

No class tuesday 11/11
Exam Friday 11/14

CS7800: Advanced Algorithms

Class 18 : Convex Optimization

- Basic Concepts
- Gradient Descent

Jonathan Ullman

November 7, 2025

Exam 2:

- Linear programming
- Reductions
 - NP-completeness/hardness
 - Applications of max flow / min cut

Convex Optimization

Linear Programming

$$\begin{aligned} \min \quad & c^T x \\ \text{s.t.} \quad & Ax \geq b \end{aligned}$$

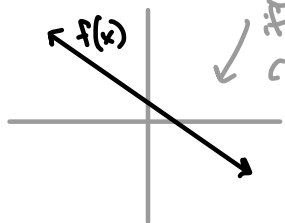
min/max interchangeable

linear function

linear constraints

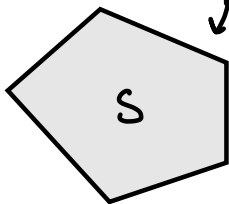
mm important

Objective



only interesting if there are constraints

Feasible region



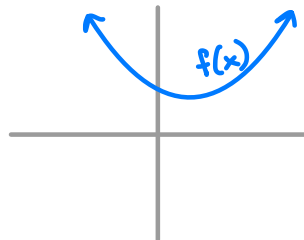
polytope

Convex Programming

$$\begin{aligned} \min \quad & f(x) \\ \text{s.t.} \quad & x \in S \end{aligned}$$

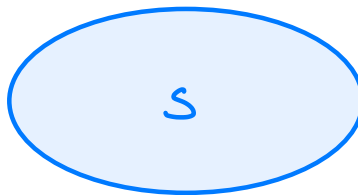
convex function

convex set



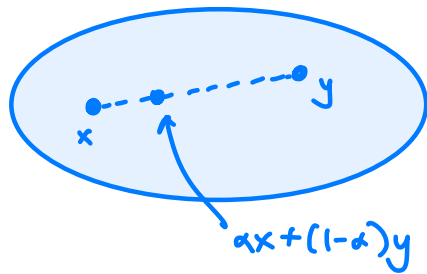
interesting even without constraints

convex set



Convex Sets

A set $S \subseteq \mathbb{R}^n$ is convex if for every $x, y \in S$ and $0 \leq \alpha \leq 1$ we have $\alpha x + (1-\alpha)y \in S$



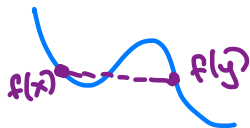
Examples:

- $S = \{x : \|x\| \leq B\}$ for any norm B feasibility easy
- $S = \{x : Ax \leq b \text{ and } x \geq 0 \text{ for any } A, b\}$ feasibility tricky

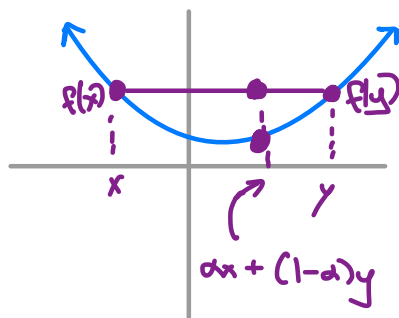
Convex Functions

A function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is convex if for every $x, y \in \mathbb{R}^n$ and $0 \leq \alpha \leq 1$ we have $f(\alpha x + (1-\alpha)y) \leq \alpha f(x) + (1-\alpha)f(y)$

not convex

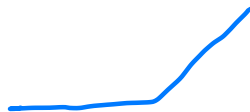


convex



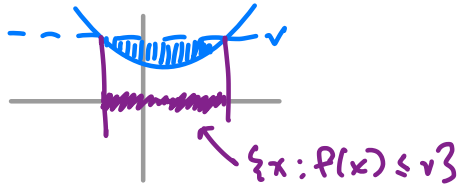
Examples:

- $f(x) = \|x\|$ for any norm
- $\|Ax - b\|_2^2 \leftarrow$ ordinary least squares regression
- Linear functions
- $\max \{0, x\} \leftarrow$ called ReLU or hinge loss

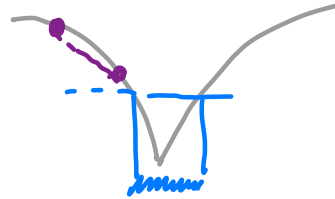


Properties of Convex Functions

If f is a convex function then $\{x: f(x) \leq v\}$ is a convex set



Note the reverse is not true!



