

CSED490Y: Optimization for Machine Learning

Week 05-2: Subgradient and projected gradient methods

Namhoon Lee

POSTECH

Spring 2022

Gradient descent example

$$f'(x) = 6x + 4 = 0$$

$$x^* = -\frac{2}{3}$$

$$1+r+r^2+\dots+r^{t-1} = \frac{1-r^t}{1-r}$$

min

Solve $f(x) = 3x^2 + 4x - 2$ using GD

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

$$= \underline{x}_t - \eta (6\underline{x}_t + 4)$$

$$= (1-6\eta) \underline{x}_t - 4\eta$$

$$= (1-6\eta) (x_{t-1} - \eta \nabla f(x_{t-1})) - 4\eta$$

$$= (1-6\eta) (\underline{x}_{t-1} - \eta (6\underline{x}_{t-1} + 4)) - 4\eta$$

$$= (1-6\eta) ((1-6\eta) x_{t-1} - 4\eta) - 4\eta$$

$$= (1-6\eta)^2 \underline{x}_{t-1} - 4\eta(1-6\eta) - 4\eta$$

$$= (1-6\eta) \left(x_1 + \frac{2}{3} \right) - \frac{2}{3}$$

$$\dots = (1-6\eta)^t x_1 - 4\eta \frac{(1-6\eta)^t - 1}{1-(1-6\eta)}$$

$$= (1-6\eta)^t x_1 - 4\eta \cdot \frac{1 - (1-6\eta)^t}{1 - (1-6\eta)}$$

$$\underline{t} \rightarrow \infty \quad |1-6\eta| < 1$$

$$= \frac{-4\eta}{1-(1-6\eta)} = \boxed{-\frac{2}{3}}$$

Gradient descent example

From the previous example

▶ solution

$$-\frac{2}{3}$$

▶ convergence

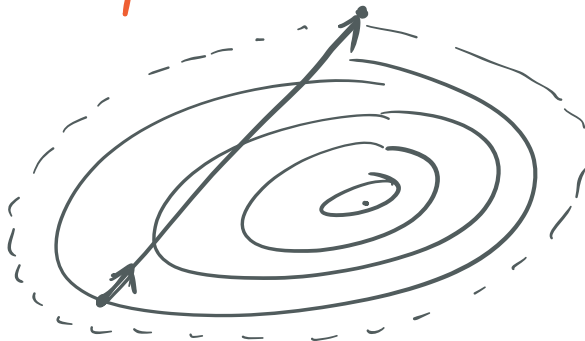
: linear : $\epsilon \sim \rho^t$ or e^{-ct}

▶ stepsize

$|1-6\eta| < 1$; stepsize small enough

▶ initial guess

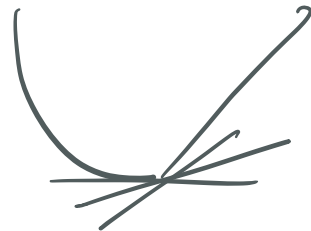
x_1



Subgradient and subdifferential

A vector g is called a subgradient of a convex function f at x if

$$\underline{f(y)} \geq f(x) + \langle \underline{g}, y - x \rangle, \quad \forall y.$$



Subgradient and subdifferential

A vector g is called a subgradient of a convex function f at x if

$$f(y) \geq f(x) + \langle g, y - x \rangle, \quad \forall y .$$

The set of subgradients of f at x is called subdifferential $\partial f(x)$.

Subgradient and subdifferential

A vector g is called a subgradient of a convex function f at x if

$$f(y) \geq f(x) + \langle g, y - x \rangle, \quad \forall y.$$

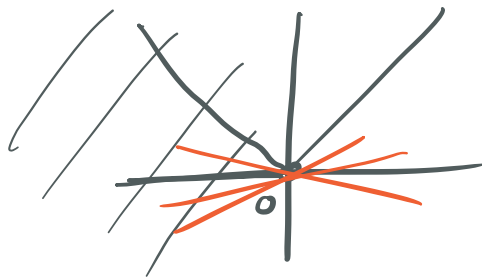
The set of subgradients of f at x is called subdifferential $\partial f(x)$.

- ▶ If a function f is differentiable at x , the gradient is the only element in the subdifferential $\partial f(x)$, *i.e.*, $g_x = \nabla f(x)$.
- ▶ The optimality condition for non-differentiable function: x^* is a global minimum if $0 \in \partial f(x^*)$.

