

Intro & Overview, Gradient Descent.

Monday, January 6, 2020 5:42 PM

What

Algorithms are discrete-time procedures to solve problems.

e.g. SMT

Does f have a SATISFYING assignment?

What should my next move be?

Translate "Vagyok boldog" to English

$$\boxed{\text{Optimization}}: \min_{x \in \Omega} f(x)$$

very general. E.g. $\exists x : f(x) = 1$

$$\Leftrightarrow \min |f(x) - 1|$$

also intractable . $f(x) = \begin{cases} 1 & x \neq x^* \\ 0 & \text{o.w.} \end{cases}$

$$\min f(x)$$

takes an infinite # evaluations of f .

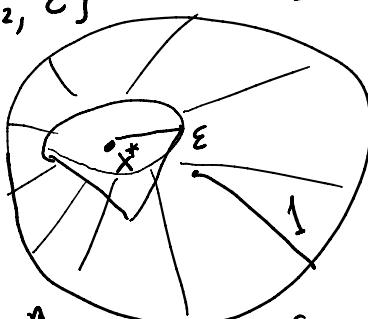
But what if f is Lipschitz, i.e. $\|\nabla f\| \leq 1$.

$$L(f) = \sup_{x,y} \frac{|f(x) - f(y)|}{\|x - y\|} \quad \text{and with bounded domain.}$$

$$f(x) = \min \{ \|x - x^*\|_2, \varepsilon \} \quad x \in B_2^n$$

is 1 -Lipschitz.

$$\begin{aligned} \text{Vol}(\{x : \|x - x^*\|_2 \leq \varepsilon\}) \\ \leq \varepsilon^n \cdot \text{Vol}(B_2^n) \end{aligned}$$



Finding x takes $\mathcal{O}(\frac{1}{\varepsilon^2})$ calls to f .

Finding x takes $\mathcal{O}\left(\frac{1}{\varepsilon}\right)^n$ calls to f .

Exercise. Show that $\mathcal{O}\left(\frac{1}{\varepsilon}\right)^n$ calls suffice.

[P2] Sampling. $f: \mathbb{R}^n \rightarrow \mathbb{R}$, sample x with density

$$p(x) \propto e^{-f(x)}.$$

$$\text{l.g. } f(x) = \frac{\|x\|^2}{2}, \quad f(x) = \begin{cases} 0 & x \in K \\ \infty & \text{o.w.} \end{cases}$$

OPT \rightarrow Sampling (binary search)

i.e. intractable!

Examples in practice: flow, matching, LP

regression, sampling Gaussian, convex body, logistic regression.

$$\frac{1}{n} \sum_{i=1}^n f(\alpha^T x_i \cdot y_i) + \lambda \|\alpha\|$$

Suppose we assume f is convex.

$$f(\lambda x + (1-\lambda)y) \leq \lambda f(x) + (1-\lambda)f(y)$$

then all kinds of wonderful things happen!

① Any local minimum $\nabla f(x) = 0 \Rightarrow$ global min $f(x^*) \leq f(x) \forall x$.

Why? Lemma. $f(y) \geq f(x) + \nabla f(x)^T (y-x)$

$$\text{Pf. } g(\lambda) = f((1-\lambda)x + \lambda y) \quad g(0) = x \quad g(1) = y$$

$$\text{by convexity } g(\lambda) \leq (1-\lambda)g(0) + \lambda g(1).$$

$$g(1) \geq g(0) + \frac{g(1) - g(0)}{\lambda}$$

$$\lambda \rightarrow 0$$

$$g(1) \geq g(0) + g'(0)$$

$$g'(0) = \nabla f(x)^T (y-x)$$

$$f(y) \geq f(x) + \nabla f(x)^T (y-x).$$

$$\text{closed : } \forall x, y \in K \Rightarrow [x, y] \subseteq K.$$

$$f(y) \geq f(x) + \nabla f(x) \cdot (y-x).$$

② Separation. K is convex, i.e. $x, y \in K \Rightarrow [x, y] \subseteq K$.

$$y \notin K \Rightarrow \exists a : a^T y > \max_{x \in K} a^T x.$$

This lets us do binary search!

Examples of convex functions.

Gradient Descent.

while $\|\nabla f(x)\| > \varepsilon$:

$$x \leftarrow x - h \nabla f(x)$$

Goal: Find x : $\|\nabla f(x)\| \leq \varepsilon$.

Q. how many iterations of GD? What h to use?

Why this goal: x is near-minimum in its neighborhood.

But even this is not true if ∇f can change quickly.

Assume. $\|\nabla^2 f\|_{op} \leq L$ i.e. $\nabla^2 f$ is L -Lipschitz.