Properties of Convex Functions

If f is a differentiable convex function then for any $x, y \in \mathbb{R}^n$

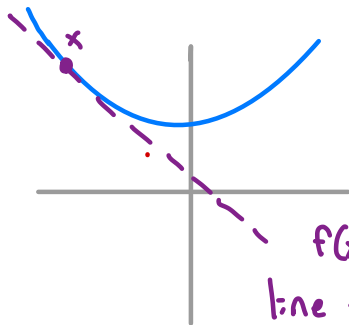
$$f(y) \geq \underbrace{f(x) + \langle \underbrace{\nabla f(x)}_{\text{gradient of } f \text{ at } x}, y - x \rangle}_{\text{tangent line to } f \text{ at } x}$$

Gradient operator

For a differentiable $f: \mathbb{R}^n \rightarrow \mathbb{R}$

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

One dimensional intuition



$$f(x) = ax^2 + bx + c$$

$$f'(x) = 2ax + b$$

$$f(x) + f'(x) \cdot (y - x) = l(y)$$

line tangent to f at x

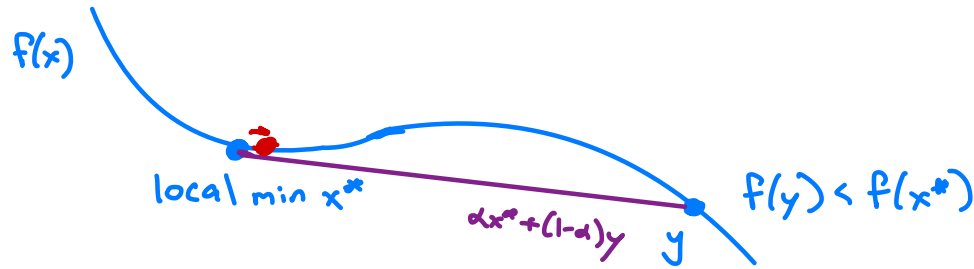
$$f(x) = x^T A x + b^T x + c$$

$$\nabla f(x) = 2Ax + b^T$$

Properties of Convex Functions

Thm: Any local minimum of a convex function is a global minimum

Proof:



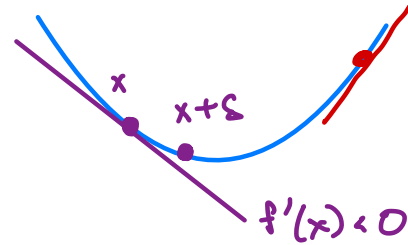
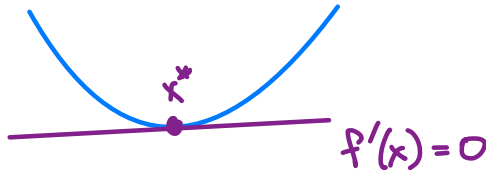
Properties of Convex Functions



Thm: If f is differentiable and convex then x^* is a global minimum if and only if $\langle \nabla f(x^*), x - x^* \rangle \geq 0$ for every $x \in S$

↑ If $S = \mathbb{R}^n$ then this is $\nabla f(x) = 0$

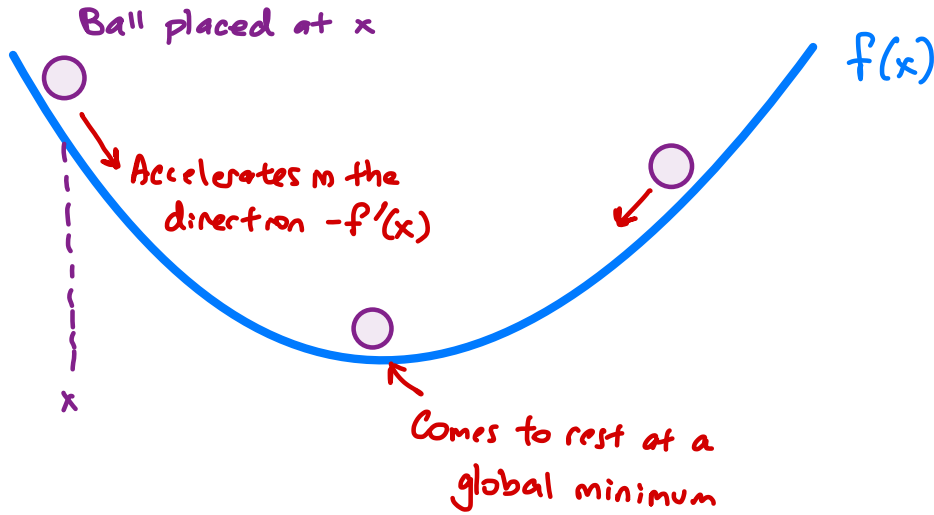
One dimensional picture



If $S = \mathbb{R}^n$ then $y = -\nabla f(x)$

Theorem \Rightarrow If x is not a global minimum then there exists a point y such that $\langle \nabla f(x), y - x \rangle < 0$, so moving in the direction $y - x$ will decrease the function

