

## More Subgradient Calculus: Function Convexity first

Following functions are again convex, but again, may not be differentiable everywhere. How does one compute their subgradients at points of non-differentiability?

- Nonnegative weighted sum:  $f = \sum_{i=1}^n \alpha_i f_i$  is convex if each  $f_i$  for  $1 \leq i \leq n$  is convex and  $\alpha_i \geq 0, 1 \leq i \leq n$ .
- Composition with affine function:  $f(Ax + b)$  is convex if  $f$  is convex. For example:
  - ▶ The log barrier for linear inequalities,  $f(x) = -\sum_{i=1}^m \log(b_i - a_i^T x)$ , is convex since  $-\log(x)$  is convex.
  - ▶ Any norm of an affine function,  $f(x) = \|Ax + b\|$ , is convex.

## More of Basic Subgradient Calculus

- Scaling:  $\partial(af) = a \cdot \partial f$  provided  $a > 0$ . The condition  $a > 0$  makes function  $f$  remain convex.
- Addition:  $\partial(f_1 + f_2) = \partial(f_1) + \partial(f_2)$
- Affine composition: if  $g(\mathbf{x}) = f(A\mathbf{x} + \mathbf{b})$ , then  $\partial g(\mathbf{x}) = A^T \partial f(A\mathbf{x} + \mathbf{b})$
- Norms: important special case,  $f(\mathbf{x}) = \|\mathbf{x}\|_p$

The derivations done in class could be used to show that if any other subgradient exists for  $g$  outside the stated set above, that could be used to construct a subgradient for  $f$  outside the stated set above as well!

## More of Basic Subgradient Calculus

- Scaling:  $\partial(af) = a \cdot \partial f$  provided  $a > 0$ . The condition  $a > 0$  makes function  $f$  remain convex.
- Addition:  $\partial(f_1 + f_2) = \partial(f_1) + \partial(f_2)$
- Affine composition: if  $g(\mathbf{x}) = f(A\mathbf{x} + \mathbf{b})$ , then  $\partial g(\mathbf{x}) = A^T \partial f(A\mathbf{x} + \mathbf{b})$
- Norms: important special case,  $f(\mathbf{x}) = \|\mathbf{x}\|_p = \max_{\|\mathbf{z}\|_q \leq 1} \mathbf{z}^T \mathbf{x}$  where  $q$  is such that

$1/p + 1/q = 1$ . Then

On the board we have used  $y$  instead of  $z$

## More of Basic Subgradient Calculus

- Scaling:  $\partial(af) = a \cdot \partial f$  provided  $a > 0$ . The condition  $a > 0$  makes function  $f$  remain convex.
- Addition:  $\partial(f_1 + f_2) = \partial(f_1) + \partial(f_2)$
- Affine composition: if  $g(\mathbf{x}) = f(A\mathbf{x} + \mathbf{b})$ , then  $\partial g(\mathbf{x}) = A^T \partial f(A\mathbf{x} + \mathbf{b})$
- Norms: important special case,  $f(\mathbf{x}) = \|\mathbf{x}\|_p = \max_{\|\mathbf{z}\|_q \leq 1} \mathbf{z}^T \mathbf{x}$  where  $q$  is such that

$1/p + 1/q = 1$ . Then

$$\partial f(\mathbf{x}) = \left\{ \mathbf{y} : \|\mathbf{y}\|_q \leq 1 \text{ and } \mathbf{y}^T \mathbf{x} = \max_{\|\mathbf{z}\|_q \leq 1} \mathbf{z}^T \mathbf{x} \right\} =$$

**y corresponds to z where the max is attained**

The part above is largely connected to previous discussion on max of convex functions

