# Lecture 12: Subgradient methods    E1 260

- Subgradients

- Subgradient methods

- Convergence analysis
    - Lipschitz Convex functions
    - Strong Convexity

# Gradient descent method:

$$\underline{x}_{t+1} = \underline{x}_t - \eta_t \nabla f(\underline{x}_t)$$

Differentiability of the objective function $f$ is essential
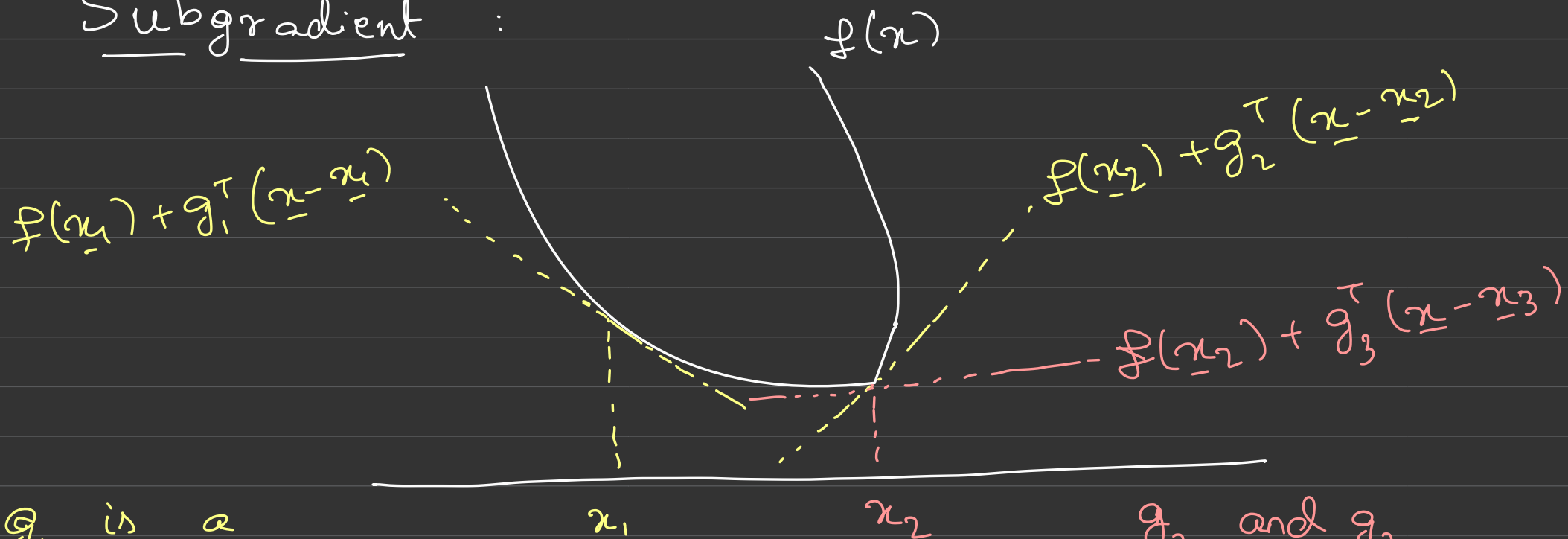
How about

minimize $\|\underline{x}\|_1$

s.t. $\|A\underline{x} - \underline{b}\| \leq \varepsilon$

minimize $\|x\|_*$

s.to $\|P_\Omega(x - y)\|_F^2 \leq \varepsilon$

# Subgradient :

$f(x)$

$f(x_1) + g_1^T (x - x_1)$

$f(x_2) + g_2^T (x - x_2)$

$f(x_2) + g_3^T (x - x_3)$

$x_1$      $x_2$

$g_1$ is a subgradient at $x_1$

$g_2$ and $g_3$ are subgradients at $x_2$

- $g$ is a subgradient of $f$ at $x$ if

$$f(y) \geqslant f(x) + g^T (y - x), \quad \forall y$$

a global linear underestimate of $f$

- Convexity is equivalent to the existence of subgradients everywhere

- if a function is convex and differentiable, $\nabla f(x)$ is a subgradient of a $f$ at $x$

- the set of subgradients of $f$ at $x$ is called the subdifferential of $f$ at $\underline{x}$
  $$\partial f(\underline{x})$$

