

## Unconstrained Minimization (Sec 9.1-9.4 + my SGD notes)

→ no dual

Aim: Solve the unconstrained min problem.

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

want: algorithm to find minimum of a point

$f: \mathbb{R}^n \rightarrow \mathbb{R}$  convex & twice continuously diff<sup>ble</sup>.

problem is solveable

Assume  $\exists x^* \in \mathbb{R}^n$  s.t.  $f(x^*) = \min_{x \in \mathbb{R}^n} f(x) = p^*$

there exist a minimizer

Necessary and sufficient condition for  $x^*$  to be optimal

System of  $n$  equations  $\nabla f(x^*) = 0$ . stationarity

⇒ Solving the unconstrained problem  $\Leftrightarrow$  solving the system of  $n$  equations in  $n$  unknowns.

↳ might be nonlinear etc.

iteration number  
 $x^{(0)}, x^{(1)}, x^{(2)}, \dots$

Iterative algorithms: Produces a sequence  $x^{(k)} \in \text{dom } f$

s.t.  $f(x^{(k)}) \rightarrow p^*$  as  $k \rightarrow \infty$ .

in practice we do not know what this is

Termination criterion: i)  $f(x^{(k)}) - p^* \leq \varepsilon$ .

ii)  $\|\nabla f(x^{(k)})\| \leq \varepsilon$ . (when gradient is sufficiently small)

Initial point  $x^{(0)} \in \text{dom } f$ .

Sublevel set  $S = \{x \in \mathbb{R}^n : f(x) \leq f(x^{(0)})\}$  is closed

quadratic      tve definite      vector

Example 1:  $\min_{x \in \mathbb{R}^n} \frac{1}{2} x^T P x + q^T x + r$      $P \in \mathbb{S}_+^n, q \in \mathbb{R}^n, r \in \mathbb{R}$ .

$$\nabla f(x^*) = 0 \Rightarrow P x^* + q = 0$$

Since  $P \in \mathbb{S}_+^n$ ,  $P^{-1}$  exists, invertible,  $x^* = -P^{-1}q$ . (no need for optimization algorithm  
↳ closed form exist)

Example 2: Unconstrained Geometric Programming → cannot

$$\min_{x \in \mathbb{R}^n} f(x) \quad f(x) = \log \left( \sum_{i=1}^m \exp(a_i^T x + b_i) \right)$$

write down  
closed form  
expression  
from  $x^*$

Stationarity condition:  $\nabla f(x^*) = 0$ .

$$\frac{\sum_{i=1}^m a_i \exp(a_i^T x + b_i)}{\sum_{i=1}^m \exp(a_i^T x + b_i)} = 0.$$

$$\cancel{\sum_{i=1}^m \exp(a_i^T x + b_i)}$$

(always  
tve)

$$\sum_{i=1}^m a_i \exp(a_i^T x + b_i) = 0.$$

↗ system of  
n equations

$\underbrace{\sum_{i=1}^m a_i}_{\mathbb{R}^n} \underbrace{\exp(a_i^T x + b_i)}_{\text{scalar}} \underbrace{0}_{\mathbb{R}^n}$

No analytical solution, so we must use an iterative algorithm to get close to  $x^*$ .

Preview:

- i) Condition. → strong convexity  
→ smoothness
- ii) Gradient Descent.

more conditions,  
the better the  
gradient descent  
(faster)

works → prove  
convergence (how fast)

- first order
- iii) Stochastic Gradient Descent.  
 iv) Steepest Descent  
 v) Newton's method. 2nd order
- 

### Strong Convexity

normal convexity  $\nabla^2 f(x) \geq 0 \quad \forall x \in S$

Assume  $f: \mathbb{R}^n \rightarrow \mathbb{R}$   $m > 0$   $m$ -strongly convex on  $S = \{x : f(x) \leq f(x^{(0)})\}$

i.e.,  $\nabla^2 f(x) \geq mI$   $\forall x \in S$ .  
 strictly larger than 0  
 (some lower bound)  
 $m$  convex  $\Rightarrow$  convex

