

## RECAP

# Demystifying the Projection Step

$$\begin{aligned} \mathbf{x}_p^{(k+1)} &= P_C(\mathbf{x}_u^{(k+1)}) \\ &= \operatorname{argmin}_{\mathbf{z} \in C} \left\| \mathbf{x}_u^{(k+1)} - \mathbf{z} \right\|_2^2 \\ &= \operatorname{argmin}_{\mathbf{z} \in C} \left\| \mathbf{x}_u^{(k+1)} - \mathbf{z} \right\|_2^2 + l_C(\mathbf{z}) \\ &= \operatorname{argmin}_{\mathbf{z} \in C} \frac{1}{2} \left\| \mathbf{x}_u^{(k+1)} - \mathbf{z} \right\|_2^2 \end{aligned}$$

- Solution set of a linear system  $\mathcal{C} = \{\mathbf{x} \in \mathbb{R}^n : A^T \mathbf{x} = \mathbf{b}\}$
- Affine images  $\mathcal{C} = \{A\mathbf{x} + \mathbf{b} : \mathbf{x} \in \mathbb{R}^n\}$
- Nonnegative orthant  $\mathcal{C} = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{x} \succeq 0\}$ . It may be hard to project on arbitrary polyhedron.
- Norm balls  $\mathcal{C} = \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\|_p \leq 1\}$ , for  $p = 1, 2, \infty$

Solution set of a linear system  $\mathcal{C} = \{\mathbf{x} \in \mathbb{R}^n : A^T \mathbf{x} = \mathbf{b}\}$

$$\mathbf{x}_p^{(k+1)} = P_{\mathcal{C}}(\mathbf{x}_u^{(k+1)}) = \arg \min_{A^T \mathbf{z} = \mathbf{b}} \frac{1}{2} \|\mathbf{x}_u^{(k+1)} - \mathbf{z}\|_2^2$$

For  $\mathbf{z}, \mathbf{x} \in \mathbb{R}^n$ ,  $A$  as an  $n \times m$  matrix,  $\mathbf{b}$  is a vector of size  $m$ , consider the slightly more general problem (58) with  $B$  as an  $n \times n$  matrix:

$$\begin{aligned} \min_{\mathbf{z} \in \mathbb{R}^n} \quad & \frac{1}{2} (\mathbf{z} - \mathbf{x})^T B (\mathbf{z} - \mathbf{x}) \\ \text{subject to} \quad & A^T \mathbf{z} = \mathbf{b} \end{aligned} \tag{58}$$

For projected gradient descent,  $B =$

# Projected Gradient Descent for Affine Constraint Set $\mathcal{C}$

RECAP

Solution set of a linear system  $\mathcal{C} = \{\mathbf{x} \in \mathbb{R}^n : A^T \mathbf{x} = \mathbf{b}\}$

$$\mathbf{x}_p^{(k+1)} = P_{\mathcal{C}}(\mathbf{x}_u^{(k+1)}) = \arg \min_{A^T \mathbf{z} = \mathbf{b}} \frac{1}{2} \|\mathbf{x}_u^{(k+1)} - \mathbf{z}\|_2^2$$

For  $\mathbf{z}, \mathbf{x} \in \mathbb{R}^n$ ,  $A$  as an  $n \times m$  matrix,  $\mathbf{b}$  is a vector of size  $m$ , consider the slightly more general problem (58) with  $B$  as an  $n \times n$  matrix:

$$\begin{aligned} \min_{\mathbf{z} \in \mathbb{R}^n} \quad & \frac{1}{2} (\mathbf{z} - \mathbf{x})^T B (\mathbf{z} - \mathbf{x}) \\ \text{subject to} \quad & A^T \mathbf{z} = \mathbf{b} \end{aligned} \tag{58}$$

For projected gradient descent,  $B = I$ . Further, if  $n = 2$  and  $m = 1$ , the minimization problem (58) amounts to finding a point  $\mathbf{y}^*$  on a line  $a_{11}z_1 + a_{12}z_2 = b$  that is closest to  $\mathbf{x}$ .

- Consider minimization of the modified objective function

$$L(\mathbf{z}, \lambda) = \frac{1}{2}(\mathbf{z} - \mathbf{x})^T B(\mathbf{z} - \mathbf{x}) + \lambda^T (A^T \mathbf{z} - \mathbf{b}).$$

$$\min_{\mathbf{z} \in \mathbb{R}^n, \lambda \in \mathbb{R}^m} \frac{1}{2}(\mathbf{z} - \mathbf{x})^T B(\mathbf{z} - \mathbf{x}) + \lambda^T (A^T \mathbf{z} - \mathbf{b}) \quad (59)$$

The function  $L(\mathbf{z}, \lambda)$  is called the lagrangian and involves the lagrange multiplier  $\lambda \in \mathbb{R}^m$ .

- A sufficient condition for optimality of  $L(\mathbf{z}^*, \lambda^*)$  is that  $\nabla L(\mathbf{z}^*, \lambda^*) = 0$  and  $\nabla^2 L(\mathbf{z}^*, \lambda^*) \succ 0$ . For this specific problem:

$$\nabla L(\mathbf{z}^*, \lambda^*) = \begin{bmatrix} B\mathbf{z}^* - \frac{1}{2}(B + B^T)\mathbf{x} + A\lambda^* \\ A^T \mathbf{z}^* - \mathbf{b} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

and

$$\nabla^2 L(\mathbf{z}^*, \lambda^*) = \begin{bmatrix} B & A \\ A^T & 0 \end{bmatrix} \succ 0$$

