

Non-smooth Optimization

Optimization Techniques (ENGG*6140)

School of Engineering,
University of Guelph, ON, Canada

Course Instructor: Benjamin Ghojogh
Winter 2023

\cup
 \cup
 \cup
 \cup
 \cup

∞
 \cup
 \rightarrow

$f \notin \cup$

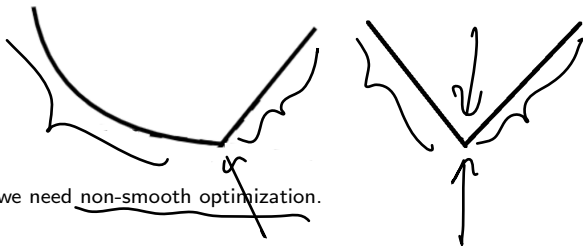
$f \in \cup$
 \downarrow
Gradient
Hessian

$f \notin \cup$

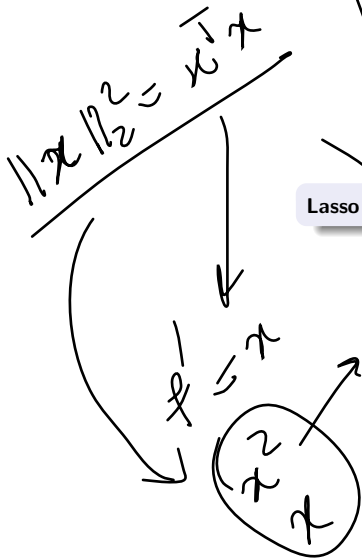
Non-smooth Function

Non-smooth function

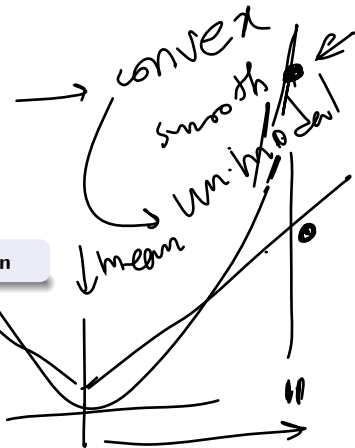
- When we have non-smooth function, the gradient is not defined at the non-smooth point(s).



- In these cases, we need non-smooth optimization.



Lasso Regularization



Lasso Regularization

- The ℓ_1 norm can be used for sparsity [1]. Sparsity is very useful and effective because of betting on sparsity principal [2] and the Occam's razor [3].
- If $\mathbf{x} = [x_1, \dots, x_d]^T$, for having sparsity, we should use subset selection for the regularization of a cost function $\Omega_0(\mathbf{x})$:

$$\underset{\mathbf{x}}{\text{minimize}} \quad \Omega(\mathbf{x}) := \Omega_0(\mathbf{x}) + \lambda \|\mathbf{x}\|_0, \quad \rightarrow \text{"}\ell_0\text{"} \quad (1)$$

where:

$$\|\mathbf{x}\|_0 := \sum_{j=1}^d \mathbb{I}(x_j \neq 0) = \begin{cases} 0 & \text{if } x_j = 0, \\ 1 & \text{if } x_j \neq 0, \end{cases} \quad (2)$$

is the " ℓ_0 " norm, which is not a norm (so we use " ℓ_0 " for it) because it does not satisfy the norm properties [4]. The " ℓ_0 " norm counts the number of non-zero elements so when we penalize it, it means that we want to have sparser solutions with many zero entries.

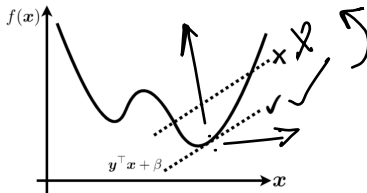
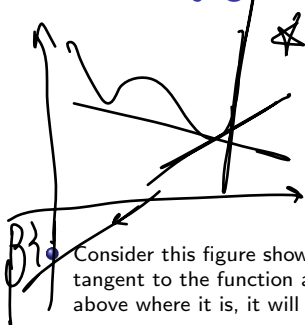
- According to [5], the convex relaxation of " ℓ_0 " norm (subset selection) is ℓ_1 norm. Therefore, we write the regularized optimization as:

$$\underset{\mathbf{x}}{\text{minimize}} \quad \Omega(\mathbf{x}) := \Omega_0(\mathbf{x}) + \lambda \|\mathbf{x}\|_1. \quad (3)$$

- The ℓ_1 regularization is also referred to as lasso (least absolute shrinkage and selection operator) regularization [6, 7]. Different methods exist for solving optimization having ℓ_1 norm, such as its approximation by Huber function [8], proximal algorithm and soft thresholding [9], coordinate descent [10, 11], and subgradients. In the following, we explain these methods.