Ex: If  $f$  is  $m$ -strongly convex, consider the new function

$$g(x) = f(x) - \underbrace{\frac{m}{2} \|x\|^2}_{\text{convex}}. \text{ Then } g \text{ is convex.}$$

Implications: By Taylor's theorem  $\forall x, y \in S, \exists z \in S$

on the line segment connecting  $x$  &  $y$  s.t.

control to hessian

$$f(y) = f(x) + (y-x)^T \nabla f(x) + \frac{1}{2} (y-x)^T \underbrace{\nabla^2 f(z)}_{\geq mI} (y-x).$$

If  $f$  is  $m$ -strongly convex,

$$f(y) \geq f(x) + (y-x)^T \nabla f(x) + \boxed{\frac{m}{2} \|y-x\|^2} \quad \forall x, y \in S - (*)$$

When  $m=0$ ,  $(*)$  becomes  $f(y) \geq f(x) + (y-x)^T \nabla f(x)$ , which is exactly convexity for diff<sup>ble</sup> functions.

$m$  is the strong convexity constant

used to characterize sub-optimality gap

PL inequality (Polyak-Łojasiewicz).

Lemma: If  $f$  is  $m$ -strongly convex,  $\| \nabla f(x) \|_2^2 \geq f(x) - f^*$

smallest possible

$$p^* \geq f(x) - \frac{1}{2m} \| \nabla f(x) \|_2^2$$



if this small then we are close to optimal (e.g.  $p^*$ )

Rank: If the gradient  $\nabla f(x)$  is small, then  $x$  is nearly optimal

Pf: Since  $f$  is  $m$ -strongly convex,

$$\text{unconstrained } f(y) \geq f(x) + (y-x)^T \nabla f(x) + \frac{m}{2} \| y-x \|^2.$$

so smallest is  $p^*$

RHS is a convex, quadratic of  $y$

Minimize both sides w.r.t.  $y$ .

$$\nabla \text{RHS}(y) = 0$$

$$\nabla f(x) + m(y-x) = 0$$

$y$  that minimizes the RHS

$$p^* \geq f(x) + (-\frac{1}{m} \nabla f(x))^T \nabla f(x)$$

$$+ \frac{m}{2} \| -\frac{1}{m} \nabla f(x) \|^2$$

$$y = x - \frac{1}{m} \nabla f(x)$$

$$= f(x) - \frac{1}{2m} \| \nabla f(x) \|^2$$

(qed).

Rank: If the gradient  $\nabla f(x)$  is small, then  $x$  is nearly optimal

$$p^* \geq f(x) - \frac{1}{2m} \| \nabla f(x) \|^2$$

$$\| \nabla f(x) \| \leq \sqrt{2m\varepsilon} \rightarrow f(x) - p^* \leq \varepsilon.$$

If norm small, very close to optimality

$x^*$  is global optimal solution  
 $x$  is the current solution

The  $m$ -strongly convex condition can also be used to derive a bound on  $\|x - x^*\|_2$ .  
 ↳ sub-optimality of  $x$  itself

$$f(y) \geq f(x) + (y-x)^T \nabla f(x) + \frac{m}{2} \|y-x\|^2 \quad \forall x, y \in S.$$

set  $y = x^*$

$$\begin{aligned} p^* = f(x^*) &\geq f(x) + (x^* - x)^T \nabla f(x) + \frac{m}{2} \|x^* - x\|^2 && \text{cauchy} \\ &\geq f(x) - \underbrace{\|x^* - x\| \|\nabla f(x)\|}_{a} + \frac{m}{2} \|x^* - x\|^2. && \left. \begin{array}{l} \text{ineq M} \\ \text{opp direction} \end{array} \right\} \end{aligned}$$

$$\|a\| \|\nabla f(x)\| \geq a^T b \geq -\|a\| \|b\|$$

Since  $p^* \leq f(x)$ ,  $\forall x$ ,

$$-\|x^* - x\| \|\nabla f(x)\| + \frac{m}{2} \|x^* - x\|^2 \leq 0.$$

$$\|x^* - x\| \leq \frac{2}{m} \|\nabla f(x)\|_2$$

gradient = 0  
characterizes the optimality of  $x^*$

Rmk: If  $\nabla f(x)$  is small, then  $x$  is close to  $x^*$ .

