

Chapter 6 Subgradient Methods

Subgradient and Subdifferential

- **Definition (Subgradient)** Let $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ be convex. A vector $g \in \mathbb{R}^n$ is a *subgradient* of f at a point $x \in \text{dom}(f)$ if $f(y) \geq f(x) + g^\top(y - x), \forall y \in \text{dom}(f)$.
 - The set of all subgradient at x is called the *subdifferential* of f at x denoted as ∂f .
- **Lemma 6.2 (uniqueness when differentiable)** If f is convex and differentiable at $x \in \text{dom}(f)$, then $\partial f = \{\nabla f(x)\}$.
 - Proof
 - Let $y = x + \epsilon d$, if $g \in \partial f$, then $f(x + \epsilon d) \geq f(x) + \epsilon g^\top d$, then $\frac{f(x + \epsilon d) - f(x)}{\epsilon} \geq g^\top d, \forall d, \forall \epsilon$.
 - Letting $\epsilon \rightarrow 0$, then we have $\nabla f(x)^\top d \geq g^\top d, \forall d$. This only holds when $g = \nabla f(x)$.
 - If f convex, we can check the definition of subgradient holds.
 - Or can use Taylor expansion and get $(g - \nabla f(x))^\top(y - x) \leq r(y - x)$, then let $y = x + \epsilon(g - \nabla f(x))$.
- **Lemma 6.3** If we don't assume convexity, we can only say $\partial f(x) \subseteq \{\nabla f(x)\}$.
- **Examples 6.4**
 - $f(x) = \frac{1}{2}x^2, \partial f(x) = x$
 - $f(x) = |x|, \partial f(x) = \begin{cases} \text{sgn}(x), x \neq 0 \\ [-1, 1], x = 0 \end{cases}$
 - $f(x) = \begin{cases} -\sqrt{x}, x \geq 0 \\ +\infty, o.w. \end{cases}, \partial f(0) = \emptyset;$
 - $f(x) = \begin{cases} 1, x = 0 \\ 0, x > 0 \\ +\infty, o.w \end{cases}, \partial f(0) = \emptyset$
- **Lemma 6.5** Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be convex, $\text{dom}(f)$ open, $B \in \mathbb{R}_+$. Then $\forall x \in \text{dom}(f), \forall g \in \partial f(x), \|g\|_{a*} \leq B \Leftrightarrow \forall x, y \in \text{dom}(f), |f(y) - f(x)| \leq B\|x - y\|_a$.
 - Proof
 - (\Rightarrow)
 - Subgradient $\rightarrow f(y) - f(x) \leq g_y^\top(y - x)$ and $f(x) - f(y) \leq g_x^\top(y - x)$
 - This means $|f(y) - f(x)| \leq \max\{g_y^\top(y - x), g_x^\top(y - x)\} \leq \max\{\|g_y\|_{a*}, \|g_x\|_{a*}\}\|y - x\|_a$
 - If $\|g\|_{a*} \leq B$ we arrive at Lipschitz
 - (\Leftarrow)
 - $f(y) \geq f(x) + g_x^\top(y - x)$, by Lipschitz, we have $B\|x - y\|_a \geq g_x^\top(y - x), \forall y$. Choosing $y - x = \arg \max_{\|y-x\|=1} g_x^\top(y - x)$, we get maximum of $\|g_x\|_{a*} \leq B$. This is true for all x .

Topological Properties of Subgradients

- **Lemma 6.6** Let $f(x)$ be a convex function and $x \in \text{dom}(f)$. Then $\partial f(x)$ is convex and closed.
 - Proof $\partial f(x) = \cap_y \{g \in \mathbb{R}^n : f(y) \geq f(x) + g^\top(y - x)\}$, solution of g for each y is convex and closed linear separation.
- **Definition 6.7 (separation of convex sets)** Let S and T be two nonempty convex sets in \mathbb{R}^n . A hyperplane $H = \{x \in \mathbb{R}^n : a^\top x = b, a \neq 0\}$ is said to separate S and T if $S \cup T \not\subset H$ and $S \subset H^- = \{x \in \mathbb{R}^n : a^\top x \leq b\}$ and $T \subset H^+ = \{x \in \mathbb{R}^n : a^\top x \geq b\}$.
 - PS: It's OK that one of S or T is a point, or part of the hyperplane.
 - *Strictly separation* if $S \subset H^{--} = \{x \in \mathbb{R}^n : a^\top x < b\}$ and $T \subset H^{++} = \{x \in \mathbb{R}^n : a^\top x > b\}$.
- **Theorem 6.8 (Hyperplane separation theorem, [Roc97, Thm 11.3], no proof)** Let S and T be two nonempty convex sets. Then S and T can be separated if and only if their (relative) interiors do not intersect, $\text{rint}(S) \cap \text{rint}(T) = \emptyset$.
 - PS The *Relative Interior* of a general non=convex set is defined to be

$$\text{relint}(C) := \{x \in C : \forall y \in C, \exists \lambda > 1, \text{ s.t. } \lambda x + (1 - \lambda)y \in C\}$$
- **Corollary 6.9** Let S be a nonempty convex set and $x_0 \in \partial S$. $\exists H = \{x : a^\top x = a^\top x_0\}, a \neq 0$ s.t. $S \subset \{x : a^\top x \leq a^\top x_0\}$, and $x_0 \in H$.
 - just let $T = \{x_0\}$.
- **Theorem 6.10 (non emptiness of subgradient for interior)** Let $f(x)$ be a convex function and $x \in \text{rint}(\text{dom}(f))$. Then $\partial f(x)$ is nonempty and bounded if $x \in \text{rint}(\text{dom}(f))$.
 - Proof
 - Non-emptiness