$f(x)$

epi $(f)$

$t \geqslant f(y) \geqslant f(x) + g^T(y-x)$
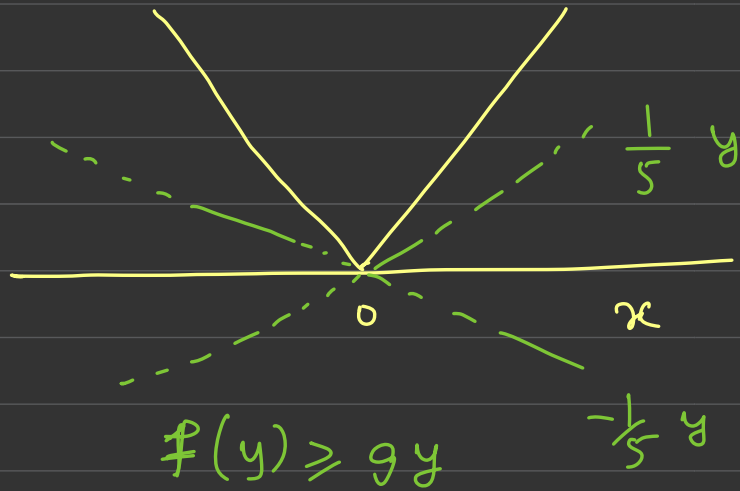
$(g, -1)$

$\underline{g}$ is a subgradient of $f$ at $x$ iff

$(\underline{g}, -1)$ defines a supporting hyperplane of epi($f$) at $(x, f(x))$

$(y, t) \in \text{epi}(f) \implies \begin{bmatrix} \underline{g} \\ -1 \end{bmatrix}^T \left( \begin{bmatrix} y \\ t \end{bmatrix} - \begin{bmatrix} x \\ f(x) \end{bmatrix} \right) \leq 0$
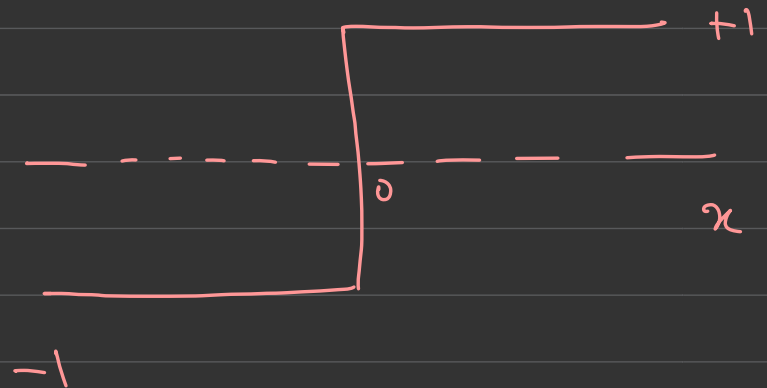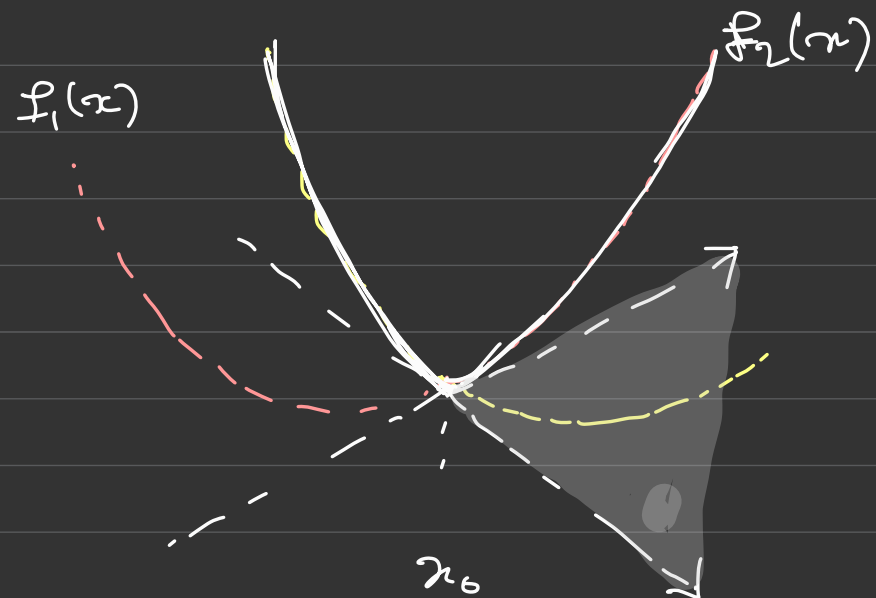
# Examples:

$$f(x) = |x|$$

$$\partial f(x) = \begin{cases} -1 & \text{, if } x < 0 \\ [-1, 1] & \text{if } x = 0 \\ 1 & \text{if } x > 0 \end{cases}$$

