

CSED490Y: Optimization for Machine Learning

Week 04-1: Gradient descent

Namhoon Lee

POSTECH

Spring 2022

Warnings based on the current progress:

1. Students with no contact
2. Single member group
3. Group without project topic
4. Group with project out of context

- o example topics? → see "S03-1"
- o TA email addresses?
 - ↳ PLMS
 - ↳ Op+ML
 - ↳ slides "S01-1"

Warnings based on the current progress:

1. Students with no contact
2. Single member group
3. Group without project topic
4. Group with project out of context

Topic W6 Monday
CWS ↓↓
28 March
11:59 PM

Team up (Due-extended: 11:59PM on Monday 21 March): WS

- ▶ Form a group of up to 3 members.
- ▶ Email TA about your group members by the due date.
- ▶ You may use the discussion board on PLMS to find your teammates.

New requirements:

- ▶ **Avoid late submission:** if you miss the due date, you will receive a penalty of 10% of the total marks. $\boxed{-4} / 40$
- ▶ **Submit a (self-)plagiarism statement:** “I certify that this project is entirely my own work, and I have not previously worked, am currently working or planning to work on any aspect of this course project ..” – TA will soon send out a form to sign.

↳ prevent re-using prev materials

↳ ensure a fair assessment.

↑
G

New requirements:

- ▶ **Avoid late submission:** if you miss the due date, you will receive a penalty of 10% of the total marks.
- ▶ **Submit a (self-)plagiarism statement:** “I certify that this project is entirely my own work, and I have not previously worked, am currently working or planning to work on any aspect of this course project ..” – TA will soon send out a form to sign.

✦ Important:

- ▶ Make sure to check important dates for project.
- ▶ This course does not require any assignments other than project.
- ▶ If there is anything you are not sure of about the project, come talk to me.

Convex function

$$f(y) \geq f(x) + \langle \sigma f(x), y-x \rangle$$

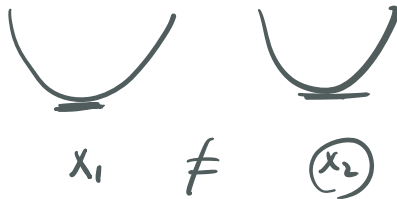
differentiable then

opt. condition

More on convex functions..

- ▶ Notice from C^1 definition that $\nabla f(x) = 0$ implies $f(y) \geq f(x)$ for all y , so x is a global minimizer; this further explains why least squares can be solved by setting the derivative equal to zero.
- ▶ (Strictly) convex function have at most one global minimum; w and v can't both be global minima if $w \neq v$; it would imply convex combinations u of w and v would have $f(u)$ below the global minimum.

\Rightarrow sol'n is unique.



Convex function

For strictly convex objective f there can be at most one global optimum.

Convex function

For strictly convex objective f there can be at most one global optimum.

Proof:



$$x^* \neq x^\#$$

1. Suppose x^* is a local minimum and also there exists another local minimum $x^\#$ ($\neq x^*$).

Convex function

For strictly convex objective f there can be at most one global optimum.

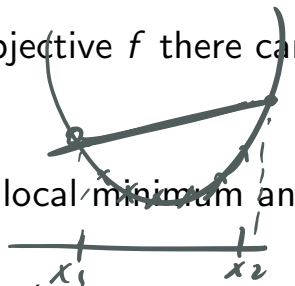
Proof: f convex \rightarrow local \rightarrow global min.

1. Suppose x^* is a local minimum and also there exists another local minimum $x^\#$ ($\neq x^*$).
2. Since f is convex (because it is strictly convex), $f(x^*)$ and $f(x^\#)$ are both global minima, and $f(x^*) = f(x^\#)$.

Convex function

For strictly convex objective f there can be at most one global optimum.

Proof:



strictly convex

$$\boxed{f(\theta x_1 + (1-\theta)x_2)} < \theta f(x_1) + (1-\theta)f(x_2)$$

1. Suppose x^* is a local minimum and also there exists another local minimum $x^\#$ ($\neq x^*$).
2. Since f is convex (because it is strictly convex), $f(x^*)$ and $f(x^\#)$ are both global minima, and $\boxed{f(x^*) = f(x^\#)}$.
3. The C^0 definition for $\underline{y = \theta x^* + (1-\theta)x^\#}$, i.e.,

$$\underline{f(y)} < \theta f(x^*) + (1-\theta)\boxed{f(x^\#)} = \theta f(x^*) + (1-\theta)\boxed{f(x^*)} = \underline{f(x^*)}$$

contradicts that x^* is a global minimum.

$$\underline{f(y)} < \underline{f(x^*)} \rightarrow$$

Convex function

For strictly convex objective f there can be at most one global optimum.