- The point  $(\mathbf{z}^*, \lambda^*)$  must therefore satisfy,  $A^T \mathbf{z}^* = \mathbf{b}$  and  $A\lambda^* = -B\mathbf{z}^* + \frac{1}{2}(B + B^T)\mathbf{x}$ .
- Recap: If  $B$  is taken to be the identity matrix,  $n = 2$  and  $m = 1$ , the minimization problem (58) amounts to finding a point  $\mathbf{y}^*$  on a line  $a_{11}z_1 + a_{12}z_2 = b$  that is closest to  $\mathbf{x}$ .
- From geometry, the point on a line closest to  $\mathbf{x}$  is the point of intersection  $\mathbf{p}^*$  of a perpendicular (or least possible<sup>8</sup> obtuse angle) from  $\mathbf{x}$  to the line. However, the solution for the minimum of (59), for these conditions coincides with  $\mathbf{p}^*$  and is given by:

$$z_1^* = x_1 - \frac{a_{11}(a_{11}x_1 + a_{12}x_2 - b)}{(a_{11})^2 + (a_{12})^2} \quad z_2^* = x_2 - \frac{a_{12}(a_{11}x_1 + a_{12}x_2 - b)}{(a_{11})^2 + (a_{12})^2}$$

That is, for  $n = 2$  and  $m = 1$ , the solution to (59) is the same as the solution to (58)

- For general  $n$  and  $m$ ,

$$\mathbf{z}^* = \mathbf{x}_p^{(k+1)} = P_{\mathcal{C}}(\mathbf{x}_u^{(k+1)}) = \arg \min_{A^T \mathbf{z} = \mathbf{b}} \frac{1}{2} \left\| \mathbf{x}_u^{(k+1)} - \mathbf{z} \right\|_2^2 = \mathbf{x}_u^{(k+1)} - A(A^T A)^{-1}(A^T \mathbf{x}_u^{(k+1)} - \mathbf{b})$$

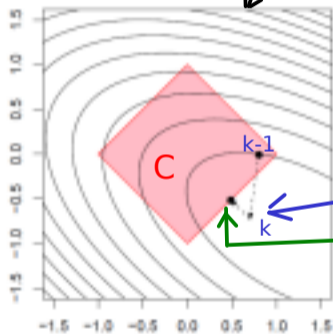
Was this an accident?

More today!

<sup>8</sup>See following slides for some elaboration on geometry of the projection

# Projected Gradient Descent: Illustrated and Summarized

## Level surfaces for the quadratic objective



- Illustration of Projected Gradient Descent on Quadratic Objective with bounded affine (Polyhedral) constraint set
- The line joining point of projection  $\mathbf{x}_p^k = P_C(\mathbf{x}_u^k)$  to  $\mathbf{x}_u^k$  forms least possible obtuse angle<sup>a</sup> with line joining  $\mathbf{x}_p^k = P_C(\mathbf{x}_u^k)$  to any point  $\mathbf{z} \in C$ .

<sup>a</sup>See following slides for some elaboration on geometry of the projection

Elaboration on the Geometry of the Projected  
Gradient Descent  
Right angle FOR Affine Set/Unbounded sets  
Least possible obtuse angle FOR  
Polyhedron/Bounded Sets

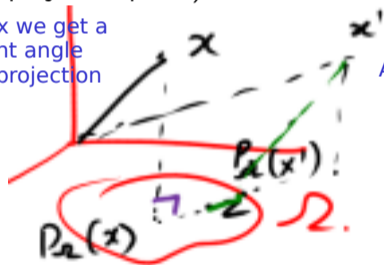


- **Claim:** If  $P_C(\mathbf{x})$  is a projection of  $\mathbf{x}$ , then

$$(\mathbf{z} - P_C(\mathbf{x}))^\top (\mathbf{x} - P_C(\mathbf{x})) \leq 0, \forall \mathbf{z} \in C$$

- That is, the angle between  $(\mathbf{z} - P_C(\mathbf{x}))$  and  $(\mathbf{x} - P_C(\mathbf{x}))$  is obtuse (or right-angled for the projected point),  $\forall \mathbf{z} \in C$

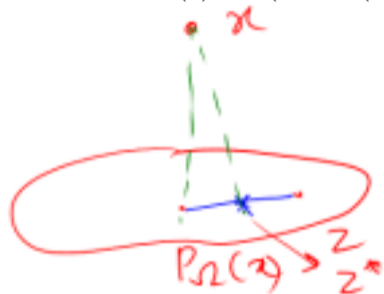
At  $\mathbf{x}$  we get a right angle at projection



At  $\mathbf{x}'$  we make an obtuse angle at projection

## Proof for $\langle z - P_C(x), x - P_C(x) \rangle \leq 0$

- To be more general, let us consider an inner product  $\langle a, b \rangle$  instead of  $a^\top b$
- Let  $z^* = (1 - \alpha)P_C(x) + \alpha z$ , for some  $\alpha \in (0, 1)$ , and  $z \in C$   
 $\implies z^* = P_C(x) + \alpha(z - P_C(x)), z^* \in C$



- Since  $P_C(x) = \operatorname{argmin}_{z \in C} \|x - z\|_2^2$ ,  
 $\|x - P_C(x)\|^2 \leq \|x - z^*\|^2$

$$\begin{aligned}
& \|x - z^*\|^2 \\
&= \left\| x - (P_C(x) + \alpha(z - P_C(x))) \right\|^2 \\
&= \|x - P_C(x)\|^2 + \alpha^2 \|z - P_C(x)\|^2 - 2\alpha \langle x - P_C(x), z - P_C(x) \rangle \\
&\geq \|x - P_C(x)\|^2
\end{aligned}$$

$$\implies \langle x - P_C(x), z - P_C(x) \rangle \leq \frac{\alpha}{2} \|z - P_C(x)\|^2, \forall \alpha \in (0, 1)$$

- Thus, the LHS can either be 0 or a negative value. Any positive value of the LHS will lead to a contradiction for some small  $\alpha \rightarrow 0$
- Hence, we proved that  $\langle z - P_C(x), x - P_C(x) \rangle \leq 0$