### Upper Bound on Hessian

If  $\exists M > 0$  s.t.  $\nabla^2 f(x) \leq MI \quad \forall x \in S$ , we say that  $f$  is  $M$ -smooth. hold if level sets are closed and bounded

Lemma: (Descent lemma)

$$f(y) \leq f(x) + \nabla f(x)^T (y-x) + \frac{M}{2} \|y-x\|^2$$

PF: By Taylor's theorem,  $\exists z \in [x, y]$  s.t.

$$f(y) = f(x) + \nabla f(x)^T (y-x) + \frac{1}{2} (y-x)^T \underbrace{\nabla^2 f(z)}_{\leq MI} (y-x)$$

$$\leq f(x) + \nabla f(x)^T (y-x) + \frac{M}{2} \|y-x\|^2 \quad \text{convex quadratic}$$

A consequence of descent lemma is

minimize over

all possible  $y$

$$p^* \leq f(x) - \frac{1}{2M} \|\nabla f(x)\|^2$$

is a counterpart to the PL inequality:  $p^* \geq f(x) - \frac{1}{2m} \|\nabla f(x)\|^2$

$f$  is both  $m$ -strongly convex &  $M$ -smooth,  $0 \leq m \leq M < \infty$   
 fast gradient descent  $MI \leq \nabla^2 f(x) \leq MI \quad \forall x \in S$ .

The ratio  $K = \frac{M}{m}$  is called the condition number of  $\nabla^2 f(x)$ .

### Descent Methods

$k$  is the iteration number

Produce a sequence  $\{x^{(k)}\}_{k \in \mathbb{N}}$  s.t.

$$x^{(k+1)} = x^{(k)} + t^{(k)} \Delta x^{(k)}$$

At iteration

$x$  want to move in the direction of the search direction

learning rate / step size

$$x^+ = x + t(\Delta x)$$

next iteration

(descent direction)

search direction.

(have many types)

most common is negative gradient)

Downhill  $\rightarrow$  we want to minimise function

next iteration derives something better

Descent methods:  $f(x^+) < f(x)$  unless  $x$  is already optimal.

From convexity,  $\nabla f(x^{(k)})^T(y - x^{(k)}) \geq 0 \Rightarrow f(y) \geq f(x^{(k)})$



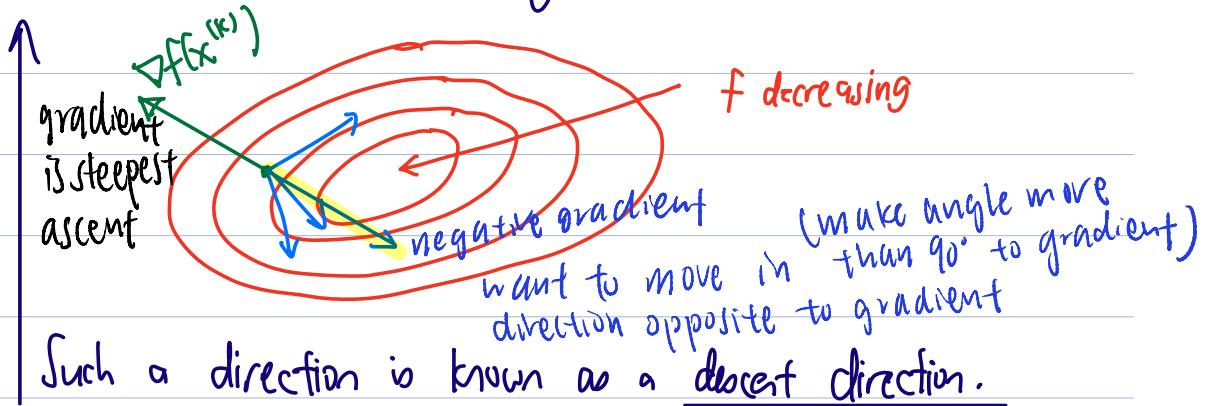
$f(y) \geq f(x^{(k)}) + \nabla f(x^{(k)})^T(y - x^{(k)})$  first order condition