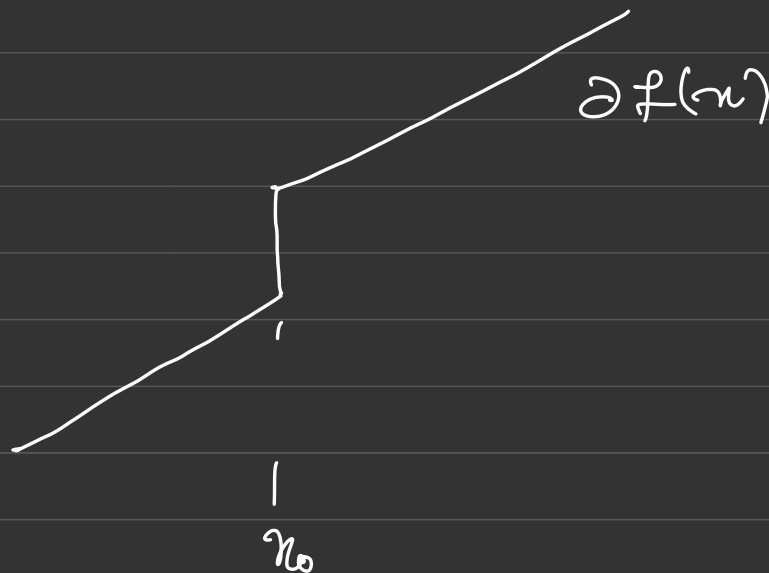$$f(y) \geqslant g y$$

$$g \in [-1, 1] \quad \text{at } 0$$

$$f(\underline{x}) = \|\underline{x}\|_1 = \sum_{i=1}^{d} |x_i| = \sum_{i=1}^{d} f_i(x)$$

$$\partial f(\underline{x}) = \begin{cases} sgn(x_i) \, \underline{e}_i & \text{if } x_i \neq 0 \\ [-1, 1] \, \underline{e}_i & \text{if } x_i = 0 \end{cases}$$

**Example:** $f(x) = \max\{f_1(x), f_2(x)\}$

$f_1(x)$

$f_2(x)$

$x_0$

$$\partial f(x) = \begin{cases} \nabla f_1(x_0), & \text{if } f_1(x_0) > f_2(x_0) \\ [\nabla f_1(x_0), \nabla f_2(x_0)], & \text{if } \\ \qquad\qquad\qquad f_1(x_0) = f_2(x_0) \\ \nabla f_2(x_0), & \text{if } f_2(x_0) > f_1(x_0) \end{cases}$$

