

# Newton Method

Second Order Approximation:

$\stackrel{\text{def}}{=} \phi(P_k)$

$$f(x_k + \alpha P_k) \cong f(x_k) + \alpha \nabla f(x_k)^T P_k + \frac{\alpha^2}{2} P_k^T \nabla^2 f(x_k) P_k$$

To find the "best" search direction, we want to find a  $P_k$  s.t.  $f(x_k + \alpha P_k)$  is minimized.

For simplicity, let's assume  $\alpha=1$ . Then

$$\frac{\partial}{\partial P_k} \phi(P_k) = \nabla f(x_k) + \nabla^2 f(x_k) P_k = 0$$

$$\Rightarrow \boxed{P_k = -[\nabla^2 f(x_k)]^{-1} \nabla f(x_k)} \quad (1)$$

Assumption:  $\nabla^2 f(x_k) > 0$ .

When will (1) be valid? (1) is valid if  $\phi(P_k)$  is a good approximation of  $f(x_k + \alpha P_k)$ , which happens only when  $\alpha P_k$  is sufficiently close to  $x_k$ .

if  $P_k = -[\nabla^2 f(x_k)]^{-1} \nabla f(x_k)$ , then one can show that

$$\begin{aligned} \nabla f(x_k)^T P_k &= -\nabla f(x_k)^T [\nabla^2 f(x_k)]^{-1} \nabla f(x_k) \\ &< 0 \quad , \quad = 0 \text{ when } \nabla f(x_k) = 0. \end{aligned}$$

Step size  $\alpha$ : We can actually fix  $\alpha=1$ .

gradient descent:  $\min_P \nabla f(x_k)^T P$  s.t.  $\|P\|_2 = 1$    
or using line search   
makes sure minimization has solution but loose length

Newton:  $\min_P \nabla f(x_k)^T P + \left(\frac{1}{2} P^T \nabla^2 f(x_k) P\right)$    
 $\nabla^2 f(x_k)$  is defining the length

## Compare Newton & descent

$$P_k = \underbrace{-[\nabla^2 f(x_k)]^{-1}}_{\text{vector step size}} \nabla f(x_k) \quad P_k = -\underbrace{[\alpha_k I]}_{\text{scalar step size}} \nabla f(x_k)$$

## Convergence of Newton

Assumption: (1)  $f$  is twice differentiable so that  $\nabla^2 f$  exists.

(2)  $\nabla^2 f(x)$  is Lipschitz continuous in a neighborhood of  $x^*$ :

$$\|\nabla^2 f(x) - \nabla^2 f(y)\| \leq L \|x - y\| \quad \forall x, y \in \text{this neighborhood}$$

## Theorem (3.7 Nocedal & Wright)

(1) if  $x_0$  is sufficiently close to  $x^*$ , then  $\{x_k\}$  converges to  $x^*$

(2) rate of convergence of  $\{x_k\}$  is quadratic:

$$\|x_{k+1} - x^*\| \leq C \|x_k - x^*\|^2, \text{ some constant } C$$

(3) rate of convergence of  $\{\nabla f(x_k)\}$  is quadratic:

$$\|\nabla f(x_{k+1}) - 0\| \leq C' \|\nabla f(x_k) - 0\|^2, \text{ some constant } C'$$

Proof: First, note that

$$\begin{aligned} x_{k+1} - x^* &= x_k + p_k - x^* \\ &= x_k - [\nabla^2 f(x_k)]^{-1} \nabla f(x_k) - x^* \\ &= [\nabla^2 f(x_k)]^{-1} \left[ \nabla^2 f(x_k) (x_k - x^*) - \nabla f(x_k) \right] \\ &= [\nabla^2 f(x_k)]^{-1} \left[ \nabla^2 f(x_k) (x_k - x^*) - (\nabla f(x_k) - \nabla f(x^*)) \right] \\ &= 0 \end{aligned}$$