go downhill

direction at this iteration

Search direction  $\Delta x^{(k)}$  must satisfy  $\nabla f(x^{(k)})^T(\Delta x^{(k)}) < 0$ .

need search direction to be different direction than gradient



General Descent Method      guess starting point

Given a starting point  $x^{(0)} \in \text{dom } f$

repeat

1. Determine a descent direction  $\Delta x^{(k)} = -\nabla f(x^{(k)})$

2. Line search: Choose a step size  $t$   $t_k$  (how much to move)

3. Update:  $x^+ = x + t \Delta x$

can tune  $t$  as a function of  $M$  and  $m$

stop  
( $100$  iterations)

or gradient until some stopping criterion is satisfied.  
sufficiently small?

Exact line search:  $t = \underset{s \geq 0}{\operatorname{arg\,min}} f(x + s \Delta x)$

optimisation problem  
in an optimisation problem

|  
already chosen  
direction

tune these hyperparameters

## Backtracking line search : (Armijo rule)

Given a descent direction  $\Delta x$ .  $\alpha \in (0, 0.5)$ ,  $\beta \in (0, 1)$ .

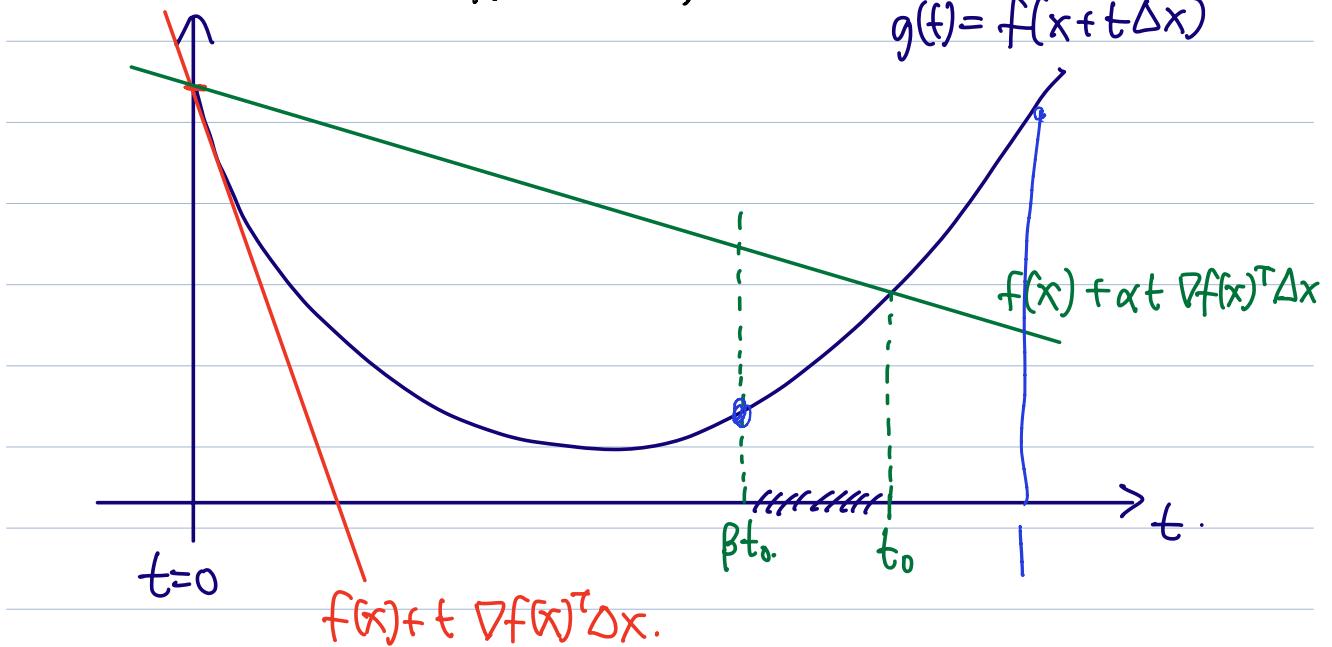
$t \leftarrow 1$  function still too big