- W.l.o.g., assume $\text{dom}(f)$ is full-dimensional and $x \in \text{int}(\text{dom}(f))$ (*interior*).
 - $\text{epi}(f)$ is convex and $(x, f(x))$ belongs to its boundary, by Thm 6.8 $\exists \alpha = (s, \beta) \neq 0$ s.t.
 $s^\top y + \beta t \geq s^\top x + \beta f(x), \forall (y, t) \in \text{epi}(f)$
 - We claim $\beta > 0$.
 - If $\beta < 0$, we can let $t \rightarrow \infty$ where $(y, t) \in \text{epi}(f)$.
 - if $\beta = 0$, this contradict to $s^\top y \geq s^\top x$, since x is the interior, we can always find opposite direction to reverse the inequality.
 - Setting $g = -\beta^{-1}s$, we have $f(y) \geq f(x) + g^\top(y - x), \forall y$
 - *Boundedness*
 - If $\partial f(x) \not\prec \infty$, $\exists \{g_k\}_{k=1}^\infty \in \partial f(x)$ s.t. $\lim_{k \rightarrow \infty} \|g_k\|_2 = \infty$.
 - Since x interior, exists ball $B(x, \delta) \subseteq \text{dom}(f)$, then $y_k = x + \delta \frac{g_k}{\|g_k\|_2} \in \text{dom}(f)$, by convexity,
 $f(y_k) \geq f(x) + g_k^\top(y_k - x) = f(x) + \delta \|g_k\|_2 \rightarrow \infty$, contradiction.
- Lemma A (non-emptiness \rightarrow convexity) If $\forall x \in \text{dom}(f), \partial f(x) \neq \emptyset$ and $\text{dom}(f)$ convex, then f convex.
- Proof
 - $z = \lambda x + (1 - \lambda)y \in \text{dom}(f)$ for $\lambda \in [0, 1]$.
 - Let $g \in \partial f(z)$, we have $f(x) \geq f(z) + g^\top(x - z), f(y) \geq f(z) + g^\top(y - z)$, add them with weight $(\lambda, 1 - \lambda)$ we get convexity.
- Remark 6.11 (Exercise 42) The subdifferential of a convex function $f(x)$ at $x \in \text{dom}(f)$ is a monotone operator, i.e., $(u - v)^\top(x - y) \geq 0, \forall x, y \in \text{dom}(f), u \in \partial f(x), v \in \partial f(y)$. *Proof similar to gradient.*

Subdifferential and directional derivative

- Def directional derivative $f'(x; d) = \lim_{\delta \rightarrow 0^+} \frac{f(x + \delta d) - f(x)}{\delta}$.
- Def differentiable f *differentiable* if $\Delta f = A\Delta x + o(\Delta x)$.
 - ∂f can exists when $d f$ not.
- Lemma 6.12 (Exercise 43) Let f convex, the ratio $\phi(\delta) := \frac{f(x + \delta d) - f(x)}{\delta}$ is non-decreasing of $\delta > 0$.
 - Proof Let $\delta_1 \geq \delta_2$, then let $\lambda := \delta_2/\delta_1$, by convexity $\phi(\delta_2) = \frac{f(x + \lambda\delta_1 d) - f(x)}{\lambda\delta_1} \leq \frac{\lambda f(x + \delta_1 d) + (1 - \lambda)f(x) - f(x)}{\lambda\delta_1} = \phi(\delta_1)$.
- Theorem 6.13 Let f be convex and $x \in \text{int}(\text{dom}(f))$, then $f'(x; d) = \max_{g \in \partial f(x)} g^\top d$.
 - Proof
 - By def of subgrad, we have $f'(x; d) \geq \max_{g \in \partial f(x)} g^\top d$.
 - Now we want to also prove the opposite.
 - Define $C_1 = \{(y, t) : f(y) < t\}$, (not **epi**(f)), and $C_2(d) = \{(y, t) : y = x + \alpha d, t = f(x) + \alpha f'(x; d), \alpha \geq 0\}$.
 - so both C_1, C_2 convex, non-empty.
 - Since $\phi(\delta)$ non-decreasing, $(t >) f(x + \alpha d) \geq f(x) + \alpha f'(x; d), \forall \alpha \geq 0$, so $C_1 \cap C_2 = \emptyset$
 - By hyperplane separation theorem $\exists (g_0, \beta) \neq 0$ s.t.
 $g_0^\top(x + \alpha d) + \beta(f(x) + \alpha f'(x; d)) \leq g_0^\top y + \beta t, \forall \alpha \geq 0, \forall t > f(y)$
 - Similar to proof of Thm 6.10, we claim $\beta > 0$. Let $\tilde{g} = \beta^{-1}g_0$
 - we get $\tilde{g}^\top(x + \alpha d) + f(x) + \alpha f'(x; d) \leq \tilde{g}^\top y + f(y), \forall \alpha \geq 0$.
 - Setting $\alpha = 0$, we get $\tilde{g}^\top x + f(x) \leq \tilde{g}^\top y + f(y) \Leftrightarrow -\tilde{g} \in \partial f(x)$
 - setting $\alpha = 1, y = x$, we get $\tilde{g}^\top d + f'(x; d) \leq 0 \Leftrightarrow f'(x; d) \leq -\tilde{g}^\top d \leq \max_{g \in \partial f(x)} g^\top d$
 - This mean $f'(x; d) = \max_{g \in \partial f(x)} g^\top d$.
- Theorem B [Roc97, Theorem 25.5] (differentiable almost everywhere) A (proper) convex function f is differentiable almost everywhere on (the interior of) $\text{dom}(f)$.
- Lemma C (optimal condition) If $0 \in \partial f(x)$, then x is a global minimum. *Proof* $f(y) \geq f(x) + g^\top(y - x) = f(x)$.

Calculus of Subgradient

- Conic combination $h(x) = \lambda f(x) + \mu g(x)$, then $\partial h(x) = \lambda \partial f(x) + \mu \partial g(x), \forall x \in \text{int}(\text{dom}(h))$.
- affine composition $h(x) = f(Ax + b)$, then $\partial h(x) = A^T \partial f(Ax + b)$.
- supremum $h(x) = \sup_{\alpha \in \mathcal{A}} f_\alpha(x)$, then $\partial h(x) \supseteq \text{conv} \{\partial f_\alpha(x) | \alpha \in \mathcal{A}\}$
- superposition $h(x) = F(f_1(x), \dots, f_m(x))$ where $F(y_1, \dots, y_m)$ non-decreasing and convex, then
 $\partial h(x) \supseteq \left\{ \sum_{i=1}^m d_i \partial f_i(x) : (d_1, \dots, d_m) \in \partial F(y_1, \dots, y_m) \right\}$.
- Example 6.14 Let $h(x) = \max_{y \in C} f(x, y)$ where $f(x, y)$ is convex in x for any y and C is closed. Then $\partial f(x, y_*(x)) \subset \partial h(x)$, where $y_*(x) = \arg\max_{y \in C} f(x, y)$
 - Proof Let $g \in \partial f(x, y_*(x))$, then $h(z) \geq f(z, y_*(x)) \geq f(x, y_*(x)) + g^\top(z - x) = h(z) + g^\top(z - x)$