If  $x^*$  is in the set  $C$ , it itself must be the projection

- We can also prove that if  $\langle x - x^*, z - x^* \rangle \leq 0, \forall z \in C$  s.t.  $z \neq x^*$ , and  $x^* \in C$ , then

$$x^* = P_C(x) = \operatorname{argmin}_{\bar{z} \in C} \|x - \bar{z}\|_2^2$$

- Consider  $\|x - z\|^2 - \|x - x^*\|^2$   
 $= \|x - x^* + (x^* - z)\|^2 - \|x - x^*\|^2$   
 $= \|x - x^*\|^2 + \|z - x^*\|^2 - 2 \langle x - x^*, z - x^* \rangle - \|x - x^*\|^2$   
 $= \|z - x^*\|^2 - 2 \langle x - x^*, z - x^* \rangle$   
 $> 0$
- $\implies \|x - z\|^2 > \|x - x^*\|^2, \forall z \in C$  s.t.  $z \neq x^*$
- This proves that  $x^* = P_C(x)$

# Lagrange Function and KKT Conditions

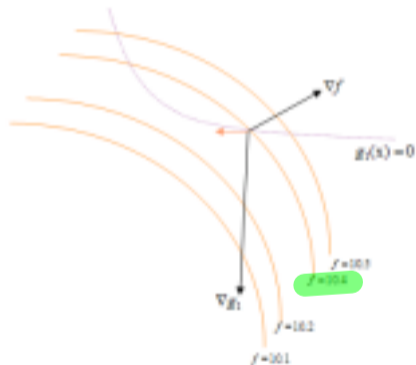
# Lagrange Function and Necessary KKT Conditions

- Can the Lagrange Multiplier construction be generalized to always find optimal solutions to a minimization problem?
- Instead of the iterative path again, assume everything can be computed analytically
- Attributed to the mathematician Lagrange (born in 1736 in Turin). Largely worked on mechanics, the calculus of variations probability, group theory, and number theory. Attributed choice of base 10 for the metric system (rather than 12).

Projected gradient descent is only one consumer for this analysis.  
There are several other results and algorithms that make use of this analysis

# Lagrange Function and Necessary KKT Conditions

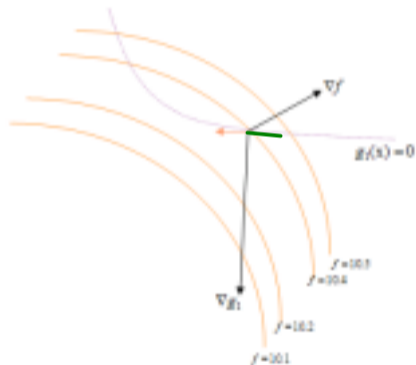
- Consider the equality constrained minimization problem (with  $\mathcal{D} \subseteq \mathbb{R}^n$ )



$$\begin{array}{ll} \min_{\mathbf{x} \in \mathcal{D}} & f(\mathbf{x}) \\ \text{subject to} & g_i(\mathbf{x}) = 0 \quad i = 1, 2, \dots, m \end{array} \quad (60)$$

- The figure shows some level curves of the function  $f$  and of a single constraint function  $g_1$  (dotted lines)
- The gradient of the constraint  $\nabla g_1$  is not parallel to the gradient  $\nabla f$  of the function at  $f = 10.4$ ; it is therefore possible to decrease  $f$  while maintaining  $g_1(\mathbf{x}) = 0$  (by moving tangential to  $g_1(\mathbf{x}) = 0$ )

# Lagrange Function and Necessary KKT Conditions



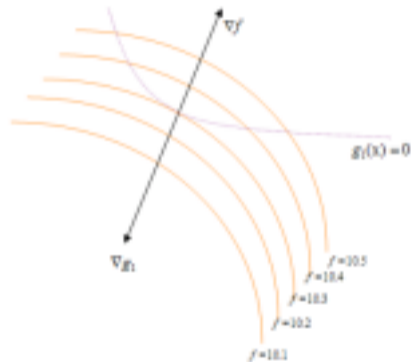
- Consider the equality constrained minimization problem (with  $\mathcal{D} \subseteq \mathbb{R}^n$ )

$$\begin{array}{ll} \min_{\mathbf{x} \in \mathcal{D}} & f(\mathbf{x}) \\ \text{subject to} & g_i(\mathbf{x}) = 0 \quad i = 1, 2, \dots, m \end{array} \quad (60)$$

- The figure shows some level curves of the function  $f$  and of a single constraint function  $g_1$  (dotted lines)
- The gradient of the constraint  $\nabla g_1$  is not parallel to the gradient  $\nabla f$  of the function at  $f = 10.4$ ; it is therefore possible to **move along the constraint surface so as to further reduce  $f$ .**



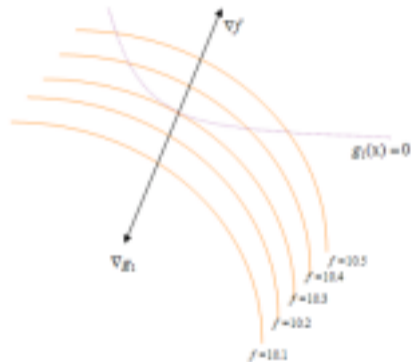
# Lagrange Function and Necessary KKT Conditions



- However,  $\nabla g_1$  and  $\nabla f$  are parallel at  $f = 10.3$ , and any motion along  $g_1(\mathbf{x}) = 0$  will not change the value of  $f(\mathbf{x})$  since gradient of  $f$  has no component perpendicular to the gradient of  $g_1(\mathbf{x}) = 0$