Convex Optimization Intuition



In \mathbb{R}^n , ball would accelerate in the direction $-\nabla f(x)$
and come to rest where $\nabla f(x) = 0$

First Order Convex Optimization

Linear Programming

$$\begin{array}{ll}\min & c^T x \\ & Ax \geq b\end{array}$$

Input is: $A \in \mathbb{R}^{n \times n}$
 $b \in \mathbb{R}^n$
 $c \in \mathbb{R}^n$

Convex Programming

$$\begin{array}{ll}\min & f(x) \\ & x \in S\end{array}$$

No general compact way
to describe $f(x)$

Input is: An oracle that takes
 $x \in \mathbb{R}^n$ and returns $\nabla f(x)$

Gradient Descent

Simplest setting

- $f: \mathbb{R}^n \rightarrow \mathbb{R}$ differentiable
- $S = \mathbb{R}^n$ (unconstrained)

Choosing step size: need to balance making progress against "overshooting" the minimizer

$$\textcircled{1} \underset{\eta}{\operatorname{argmin}} f(x^{(t-1)} - \eta \nabla f(x^{(t-1)}))$$

$$\textcircled{2} \eta_t \approx 1/t$$

$$\textcircled{3} \underline{\underline{\eta_t = 1/\sqrt{t}}}$$

Gradient descent

Initialize $x^{(0)}$

For $t = 1, \dots, T$:

Choose step size η_t

$$\text{let } x^{(t)} = x^{(t-1)} - \eta_t \nabla f(x^{(t-1)})$$

$$\text{Return } \frac{1}{T} \sum_{t=1}^T x^{(t)}$$

In practice return $x^{(\tau)}$

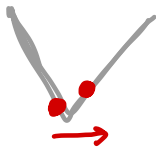
Analyzing Gradient Descent

→ A potential function

How do we keep track of progress? Two natural choices:

- ① (Function value) $\Phi_t = f(x^{(t)}) - f(x^*)$
 - ② (Distance to optimality) $\Phi_t = \|x^{(t)} - x^*\|_2^2$
- } Also combinations thereof

Analysis outline:



- Bound the decrease in potential $\Phi_{t-1} - \Phi_t \geq B_t$
- Telescoping sum gives $\Phi_T \leq \Phi_0 - \sum_{t=1}^T B_t$

Analyzing Gradient Descent

$$\Phi_t - \Phi_{t-1}$$

$$\|v\|^2 = \langle v, v \rangle$$

$$= \frac{1}{2\eta} \left[\|x^{(t)} - x^*\|^2 - \|x^{(t-1)} - x^*\|^2 \right]$$

$$= \frac{1}{2\eta} \left[\langle x^{(t)} - x^*, x^{(t)} - x^* \rangle - \langle x^{(t-1)} - x^*, x^{(t-1)} - x^* \rangle \right]$$

$$= \frac{1}{2\eta} \left[\langle x^{(t)} - x^{(t-1)}, x^{(t)} + x^{(t-1)} - 2x^* \rangle \right] \quad x^{(t)} - x^{(t-1)} = -\eta \nabla f(x^{(t-1)})$$

$$= \frac{1}{2\eta} \left[\langle -\eta \nabla f(x^{(t-1)}), -\eta \nabla f(x^{(t-1)}) + 2x^{(t-1)} - 2x^* \rangle \right]$$

$$= \frac{1}{2\eta} \left[\eta^2 \|\nabla f(x^{(t-1)})\|^2 + 2\eta \langle \nabla f(x^{(t-1)}), x^* - x^{(t-1)} \rangle \right]$$

Function f

$$\text{Iterates } x^{(t)} = x^{(t-1)} - \eta \nabla f(x^{(t-1)})$$

$$\text{Potential } \underline{\Phi_t} = \frac{1}{2\eta} \underline{\|x^{(t)} - x^*\|_2^2}$$