## More of Basic Subgradient Calculus

- Scaling:  $\partial(af) = a \cdot \partial f$  provided  $a > 0$ . The condition  $a > 0$  makes function  $f$  remain convex.
- Addition:  $\partial(f_1 + f_2) = \partial(f_1) + \partial(f_2)$
- Affine composition: if  $g(\mathbf{x}) = f(A\mathbf{x} + \mathbf{b})$ , then  $\partial g(\mathbf{x}) = A^T \partial f(A\mathbf{x} + \mathbf{b})$
- Norms: important special case,  $f(\mathbf{x}) = \|\mathbf{x}\|_p = \max_{\|\mathbf{z}\|_q \leq 1} \mathbf{z}^T \mathbf{x}$  where  $q$  is such that

$1/p + 1/q = 1$ . Then

$$\partial f(\mathbf{x}) = \left\{ \mathbf{y} : \|\mathbf{y}\|_q \leq 1 \text{ and } \mathbf{y}^T \mathbf{x} = \max_{\|\mathbf{z}\|_q \leq 1} \mathbf{z}^T \mathbf{x} \right\} = \left\{ \mathbf{y} : \|\mathbf{y}\|_q \leq 1 \text{ and } \mathbf{y}^T \mathbf{x} = \|\mathbf{x}\|_p \right\}$$

This is derived in  
class

Why  $\|\mathbf{y}\|_q \leq 1$   
is because of Minkowski's  
inequality

# Subgradients for the 'Lasso' Problem in Machine Learning

We use Lasso ( $\min_{\mathbf{x}} f(\mathbf{x})$ ) as an example to illustrate subgradients of affine composition:

$$f(\mathbf{x}) = \frac{1}{2} \|\mathbf{y} - \mathbf{x}\|^2 + \lambda \|\mathbf{x}\|_1$$

The subgradients of  $f(\mathbf{x})$  are

$$\mathbf{x} - \mathbf{y} + \lambda \mathbf{s}$$

Where  $\mathbf{s} = \{+1, -1\}^n$   
such that  $\|\mathbf{x}\|_1 = \mathbf{s}^T \mathbf{x}$

# Subgradients for the 'Lasso' Problem in Machine Learning

We use Lasso ( $\min_{\mathbf{x}} f(\mathbf{x})$ ) as an example to illustrate subgradients of affine composition:

$$f(\mathbf{x}) = \frac{1}{2} \|\mathbf{y} - \mathbf{x}\|^2 + \lambda \|\mathbf{x}\|_1$$

The subgradients of  $f(\mathbf{x})$  are

$$\mathbf{h} = \mathbf{x} - \mathbf{y} + \lambda \mathbf{s},$$

where  $\mathbf{s}_i = \text{sign}(x_i)$  if  $x_i \neq 0$  and  $\mathbf{s}_i \in [-1, 1]$  if  $x_i = 0$ .

Second component is a result of the convex hull

## More Subgradient Calculus: Composition

Following functions, though convex, may not be differentiable everywhere. How does one compute their subgradients? (what holds for subgradient also holds for gradient)

- **Composition with functions:** Let  $p : \mathbb{R}^k \rightarrow \mathbb{R}$  with  $q(x) = \infty, \forall \mathbf{x} \notin \text{dom } h$  and  $q : \mathbb{R}^n \rightarrow \mathbb{R}^k$ . Define  $f(\mathbf{x}) = p(q(\mathbf{x}))$ .  $f$  is convex if
  - ▶  $q_i$  is convex,  $p$  is convex and nondecreasing in each argument
  - ▶ or  $q_i$  is concave,  $p$  is convex and nonincreasing in each argument

We will consider only the first case



## More Subgradient Calculus: Composition

Following functions, though convex, may not be differentiable everywhere. How does one compute their subgradients? (what holds for subgradient also holds for gradient)

- **Composition with functions:** Let  $p: \mathbb{R}^k \rightarrow \mathbb{R}$  with  $q(x) = \infty, \forall \mathbf{x} \notin \text{dom } h$  and  $q: \mathbb{R}^n \rightarrow \mathbb{R}^k$ . Define  $f(\mathbf{x}) = p(q(\mathbf{x}))$ .  $f$  is convex if **In both conditions, composition will be concave if  $p$  is concave**
  - ▶  $q_i$  is convex,  $p$  is convex and nondecreasing in each argument
  - ▶ or  $q_i$  is concave,  $p$  is convex and nonincreasing in each argument

