

## Gradient Descent:

①

• Convex set:  $K$  st  $\forall x, y \in K, \lambda x + (1-\lambda)y \in K \quad \forall \lambda \in [0, 1]$ .

• Convex function:  $f: X \rightarrow \mathbb{R}$  is convex if

(i)  $f(\lambda x + (1-\lambda)y) \leq \lambda f(x) + (1-\lambda)f(y), \quad \forall x, y \in X.$

or (ii) for all lines in  $X$ ,  $f(x)$  restricted to the line is convex like we think of 1-d convexity  $\checkmark$ .

• if  $f$  is differentiable then define the gradient of  $f$  at  $x$ :

$$\nabla f(x) = \left( \frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right)$$

such  $f$  is convex iff  $f(y) \geq f(x) + \langle \nabla f(x), y-x \rangle$   
 $\uparrow$  the linear approximation of the function  $f$  at pt  $x$ .

Eg:  $f(x) = x^2$  then  $\nabla f(x) = 2x$  and linear approximation is  
 $y^2 \stackrel{?}{\geq} f(x) + \langle \nabla f(x), y-x \rangle$   
 $= x^2 + 2x(y-x)$   
 $\Rightarrow y^2 - 2xy + x^2 \stackrel{?}{\geq} 0$  which is true.

---

if  $f$  is twice differentiable then Hessian( $f$ ) =  $H(f) = \nabla^2 f$ .

such  $f$  is convex iff  $H(f)$  is positive semidefinite (all its eigenvalues  $\geq 0$ ).  
(analog of 2nd derivative being positive).

---

## Problems today:

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

$\min_{x \in K} f(x) \leftarrow$  constrained fn. minimization

① we'll give approximate solutions: - additive approximations.

i.e. if  $x^*$  is the minimizer then find  $\hat{x}$  st  $f(\hat{x}) - f(x^*) \leq \epsilon$ .

② Need to assume something about  $f$ , and about the distance b/w  $x_{\text{start}}$  &  $x^*$ .

Simplest case [minimal assumptions]  $\leftrightarrow$  can we assume even less?

(2)

•  $f$  is differentiable, convex

•  $f$  is  $G$ -Lipschitz (e.  $|f(x) - f(y)| \leq L \|x - y\|$   $\forall x, y \in K$ )  $\leftarrow$  assume Euclidean norm, say, for now.

$$\Leftrightarrow \|\nabla f(x)\| \leq G \quad \forall x \in K.$$

Btw: let's figure out some properties of minimizer  $x^*$  of  $f(x)$ . (convex)

• unconstrained:  $x^*$  is a <sup>global</sup> minimizer of  $f$  in  $\mathbb{R}^n$   
 $\Leftrightarrow x^*$  is a local min of  $f$  in  $\mathbb{R}^n$   
 $\Leftrightarrow \nabla f(x^*) = \vec{0}$

• constrained ( $x \in K$ ):  $x^*$  is global min  $\Leftrightarrow x^*$  is local min in  $K$   
 $\Leftrightarrow \langle \nabla f(x^*), y - x^* \rangle \geq 0 \quad \forall y \in K.$

• Let's do unconstrained minimization  $\min_{x \in \mathbb{R}^n} f(x)$   
(diff, convex,  $G$ -Lipschitz).

Assume that  $x_0$  is known, such that  $\|x_0 - x^*\| \leq D \leftarrow$  diameter.

Will find  $\epsilon$ -approximate minimizer; i.e.  $\exists \bar{x}$  st  $f(\bar{x}) - f(x^*) \leq \epsilon$ .  
in time poly( $G, D, \epsilon$ ).

Algorithm: Gradient Descent

Start at  $\bar{x}_0 \in \mathbb{R}^n$

At each timestep  $\bar{x}_{t+1} \leftarrow \bar{x}_t - \eta_t \cdot \nabla f(\bar{x}_t)$

Repeat for  $T$  steps.

Intuition:  
 $\min_x \frac{1}{2} \|y - x\|^2 + \eta \langle \nabla f(x), y - x \rangle$

Q: What is  $\eta_t$ ? For this algo, we'll set  $\eta_t = \eta$  for some  $\eta$ .