Convex Conjugate and the Huber function

Convex Conjugate



f
support

- Consider this figure showing a line which supports the function f meaning that it is tangent to the function and the function upper-bounds it. In other words, if the line goes above where it is, it will intersect the function in more than a point.
- Now let the support line be multi-dimensional to be a support hyperplane. For having this tangent support hyperplane with slope $y \in \mathbb{R}^d$ and intercept $\beta \in \mathbb{R}$, we should have:

$$y^\top x + \beta = f(x) \implies \beta = f(x) - y^\top x.$$

- We want the smallest intercept for the support hyperplane:

$$\beta^* = \min_{x \in \mathbb{R}^d} (f(x) - y^\top x) = - \max_{x \in \mathbb{R}^d} (y^\top x - f(x)).$$

Convex Conjugate

- We found:

$$\beta^* = \min_{\mathbf{x} \in \mathbb{R}^d} (f(\mathbf{x}) - \mathbf{y}^\top \mathbf{x}) = - \max_{\mathbf{x} \in \mathbb{R}^d} (\mathbf{y}^\top \mathbf{x} - f(\mathbf{x})).$$

Definition (Convex conjugate of function)

The conjugate gradient of function $f(\cdot)$ is defined as:

$$f^*(\mathbf{y}) := \sup_{\mathbf{x} \in \mathbb{R}^d} (\mathbf{y}^\top \mathbf{x} - f(\mathbf{x})).$$

(4)

- Therefore, we define $f^*(\mathbf{y}) := -\beta^*$ to have the convex conjugate defined above.
- The convex conjugate of a function is always convex, even if the function itself is not convex, because it is point-wise maximum of affine functions.

Convex Conjugate

Recall the convex conjugate: $f^*(\mathbf{y}) := \sup_{\mathbf{x} \in \mathbb{R}^d} (\mathbf{y}^\top \mathbf{x} - f(\mathbf{x}))$. ✱

Lemma (Conjugate of convex conjugate)

The conjugate of convex conjugate of a function is:

$$\underbrace{f^{**}(\mathbf{x})}_{\text{conjugate of } f^*} = \sup_{\mathbf{y} \in \text{dom}(f^*)} \underbrace{(\mathbf{x}^\top \mathbf{y} - f^*(\mathbf{y}))}_{\text{conjugate of } f}. \quad (5)$$

It is always a lower-bound for the function, i.e., $f^{**}(\mathbf{x}) \leq f(\mathbf{x})$.

If the function $f(\cdot)$ is convex, we have $\underline{f^{**}(\mathbf{x}) = f(\mathbf{x})}$; hence, for a convex function, we have:

$$\star \quad \boxed{f(\mathbf{x}) = \sup_{\mathbf{y} \in \text{dom}(f^*)} (\mathbf{x}^\top \mathbf{y} - f^*(\mathbf{y}))}. \quad (6)$$

$f^{**}(\mathbf{x})$

Huber Function: Smoothing L1 Norm by Convex Conjugate

Lemma (The convex conjugate of ℓ_1 norm)

The convex conjugate of $f(\cdot) = \|\cdot\|_1$ is:

$$f(x) = \|x\|_1 \quad \rightarrow \quad f^*(y) = \begin{cases} 0 & \text{if } \|y\|_\infty \leq 1 \\ \infty & \text{Otherwise.} \end{cases} \quad \star \star$$

$$x = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad \star$$

Proof.

We can write ℓ_1 norm as:

$$\star \star \quad f(x) = \|x\|_1 = \max_{\|z\|_\infty \leq 1} x^T z. \quad \star \star$$



Using this in Eq. (4), $f^*(y) := \sup_{x \in \mathbb{R}^d} (y^T x - f(x))$, results in Eq. (7) because:

$$f^*(y) := \sup_{x \in \mathbb{R}^d} (y^T x - \underbrace{\max_{\|z\|_\infty \leq 1} \underbrace{x^T z}_{f(x)}}_{\star \star}) = \sup_{x \in \mathbb{R}^d} (y^T x - \max_{\|z\|_\infty \leq 1} \underbrace{z^T x}_{\star \star}) = \begin{cases} 0 & \text{if } \|y\|_\infty \leq 1 \\ \infty & \text{Otherwise.} \end{cases}$$

$\star \star \quad \|y\|_\infty \leq 1$

□

Huber Function: Smoothing L1 Norm by Convex Conjugate

Lemma (Gradient in terms of convex conjugate)

For any function $f(\cdot)$, we have:

★

$$\nabla f(\mathbf{x}) = \arg \max_{\mathbf{y} \in \text{dom}(f^*)} (\mathbf{x}^\top \mathbf{y} - f^*(\mathbf{y})).$$

(8)

- We saw that the convex conjugate of $f(\cdot) = \|\cdot\|_1$ is:

★

$$f^*(\mathbf{y}) = \begin{cases} 0 & \text{if } \|\mathbf{y}\|_\infty \leq 1 \\ \infty & \text{Otherwise.} \end{cases}$$

- According to Eq. (8), we have $\nabla f(\mathbf{x}) = \arg \max_{\|\mathbf{y}\|_\infty \leq 1} \mathbf{x}^\top \mathbf{y}$ for $f(\cdot) = \|\cdot\|_1$.
- For $\mathbf{x} = \mathbf{0}$, we have $\nabla f(\mathbf{x}) = \arg \max_{\|\mathbf{y}\|_\infty \leq 1} \mathbf{0}$ which has many solutions. Therefore, at $\mathbf{x} = \mathbf{0}$, the function $\|\cdot\|_1$ norm is not differentiable and not smooth because the gradient at that point is not unique.