Mean-Value theorem:

if  $f$  is continuous on  $[a, b]$ , then  $\exists c$  s.t.  $a \leq c \leq b$ ,

$$f(c)(b-a) = \int_a^b f(x) dx$$

Multi-variate: (see Wikipedia)

$$f(x+p) - f(x) = \int_0^1 \nabla f(x+tp)^T p dt$$

Fund. theorem of Calculus:

$$f(x+p) - f(x) = \int_x^{x+p} f'(u) du = \int_0^1 f'(x+tp) dt p$$

Therefore,

$$\begin{aligned} \nabla f(x_k) - \nabla f(x^*) &= \int_0^1 \nabla^2 f(x_k + t(x^* - x_k)) dt (x_k - x^*) \\ &= \int_0^1 \nabla^2 f(x_k + t(x^* - x_k)) \cdot (x_k - x^*) dt. \end{aligned}$$

So,

$$\begin{aligned} &\| \nabla^2 f(x_k)(x_k - x^*) - (\nabla f(x_k) - \nabla f(x^*)) \| \\ &= \| \nabla^2 f(x_k)(x_k - x^*) - \int_0^1 \nabla^2 f(x_k + t(x^* - x_k)) (x_k - x^*) dt \| \\ &= \| \int_0^1 [\nabla^2 f(x_k) - \nabla^2 f(x_k + t(x^* - x_k))] (x_k - x^*) dt \| \\ &\leq \int_0^1 \| \nabla^2 f(x_k) - \nabla^2 f(x_k + t(x^* - x_k)) \| \| x_k - x^* \| dt \\ &\leq \int_0^1 L \| x_k - (x_k + t(x^* - x_k)) \| \| x_k - x^* \| dt \\ &= \frac{1}{2} \| x_k - x^* \|^2 \int_0^1 L t dt = \frac{L}{2} \| x_k - x^* \|^2 \end{aligned}$$

### Banach Perturbation Lemma:

$A$ : non-singular,  $\|A^{-1}\| \leq \mu$ . If  $\|A - \bar{A}\| \leq \epsilon$ , where  $\epsilon\mu < 1$ , then  $\bar{A}$  is non-singular and

$$\|\bar{A}^{-1}\| \leq \frac{\mu}{1 - \mu\epsilon} \quad \epsilon = \frac{1}{2\mu}$$

Corollary:  $\nabla^2 f(x^*) > 0$ . Let  $\mu = \|\nabla^2 f(x^*)^{-1}\|$ . Then there is a  $B(x^*, \delta)$  s.t.  $\forall x \in B(x^*, \delta)$ ,

$$\|\nabla^2 f(x)^{-1}\| \leq 2\mu = 2\|\nabla^2 f(x^*)^{-1}\|.$$

Then,

$$\|x_{k+1} - x^*\| \leq \underbrace{2\|\nabla^2 f(x^*)^{-1}\|}_{\text{this shows (a) and (b)}} \|x_k - x^*\|^2.$$

So  $\{x_k\}$  converges  $\subset$  quadratically to  $x^*$ .

For (c),

$$\begin{aligned} \|\nabla f(x_{k+1}) - 0\| &= \|\nabla f(x_{k+1})\| \quad \xrightarrow{=0 \text{ by Newton method definition}} \\ &= \|\nabla f(x_{k+1}) - (\nabla f(x_k) + \nabla^2 f(x_k) P_k)\| \\ &= \|\int_0^1 \nabla^2 f(x_k + tP_k)(x_{k+1} - x_k) dt - \nabla^2 f(x_k) P_k\| \\ &= \|\int_0^1 [\nabla^2 f(x_k + tP_k) - \nabla^2 f(x_k)] P_k dt\| \\ &\leq \int_0^1 \|\nabla^2 f(x_k + tP_k) - \nabla^2 f(x_k)\| \|P_k\| dt \\ &\leq \frac{1}{2} L \|P_k\|^2 \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{2} L \left\| \left[ \nabla^2 f(x_k) \right]^{-1} \nabla f(x_k) \right\|^2 \\
&\leq \frac{1}{2} L \left\| \nabla^2 f(x_k)^{-1} \right\|^2 \left\| \nabla f(x_k) \right\|^2 \\
&\leq \frac{1}{2} L \cdot 4 \left\| \nabla^2 f(x^*)^{-1} \right\|^2 \left\| \nabla f(x_k) \right\|^2 \\
&= 2L \left\| \nabla^2 f(x^*)^{-1} \right\|^2 \left\| \nabla f(x_k) \right\|^2.
\end{aligned}$$

So we conclude that

$$\left\| \nabla f(x_{k+1}) \right\| \leq \underbrace{2L \left\| \nabla^2 f(x^*)^{-1} \right\|^2}_{C'} \left\| \nabla f(x_k) \right\|^2.$$