- Lemma 6.15 (Exercise 44) Consider the function $f(x) = \|x\|_a$. Then $\partial f(x) = \{g : g^T x = \|x\|_a \wedge \|g\|_{a^*} \leq 1\}$. (In particular $\partial f(0) = \{g : \|g\|_{a^*} \leq 1\}$.
 - Proof
 - $g \in \partial f(x) \Leftrightarrow \forall \delta, \|x + \delta\|_a \geq \|x\|_a + g^\top \delta$
 - setting $\delta = -x$, we get $0 = \|0\|_a = \|x\|_a - g^\top x$.
 - By triangular ineq, $\|x\|_a + \|\delta\|_a \geq \|x + \delta\|_a \geq \|x\|_a + g^\top \delta \Rightarrow \|\delta\|_a \geq g^\top \delta, \forall \delta$.
 - subjecting to $\|\delta\|_a = 1$ we get $\|g\|_{a^*} \leq 1$.

Subgradient Methods

- Remark 6.16 Negative subgradient may not be a descending direction.
 - Example $f(x_1, x_2) = |x_1| + 2|x_2|$ at $\mathbf{x} = (1, 0)$, $\partial f(\mathbf{x}) = \{(1, a) : a \in [-2, 2]\}$
 - $\mathbf{g} = (1, 0), \mathbf{d} = -\mathbf{g}$ is descending direction, while $\mathbf{g} = (1, 2), \mathbf{d} = -\mathbf{g}$ is not.
- Problem $\min f(x)$ s.t. $x \in X$. Denote:
 - $R := \sqrt{\max_{x,y \in X} \|x - y\|_2^2}$ the diameter of X .
 - $B := \sup_{x,y \in X} \frac{|f(x) - f(y)|}{\|x - y\|_2} < +\infty$ the Lipschitz constant under ℓ_2 norm.

Subgradient Descent

- Algorithm $\mathbf{x}_{t+1} = \Pi_X(\mathbf{x}_t - \gamma_t \mathbf{g}(\mathbf{x}_t))$, where $\mathbf{g}(\mathbf{x}_t) \in \partial f(\mathbf{x}_t)$.
- Lemma D(Descent Lemma) f convex, for any optimal $\mathbf{x}^* \in X^*$, $\|\mathbf{x}_{t+1} - \mathbf{x}^*\|_2^2 \leq \|\mathbf{x}_t - \mathbf{x}^*\|_2^2 - 2\gamma_t(f(\mathbf{x}_t) - f^*) + \gamma_t^2 \|\mathbf{g}_t\|_2^2$
 - Proof
 - $\|\mathbf{x}_{t+1} - \mathbf{x}^*\|_2^2 = \|\Pi_X(\mathbf{x}_t - \gamma_t \mathbf{g}_t) - \mathbf{x}^*\|_2^2$,
 - by Fact 4.1 (ii) RHS $\leq \|\mathbf{x}_t - \gamma_t \mathbf{g}_t - \mathbf{x}^*\|_2^2 = \|\mathbf{x}_t - \mathbf{x}^*\|_2^2 - 2\gamma_t \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle + \gamma_t^2 \|\mathbf{g}_t\|_2^2$
 - By convexity $\langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}^* \rangle \leq f(\mathbf{x}_t) - f^*$, then we arrive at conclusion.

$\mathcal{O}(1/\sqrt{t})$ convergence for B -Lipschitz convex functions

- Theorem 6.17 f convex, then subgradient descent gives
$$\max\{\min_{1 \leq t \leq T} f(\mathbf{x}_t), f(\hat{\mathbf{x}}_T)\} - f^* \leq \frac{\|\mathbf{x}_1 - \mathbf{x}^*\|_2^2 + \sum_{t=1}^T \gamma_t^2 \|\mathbf{g}_t\|_2^2}{2 \sum_{t=1}^T \gamma_t} \leq \frac{R^2 + \sum_{t=1}^T \gamma_t^2 B^2}{2 \sum_{t=1}^T \gamma_t}.$$
 - Proof
 - By descent lemma $\gamma_t(f(\mathbf{x}_t) - f^*) \leq \frac{1}{2} (\|\mathbf{x}_t - \mathbf{x}^*\|_2^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|_2^2 + \gamma_t^2 \|\mathbf{g}(\mathbf{x}_t)\|_2^2)$
 - sum over $t = [1 : T]$ and divided by $\sum_t \gamma_t$ gives upper bound of the average $\mathbb{E} f_t$, after that is straightforward.
- Step size is crucial for convergence behavior, unlike GD

- Constant Stepsize $\gamma_t \equiv \gamma, \epsilon_T \leq \frac{R^2}{2T} \cdot \frac{1}{\gamma} + \frac{B^2}{2} \gamma \xrightarrow{T \rightarrow \infty} \frac{B^2}{2} \gamma$,
 - choosing $\gamma_* = \frac{R}{B\sqrt{T}}$ yields $\epsilon_T \leq \frac{RB}{\sqrt{T}}$.
- Scaled stepsize $\gamma_t = \frac{\gamma}{\|\mathbf{g}_t\|_2}, \epsilon_T \xrightarrow{T \rightarrow \infty} B\gamma/2$
- Non-summable but diminishing stepsize $\sum_{t=1}^{\infty} \gamma_t = \infty, \lim_{t \rightarrow \infty} \gamma_t = 0$
 - Since γ_t^2 approach zero faster than $\gamma_t, \epsilon_T \xrightarrow{T \rightarrow \infty} 0$.
 - Example $\gamma_t = O(t^{-q})$ with $q \in (0, 1]$, e.g. $q = 1/2$, then $\epsilon_T \leq O\left(\frac{BR \ln(T)}{\sqrt{T}}\right)$
 - If we average among $[T/2 : T]$, then $\min_{\lfloor \frac{T}{2} \rfloor \leq t \leq T} f(\mathbf{x}_t) - f^* \leq O\left(\frac{BR}{\sqrt{T}}\right)$
- Non-summable but square-summable (Robbins-Monro) stepsize $\sum_{t=1}^{\infty} \gamma_t = \infty, \sum_{t=1}^{\infty} \gamma_t^2 < \infty$, e.g. $\gamma_t = t^{-1}, \epsilon_T \xrightarrow{T \rightarrow \infty} 0$.
- Polyak stepsize Assume $f^* = f(\mathbf{x}^*)$ is known and $\gamma_t = \frac{f(\mathbf{x}_t) - f^*}{\|\mathbf{g}(\mathbf{x}_t)\|_2^2}$
 - Using descent lemma, we get $\|\mathbf{x}_{t+1} - \mathbf{x}^*\|_2^2 \leq \|\mathbf{x}_t - \mathbf{x}^*\|_2^2 - \frac{(f(\mathbf{x}_t) - f^*)^2}{\|\mathbf{g}(\mathbf{x}_t)\|_2^2}$
 - then $(f(\mathbf{x}_t) - f(\mathbf{x}^*))^2 \leq \|\mathbf{g}_t\|^2 (\|\mathbf{x}_t - \mathbf{x}^*\|_2^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|_2^2) \leq B^2 (\|\mathbf{x}_t - \mathbf{x}^*\|_2^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|_2^2)$
 - then $\sum_{t=1}^T (f(\mathbf{x}_t) - f^*)^2 \leq R^2 \cdot M < \infty$, this means $f(\mathbf{x}_t) \rightarrow f^*$.
 - Note this is last iterate convergence, not average convergence!
 - can also show that $\exists \mathbf{x}^* \in X^*$ s.t. $\limsup_{t \rightarrow \infty} \|\mathbf{x}_t - \mathbf{x}^*\|_2 \rightarrow 0$, converge to one of the optimal point.

