



262B-Lecture 10

Date created: 2021.02.18
N. of Pages: 12

Conjugate gradient method:

$$\min_{x \in \mathbb{R}^n} f(x) \rightarrow \frac{1}{2} x^T Q x - b^T x$$

$$\text{Algorithm: } x^{(k+1)} = x^{(k)} + \alpha^{(k)} d^{(k)} \quad k = 0, 1, \dots$$

Instead of $d^{(k)} = -\nabla f(x^{(k)})$, we generate

$d^{(0)}, d^{(1)}, \dots, d^{(n-1)}$ based on $-\nabla f(x^{(0)}), -\nabla f(x^{(1)})$

, ..., $-\nabla f(x^{(n-1)})$ in such a way that the spanning

rule on subspaces is preserved.

Formula:

$$\begin{cases} d^{(k)} = -\nabla f(x^{(k)}) + \frac{\sum_{i=0}^{k-1} (d^{(i)})^T Q \nabla f(x^{(k)})}{(d^{(i)})^T Q (d^{(i)})} d^{(i)} \\ d^{(0)} = -\nabla f(x^{(0)}) \end{cases}$$

Also, since $\alpha^{(k)}$ is optimal: $(d^{(i)})^T \nabla f(x^{(k+1)}) = 0$

$$i = 0, \dots, k$$

Also, span of $d^{(0)}, \dots, d^{(k)}$
= span of $\nabla f(x^{(0)}), \dots, \nabla f(x^{(k)})$

$\Rightarrow \nabla f(x^{(k+1)})$ is orthogonal to $\nabla f(x^{(0)}), \dots, \nabla f(x^{(k)})$

Simplify:

$$\left\{ \begin{array}{l} x^{(k+1)} = x^{(k)} + \alpha^{(k)} d^{(k)} \\ d^{(k)} = -\nabla f(x^{(k)}) + \frac{\nabla f(x^{(k)})^T \nabla f(x^{(k-1)})}{\nabla f(x^{(k-1)})^T \nabla f(x^{(k-1)})} d^{(k-1)} \\ d^{(0)} = -\nabla f(x^{(0)}) \end{array} \right.$$

And $\alpha^{(k)} : \min_{\alpha} f(x^{(k)} + \alpha d^{(k)})$

Thm: If $f(x)$ is quadratic, then conjugate gradient stops after at most n steps.

→ finite-time convergence

Note: algorithm can be used for arbitrary $f(x)$, but convergence is not finite.

$$\min \frac{1}{2} x^T Q x + b^T x$$

Conjugate gradient = Quasi-newton with exact line search & $D^{(0)} = I$
(true for arbitrary $f(x)$)

Convergence rate: $\min \frac{1}{2} x^T Q x$ (no b)
 λ_0

$$\Rightarrow \begin{cases} x^{(1)} \in x^{(0)} + \text{span} \{ \nabla f(x^{(0)}) \} \\ x^{(2)} \in x^{(0)} + \text{span} \{ \underbrace{\nabla f(x^{(0)})}_{Qx^{(0)}}, \underbrace{\nabla f(x^{(1)})}_{Qx^{(1)}} \} \end{cases}$$

$$\Rightarrow x^{(2)} \in x^{(0)} + \text{span} \{ Qx^{(0)}, Q^2x^{(0)} \}$$

⋮

$$x^{(k+1)} \in x^{(0)} + \text{span} \{ Qx^{(0)}, Q^2x^{(0)}, \dots, Q^{k+1}x^{(0)} \}$$

$$x^{(k+1)} = x^{(0)} + Q \left(\sum_{j=0}^k \beta_j Q^j \right) x^{(0)}$$

$P^k(Q)$: polynomial of

$$\Rightarrow x^{(k+1)} = \left(I + Q \underbrace{P^k(Q)}_{\text{degree } k} \right) x^{(0)}$$

previously: coefficients are optimal

$$\left(x^{(k+1)} = \arg \min_{x \in M^{(k)}} f(x) \right)$$

$$\Rightarrow \left(f(x^{(k+1)}) = \min_{P^k(\cdot)} \frac{1}{2} (x^{(0)})^T Q (I + Q P^k(Q))^{-2} x^{(0)} \right)$$