Some examples illustrating this property are:

- ▶  $\exp q(\mathbf{x})$  is convex if  $q$  is convex **exp is a monotonic and convex  $p$**
- ▶  $\sum_{i=1}^m \log q_i(\mathbf{x})$  is concave if  $q_i$  are concave and positive  **$p$  is concave and hence the composition is concave**
- ▶  $\log \sum_{i=1}^m \exp q_i(\mathbf{x})$  is convex if  $q_i$  are convex
- ▶  $1/q(\mathbf{x})$  is convex if  $q$  is concave and positive

## More Subgradient Calculus: Composition (contd)

- **Composition with functions:** Let  $p: \mathbb{R}^k \rightarrow \mathbb{R}$  with  $q(x) = \infty, \forall \mathbf{x} \notin \text{dom } h$  and  $q: \mathbb{R}^n \rightarrow \mathbb{R}^k$ . Define  $f(\mathbf{x}) = p(q(\mathbf{x}))$ .  $f$  is convex if
  - ▶  $q_i$  is convex,  $p$  is convex and nondecreasing in each argument
  - ▶ or  $q_i$  is concave,  $p$  is convex and nonincreasing in each argument
- Subgradients for the first case (second one is homework):

## More Subgradient Calculus: Composition (contd)

- **Composition with functions:** Let  $p : \mathbb{R}^k \rightarrow \mathbb{R}$  with  $q(x) = \infty, \forall \mathbf{x} \notin \text{dom } h$  and  $q : \mathbb{R}^n \rightarrow \mathbb{R}^k$ . Define  $f(\mathbf{x}) = p(q(\mathbf{x}))$ .  $f$  is convex if
  - ▶  $q_i$  is convex,  $p$  is convex and nondecreasing in each argument
  - ▶ or  $q_i$  is concave,  $p$  is convex and nonincreasing in each argument
- Subgradients for the first case (second one is homework):
  - ▶  $f(\mathbf{y}) = p(q_1(\mathbf{y}), \dots, q_k(\mathbf{y})) \geq p(q_1(\mathbf{x}) + \mathbf{h}_{q_1}^T(\mathbf{y} - \mathbf{x}), \dots, q_k(\mathbf{x}) + \mathbf{h}_{q_k}^T(\mathbf{y} - \mathbf{x}))$   
Where  $\mathbf{h}_{q_i} \in \partial q_i(\mathbf{x})$  for  $i = 1..k$  and since  $p(\cdot)$  is non-decreasing in each argument.

$p$  applied to  $q_i(x)$  is  $\geq$   $p$  applied to the lower bounds on  $q_i(x)$

## More Subgradient Calculus: Composition (contd)

- **Composition with functions:** Let  $p: \mathbb{R}^k \rightarrow \mathbb{R}$  with  $q(x) = \infty, \forall \mathbf{x} \notin \text{dom } h$  and  $q: \mathbb{R}^n \rightarrow \mathbb{R}^k$ . Define  $f(\mathbf{x}) = p(q(\mathbf{x}))$ .  $f$  is convex if
  - ▶  $q_i$  is convex,  $p$  is convex and nondecreasing in each argument
  - ▶ or  $q_i$  is concave,  $p$  is convex and nonincreasing in each argument
- Subgradients for the first case (second one is homework):

$$\text{▶ } f(\mathbf{y}) = p(q_1(\mathbf{y}), \dots, q_k(\mathbf{y})) \geq p(q_1(\mathbf{x}) + \mathbf{h}_{q_1}^T(\mathbf{y} - \mathbf{x}), \dots, q_k(\mathbf{x}) + \mathbf{h}_{q_k}^T(\mathbf{y} - \mathbf{x}))$$