Proof:

1. Suppose x^* is a local minimum and also there exists another local minimum $x^\#$ ($\neq x^*$).
2. Since f is convex (because it is strictly convex), $f(x^*)$ and $f(x^\#)$ are both global minima, and $f(x^*) = f(x^\#)$.
3. The C^0 definition for $y = \theta x^* + (1 - \theta)x^\#$, i.e.,

$$f(y) < \theta f(x^*) + (1 - \theta)f(x^\#) = \theta f(x^*) + (1 - \theta)f(x^*) = f(x^*)$$

contradicts that x^* is a global minimum.

- ✓ 4. This means that for $x^\#$ to be a local minimum, it must be that $x^\# = x^*$.

Gradient descent

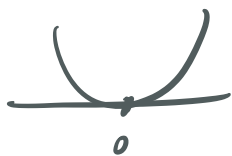
W2-2

prob. convex

Least squares :

$$\min_{w \in \mathbb{R}^d} f(w) := \frac{1}{2} \|Xw - y\|^2$$

$$f(x) = x^2$$



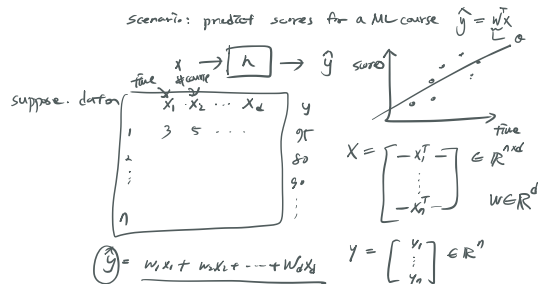
$$f'(x) = 2x = 0$$

$$x = 0$$

$\Rightarrow w^*$ s.t. $\nabla f(w^*) = 0$ global min.

$$\nabla f(w) = X^T(Xw - y) = 0 \Rightarrow w^* = (X^T X)^{-1} X^T y$$

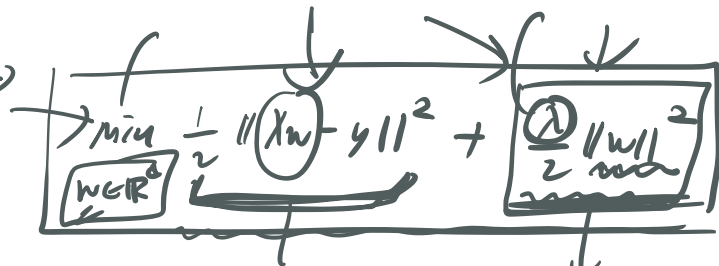
$$\nabla^2 f(w) = X^T X \succeq 0 \quad \text{p.s.d} \Rightarrow \text{by } \epsilon \text{ det.} \Rightarrow \underline{f \text{ convex}}$$



Gradient descent

⊕ Cross-validation

(L2-regularized) Least squares



$$\nabla f(w) = X^T(Xw - y) + \lambda w$$

$$\nabla^2 f(w) = X^T X + \lambda I \succeq 0$$

$$\begin{aligned} \underbrace{v^T (X^T X + \lambda I) v}_{\forall v \in \mathbb{R}^d, v \neq 0} &= v^T X^T X v + \lambda v^T v \\ &= \underbrace{\|Xv\|^2}_{\geq 0} + \lambda \underbrace{\|v\|^2}_{> 0} > 0 \end{aligned}$$

unique solution

Gradient descent

Linear system

Cost of solving (regularized) least squares \Rightarrow

$$\frac{O(nd^2 + \underline{d^3})}{O(d^3)}$$

$$\nabla f(w) = X^T(Xw - y) + \lambda w = 0$$

$$w^* = \underbrace{(X^T X + \lambda I)^{-1}}_{nd^2} X^T y$$

nd^2

pseudo inverse

502-1

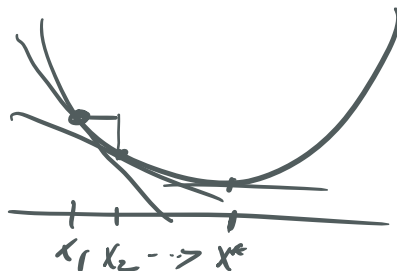
Gradient descent

$$x_{t+1} = x_t - \eta \nabla f(x_t)$$

(GTD)

Gradient descent algorithm

- ↳ for finding a minimum $f(x)$
- ↳ iteratively
- ↳ x_1 , initial random guess \rightarrow $f(x_1)$?



