

# Outline

- 1 Optimization without constraints
  - Optimization scheme
  - Linear search methods
  - Gradient descent
  - Conjugate gradient
  - Newton method
  - Quasi-Newton methods
- 2 Optimization under constraints
  - Lagrange
  - Equality constraints
  - Inequality constraints
  - Dual problem - Resolution by duality
  - Numerical methods
    - Penalty functions
    - Projected gradient: equality constraints
    - Projected gradient: inequality constraints
- 3 Conclusion

# Outline

- 1 Optimization without constraints
  - Optimization scheme
  - Linear search methods
  - Gradient descent
  - Conjugate gradient
  - Newton method
  - Quasi-Newton methods
- 2 Optimization under constraints
  - Lagrange
  - Equality constraints
  - Inequality constraints
  - Dual problem - Resolution by duality
  - Numerical methods
    - Penalty functions
    - Projected gradient: equality constraints
    - Projected gradient: inequality constraints
- 3 Conclusion

# General numerical optimization scheme

$$\min_{x \in \mathbb{R}^d} f(x)$$

- ① Init: by  $x^0$ , an initial guess of a minimum  $x^*$ .
- ② Recursion: until a convergence criterion is satisfied at  $x^n$ 
  - at  $x^n$ , determine of a **search direction**  $d^n \in \mathbb{R}^d$ ,
  - **linear search**: find  $x^{n+1}$  along the semi-line  $x^n + t d^n$ ,  $t \in \mathbb{R}^+$  ; amounts to minimizing  $\phi(t) = f(x^n + t d^n)$  in  $t > 0$ .



## Remarks :

- It is important to notice that  $f$  can't be plotted ( $d$  large).  
An optimization scheme is **short-sighted** *i.e.* only has access to **local knowledge** on  $f$ .
- The information available will determine the method to use :
  - order 0 : only  $f(x^n)$  available,
  - 1st order :  $\nabla f(x^n)$  also known
  - 2nd order :  $\nabla^2 f(x^n)$  known (or estimated)
- Stop criteria are on  $\|\nabla f(x^n)\|$ , on the relative norm of the last step  $\frac{\|x^{n+1} - x^n\|}{\|x^n\|}$ , etc.
- The linear search may simply be an approximate minimization.



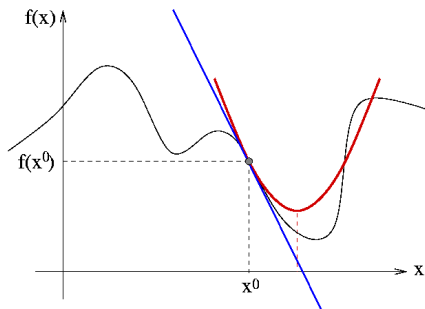
# Newton-Raphson method

Assumes  $f$  is  $C^2$ . Principle:

- approximate  $f$  by its second order expansion around  $x^n$ ,  

$$f(x) = f(x^n) + f'(x^n)(x - x^n) + \frac{1}{2}f''(x^n)(x - x^n)^2 + o(x - x^n)^2$$
- take as  $x^{n+1}$  the min of the quadratic approx. of  $f$ .

$$x^{n+1} = x^n - \frac{f'(x^n)}{f''(x^n)}$$









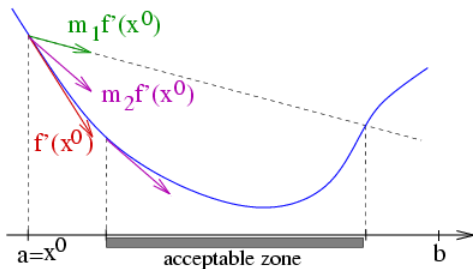


Let  $0 < m_1 < \frac{1}{2} < m_2 < 1$  be two parameters.