Q: What is  $T$ ?

And what guarantees can we set on final fn value?

Thm 1: Given  $\epsilon, x_0$  s.t.  $\|x_0 - x^*\| \leq D$ , function  $f$   $\mathbb{R}^n$  diff, convex,  $G$ -Lipschitz,

the gradient descent algorithm produces  $x_1, x_2, \dots, x_T$

such that  $f(x_T) - f(x^*) \leq \epsilon$  if  $T \leq \left(\frac{GD}{\epsilon}\right)^2$ .

[N.b. if  $G, D$  are constants, then this is like  $T = O(1/\epsilon^2)$ .]

[There are examples s.t. this guarantee is tight].

Pf:  $f(x_t) - f(x^*) \leq \langle \nabla f(x_t), x_t - x^* \rangle$

$= \langle \frac{1}{\eta}(x_t - x_{t+1}), x_t - x^* \rangle$   $\langle \bar{a}, \bar{b} \rangle = \left( \|\bar{a}\|^2 + \|\bar{b}\|^2 - 2\|\bar{a}\|\|\bar{b}\|\cos\theta \right)^{\frac{1}{2}}$

$= \frac{1}{\eta} \cdot \frac{1}{2} (\|x_t - x_{t+1}\|^2 + \|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2)$

$= \frac{1}{2\eta} (\|\eta \nabla f(x_t)\|^2 + \|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2)$

$\uparrow \leq \eta^2 G^2$   $\uparrow$  will telescope

$\Rightarrow \frac{1}{T} \sum_{t=1}^T (f(x_t) - f(x^*)) \leq \frac{1}{2\eta} \left[ \eta^2 G^2 + \frac{\|x_0 - x^*\|^2}{T} \right]$  drop negative term.

$= \frac{1}{2\eta} \left[ \eta^2 G^2 + \frac{D^2}{T} \right].$

[apply Jensen's, and define  $\hat{x} = \frac{1}{T} \sum x_t$ ]

$f(\hat{x}) = f\left(\frac{1}{T} \sum x_t\right) \leq \frac{1}{T} \sum f(x_t).$

$\Rightarrow f(\hat{x}) - f(x^*) \leq \frac{1}{2} \left[ \eta G^2 + \frac{D^2}{\eta T} \right]$

Set  $\eta = \frac{D}{G\sqrt{T}}$  gives  $\leq \frac{DG}{\sqrt{T}} \leq \epsilon$

if we set  $T \geq \left(\frac{DG}{\epsilon}\right)^2$ .

Btw: if we set  $\eta_t = \frac{D}{\sqrt{t} \cdot G}$  then base a small log T term but now not dep. on length.

Not a polytime algo, since rate of convergence is slow.

Basically get one more bit of info ( $\epsilon \rightarrow \epsilon/2$ ) need 4 times more time.

$\Rightarrow$  for  $n$  bits will need exponential time. Want better.

But I tight examples if we don't assume more  $\leftarrow$  (soon).

But it's a great start: -

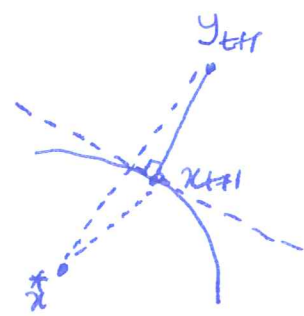
① Constrained optimization  $\min_{x \in K} f(x)$ .

Same guarantees. Almost same algorithm; change thus: -

$$y_{t+1} \leftarrow x_t - \eta \nabla f(x_t).$$

but  $y_{t+1}$  may be outside  $K$ .

so def  $x_{t+1} \leftarrow \Pi_K y_{t+1}$  = nearest point in  $K$  to  $y_{t+1}$   
= "projection" of  $y_{t+1}$  onto  $K$ .



Lemma:  $\|y_{t+1} - x^*\|^2 \geq \|x_{t+1} - x^*\|^2$

B/c the angle between  $x^*$  and  $y_{t+1}$  must be obtuse, since  $x^*, x_{t+1} \in K$  and  $K$  is convex.