Huber Function: Smoothing L1 Norm by Convex Conjugate

- We can smooth the ℓ_1 norm at $\mathbf{x} = \mathbf{0}$ using convex conjugate.
- We saw that the convex conjugate of $f(\cdot) = \|\cdot\|_1$ is:

$$f^*(\mathbf{y}) = \begin{cases} 0 & \text{if } \|\mathbf{y}\|_\infty \leq 1 \\ \infty & \text{Otherwise.} \end{cases}$$

- Let $\mathbf{x} = [x_1, \dots, x_d]^\top$. As we have $f(\mathbf{x}) = \|\mathbf{x}\|_1 = \sum_{j=1}^d |x_j|$, we can use the convex conjugate for every dimension $f(x_j) = |x_j|$:

$$f^*(y_j) = \begin{cases} 0 & \text{if } |y_j| \leq 1 \\ \infty & \text{Otherwise.} \end{cases} \quad (9)$$

- According to Eq. (6), $f(\mathbf{x}) = \sup_{\mathbf{y} \in \text{dom}(f^*)} (\mathbf{x}^\top \mathbf{y} - f^*(\mathbf{y}))$, we have:

$$|x_j| = \sup_{y \in \mathbb{R}} (x_j y - f^*(y_j)) \stackrel{(9)}{=} \max_{|y_j| \leq 1} x_j y_j = \max_{|y_j| \leq 1} x_j y_j$$

- This is not unique for $x_j = 0$. Hence, we add a μ -strongly convex function to the above equation to make the solution unique at $x_j = 0$ also.
- This added term is named the **proximity function** defined below.

Huber Function: Smoothing L1 Norm by Convex Conjugate

Definition (Proximity function [12])

A proximity function $p(y)$ for a closed convex set $S \in \text{dom}(p)$ is a function which is continuous and strongly convex. We can change Eq. (6), $f(x) = \sup_{y \in \text{dom}(f^*)} (x^\top y - f^*(y))$, to:

$$\star \quad f(x) \approx f_\mu(x) := \sup_{y \in \text{dom}(f^*)} \underbrace{(x^\top y - f^*(y))}_{\star} - \underbrace{\mu p(y)}_{\star}, \quad (10)$$

where $\mu > 0$.

- Recall Eq. (9): $f^*(y_j) = \begin{cases} 0 & \text{if } |y_j| \leq 1 \\ \infty & \text{Otherwise.} \end{cases} \star$
- Using Eq. (10), we can have:

$$\star \quad |x_j| \approx \sup_{y \in \mathbb{R}} (x_j y_j - f^*(y_j) - \frac{\mu}{2} y_j^2) \stackrel{(9)}{=} \max_{|y_j| \leq 1} (x_j y_j - \frac{\mu}{2} y_j^2) = \begin{cases} \frac{x_j^2}{2\mu} & \text{if } |x_j| \leq \mu \\ |x_j| - \frac{\mu}{2} & \text{if } |x_j| > \mu. \end{cases} \star$$

- This approximation to ℓ_1 norm, which is differentiable everywhere, including at $x_j = 0$, is named the Huber function.
- Note that the Huber function is the Moreau envelope of absolute value.

Huber Function: Smoothing L1 Norm by Convex Conjugate

Definition (Huber and pseudo-Huber functions (1992) [8])

The Huber function and pseudo-Huber functions are:

$$\star \quad h_{\mu}(x) = \begin{cases} \frac{x^2}{2\mu} & \text{if } |x| \leq \mu \\ |x| - \frac{\mu}{2} & \text{if } |x| > \mu, \end{cases} \quad (11)$$

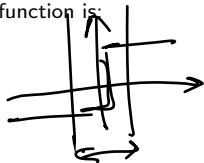
$$\star \quad \hat{h}_{\mu}(x) = \sqrt{\left(\frac{x}{\mu}\right)^2 + 1} - 1, \quad (12)$$

respectively, where $\mu > 0$.

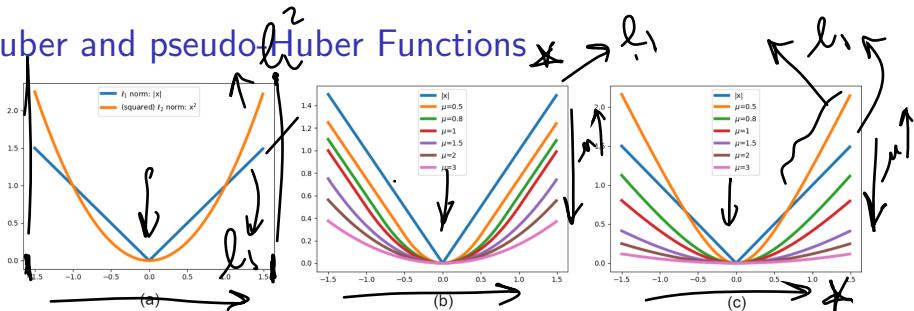
The derivative of these functions is easily calculated. For example, the derivative of Huber function is:

$$\nabla h_{\mu}(x) = \begin{cases} \frac{x}{\mu} & \text{if } |x| \leq \mu \\ \text{sign}(x) & \text{if } |x| > \mu. \end{cases}$$

Handwritten notes: A box labeled $\mu \sim 0$ has an arrow pointing to the $|x| \leq \mu$ case of the derivative formula. Another arrow points from the $|x| > \mu$ case to the right.