## Quasi-Newton Method

Recall Mean-Value Theorem:

$$\begin{aligned}
f(x+p) - f(x) &= \int_0^1 \nabla f(x+tp)^T p \, dt \\
\Rightarrow \nabla f(x+p) - \nabla f(x) &= \int_0^1 \nabla^2 f(x+tp) p \, dt \\
\Rightarrow \nabla f(x_{k+1}) - \nabla f(x_k) &= \int_0^1 \nabla^2 f(x_{k+1}) p_k \, dt \\
\Rightarrow \nabla f(x_{k+1}) - \nabla f(x_k) - \nabla^2 f(x_k) p_k &= \int_0^1 \left[ \nabla^2 f(x_{k+1}) - \nabla^2 f(x_k) \right] p_k \, dt \\
&\leq \frac{L}{2} \|p_k\|^2 \\
&= \frac{L}{2} \|x_{k+1} - x_k\|^2.
\end{aligned}$$

So if  $x_{k+1} \approx x^*$  and  $x_k \approx x^*$  (near  $x^*$ ), then

$$\nabla f(x_{k+1}) \approx \nabla f(x_k) + \nabla^2 f(x_k) p_k$$

or, in other words,

$$\underbrace{\nabla f(x_{k+1}) - \nabla f(x_k)}_{\text{def } y_k} \cong \underbrace{\nabla^2 f(x_k)}_{\text{def } B_{k+1}} \underbrace{(x_{k+1} - x_k)}_{\text{def } s_k}$$

Therefore, we can replace Newton Equation

$$P_k = -[\nabla^2 f(x_k)]^{-1} \nabla f(x_k) \longrightarrow \boxed{P_k = -[B_{k+1}]^{-1} \nabla f(x_k)}$$

Update Rules for  $B_{k+1}$ :

(Because we want to find cheap ways to compute  $B_{k+1}$ )

(i) Symmetric Rank-1 Update

$$B_{k+1} = B_k + \frac{(y_k - B_k s_k)(y_k - B_k s_k)^T}{(y_k - B_k s_k)^T s_k}$$

(ii) BFGS:

$$B_{k+1} = B_k - \frac{B_k s_k s_k^T B_k}{s_k^T B_k s_k} + \frac{y_k y_k^T}{y_k^T s_k}$$

Why these updates?

- Maintain Symmetry
- Prefer low-rank update
- satisfies  $B_{k+1} s_k = y_k$ .

## Barzilai-Borwein Method

- Steepest descent with a special step size
- Motivated from Newton but does not require Hessian.

Recall Newton method:

$$\nabla^2 f(x_k) p_k = -\nabla f(x_k)$$

$$\Rightarrow p_k = -[\nabla^2 f(x_k)]^{-1} \nabla f(x_k)$$

Idea: Replace  $-\left[\nabla^2 f(x_k)\right]^{-1}$  by  $-\alpha_k I$ .

So we want

$$\nabla^2 f(x_k) \leftarrow \alpha_k^{-1} I.$$

Let's also recall Quasi-Newton Method:

$$B_{k+1} s_k = y_k \quad \begin{cases} s_k = x_{k+1} - x_k \\ y_k = \nabla f(x_{k+1}) - \nabla f(x_k) \end{cases}$$

So, we can choose  $\alpha_k$  s.t.

$$\boxed{-\left(\alpha_k^{-1} I\right) s_{k-1} = y_{k-1}.}$$

How to choose  $\alpha_k$ ?

$$(i) \quad \alpha_k^{-1} = \underset{\beta}{\operatorname{argmin}} \left\| s_{k-1} \beta - y_{k-1} \right\|^2$$

$$\Rightarrow \alpha_k = \frac{s_{k-1}^T s_{k-1}}{s_{k-1}^T y_{k-1}}$$

$$(ii) \quad \alpha_k = \underset{\alpha}{\operatorname{argmin}} \left\| s_{k-1} - y_{k-1} \alpha \right\|^2$$

$$\Rightarrow \alpha_k = \frac{s_{k-1}^T y_{k-1}}{y_{k-1}^T y_{k-1}}$$

## Gradient Projection Method

Recall steepest descent:

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

Claim:

$$x_{k+1} = \operatorname{argmin}_x \left\{ f(x_k) + \nabla f(x_k)^T (x - x_k) + \frac{1}{2\alpha_k} \|x - x_k\|^2 \right\}$$

Proof: Take derivative:

$$\frac{\partial}{\partial x}(\cdot) = \nabla f(x_k) + \frac{1}{\alpha_k} (x - x_k) = 0$$

$$\Rightarrow x = x_k - \alpha_k \nabla f(x_k).$$

Implication:

Gradient descent is a minimizer of a local linear model added with a quadratic penalty. (or think as we set  $\nabla^2 f(x_k) \cong \frac{1}{\alpha_k} I$ )

Gradient Projection:

if we assume that  $x \in \Omega$ , say  $l \leq x \leq u$ , then we can modify the min. as

$$x_{k+1} = \operatorname{argmin}_{x \in \Omega} \left\{ f(x_k) + \nabla f(x_k)^T (x - x_k) + \frac{1}{2\alpha_k} \|x - x_k\|^2 \right\}.$$

$$\Rightarrow \boxed{x_{k+1} = P_{\Omega}(x_k - \alpha_k \nabla f(x_k))}$$

where  $P_{\Omega}$  is a projection operator that projects its input to  $\Omega$ .



Example: GPSR algorithm (Figueiredo, Nowak, Wright)

GPSR tries to solve this problem

$$\min_x \frac{1}{2} \|y - Ax\|^2 + \lambda \|x\|_1$$

Its idea is to turn the problem into

$$\min_{u,v} \frac{1}{2} \|y - A(u-v)\|^2 + \lambda (\mathbf{1}^T u + \mathbf{1}^T v)$$

$$\text{s.t. } u \geq 0, v \geq 0.$$

$$\Leftrightarrow \boxed{\begin{array}{l} \min_z f(z) = \frac{1}{2} z^T B z + c^T z \\ \text{s.t. } z \geq 0 \end{array}}$$

$$\text{where } z = \begin{bmatrix} u \\ v \end{bmatrix}, b = A^T y, c = \lambda \mathbf{1} + \begin{bmatrix} -b \\ b \end{bmatrix},$$

$$B = \begin{bmatrix} A^T A & -A^T A \\ -A^T A & A^T A \end{bmatrix}.$$

• The gradient projection is the following:

$$z_{k+1} = \left[ z_k - \alpha_k \nabla f(z_k) \right]_+,$$

$$\text{where } \nabla f(z_k) = B z_k + c,$$

and  $[\cdot]_+$  returns the positive entries only.

$$\text{e.g. } [x_i]_+ = \begin{cases} x_i & , \text{ if } x_i \geq 0 \\ 0 & , \text{ if } x_i < 0 \end{cases}.$$

The step size  $\alpha_k$  is found using the line search.

• Variation: (Barzilai-Borwein)

$$z_{k+1} = z_k - \lambda_k P_k, \quad \lambda_k \text{ found by line search}$$

$$\text{where } P_k = z_k - z_{k-1} = \left[ z_{k-1} - \alpha_k \nabla f(z_{k-1}) \right]_+ - z_{k-1}, \quad \alpha_k = \frac{s_k^T s_k}{s_k^T B_k s_k} \text{ (BB step)}$$

## Proximal Gradient Methods

The proximal mapping of a convex function  $h$  is defined as

$$\text{prox}_h(x) = \underset{v}{\text{argmin}} \left\{ h(v) + \frac{1}{2} \|v - x\|^2 \right\}$$

For  $\lambda h(x)$ , it is easy to show that

$$\text{prox}_{\lambda h}(x) = \underset{v}{\text{argmin}} \left\{ h(v) + \frac{1}{2\lambda} \|v - x\|^2 \right\}.$$