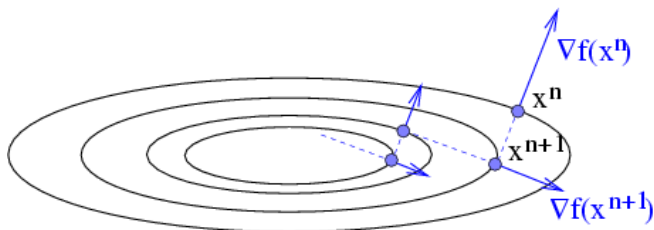
The point  $x$  is acceptable as  $x^{n+1}$  iff

$$f(x) \leq f(x^0) + m_1(x - x^0)f'(x_0)$$

$$f'(x) \geq m_2f'(x^0)$$





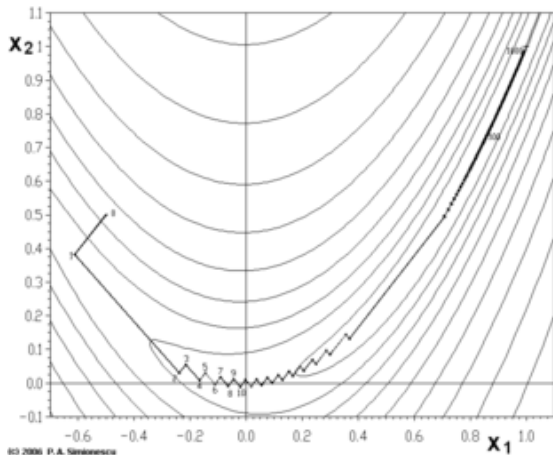


Properties :

- + Easy to implement, requires only first order information on  $f$ .
- - Slow convergence. Performs poorly even on simple functions like quadratic forms !
- In practice, a fast suboptimal descent step is preferred.

Example: on "Rosenbrock's banana"

$$f(x) = 100(x_2 - x_1^2)^2 + (1 - x_1)^2$$







We assume  $A$  is a *positive* symmetric matrix and

$$f(x) = \frac{1}{2}x^tAx + b^tx \quad \text{and} \quad \nabla f(x) = Ax + b$$

Principle :

- start at  $x^0$ ,  $d^0 = -\nabla f(x^0) \triangleq -g^0$ ,
- at  $x^n$ , instead of  $d^n = -\nabla f(x^n) \triangleq -g^n$ ,  
look for the minimum of  $f$  in the affine space

$$\mathcal{W}_{n+1} = x^0 + \text{sp}\{d^0, d^1, \dots, d^{n-1}, g^n\}$$

### Lemma

$$x^{n+1} \text{ is the min of } f \text{ in } \mathcal{W}_{n+1} \Rightarrow g^{n+1} \triangleq \nabla f(x^{n+1}) \perp \mathcal{W}_{n+1}$$











**Question :** How to find direction  $d^n$ , conjugate to all previous  $d^i$  ?

Notice  $g^{i+1} - g^i = A(x^{i+1} - x^i) \propto Ad^i$ , so

$$(d^n)^t A d^i = 0 \Rightarrow (d^n)^t g^{i+1} = (d^n)^t g^i = cst$$

Since the  $g^i$  form an orthogonal family, one has

$$d^n \propto \sum_{i=0}^n \frac{g^i}{\|g^i\|^2} \Rightarrow d^n = -g^n + c_n d^{n-1}$$

**Answer :** steepest slope, slightly corrected by previous descent direction.



Expressions of the correction coefficient  $c_n$  :

- $c_n = \frac{\|g^n\|^2}{\|g^{n-1}\|^2}$  Fletcher & Reeves (1964)
- $c_n = \frac{(g^n - g^{n-1})^t g^n}{\|g^{n-1}\|^2}$  Polak & Ribière (1971)
- $c_n = \frac{(g^n)^t A d^{n-1}}{\|d^{n-1}\|_A^2}$

Properties :

- converges in  $d$  steps for a quadratic form  $f : \mathbb{R}^d \rightarrow \mathbb{R}$
- same complexity as the gradient method !
- Works well on non quadratic forms if the Hessian doesn't change much between  $x^n$  and  $x^{n+1}$
- Caution :  $d^n$  may not be a descent direction... In this case, reset to  $-g^n$ .