- Since series $\|\mathbf{x}_t - \mathbf{x}^*\|_2$ is bounded and non-increasing, there exists a subsequence $\{\mathbf{x}_{t_k}\}$ with accumulation point $\hat{\mathbf{x}}$.
- Since we already have $f(\mathbf{x}_t) \rightarrow f^*$, then $\hat{\mathbf{x}} \in X^*$.
- So $\|\mathbf{x}_{t+1} - \hat{\mathbf{x}}\|_2^2 \leq \|\mathbf{x}_{t_k} - \hat{\mathbf{x}}\|_2^2 - \sum_{j=t_k}^t \frac{(f(\mathbf{x}_j) - f^*)^2}{\|g(\mathbf{x}_j)\|_2^2}$,
- taking \limsup we have $\limsup_{t \rightarrow \infty} \|\mathbf{x}_t - \hat{\mathbf{x}}\|_2^2 \leq \|\mathbf{x}_{t_k} - \hat{\mathbf{x}}\|_2^2 - \sum_{j=t_k}^\infty \frac{(f(\mathbf{x}_j) - f^*)^2}{\|g(\mathbf{x}_j)\|_2^2}$
- then take $k \rightarrow \infty$, we get $\limsup_{t \rightarrow \infty} \|\mathbf{x}_t - \hat{\mathbf{x}}\|_2^2 \leq \lim_{k \rightarrow \infty} \|\mathbf{x}_{t_k} - \hat{\mathbf{x}}\|_2^2 = 0$.

Similar result under Polyak's step size (Exercise Prob 5)

- For non-differentiable convex f , μ -strongly convex w.r.t. 2-norm is defined to be $f_\mu(\mathbf{x}) := f(\mathbf{x}) - \frac{\mu}{2} \|\mathbf{x}\|^2$ is convex. It can also be shown that this means $f(y) \geq f(x) + g^\top(y - x) + \frac{\mu}{2} \|x - y\|^2$.
- Lemma I If f convex and B -Lipschitz, optimization in whole space $\mathbf{x} \in \mathbb{R}^d$ with Polyak step size gives $\min_{1 \leq t \leq T} f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \frac{B\|\mathbf{x}_1 - \mathbf{x}^*\|_2}{\sqrt{T}}$.

◦ Proof

- Denote $h_t := f(\mathbf{x}_t) - f(\mathbf{x}^*)$ and $d_t := \|\mathbf{x}_t - \mathbf{x}^*\|_2$, update gives $d_{t+1}^2 \leq d_t^2 - \frac{h_t^2}{\|g_t\|_2^2}$
- equivalently $h_t^2 \leq \|g_t\|_2^2 (d_t^2 - d_{t+1}^2) \leq B^2 (d_t^2 - d_{t+1}^2)$
- sum over $t = [1 : T]$ gives the result.

- Lemma J If f μ -strongly convex and B -Lipschitz, optimization in whole space $\mathbf{x} \in \mathbb{R}^d$ with Polyak step size gives $\min_{1 \leq t \leq T} f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \frac{4B^2}{\mu T}$

◦ Proof

- Similarly, $d_{t+1}^2 - d_t^2 \leq -\frac{h_t^2}{\|g_t\|_2^2} \leq -\frac{h_t^2}{B^2}$
- By strong convexity w.r.t. optimial point (since optimization on whole space, $0 \in \partial f(\mathbf{x}^*)$, $h_t = f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \mu d_t^2 / 2$)
 - this means $d_{t+1}^2 - d_t^2 \leq -\frac{h_t^2}{B^2} \leq -\frac{\mu^2}{4B^2} d_t^4$. Denote $a_t = \mu^2 \|\mathbf{x}_t - \mathbf{x}^*\|_2^2 / (4B^2)$, we have $a_{t+1} \leq a_t(1 - a_t)$
- By strong convexity w.r.t. \mathbf{x}_0 and \mathbf{x}^* , we have $f^* \geq f(\mathbf{x}_0) + g_0^\top(\mathbf{x}^* - \mathbf{x}_0) + \mu \|\Delta \mathbf{x}_0\|_2^2 / 2$
 - by sufficient descent of Polyak step size, $0 \geq -h_0 \geq \mu d_0^2 / 2 - g_0^\top \Delta \mathbf{x}_0$,
 - this means $\mu d_0^2 / 2 \geq g_0^\top \Delta \mathbf{x}_0 \leq B d_0$, take square we get $a_0 \leq 1$.
- We claim $a_t \leq 1/(t+1)$, by induction, we have $a_{t+1} \leq \frac{1}{t+2} \frac{t(t+2)}{(t+1)^2}$ (by monotonicity of $x(1-x)$ in $[0, 1/2]$)
 - Since $b_t := t/t+1$ is increasing, $b_t/b_{t+1} = \frac{t(t+2)}{(t+1)^2} \leq 1$, this means $a_{t+1} \leq \frac{1}{t+2}$, induction succeeds.
 - This means $d_t^2 \leq \frac{4B^2}{\mu^2} \frac{1}{t+1} \leq \frac{4B^2}{\mu^2} \frac{1}{t}$
- Again, by $d_{t+1}^2 - d_t^2 \leq -\frac{h_t^2}{B^2}$, sum over $t \in [T/2, T]$, we get $\sum_{t=T/2}^T h_t^2 \leq B^2 (d_{T/2}^2 - d_T^2) \leq B^2 d_{T/2}^2 \leq \frac{8B^4}{\mu^2} \frac{1}{T}$
 - this means $\min_{t \in [T/2:T]} h_t^2 \leq \frac{16B^4}{\mu^2} \frac{1}{T^2}$, taking square root and we finish our proof.