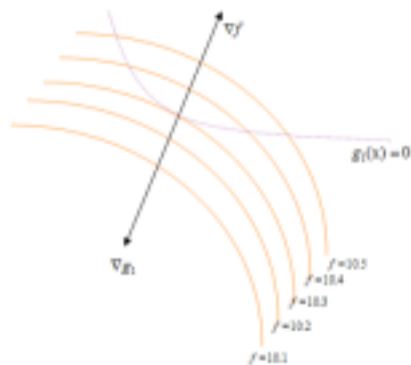
At  $\mathbf{x}$  s.t  $f(\mathbf{x}) = 10.3$ ,  
gradient of  $f = \lambda$  gradient of  $g_1$   
sign of  $\lambda$  does not matter

# Lagrange Function and Necessary KKT Conditions



- However,  $\nabla g_1$  and  $\nabla f$  are parallel at  $f = 10.3$ , and any motion along  $g_1(\mathbf{x}) = 0$  will **increase  $f$ , or leave it unchanged**.
- Hence, at the solution  $\mathbf{x}^*$ ,

# Lagrange Function and Necessary KKT Conditions



- However,  $\nabla g_1$  and  $\nabla f$  are parallel at  $f = 10.3$ , and any motion along  $g_1(\mathbf{x}) = 0$  will **increase  $f$ , or leave it unchanged**.
- Hence, at the solution  $\mathbf{x}^*$ ,  $\nabla f(\mathbf{x}^*)$  must be proportional to  $-\nabla g_1(\mathbf{x}^*)$ , yielding,  $\nabla f(\mathbf{x}^*) = -\lambda \nabla g_1(\mathbf{x}^*)$ , for some constant  $\lambda \in \mathfrak{R}$ ;  $\lambda$  is called a *Lagrange multiplier*.
- Often  $\lambda$  itself need never be computed and therefore often qualified as the *undetermined* lagrange multiplier.

## Lagrange Function and Necessary KKT Conditions

- The necessary condition for an optimum at  $\mathbf{x}^*$  for the optimization problem in (60) with  $m = 1$  can be stated as in (61); the gradient is now in

## Lagrange Function and Necessary KKT Conditions

- The necessary condition for an optimum at  $\mathbf{x}^*$  for the optimization problem in (60) with  $m = 1$  can be stated as in (61); the gradient is now in  $\Re^{n+1}$  with its last component being a partial derivative with respect to  $\lambda$ .

$$\nabla L(\mathbf{x}^*, \lambda^*) = \nabla f(\mathbf{x}^*) + \lambda^* \nabla g_1(\mathbf{x}^*) = 0 \quad (61)$$

- The solutions to (61) are the stationary points of the Lagrangian  $L$ ; they are not necessarily local extrema of  $L$ .  $L$  is unbounded: given a point  $\mathbf{x}$  that doesn't lie on the constraint, letting  $\lambda \rightarrow \pm\infty$  makes  $L$  arbitrarily large or small. However, under certain stronger assumptions, if the *strong Lagrangian principle* holds, the minima of  $f$  minimize the Lagrangian globally.

## Lagrange Function and Necessary KKT Conditions

- Let us extend the necessary condition for optimality of a minimization problem with single constraint to minimization problems with multiple equality constraints (*i.e.*,  $m > 1$ . in (60)).
- Let  $\mathcal{S}$  be the subspace spanned by  $\nabla g_i(\mathbf{x})$  at any point  $\mathbf{x}$  and let  $\mathcal{S}_\perp$  be its orthogonal complement. Let  $(\nabla f)_\perp$  be the component of  $\nabla f$  in the subspace  $\mathcal{S}_\perp$ .

There is no component of gradient  $f$  perpendicular to  $S$   
SAME AS  
gradient of  $f$  lies in  $S$

## Lagrange Function and Necessary KKT Conditions

- Let us extend the necessary condition for optimality of a minimization problem with single constraint to minimization problems with multiple equality constraints (*i.e.*,  $m > 1$ . in (60)).
- Let  $\mathcal{S}$  be the subspace spanned by  $\nabla g_i(\mathbf{x})$  at any point  $\mathbf{x}$  and let  $\mathcal{S}_\perp$  be its orthogonal complement. Let  $(\nabla f)_\perp$  be the component of  $\nabla f$  in the subspace  $\mathcal{S}_\perp$ .
- At any solution  $\mathbf{x}^*$ , it must be true that the gradient of  $f$  has  $(\nabla f)_\perp = 0$  (*i.e.*, no components that are perpendicular to all of the  $\nabla g_i$ ), because otherwise you could move  $\mathbf{x}^*$  a little in that direction (or in the opposite direction) to increase (decrease)  $f$  without changing any of the  $g_i$ , *i.e.* without violating any constraints.
- Hence for multiple equality constraints, it must be true that at the solution  $\mathbf{x}^*$ , the space  $\mathcal{S}$  contains the vector  $\nabla f$ , *i.e.*, there are some constants  $\lambda_i$  such that  $\nabla f(\mathbf{x}^*) = \lambda_i \nabla g_i(\mathbf{x}^*)$ .