Thm. If $\nabla^2 f$ is L -Lip, $\forall \varepsilon > 0$, starting at x_0 , we

reach x with $\|\nabla f(x)\| \leq \varepsilon$ in at most

$$2 \frac{L}{\varepsilon} (f(x_0) - f(x^*)) \text{ steps}$$

Lem 1. $f(x - \frac{1}{L} \nabla f(x)) \leq f(x) - \frac{1}{2L} \|\nabla f(x)\|^2$

If $f: \mathbb{R}^d \rightarrow \mathbb{R}$ and x, y

P.
Thm [Taylor]. If $k+1$ times differentiable $g: \mathbb{R} \rightarrow \mathbb{R}$, x, y

$$g(y) = g(x) + \sum_{i=1}^k \frac{(y-x)^i}{i!} g^{(i)}(x) + \frac{(y-x)^{k+1}}{(k+1)!} g^{(k+1)}(\zeta) \quad \zeta \in [x, y]$$

$$g(z) = g(0) + g'(0) + \frac{1}{2} g''(z) \quad z \in [0, 1].$$

$$g(t) = f((1-t)x + t y)$$

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

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

$$\begin{aligned} f\left(x - \frac{1}{L} \nabla f(x)\right) &\leq f(x) - \frac{1}{L} \|\nabla f(x)\|^2 + \frac{1}{2} L \cdot \frac{\|\nabla f(x)\|^2}{L^2} \\ &\leq f(x) - \frac{1}{2L} \|\nabla f(x)\|^2. \end{aligned}$$

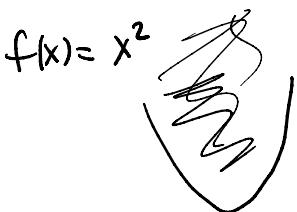
Start at $f(x^0)$, min in $f(x^*)$. each step decrease is

$$\geq \frac{\varepsilon^2}{2L}.$$

$$\Rightarrow \# \text{steps} \leq \frac{2L}{\varepsilon^2} (f(x^0) - f(x^*)).$$

f convex

$$\text{Epi}(f) = \{(x, t) : f(x) \leq t\}$$



Exercise: $\text{Epi}(f)$ is convex.

$$\min f(x) = \arg \min_t \text{Epi}(f).$$

\min convex functions \Leftrightarrow minimizing Linear functions
on convex sets.

Next time: 6D for convex f .

Checking convexity:

Lemma: Equivalent: (1) f is convex
 (2) $\forall x, y \quad f(y) \geq f(x) + \nabla f(x)^T (y-x)$
 (3) $\forall x \quad \nabla^2 f(x) \succcurlyeq 0$.

Pf: (1) \Rightarrow (2) earlier.

$$(2) \Rightarrow (3) \quad \text{By Taylor} \quad f(x+th) = f(x) + t \nabla f(x)^T h + \frac{t^2}{2} h^T \nabla^2 f(z) h \quad z \in [x, x+th]$$

$$\Rightarrow h^T \nabla^2 f(z) h \geq 0 \quad \forall h$$

$$t \rightarrow 0 \Rightarrow \nabla^2 f(x) \succeq 0.$$

$$(3) \Rightarrow (1) \quad g(\lambda) = f(\lambda x + (1-\lambda)y) - \lambda f(x) - (1-\lambda)f(y)$$

$$g'(\lambda) = \nabla f(\quad)(x-y) - f(x) + f(y)$$

$$g''(\lambda) = (x-y)^T \nabla^2 f(\quad)(x-y)$$

$$\text{Let } \lambda^* = \arg\max g(\lambda). \quad \text{If } \lambda^* = 0 \text{ or } 1, \Rightarrow g(\lambda) \leq 0 \quad \checkmark$$

$$g(1) = g(\lambda^*) + g'(\lambda^*) + \frac{1}{2} g''(z)(1-\lambda^*)^2 \quad z \in [\lambda^*, 1]$$

$$0 = g(1) \Rightarrow g(\lambda^*) \geq g(\lambda). \quad \checkmark.$$

what about (1) applied to convex f ?

$$\text{Th.} \quad f(x^k) - f(x^*) \leq \frac{2LR^2}{k+4}$$

$$\|\nabla^2 f\|_{op} \leq L$$

$$\max \|x - x^0\|_2 \leq R$$

$$f(x) \leq f(x^0)$$

\dots, x^1, x^k, x^*

$$f(x) \leq f(x^*)$$

Pf.

$$\begin{aligned} \varepsilon_k &= f(x^k) - f(x^*) \leq \nabla f(x^k)(x^k - x^*) \\ &\leq \|\nabla f(x^*)\|_2 \|x^k - x^*\|_2 \\ \frac{\varepsilon_k}{R} &\leq \|\nabla f(x^*)\|. \end{aligned}$$