$\mathcal{O}(1/t)$ convergence for strong convex functions

- Lemma E(Descent Lemma under strong convexity) f is μ -strongly convex, then $\|\mathbf{x}_{t+1} - \mathbf{x}^*\|_2^2 \leq (1 - \mu\gamma_t) \|\mathbf{x}_t - \mathbf{x}^*\|_2^2 - 2\gamma_t (f(\mathbf{x}_t) - f^*) + \gamma_t^2 \|g_t\|_2^2$
 - Proof Replace $\langle g_t, \mathbf{x}_t - \mathbf{x}^* \rangle \leq f(\mathbf{x}_t) - f^*$ by $\langle g_t, \mathbf{x}_t - \mathbf{x}^* \rangle \leq f(\mathbf{x}_t) - f^* - (\mu/2) \|\mathbf{x}_t - \mathbf{x}^*\|_2^2$.
- There are two choice of step size $1/\mu t$ or $2/\mu(t+1)$, which give different convergence rate.
- Theorem 6.18 ($\mathcal{O}(\ln T/T)$) f is μ -strongly convex and B -Lipschitz, then subgradient descent of step size $\gamma_t = 1/\mu t$ gives $\max\{\min_{1 \leq t \leq T} f(\mathbf{x}_t), f(\hat{\mathbf{x}}_T)\} - f^* \leq \frac{B^2(\ln(T)+1)}{2\mu T}$ where $\hat{\mathbf{x}}_T := \frac{1}{T} \sum_{t=1}^T \mathbf{x}_t$.
 - Proof
 - Plug in $\gamma = 1/\mu t$ we have $f(\mathbf{x}_t) - f^* \leq \left(\frac{\mu}{2}(t-1) \|\mathbf{x}_t - \mathbf{x}^*\|_2^2 - \frac{\mu}{2}t \|\mathbf{x}_{t+1} - \mathbf{x}^*\|_2^2 + \frac{1}{2\mu t} \|g(\mathbf{x}_t)\|_2^2 \right)$,
 - sum over $t = [1 : T]$ we get $\sum_{t=1}^T [f(\mathbf{x}_t) - f^*] \leq \sum_{t=1}^T \frac{1}{2\mu t} \|g(\mathbf{x}_t)\|_2^2 \leq \frac{B^2}{2\mu} \sum_{t=1}^T \frac{1}{t} \leq \frac{B^2}{2\mu} (\ln(T) + 1)$
- Theorem F ($\mathcal{O}(1/T)$) f is μ -strongly convex and B -Lipschitz, then subgradient descent of step size $\gamma_t = 2/\mu(t+1)$ gives $\max\{\min_{1 \leq t \leq T} f(\mathbf{x}_t), f(\hat{\mathbf{x}}_T)\} - f^* \leq \frac{2B^2}{\mu \cdot (T+1)}$ where $\hat{\mathbf{x}}_T := \sum_{t=1}^T t \mathbf{x}_t / \sum_{t=1}^T t$.
 - Proof
 - By descent lemma $(f(\mathbf{x}_t) - f^*) \leq \frac{1-\mu\gamma_t}{2\gamma_t} \|\mathbf{x}_t - \mathbf{x}^*\|_2^2 - \frac{1}{2\gamma_t} \|\mathbf{x}_{t+1} - \mathbf{x}^*\|_2^2 + \frac{\gamma_t}{2} \|g_t\|_2^2$
 - Plug in $\gamma_t = 2/\mu(t+1)$ we get RHS = $\frac{\mu(t-1)}{4} \|\mathbf{x}_t - \mathbf{x}^*\|_2^2 - \frac{\mu(t+1)}{4} \|\mathbf{x}_{t+1} - \mathbf{x}^*\|_2^2 + \frac{1}{\mu(t+1)} \|g_t\|_2^2$
 - Times t on each side and we get $t(f(\mathbf{x}_t) - f^*) \leq \frac{\mu(t-1)}{4} \|\mathbf{x}_t - \mathbf{x}^*\|_2^2 - \frac{\mu(t+1)t}{4} \|\mathbf{x}_{t+1} - \mathbf{x}^*\|_2^2 + \frac{t}{\mu(t+1)} \|g_t\|_2^2$
 - Since $\frac{t}{\mu(t+1)} \|g_t\|_2^2 \leq \frac{1}{\mu} B^2$, sum over $t = [1 : T]$ and get $\sum_{t=1}^T t(f(\mathbf{x}_t) - f^*) \leq -\frac{\mu T(T+1)}{4} \|\mathbf{x}_{T+1} - \mathbf{x}^*\|_2^2 + \frac{T}{\mu} B^2$
 - else is straightforward.

Lower Bound Complexity

We can show that, the above $\mathcal{O}(1/\sqrt{t})$ and $\mathcal{O}(1/t)$ can not be improved. The worst case function is given by the piecewise-linear function similar to $f(\mathbf{x}) = \max_{1 \leq i \leq n} x_i$.