## Comments :

- + faster convergence (1 step for quadratic functions!), but expensive: requires second order information on  $f$
- yields a stationary point of  $f$ : one still has to check that it is a minimum
- in practice, try  $d^n = -[\nabla^2 f(x^n)]^{-1} \nabla f(x^n)$  as descent direction, and perform a linear search
- - no guarantee that  $d^n$  is an admissible descent direction...
- - no guarantee that  $x^{n+1}$  is a better point than  $x^n$ ...
- -  $\nabla^2 f(x^n)$  may be singular, or badly conditioned...
- the Levenberg-Marquardt **regularization** suggests to solve

$$[\nabla^2 f(x^n) + \mu \mathbf{1}] d^n = -\nabla f(x^n)$$



# Quasi-Newton methods

## Principle :

- Take advantage of the efficiency of the Newton method...
- ... when the Hessian  $\nabla^2 f(x)$  is unavailable !
- **Idea :** *approximate*  $[\nabla^2 f(x^n)]^{-1}$  by matrix  $K_n$  in

$$x^{n+1} = x^n - [\nabla^2 f(x^n)]^{-1} \nabla f(x^n)$$

- More precisely, explore direction  $d^n = -K_n \nabla f(x^n)$  from  $x^n$ .

# Quasi-Newton equation

- Consider the second order Taylor expansion of  $f$  at  $x^n$

$$f(x) = f(x^n) + \nabla f(x^n)^t (x - x^n) + \frac{1}{2}(x - x^n)^t \nabla^2 f(x^n) (x - x^n) + o(\|x - x^n\|^2)$$

$$\nabla f(x) = \nabla f(x^n) + \nabla^2 f(x^n) (x - x^n) + o(\|x - x^n\|)$$

- The estimate  $K_n$  of the inverse Hessian must satisfy the **quasi-Newton equation (QNE)**

$$x^{n+1} - x^n = K_{n+1} [\nabla f(x^{n+1}) - \nabla f(x^n)]$$

- Notice that this should be  $K_n$ ... but  $K_n$  is used to find  $x^{n+1}$ , so we impose the relation be satisfied at the next step.



# Quasi-Newton equation

- Consider the second order Taylor expansion of  $f$  at  $x^n$

$$f(x) = f(x^n) + \nabla f(x^n)^t (x - x^n) + \frac{1}{2}(x - x^n)^t \nabla^2 f(x^n) (x - x^n) + o(\|x - x^n\|^2)$$

$$\nabla f(x) = \nabla f(x^n) + \nabla^2 f(x^n) (x - x^n) + o(\|x - x^n\|)$$

- The estimate  $K_n$  of the inverse Hessian must satisfy the **quasi-Newton equation (QNE)**

$$x^{n+1} - x^n = K_{n+1} [\nabla f(x^{n+1}) - \nabla f(x^n)]$$

- Notice that this should be  $K_n$ ... but  $K_n$  is used to find  $x^{n+1}$ , so we impose the relation be satisfied at the next step.

All quasi-Newton Methods recursively build the  $K_n$  by

$$K_{n+1} = K_n + C_n$$

where the correction  $C_n$  is adjusted to satisfy the QNE.

Notations:

$$u^n = x^n - x^{n-1}$$

$$v^n = g^n - g^{n-1}$$

$$\text{QNE : } u^{n+1} = K_{n+1} v^{n+1}$$

- Correction  $C_n$  of rank 1

$$K_{n+1} = K_n + \frac{w^n (w^n)^t}{(w^n)^t v^{n+1}} \quad \text{where} \quad w^n = u^{n+1} - K_n v^{n+1}$$

- If initialized with  $K_0 = \mathbf{I}$ ,  $K_n$  converges in  $d$  steps to the true  $A^{-1}$  for a quadratic form.
- **DFP** (Davidon, Fletcher, Powell) correction of rank 2

$$K_{n+1} = K_n + \frac{u^{n+1} (u^{n+1})^t}{(u^{n+1})^t v^{n+1}} - \frac{K_n v^{n+1} (v^{n+1})^t K_n}{(v^{n+1})^t K_n v^{n+1}}$$

- Converges in  $d$  steps to the true  $A^{-1}$  for a quadratic form.
- Descent directions are **conjugate** w.r.t.  $A$ .
- Coincides with the conjugate gradient method.



