# Optimization:

$$\text{minimize} \quad f(\underline{x})$$
$$\underline{x} \in X$$

$$f : \mathbb{R}^D \longrightarrow \mathbb{R}$$

$\underline{x} \in \mathbb{R}^D$ is the design variable

$$X \subseteq \mathbb{R}^D \quad \text{Constraint set.}$$



local minima          Strict minima (local)          global minimum

- $\underline{x}^* \in X$ is a local minima if $\exists \, \epsilon > 0$

  s.t. $f(\underline{x}) \geqslant f(\underline{x}^*)$ $\forall \, \underline{x} \in X$ with
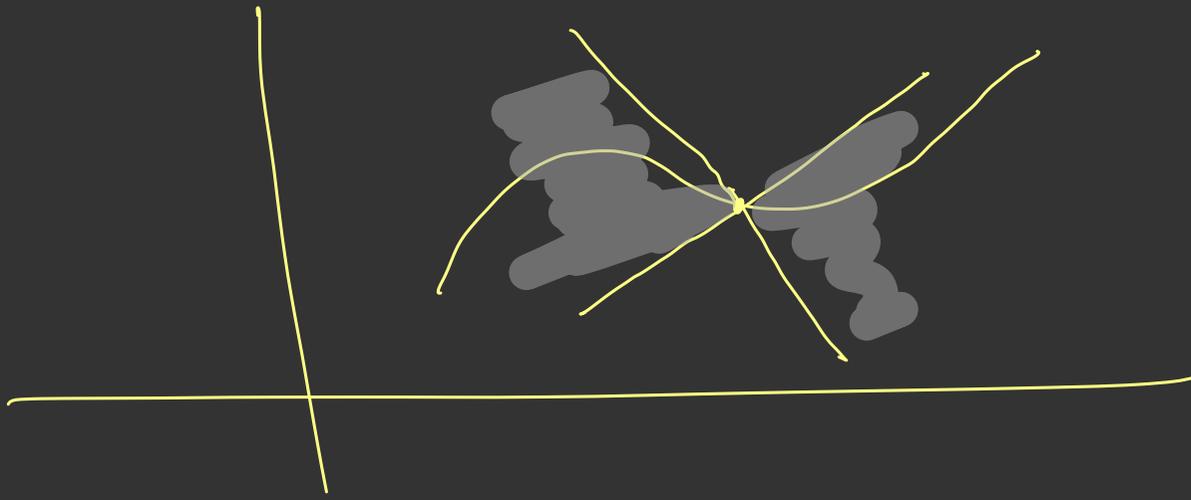
  $\| \underline{x} - \underline{x}^* \| \leq \epsilon$

- $\underline{x}^* \in X$ is a global minimum if $\exists \, \underline{x}^* \in X$

  s.t. $f(\underline{x}) \geqslant f(\underline{x}^*)$ $\forall \, \underline{x} \in X$
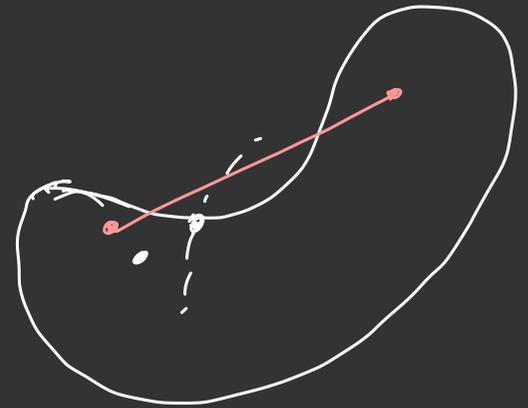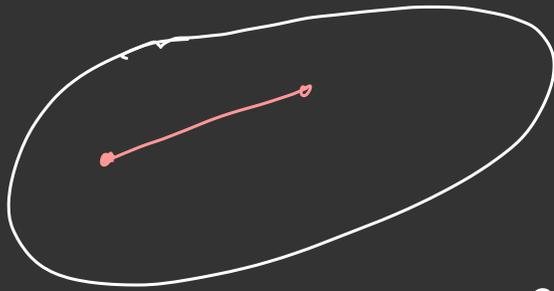
## Lipschitz - continuous function:

$f : \mathbb{R}^D \rightarrow \mathbb{R}$ is Lipschitz continuous

if $\exists \, L \geqslant 0$

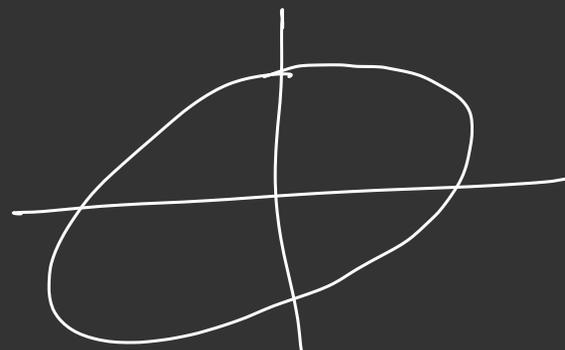$$\| f(\underline{x}) - f(\underline{y}) \| \leq L \, \| x - y \|$$

# Convex Sets:

A subset $x \subseteq \mathbb{R}^d$ is convex if

$$\theta \underline{x}_1 + (1-\theta) \underline{x}_2 \in x \quad \forall \; \underline{x}_1, \underline{x}_2 \in x$$

$$\theta \in [0, 1]$$

E.g.     norm ball

$$\chi = \left\{ \underline{x} \mid \underline{x}^{T} \overset{\Sigma}{\underset{\downarrow}{\chi}} \underline{x} \leqslant 1 \right\}$$

Convex function:



$f$ is convex if : $\forall \underline{x}, \underline{y} \in \chi$ , $\theta \in [0, 1]$

$$f(\theta \underline{x} + (1-\theta) y) \leq \theta f(x) + (1-\theta) f(y)$$

OR,

$$f(\underline{y}) \geq f(\underline{x}) + \nabla f^T(\underline{x})(\underline{y} - \underline{x})$$

$$\forall \underline{x}, y \in X$$

OR,

$$\nabla^2 f(x) \succeq 0$$

# Auto-diff:

## Reverse mode auto-diff:

$$f(x, y) = x^2y + y + 2$$



$\partial f/\partial n_7 = 1$

$n_7$ : $+$ 42

$\partial f/\partial n_5 = \partial f/\partial n_7 \times \partial n_7/\partial n_5$
$= 1 \times 1 = 1$

$\partial f/\partial n_6 = \partial f/\partial n_7 \times \partial n_7/\partial n_6$
$= 1 \times 1 = 1$

$n_5$ : $\times$ 36

$n_6$ : $+$ 6

$\partial f/\partial n_4 = \partial f/\partial n_5 \times \partial n_5/\partial n_4$
$= 1 \times n_2 = 4$

(1)
$\partial f/\partial n_2 = \partial f/\partial n_5 \times \partial n_5/\partial n_2$
$= 1 \times n_4 = 9$

$n_4$ : $\times$ 9

(2)
$\partial f/\partial n_2 = 1$

(1)    (2)

$n_2$ : y 4

$n_3$ : 2 2

$n_1$ : x 3

(1)    (2)
$\partial f/\partial y = \partial f/\partial n_2 + \partial f/\partial n_2 = 9 + 1 = 10$

(1)    (2)
$\partial f/\partial x = n_1 \times 4 + n_1 \times 4 = 24$

# Subgradients :

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

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

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

$$f(x_3) + g_4^T(x - x_3)$$

$x_1$     $x_3$   $x_2$

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

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

## Example:

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

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

# Iterative descent methods:

$$\min_{\underline{x} \in \mathbb{R}^d} f(x)$$



$$f(x_{k+1}) < f(x_k)$$

$$k = 0, 1, 2 \ldots$$

$$\underline{x}_{k+1} = x_k + \eta_k \underline{d}_k$$

$$f(\underline{x}_k + d_k) \approx f(x_k) + \nabla f(x_k)^T d_k$$
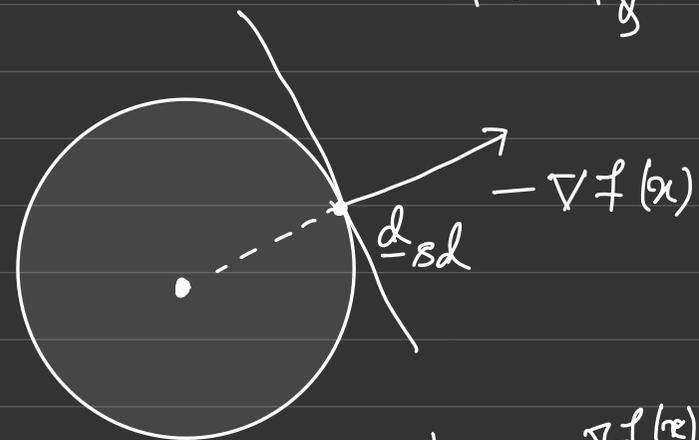
# Steepest descent method:

$\rightarrow$ descent direction that is as negative as possible
~~that~~ that yields the greatest rate of objective value
improvement

$$\underline{d}_{sd} = \underset{\underline{d}}{\arg \min} \left\{ \nabla f^{\top}(\underline{x}) \, \underline{d} \; : \; \| \underline{d} \| \leq 1 \right\}$$

$$\| \underline{d} \|_{g} \leq 1$$

## For Euclidean norm :

$$\underline{d}_{sd} = - \nabla f(\underline{x})$$

$-\nabla f(x)$

$\underline{d}_{sd}$

## For $l_1$ - norm :

$$\underline{d}_{sd} = \underset{\underline{d}}{\arg \min} \left\{ \nabla f^{\top}(\underline{x}) \, \underline{d} \; : \; \| \underline{d} \|_{1} \leq 1 \right\}$$

$-\nabla f(x)$

$\underline{d}_{sd}$

Let $i$ be any index s.t. $\| \nabla f(\underline{x}) \|_{\infty} = \left| [\nabla f(x)]_{i} \right|$

Then $\underline{d}_{sd} = - \text{sign} \left( \dfrac{\partial f(x)}{\partial x_i} \right) \underline{e}_i \; : \; \text{coordinate descent.}$

$$\nabla f^{\top}(x_k) \, d_k < 0$$



$$d_k = -\nabla f(x_k)$$

$$\nabla f^{\top}(x_k) \, d_k = -\|\nabla f(x_k)\| < 0$$

Gradient descent:

$$x_{k+1} = x_k - \eta_k \nabla f(x_k)$$

Vanilla analysis:
$$\underline{x}_{t+1} = \underline{x}_t - \eta \, \nabla f(\underline{x}_t)$$
$$t = 0, \ldots, T-1$$

Define:
$$\underline{g}_t = \nabla f(\underline{x}_t)$$

So
$$\underline{g}_t = (x_{t+1} - x_t)/\eta$$

Let us relate $\underline{x}_t$ vector to our current direction from an optimum $\underline{x}^*$ : $x_t - \underline{x}^*$

$$\underline{g}_t^T (\underline{x}_t - \underline{x}^*) = \frac{1}{\eta} (\underline{x}_{t+1} - x_t)^T (\underline{x}_t - x^*)$$

Cosine theorem:
$$2 \underline{v}^T \underline{w} = \|\underline{v}\|^2 + \|\underline{w}\|^2 - \|\underline{v} - \underline{w}\|_2^2$$

$$\underline{g}_t^T (\underline{x}_t - \underline{x}^*) = \frac{1}{2\eta} \left[ \|\underline{x}_{t+1} - \underline{x}_t\|_2^2 + \|\underline{x}_t - \underline{x}^*\|_2^2 - \|\underline{x}_{t+1} - x^*\|_2^2 \right]$$

(✱)

$$= \frac{1}{2\eta} \left[ \eta^2 \|\underline{g}_t\|_2^2 + \|\underline{x}_t - x^*\|_2^2 - \|\underline{x}_{t+1} - \underline{x}^*\|_2^2 \right]$$

$$= \frac{\eta}{2} \|\underline{g}_t\|^2 + \frac{1}{2\eta} \left[ \|\underline{x}_t - x^*\|_2^2 - \|\underline{x}_{t+1} - \underline{x}^*\|^2 \right]$$

(✱) Sum over the iteration $t$.

$$\sum_{t=0}^{T-1} \underline{g}_t^T (\underline{x}_t - \underline{x}^*) = \frac{\eta}{2} \sum_{t=0}^{T-1} \| \underline{g}_t \|^2$$

$$+ \frac{1}{2\eta} \left[ \| \underline{x}_0 - \underline{x}^* \|_2^2 - \| \underline{x}_T - \underline{x}^* \|_2^2 \right]$$