$\frac{1}{2} x^{(k+1)T} Q x^{(k+1)}$

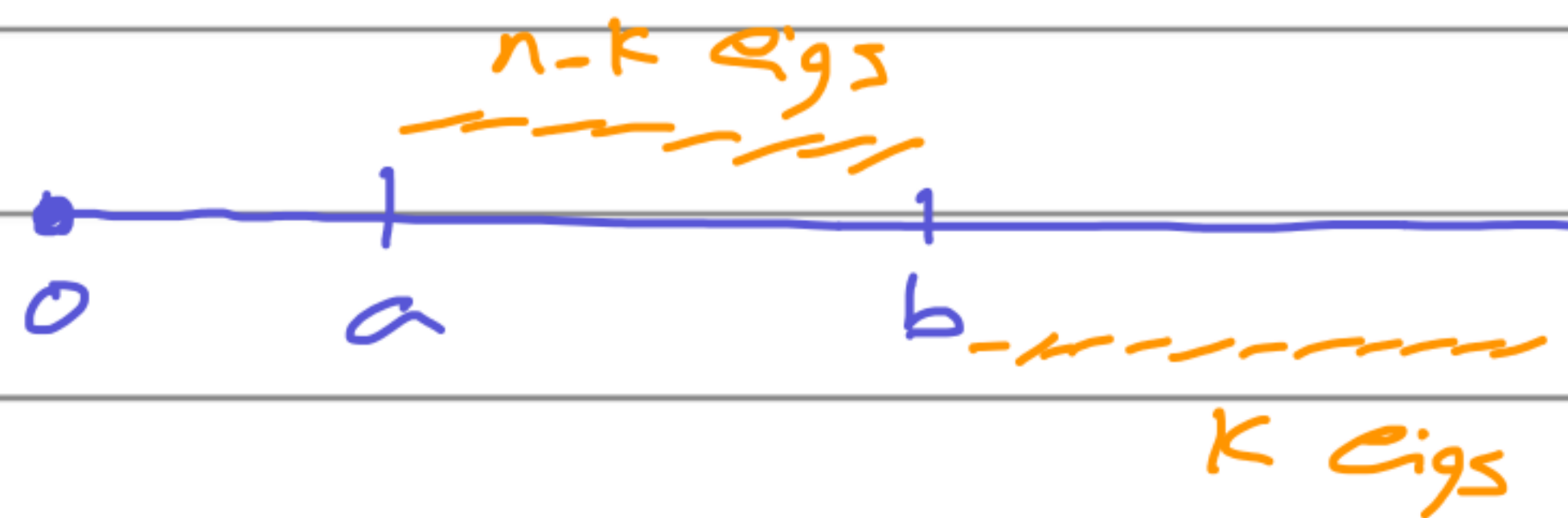
If we denote eigenvalues of Q with $\lambda_1, \dots, \lambda_n$ and eigenvectors as v_1, \dots, v_n , then

$$x^{(0)} = \underbrace{\mu_1}_{\in \mathbb{R}} v_1 + \underbrace{\mu_2}_{\in \mathbb{R}} v_2 + \dots + \underbrace{\mu_n}_{\in \mathbb{R}} v_n$$

$$\textcircled{*} \Rightarrow f(x^{(k+1)}) \leq \max_{1 \leq i \leq n} (1 + \lambda_i P^k(\lambda_i))^2 f(x^{(0)})$$

$\forall k, \forall P^k(\cdot)$

Thm: Assume Q has $n-k$ eigenvalues in the interval $[a, b]$ where $a > 0$ and k eigenvalues in the interval (b, ∞) .