- **BFGS** (Broyden, Fletcher, Goldfarb, Shanno, 1970),  
correction of rank 3

$$K_{n+1} = K_n - \frac{u^{n+1}(v^{n+1})^t K_n + K_n v^{n+1}(u^{n+1})^t}{(u^{n+1})^t v^{n+1}} \\ + \left(1 + \frac{(v^{n+1})^t K_n v^{n+1}}{(u^{n+1})^t v^{n+1}}\right) \frac{u^{n+1}(u^{n+1})^t}{(u^{n+1})^t v^{n+1}}$$

Considered as the best Quasi-Newton method.

- In practice, one should check that  $-K_n g^n$  is a descent direction, i.e.  $-(g^n)^t K_n g^n < 0$ , otherwise reinitialize by  $K_n = \mathbf{1}$ .







# Outline

- 1 Optimization without constraints
  - Optimization scheme
  - Linear search methods
  - Gradient descent
  - Conjugate gradient
  - Newton method
  - Quasi-Newton methods
- 2 Optimization under constraints
  - Lagrange
  - Equality constraints
  - Inequality constraints
  - Dual problem - Resolution by duality
  - Numerical methods
    - Penalty functions
    - Projected gradient: equality constraints
    - Projected gradient: inequality constraints
- 3 Conclusion

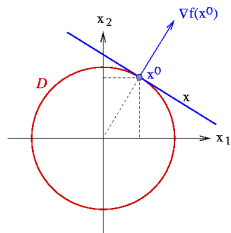




# Equality constraints

$$\min_x f(x) \quad \text{s.t.} \quad \theta_j(x) = 0, \quad 1 \leq j \leq m$$

- $\mathcal{D} : \theta(x) = [\theta_1(x), \dots, \theta_m(x)]^t = 0$   
defines a **manifold** of dimension  $d - m$  in  $\mathbb{R}^d$
- $\nabla \theta_j(x^0)^t (x - x^0) = 0$  : **tangent hyperplane** to  $\theta_j(x) = 0$  at point  $x^0$



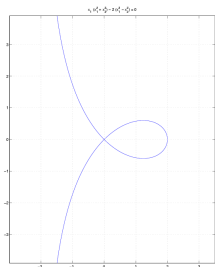
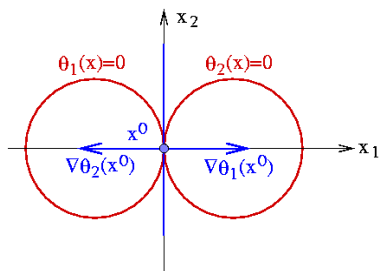
- $\nabla \theta(x^0)^t (x - x^0) = 0$  : **tangent space** to  $\mathcal{D}$  at  $x^0$

## Definition

In domain  $\mathcal{D} = \{x \in \mathbb{R}^d : \theta(x) = 0\}$ , the point  $x^0$  is **regular** iff the gradients  $\nabla\theta_j(x^0)$  of the  $m$  constraints are linearly independent.

## Lemma

If  $x^0$  is regular, every (unit) direction  $d$  in the tangent space is **admissible**, i.e. can be obtained as the limit of  $\frac{x^n - x^0}{\|x^n - x^0\|}$ , with  $\lim_n x^n = x^0$  and  $x^n \in \mathcal{D}$ .





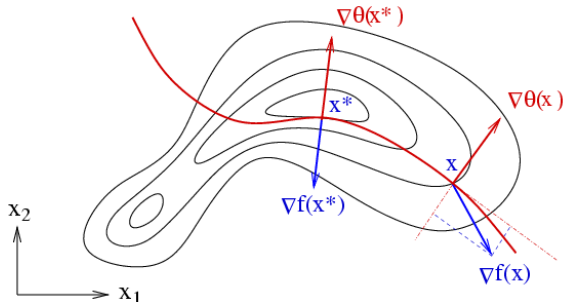


Proof:

- Project  $\nabla f(x^*)$  on  $sp\{\nabla\theta_1(x^*), \dots, \nabla\theta_m(x^*)\}$

$$\nabla f(x^*) = \sum_{j=1}^m -\lambda_j^* \nabla\theta_j(x^*) + u$$

- $u$  belongs to the tangent space to  $\mathcal{D}$  at  $x^*$
- progressing along  $-u$  decreases  $f$  and doesn't change  $\theta$



# Application

To solve  $\min_x f(x)$  s.t.  $\theta(x) = 0$ ,

- 1 build the **Lagrangian**

$$L(x, \lambda) = f(x) + \sum_j \lambda_j \theta_j(x)$$

- 2 find a stationary point  $(x^*, \lambda^*)$  of the Lagrangian, *i.e.* a zero of  $\nabla L(x, \lambda)$

$$\nabla_x L(x, \lambda) = \nabla f(x) + \sum_j \lambda_j \nabla \theta_j(x)$$

$$\nabla_\lambda L(x, \lambda) = \theta(x)$$

*i.e.*  $d + m$  (non-linear) equations, with  $d + m$  unknowns.





Solution : Lagrangian  $L(x, \lambda) = f(x) + \lambda\theta(x)$

$$\frac{\partial L(x, \lambda)}{\partial x_1} = 2\pi x_1 + 2\pi x_2 + \lambda 2\pi x_1 x_2 = 0$$

$$\frac{\partial L(x, \lambda)}{\partial x_2} = 2\pi x_1 + \lambda \pi x_1^2 = 0$$

We obtain  $x_1^* = x_2^* = -\frac{2}{\lambda}$ .

Finally,  $\theta(x^*) = 0$  gives the value of  $\lambda$  to plug :

$\lambda^* = -\frac{\pi^{1/3}}{2}$ , so  $x_1^* = x_2^* = \pi^{-1/3}$ .

## Another interpretation

- Consider the unconstrained problem, where  $\lambda$  is fixed

$$\min_x L(x, \lambda) = f(x) + \lambda\theta(x)$$

- $f$  and  $L$  have the same local minima in  $\mathcal{D} = \{x : \theta(x) = 0\}$ .
- Let  $x^*(\lambda)$  be a local minimum of  $L(x, \lambda)$  in  $\mathbb{R}^d$ .
- If  $x^*(\lambda) \in \mathcal{D}$ , then it is also a local min of  $f$ .
- So one just has to adjust  $\lambda$  to get this property.

## Second order conditions

### Theorem

Let  $(x^*, \lambda^*)$  be a stationary point of  $L(x, \lambda)$ , and consider the Hessian of the Lagrangian

$$\nabla_x^2 L(x^*, \lambda^*) = \nabla^2 f(x^*) + \sum_{j=1}^m \lambda_j^* \nabla^2 \theta_j(x^*)$$

- NC:  $x^*$  is a local min of  $f$  on  $\mathcal{D} \Rightarrow \nabla_x^2 L(x^*, \lambda^*)$  is a positive quadratic form on the tangent space at  $x^*$ , *i.e.* the kernel of matrix  $\nabla \theta(x^*)^t$ .
- SC:  $\nabla_x^2 L(x^*, \lambda^*)$  is strictly positive on the tangent space  $\Rightarrow x^*$  is a local min of  $f$  on  $\mathcal{D}$





# Inequality constraints

$$\min_x f(x) \quad \text{s.t.} \quad \theta_j(x) \leq 0, \quad 1 \leq j \leq m$$

- $\mathcal{D} : \theta(x) = [\theta_1(x), \dots, \theta_m(x)]^t \leq 0$  defines a **volume** in  $\mathbb{R}^d$  limited by  $m$  **manifolds** of dimension  $d - 1$
- At point  $x$ , constraint  $\theta_j$  is **active** iff  $\theta_j(x) = 0$ .  
 $\mathcal{A}(x) = \{j : \theta_j(x) = 0\} =$  active set at  $x$ .
- One could have simultaneously equality and inequality constraints (not done here for a of matter clarity).  
 Equality constraints are always active.
- $\bigcap_{j \in \mathcal{A}(x^0)} \{x : \nabla \theta_j(x^0)^t (x - x^0) = 0\}$   
 defines the **tangent space** to  $\mathcal{D}$  at  $x^0$