- **Theorem 6.20 (Nemirovski & Yudin 1979).** For any $1 \leq t \leq n, \mathbf{x}_1 \in \mathbb{R}_n$,
 - (i) **B -Lipschitz** there exists a B -Lipschitz continuous function f and a convex set X with diameter R , such that for any first-order method that generates $\mathbf{x}_t \in \mathbf{x}_1 + \text{span}(\mathbf{g}_1, \dots, \mathbf{g}_{t-1})$, where $\mathbf{g}_i \in \partial f(\mathbf{x}_i)$, we have $\min_{1 \leq s \leq t} f(\mathbf{x}_s) - f^* \geq \frac{B \cdot R}{4(1+\sqrt{t})}$.
 - (ii) **B -Lipschitz and μ -convex** there exists a μ -strongly convex, B -Lipschitz continuous function f and a convex set X with diameter R , for any first-order method as described above, we always have $\min_{1 \leq s \leq t} f(\mathbf{x}_s) - f^* \geq \frac{B^2}{8\mu t}$.
 - **Proof**
 - The function we construct is $f(\mathbf{x}) = C \cdot \max_{1 \leq i \leq t} x_i + \frac{\mu}{2} \|\mathbf{x}\|_2^2$, and we restrict ourself within region $X = \{\mathbf{x} \in \mathbb{R}^n, \|\mathbf{x}\|_2 \leq R\}$.
 - The subgradient is $\partial f(\mathbf{x}) = \mu \mathbf{x} + C \cdot \text{conv} \{e_i : i \text{ such that } x_i = \max_{1 \leq j \leq t} x_j\}$.
 - The optimal solution is $\mathbf{x}_i^* = \begin{cases} -\frac{C}{\mu t} & 1 \leq i \leq t \\ 0 & t < i \leq n \end{cases}$ and $f^* = -\frac{C^2}{2\mu t}$, can be verified by $0 \in \partial f(\mathbf{x})$.
 - Since all subgradient method controls γ_T , but not how to choose subgradient, we can design a process that gives worst case. The update is design to have $g(\mathbf{x}) = C \cdot e_i + \mu \mathbf{x}$, only one pure direction of $x_i = \max_{1 \leq j \leq t} x_j$, but not its combination.
 - To simplify, we always choose the smallest one.
 - We claim that if we start from $\mathbf{x}_{s=1} = \mathbf{0}$, then $\mathbf{x}_s \in \text{span}(e_1, \dots, e_{s-1})$, always one dimension less than s .
 - When $s = 2$, all coordinates are equal and we choose e_1 to update, so $\mathbf{x}_{s=2} \in \text{span}\{e_1\}$.
 - By our update rule $g(\mathbf{x}) = C \cdot e_i + \mu \mathbf{x}$, the span can only increase at most one dimension at each update, so induction still holds.
 - This means for $\mathbf{x}_s, 1 \leq s \leq t$, at least one coordinate is unchanged as 0. So $\max_{1 \leq i \leq t} x_i \rightarrow f(\mathbf{x}_s) \geq 0$.
 - This gives $\min_{1 \leq s \leq t} f(\mathbf{x}_s) - f^* \geq \frac{C^2}{2\mu t}$
 - Choosing $C = \frac{B\sqrt{t}}{1+\sqrt{t}}, \mu = \frac{2B}{R(1+\sqrt{t})}$, we have $\|\partial f(\mathbf{x})\|_2 \leq C + \mu \|\mathbf{x}\|_2 \leq C + \mu R =: M$, this is the M -Lipschitz case, we have $\min_{1 \leq s \leq t} f(\mathbf{x}_s) - f^* \geq \frac{C^2}{2\mu t} = \frac{B \cdot R}{2(1+\sqrt{t})}$.
 - Choosing $C = \frac{B}{2}, \mu = \frac{B}{R}$, we have $\|\partial f(\mathbf{x})\|_2 \leq C + \mu R =: M$, this is M -Lispchitz and μ -strongly convex case, we have $\min_{1 \leq s \leq t} f(\mathbf{x}_s) - f^* \geq \frac{C^2}{2\mu t} = \frac{B^2}{8\mu t}$
 - **Comment** This theorem is so strange that number of updates is related to dimension.

Mirror Descent

- **Key Idea** Update can be viewed as $\mathbf{x}_{t+1} = \underset{\mathbf{x} \in X}{\text{argmin}} \left\{ \frac{1}{2} \|\mathbf{x} - \mathbf{x}_t\|_2^2 + \langle \gamma_t g(\mathbf{x}_t), \mathbf{x} \rangle \right\}$, how about we change $\|\cdot\|_2^2$ to something else.

Bregman Divergence

- **Definition 6.22** Let $\omega(\mathbf{x}) : X \rightarrow \mathbb{R}$ be a function that is *strictly* convex, continuously differentiable on a closed convex set X . The *Bregman divergence* is defined as $V_\omega(\mathbf{x}, \mathbf{y}) = \omega(\mathbf{x}) - \omega(\mathbf{y}) - \nabla \omega(\mathbf{y})^T (\mathbf{x} - \mathbf{y}), \forall \mathbf{x}, \mathbf{y} \in X$.
 - Asymmetric $V_\omega(\mathbf{x}, \mathbf{y}) \neq V_\omega(\mathbf{y}, \mathbf{x})$, not a valid distance, triangle inequality may not hold.
 - $\omega(\cdot)$ is called *distance-generating function*.
 - If $\omega(\cdot)$ is also σ -strongly convex w.r.t. norm $\|\cdot\|_a$, then $V_\omega(\mathbf{x}, \mathbf{y}) \geq \sigma \|\mathbf{x} - \mathbf{y}\|_a^2 / 2$.
- **Example**
 - *Euclidean Distance* $\Omega = \mathbb{R}^d, \omega(\mathbf{x}) = \frac{1}{2} \|\mathbf{x}\|_2^2$, ω is 1 strongly convex w.r.t ℓ_2 norm,
 - $V_\omega(\mathbf{x}, \mathbf{y}) = \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2$.
 - *Mahalanobis distance* $\Omega = \mathbb{R}^d, \omega(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T Q \mathbf{x}$ where $Q \succeq I$, ω is 1 strongly convex w.r.t ℓ_2 norm,
 - $V_\omega(\mathbf{x}, \mathbf{y}) = \frac{1}{2} (\mathbf{x} - \mathbf{y})^T Q (\mathbf{x} - \mathbf{y})$.
 - *Kullback-Leibler divergence* $\Omega = \Delta_d, \omega(\mathbf{x}) = \sum_{i=1}^d x_i \ln x_i, \nabla \omega(\mathbf{x}) = \ln(\mathbf{x}) - 1$,
 - $V_\omega(\mathbf{x}, \mathbf{y}) = \sum_i x_i \ln x_i - y_i \ln y_i - (\ln(y_i) - 1)(x_i - y_i) = \sum_i x_i \ln \frac{x_i}{y_i} = \text{KL}(\mathbf{x} \parallel \mathbf{y})$, since $\sum_i x_i = \sum_i y_i = 1$.
 - The proof for $V_\omega(\mathbf{x}, \mathbf{y}) \geq \|\mathbf{x} - \mathbf{y}\|_1^2 / \ln 4$ is complicated, see [this](#).
- **Lemma 6.23 (Generalized Pythagorean Theorem, Exercise 46)** If \mathbf{x}^* is the Bregman projection of \mathbf{x}_0 onto a convex set

$C \subset X : x^* \operatorname{argmin}_{x \in C} V_\omega(x, x_0)$. Then $\forall y \in C, V_\omega(y, x_0) \geq V_\omega(y, x^*) + V_\omega(x^*, x_0)$.

