ECE 595: Machine Learning I Lecture 04 Intro to Optimization

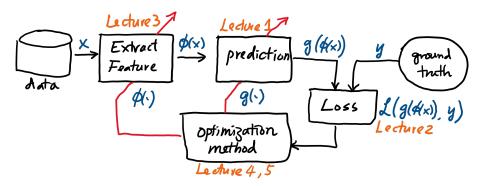
Spring 2020

Stanley Chan

School of Electrical and Computer Engineering Purdue University



Outline



Outline

Mathematical Background

- Lecture 4: Intro to Optimization
- Lecture 5: Gradient Descent

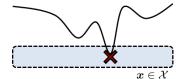
Lecture 4: Intro to Optimization

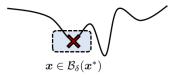
- Unconstrained Optimization
 - First Order Optimality
 - Second Order Optimality
- Convexity
 - What is convexity?
 - Convex optimization
- Constrained Optimization
 - Lagrangian
 - Examples

Unconstrained Optimization

$$\underset{\mathbf{x} \in \mathcal{X}}{\text{minimize}} f(\mathbf{x})$$

- $x^* \in \mathcal{X}$ is a global minimizer if
 - $f(x^*) \le f(x)$ for any $x \in \mathcal{X}$
- $x^* \in \mathcal{X}$ is a **local minimizer** if
 - $f(\mathbf{x}^*) \leq f(\mathbf{x})$, for any \mathbf{x} in a neighborhood $\mathcal{B}_{\delta}(\mathbf{x}^*)$
 - $\mathcal{B}_{\delta}(\mathbf{x}^*) = \{\mathbf{x} \mid ||\mathbf{x} \mathbf{x}^*||_2 \leq \delta\}$





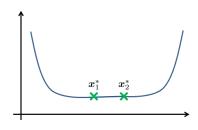
Uniqueness of Global Minimizer

If x^* is global minimizer, then

- Objective value $f(x^*)$ is unique
- Solution **x*** is not necessarily unique

Therefore:

- Suppose $f(\mathbf{x}) = g(\mathbf{x}) + \lambda ||\mathbf{x}||_1$ for some convex g.
- "minimize f(x)" has a global optimal $f(x^*)$.
- But there could be multiple x*'s.
- Some x^* maybe better, but not in the sense of f(x).



First and Second Order Optimality

$$\underbrace{\nabla f(\mathbf{x}^*) = \mathbf{0}}_{\text{First order condition}} \quad \text{and} \quad$$

$$\nabla^2 f(\mathbf{x}^*) \succeq 0$$

Necessary Condition:

If x^* is a global (or local) minimizer, then

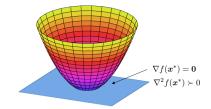
- $\nabla f(\mathbf{x}^*) = \mathbf{0}$.
- $\nabla^2 f(\mathbf{x}^*) \succeq 0$.

Sufficient Condition:

If x* satisfies

- $\nabla^2 f(\mathbf{x}^*) \succ 0$.

then x^* is a global (or local) minimizer.



Why? First Order

- Why is $\nabla f(\mathbf{x}^*) = \mathbf{0}$ necessary?
- Suppose **x*** is the minimizer.
- ullet Pick any direction $oldsymbol{d}$, and any step size ϵ . Then

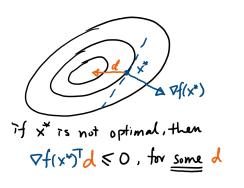
$$f(\mathbf{x}^* + \epsilon \mathbf{d}) = f(\mathbf{x}^*) + \epsilon \nabla f(\mathbf{x}^*)^T \mathbf{d} + \mathcal{O}(\epsilon^2).$$

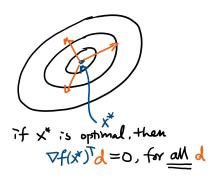
• Rearranging the terms yields

$$\underbrace{\lim_{\epsilon \to 0} \left\{ \frac{f(\mathbf{x}^* + \epsilon \mathbf{d}) - f(\mathbf{x}^*)}{\epsilon} \right\}}_{>0, \forall \mathbf{d}} = \nabla f(\mathbf{x}^*)^T \mathbf{d}.$$

• So $\nabla f(\mathbf{x}^*)^T \mathbf{d} \geq 0$ for all \mathbf{d} . True only when $\nabla f(\mathbf{x}^*) = \mathbf{0}$.