$\partial f(x)$

$x_0$

# Optimality Condition

- Recall for convex and differentiable $f$

$$f(\underline{x}^*) = \inf_x f(x) \iff 0 = \nabla f(\underline{x}^*)$$

- For differentiable $f$, $\nabla f(x) = 0$ we can only say $\underline{x}$ is a critical point.
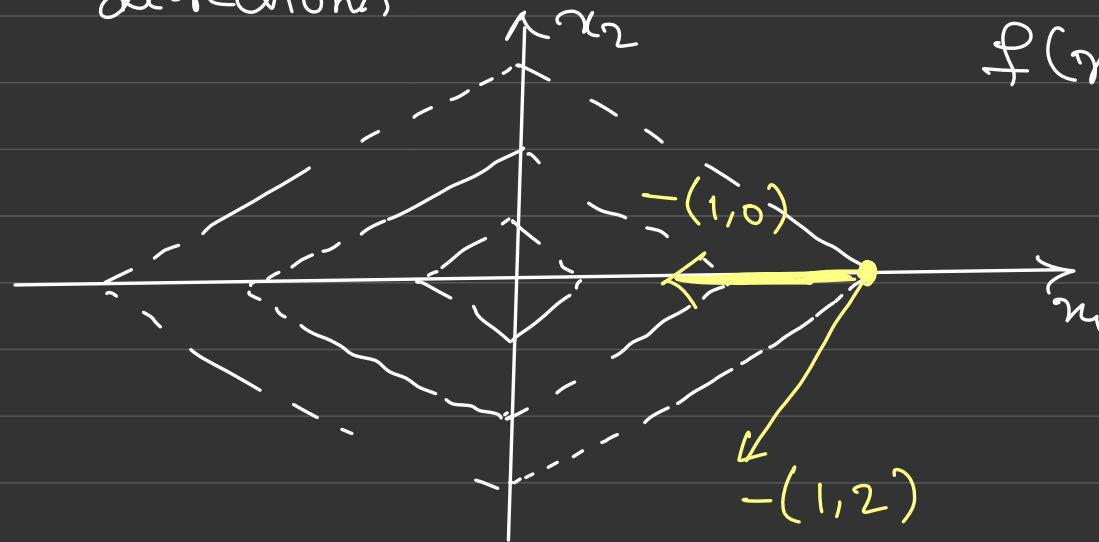
- Suppose $f : \text{dom}(f) \longrightarrow \mathbb{R}$ and $\underline{x} \in \text{dom}(f)$. if $\underline{0} \in \partial f(\underline{x})$, then $\underline{x}$ is "a global minimum"

If $\underline{g} = \underline{0} \in \partial f(x)$, $f(y) \geq f(x) + \underline{g}^T(y - x) = f(\underline{x})$

$$\forall\, y \in \text{dom}(f)$$

# Descent direction:

- Negative subgradients are not necessarily descent directions

$$f(x) = |x_1| + 2|x_2|$$



at $\underline{x} = (1,0)$

- $\underline{g}_1 = (1,0) \in \partial f(x)$ and $-g_1$ is a descent direction

- $g_2 = (1,2) \in \partial f(x)$, but $-g_2$ is not

# The algorithm :

$$\implies \quad x_{t+1} = x_t - \eta_t g_t \quad ; \quad g_t \in \partial f(x)$$

$$\implies \quad x_{t+1} = P_c \left( x_t - \eta_t g_t \right) \; ; \; g_t \in \partial f(x)$$

- Since $f(x_{t+1})$ is not necessarily monotone, we also keep track of

$$f_t^{best} := \min_{1 \le i \le t} f(x_i)$$

- Define $f^{opt} = \min_{x} f(x)$

# Fundamental inequality for Projected Subgradient methods:

$$\| x_{t+1} - x^* \|_2^2 \leq \| x_t - x^* \|_2^2 - 2\eta_t \left( f(x_t) - f^{opt} \right)$$

$$+ \underbrace{\eta_t^2 \| g_t \|_2^2}_{\text{majorizing function}}$$

We wish to optimize $\| x_{t+1} - x^* \|_2^2$, but without access to $x^*$ we optimize by finding another function that majorizes it.

$$\| x_{t+1} - x^* \|_2^2 = \| P_C ( x_t - \eta_t g_t ) - P_C (x^*) \|_2^2$$

From non expansiveness

$$\leq \| x_t - \eta_t g_t - x^* \|_2^2$$

$$= \| x_t - x^* \|_2^2 - 2 \eta_t g_t^T ( x_t - x^* ) + \eta_t^2 \| g_t \|_2^2$$

$$\leq \| x_t - x^* \|_2^2 - 2 \eta_t [ f ( x_t ) - f ( x^* ) ] + \eta_t^2 \| g_t \|_2^2$$

as

$$f ( x^* ) - f ( x_t ) \geq g_t ( x^* - x_t )$$

# Polyak step size :

$$\eta_t = \frac{f(x_t) - f^{opt}}{\| g_t \|^2}$$

- we get an error reduction

$$\| x_{t+1} - x^* \|_2^2 \leq \| x_t - x^* \|_2^2 - \frac{(f(x_t) - f(x^*))^2}{\| g_t \|_2^2}$$

- But needs $f^{opt}$ to be known

Suppose $f$ is convex and $B$ Lipschitz
Continuous

Then

① $\|g\| \leq B$ $\quad\quad$ $\forall$ $\underline{g} \in \partial f(\underline{x})$