$5 > 3$

$-5 < -3$

$$\leq \frac{\eta}{2} \sum_{t=0}^{T-1} \| \underline{g}_t \|^2 + \frac{1}{2\eta} \| \underline{x}_0 - \underline{x}^* \|_2^2$$

Now, suppose $f$ is convex : $f(\underline{y}) \geq f(x) + \nabla f^T(x)(\underline{y} - x)$

$\left. \begin{array}{c} y = \underline{x}^* \\ x = \underline{x}_t \end{array} \right\}$ $\Rightarrow$ $f(x_t) - f(x^*) \leq \nabla f^T(\underline{x}_t)(\underline{x}_t - x^*)$

$$= \underline{g}_t^T (\underline{x}_t - \underline{x}^*)$$

• Upper bound on the average error: $f(x_t) - f(x^*)$

$$\sum_{t=0}^{T-1} f(x_t) - f(x^*) \leq \frac{\eta}{2} \sum \|\underline{g}_t\|^2 + \frac{1}{2\eta} \|\underline{x}_0 - \underline{x}^*\|_2^2$$

• Last iterate is not neccesarily the best one as "fixed" step size can make steps overshoot & increase function value

$f:$ $X \longrightarrow \mathbb{R}^n$ is called Lipschitz
continuous if there exists $B \geqslant 0$

$$\|f(x) - f(y)\| \leq B \|x - y\|$$
$$\forall \ x, y \in X$$

$$\iff \quad \|\nabla f(x)\| \leq B \qquad \forall \ x \in X$$

$$f: \mathbb{R}^{n^{\sigma}} \rightarrow \mathbb{R}$$

• $$\|x_0 - x^*\| \leq R$$

Then,

$$\sum_{t=0}^{T-1} \left( f(x_t) - f(x^*) \right) \leq \frac{\eta}{2} B^2 T + \frac{1}{2\eta} R^2 \qquad \text{(**)}$$

So, choose $\eta$ such that $q(\eta) = \frac{\eta}{2} B^2 T + \frac{1}{2\eta} R^2$ is minimized.

$$\frac{d}{d\eta} q(\eta) = 0 \implies \frac{1}{2} B^2 T - \frac{1}{2\eta^2} R^2 = 0$$

$$\implies \quad \eta = \frac{R^2}{B^2 \sqrt{T}} \qquad \text{and} \quad q(\eta) = R B \sqrt{T}$$

(**) $\div T$

$$\frac{1}{T} \sum_{t=0}^{T-1} \left( f(x_t) - f(x^*) \right) \leq \frac{RB}{\sqrt{T}} \approx O\left(\frac{1}{\sqrt{T}}\right)$$

independent of $n$ but $R$ & $B$ does depend on $n$

So to obtain $\min_{t=0 \dots T-1} f(x_t) - f(x^*) \leq \varepsilon$

we need $T \geqslant \frac{R^2 B^2}{\varepsilon^2}$ · $\varepsilon o s.$ $\varepsilon = 10^{-8}$

$R$ and $B \approx O(\sqrt{9} ??)$

Second-order methods (Newton's):

$$f(x_1, x_2) = x_1^2 \cdot \frac{1}{100} + x_2^2 \cdot 100$$

$$x_k = \begin{bmatrix} -10 \\ -0.1 \end{bmatrix}$$

$$\nabla f(x_t) = \begin{bmatrix} -\frac{1}{5} \\ -20 \end{bmatrix}$$

Check Newton's method with $B = H^{-1}$?
(Hessian)

$$x_{k+1} = x_k - \eta_k H_k^{-1} \nabla f(x_k)$$

# Linear regression:

$$f(\underline{w}) = \sum_{n=1}^{N} \left( y_n - \underline{w}^T \underline{x}_n \right)^2$$

$$\nabla f(\underline{w}) = \sum_{n=1}^{N} \nabla f_n(\underline{w})$$

$$= - \sum_{n=1}^{N} 2\underline{x}_n \left( y_n - \underline{w}^T \underline{x}_n \right)$$