First Order Condition Illustrated





Why? Second Order

Do third order approximation:

$$f(\mathbf{x}^* + \epsilon \mathbf{d}) = f(\mathbf{x}^*) + \epsilon \underbrace{\nabla f(\mathbf{x}^*)^T \mathbf{d}}_{=0} + \frac{\epsilon^2}{2} \mathbf{d}^T \nabla^2 f(\mathbf{x}^*) \mathbf{d} + \frac{\epsilon^3}{6} \mathcal{O}(\|\mathbf{d}\|^3)$$

Therefore,

$$\frac{1}{\epsilon^2} \left[f(\mathbf{x}^* + \epsilon \mathbf{d}) - f(\mathbf{x}^*) \right] = \frac{1}{2} \mathbf{d}^T \nabla^2 f(\mathbf{x}^*) \mathbf{d} + \left[\frac{\epsilon}{6} \mathcal{O}(\|\mathbf{d}\|^3) \right]$$

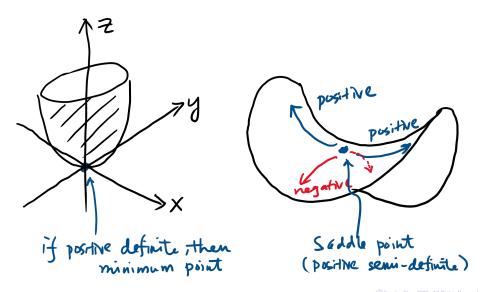
$$\lim_{\epsilon \to 0} \frac{1}{\epsilon^2} \underbrace{\left[f(\mathbf{x}^* + \epsilon \mathbf{d}) - f(\mathbf{x}^*) \right]}_{\geq 0} = \frac{1}{2} \mathbf{d}^T \nabla^2 f(\mathbf{x}^*) \mathbf{d} + \lim_{\epsilon \to 0} \left[\frac{\epsilon}{6} \mathcal{O}(\|\mathbf{d}\|^3) \right],$$

Hence,

$$\frac{1}{2}\boldsymbol{d}^{\mathsf{T}}\nabla^{2}f(\boldsymbol{x}^{*})\boldsymbol{d}\geq0,\qquad\forall\boldsymbol{d}$$

⇒ positive semi-definite!

Second Order Condition Illustrated



Outline

Mathematical Background

- Lecture 4: Intro to Optimization
- Lecture 5: Gradient Descent

Lecture 4: Intro to Optimization

- Unconstrained Optimization
 - First Order Optimality
 - Second Order Optimality
- Convexity
 - What is convexity?
 - Convex optimization
- Constrained Optimization
 - Lagrangian
 - Examples

Most Optimization Problems are Not Easy

Minimize the log-sum-exp function:

$$f(\mathbf{x}) = \log \left(\sum_{i=1}^{m} \exp(\mathbf{a}_{i}^{T} \mathbf{x} + b_{i}) \right)$$

• Gradient is (exercise)

$$\nabla f(\mathbf{x}^*) = \frac{1}{\sum_{j=1}^m \exp(\mathbf{a}_i^T \mathbf{x}^* + b_j)} \sum_{i=1}^m \exp(\mathbf{a}_i^T \mathbf{x}^* + b_i) \mathbf{a}_i.$$

- Non-linear equation. No closed-form solution.
- Need iterative algorithms, e.g., gradient descent.
- Or off-the-shelf optimization solver, e.g., CVX.

CVX Demonstration

- Disciplined optimization: It translates the problem for you.
- Developed by S. Boyd and colleagues (Stanford).
- E.g., Minimize $f(\mathbf{x}) = \log \left(\sum_{i=1}^{n} \exp(\mathbf{a}_{i}^{T} \mathbf{x} + b_{i}) \right) + \lambda \|\mathbf{x}\|^{2}$.

```
import cvxpy as cp
import numpy as np
n = 100
d = 3
A = np.random.randn(n, d)
b = np.random.randn(n)
lambda = 0.1
x = cp.Variable(d)
objective = cp.Minimize(cp.log_sum_exp(A*x - b) + lambda_*cp.sum_squares(x))
constraints = []
prob = cp.Problem(objective, constraints)
optimal_objective_value = prob.solve()
print(optimal_objective_value)
print(x.value)
```