# Admissible directions

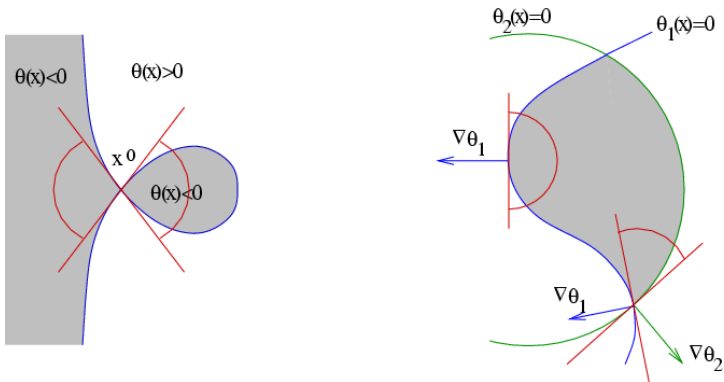
Let  $x^0 \in \mathcal{D}$ , we look for directions  $d \in \mathbb{R}^d$  that keep us inside domain  $\mathcal{D}$ :  $x^0 + \epsilon \cdot d \in \mathcal{D}$ .

## Definition

Direction  $d$  is **admissible** from  $x^0$  iff  $\exists (x^n)_{n>0}$  in  $\mathcal{D}$  such that

$$\lim_n x^n = x^0 \quad \text{and} \quad \lim_n \frac{x^n - x^0}{\|x^n - x^0\|} = \frac{d}{\|d\|}$$

- Admissible directions at  $x^0$  form a **cone**  $\mathcal{C}(x^0)$ .
- This cone is not necessarily convex...



- $\mathcal{C}(x^0)$  can be determined from the  $\nabla\theta_j(x^0)$  of the active constraints.

## Theorem

If  $x^0$  is a **regular** point, *i.e.* the gradients of the active constraints at  $x^0$  are linearly independent, then  $\mathcal{C}(x^0)$  is the *convex* cone given by

$$\mathcal{C}(x^0) = \{u \in \mathbb{R}^d : \nabla \theta_j(x^0)^t u \leq 0, j \in \mathcal{A}(x^0)\}$$

Interpretation: an admissible displacement must not increase the value of  $\theta_j(x^0)$  for an already active constraint, it can only decrease it or leave it unchanged.

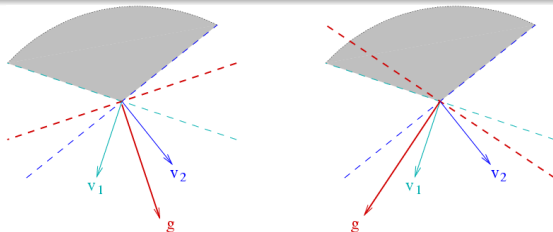
# Dual cone and Farkas lemma

For  $v_1, \dots, v_J \in \mathbb{R}^d$ , consider cone  $\mathcal{C} = \{u : u^t v_1 \leq 0, \dots, u^t v_j \leq 0\}$ .

## Farkas-Minkowski lemma

Let  $g \in \mathbb{R}^d$ , one has the equivalence

- $\forall u \in \mathcal{C}, g^t u \leq 0$ ,
- $\mathcal{C}$  is included in the half-space  $\{u : g^t u \leq 0\}$ ,
- $g$  belongs to the **dual cone**  $\mathcal{C}' = \{w : \forall u \in \mathcal{C}, w^t u \leq 0\}$
- $g = \sum_{j=1}^J \alpha_j v_j$  where  $\alpha_j \geq 0$  for all  $j$



# 1st order optimality conditions

## Theorem (Karush-Kuhn-Tucker conditions)

Let  $x^*$  be a regular point of domain  $\mathcal{D}$ . If  $x^*$  is a local minimum of  $f$  in  $\mathcal{D}$ , there exists a unique set of **generalized Lagrange multipliers**  $\lambda_j^*$  for  $j \in \mathcal{A}(x^*)$  such that

$$\nabla f(x^*) + \sum_{j \in \mathcal{A}(x^*)} \lambda_j^* \nabla \theta_j(x^*) = 0 \quad \text{and} \quad \lambda_j^* \geq 0, \quad j \in \mathcal{A}(x^*)$$

Remarks :

- Similar to the case of equality constraints: here only *active* constraints are considered.
- The positivity condition is new: translates the fact that one side of the manifold is permitted.