Subdifferential example

An absolute value function: $f(x) = |x|$

$$\partial |x| = \begin{cases} -1, & x < 0 \\ [-1, 1], & x = 0 \\ +1, & x > 0 \end{cases}$$



Non-smooth problem

X, y : data set

$$\hat{y} = w^T x = \underline{w_1} x_1 + \overset{0}{\underline{w_2} x_2} + \dots + \overset{0}{\underline{w_d} x_d}$$

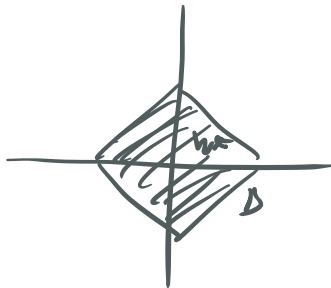
$$w^* = (\overset{d}{x \ x \ 0 \ 0 \ 0 \ x \ \dots \ 0})$$

L1 regularized least squares

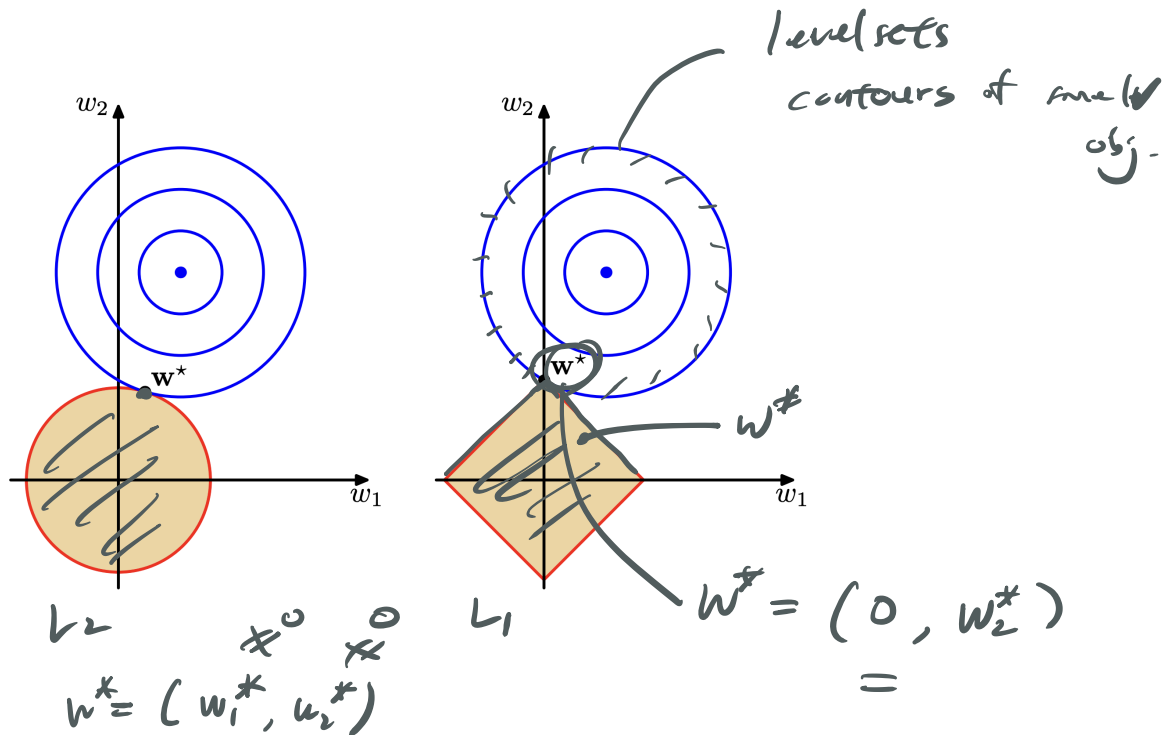
$$\min_{w \in \mathbb{R}^d} \frac{1}{2} \|Xw - y\|^2 + \lambda \|w\|_1$$

$$\Leftrightarrow \begin{cases} \min & \frac{1}{2} \|Xw - y\|^2 \\ \text{s.t.} & w \in \mathbb{R}^d \\ & \|w\|_1 \leq B \end{cases} \quad w^*$$

- ▶ You can obtain a sparse solution.
- ▶ It can be used for feature selection.