Where  $\mathbf{h}_{q_i} \in \partial q_i(\mathbf{x})$  for  $i = 1..k$  and since  $p(\cdot)$  is non-decreasing in each argument.

$$\text{▶ } p(q_1(\mathbf{x}) + \mathbf{h}_{q_1}^T(\mathbf{y} - \mathbf{x}), \dots, q_k(\mathbf{x}) + \mathbf{h}_{q_k}^T(\mathbf{y} - \mathbf{x})) \geq p(q_1(\mathbf{x}), \dots, q_k(\mathbf{x})) + \mathbf{h}_p^T(\mathbf{h}_{q_1}^T(\mathbf{y} - \mathbf{x}), \dots, \mathbf{h}_{q_k}^T(\mathbf{y} - \mathbf{x}))$$

Where  $\mathbf{h}_p \in \partial p(q_1(\mathbf{x}), \dots, q_k(\mathbf{x}))$

All we need to do next is club together  $\mathbf{h}_p$  and  $\mathbf{h}_q$  and leave only  $(\mathbf{y}-\mathbf{x})$  in the second component

## More Subgradient Calculus: Composition (contd)

- **Composition with functions:** Let  $p: \mathbb{R}^k \rightarrow \mathbb{R}$  with  $q(x) = \infty, \forall \mathbf{x} \notin \text{dom } h$  and  $q: \mathbb{R}^n \rightarrow \mathbb{R}^k$ . Define  $f(\mathbf{x}) = p(q(\mathbf{x}))$ .  $f$  is convex if

- ▶  $q_i$  is convex,  $p$  is convex and nondecreasing in each argument
- ▶ or  $q_i$  is concave,  $p$  is convex and nonincreasing in each argument

- Subgradients for the first case (second one is homework):

- ▶  $f(\mathbf{y}) = p(q_1(\mathbf{y}), \dots, q_k(\mathbf{y})) \geq p(q_1(\mathbf{x}) + \mathbf{h}_{q_1}^T(\mathbf{y} - \mathbf{x}), \dots, q_k(\mathbf{x}) + \mathbf{h}_{q_k}^T(\mathbf{y} - \mathbf{x}))$

Where  $\mathbf{h}_{q_i} \in \partial q_i(\mathbf{x})$  for  $i = 1..k$  and since  $p(\cdot)$  is non-decreasing in each argument.

- ▶  $p(q_1(\mathbf{x}) + \mathbf{h}_{q_1}^T(\mathbf{y} - \mathbf{x}), \dots, q_k(\mathbf{x}) + \mathbf{h}_{q_k}^T(\mathbf{y} - \mathbf{x})) \geq$

$$p(q_1(\mathbf{x}), \dots, q_k(\mathbf{x})) + \mathbf{h}_p^T(\mathbf{h}_{q_1}^T(\mathbf{y} - \mathbf{x}), \dots, \mathbf{h}_{q_k}^T(\mathbf{y} - \mathbf{x}))$$

Where  $\mathbf{h}_p \in \partial p(q_1(\mathbf{x}), \dots, q_k(\mathbf{x}))$

- ▶  $p(q_1(\mathbf{x}), \dots, q_k(\mathbf{x})) + \mathbf{h}_p^T(\mathbf{h}_{q_1}^T(\mathbf{y} - \mathbf{x}), \dots, \mathbf{h}_{q_k}^T(\mathbf{y} - \mathbf{x})) = f(\mathbf{x}) + \sum_{i=1}^k (h_p)_i \mathbf{h}_{q_i}^T(\mathbf{y} - \mathbf{x})$

That is,  $\sum_{i=1}^k (h_p)_i \mathbf{h}_{q_i}$  is a subgradient of the composite function at  $\mathbf{x}$ .

H/W: Derive the subdifferentials to example functions on previous slide

## More Subgradient Calculus: Proximal Operator

Following functions are again convex, but again, may not be differentiable everywhere. How does one compute their subgradients at points of non-differentiability?