Analyzing Gradient Descent

$$\begin{aligned} & \Phi_t - \Phi_{t-1} \\ &= \frac{1}{2\eta} \left[\eta^2 \|\nabla f(x^{(t-1)})\|^2 + 2\eta \langle \nabla f(x^{(t-1)}), x^* - x^{(t-1)} \rangle \right] \\ &= \frac{\eta}{2} \|\nabla f(x^{(t-1)})\|^2 + \underbrace{\langle \nabla f(x^{(t-1)}), x^* - x^{(t-1)} \rangle}_{f(x^*) - f(x^{(t-1)}) \geq \langle \nabla f(x^{(t-1)}), x^* - x^{(t-1)} \rangle} \\ &\leq \underbrace{\frac{\eta}{2} \|\nabla f(x^{(t-1)})\|^2}_{\text{Assume } \|\nabla f\|^2 \leq G^2} + \underbrace{f(x^*) - f(x^{(t-1)})}_{\leq 0} \end{aligned}$$

$$\begin{aligned} \Phi_T - \Phi_0 &= (\Phi_T - \Phi_{T-1}) + (\Phi_{T-1} - \Phi_{T-2}) + \dots \\ &\leq \frac{T\eta G^2}{2} + \sum_{t=1}^T f(x^*) - f(x^{(t-1)}) \end{aligned}$$

Analyzing Gradient Descent

$$\begin{aligned}\Phi_T - \Phi_0 &= (\Phi_T - \Phi_{T-1}) + (\Phi_{T-1} - \Phi_{T-2}) + \dots \\ &\leq \frac{\eta G^2}{2} + \sum_{t=1}^T f(x^*) - f(x^{(t-1)})\end{aligned}$$

$$\bar{x} = \frac{1}{T} \sum_{t=1}^T x^{(t)}$$

convexity

$$f(\bar{x}) - f(x^*) \leq \frac{1}{T} \sum_{t=1}^T f(x^{(t-1)}) - f(x^*) \leq \frac{\eta G^2}{2} + \frac{\Phi_0}{T} - \frac{\Phi_T}{T} \leq \frac{\eta G^2}{2} + \frac{\Phi_0}{T}$$

$$f(\bar{x}) - f(x^*) \leq \frac{\eta G^2}{2} + \frac{\|x^{(0)} - x^*\|^2}{2\eta T} \quad \text{Assume } R^2 \geq \|x^{(0)} - x^*\|^2$$

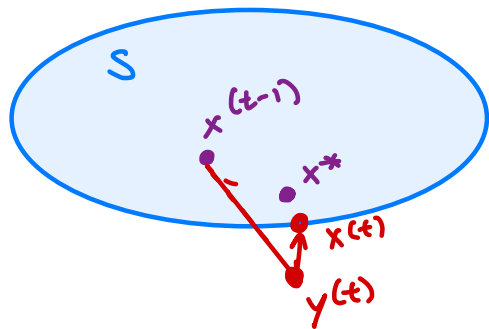
$$\leq \frac{\eta G^2}{2} + \frac{R^2}{2\eta T} \quad \text{set } \eta \text{ optimally}$$

$$f(\bar{x}) - f(x^*) \leq \frac{RG}{\sqrt{T}}$$

Constrained Optimization

Simplest setting

- $f: \mathbb{R}^n \rightarrow \mathbb{R}$ differentiable
- $S \subseteq \mathbb{R}^n$



Ex: $S = \{x: \|x\|_2 = 1\}$

$$x^{(t)} = \frac{y^{(t)}}{\|y^{(t)}\|}$$

Thm: $\|x^{(t)} - x^*\|_2 \leq \|y^{(t)} - x^*\|_2$
"Projection decreases distance"

Projected Gradient Descent

Initialize $x^{(0)}$

For $t = 1, \dots, T$:

choose step size η_t

$$\text{let } y^{(t)} = x^{(t-1)} - \eta_t \nabla f(x^{(t-1)})$$

$$\text{let } x^{(t)} = \underset{x \in S}{\operatorname{argmin}} \|x - y^{(t)}\|_2$$

Return $\frac{1}{T} \sum_{t=1}^T x^{(t)}$

"Projection" can often
be computed efficiently

Stochastic Optimization

Often in machine learning and statistics

$$f(x) = \frac{1}{n} \sum_{i=1}^n \underbrace{\ell_{\tilde{z}_i}}_{\substack{\text{error of model} \\ \text{on example } \tilde{z}_i}}(\tilde{x})$$

model parameters

Computing $\nabla f(x)$ is expensive, so we use
a single $\nabla \ell_{\tilde{z}}(x)$ for one random \tilde{z}
a stochastic gradient

Fact:

$$\mathbb{E}_{\tilde{z}}[\nabla \ell_{\tilde{z}}(x)] = \nabla f(x)$$