while  $f(x + t \Delta x) > f(x) + \alpha t \nabla f(x)^T \Delta x$

$t \leftarrow \beta t$  (set  $t$  to be smaller)

restricting to a certain direction: still convex

$$g(t) = f(x + t \Delta x)$$



Backtracking line search stops with a step size of  $t$  that is  $t=1$  or  $t \in (\beta t_0, t_0)$

$\alpha$ : fraction of the  $\downarrow$  in  $f$  predicted by linear extrapolation that we will accept.

If  $x$  large, then loop have to go many times

↑ make obtuse angle  
with gradient

## Gradient Descent

$$\Delta x = -\nabla f(x).$$

Stopping criterion:  $\|\nabla f(x)\| \leq \eta$   $\eta > 0.$

(sufficiently small e.g.  $10^{-3}$ )

Thm: If GD is implemented on a function  $f$  that

i)  $m$ -strongly convex &  $M$ -smooth (i.e.,  $MI \leq \nabla^2 f(x) \leq MI$  for all  $x \in S$ ), with exact line search,

after  $k$  iterations  $f(x^{(k)}) - p^* \leq c^k (f(x^{(0)}) - p^*)$  linear convergence

(less than 1)  $c = 1 - \left(\frac{m}{M}\right)$  (contraction factor) every iteration contract by  $c$  (want  $c$  to be small)  
 $\Rightarrow$  want this close to 0

Pf: From the  $M$ -smoothness

$$x^* \quad f(y) \leq f(x) + \nabla f(x)^T (y-x) + \frac{M}{2} \|y-x\|^2$$

Let  $y = x - t \nabla f(x)$   $\tilde{f}(t) = f(x - t \nabla f(x))$ .  
(next iteration)

$$f(x - t \nabla f(x)) \leq f(x) + \nabla f(x)^T (-t \nabla f(x)) + \frac{M}{2} \| -t \nabla f(x) \|^2$$

||

$$\tilde{f}(t) \leq f(x) + \nabla f(x)^T (-t \nabla f(x)) + \frac{M}{2} \| -t \nabla f(x) \|^2$$

$$= f(x) - t \|\nabla f(x)\|^2 + \frac{Mt^2}{2} \|\nabla f(x)\|^2$$

min over all  $t$ .

Since we use exact line search, min LHS over all  $t$

$$t_{exact} = \underset{s \geq 0}{\operatorname{argmin}} f(x - s \nabla f(x))$$

choose step length to minimize this

$$f(x^t) = f(t_{\text{exact}}) \leq f(x) - \frac{1}{2M} \|\nabla f(x)\|^2$$

diff w.r.t t  
set to 0.

$$t^* = \frac{1}{M}$$

Subtract  $p^*$  on both sides.

$$f(x^t) - p^* \leq f(x) - p^* - \frac{1}{2M} \|\nabla f(x)\|^2.$$

Since  $f$  is  $m$ -strongly convex, PL ineq. holds

$$\|\nabla f(x)\|^2 \geq 2m(f(x) - p^*)$$

replace gradient  
with PL ineq

$$f(x^t) - p^* \leq f(x) - p^* - \frac{1}{2M} (2m(f(x) - p^*))$$

$$= \left(1 - \frac{m}{M}\right) (f(x) - p^*)$$

amount of improvement in 1 iteration

Applying this  $k$  times

$$f(x^{(k)}) - p^* \leq \left(1 - \frac{m}{M}\right)^k (f(x^{(0)}) - p^*) \quad \text{////.}$$

Cor: We have  $f(x^{(k)}) - p^* \leq \varepsilon$  if we use

$$k = \left\lceil \frac{\log((f(x^{(0)}) - p^*)/\varepsilon)}{\log(\frac{1}{c})} \right\rceil$$

condition # ↑  
denominator small

ratio of initial to final  
suboptimality.