$\Rightarrow$  do the same analysis: -

$$f(x_t) - f(x^*) \leq \dots \leq \frac{1}{2\eta} [\eta^2 G^2 + \|x_t - x^*\|^2 - \|y_{t+1} - x^*\|^2]$$

$$\leq \frac{1}{2\eta} [\eta^2 G^2 + \|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2] \dots$$

and same follows.

$\Rightarrow$  no change.

Can do other things: eg [Frank Wolfe algorithm]: find the optimizer  $\hat{z}$  in the

direction of  $\nabla f(x_t)$ , say  $\hat{z}_{t+1} = \min_{z \in K} \langle \nabla f(x_t), z \rangle \leftarrow$  linear optimization

and take a small step in the direction of  $\hat{z}_{t+1}$  to get  $x_{t+1}$ .

Guarantees?

② Online convex optimization:

Each time algorithm picks point  $x_t$

Then adversary picks function (convex)  $f_t(\cdot)$ .

sps  $\|\nabla f_t(x_t)\| \leq G$   
 $\forall t \forall x \in K$ .

Cost =  $f_t(x_t)$ . want to minimize regret.

Algo: same. just use  $x_{t+1} \leftarrow x_t - \eta \nabla f_t(x_t)$ .

Analysis: shows  $\sum_{t=1}^T (f_t(x_t) - f_t(x^*)) \leq \frac{DG}{\sqrt{T}} \leq \epsilon$  if  $T \geq \frac{(DG)^2}{\epsilon^2}$ .

- Same  $\frac{1}{\epsilon^2}$  dependence, but now depends on  $D, G$  not on  $\lg N$  (#experts).

- However, holds for all convex bodies  $K$ , not just for the probability simplex.

What does it give for  $K = \Delta_N$ ? and say loss function  $l_t(x) = \sum_i l_i^t x_i$   
 $\|\nabla l_t\|_2 \leq \|\mathbf{1}\|_2 \leq \sqrt{N}$  since each  $l_i^t \in [0, 1]$ .

$\text{diam}(K) = 1$ .

$\Rightarrow$  ~~get~~ need  $T \geq \frac{N}{\epsilon^2}$  ~~needed~~ for  $\epsilon$ -regret-on-average. much worse. (but we'll do better).

③ What if not differentiable?

Use "subgradients". need a vector  $\nabla f(x_t)$  such that

$f(y) \geq f(x_t) + \langle \nabla f(x_t), y - x_t \rangle$ .

If not diff, may have many such vectors. [Set called  $\partial f(x_t)$ ]

Take any of those. [also if approximate subgradients, can lose a little  $\delta$  and still be OK etc.]

N.b. no dependence on the dimension  $n$ . [not explicitly, anyways].

• Could choose a single coordinate to move along — at random

$x_{t+1} \leftarrow x_t + \eta_t \langle \nabla f(x_t), e_i \rangle \cdot e_i$  [Coordinate descent]

Naturally slows down things by factor of  $n \leftarrow$  dimensions.

How do you get the gradient?

- Compute it explicitly
- estimate it explicitly.

- or "stochastic gradient" descent: -  $\nabla f(\bar{x}) = \lim_{\delta \rightarrow 0} \mathbb{E}_{\bar{u} \in B(0, \delta)} \left[ \frac{f(\bar{x} + \bar{u}) - f(\bar{x})}{\delta} \cdot \bar{u} \right]$

So can use a single sample to move along.

Q: how noisy is this estimator  $\leftarrow$  depends, see lots of work on this

But now: if we assume more, we get more.

Def:  $f$  is  $\ell$ -strongly convex if

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\ell}{2} \|y - x\|^2.$$

*locally "grows" at least as fast as a quadratic*

Fact:  $f$  is  $\ell$ -strongly convex  $\Leftrightarrow H(f) \succeq \ell I$ . (min eigenvalue  $\geq \ell$ ).

For this case: almost same analysis works.

choose a step size  $\eta_t \leq \frac{1}{L t}$  or so (constants dep on  $G, D$  dropped).

get that  $T = \frac{O(G^2 / \epsilon T)}{\ell \epsilon}$  is enough [HW]

Even faster?