Huber and pseudo-Huber Functions



- (a) Comparison of ℓ_1 and ℓ_2 norms in \mathbb{R}^1 , (b) comparison of ℓ_1 norm (i.e., absolute value in \mathbb{R}^1) and the Huber function, and (c) comparison of ℓ_1 norm (i.e., absolute value in \mathbb{R}^1) and the pseudo-Huber function.
- In contrast to ℓ_1 norm or absolute value, these two functions are smooth so they approximate the ℓ_1 norm smoothly. This figure also shows that the Huber function is always upper-bounded by absolute value (ℓ_1 norm); however, this does not hold for pseudo-Huber function.
- We can also see that the approximation of Huber function is better than the approximation of pseudo-Huber function; however, its calculation is harder than pseudo-Huber function because it is a piece-wise function (compare Eqs. (11) and (12)).
- Moreover, the figure shows a smaller positive value μ give better approximations, although it makes calculation of the Huber and pseudo-Huber functions harder.

Soft-thresholding and Proximal Methods

Soft-thresholding and Proximal Methods

- Proximal mapping was introduced before:

$$\text{prox}_{\lambda g}(\mathbf{x}) := \arg \min_{\mathbf{u}} \left(g(\mathbf{u}) + \frac{1}{2\lambda} \|\mathbf{u} - \mathbf{x}\|_2^2 \right). \quad (13)$$

- We can use proximal mapping of non-smooth functions to solve non-smooth optimization by proximal methods introduced before.
- For example, we can solve an optimization problem containing ℓ_1 norm in its objective function using the proximal mapping of ℓ_1 norm (soft-thresholding).

$$[\text{prox}_{\lambda \|\cdot\|_1}(\mathbf{x})]_j = \max(0, |x_j| - \lambda) \text{sign}(x_j) = s_\lambda(x_j) := \begin{cases} x_j - \lambda & \text{if } x_j \geq \lambda \\ 0 & \text{if } |x_j| < \lambda \\ x_j + \lambda & \text{if } x_j \leq -\lambda \end{cases} \quad (14)$$

- Then, we can use any of the proximal methods such as proximal point method and proximal gradient method.
- For solving the regularized problem (3), minimize_x $\Omega(\mathbf{x}) := \Omega_0(\mathbf{x}) + \lambda \|\mathbf{x}\|_1$ which is optimizing a composite function, we can use the proximal gradient method introduced before.

Coordinate Descent

Coordinate Method

$$P(x_1, \dots, x_d)$$

- Assume $\mathbf{x} = [x_1, \dots, x_d]^T$. For solving $\min_{\mathbf{x}} f(\mathbf{x})$, coordinate method [10] updates the dimensions (coordinates) of solution one-by-one and not all dimensions together at once:

$$\left. \begin{aligned} x_1^{(k+1)} &:= \arg \min_{x_1} f(x_1, x_2^{(k)}, x_3^{(k)}, \dots, x_d^{(k)}), \\ x_2^{(k+1)} &:= \arg \min_{x_2} f(x_1^{(k+1)}, x_2, x_3^{(k)}, \dots, x_d^{(k)}), \\ &\vdots \\ x_d^{(k+1)} &:= \arg \min_{x_d} f(x_1^{(k+1)}, x_2^{(k+1)}, x_3^{(k+1)}, \dots, x_d), \end{aligned} \right\} (15)$$

until convergence of all dimensions of solution.

- Note that the update of every dimension uses the latest update of previously updated dimensions.
- The order of updates for the dimensions does not matter.
- The idea of coordinate descent algorithm is similar to the idea of Gibbs sampling [13, 14] where we work on the dimensions of the variable one by one.

*

Coordinate Descent

- If we use a step of gradient descent, $\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} - \eta \nabla f(\mathbf{x}^{(k)})$, for every of the above updates, the method is named coordinate descent.
- If we use proximal gradient method, $\mathbf{x}^{(k+1)} := \text{prox}_{\eta^{(k)}g}(\mathbf{x}^{(k)} - \eta^{(k)} \nabla f(\mathbf{x}^{(k)}))$, for every update in coordinate method, the method is named the proximal coordinate descent.
- We can group some of the dimensions (features) together and alternate between updating the blocks (groups) of features. That method is named block coordinate descent.
- The convergence analysis of coordinate descent and block coordinate descent methods can be found in [15, 16] and [17], respectively. They show that if the function $f(\cdot)$ is continuous, proper, and closed, the coordinate descent method converges to a stationary point.
- There exist some other faster variants of coordinate descent named accelerated coordinate descent [18, 19].
- Similar to SGD, the full gradient is not available in coordinate descent to use for checking convergence. So, we should use other convergence criteria such as maximum number of iterations or checking convergence for each of the variables.
- Although coordinate descent methods are very simple and have shown to work properly for ℓ_1 norm optimization [11], they have not sufficiently attracted the attention of researchers in the field of optimization [10].