After $k+1$ steps of conjugate gradient:

$$f(x^{(k+1)}) \leq \left(\frac{b-a}{b+a} \right)^2 f(x^{(0)})$$

Compare it with gradient method:
(with optimal stepsize)

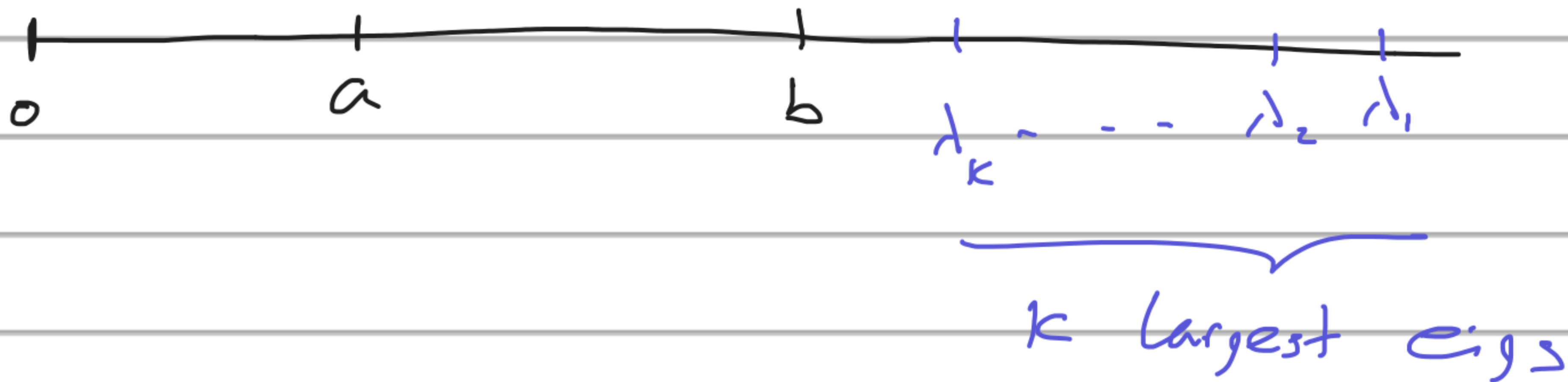
$$f(x^{(k+1)}) \leq \left(\frac{\underbrace{\text{c.d.}(Q) - 1}_{\lambda_{\max} - \lambda_{\min}}}{\text{c.d.}(Q) + 1}_{\lambda_{\max} + \lambda_{\min}} \right)^2 f(x^{(0)})$$

So, after $k+1$ iterations of conjugate gradient the k largest eigs drop and condition number changes from $\frac{\lambda_{\max}}{\lambda_{\min}}$ to $\frac{b}{a}$.

Ex: $\underbrace{\text{xxxx}}_{\sigma_{\text{round}} \perp} \quad \frac{x}{10^4} \Rightarrow \text{c.d.} = \text{huge}$
 \downarrow
gradient alg = slow

If you use conjugate gradient \Rightarrow c.d. after one iteration ≈ 1

proof:



pick $p^k(\cdot)$:

$$1 + \lambda p^k(\lambda) = \frac{2}{(a+b)\lambda_1 \dots \lambda_k} \left(\frac{a+b}{2} - \lambda \right) (\lambda_1 - \lambda) \dots (\lambda_k - \lambda)$$

degree = $k+1$

at $\lambda=0$: value = 1

satisfies both properties

$$\textcircled{1}: 1 + \lambda_i p^k(\lambda_i) = 0 \quad i=1, \dots, k$$

$$\textcircled{2}: \lambda \leq b \Rightarrow \frac{\lambda_i - \lambda}{\lambda_i} \leq 1 \quad i=1, \dots, k$$

$$\Rightarrow f(x^{(k+1)}) \leq \max_{1 \leq i \leq n} (1 + \lambda_i p^k(\lambda_i))^2 f(x^{(0)})$$

$$\textcircled{1}, \textcircled{2} \Rightarrow f(x^{(k+1)}) \leq \max_{a \leq \lambda \leq b} \left(\frac{\lambda - \frac{a+b}{2}}{\frac{a+b}{2}} \right)^2 f(x^{(0)})$$