Rmk: If  $\frac{M}{m}$  is large

want condition number to be small

if condition number is bad  $\Rightarrow$  we need many iterations

$$\log\left(\frac{1}{c}\right) = \log\left(\frac{1}{1 - \frac{m}{M}}\right) = -\log\left(1 - \frac{m}{M}\right) \approx \frac{m}{M}$$

$$\log(1-x) \approx -x$$

↗  
small.

### Analysis of backtracking line search.

Thm:  $f$  is  $m$ -strongly convex &  $M$ -smooth. Use Armijo or backtracking line search with  $(\alpha, \beta)$ . Then

$$f(x^{(k)}) - p^* \leq c^k (f(x^{(0)}) - p^*)$$

$$c = 1 - \min\left\{2m\alpha, 2\beta\alpha \frac{m}{M}\right\} < 1.$$

### Steepest Descent

$$f(x+v) \approx f(x) + \nabla f(x)^T v$$

$\nabla f(x)^T v$  is the directional derivative of  $f$  at  $x$  in the direction  $v$ .

↗ can derive  $f(x+v)$  function by a lot

Choose  $v$  to make  $\nabla f(x)^T v$  as negative as possible.

$$\Delta x_{nsd} = \arg \min \left\{ \nabla f(x)^T v : \|v\|=1 \right\}$$

↑

normalized steepest descent

$\|\cdot\|$  suitably chosen norm.

so cannot choose  $v$  to be very long vector

- choose different norm  
 $\Rightarrow$  get diff type of direction

Unnormalized version  $\Delta x_{sd} = \|\nabla f(x)\|_* \Delta x_{nsd}$

$$\|x\|_* = \sup \{ y^T x : \|y\| \leq 1 \}. \text{ dual norm}$$

$$\begin{aligned} \text{Rmk: } \nabla f(x)^T \Delta x_{sd} &= \nabla f(x)^T (\|\nabla f(x)\|_* \Delta x_{nsd}) \\ &= -\|\nabla f(x)\|_*^2 \text{ by definition of dual norm} \end{aligned}$$

Ex: Euclidean norm:  $\|\cdot\| = \|\cdot\|_2$

$$\begin{aligned} \Delta x_{sd} &= \arg \min \left\{ \nabla f(x)^T v : \|v\|_2 = 1 \right\} \\ &= -\nabla f(x) \underbrace{a}_{\text{(b constrained to be norm 1)}} \underbrace{b}_{\text{achieved when b is a?}} \end{aligned}$$

Steepest descent method with Euclidean norm is nothing but gradient descent!

Quadratic norm:  $P \in S_{ff}^n \quad \|z\|_P = \sqrt{z^T P z}$

$$\Delta x_{nsd} = -(\nabla f(x)^T P^{-1} \nabla f(x))^{-\frac{1}{2}} P^{-1} \nabla f(x).$$

$$\Delta x_{sd} = -P^{-1} \nabla f(x)$$

Can be regarded as a gradient method applied to the problem after a change of coordinates:  $\bar{x} = P^{\frac{1}{2}} x$ .

## $l_1$ -norm

$$\Delta x_{\text{nsd}} = \arg \min \left\{ \nabla f(x)^T v : \|v\|_1 = 1 \right\}.$$

$$= -\text{sign}\left(\frac{\partial f(x)}{\partial x_i}\right) e_i$$

$e_i$ :  $i^{\text{th}}$  standard basis vector.

$$\Delta x_{\text{sd}} = \Delta x_{\text{nsd}} \|\nabla f(x)\|_\infty = -\frac{\partial f}{\partial x_i} e_i$$

*vector*

At each iteration, choose a component of  $\nabla f(x)$  that has the largest abs value, decrease or increase the corresponding comp. of  $x$  by  $\text{sign}((\nabla f(x))_i)$ . only move in that coordinate rather than full gradient

$\Rightarrow$  coordinate descent algorithm.