$\uparrow$

Can we update using
one or few examples
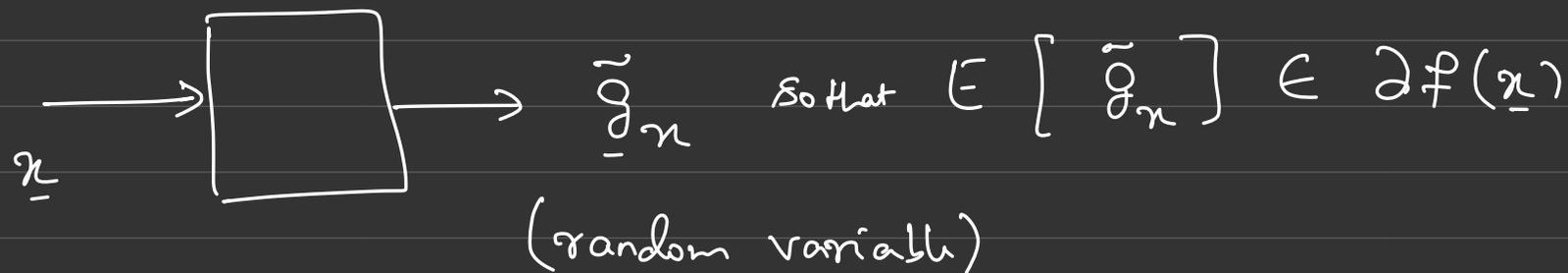in each iteration?

# Stochastic gradient descent:

$$x \longrightarrow \boxed{\phantom{xxxx}} \longrightarrow \widetilde{\nabla} f(x) = \begin{bmatrix} 0 \\ \vdots \\ [\nabla f(x)]_j \\ \vdots \\ 0 \end{bmatrix}$$

$$\underbrace{j \sim \text{unib}(1,\dots d) \quad \text{r. v.}}$$

$$\underline{x}_{k+1} = \underline{x}_{k+1} - \eta_k \, \widetilde{\nabla} f(x_k)$$

$$\widetilde{\nabla} f(x_k) \approx \nabla f(x_k) + \omega$$

Stochastic gradient:

$$\boxed{\mathbb{E}\left[ \widetilde{\nabla} f(x_k) \right] = \nabla f(x_k)}$$

$$\tilde{g}_n \quad \text{so that} \quad E\left[\tilde{g}_n\right] \in \partial f(\underline{x})$$

(random variable)

## Example:

① 
$$\tilde{g}_n = \nabla f(x) + \underline{w} \quad ; \quad \underline{w} \text{ is zero-mean noise}$$

$$E\left[\tilde{g}_n\right] = E\left[\nabla f(x) + \underline{w}\right] = \nabla f(\underline{x})$$

② Random coordinate descent:

$$x \longrightarrow \boxed{\phantom{xx}} \longrightarrow \tilde{g}_n = \begin{bmatrix} 0 \\ \vdots \\ \partial f/\partial x_j \\ \vdots \\ 0 \end{bmatrix} \cdot d \quad ; \quad \underline{x} \in \mathbb{R}^d$$

$$j \sim \text{Unif}(1, \ldots, d)$$

$$\nabla f(\underline{x}) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_d} \end{bmatrix}$$

$$E_j\left[\tilde{g}_n\right] = \sum_{i=1}^{d} d \cdot \begin{bmatrix} 0 \\ \vdots \\ \partial f/\partial x_j \\ \vdots \\ 0 \end{bmatrix} \cdot \frac{1}{d} = \nabla f(\underline{x})$$

# Unbiasedness and the vanilla analysis:

Recall: In gradient descent, we could lower bound

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

but now we cannot as $\tilde{g}_t$ may be far from being the true gradient.

- So inequality

$$f(x_t) - f(x^*) \leq \tilde{g}_t^T (x - x^*)$$

(from convexity) may not hold.

We have

$$\mathbb{E} \left[ \underbrace{g_t \mid x_t = x} \right] = \frac{1}{n} \sum_{j=1}^{n} \nabla f_j(x) = \nabla f(x) \qquad \forall x \in \mathbb{R}^d$$

Conditional expectation of $g_t$ given the event $\{x = x_t\}$.

$$\Rightarrow \mathbb{E}\left[g_t^T(\underline{x} - \underline{x}^*)\,\middle|\, \underline{x}_t = \underline{x}\right] =$$

$$\mathbb{E}\left[g_t^T\,\middle|\, \underline{x}_t = \underline{x}\right](\underline{x} - \underline{x}^*) = \nabla f^T(\underline{x})(\underline{x} - \underline{x}^*)$$