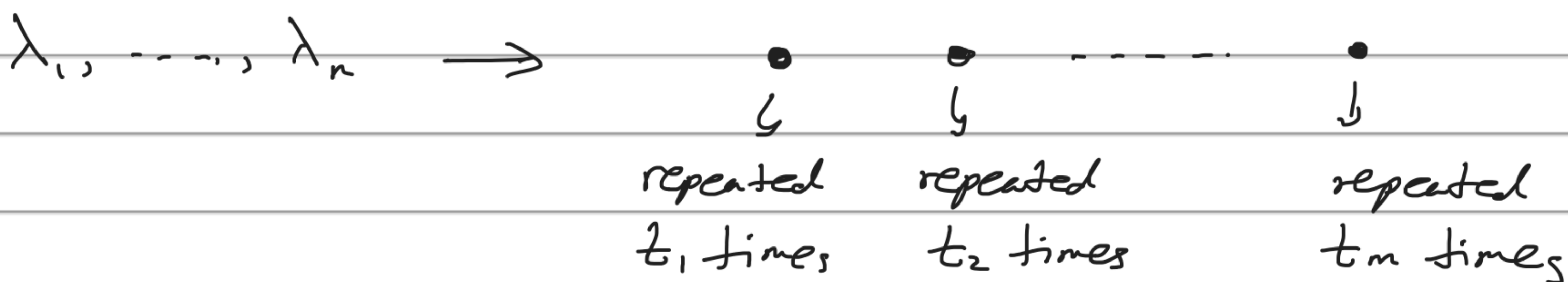
$$= \left(\frac{b-a}{b+a} \right)^2 f(x^{(0)})$$

Homework: If Q is low rank ($\text{rank} = m$)

\Rightarrow Conjugate gradient solves the problem

in at most m steps.

If Q has clustered eigs, $Q \succ 0$



$\rightarrow m$: number of clusters

Then, conjugate gradient solves the problem

in at most $m+1$ steps.

Newton's method but use conjugate gradient

to find Newton's step:

$$x^{(k+1)} = x^{(k)} + \alpha^{(k)} \Delta x^{(k)} \quad ; \quad \underbrace{\underbrace{\underbrace{\nabla f(x^{(k)})^2}_{\text{Hessian}}}_{\text{CG}}}_{\text{CG}} \Delta x^{(k)} = -\nabla f(x^{(k)})$$

Then instead of complexity of $O(n^3)$,

it could be much lower if $\nabla^2 f(x_*)$

is structured and $x^{(0)}$ is close to x_* .

Ex: $\nabla^2 f(x_*)$: ① ill-conditioned \rightarrow

conjugate gradient

gets rid of

bad eigs

② almost low-rank

③ eigs are almost clustered

⋮

Coordinate descent :

Goal : low-complexity algorithm

Here, instead of iterating over n entries,

what if we iterate over a single entry ?

$$x^{(k+1)} = x^{(k)} - \alpha^{(k)} \nabla f(x^{(k)})$$

instead, pick entry $i \in \{1, \dots, n\}$

$$\textcircled{1} \quad x^{(k+1)} = x^{(k)} - \alpha^{(k)} \frac{\partial f(x^{(k)})}{\partial x_i} \begin{pmatrix} \vdots \\ 1 \\ \vdots \end{pmatrix} \quad \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \quad \begin{matrix} \text{=} i\text{-th} \\ \text{coordi} \\ \text{rate} \end{matrix}$$

or

$$\textcircled{2} \quad x^{(k+1)} = \arg \min_{x_i \in \mathbb{R}} f(x_1^{(k)}, \dots, x_{i-1}^{(k)}, x_i, x_{i+1}^{(k)}, \dots, x_n^{(k)})$$

Known
variable
Known

Strategy $\textcircled{1}$: one improvement step for entry i

Strategy $\textcircled{2}$: univariate opt \Rightarrow low-complexity ||c||

later on, we will study convergence in a general setting subject to constraints.