### Examples

(i)  $h(x) = 0$ . Then  $\text{prox}_h(x) = x$

(ii)  $h(x) = \begin{cases} 0, & \text{if } x \in \Omega \\ +\infty, & \text{if } x \notin \Omega. \end{cases}$

then

$$\begin{aligned} \text{prox}_h(x) &= \underset{v \in \Omega}{\text{argmin}} \frac{1}{2} \|v - x\|^2 \\ &= P_\Omega(x). \end{aligned}$$

if  $\Omega = \{x \mid Ax = b\}$   
and if  $A$  full rank,  
then,

$$P_\Omega(x) = x - A^\dagger (Ax - b)$$

Pseudo-inverse

(iii)  $h(x) = \|x\|_1$ . Then

$$\text{prox}_h(x) = \underset{v}{\text{argmin}} \left\{ \|x\|_1 + \frac{1}{2} \|v - x\|^2 \right\}$$

in general,

$$\begin{aligned} \text{prox}_{\lambda h}(x) &= \underset{v}{\text{argmin}} \left\{ \|x\|_1 + \frac{1}{2\lambda} \|v - x\|^2 \right\} \\ &= S_\lambda(x) = \max\{|x| - \lambda, 0\} \text{sign}(x). \end{aligned}$$

(iv)  $h(x) = -\log x$  (appears for Poisson regression) ↙ scalar  $x$

$$\text{prox}_{\lambda h}(x) = \underset{v}{\text{argmin}} \left\{ -\log v + \frac{1}{2\lambda}(v-x)^2 \right\}$$

$$\frac{\partial}{\partial v}(\cdot) = \frac{-1}{v} + \frac{1}{\lambda}(v-x) = 0$$

$$\Rightarrow -1 + v^2 - vx = 0$$

$$\Rightarrow v = \frac{x + \sqrt{x^2 + 4\lambda}}{2}$$

This technique applies to vector case if

$$h(x) = -\sum_{i=1}^n \log x_i$$

$$(v) \quad h(x) = \begin{cases} 0 & , \quad x \in \Omega \\ +\infty & , \quad x \notin \Omega \end{cases}$$

$$\Omega = \left\{ x \mid x \geq 0, \quad x^T \mathbf{1} = 1 \right\}$$

happens for  
probability-type  
of problems

$$\text{prox}_{\lambda h}(x) = \underset{v}{\text{argmin}} \left\{ \frac{1}{2} \|v - x\|^2 \right\}$$

$$\text{s.t. } x \geq 0, \quad x^T \mathbf{1} = 1.$$

$$= (x - \lambda \mathbf{1})_+, \quad \text{where } \lambda \text{ is the solution of}$$

$$\mathbf{1}^T (x - \lambda \mathbf{1}) = 1, \quad \text{or}$$

$$\text{i.e. } \sum_{i=1}^n \max\{x_i - \lambda, 0\} = 1.$$

(This is a constrained problem. We will discuss it after we study the KKT condition & Lagrange multiplier.)

Consider an unconstrained optimization with two functions in the objective

$$\min_x f(x) + h(x)$$

- $f(x)$  is convex and differentiable (eg.  $\|Ax-y\|^2$ )
- $h(x)$  is convex (eg.  $\lambda\|x\|_1$ ).

The proximal gradient method is

$$x_{k+1} = \text{prox}_{\alpha_k h} \left( x_k - \alpha_k \nabla f(x_k) \right)$$

↑  
step size (constant or line search)

What is it doing?

$$x_{k+1} = \text{prox}_{\alpha_k h} \left( x_k - \alpha_k \nabla f(x_k) \right)$$

$$= \arg\min_v \left\{ h(v) + \frac{1}{2\alpha_k} \|v - (x_k - \alpha_k \nabla f(x_k))\|^2 \right\}$$

$$= \arg\min_v \left\{ h(v) + \frac{1}{2\alpha_k} \left( \|v - x_k\|^2 + 2\alpha_k \nabla f(x_k)^T (v - x_k) + \alpha_k^2 \|\nabla f(x_k)\|^2 \right) \right\}$$

$$= \arg\min_v \left\{ h(v) + \cancel{2\alpha_k} \nabla f(x_k)^T (v - x_k) + \frac{1}{2\alpha_k} \|v - x_k\|^2 \right\}$$

$$= \arg\min_v \left\{ h(v) + \underbrace{f(x_k) + \nabla f(x_k)^T (v - x_k)}_{\text{local quadratic model around } x_k} + \frac{1}{2\alpha_k} \|v - x_k\|^2 \right\}$$

local quadratic model around  $x_k$

Example:

(i)  $h(x) = 0$ . Then

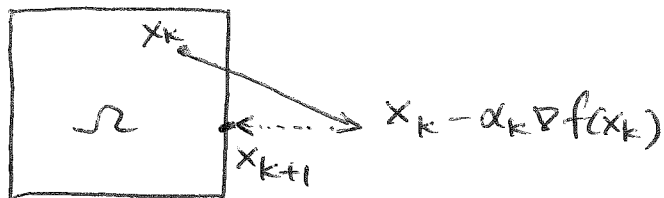
$$\begin{aligned}x_{k+1} &= \text{prox}_{\alpha_k h}(x_k - \alpha_k \nabla f(x_k)) \\ &= \underset{v}{\text{argmin}} \left\{ 0 - \frac{1}{2\alpha_k} \|v - (x_k - \alpha_k \nabla f(x_k))\|^2 \right\} \\ &= x_k - \alpha_k \nabla f(x_k) \quad \text{See previous example (i)}\end{aligned}$$

So the proximal mapping is the classical steepest descent.

(ii)  $h(x) = \begin{cases} 0 & , x \in \Omega \\ +\infty & , x \notin \Omega \end{cases}$ . Then

$$\begin{aligned}x_{k+1} &= \text{prox}_{\alpha_k h}(x_k - \alpha_k \nabla f(x_k)) \\ &= P_{\Omega}(x_k - \alpha_k \nabla f(x_k))\end{aligned}$$

So the proximal mapping is the project gradient method.



(iii)  $h(x) = \|x\|_1$ . Then

$$\begin{aligned}x_{k+1} &= \text{prox}_{\alpha_k h}(x_k - \alpha_k \nabla f(x_k)) \\ &= S_{\alpha_k}(x_k - \alpha_k \nabla f(x_k)) \\ &= \max \left\{ |x_k - \alpha_k \nabla f(x_k)| - \alpha_k, 0 \right\} \text{sign} \left\{ x_k - \alpha_k \nabla f(x_k) \right\}.\end{aligned}$$

## Application: ISTA

$\min_x f(x) + h(x)$ , where

$$\begin{cases} f(x) = \frac{1}{2} \|Ax - y\|^2 \\ h(x) = \lambda \|x\|_1. \end{cases}$$

The proximal gradient method defines

$$x_{k+1} = \text{prox}_{\alpha_k h}(x_k - \alpha_k \nabla f(x_k)).$$

$$\nabla f(x_k) = A^T(Ax_k - y)$$

$$\begin{aligned} \text{prox}_{\alpha_k h}(x) &= \underset{v}{\text{argmin}} \left\{ \lambda \|x\|_1 + \frac{1}{2\alpha_k} \|v - x\|^2 \right\} \\ &= \underset{v}{\text{argmin}} \left\{ \|x\|_1 + \frac{1}{2\lambda\alpha_k} \|v - x\|^2 \right\} \\ &= S_{\lambda\alpha_k}(x). \end{aligned}$$

## Convergence Rate of ISTA:

Assume  $f$  is differentiable & convex, and has Lipschitz continuous gradient:

$$\|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\| \quad \forall x, y \in \mathbb{R}^n.$$

e.g. if  $f(x) = \frac{1}{2} \|Ax - y\|^2$ , then

$$\nabla f(x) = A^T(Ax - y)$$

$$\begin{aligned} \|\nabla f(x) - \nabla f(x')\| &= \|A^T(Ax - y) - A^T(Ax' - y)\| \\ &= \|A^T A(x - x')\| \\ &\leq \underbrace{\lambda_{\max}(A^T A)}_L \|x - x'\| \end{aligned}$$

### Theorem (Beck & Teboulle)

Let  $\{x_k\}$  be a sequence generated by ISTA with either a constant step size or a line search, then for any  $k \geq 1$ , it holds that

$$F(x_k) - F(x^*) \leq \frac{\gamma L \|x_0 - x^*\|^2}{2k},$$

where  $\gamma = 1$  for constant step size,

$\gamma = \gamma_s$  for line search.

$F(x) \stackrel{\text{def}}{=} f(x) + h(x)$ , is the overall objective function.

If we want

$$F(x_k) - F(x^*) < \epsilon,$$

then Theorem tells us that we need

$$\frac{\gamma L \|x_0 - x^*\|^2}{2k} < \epsilon$$

$$\Rightarrow k > \frac{\gamma L \|x_0 - x^*\|^2}{2\epsilon} \leftarrow O\left(\frac{1}{\epsilon}\right).$$

Operations.