

1. Property of Smooth Functions

$$f \in C^1(\mathbb{R}^d) \quad \|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\|, \quad x, y \in \mathbb{R}^d$$

$$\text{Prove} \quad f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{L}{2} \|y - x\|^2$$

Proof. Let $g(t) = f(x + t(y - x))$, $t \in [0, 1]$

$$g(0) = f(x), \quad g(1) = f(y)$$

$$g'(t) = \langle \nabla f(x + t(y - x)), y - x \rangle \quad \text{by chain rule.}$$

Thus $g \in C[0, 1]$ and $g' \in C(0, 1)$.

We need to prove

$$g(1) - g(0) \leq g'(0) + \frac{L}{2} \|y - x\|^2$$

$$\text{In fact} \quad g(1) - g(0) = \int_0^1 g'(t) dt = \int_0^1 \langle \nabla f(x + t(y - x)), y - x \rangle dt$$

$$g(1) - g(0) - g'(0) = \int_0^1 \langle \nabla f(x_t) - \nabla f(x), y - x \rangle dt, \quad x_t = x + t(y - x)$$

$$\leq \int_0^1 \|\nabla f(x_t) - \nabla f(x)\| \|y - x\| dt \quad (\text{Cauchy inequality})$$

$$\leq L \int_0^1 t \|y - x\|^2 dt = \frac{L \|y - x\|^2}{2}$$

We are done. #

2 Convergence of GF under KL condition

$$\inf_x f(x) = 0, \quad \|\nabla f(x)\|^2 \geq \mu f(x)^\alpha \quad (x_t)_{t \geq 0} \text{ be GF solution}$$

(a) $\alpha > 1$,

$$f(x_t) \leq \frac{1}{(f(x_0)^{1-\alpha} + \mu(\alpha-1)t)^{\frac{1}{\alpha-1}}} \sim t^{-\frac{1}{\alpha-1}}$$

Proof. If there exist $0 \leq s \leq t$ s.t. $x_s = 0$

Then from $\dot{x}_t = -\nabla f(x_t)$, $x_t \equiv 0$ for $t \geq s$ we are done.

Otherwise, $x_s > 0$, $0 \leq s \leq t$

We have $\frac{df(x_t)}{dt} = -\|\nabla f(x_t)\|^2$

$$f(x_t)^{1-\alpha} - f(x_0)^{1-\alpha} = \int_0^t (f(x_r)^{1-\alpha})'_r dr$$

$$= \int_0^t (\alpha-1) f(x_r)^{-\alpha} \|\nabla f(x_r)\|^2 dt$$

$$\geq \int_0^t \mu(\alpha-1) dt = \mu(\alpha-1)t$$

Therefore $f(x_t)^{1-\alpha} \geq f(x_0)^{1-\alpha} + \mu(\alpha-1)t$

$$f(x_t) \leq \frac{1}{(f(x_0)^{1-\alpha} + \mu(\alpha-1)t)^{\frac{1}{\alpha-1}}} \quad (\text{As } \alpha > 1)$$