Machine learning application:

$$\min_x g(x) + \underbrace{|x|}_{\text{regularizer}} \quad \longrightarrow \quad \min_{x_i} g(x_i) + |x_i|$$

non-smooth
1-d non-smooth \rightarrow closed-form

Convergence rate for constant stepsize:

Assumptions:

$$- \left| \frac{\partial f(x+te_i)}{\partial x_i} - \frac{\partial f(x)}{\partial x_i} \right| \leq L|t|$$

$$\forall x \in \mathbb{R}^n, \forall t \in \mathbb{R}$$

$$\forall i \in \{1, \dots, n\}$$

Lipschitz continuity of gradient

partial derivatives

$$- \nabla^2 f(x) \succeq \frac{m}{r} I \quad \forall x$$

similar result holds for non-convex case

$$- \text{pick stepsize } \alpha^{(k)} = \alpha < \frac{1}{L}$$

$$\text{Alg: } x^{(k+1)} = x^{(k)} - \alpha \frac{\partial f(x^{(k)})}{\partial x_i} e_i \rightarrow y$$

$$\begin{aligned} f(x^{(k+1)}) &\leq f(x^{(k)}) \\ &+ \nabla f(x^{(k)})^T \left(-\alpha \frac{\partial f(x^{(k)})}{\partial x_i} \right) e_i \\ &+ \frac{L}{2} \left\| -\alpha \frac{\partial f(x^{(k)})}{\partial x_i} \right\|^2 \end{aligned}$$

Thm:
 $f(x+y) \leq f(x) + \nabla f(x)^T y + \frac{L}{2} \|y\|^2 \rightarrow$ redo for partial derivatives

*

$$\textcircled{*} \Rightarrow f(x^{(k+1)}) \leq f(x^{(k)}) - \left(\frac{\partial f(x^{(k)})}{\partial x_i} \right)^2$$

$$\times \alpha \left(1 - \frac{\alpha L}{2} \right)$$

How to pick \underline{i} :

1-cyclic : 1, 2, ..., n, 1, 2, ..., n, 1, 2, ..., n, ...

2-random : pick a coordinate in a uniform way

call them $i_1, i_2, \dots, i_k, \dots$

let's analyze case $\underline{2}$:

Take expected value with respect to i_k from

$\textcircled{**}$:

$$\Rightarrow \mathbb{E}_{i_k} (f(x^{(k+1)}) - f_*) \leq (f(x^{(k)}) - f_*)$$

$$- \alpha \left(1 - \frac{\alpha L}{2} \right) \times \frac{1}{n} \sum_{i=1}^n \left(\frac{\partial f(x^{(k)})}{\partial x_i} \right)^2$$

$$2n (f(x^{(k)}) - f_*) \leq \left(\|\nabla f(x^{(k)})\| \right)^2$$

$$\Rightarrow \sum_{i_k} (f(x^{(k+1)}) - f_*) \leq$$

$$(f(x^{(k)}) - f_*) \times \left(1 - \frac{m}{n} \alpha (2 - \alpha L)\right)$$

Take expected value with respect to i_0, \dots, i_{k-1}

$$\Rightarrow (\phi_{k+1} - f_*) \leq (\phi_k - f_*) \underbrace{\left(1 - \frac{m}{n} \alpha (2 - \alpha L)\right)}$$

where $\phi_k = \mathbb{E} (f(x^{(k)}))$
 $\{i_0, i_1, \dots, i_{k-1}\}$

if α : small
 then this term
 is in $(0, 1)$

\Rightarrow Linear convergence on $\phi_k - f_*$

\Rightarrow Linear convergence on expected value
 of $f(x)$.