- ↳ update $x_1 \rightarrow x_2$ by taking a step into the neg. direction of the slope w/ step size η
- ↳ repeat until find soln.

Gradient descent

$$f(w) = \|xw - y\|^2$$

Numerical

$$f(x) = \|Ax - b\|^2 \quad (0)$$

Cost of solving (regularized) least squares with GD

$$x_{t+1} = x_t - \eta \nabla f(x_t)$$

$$= x_t - \eta \left(\cancel{x_t^T(x_t w - y)} + \lambda w \right) \\ \left(A^T(Ax_t - b) + \lambda x_t \right)$$

cost (1 step update) : $O(nd)$

cost for t steps : $O(nd \cdot t) + \lambda w$

Lin. alg. $\left\{ \begin{array}{l} - \text{Gauss} \\ - \text{QR} \end{array} \right.$

Analytical (x)

$$O(nd^2 + \boxed{d^3})$$

$$t < \max \left\{ d, \frac{d^2}{\eta} \right\}$$

\hookrightarrow GD is fast enough

Q: How fast is GD?

Lipschitz continuous objective gradients

A differentiable function f is called L -smooth if there exists an $L > 0$ such that the following satisfies:

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\| \quad \forall \{x, y\} .$$

Lipschitz continuous objective gradients

A differentiable function f is called L -smooth if there exists an $L > 0$ such that the following satisfies:

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\| \quad \forall \{x, y\} .$$

- ▶ Gradient does not change arbitrarily quickly.
- ▶ Intuitively, without it the gradient would not be useful to decrease f .
- ▶ It is essential for convergence analyses of most gradient based methods.
- ▶ This is a fairly weak assumption and holds true for most ML models (including neural networks with smooth activations).

An important consequence of Lipschitz continuous objective gradient:

$$f(y) \leq f(x) + \nabla f(x)^\top (y - x) + \frac{L}{2} \|y - x\|^2 .$$

Proof (recall ftc):

An important consequence of Lipschitz continuous objective gradient:

$$f(y) \leq f(x) + \nabla f(x)^\top (y - x) + \frac{L}{2} \|y - x\|^2 .$$

Illustration:

Gradient descent progress bound

Under the quadratic upper bound, we are interested in how much progress gradient descent can make at each step.

Gradient descent progress bound

Under the quadratic upper bound, we are interested in how much progress gradient descent can make at each step.

Consider gradient descent with $\eta = 1/L$.

$$x_{t+1} = x_t - \frac{1}{L} \nabla f(x_t) .$$

Gradient descent progress bound

Under the quadratic upper bound, we are interested in how much progress gradient descent can make at each step.

Consider gradient descent with $\eta = 1/L$.

$$x_{t+1} = x_t - \frac{1}{L} \nabla f(x_t) .$$

Plugging this into the bound gives

Convergence of gradient descent

From the C^1 definition of convex function, we can get

$$f(x_t) \leq f(x^*) + \nabla f(x_t)^\top (x_t - x^*)$$

Convergence of gradient descent

From the C^1 definition of convex function, we can get

$$f(x_t) \leq f(x^*) + \nabla f(x_t)^\top (x_t - x^*)$$

Plugging this into the progress bound we derived previously gives

Convergence of gradient descent

Thank you

Any questions?

A lot of material in this course is borrowed or derived from the following:

- ▶ Numerical Optimization, Jorge Nocedal and Stephen J. Wright.
- ▶ Convex Optimization, Stephen Boyd and Lieven Vandenberghe.
- ▶ Convex Optimization, Ryan Tibshirani.
- ▶ Optimization for Machine Learning, Martin Jaggi and Nicolas Flammarion.
- ▶ Optimization Algorithms, Constantine Caramanis.
- ▶ Advanced Machine Learning, Mark Schmidt.