*"locally grows no faster than a quadratic"*

equivalent to  $\| \nabla f(x) - \nabla f(y) \| \leq \beta \|x - y\|$   
 $\Rightarrow$  gradients are also  $\beta$ -Lipschitz.

Def:  $f$  is  $\beta$ -smooth if

$$f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\beta}{2} \|y - x\|^2.$$

Fact:  $f$  is  $\beta$ -smooth  $\Leftrightarrow H(f) \preceq \beta I$ .

For this case too, slightly different analysis shows that  $T = \Theta(\frac{1}{\epsilon})$  enough

What if we assume both?  $\ell I \preceq H(f) \preceq \beta I$

$f$  is "well conditioned"

then gradient descent gives a  $\log(\frac{1}{\epsilon})$  convergence.  $\leftarrow !!$

Also called "linear convergence"  $\leftarrow$  one bit of accuracy per constant # of steps.

Assume that  $f$  is both  $\beta$ -smooth  
and  $l$ -strongly-convex

$$lI \leq H(f) \leq \beta I.$$

Fact 1:  $f(x_t) - f(x^*) \leq \frac{\|\nabla f(x_t)\|^2}{2l}$

Pf:  $f(y) \geq f(x) + \langle \nabla f(x), y-x \rangle + \frac{l}{2} \|y-x\|^2$  by  $l$ -strong-convexity

$$\geq \min_z \left\{ f(x) + \langle \nabla f(x), z-x \rangle + \frac{l}{2} \|z-x\|^2 \right\}$$

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

$$\Rightarrow f(x_t) - f(x^*) \leq \frac{1}{2l} \|\nabla f(x_t)\|^2$$

$$\begin{aligned} \nabla f(x) &= l \cdot (x-z) \\ \Rightarrow z &= x - \frac{\nabla f(x)}{l} \\ \text{plug back in} \end{aligned}$$

Claim:  $f(x_{t/\beta}) - f(x^*) \leq \exp\left(-\frac{t}{\beta l}\right) (f(x_0) - f(x^*))$

Pf: let  $\Delta_t = f(x_t) - f(x^*)$ .

condition #

$$\Rightarrow \Delta_t - \Delta_{t-1} = f(x_t) - f(x_{t-1})$$

$$\leq \langle \nabla f(x_{t-1}), x_t - x_{t-1} \rangle + \frac{\beta}{2} \|x_t - x_{t-1}\|^2$$

$$= -\eta_t \|\nabla f(x_{t-1})\|^2 + \frac{\beta}{2} \eta_t^2 \|\nabla f(x_{t-1})\|^2$$

set  $\eta_t = \frac{1}{\beta}$

$$= -\frac{1}{2\beta} \|\nabla f(x_{t-1})\|^2 \leq -\frac{1}{2\beta} \cdot 2l \Delta_{t-1}$$

from fact 1.

$$= -\left(\frac{l}{\beta}\right) \Delta_{t-1}$$

$$\Rightarrow \Delta_t \leq \left(1 - \frac{l}{\beta}\right) \Delta_{t-1}$$

$$\leq \Delta_0 e^{-(l/\beta)t}$$

Can also show convergence  $x_t \rightarrow x^*$  (see the other proof)

(7)

Theorem: Suppose  $f$  is both  $\beta$ -smooth &  $l$ -strongly convex.

then if  $x^* \leftarrow \min_{x \in \mathbb{R}^n} f(x)$  [unconstrained]  $L \leq H(f) \leq \beta L$ .

we get (1)  $\|x_t - x^*\|^2 \leq e^{-\delta t} \|x_0 - x^*\|^2$  for some  $\delta = \delta(\beta, l)$

(2)  $|f(x_t) - f(x^*)| \leq \frac{\beta}{2} \|x_0 - x^*\|^2 \exp(-\frac{4t}{\beta/l})$

if we set  $\eta = \frac{2}{l+\beta}$

$\frac{4t}{\beta/l}$  = condition # of the function.

Proof: [will skip in class].