L1 Norm Optimization

- Coordinate descent method can be used for ℓ_1 norm (lasso) optimization [11] because every coordinate of the ℓ_1 norm is an absolute value ($\|\mathbf{x}\|_1 = \sum_{j=1}^d |x_j|$ for $\mathbf{x} = [x_1, \dots, x_d]^\top$) and the derivative of absolute value is a simple sign function.
- One of the well-known ℓ_1 optimization methods is the lasso regression (1996) [6, 2, 7]:

$$\star \quad \underset{\beta}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_1, \quad (16)$$

where $\mathbf{y} \in \mathbb{R}^n$ are the labels, $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_d] \in \mathbb{R}^{n \times d}$ are the observations, $\beta = [\beta_1, \dots, \beta_d]^\top \in \mathbb{R}^d$ are the regression coefficients, and λ is the regularization parameter. The lasso regression is sparse which is effective as explained before.

- Let c denote the objective function in Eq. (16). The objective function can be simplified as:

$$c = 0.5(\mathbf{y}^\top \mathbf{y} - 2\beta^\top \mathbf{X}^\top \mathbf{y} + \beta^\top \mathbf{X}^\top \mathbf{X} \beta) + \lambda \|\beta\|_1.$$

$\|\mathbf{x}\|_2^2 = \mathbf{x}^\top \mathbf{x}$

- We can write the j -th element of this objective, denoted by c_j , as:

$$\star \quad c_j = \frac{1}{2}(\mathbf{y}^\top \mathbf{y} - 2\mathbf{x}_j^\top \mathbf{y} \beta_j + \beta_j \mathbf{x}_j^\top \mathbf{x}_j \beta_j + \beta_j \mathbf{x}_j^\top \mathbf{X}_{-j} \beta_{-j}) + \lambda |\beta_j|,$$

where $\mathbf{X}_{-j} := [\mathbf{x}_1, \dots, \mathbf{x}_{j-1}, \mathbf{x}_{j+1}, \dots, \mathbf{x}_d]$ and $\beta_{-j} := [\beta_1, \dots, \beta_{j-1}, \beta_{j+1}, \dots, \beta_d]^\top$.

L1 Norm Optimization

- We wrote the j th element of this objective, denoted by c_j , as:
 $c_j = \frac{1}{2}(\mathbf{y}^\top \mathbf{y} - 2\mathbf{x}_j^\top \mathbf{y} \beta_j + \beta_j \mathbf{x}_j^\top \mathbf{x}_j \beta_j + \beta_j \mathbf{x}_j^\top \mathbf{X}_{-j} \beta_{-j}) + \lambda |\beta_j|$
- For coordinate descent, we need gradient of objective function w.r.t. every coordinate.
 The derivatives of other coordinates of objective w.r.t. β_j are zero so we need c_j for derivative w.r.t. β_j .
- Taking derivative of c_j w.r.t. β_j and setting it to zero gives:

$$\frac{\partial c}{\partial \beta_j} = \frac{\partial c_j}{\partial \beta_j} = \mathbf{x}_j^\top \mathbf{x}_j \beta_j + \mathbf{x}_j^\top (\mathbf{X}_{-j} \beta_{-j} - \mathbf{y}) + \lambda \text{sign}(\beta_j) \stackrel{\text{set}}{=} 0$$

$$\Rightarrow \beta_j = s_{\frac{\lambda}{\|\mathbf{x}_j\|_2^2}} \left(\frac{\mathbf{x}_j^\top (\mathbf{y} - \mathbf{X}_{-j} \beta_{-j})}{\mathbf{x}_j^\top \mathbf{x}_j} \right)$$

$$\star = \begin{cases} \frac{\mathbf{x}_j^\top (\mathbf{y} - \mathbf{X}_{-j} \beta_{-j})}{\|\mathbf{x}_j\|_2^2} - \frac{\lambda}{\|\mathbf{x}_j\|_2^2} & \text{if } \frac{\mathbf{x}_j^\top (\mathbf{y} - \mathbf{X}_{-j} \beta_{-j})}{\mathbf{x}_j^\top \mathbf{x}_j} \geq \frac{\lambda}{\|\mathbf{x}_j\|_2^2} \\ 0 & \text{if } \left| \frac{\mathbf{x}_j^\top (\mathbf{y} - \mathbf{X}_{-j} \beta_{-j})}{\mathbf{x}_j^\top \mathbf{x}_j} \right| < \frac{\lambda}{\|\mathbf{x}_j\|_2^2} \\ \frac{\mathbf{x}_j^\top (\mathbf{y} - \mathbf{X}_{-j} \beta_{-j})}{\|\mathbf{x}_j\|_2^2} + \frac{\lambda}{\|\mathbf{x}_j\|_2^2} & \text{if } \frac{\mathbf{x}_j^\top (\mathbf{y} - \mathbf{X}_{-j} \beta_{-j})}{\mathbf{x}_j^\top \mathbf{x}_j} \leq -\frac{\lambda}{\|\mathbf{x}_j\|_2^2} \end{cases}$$

which is a soft-thresholding function (see Eq. (14)).

- Therefore, coordinate descent for ℓ_1 optimization finds the soft-thresholding solution, the same as the proximal mapping. We can use this soft-thresholding in coordinate descent where we use β_j 's in Eq. (15) rather than \mathbf{x}_j 's.