L1 vs L2 regularization



L1 vs L2 regularization

L2 regularized least squares $f(w)$

$$\min_w \frac{1}{2} \|Xw - y\|^2 + \frac{\lambda}{2} \|w\|_2^2$$

$$w_j^* = 0$$

$$\underline{\underline{w^T x}}$$

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

$$\nabla_j f(w) = \underbrace{X_j^T}_{\text{data features}} (\underbrace{Xw^* - y}_{\text{error, residual}}) + \lambda \underbrace{w_j^*}_{=0} = 0$$

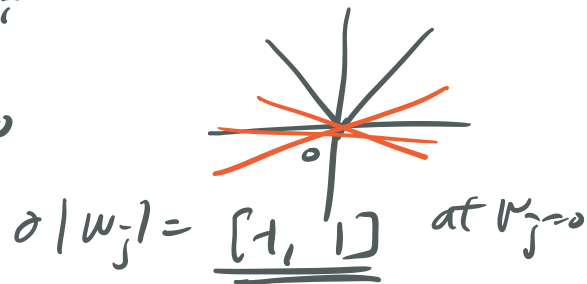
data features ? error, residual

L1 vs L2 regularization

L2 regularized least squares

$$\|w\|_1 = \sum_i |w_i|$$

$$w_j^* = 0$$



L1 regularized least squares

$$\min_w f(w) = \frac{1}{2} \|Xw - y\|^2 + \lambda \|w\|_1$$

$$\partial f(w) = X^T(Xw - y) + \lambda \partial \|w\|_1$$

$$0 \in \partial_j f(w) = \underline{X_j^T(Xw^* - y)} + \lambda [-1, 1]$$

$$\underline{-X_j^T(Xw^* - y)} \in \underline{\lambda [-1, 1]}$$

Subgradient method

Subgradient method:

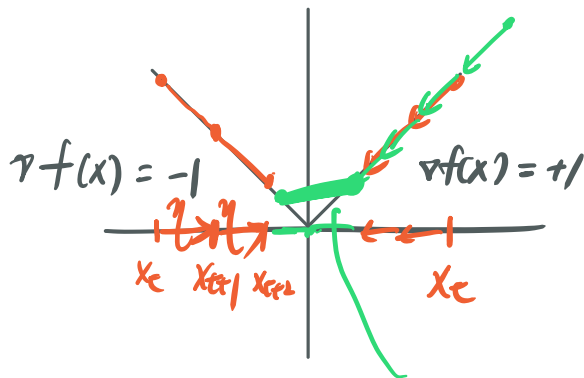
$$x_{t+1} = x_t - \eta \underline{g}_{x_t}$$

where $\underline{g}_{x_t} \in \partial f(x_t)$.

- ▶ Gradient descent for non-differentiable cases.
- ▶ Applicable to the previous example of absolute value function.

Subgradient method example

$$f(x) = |x|$$



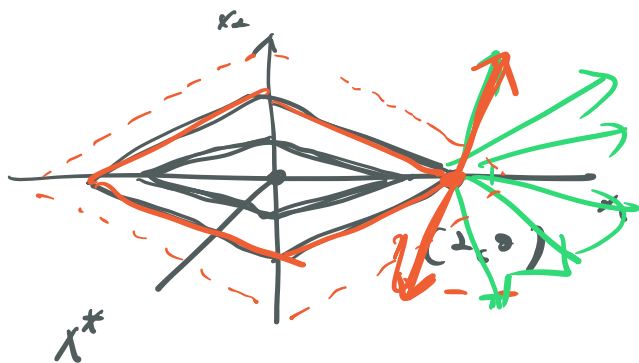
$$\underline{x} < 0 : \quad x_{t+1} = x_t + \eta$$

$$x > 0 : \quad x_{t+1} = x_t - \eta$$

Q: what happens when you get close to the min?

Subgradient method example

$$f(x_1, x_2) = |x_1| + 2|x_2|$$



$$\partial f(x_1, x_2) = \boxed{(1, 4[-1, 1])^T} \quad \underline{x_1 > 0, x_2 = 0}$$



Gradient descent vs subgradient method

Differences between gradient descent and subgradient method:

- ▶ Gradient descent improves at every iteration, unlike sub-gradient method.
- ▶ Gradient descent can take a big step size: self-tuning property.
- ▶ Gradient descent takes bigger steps when far away.

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.