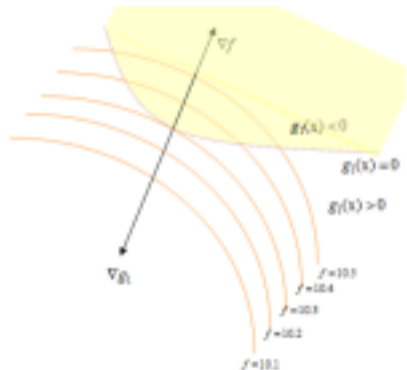
## Lagrange Multipliers with Inequality Constraints

- We also need to impose that the solution is on the correct constraint surface (*i.e.*,  $g_i = 0, \forall i$ ). In the same manner as in the case of  $m = 1$ , this can be encapsulated by introducing the Lagrangian  $L(\mathbf{x}, \lambda) = f(\mathbf{x}) - \sum_{i=1}^m \lambda_i g_i(\mathbf{x})$ , whose gradient with respect to both  $\mathbf{x}$ , and  $\lambda$  vanishes at the solution.
- This gives us the following necessary condition for optimality of (60):

$$\nabla L(\mathbf{x}^*, \lambda^*) = \nabla \left( f(\mathbf{x}) - \sum_{i=1}^m \lambda_i g_i(\mathbf{x}) \right) = 0 \quad (62)$$



# Lagrange Multipliers with Inequality Constraints



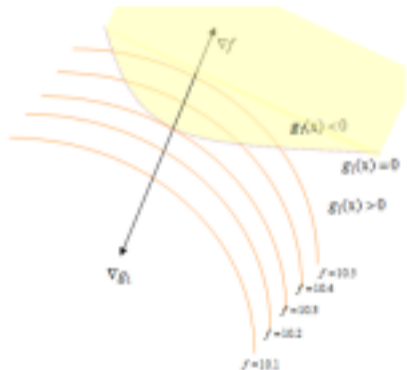
- Single equality constraint  $g_1(\mathbf{x}) = 0$ , replaced with a single inequality constraint  $g_1(\mathbf{x}) \leq 0$ . The entire region labeled  $g_1(\mathbf{x}) \leq 0$  in the Figure becomes feasible.
- At the solution  $\mathbf{x}^*$ , if  $g_1(\mathbf{x}^*) = 0$ , *i.e.*, if the constraint is active, we must have

gradient of  $f(\mathbf{x}^*)$  has no component  
perpendicular to gradient  $g_1(\mathbf{x}^*)$   
AND

- gradient of  $f(\mathbf{x}^*)$  is not along direction  
of - gradient of  $g_1(\mathbf{x}^*)$

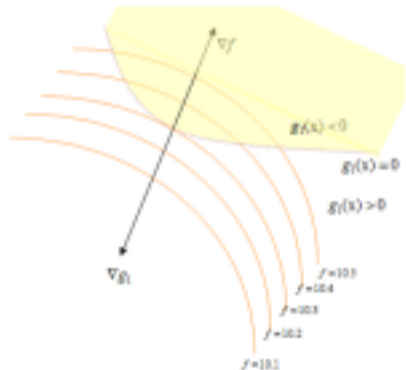
THAT IS, the two gradients MUST be in  
opposite directions

# Lagrange Multipliers with Inequality Constraints



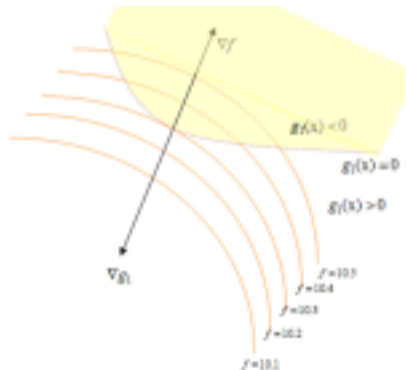
- Single equality constraint  $g_1(\mathbf{x}) = 0$ , replaced with a single inequality constraint  $g_1(\mathbf{x}) \leq 0$ . The entire region labeled  $g_1(\mathbf{x}) \leq 0$  in the Figure becomes feasible.
- At the solution  $\mathbf{x}^*$ , if  $g_1(\mathbf{x}^*) = 0$ , *i.e.*, if the constraint is active, we must have (as in the case of a single equality constraint) that  $\nabla f$  is parallel to  $\nabla g_1$ , by the same argument as before.
- Additionally, necessary for the two gradients to point in **opposite directions!**

## Lagrange Multipliers with Inequality Constraints



- Single equality constraint  $g_1(\mathbf{x}) = 0$ , replaced with a single inequality constraint  $g_1(\mathbf{x}) \leq 0$ . The entire region labeled  $g_1(\mathbf{x}) \leq 0$  in the Figure becomes feasible.
- At the solution  $\mathbf{x}^*$ , if  $g_1(\mathbf{x}^*) = 0$ , *i.e.*, if the constraint is active, we must have (as in the case of a single equality constraint) that  $\nabla f$  is parallel to  $\nabla g_1$ , by the same argument as before.
- Additionally, necessary for the two gradients to point in opposite directions; else a move away from the surface  $g_1 = 0$  and into the feasible region would further reduce  $f$ .
- With Lagrangian  $L = f + \lambda g_1$ , an additional constraint is that

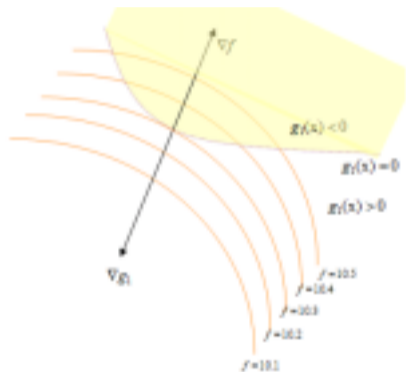
## Lagrange Multipliers with Inequality Constraints



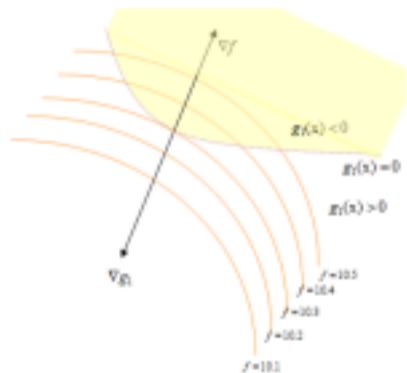
- Single equality constraint  $g_1(\mathbf{x}) = 0$ , replaced with a single inequality constraint  $g_1(\mathbf{x}) \leq 0$ . The entire region labeled  $g_1(\mathbf{x}) \leq 0$  in the Figure becomes feasible.
- At the solution  $\mathbf{x}^*$ , if  $g_1(\mathbf{x}^*) = 0$ , *i.e.*, if the constraint is active, we must have (as in the case of a single equality constraint) that  $\nabla f$  is parallel to  $\nabla g_1$ , by the same argument as before.
- Additionally, necessary for the two gradients to point in opposite directions; else a move away from the surface  $g_1 = 0$  and into the feasible region would further reduce  $f$ .
- With Lagrangian  $L = f + \lambda g_1$ , an additional constraint is that  $\lambda \geq 0$