#

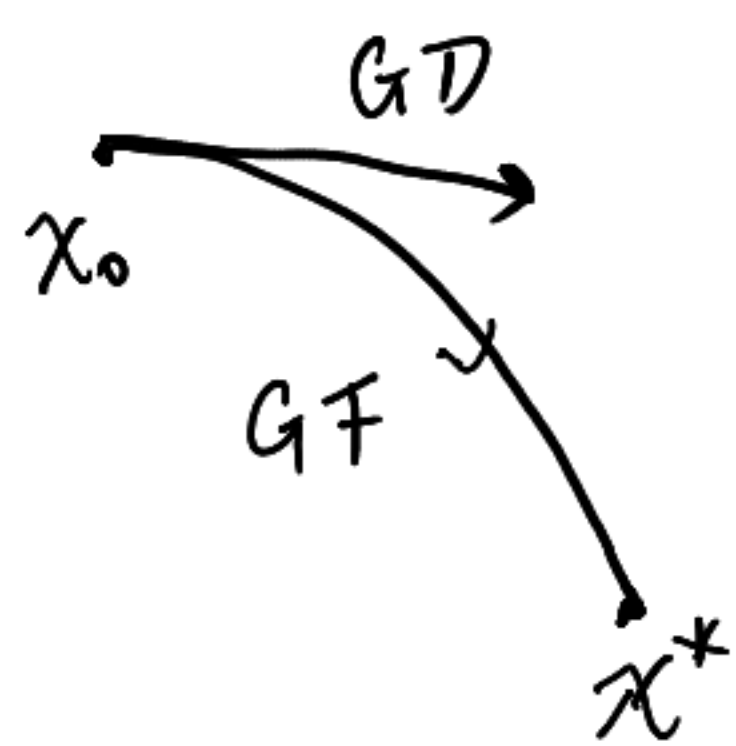
$$(b) \alpha < 1, \quad f(x_t) \leq \left(f(x_0)^{1-\alpha} - \mu(1-\alpha)t \right)^{\frac{1}{1-\alpha}}, \quad \forall t < \frac{f(x_0)^{1-\alpha}}{\mu(1-\alpha)}$$

Proof. If there exists $s \in S \leq t$, $f(x_s) = 0$
 then $\dot{x}_t = -\nabla f(x_t)$ implies $f(x_t) = 0$
 with $t < \frac{f(x_0)^{1-\alpha}}{\mu(1-\alpha)}$ we are done. Otherwise $f(x_s) > 0$, $0 \leq s \leq t$,

we have

$$\begin{aligned} & f(x_t)^{1-\alpha} - f(x_0)^{1-\alpha} \\ &= \int_0^t \left(f(x_r)^{1-\alpha} \right)'_r dt \\ &= - \int_0^t (1-\alpha) f(x_r)^{-\alpha} \|\nabla f(x_r)\|^2 dt \\ &\leq - \int_0^t (1-\alpha) \mu dt = -\mu(1-\alpha)t \\ & f(x_t)^{1-\alpha} \leq f(x_0)^{1-\alpha} - \mu(1-\alpha)t \\ & f(x_t) \leq \left(f(x_0)^{1-\alpha} - \mu(1-\alpha)t \right)^{\frac{1}{1-\alpha}} \quad \# \end{aligned}$$

(c) Because finite-time convergence of GF means after a finite time T , going on the trajectory of GF will reach minima



But finite-step GD goes on the tangent line of trajectory of GF, will probably reach different point after T .

From another perspective, suppose learning rate is $\{\eta_t\}_0^{T_1}$

To let GF approximate GD, $\eta_t < 1$

but $T \approx \sum_{t=0}^{T_1} \eta_t$ This implies $T_1 \rightarrow \infty$

3 Implicit bias of GD for linear regression

$$\hat{R}(\beta) = \frac{1}{2} \|X\beta - y\|_2^2$$

$$X \in \mathbb{R}^{n \times d} \quad \beta \in \mathbb{R}^d \quad y \in \mathbb{R}^n \quad d > n, \quad \min_{\beta} \hat{R}(\beta) = 0$$

Optimize with GD

$$\beta_0 = 0$$

$$\beta_{t+1} = \beta_t - \eta \nabla \hat{R}(\beta_t)$$

Prove $\lim_{t \rightarrow \infty} \beta_t = \bar{\beta}$, $\bar{\beta} := \arg\min_{\beta} \|\beta\|_2 \text{ s.t. } X\beta = y$.