Next

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

$$\varepsilon_{k+1} \leq \varepsilon_k - \frac{1}{2L} \frac{\varepsilon_k^2}{R^2}$$

$$\frac{1}{\varepsilon_{k+1}} - \frac{1}{\varepsilon_k} = \frac{\varepsilon_k - \varepsilon_{k+1}}{\varepsilon_{k+1}\varepsilon_k} \geq \frac{1}{2LR^2}.$$

Also,

$$\begin{aligned} \varepsilon_0 &= f(x^0) - f(x^*) = \nabla f(x^*)(x^* - x^0) + \frac{1}{2} (x^* - x^0)^T \nabla f^2(x^*) \\ &\quad z \in [x^0, x^*] \\ &\leq 0 + \frac{1}{2} L \cdot \|x^* - x^0\|^2 \\ \varepsilon_0 &\leq \frac{1}{2} LR^2 \end{aligned}$$

$$\therefore \frac{1}{\varepsilon_k} \geq \frac{k}{2LR^2} + \frac{1}{\varepsilon_0} \geq \frac{k+4}{2LR^2}$$

or

$$\varepsilon_k \leq \frac{2LR^2}{k+4}$$

$K+4$

to get error $\propto \varepsilon$, takes $O\left(\frac{1}{\varepsilon}\right)$ steps

How to check convexity?

- sometimes original definition
 - usually $\nabla^2 f \succcurlyeq 0$.
 - sometimes separation!
-

e.g. Logistic regression.

$$\min \frac{1}{m} \sum_{i=1}^m \begin{cases} 1 & \text{if } (\theta^T x^i) y^i < 0 \\ 0 & \text{otherwise} \end{cases} + \gamma \|\theta\|^2$$

First some examples: $\{x : f(x) \leq t\}$

$$x: Ax \geq b \quad \left| \begin{array}{l} \forall v \quad v^T A v \geq 0 \\ \lambda_{\min}(A) \geq 0 \end{array} \right.$$

$$x: x^T A x \leq 1 \quad A \succcurlyeq 0 \quad \left| \begin{array}{l} \det(\text{principal minor}) \geq 0 \end{array} \right.$$

$$\|x\|_p \leq 1 \quad p \geq 1$$

$$\|x\|_K = \{\inf t : t x \in K\}$$

$$x \in K \quad \left| \begin{array}{l} \|x\|_p \\ e^x, x^\alpha \quad \alpha > 1 \\ -\log x, x \log x \end{array} \right.$$

Log. Reg.

$\|x\|_p$ is convex.

$$F(a) = \sum_{i=1}^m f(a^T x^i)$$

$$\nabla_a F(a) = \sum_{i=1}^m f'(a^T x^i) a$$

$$\nabla^2 F(a) = \sum_{i=1}^m f''(a^T x^i) a a^T$$

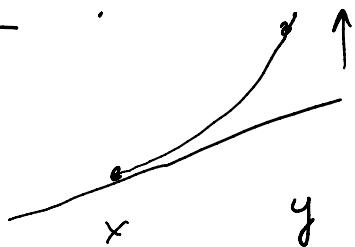
So it suffices to show $f'' \geq 0$.

$$f(z) = \log(1 + e^z)$$

$$f'(z) = \frac{e^z}{1 + e^z} = 1 - \frac{1}{1 + e^z} \quad f''(z) = \frac{1 \cdot e^z}{(1 + e^z)^2} \geq 0.$$

convex

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



Strongly convex

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

$$\text{Thm. } f(x^k) - f(x^*) \leq \left(1 - \frac{\mu}{L}\right)^k (f(x^0) - f(x^*))$$

PF

As before

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

also

$$f(x^*) - f(x^k) \geq \nabla f(x^k)^T (x^* - x^k) + \frac{\gamma}{2} \|x^* - x^k\|^2$$

or

$$f(x^k) - f(x^*) \leq \nabla f(x^k)^T (x^k - x^*) - \frac{\gamma}{2} \|x^* - x^k\|^2$$

arg $\max_{\Delta} \nabla f(x^k)^T \Delta - \frac{\gamma}{2} \|\Delta\|^2 \Rightarrow \Delta = \frac{\nabla f(x^k)}{\gamma}$.

and

$$\leq \frac{1}{2\gamma} \|\nabla f(x^k)\|^2.$$

$$f(x^{k+1}) - f(x^*) \leq f(x^k) - f(x^*) \left(1 - \frac{\gamma}{L}\right).$$

Cor. After $T = \frac{L}{\gamma} \log \left(\frac{f(x^0) - f(x^*)}{\epsilon} \right)$,

$$f(x^T) \leq f(x^*) + \epsilon$$

poly $(\frac{1}{\epsilon}) \rightarrow \log \frac{1}{\epsilon}$.