- **Infimum:** If  $c(x, y)$  is convex in  $(x, y)$  and  $\mathcal{C}$  is a convex set, then  $d(x) = \inf_{y \in \mathcal{C}} c(x, y)$  is

convex. For example:

- ▶ Let  $d(\mathbf{x}, \mathcal{C})$  that returns the distance of a point  $\mathbf{x}$  to a convex set  $\mathcal{C}$ . That is  $d(\mathbf{x}, \mathcal{C}) = \inf_{y \in \mathcal{C}} \|\mathbf{x} - \mathbf{y}\| = \|\mathbf{x} - P_{\mathcal{C}}(\mathbf{x})\|$ , where,  $P_{\mathcal{C}}(\mathbf{x}) = \operatorname{argmin} d(\mathbf{x}, \mathcal{C})$ . Then  $d(\mathbf{x}, \mathcal{C})$  is a

convex function and  $\nabla d(\mathbf{x}, \mathcal{C}) = \frac{\mathbf{x} - P_{\mathcal{C}}(\mathbf{x})}{\|\mathbf{x} - P_{\mathcal{C}}(\mathbf{x})\|}$

H/w: Prove that  $d$  is convex if  $c$  is a convex function  
and if  $\mathcal{C}$  is a convex set

## More Subgradient Calculus: Proximal Operator

Following functions are again convex, but again, may not be differentiable everywhere. How does one compute their subgradients at points of non-differentiability?

- **Infimum:** If  $c(x, y)$  is convex in  $(x, y)$  and  $\mathcal{C}$  is a convex set, then  $d(x) = \inf_{y \in \mathcal{C}} c(x, y)$  is

convex. For example:

- ▶ Let  $d(\mathbf{x}, \mathcal{C})$  that returns the distance of a point  $\mathbf{x}$  to a convex set  $\mathcal{C}$ . That is  $d(\mathbf{x}, \mathcal{C}) = \inf_{\mathbf{y} \in \mathcal{C}} \|\mathbf{x} - \mathbf{y}\| = \|\mathbf{x} - P_{\mathcal{C}}(\mathbf{x})\|$ , where,  $P_{\mathcal{C}}(\mathbf{x}) = \operatorname{argmin}_{\mathbf{y} \in \mathcal{C}} \|\mathbf{x} - \mathbf{y}\|$ . Then  $d(\mathbf{x}, \mathcal{C})$  is a

convex function and  $\nabla d(\mathbf{x}, \mathcal{C}) = \frac{\mathbf{x} - P_{\mathcal{C}}(\mathbf{x})}{\|\mathbf{x} - P_{\mathcal{C}}(\mathbf{x})\|}$  .... The point of intersection of convex sets  $\mathcal{C}_1, \mathcal{C}_2, \dots, \mathcal{C}_m$  by minimizing...

## More Subgradient Calculus: Proximal Operator

Following functions are again convex, but again, may not be differentiable everywhere. How does one compute their subgradients at points of non-differentiability?

- **Infimum:** If  $c(x, y)$  is convex in  $(x, y)$  and  $\mathcal{C}$  is a convex set, then  $d(x) = \inf_{y \in \mathcal{C}} c(x, y)$  is

convex. For example:

- ▶ Let  $d(\mathbf{x}, \mathcal{C})$  that returns the distance of a point  $\mathbf{x}$  to a convex set  $\mathcal{C}$ . That is  $d(\mathbf{x}, \mathcal{C}) = \inf_{y \in \mathcal{C}} \|\mathbf{x} - \mathbf{y}\| = \|\mathbf{x} - P_{\mathcal{C}}(\mathbf{x})\|$ , where,  $P_{\mathcal{C}}(\mathbf{x}) = \operatorname{argmin}_{y \in \mathcal{C}} d(\mathbf{x}, \mathcal{C})$ . Then  $d(\mathbf{x}, \mathcal{C})$  is a convex function and  $\nabla d(\mathbf{x}, \mathcal{C}) = \frac{\mathbf{x} - P_{\mathcal{C}}(\mathbf{x})}{\|\mathbf{x} - P_{\mathcal{C}}(\mathbf{x})\|}$  ....The point of intersection of convex sets  $\mathcal{C}_1, \mathcal{C}_2, \dots, \mathcal{C}_m$  by minimizing... (Subgradients and Alternating Projections)
- ▶  $\operatorname{argmin}_{y \in \mathcal{C}} d(\mathbf{x}, \mathcal{C})$  is a special case of the proximity operator:  $\operatorname{prox}_c(\mathbf{x}) = \operatorname{argmin}_y PROX_c(\mathbf{x})$  of a convex function  $c(\mathbf{x})$ . Here,  $PROX_c(\mathbf{x}) = c(\mathbf{y}) + \frac{1}{2}\|\mathbf{x} - \mathbf{y}\|^2$  The special case is when

$c(\mathbf{x})$  is the indicator function over  $\mathcal{C}$



## More Subgradient Calculus: Proximal Operator

Following functions are again convex, but again, may not be differentiable everywhere. How does one compute their subgradients at points of non-differentiability?

- **Infimum:** If  $c(x, y)$  is convex in  $(x, y)$  and  $\mathcal{C}$  is a convex set, then  $d(x) = \inf_{y \in \mathcal{C}} c(x, y)$  is

convex. For example:

- ▶ Let  $d(\mathbf{x}, \mathcal{C})$  that returns the distance of a point  $\mathbf{x}$  to a convex set  $\mathcal{C}$ . That is  $d(\mathbf{x}, \mathcal{C}) = \inf_{y \in \mathcal{C}} \|\mathbf{x} - \mathbf{y}\| = \|\mathbf{x} - P_{\mathcal{C}}(\mathbf{x})\|$ , where,  $P_{\mathcal{C}}(\mathbf{x}) = \operatorname{argmin}_{y \in \mathcal{C}} d(\mathbf{x}, \mathcal{C})$ . Then  $d(\mathbf{x}, \mathcal{C})$  is a

convex function and  $\nabla d(\mathbf{x}, \mathcal{C}) = \frac{\mathbf{x} - P_{\mathcal{C}}(\mathbf{x})}{\|\mathbf{x} - P_{\mathcal{C}}(\mathbf{x})\|}$  ....The point of intersection of convex sets

$\mathcal{C}_1, \mathcal{C}_2, \dots, \mathcal{C}_m$  by minimizing... (Subgradients and Alternating Projections)

- ▶  $\operatorname{argmin}_{y \in \mathcal{C}} d(\mathbf{x}, \mathcal{C})$  is a special case of the proximity operator:  $\operatorname{prox}_{c(\cdot)}(\mathbf{x}) = \operatorname{argmin}_y \operatorname{PROX}_c(\mathbf{x})$  of a

convex function  $c(\mathbf{x})$ . Here,  $\operatorname{PROX}_c(\mathbf{x}) = c(\mathbf{y}) + \frac{1}{2}\|\mathbf{x} - \mathbf{y}\|^2$  The special case is when  $c(\mathbf{y})$  is

the indicator function  $I_{\mathcal{C}}(\mathbf{y})$  introduced earlier to eliminate the constraints of an optimization problem.

- ★ Recall that  $\partial I_{\mathcal{C}}(\mathbf{y}) = N_{\mathcal{C}}(\mathbf{y}) = \{\mathbf{h} \in \mathbb{R}^n : \mathbf{h}^T \mathbf{y} \geq \mathbf{h}^T \mathbf{z} \text{ for any } \mathbf{z} \in \mathcal{C}\}$
- ★ The subdifferential  $\partial \operatorname{PROX}_c(\mathbf{x}) = \partial c(\mathbf{y}) + \mathbf{y} - \mathbf{x}$  which can now be obtained for the special case  $c(\mathbf{y}) = I_{\mathcal{C}}(\mathbf{y})$ .

Proximal  
will be done  
in details  
later