Proof.
$$\begin{aligned} \hat{R}(\beta) &= \frac{1}{2} (\beta^T X^T - y^T)(X\beta - y) \\ &= \frac{1}{2} \beta^T X^T X \beta - (X^T y)^T \beta + \frac{1}{2} y^T y \end{aligned}$$

$$\nabla_{\beta} \hat{R}(\beta) = X^T X \beta - X^T y$$

Dynamics is
$$\beta_{t+1} = \beta_t - \eta (X^T X \beta_t - X^T y)$$

$$\beta_{t+1} = (I_d - \eta X^T X) \beta_t + \eta X^T y$$

$$\beta_{t+1} - \bar{\beta} = (I_d - \eta X^T X) (\beta_t - \bar{\beta})$$

$$\beta_t - \bar{\beta} = -(I_d - \eta X^T X)^t \bar{\beta}$$

We analyze spectrum of $I_d - \eta X^T X$.

Consider SVD decomposition of X

$$X = T \Sigma S, T \in O(n), S \in O(d), \Sigma = \begin{pmatrix} \sigma_1 & \dots & \sigma_r & 0 & \dots & 0 \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \end{pmatrix}, \sigma_i > 0, r \leq n$$

$$X^T X = S^T \Sigma^T T^T T \Sigma S = S^T \Sigma^T \Sigma S$$

$$= S^T \begin{pmatrix} \sigma_1^2 & \dots & \sigma_r^2 & 0 & \dots & 0 \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \end{pmatrix} S$$

$$I_d - \gamma X^T X = S^T I_d S - \gamma S^T \begin{pmatrix} \sigma_1^2 & \dots & \sigma_r^2 & 0 & \dots & 0 \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \end{pmatrix} S$$

$$= S^T \begin{pmatrix} 1 - \gamma \sigma_1^2 & \dots & 1 - \gamma \sigma_r^2 & 0 & \dots & 0 \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \end{pmatrix} S$$

$$\beta_t - \bar{\beta} = -S^T \begin{pmatrix} (1 - \gamma \sigma_1^2)^t & \dots & (1 - \gamma \sigma_r^2)^t & 0 & \dots & 0 \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \end{pmatrix} S \bar{\beta}$$

$$\text{For } 0 < \gamma < \min_i \frac{1}{\sigma_i^2},$$

$$\lim_{t \rightarrow \infty} \beta_t - \bar{\beta} = -S^T \begin{pmatrix} 0 & \dots & 0 & 0 & \dots & 0 \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \end{pmatrix} S \bar{\beta}$$

It is well known $\bar{\beta}$ admits the form $\bar{\beta} = S^T \begin{pmatrix} \frac{1}{\sigma_1} & \dots & \frac{1}{\sigma_r} & 0 & \dots & 0 \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \end{pmatrix} T^T y$ from linear algebra

Thus

$$-S^T \begin{pmatrix} 0 & \dots & 0 & 0 & \dots & 0 \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \end{pmatrix} S \bar{\beta} = -S^T \begin{pmatrix} 0 & \dots & 0 & 0 & \dots & 0 \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \end{pmatrix} \begin{pmatrix} \frac{1}{\sigma_1} & \dots & \frac{1}{\sigma_r} & 0 & \dots & 0 \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \end{pmatrix} T^T y$$

$$= -S^T 0 T^T y = 0$$

We are done #

4 The dynamic behavior of HB method

(a) $f(x) = \frac{1}{2}hx^2 \Rightarrow f'(x) = hx$

Dynamics is $\dot{x}_t = -\gamma \dot{x}_t - hx_t$ This is second-order constant coefficient linear ODE.

characteristic equation

$$\lambda^2 + \gamma\lambda + h = 0$$

Case 1. $\gamma^2 > 4h$
$$x_t = A_1 e^{\frac{-\gamma + \sqrt{\gamma^2 - 4h}}{2} t} + A_2 e^{\frac{-\gamma - \sqrt{\gamma^2 - 4h}}{2} t}$$

for $\gamma > 0, h > 0$, $x_t \rightarrow 0$ with rate $O(e^{\frac{-\gamma + \sqrt{\gamma^2 - 4h}}{2} t})$
Other (γ, h) leads to diverges ($x_t \rightarrow \infty$)