◦ Proof

- $\nabla_x V_\omega(x, y) = \nabla \omega(x) - \nabla \omega(y)$, optimal condition means $\forall y \in C, (\nabla \omega(x^*) - \nabla \omega(x_0))^\top (y - x^*) \geq 0$.
- equivalently, $-\nabla \omega(x_0)^\top (y - x_0 + x_0 - x^*) \geq -\nabla \omega(x^*)^\top (y - x^*)$
- or $-\nabla \omega(x_0)^\top (y - x_0) \geq -\nabla \omega(x^*)^\top (y - x^*) - \nabla \omega(x_0)^\top (x_0 - x^*)$
- then $\omega(x_0) - \omega(y) - \nabla \omega(x_0)^\top (y - x_0) \geq \omega(x_0) - \omega(x^*) + \omega(x^*) - \omega(y) - \nabla \omega(x^*)^\top (y - x^*) + -\nabla \omega(x_0)^\top (x_0 - x^*)$, QED.

Mirror Descent Algorithm

- Prox-mapping $\operatorname{Prox}_x(\xi) = \operatorname{argmin}_{u \in X} \{V_\omega(u, x) + \langle \xi, u \rangle\}$, and suppose ω is 1-strongly convex on norm $\|\cdot\|_a$.
- Algorithm $x_{t+1} = \operatorname{Prox}_{x_t}(\gamma_t g(x_t))$, where $g(x_t) \in \partial f(x_t)$.
- Example 6.25 Under KL divergence, the prox-mapping becomes $\operatorname{Prox}_x(\xi) = (\sum_{i=1}^n x_i e^{-\xi_i})^{-1} \begin{bmatrix} x_1 e^{-\xi_1} \\ \dots \\ x_n e^{-\xi_n} \end{bmatrix}$
- Proof
 - $V_\omega(u, x) + \langle \xi, u \rangle = \sum_i u_i \ln(u_i/x_i) + u_i \xi_i$, optimal condition under constraint is $0 = \ln(u_i/x_i) + 1 + \lambda + \xi_i$
 - So solution is $u_i = x_i e^{-\xi_i + \alpha}$ where $\sum_i u_i = 1 \Rightarrow e^{-\alpha} = \sum_i x_i e^{-\xi_i}$ QED.

$\mathcal{O}(1/\sqrt{t})$ convergence for convex functions

- Lemma 6.26 (Three point identity) $\forall x, y, z \in \operatorname{dom}(\omega), V_\omega(x, z) = V_\omega(x, y) + V_\omega(y, z) - \langle \nabla \omega(z) - \nabla \omega(y), x - y \rangle$
- Proof
 - $V_\omega(x, y) + V_\omega(y, z) = \omega(x) - \omega(y) + \omega(y) - \omega(z) - \langle \nabla \omega(y), x - y \rangle - \langle \nabla \omega(z), y - z \rangle$
 - RHS = $V_\omega(x, z) + \langle \nabla \omega(z), x - z \rangle - \langle \nabla \omega(y), x - y \rangle - \langle \nabla \omega(z), y - z \rangle$
 - RHS = $V_\omega(x, z) + \langle \nabla \omega(z) - \nabla \omega(y), x - y \rangle$
- When $\omega = \|\cdot\|_2^2$, this becomes law of cosines.
- Lemma G (Descent Lemma) $\gamma_t (f(x_t) - f^*) \leq V_\omega(x^*, x_t) - V_\omega(x^*, x_{t+1}) + \frac{\gamma_t^2}{2} \|g_t\|_{a^*}^2$
- Proof
 - By 3-point identiy, $V_\omega(x^*, x_t) = V_\omega(x^*, x_{t+1}) + V_\omega(x_{t+1}, x_t) - \langle \nabla \omega(x_t) - \nabla \omega(x_{t+1}), x^* - x_{t+1} \rangle$, equivalently $V_\omega(x^*, x_t) - V_\omega(x^*, x_{t+1}) + \langle \nabla \omega(x_t) - \nabla \omega(x_{t+1}), x^* - x_{t+1} \rangle = V_\omega(x_{t+1}, x_t)$
 - by def of algo, the optimial condition is $\forall x \in X$, especially x^*
 $\langle \nabla_{x_{t+1}} V_\omega(x_{t+1}, x_t) + \gamma_t g_t, x^* - x_{t+1} \rangle = \langle \nabla \omega(x_{t+1}) - \nabla \omega(x_t) + \gamma_t g_t, x^* - x_{t+1} \rangle \geq 0$, this means
 $V_\omega(x^*, x_t) - V_\omega(x^*, x_{t+1}) + \langle \gamma_t g_t, x^* - x_{t+1} \rangle \geq V_\omega(x_{t+1}, x_t)$
 - By convexity, $\langle g_t, x^* - x_t \rangle \leq f(x^*) - f(x_t)$, so
 $V_\omega(x^*, x_t) - V_\omega(x^*, x_{t+1}) + \gamma_t (f(x^*) - f(x_t)) \geq V_\omega(x_{t+1}, x_t) + \langle \gamma_t g_t, x_{t+1} - x_t \rangle$
 - By Young's inequality, $\langle a, b \rangle \leq \|a\|_a \cdot \|b\|_{a^*} \leq \|a\|_a^2/2 + \|b\|_{a^*}^2/2$, therefore
 $\langle \gamma_t g_t, x_t - x_{t+1} \rangle \leq \gamma_t^2 \|g_t\|_{a^*}^2/2 + \|x_t - x_{t+1}\|_a^2/2$
 - By the assumption of ω is 1-strongly convex, $V_\omega(x_{t+1}, x_t) \geq \|x_t - x_{t+1}\|_a^2/2$
 - Combine the above two,
 $V_\omega(x_{t+1}, x_t) + \langle \gamma_t g_t, x_{t+1} - x_t \rangle \geq \|x_t - x_{t+1}\|_a^2/2 - \gamma_t^2 \|g_t\|_{a^*}^2/2 - \|x_t - x_{t+1}\|_a^2/2 = -\gamma_t^2 \|g_t\|_{a^*}^2/2$
 - Combine above, we arrive at the descent lemma.

