}_{n+1} X V_{n+1}$$

$$= y - X (\tilde{z}_{n+1} V_{n+1} + (1 - \tilde{z}_{n+1}) \beta_{n+1})$$

$$= y - X \beta_{n+2}$$

The iteration rule of the modified forward stage regression is equivalent to that of the Frank-Wolfe algorithm applied to the Lasso problem #.

Q2 = $f^* = \min_x \{ f(x) \mid x \in \mathbb{R}^{d \times d}, x \geq 0, \text{Tr}(x) = 1 \}$, $f \stackrel{L}{\sim}$ smooth with Frobenius $\rightarrow \mathcal{X}$

consider using Frank-wolfe algorithm with

by assumption, the curvature

$$C_f \leq L \max_{x,y} \{ \|x-y\|_F^2 \mid x,y \in \mathcal{X} \}$$

$$\begin{aligned} \|x-y\|_F^2 &= \text{trace}((x-y)^T(x-y)) \\ &= \text{trace}(x^T x) - 2\text{trace}(x^T y) + \text{trace}(y^T y) \end{aligned}$$

note that: $\text{trace}(x) = \sum_{i=1}^d \lambda_i$, where x is PSD.
 $\lambda_d \geq 0, \lambda_d$ eigenvalue

$$\text{trace}(x^T x) = \sum_{i=1}^d \lambda_i^2 \leq 1 \quad \text{of } \mathcal{X}.$$

(each $\lambda_i^2 \leq \lambda_i$).

$$\text{trace}(y^T y) \leq 1$$

notice that $\text{tr}(x, y)$ defines inner-product $\langle x, y \rangle$,

by cauchy $\text{tr}(x^T y) \leq \sqrt{\text{tr}(x^T x)} \sqrt{\text{tr}(y^T y)} \leq 1$

It shows that $C_f \leq 4L$

C_f exists.

Set $X_0 = u_{-1}u_{-1}^T$, where $u_{-1} = e_1 \in \mathbb{R}^d$.

$$X_0 \in \mathcal{X}, \quad \tau_t = \frac{2}{t+2}$$

using Frank-wolfe =

for $t=0, 1, \dots, T$ do.

$$V_t = \underset{v}{\operatorname{argmin}} \{ \langle \nabla f(X_t), v \rangle \mid v \in \mathcal{X} \}$$

$$= u_t u_t^T \quad (\text{by proposition 8-19})$$

(

$$X_{t+1} = (1 - \tau_t) X_t + \tau_t V_t$$

end

by convergence guarantee of Frank-wolfe,

$$\text{take } T > \frac{2C_f}{\varepsilon} - 2.$$

then $\nabla f(X_T) - f^* \leq \varepsilon$, X_T is a ε -approximation solution, since $\operatorname{rank}(A+B) \leq \operatorname{rank}(A) + \operatorname{rank}(B)$,

$$\operatorname{rank}(X_T) \leq T, \quad \operatorname{rank}(X_T) = O\left(\frac{1}{\varepsilon}\right) \quad \#$$

Remark =

We don't have ^{any} constraint to
force $\nabla f(x)$ p.s.d $\forall x \in \mathcal{X}$,

so it is not true that $V_t = U_t U_t^T$.

Suppose we have the linear minimization
oracle, let \tilde{V}_t be the _{output} $\tilde{V}_t \in \mathcal{X}$, $\text{trace}(\tilde{V}_t) = 1$.

$$\tilde{V}_t = \sum \lambda_k b_k b_k^T \quad (\text{by eigen-decomposition}) \quad \tilde{V}_t \text{ is now p.s.d.}$$

$$\text{then } \langle \nabla f(x_t), \tilde{V}_t \rangle = \langle \nabla f(x_t), b_k b_k^T \rangle \\ \forall k = 1, 2, \dots, \text{rank}(\tilde{V}_t)$$

we can choose our V_t as $b_k b_k^T$ for
some $k \neq$.