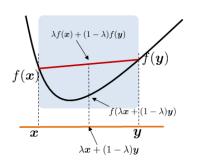
Convex Function

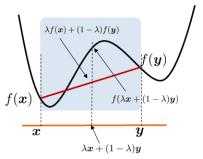
Definition

Let $x \in \mathcal{X}$ and $y \in \mathcal{X}$. Let $0 \le \lambda \le 1$. A function $f : \mathbb{R}^n \to \mathbb{R}$ is **convex** over \mathcal{X} if

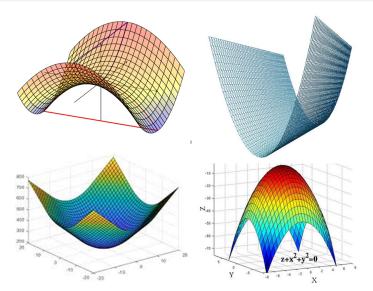
$$f(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda)f(\mathbf{y}).$$

The function is called strictly convex if " \leq " is replaced by "<".





Example: Which one is convex?



Verifying Convexity

Any of the following conditions is **necessary** and **sufficient** for convexity:

By definition:

$$f(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y}) \le \lambda f(\mathbf{x}) + (1 - \lambda)f(\mathbf{y}).$$

- Function value is lower than the line.
- First Order Convexity:

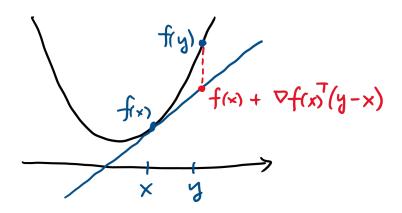
$$f(\mathbf{y}) \ge f(\mathbf{x}) + \nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x}), \quad \forall \mathbf{x}, \mathbf{y} \in \mathcal{X}.$$

- Tangent line is always lower than the function
- **3** Second Order Convexity: f is convex over \mathcal{X} if and only if

$$\nabla^2 f(\mathbf{x}) \succeq 0 \quad \forall \mathbf{x} \in \mathcal{X}.$$

• Curvature is positive.

Tangent Line Condition Illustrated



Outline

Mathematical Background

- Lecture 4: Intro to Optimization
- Lecture 5: Gradient Descent

Lecture 4: Intro to Optimization

- Unconstrained Optimization
 - First Order Optimality
 - Second Order Optimality
- Convexity
 - What is convexity?
 - Convex optimization
- Constrained Optimization
 - Lagrangian
 - Examples

Constrained Optimization

Equality Constrained Optimization:

Requires a function: Lagrangian function

$$\mathcal{L}(\mathbf{x}, \mathbf{\nu}) \stackrel{\mathsf{def}}{=} f(\mathbf{x}) - \sum_{j=1}^{k} \nu_j h_j(\mathbf{x}).$$

 $\nu = [\nu_1, \dots, \nu_k]$: Lagrange multipliers or the dual variables.

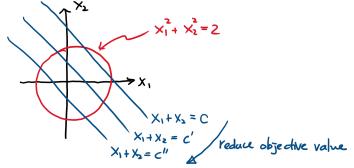
Solution $(\mathbf{x}^*, \mathbf{\nu}^*)$ satisfies

$$abla_{\mathbf{x}}\mathcal{L}(\mathbf{x}^*, \mathbf{\nu}^*) = \mathbf{0},
abla_{\mathbf{\nu}}\mathcal{L}(\mathbf{x}^*, \mathbf{\nu}^*) = \mathbf{0}.
abla_{\mathbf{\nu}}\mathcal{L}(\mathbf{x}^*, \mathbf{\nu}^*) = \mathbf{0}.$$

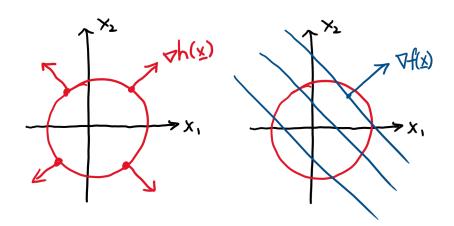
Example: Illustrating Lagrangian

Consider the problem

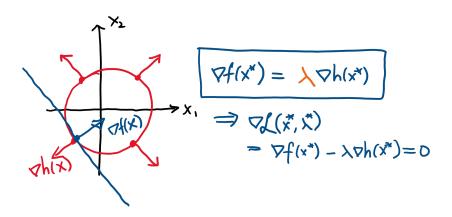
• Minimizer is x = (-1, -1).



Example: Illustrating Lagrangian



Example: Illustrating Lagrangian



Example: ℓ_2 -minimization with constraint

minimize
$$\frac{1}{2} \| \mathbf{x} - \mathbf{x}_0 \|^2$$
, subject to $\mathbf{A}\mathbf{x} = \mathbf{y}$.