# 1st order optimality conditions

## Theorem (Karush-Kuhn-Tucker conditions)

Let  $x^*$  be a regular point of domain  $\mathcal{D}$ . If  $x^*$  is a local minimum of  $f$  in  $\mathcal{D}$ , there exists a unique set of **generalized Lagrange multipliers**  $\lambda_j^*$  for  $j \in \mathcal{A}(x^*)$  such that

$$\nabla f(x^*) + \sum_{j \in \mathcal{A}(x^*)} \lambda_j^* \nabla \theta_j(x^*) = 0 \quad \text{and} \quad \lambda_j^* \geq 0, \quad j \in \mathcal{A}(x^*)$$

Remarks :

- Similar to the case of equality constraints: here only *active* constraints are considered.
- The positivity condition is new: translates the fact that one side of the manifold is permitted.

Proof:

- take any admissible direction :  
 $d \in \mathcal{C}(x^*) = \{u : u^t \nabla \theta_j(x^*) \leq 0, j \in \mathcal{A}(x^*)\}$
- progressing along  $d$  doesn't decrease  $f$  :  $[-\nabla f(x^*)]^t d \leq 0$
- this means that  $g = -\nabla f(x^*)$  belongs to the **dual cone**  $\mathcal{C}(x^*)'$ , so by Farkas lemma

$$-\nabla f(x^*) = \sum_{j \in \mathcal{A}(x^*)} \lambda_j^* \nabla \theta_j(x^*) \quad \text{and} \quad \lambda_j^* \geq 0$$



## Corollary

The Karush-Kuhn-Tucker conditions are equivalent to

$$\nabla f(x^*) + \sum_{j=1}^m \lambda_j^* \nabla \theta_j(x^*) = 0 \quad \text{and} \quad \lambda_j^* \geq 0, \quad 1 \leq j \leq m$$

with the extra **complementarity condition**

$$\sum_{j=1}^m \lambda_j^* \theta_j(x^*) = 0$$

- This entails  $\lambda_j^* = 0$  for an inactive constraint  $\theta_j$  at  $x^*$ .
- To be usable, requires to know/guess the set of active constraints at the optimum.
- $\mathcal{A}(x^*)$  known, leaves a set of non-linear equations + positivity constraints.



Objective : cancel the gradient of the Lagrangian

$$\frac{\partial L(x, \lambda)}{\partial x_1} = x_1 - 1 + \lambda_1 + \lambda_2 - \lambda_3 = 0$$

$$\frac{\partial L(x, \lambda)}{\partial x_2} = x_2 - 2 - \lambda_1 + \lambda_2 - \lambda_4 = 0$$

equality  $\frac{\partial L(x, \lambda)}{\partial \lambda_1} = x_1 - x_2 - 1 = 0$

inequalities  $\lambda_2(x_1 + x_2 - 2) = 0, \quad \lambda_2 \geq 0$

$$-\lambda_3 x_1 = 0, \quad \lambda_3 \geq 0$$

$$-\lambda_4 x_2 = 0, \quad \lambda_4 \geq 0$$

**1st guess:**  $\mathcal{A}(x^*) = \{1\}$ , i.e. only  $\theta_1$  active at the optimum.

Complementarity  $\Rightarrow \lambda_2^* = \lambda_3^* = \lambda_4^* = 0$ .

This yields  $x^* = (2, 1)$  which violates  $\theta_2(x) \leq 0$ .