- $\{\underline{x}_t = \underline{x}\}$ can occur only for $\underline{x}$ in finite set $X$

$$\mathbb{E}\left[g_t^T(\underline{x}_t - x^*)\right]$$

$$= \sum_{\underline{x} \in X} \mathbb{E}\left[g_t^T(\underline{x} - \underline{x}^*)\,\middle|\, \underline{x}_t = \underline{x}\right]\text{prob}(\underline{x}_t = \underline{x})$$

$$= \sum_{\underline{x} \in X} \nabla f^T(\underline{x})(\underline{x} - \underline{x}^*)\,\text{prob}(\underline{x}_t = \underline{x})$$

$$= \mathbb{E}\left[\nabla f^T(\underline{x}_t)(\underline{x}_t - \underline{x}^*)\right]$$

$$\Rightarrow \quad \mathbb{E}\left[ \tilde{g}_t^\top (x_t - x^*) \right] = \mathbb{E}\left[ \nabla f^\top (x_t)(x_t - x^*) \right]$$

$$\geq \mathbb{E}\left[ f(x_t) - f(x^*) \right]$$

So the lower bound holds in expectation.

# Bounded stochastic gradients:

- Same convergence rate as gradient descent method

**Claim:** Let $f : \mathbb{R}^d \to \mathbb{R}$ be a convex and differentiable function, $\underline{x}^*$ be a global minimum of $f$;

$$\| \underline{x}_0 - \underline{x}^* \| \leq R \quad \text{and that} \quad \boxed{\mathbb{E}\left[ \| \underline{g}_t \|^2 \right] \leq B^2} \quad \forall t$$

Then Stochastic gradient descent with

constant step size $\eta = \dfrac{R}{B\sqrt{T}}$ yields

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}\left[ f(\underline{x}_t) \right] - f(\underline{x}^*) \leq \frac{RB}{\sqrt{T}}$$

Iteration complexity : $O\left( \dfrac{1}{\varepsilon^2} \right)$ $\qquad O\left( \dfrac{1}{\sqrt{T}} \right)$

Recall our vanilla analysis:

$$\underline{g}_t^T (\underline{x}_t - \underline{x}^*) = \frac{\eta}{2} \| g_t \|^2 + \frac{1}{2\eta} \left( \| \underline{x}_t - \underline{x}^* \|^2 - \| \underline{x}_{t+1} - \underline{x}^* \|^2 \right)$$

Telescoping sum :

$$\sum_{t=0}^{T-1} \underline{g}_t^T (\underline{x}_t - \underline{x}^*) = \frac{\eta}{2} \sum_{t=0}^{T-1} \| \underline{g}_t \|^2 + \frac{1}{2\eta} \left( \| \underline{x}_0 - \underline{x}^* \|^2 - \| \underline{x}_T - \underline{x}^* \|^2 \right)$$

$$\leq \frac{\eta}{2} \sum_{t=0}^{T-1} \| \underline{g}_t \|^2 + \frac{1}{2\eta} \| \underline{x}_0 - \underline{x}^* \|^2$$

Taking expectation on both sides

$$\sum_{t=0}^{T-1} \mathbb{E}\left[ \tilde{\underline{g}}_t^T (\underline{x}_t - \underline{x}^*) \right] \leq \frac{\eta}{2} \sum_{t=0}^{T-1} \underbrace{\mathbb{E}\left[ \| \tilde{\underline{g}}_t \|^2 \right]}_{\leq B^2} + \frac{1}{2\eta} \underbrace{\| \underline{x}_0 - \underline{x}^* \|^2}_{\leq R^2}$$

We have the lower bound :

$$\mathbb{E}\left[ \tilde{\underline{g}}_t^T (\underline{x}_t - \underline{x}^*) \right] \geq \mathbb{E}\left[ f(\underline{x}_t) - f(\tilde{\underline{x}}^*) \right]$$

$$\sum_{t=0}^{T-1} \mathbb{E}\left[f(x_t) - f(x^*)\right] \leq \frac{\eta}{2} B^2 T + \frac{1}{2\eta} R^2$$

$$= q(\eta)$$

Choose $\eta$ that minimize the upper bound:

$$\frac{1}{2} B^2 T - \frac{1}{2\eta^2} R^2 = 0$$

$$\eta = \frac{R}{B\sqrt{T}}$$

for which we have $O\left(\frac{1}{\sqrt{T}}\right)$