## Stochastic Gradient Descent

$\leftarrow$  expensive to compute gradient

each

$$(x_i, y_i) \quad \text{Loss}(\theta; D) = \frac{1}{m} \sum_{i=1}^m \text{Loss}(\theta; x_i, y_i)$$

$i=1, \dots, m$

finite sum.

expensive to calculate full gradient not assuming m-s.c.

$$\min_{x \in \mathbb{R}^n} f(x) = \frac{1}{m} \sum_{i=1}^m f_i(x)$$

(a lot of training samples)

not using full gradient.

pick 1 training sample from  $m$

SGD:  $\rightarrow i_k \in [m]$  randomly.

$$\Pr(i_k=i) = \frac{1}{m} \quad \forall i \in [m]$$

$$\rightarrow x^{(k+1)} = x^{(k)} - t_k$$

$$\boxed{\nabla f_{i_k}(x^{(k)})}$$

$$\mathbb{E}[\nabla f_i(x)] = \nabla f(x)$$

one sample gradient

unbiased

one sample representative of whole dataset

Assumption: 1)  $p^* = \inf_{x \in \mathbb{R}^n} f(x)$  <sup>finite</sup> > -∞       $|M| - J^2 f(x)$   
 J is positive semi-definite  
 2) f: M-smoothness       $\nabla^2 f(x) \leq M I \quad \forall x$   
 3) Variance is bounded  
 $\sup_x \mathbb{E}[\|\nabla f_i(x)\|^2] \leq \sigma^2$   
 only the i (index) is random  
 tk step size

Thm: Assume 1), 2), 3), run SGD for T iterations.

find 1 iteration such that  $\min_{k=0,1,\dots,T} \mathbb{E}[\|\nabla f(x^{(k)})\|^2] \leq \frac{f(x^{(0)}) - p^*}{\sum_{k=0}^{T-1} t_k} + \frac{M\sigma^2}{2} \frac{\sum_{k=0}^{T-1} t_k^2}{\sum_{k=0}^{T-1} t_k}$

want this to be small

(or: 1)  $t_k = t > 0$ ,  $\sigma^2 = 0$  (no noise essentially doing gradient descent)  
 constant  $\sum_{k=0}^{T-1} t_k = tT$ . First term gives  $\frac{f(x^{(0)}) - p^*}{tT}$ .

Convergence rate of  $O(\frac{1}{T})$ .

If assume M-smoothness (no m-s.c.), GD converges at a rate of  $O(\frac{1}{T})$ , much slower than  $(C^T)$  exponential

2) SGD:  $t_k = \frac{1}{\sqrt{k}}$        $\sum_{k=1}^T \frac{1}{\sqrt{k}} \sim \sqrt{T}$ .       $\sum_{k=1}^T t_k^2 = \sum_{k=1}^T \frac{1}{k}$   
 (approximation)       $= O(\log T)$ .

$\min_{k=0,1,\dots,T} \mathbb{E}[\|\nabla f(x^{(k)})\|^2] \leq O\left(\frac{\log T}{\sqrt{T}}\right)$

slower than  $(\frac{1}{T})$   
 worse than gradient descent without s-c assumption

minibatch set

$N_{\text{it}}$

$\rightarrow B_k \subseteq [n]$  randomly.