② $|f(x) - f(y)| \leq B \|\underline{x} - \underline{y}\|$ $\quad\quad$ $\forall \underline{x}, y \in \partial nf$

Claim: The projected subgradient descent with

$\quad\quad$ Polyak's step size rule satisfies

$$f_T^{best} - f^{opt} \leq \frac{B}{\sqrt{T}} \|\underline{x}_0 - \underline{x}^*\|_2$$

- Sublinear convergence rate of $O\left(\frac{1}{\sqrt{T}}\right)$

We had for Polyak's step size rule:

$$\|\underline{x}_{t+1} - \underline{x}_t^*\|_2^2 \leq \|\underline{x}_t - x^*\|_2^2 - \frac{(f(\underline{x}_t) - f(x^*))^2}{\|g_t\|_2^2}$$

$$\Rightarrow \left(f(\underline{x}_t) - f(\underline{x}^*)\right)^2 \leq \Big[\|\underline{x}_t - \underline{x}^*\|_2^2 -$$

$$\|\underline{x}_{t+1} - \underline{x}^*\|_2^2\Big] B^2$$

Apply recursively and sum over iterations

$$t = 0 \text{ to } T-1$$

$$\sum_{t=0}^{T-1} \left(f(\underline{x}_t) - f(\underline{x}^*)\right)^2 \leq B^2 \Big[\|\underline{x}_0 - x^*\|_2^2 - \|\underline{x}_T - \underline{x}^*\|_2^2\Big]$$

$$\Rightarrow T\left(f_{T-1}^{best} - f(\underline{x}^*)\right)^2 \leq B^2 \|\underline{x}_0 - x^*\|_2^2$$

$$\Rightarrow f_t^{best} - f^{opt} \leq \frac{B}{\sqrt{T}} \|\underline{x}_0 - \underline{x}^*\|$$

How about other step sizes (diminishing ?)

Claim:

Suppose $f$ is convex and $B$ Lipschitz continuous.
The projected subgradient method

$$f_T^{best} - f^{opt} \leq \frac{\| x_0 - x^* \| + B^2 \sum_{t=0}^{T-1} \eta_t^2}{2 \sum_{t=0}^{T-1} \eta_t}$$

$$\|x_{t+1} - x^*\|_2^2 \leq \|x_t - x^*\|_2^2 - 2\eta_t \left[ f(x_t) - f(x^*) \right] + \eta_t^2 \|g_t\|_2^2$$

$$\|x_T - x^*\|_2^2 \leq \|x_0 - x^*\|_2^2 - 2 \sum_{t=0}^{T-1} \eta_t \left( f(x_t) - f^{opt} \right)$$
$$+ \sum_{t=0}^{T-1} \eta_t^2 \|g_t\|_2^2$$

$$\Rightarrow$$

$$2 \sum_{t=0}^{T-1} \eta_t \left( f(x_t) - f^{opt} \right) \leq \|x_0 - x^*\|_2^2 - \|x_T - x^*\|_2^2$$
$$+ \beta^2 \sum_{t=0}^{T-1} \eta_t^2$$

But
$$2 \sum_{t=0}^{T-1} \eta_t \left( f_t^{best} - f^{opt} \right) \leq 2 \sum_{t=0}^{T-1} \eta_t \left( f(x_t) - f^{opt} \right)$$

or
$$f_T^{best} - f^{opt} \leq \sum_{t=0}^{T-1} \eta_t \left( f(x_t) - f^* \right) \Big/ \sum_{t=0}^{T-1} \eta_t$$

$$f_T^{best} - f^{opt} \leq \frac{\|x_0 - x^*\|_2^2 + B^2 \sum_{t=0}^{T-1} \eta_t^2}{2 \sum_{t=0}^{T-1} \eta_t}$$

- Constant step-size: $\eta_t = \eta$

$$\lim_{T \to \infty} f_T^{best} \leq \frac{B^2 \eta}{2}$$

may not converge to optimal points

- Diminishing step size: $\sum_{t=0}^{T-1} \eta_t^2 < \infty$ and $\sum_{t=0}^{T-1} \eta_t \to \infty$

$$\lim_{T \to \infty} f_T^{best} = 0 \qquad\qquad e.g. \ \frac{a}{b+t}$$

$$a > 0 \ ; \ b \geq 0$$

Converges to optimal points.