# Lagrange Multipliers with Inequality Constraints

- If the constraint is not active at the solution  $\nabla f(\mathbf{x}^*) = 0$ , then removing  $g_1$

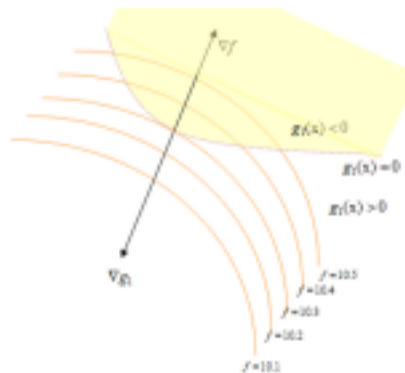


# Lagrange Multipliers with Inequality Constraints



- If the constraint is not active at the solution  $\nabla f(\mathbf{x}^*) = 0$ , then removing  $g_1$  makes no difference and we can drop it from  $L = f + \lambda g_1$ ,
- This is equivalent to setting

# Lagrange Multipliers with Inequality Constraints



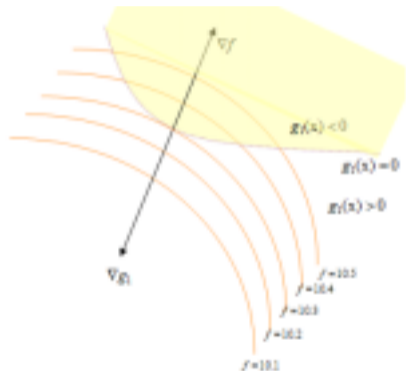
- If the constraint is not active at the solution  $\nabla f(\mathbf{x}^*) = 0$ , then removing  $g_1$  makes no difference and we can drop it from  $L = f + \lambda g_1$ ,
- This is equivalent to setting  $\lambda = 0$ .

- Thus, whether or not the constraints  $g_1 = 0$  are active, we can find the solution by requiring that
  - 1 the gradients of the Lagrangian vanish, and
  - 2  $\lambda g_1(\mathbf{x}^*) = 0$ .

This latter condition is one of the important Karush-Kuhn-Tucker conditions of convex optimization theory that can facilitate the search for the solution and will be more formally discussed subsequently.

# Lagrange Multipliers with Inequality Constraints

- Now consider the general inequality constrained minimization problem



$$\begin{aligned} \min_{\mathbf{x} \in \mathcal{D}} \quad & f(\mathbf{x}) \\ \text{subject to} \quad & g_i(\mathbf{x}) \leq 0 \quad i = 1, 2, \dots, m \end{aligned} \quad (63)$$

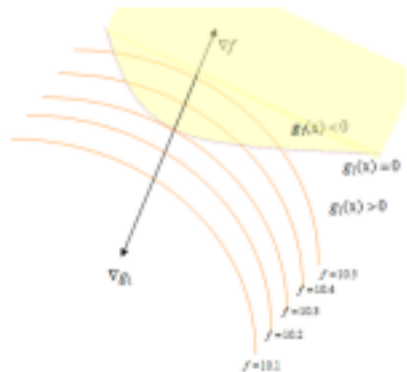
- With multiple inequality constraints, for constraints that are active, (as in the case of multiple equality constraints),

- 1  $\nabla f$  must lie in the space spanned by the  $\nabla g_i$ 's,
- 2 if the Lagrangian is  $L = f + \sum_{i=1}^m \lambda_i g_i$ , then we must also have  $\lambda_i \geq 0, \forall i$  (since otherwise  $f$  could be reduced by moving into the feasible region).



## Lagrange Multipliers with Inequality Constraints

- As for an inactive constraint  $g_j$  ( $g_j < 0$ ), removing  $g_j$  from  $L$  makes no difference and we can drop  $\nabla g_j$  from  $\nabla f = -\sum_{i=1}^m \lambda_i \nabla g_i$  or equivalently set  $\lambda_j = 0$ .
- Thus, the foregoing KKT condition generalizes to  $\lambda_i g_i(\mathbf{x}^*) = 0, \forall i$ .
- The necessary condition for optimality of (67) is summarized as:



$$\nabla L(\mathbf{x}^*, \lambda^*) = \nabla \left( f(\mathbf{x}) - \sum_{i=1}^m \lambda_i g_i(\mathbf{x}) \right) = 0$$

$$\forall i \quad \lambda_i g_i(\mathbf{x}) = 0 \quad (64)$$

# Some Algebraic Justification: Lagrange Multipliers with Inequality Constraints

# Algebraic Justification: Lagrange Multipliers with Inequality Constraints

- For the constrained optimization problem

$$\begin{array}{ll} \min_{\mathbf{x} \in \mathcal{D}} & f(\mathbf{x}) \\ \text{subject to} & \mathbf{x} \in \mathcal{C} \end{array} \quad (65)$$

$$\mathbf{x}^* = \operatorname{argmin}_{\mathbf{x} \in \mathcal{C}} f(\mathbf{x}) \iff \operatorname{argmin}_{\mathbf{x}} f(\mathbf{x}) + I_{\mathcal{C}}(\mathbf{x}), \text{ where } I_{\mathcal{C}}(\mathbf{x}) = I\{\mathbf{x} \in \mathcal{C}\} = \begin{cases} 0 & \text{if } \mathbf{x} \in \mathcal{C} \\ \infty & \text{if } \mathbf{x} \notin \mathcal{C} \end{cases}$$

$$N_{\mathcal{C}}(\mathbf{x}) = \partial I_{\mathcal{C}}(\mathbf{x}) = \left\{ \mathbf{h} \in \mathbb{R}^n \mid \mathbf{h}^T \mathbf{x} \geq \mathbf{h}^T \mathbf{z} \text{ for any } \mathbf{z} \in \mathcal{C} \right\} = \left\{ \mathbf{h} \in \mathbb{R}^n \mid \mathbf{h}^T (\mathbf{x} - \mathbf{z}) \geq 0 \text{ for an} \right.$$