The Lagrangian function of the problem is

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\nu}) = \frac{1}{2} \|\mathbf{x} - \mathbf{x}_0\|^2 - \boldsymbol{\nu}^T (\mathbf{A}\mathbf{x} - \mathbf{y}).$$

The first order optimality condition requires

$$abla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \mathbf{\nu}) = (\mathbf{x} - \mathbf{x}_0) - \mathbf{A}^T \mathbf{\nu} = \mathbf{0}$$

$$abla_{\mathbf{\nu}} \mathcal{L}(\mathbf{x}, \mathbf{\nu}) = \mathbf{A}\mathbf{x} - \mathbf{y} = \mathbf{0}.$$

Multiply the first equation by **A** on both sides:

$$\begin{array}{cccc}
A(x - x_0) - AA^T \nu &= 0 \\
\Rightarrow & \underbrace{Ax} - Ax_0 &= AA^T \nu \\
\Rightarrow & y - Ax_0 &= AA^T \nu \\
\Rightarrow & (AA^T)^{-1} (y - Ax_0) &= \nu
\end{array}$$

Example: ℓ_2 -minimization with constraint

minimize
$$\frac{1}{2} \| \mathbf{x} - \mathbf{x}_0 \|^2$$
, subject to $\mathbf{A}\mathbf{x} = \mathbf{y}$.

The first order optimality condition requires

$$egin{aligned}
abla_{x}\mathcal{L}(x,
u) &= (x-x_0) - \mathbf{A}^{T}
u = \mathbf{0} \\
abla_{
u}\mathcal{L}(x,
u) &= \mathbf{A}x - \mathbf{y} = \mathbf{0}. \end{aligned}$$

We just showed: $\nu = (\mathbf{A}\mathbf{A}^T)^{-1}(\mathbf{y} - \mathbf{A}\mathbf{x}_0)$. Substituting this result into the first order optimality yields

$$\mathbf{x} = \mathbf{x}_0 + \mathbf{A}^T \mathbf{\nu}$$
$$= \mathbf{x}_0 + \mathbf{A}^T (\mathbf{A} \mathbf{A}^T)^{-1} (\mathbf{y} - \mathbf{A} \mathbf{x}_0)$$

Therefore, the solution is $\mathbf{x} = \mathbf{x}_0 + \mathbf{A}^T (\mathbf{A} \mathbf{A}^T)^{-1} (\mathbf{y} - \mathbf{A} \mathbf{x}_0)$.

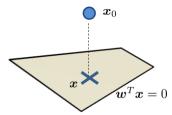
Special Case

$$\underset{\boldsymbol{x} \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{2} \|\boldsymbol{x} - \boldsymbol{x}_0\|^2, \quad \text{subject to} \quad \boldsymbol{A}\boldsymbol{x} = \boldsymbol{y}.$$

Special case: When $\mathbf{A}\mathbf{x} = \mathbf{y}$ is simplified to $\mathbf{w}^T \mathbf{x} = 0$.

- $\mathbf{w}^T \mathbf{x} = 0$ is a line.
- Find a point x on the line that is closest to x_0 .
- Solution is

$$x = x_0 + w(w^T w)^{-1}(0 - w^T x_0)$$
$$= x_0 - \left(\frac{w^T x_0}{\|w\|^2}\right)^T w.$$



In practice ...

- Use CVX to solve problem
- Here is a MATLAB code
- Exercise: Turn it into Python.

Reading List

Unconstrained Optimality Conditions

- Nocedal-Wright, Numerical Optimization. (Chapter 2.1)
- Boyd-Vandenberghe, Convex Optimization. (Chapter 9.1)

Convexity

- Nocedal-Wright, Numerical Optimization. (Chapter 1)
- Boyd-Vandenberghe, Convex Optimization. (Chapter 2 and 3)
- CMU, Convex Optimization (Lecture 2 and 4)
 https://www.stat.cmu.edu/~ryantibs/convexopt-F18/
- Stanford CS 229 (Tutorial)
 http://cs229.stanford.edu/section/cs229-cvxopt.pdf
- UCSD ECE 273 (Tutorial)
 http://eceweb.ucsd.edu/~gert/ECE273/CvxOptTutPaper.pdf

Constrained Optimization

Nocedal-Wright, Numerical Optimization. (Chapter 12.1)