Case 2. $\gamma^2 = 4h$
$$x_t = (A_1 t + A_2) e^{-\frac{\gamma}{2} t}$$

for $\gamma > 0$, $x_t \rightarrow 0$ with rate $O(t e^{-\frac{\gamma}{2} t})$, not monotonic
for $\gamma \leq 0$ diverges ($x_t \rightarrow \infty$)

Case 3. $\gamma^2 < 4h$
$$x_t = A_1 e^{-\frac{\gamma}{2} t} \cos\left(\frac{\sqrt{4h - \gamma^2}}{2} t\right) + A_2 e^{-\frac{\gamma}{2} t} \sin\left(\frac{\sqrt{4h - \gamma^2}}{2} t\right)$$

for $\gamma > 0$, $x_t \rightarrow 0$ with rate $O(e^{-\frac{\gamma}{2} t})$, not monotonic
for $\gamma \leq 0$ diverges

where above, A_1, A_2 are constants depending on initializations.

(b) See last page of this PDF

5 Convergence of HB method for convex problems

$$\ddot{\chi}_t = -\dot{\chi}_t - \nabla f(\chi_t), \quad f \text{ convex}, \quad \chi^* = \min_{\chi} f(\chi)$$

$$V(t) = f(\chi_t) - f(\chi^*) + \frac{1}{2} \|\chi_t - \chi^* + \dot{\chi}_t\|^2$$

$$(a) \quad \dot{V}(t) \leq f(\chi^*) - f(\chi_t)$$

Proof.

$$\begin{aligned} \dot{V}(t) &= \langle \nabla f(\chi_t), \dot{\chi}_t \rangle + \langle \chi_t - \chi^* + \dot{\chi}_t, \dot{\chi}_t + \ddot{\chi}_t \rangle \\ &= \langle \nabla f(\chi_t), \dot{\chi}_t \rangle + \langle -\nabla f(\chi_t), \chi_t - \chi^* + \dot{\chi}_t \rangle \\ &= \langle \nabla f(\chi_t), \chi^* - \chi_t \rangle \\ &\leq f(\chi^*) - f(\chi_t) \quad \text{by convexity of } f \quad \# \end{aligned}$$

$$(b) \quad \int_0^T (f(\chi_t) - f(\chi^*)) dt \leq V(\omega)$$

Proof. Because of (a),

$$\begin{aligned} \int_0^T (f(\chi_t) - f(\chi^*)) dt &\leq - \int_0^T \dot{V}(t) dt = V(\omega) - V(T) \\ &\leq V(\omega) \quad \left(\text{as } V(T) = f(\chi_T) - \min_{\chi} f(\chi) + \frac{1}{2} \|\chi_T - \chi^* + \dot{\chi}_T\|^2 \geq 0 \right) \end{aligned}$$

$$(c) \quad \bar{\chi}_T := \frac{1}{T} \int_0^T \chi_t dt \quad \text{Show } f(\bar{\chi}_T) - f(\chi^*) \leq \frac{V(\omega)}{T}$$

$$\text{Proof} \quad f(\bar{\chi}_T) = f\left(\frac{1}{T} \int_0^T \chi_t dt\right) \leq \frac{1}{T} \int_0^T f(\chi_t) dt \quad (*)$$

$$\leq \frac{1}{T} \left(\int_0^T f(\chi^*) dt + V(\omega) \right) \quad (\text{because of (b)})$$

$$\text{i.e. } f(\bar{\chi}_T) - f(\chi^*) \leq \frac{V(\omega)}{T} \quad \#$$

(*) : Random vector X realized by uniformly pick $t \in [0, T]$, and take χ_t
 $\mathbb{E} X = \frac{1}{T} \int_0^T \chi_t dt$ by Jensen inequality $f(\mathbb{E} X) \leq \mathbb{E} f(\chi_t)$.