**Strongly convex** : $O\left(\frac{1}{\varepsilon}\right)$

better than $O\left(\frac{1}{\varepsilon^2}\right)$, worse than $O\left(\log\left(\frac{1}{\varepsilon}\right)\right)$

**Claim:**

$f: \mathbb{R}^d \rightarrow \mathbb{R}$ is $\mu$-strongly convex and $\underline{x}^*$ be

the unique minimizer of $f$. With $\eta_t = \dfrac{2}{\mu(t+1)}$.

Then Subgradient method, yields

$$f\left(\frac{2}{T(T+1)} \sum_{t=1}^{T} t \cdot x_t\right) - f(\underline{x}^*) \leq \frac{2\beta^2}{\mu(T+1)}$$

where $\qquad B = \max_{1 \leq t \leq T} \|g_t\|$.

Recall :

$$g_t^T (\underline{x}_t - \underline{x}^*) = \frac{n_t}{2} \|g_t\|^2 + \frac{1}{2n_t} \left[ \|\underline{x}_t - \underline{x}^*\|^2 \right.$$

$$\overbrace{}^{\leq B^2}$$

$$\left. - \|\underline{x}_{t+1} - \underline{x}^*\|^2 \right]$$

use the quadratic lower bound :

$$g_t^T (\underline{x}_t - \underline{x}^*) \geq f(\underline{x}_t) - f(\underline{x}^\circ) + \frac{\mu}{2} \|\underline{x}_t - \underline{x}^*\|^2$$

$$\Rightarrow f(\underline{x}_t) - f(\underline{x}^*) \leq \frac{n_t}{2} B^2 + \frac{(n_t^{-1} - \mu)}{2} \|\underline{x}_t - \underline{x}^*\|^2$$

$$- \frac{n_t^{-1}}{2} \|\underline{x}_{t+1} - \underline{x}^*\|^2$$

- Unlike gradient descent with fixed step size
  we cannot telescope anymore when we
  sum over iterations

- To get a telescopic sum

$$\eta_t^{-1} = \eta_{t+1}^{-1} - \mu$$

One choice of $\eta_t$ that satisfies this is:

$$\eta_t^{-1} = \mu(t+1)$$

Actually our choice

$$\eta_t^{-1} = \mu(t+1)/2 \quad \text{does not}$$

Homework 2: Check what happens if we proceed

with $$\eta_t^{-1} = \mu(t+1)$$

$$t\left(f(x_t) - f(x^*)\right) \leq \frac{B^2 t}{\mu(t+1)} + \frac{\mu}{4}\left\{t(t-1)\|x_t - x^*\|^2\right.$$

$$\left. - (t+1)t \|x_{t+1} - x^*\|^2\right\}$$

$$\leq \frac{B^2}{\mu} + \frac{\mu}{4}\left\{t(t-1)\|x_t - x^*\|^2\right.$$

$$\left. - (t+1)t \|x_{t+1} - x^*\|^2\right\}$$

Summing over $t = 1$ to $T$:

$$\sum_{t=1}^{T} t\left(f(x_t) - f(x^*)\right) \leq T\frac{B^2}{\mu} + \frac{\mu}{4}\left[0 - x^*\|^2\right.$$

$$\left. T(T+1) \|x_T - x^*\|^2\right\} \leq \frac{TB^2}{\mu}$$

Since $\quad \dfrac{2}{T(T+1)} \displaystyle\sum_{t=1}^{T} t = 1$

and $\quad f(\cdot) \quad$ is convex $\quad$ ( Jensen's inequality)

$$f\left( \dfrac{2}{T(T+1)} \displaystyle\sum_{t=1}^{T} t\, \underline{x}_t \right) - f(x^*) \leq \dfrac{2}{T(T+1)} \displaystyle\sum_{t=1}^{T} t \cdot \left[ f(x_t) \right.$$

$$\left. - f(\underline{x}^*) \right]$$

$$\Rightarrow \quad f\left( \dfrac{2}{T(T+1)} \displaystyle\sum_{t=1}^{T} t \cdot x_t \right) - f(\underline{x}^*)$$

$$\leq \dfrac{2 B^2}{\mu (T+1)}$$