$|B_k| = \# \text{ machines}$

$\rightarrow x^{(k+1)} = x^{(k)} - t_k \frac{1}{|B_k|} \sum_{i \in B_k} \nabla f_i(x^{(k)})$  batch gradient

Pf: Only assume  $f$  is  $M$ -smooth.  $\Rightarrow$  Descent lemma.

$$f(x^{(k+1)}) \leq f(x^{(k)}) + \nabla f(x^{(k)})^T (x^{(k+1)} - x^{(k)}) + \frac{M}{2} \|x^{(k+1)} - x^{(k)}\|^2$$

only thing random is sample chosen

SGD does:

$$x^{(k+1)} = x^{(k)} - t_k \nabla f_{i_k}(x^{(k)})$$

at iteration  $k$ , pick an index  $i_k$

$$f(x^{(k+1)}) \leq f(x^{(k)}) - t_k \nabla f(x^{(k)})^T \nabla f_{i_k}(x^{(k)}) + \frac{Mt_k^2}{2} \|\nabla f_{i_k}(x^{(k)})\|^2$$

Take expectation

on both sides

$$\mathbb{E}[f(x^{(k+1)})] \leq f(x^{(k)}) - t_k \nabla f(x^{(k)})^T \mathbb{E}[\nabla f_{i_k}(x^{(k)})]$$

every time see  $i_k$ , it is random and we take  $E$

$$\text{unbiased} \quad + \frac{Mt_k^2}{2} \mathbb{E}[\|\nabla f_{i_k}(x^{(k)})\|^2].$$

$$\leq f(x^{(k)}) - t_k \|\nabla f(x^{(k)})\|^2 + \frac{Mt_k^2}{2} \sigma^2$$

$$\Rightarrow t_k \|\nabla f(x^{(k)})\|^2 \leq f(x^{(k)}) - \mathbb{E}[f(x^{(k+1)})] + \frac{Mt_k^2 \sigma^2}{2}$$

Run SGD for  $T$  iterations and sum up both sides.

take expectation  
on both sides, then sum up

$$\sum_{k=1}^T t_{k-1} \mathbb{E}[\|\nabla f(x^{(k)})\|^2] \leq \sum_{k=1}^T [\underbrace{\mathbb{E}[f(x^{(k)}) - f(x^{(k+1)})]}_{c_k} + \underbrace{\frac{M\sigma^2}{2} t_k^2}_{c_{k+1}}]$$

$$+ \frac{M\sigma^2}{2} \sum_{k=1}^T t_k^2$$

$$(c_1 - c_2) + (c_2 - c_3) + (c_3 - c_4) + \dots + (c_T - c_{T+1})$$

telescoping

$$\sum_{k=1}^T t_{k-1} \mathbb{E}[\|\nabla f(x^{(k)})\|^2] \leq f(x^{(0)}) - \mathbb{E}[f(x^{(T)})] + \frac{M\sigma^2}{2} \sum_{k=1}^T t_k^2$$

\geq p^\* \text{ (finite)}

**Lower Bound** LHS by  $\min_{k=0,1,\dots,T-1} \mathbb{E}[\|\nabla f(x^{(k)})\|^2]$ .

$$\min_{k=0,1,\dots,T-1} \mathbb{E}[\|\nabla f(x^{(k)})\|^2] \leq \frac{f(x^{(0)}) - p^*}{\sum_{k=0}^{T-1} t_k} + \frac{M\sigma^2}{2} \frac{\sum_{k=0}^{T-1} t_k^2}{\sum_{k=0}^{T-1} t_k}$$

$$t_k = \frac{1}{\sqrt{k}} \quad (\text{step size})$$

$$\sum_{k=1}^T \frac{1}{k} \approx \int_1^T \frac{1}{u} du = \log T = O\left(\frac{\log T}{\sqrt{T}}\right).$$

$$\sum_{k=1}^T \frac{1}{\sqrt{k}} = \int_1^T \frac{1}{\sqrt{u}} du \approx \sqrt{T}$$

- Summary:
- 1) Conditions on function ( $m$ -strongly convex,  $M$ -smooth)
  - 2) Gradient descent has linear convergence  $\uparrow$

3) Steepest descent

4) Stochastic GD ( $M$ -smoothness)

assume less about function



get worse results