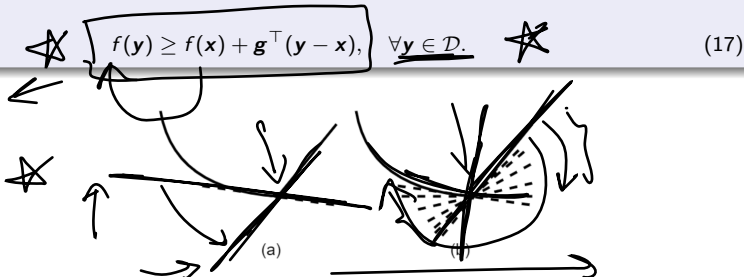
Subgradient Methods

Subgradient

- We know that the convex conjugate $f^*(y) := \sup_{x \in \mathbb{R}^d} (y^\top x - f(x))$ is always convex.
- If the convex conjugate $f^*(y)$ is strongly convex, then we have only one gradient according to Eq. (8), $\nabla f(x) = \arg \max_{y \in \text{dom}(f^*)} (x^\top y - f^*(y))$ ~~✗~~
- However, if the convex conjugate is only convex and not strongly convex, Eq. (8) might have several solutions so the gradient may not be unique.
- For the points in which the function does not have a unique gradient, we can have a set of subgradients, defined below.

Definition (Subgradient)

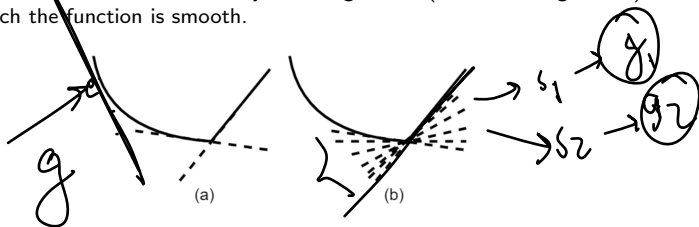
Consider a convex function $f(\cdot)$ with domain \mathcal{D} . The vector $g \in \mathbb{R}^d$ is a subgradient of $f(\cdot)$ at $x \in \mathcal{D}$ if it satisfies:

$$f(y) \geq f(x) + g^\top (y - x), \quad \forall y \in \mathcal{D}. \quad (17)$$


The diagram illustrates the definition of a subgradient. It shows a convex function $f(y)$ and a supporting hyperplane at point x . The hyperplane is defined by the equation $f(y) \geq f(x) + g^\top (y - x)$. The diagram includes a 2D plot (a) and a 3D plot (b) showing the supporting hyperplane and the subgradient vector g .

Subdifferential

- As this figure shows, if the function is not smooth at a point, it has multiple subgradients at that point. This is while there is only one subgradient (which is the gradient) for a point at which the function is smooth.



Definition (subdifferential)

The subdifferential of a convex function $f(\cdot)$, with domain \mathcal{D} , at a point $x \in \mathcal{D}$ is the set of all subgradients at that point:

$$\star \quad \partial f(x) = \{ \mathbf{g} \mid \mathbf{g}^\top (\mathbf{y} - \mathbf{x}) \leq f(\mathbf{y}) - f(\mathbf{x}), \forall \mathbf{y} \in \mathcal{D} \}. \quad \star \quad (18)$$

- The subdifferential is a closed convex set.
- Every subgradient is a member of the subdifferential, i.e., $\mathbf{g} \in \partial f(\mathbf{x})$. An example subdifferential is shown in the above figure.

Subdifferential for ℓ_1 norm

- An example of subgradient is the subdifferential of absolute value, $f(\cdot) = |\cdot|$:

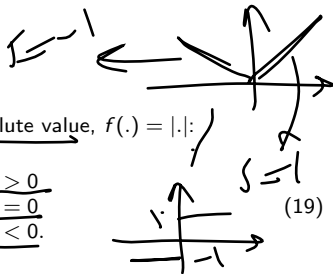
$$\partial f(x) = \begin{cases} 1 & \text{if } x > 0 \\ \in [-1, 1] & \text{if } x = 0 \\ -1 & \text{if } x < 0. \end{cases}$$

- The subgradient of absolute value is equal to the gradient for $x < 0$ and $x > 0$ while there exists a set of subgradients at $x = 0$ because absolute value is not smooth at that point.
- We can also compute the subgradient of ℓ_1 norm because we have:

$$f(\mathbf{x}) = \|\mathbf{x}\|_1 = \sum_{i=1}^d |x_i| = \sum_{i=1}^d f_i(\mathbf{x}_i),$$

for $\mathbf{x} = [x_1, \dots, x_d]^\top$.

- We take Eq. (19) as the subdifferential of the i -th dimension, denoted by $\partial f_i(x_i)$. Hence, for $f(\mathbf{x}) = \|\mathbf{x}\|_1$, we have $\partial f(\mathbf{x}) = \partial f_1(x_1) \times \dots \times \partial f_d(x_d)$ where \times denotes the Cartesian product of sets.



Subgradient

- ★ • We can have the first-order optimality condition using subgradients by generalizing $\nabla f(\mathbf{x}^*) = \mathbf{0}$ as follows.

Lemma (First-order optimality condition with subgradient)

If \mathbf{x}^* is a local minimizer for a function $f(\cdot)$, then:

★ ★ $\boxed{\mathbf{0} \in \partial f(\mathbf{x}^*)}.$



(20)

Note that if $f(\cdot)$ is convex, this equation is a necessary and sufficient condition for a minimizer.

Proof.