- Necessary condition for optimality at  $\mathbf{x}^*$ :  $0 \in \{ \mathbf{x}^* \mid \nabla f(\mathbf{x}^*) + N_{\mathcal{C}}(\mathbf{x}^*) \}$ , that is,  $\nabla f(\mathbf{x}^*) = -N_{\mathcal{C}}(\mathbf{x}^*) = 0$  and therefore

zero belongs to the subdifferential

Negative of gradient of  $f$  at  $\mathbf{x}^*$   
must lie in normal cone

$$\underline{\nabla^T f(\mathbf{x}^*)(\mathbf{z} - \mathbf{x}^*) \geq 0 \quad \text{for any } \mathbf{z} \in \mathcal{C}} \quad (66)$$

## Algebraic Justification: Lagrange Multipliers with Inequality Constraints(contd.)

- Specifically, let  $C = \{ \mathbf{x} \in \mathbb{R}^n \mid g_i(\mathbf{x}) \leq 0 \forall i = 1, 2, \dots, m \}$

$$\begin{array}{ll} \min_{\mathbf{x} \in \mathcal{D}} & f(\mathbf{x}) \\ \text{subject to} & g_i(\mathbf{x}) \leq 0 \quad i = 1, 2, \dots, m \end{array} \quad (67)$$

Assume that each  $g_i$  is convex and is differentiable. Then, we must have, for each  $i$ ,

$$\nabla^T g_i(\mathbf{x}^*)(\mathbf{z} - \mathbf{x}^*) + g_i(\mathbf{x}^*) \leq g_i(\mathbf{z}) \quad \text{for any } \mathbf{z} \in C \quad \begin{array}{l} \text{First order condition for convexity} \\ \text{of } g_i \end{array} \quad (68)$$

- Since  $g_i(\mathbf{z}) \leq 0$  whenever  $\mathbf{z} \in C$ ,

## Algebraic Justification: Lagrange Multipliers with Inequality Constraints(contd.)

- Specifically, let  $C = \{\mathbf{x} \in \mathbb{R}^n \mid g_i(\mathbf{x}) \leq 0 \forall i = 1, 2, \dots, m\}$

$$\begin{array}{ll} \min_{\mathbf{x} \in \mathcal{D}} & f(\mathbf{x}) \\ \text{subject to} & g_i(\mathbf{x}) \leq 0 \quad i = 1, 2, \dots, m \end{array} \quad (67)$$

Assume that each  $g_i$  is convex and is differentiable. Then, we must have, for each  $i$ ,

$$\nabla^T g_i(\mathbf{x}^*)(\mathbf{z} - \mathbf{x}^*) + g_i(\mathbf{x}^*) \leq g_i(\mathbf{z}) \quad \text{for any } \mathbf{z} \in C \quad (68)$$

- Since  $g_i(\mathbf{z}) \leq 0$  whenever  $\mathbf{z} \in C$ ,

$$\Rightarrow \begin{array}{ll} \nabla^T g_i(\mathbf{x}^*)(\mathbf{z} - \mathbf{x}^*) + g_i(\mathbf{x}^*) \leq 0 & \text{for any } \mathbf{z} \in C \\ -\nabla^T g_i(\mathbf{x}^*)(\mathbf{z} - \mathbf{x}^*) - g_i(\mathbf{x}^*) \geq 0 & \text{for any } \mathbf{z} \in C \end{array} \quad (69)$$

## Algebraic Justification: Lagrange Multipliers with Inequality Constraints(contd.)

- Since any non-negative scalar (such as in (66)) is a linear combination of non-negative scalars (such as in (69)) with non-negative weights, there exists scalar (vector)  $\lambda \in \mathfrak{R}_+^m$  such that

$$\nabla^T f(\mathbf{x}^*)(\mathbf{z} - \mathbf{x}^*) = \sum_{i=1}^m -\lambda_i \nabla^T g_i(\mathbf{x}^*)(\mathbf{z} - \mathbf{x}^*) - \lambda_i g_i(\mathbf{x}^*) \quad \text{for any } \mathbf{z} \in C \quad (70)$$

- Since (70) must hold for any  $\mathbf{z} \in C$  and since  $\mathbf{x}^* \in C$ , we should have  $\lambda_i g_i(\mathbf{x}^*) = 0$ . Since the equality (70) should also continuously hold on the convex set  $C$ , we must also have

$$\nabla f(\mathbf{x}^*) = \sum_{i=1}^m -\lambda_i \nabla g_i(\mathbf{x}^*), \text{ that is } \nabla f(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i \nabla g_i(\mathbf{x}^*) = 0$$

- Since any equality constraint  $h_j(\mathbf{x}) = 0$  can be expressed as two inequality constraints:  $h_j(\mathbf{x}) \geq 0$  and  $-h_j(\mathbf{x}) \geq 0$ , the corresponding lagrange multiplier  $\mu_j$  will have no sign constraints. Additionally we require  $-h$  and  $h$  to be both convex  $\implies h$  is affine

# Duality Theory for Constrained Optimization

A tricky thing in duality theory is to decide what we call the domain or *ground set*  $\mathcal{D}$  and what we call the constraints  $g_i$ 's or  $h_j$ 's. Based on whether constraints are explicitly stated or implicitly stated in the form of the ground set, the dual problem could be very different. Thus, many duals are possible for the given primal.