- Theorem 6.28 f convex, then $\max \left\{ \min_{1 \leq t \leq T} f(x_t), f(\hat{x}_T) \right\} - f^* \leq \frac{V_\omega(x^*, x_1) + \frac{1}{2} \sum_{t=1}^T \gamma_t^2 \|g(x_t)\|_*^2}{\sum_{t=1}^T \gamma_t}$, where
 $\hat{x}_T = \frac{\sum_{t=1}^T \gamma_t x_t}{\sum_{t=1}^T \gamma_t}$.

- Proof Similar to Theorem 6.17, sum over all $t \in [1 : T]$ we get the result.

- Convergence similar to subgradient case $\min_{1 \leq t \leq T} f(x_t) - f^* = O\left(\frac{BR}{\sqrt{T}}\right)$, where $R = \sqrt{\max_{\mathbf{x} \in X} V_\omega(\mathbf{x}, \mathbf{x}_1)}$ and $B := \sup_{\mathbf{x} \in X} \frac{|f(\mathbf{x}) - f(\mathbf{y})|}{\|\mathbf{x} - \mathbf{y}\|}$.

- The constant may be different. Example of Simplex of $X = \left\{ \mathbf{x} \in \mathbb{R}^d : \mathbf{x}_i \geq 0, \sum_{i=1}^d \mathbf{x}_i = 1 \right\}$, assume $\|\mathbf{g}\|_\infty \leq 1$,
 - In normal case of $\omega(\mathbf{x}) = \frac{1}{2} \|\mathbf{x}\|_2^2$, $R^2 \leq 2 = O(1)$, $B \sim \|\mathbf{g}\|_2 = O(\sqrt{d})$, overall it's $O(\frac{\sqrt{d}}{\sqrt{T}})$.
 - If we choose $\omega(\mathbf{x}) = \sum_{i=1}^d \mathbf{x}_i \ln \mathbf{x}_i$ and starting point to be $\mathbf{x}_1 = \operatorname{argmin}_{\mathbf{x} \in X} \omega(\mathbf{x})$, then $\Omega \leq \max_{\mathbf{x} \in X} \omega(\mathbf{x}) - \min_{\mathbf{x} \in X} \omega(\mathbf{x}) = 0 - (-\ln d) = \ln(d)$, overall it's $O(\frac{\ln d}{\sqrt{T}})$
 - The ratio of efficiency is then $O\left(\frac{1}{\ln(d)} \cdot \frac{\max_{\mathbf{x} \in X} \|g(\mathbf{x})\|_2}{\max_{\mathbf{x} \in X} \|g(\mathbf{x})\|_\infty}\right)$, since $\|g(\mathbf{x})\|_\infty \leq \|g(\mathbf{x})\|_2 \leq \sqrt{d} \|g(\mathbf{x})\|_\infty$, in worst case Mirror Descent is $O(\sqrt{d})$ faster than norm subgradient descent.

Mirror Descent under Smoothness (Exercise 47)

- Lemma H If f convex and gradient Lipschitz w.r.t. some norm, $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_{a*} \leq L \|\mathbf{x} - \mathbf{y}\|_a$, then setting $\gamma_t = 1/L$ will give $\min_{1 \leq t \leq T} f(\mathbf{x}_t) - f^* \leq L \cdot V_\omega(\mathbf{x}^*, \mathbf{x}_1)/T$

- Proof

- By equivalence of smoothness and Lipschitz, $f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) + \mathbf{g}_t^\top (\mathbf{x}_{t+1} - \mathbf{x}_t) + \frac{L}{2} \|\mathbf{x}_t - \mathbf{x}_{t+1}\|_a^2$
 - Using the fact of ω is 1-strongly convex, we have $\|\mathbf{x}_t - \mathbf{x}_{t+1}\|_a^2/2 \leq V_\omega(\mathbf{x}_{t+1}, \mathbf{x}_t)$
 - This makes $\frac{1}{\gamma_t} V_\omega(\mathbf{x}_{t+1}, \mathbf{x}_t) + \langle \mathbf{g}_t, \mathbf{x}_{t+1} - \mathbf{x}_t \rangle \geq f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t)$
 - By the optimal condition of prox-mapping, we can show that $V_\omega(\mathbf{x}_{t+1}, \mathbf{x}_t) + \langle \gamma_t \mathbf{g}_t, \mathbf{x}_{t+1} \rangle \leq V_\omega(\mathbf{x}_t, \mathbf{x}_t) + \langle \gamma_t \mathbf{g}_t, \mathbf{x}_t \rangle = 0$
 - This means $f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t) \leq 0$, always non-increasing.
- In the third step of descent lemma, by convexity we have $V_\omega(\mathbf{x}^*, \mathbf{x}_t) - V_\omega(\mathbf{x}^*, \mathbf{x}_{t+1}) + \gamma_t(f(\mathbf{x}^*) - f(\mathbf{x}_t)) \geq V_\omega(\mathbf{x}_{t+1}, \mathbf{x}_t) + \langle \gamma_t \mathbf{g}_t, \mathbf{x}_{t+1} - \mathbf{x}_t \rangle$
 - take the above into it we get $V_\omega(\mathbf{x}^*, \mathbf{x}_t) - V_\omega(\mathbf{x}^*, \mathbf{x}_{t+1}) + \gamma_t(f(\mathbf{x}^*) - f(\mathbf{x}_t)) \geq \gamma_t(f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t))$
- take $\gamma_t = L^{-1}$ and sum over $t = [0 : T - 1]$, we have
 - $\sum_{i=1}^T f(\mathbf{x}_i) - f(\mathbf{x}^*) \leq L(V_\omega(\mathbf{x}^*, \mathbf{x}_1) - V_\omega(\mathbf{x}^*, \mathbf{x}_T))$, QED