According to Eq. (17) in the definition of subgradient, we have:

★ $f(\mathbf{y}) \geq f(\mathbf{x}^*) + \underbrace{\mathbf{g}^\top}_{\mathbf{0}} (\mathbf{y} - \mathbf{x}^*), \forall \mathbf{y}.$

If we have $\mathbf{g} = \mathbf{0} \in \partial f(\mathbf{x}^*)$, we have:

$f(\mathbf{y}) \geq f(\mathbf{x}^*) + \mathbf{0}^\top (\mathbf{y} - \mathbf{x}^*) = f(\mathbf{x}^*),$

which means that \mathbf{x}^* is a minimizer.

□

Subgradient in Some Example Functions

- The following lemma can be useful for calculation of subdifferential of functions.

Lemma

Some useful properties for calculation of subdifferential of functions:

- For a smooth function or at points where the function is smooth, subdifferential has only one member which is the gradient: $\partial f(\mathbf{x}) = \{\nabla f(\mathbf{x})\}$.
- Linear combination: If $f(\mathbf{x}) = \sum_{i=1}^n a_i f_i(\mathbf{x})$ with $a_i \geq 0$, then $\partial f(\mathbf{x}) = \sum_{i=1}^n a_i \partial f_i(\mathbf{x})$.
- Affine transformation: If $f(\mathbf{x}) = f_0(\mathbf{Ax} + \mathbf{b})$, then $\partial f(\mathbf{x}) = \mathbf{A}^\top \partial f_0(\mathbf{Ax} + \mathbf{b})$.
- Point-wise maximum: Suppose $f(\mathbf{x}) = \max\{f_1(\mathbf{x}), \dots, f_n(\mathbf{x})\}$ where f_i 's are differentiable. Let $I(\mathbf{x}) := \{i | f_i = f(\mathbf{x})\}$ states which function has the maximum value for the point \mathbf{x} . At any point other than the intersection point of functions (which is smooth), the subgradient is $\mathbf{g} = \nabla f_i(\mathbf{x})$ for $i \in I(\mathbf{x})$. At the intersection point of two functions (which is not smooth), e.g. $f_i(\mathbf{x}) = f_{i+1}(\mathbf{x})$, we have:

$$\partial f(\mathbf{x}) = \{\mathbf{g} \mid t \nabla f_i(\mathbf{x}) + (1-t) \nabla f_{i+1}(\mathbf{x}), \forall t \in [0, 1]\}.$$



Subgradient Method

- The subgradient method, first proposed in [20], is used for solving the unconstrained optimization problem, minimize_x f(x), where the function f(.) is not smooth, i.e., not differentiable, everywhere in its domain.
- It iteratively updates the solution as:

$$\star \boxed{\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} - \eta^{(k)} \mathbf{g}^{(k)},} \star \quad (21)$$

$\nearrow \nabla f$

where $\mathbf{g}^{(k)}$ is any subgradient of function $f(.)$ in point \mathbf{x} at iteration k , i.e. $\mathbf{g}^{(k)} \in \partial f(\mathbf{x}^{(k)})$, and $\eta^{(k)}$ is the step size at iteration k .

- Comparing this update with the update in gradient descent:

$$\boxed{\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} - \eta \nabla f(\mathbf{x}^{(k)})},$$

shows that gradient descent is a special case of the subgradient method because for a smooth function, gradient is the only member of the subdifferential set (see Lemma in the previous slide); hence, the only subgradient is the gradient.

Stochastic Subgradient Method

- Consider the optimization problem $\text{minimize}_{\mathbf{x}} f(\mathbf{x})$ where at least one of the $f_i(\cdot)$ functions is not smooth.
- Stochastic subgradient method** [21] randomly samples one of the points to update the solution in every iteration:

$$\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} - \eta^{(k)} \mathbf{g}_i^{(k)},$$

$$\nabla f_i \quad (22)$$

where $\mathbf{g}_i^{(k)} \in \partial f_i(\mathbf{x}^{(k)})$.

- Comparing this with the update in stochastic gradient descent (SGD):

$$\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} - \eta^{(k)} \nabla f_i(\mathbf{x}^{(k)}),$$

shows that stochastic gradient descent is a special case of stochastic gradient descent because for a smooth function, gradient is the only member of the subdifferential set (see Lemma in two previous slides).

- We can have **mini-batch stochastic subgradient method** which is a generalization of mini-batch SGD for non-smooth functions. In this case, the update of solution is:

$$\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} - \eta^{(k)} \frac{1}{b} \sum_{i \in \mathcal{B}_{k'}} \mathbf{g}_i^{(k)}.$$

(23)

- If the function is not smooth, we can also use subgradient instead of gradient in other stochastic methods such as SAG and SVRG, which were introduced before. For this, we need to use $\mathbf{g}_i^{(k)}$ rather than $\nabla f(\mathbf{x}^{(k)})$ in these methods.

Projected Subgradient Method

- Consider the constrained optimization problem:

$$\begin{array}{ll} \underset{\mathbf{x}}{\text{minimize}} & f(\mathbf{x}) \\ \text{subject to} & \mathbf{x} \in \mathcal{S}, \end{array} \quad (24)$$

where \mathcal{S} is the feasible set of constraints.