For the rest of the discussion  $\mathcal{D}$  will mostly mean  $\mathbb{R}^n$

## Formally: The Dual Theory for Constrained Optimization

Consider the general constrained minimization problem

$$\begin{aligned} \min_{\mathbf{x} \in \mathcal{D}} \quad & f(\mathbf{x}) \\ \text{subject to} \quad & g_i(\mathbf{x}) \leq 0, \quad i = 1, 2, \dots, m \\ \text{subject to} \quad & h_j(\mathbf{x}) = 0, \quad j = 1, 2, \dots, n \end{aligned} \tag{71}$$

- Consider forming the Lagrange function by associating prices (called Lagrange multipliers)  $\lambda_i$  and  $\mu_j$ , with constraints involving  $g_i$  and  $h_j$  respectively.

$$L(\mathbf{x}, \lambda, \mu) = f(\mathbf{x}) + \sum_{i=1}^n \lambda_i g_i(\mathbf{x}) + \sum_{j=1}^n \mu_j h_j(\mathbf{x}) = f(\mathbf{x}) + \lambda^T \mathbf{g}(\mathbf{x}) + \mu^T \mathbf{h}(\mathbf{x})$$

- At each **feasible**  $\mathbf{x}$ , for fixed  $\lambda_i \geq 0 \forall i \in \{1..m\}$ ,

$f(\mathbf{x})$  is lower bounded by the value of the Lagrange function for all primal feasible  $\mathbf{x}$  and dual feasible  $\lambda$



## Formally: The Dual Theory for Constrained Optimization

Consider the general constrained minimization problem

$$\begin{aligned} \min_{\mathbf{x} \in \mathcal{D}} \quad & f(\mathbf{x}) \\ \text{subject to} \quad & g_i(\mathbf{x}) \leq 0, \quad i = 1, 2, \dots, m \\ \text{subject to} \quad & h_j(\mathbf{x}) = 0, \quad j = 1, 2, \dots, n \end{aligned} \tag{71}$$

- Consider forming the lagrange function by associating prices (called lagrange multipliers)  $\lambda_i$  and  $\mu_j$ , with constraints involving  $g_i$  and  $h_j$  respectively.

$$L(\mathbf{x}, \lambda, \mu) = f(\mathbf{x}) + \sum_{i=1}^m \lambda_i g_i(\mathbf{x}) + \sum_{j=1}^n \mu_j h_j(\mathbf{x}) = f(\mathbf{x}) + \lambda^T \mathbf{g}(\mathbf{x}) + \mu^T \mathbf{h}(\mathbf{x})$$

- At each **feasible**  $\mathbf{x}$ , for fixed  $\lambda_i \geq 0 \forall i \in \{1..m\}$ ,

$$f(\mathbf{x}) \geq L(\mathbf{x}, \lambda, \mu) \quad \text{if } g_i(\mathbf{x}) \leq 0 \text{ \& } h_j(\mathbf{x}) = 0 \tag{72}$$

## Formally: The Dual Theory for Constrained Optimization

- For  $\lambda_i \geq 0 \forall i \in \{1..m\}$  and  $\mu_j$ , minimizing the right hand side of (72) over all **feasible  $\mathbf{x}$**

$$f(\mathbf{x}) \geq \min_{\mathbf{x} \text{ s.t. } g_i(\mathbf{x}) \leq 0, h_j(\mathbf{x}) = 0} L(\mathbf{x}, \lambda, \mu) \triangleq L^*(\lambda, \mu) \quad (73)$$

- $L^*(\lambda, \mu)$  is a pointwise (w.r.t  $\mathbf{x} \in g_i(\mathbf{x}) \leq 0, h_j(\mathbf{x}) = 0$ ) minimum of linear functions ( $L(\mathbf{x}, \lambda, \mu)$ ) and is therefore always a **concave**

$L(\dots)$  for a fixed  $\mathbf{x}$  is affine function of lambda and mu

RECAP: Pointwise max/supremum of affine functions is always convex

## Formally: The Dual Theory for Constrained Optimization

- For  $\lambda_i \geq 0 \forall i \in \{1..m\}$  and  $\mu_j$ , minimizing the right hand side of (72) over all **feasible  $\mathbf{x}$**

$$f(\mathbf{x}) \geq \min_{\mathbf{x} \text{ s.t. } g_i(\mathbf{x}) \leq 0, h_j(\mathbf{x}) = 0} L(\mathbf{x}, \lambda, \mu) \triangleq L^*(\lambda, \mu) \quad (73)$$

- $L^*(\lambda, \mu)$  is a pointwise (w.r.t  $\mathbf{x} \in g_i(\mathbf{x}) \leq 0, h_j(\mathbf{x}) = 0$ ) minimum of linear functions ( $L(\mathbf{x}, \lambda, \mu)$ ) and is therefore always a concave function.
- Since  $f(\mathbf{x}) \geq L^*(\lambda, \mu)$  for all **primal feasible  $\mathbf{x}$**  and **dual feasible i.e.,  $\lambda_i \geq 0$  and  $\mu_j$** , we can maximize the lower bound  $L^*(\lambda, \mu)$  to give the following **Dual Problem**

$$\begin{aligned} & \max_{\lambda \in \mathbb{R}^m, \mu \in \mathbb{R}^p} L^*(\lambda, \mu) \\ & \text{subject to } \lambda \geq \mathbf{0} \end{aligned} \quad (74)$$

### Theorem

(i) The dual function  $L^*(\lambda, \mu)$  is always concave. (ii) If  $p^*$  is solution of (71) and  $d^*$  of (74) then  $p^* \geq d^*$