- If the function $f(\cdot)$ is not smooth, we can use the projected subgradient method [22] which generalizes the projected gradient method introduced before.
- Similar to the update in projected gradient method:

$$\mathbf{x}^{(k+1)} := \Pi_{\mathcal{S}}(\mathbf{x}^{(k)} - \eta^{(k)} \nabla f(\mathbf{x}^{(k)})),$$

projected subgradient method iteratively updates the solution as:

$$\mathbf{x}^{(k+1)} = \Pi_{\mathcal{S}}(\mathbf{x}^{(k)} - \eta^{(k)} \mathbf{g}^{(k)}), \quad (25)$$

until convergence of the solution.

Acknowledgement

- Some slides of this slide deck are inspired by the lectures of Prof. Stephen Boyd at the Stanford University.
- Some slides of this slide deck are inspired by the teaching of Prof. Mu Zhu at the University of Waterloo.
- Our tutorial also has the materials of this slide deck: [23]

References

- [1] B. Ghogh and M. Crowley, "The theory behind overfitting, cross validation, regularization, bagging, and boosting: tutorial," *arXiv preprint arXiv:1905.12787*, 2019.
- [2] J. Friedman, T. Hastie, R. Tibshirani, et al., *The elements of statistical learning*, vol. 1. Springer series in statistics New York, 2001.
- [3] P. Domingos, "The role of Occam's razor in knowledge discovery," *Data mining and knowledge discovery*, vol. 3, no. 4, pp. 409–425, 1999.
- [4] S. Boyd and L. Vandenberghe, *Convex optimization.* Cambridge university press, 2004.
- [5] D. L. Donoho, "For most large underdetermined systems of linear equations the minimal ℓ_1 -norm solution is also the sparsest solution," *Communications on Pure and Applied Mathematics: A Journal Issued by the Courant Institute of Mathematical Sciences*, vol. 59, no. 6, pp. 797–829, 2006.
- [6] R. Tibshirani, "Regression shrinkage and selection via the lasso," *Journal of the Royal Statistical Society: Series B (Methodological)*, vol. 58, no. 1, pp. 267–288, 1996.
- [7] T. Hastie, R. Tibshirani, and M. Wainwright, *Statistical learning with sparsity: the lasso and generalizations.* Chapman and Hall/CRC, 2019.

References (cont.)

- [8] P. J. Huber, "Robust estimation of a location parameter," in *Breakthroughs in statistics*, pp. 492–518, Springer, 1992.
- [9] N. Parikh and S. Boyd, "Proximal algorithms," *Foundations and Trends in optimization*, vol. 1, no. 3, pp. 127–239, 2014.
- [10] S. J. Wright, "Coordinate descent algorithms," *Mathematical Programming*, vol. 151, no. 1, pp. 3–34, 2015.
- [11] T. T. Wu and K. Lange, "Coordinate descent algorithms for lasso penalized regression," *The Annals of Applied Statistics*, vol. 2, no. 1, pp. 224–244, 2008.
- [12] B. Banaschewski and J.-M. Maranda, "Proximity functions," *Mathematische Nachrichten*, vol. 23, no. 1, pp. 1–37, 1961.
- [13] S. Geman and D. Geman, "Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images," *IEEE Transactions on pattern analysis and machine intelligence*, no. 6, pp. 721–741, 1984.
- [14] B. Ghojogh, H. Nekoie, A. Ghojogh, F. Karay, and M. Crowley, "Sampling algorithms, from survey sampling to Monte Carlo methods: Tutorial and literature review," *arXiv preprint arXiv:2011.00901*, 2020.
- [15] Z.-Q. Luo and P. Tseng, "On the convergence of the coordinate descent method for convex differentiable minimization," *Journal of Optimization Theory and Applications*, vol. 72, no. 1, pp. 7–35, 1992.

References (cont.)

- [16] Z.-Q. Luo and P. Tseng, "Error bounds and convergence analysis of feasible descent methods: a general approach," *Annals of Operations Research*, vol. 46, no. 1, pp. 157–178, 1993.
- [17] P. Tseng, "Convergence of a block coordinate descent method for nondifferentiable minimization," *Journal of optimization theory and applications*, vol. 109, no. 3, pp. 475–494, 2001.
- [18] Y. T. Lee and A. Sidford, "Efficient accelerated coordinate descent methods and faster algorithms for solving linear systems," in *2013 IEEE 54th annual symposium on foundations of computer science*, pp. 147–156, IEEE, 2013.
- [19] O. Fercoq and P. Richtárik, "Accelerated, parallel, and proximal coordinate descent," *SIAM Journal on Optimization*, vol. 25, no. 4, pp. 1997–2023, 2015.
- [20] N. Z. Shor, *Minimization methods for non-differentiable functions*, vol. 3. Springer Science & Business Media, 2012.
- [21] N. Z. Shor, *Nondifferentiable optimization and polynomial problems*, vol. 24. Springer Science & Business Media, 1998.
- [22] Y. I. Alber, A. N. Iusem, and M. V. Solodov, "On the projected subgradient method for nonsmooth convex optimization in a Hilbert space," *Mathematical Programming*, vol. 81, no. 1, pp. 23–35, 1998.

References (cont.)

- * [23] B. Ghojogh, A. Ghodsi, F. Karray, and M. Crowley, "KKT conditions, first-order and second-order optimization, and distributed optimization: Tutorial and survey," *arXiv preprint arXiv:2110.01858*, 2021.