Objective : cancel the gradient of the Lagrangian

$$\frac{\partial L(x, \lambda)}{\partial x_1} = x_1 - 1 + \lambda_1 + \lambda_2 - \lambda_3 = 0$$

$$\frac{\partial L(x, \lambda)}{\partial x_2} = x_2 - 2 - \lambda_1 + \lambda_2 - \lambda_4 = 0$$

equality  $\frac{\partial L(x, \lambda)}{\partial \lambda_1} = x_1 - x_2 - 1 = 0$

inequalities  $\lambda_2(x_1 + x_2 - 2) = 0, \quad \lambda_2 \geq 0$

$$-\lambda_3 x_1 = 0, \quad \lambda_3 \geq 0$$

$$-\lambda_4 x_2 = 0, \quad \lambda_4 \geq 0$$

**2nd guess :**  $\mathcal{A}(x^*) = \{1, 2\}$ , i.e.  $\theta_2$  is added to the active set.

Complementarity  $\Rightarrow \lambda_3^* = \lambda_4^* = 0$ .

This yields  $x^* = (\frac{3}{2}, \frac{1}{2})$  which belongs to  $\mathcal{D}$ .

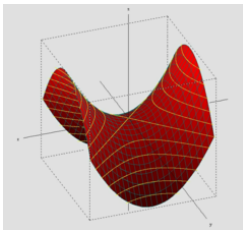
# Dual problem - Resolution by duality

[ For simplicity we consider the case of inequality constraints. ]

**Idea:** under some conditions, a stationary point  $(x^*, \lambda^*)$  of the Lagrangian, i.e.  $\nabla L(x^*, \lambda^*) = \nabla f(x^*) + \sum_i \lambda_i^* \nabla \theta(x^*) = 0$  corresponds to a **saddle point** of the Lagrangian, i.e.

$$\inf_x L(x, \lambda^*) = L(x^*, \lambda^*) = \sup_{\lambda} L(x^*, \lambda)$$

So the resolution amounts to finding such saddle points, and then extract  $x^*$ .

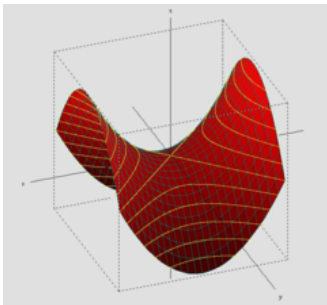


# Saddle points

## Definition

$(x^*, \lambda^*)$  is a **saddle point** of  $L$  in  $\mathcal{D}_x \times \mathcal{D}_\lambda$  iff

$$\sup_{\lambda \in \mathcal{D}_\lambda} L(x^*, \lambda) = L(x^*, \lambda^*) = \inf_{x \in \mathcal{D}_x} L(x, \lambda^*)$$



## Lemma

If  $(x^*, \lambda^*)$  is a saddle point of  $L$  in  $\mathcal{D}_x \times \mathcal{D}_\lambda$ , then

$$\sup_{\lambda \in \mathcal{D}_\lambda} \inf_{x \in \mathcal{D}_x} L(x, \lambda) = L(x^*, \lambda^*) = \inf_{x \in \mathcal{D}_x} \sup_{\lambda \in \mathcal{D}_\lambda} L(x, \lambda)$$

## Proof

- one always has  $\sup_{\lambda} \inf_x L(x, \lambda) \leq \inf_x \sup_{\lambda} L(x, \lambda)$   
the difference is called the *duality gap*, generally  $> 0$
- from the def. of a saddle point, one has

$$\sup_{\lambda} L(x^*, \lambda) = L(x^*, \lambda^*) = \inf_x L(x, \lambda^*)$$

then

$$\begin{aligned} \inf_x [\sup_{\lambda} L(x, \lambda)] &\leq \sup_{\lambda} L(x^*, \lambda) \\ \inf_x L(x, \lambda^*) &\leq \sup_{\lambda} [\inf_x L(x, \lambda)] \end{aligned}$$

- Morality: one can look for  $\lambda^*$  first, and then for  $x^*$ ...

## Lemma

If  $(x^*, \lambda^*)$  is a saddle point of  $L$  in  $\mathcal{D}_x \times \mathcal{D}_\lambda$ , then

$$\sup_{\lambda \in \mathcal{D}_\lambda} \inf_{x \in \mathcal{D}_x} L(x, \lambda) = L(x^*, \lambda^*) = \inf_{x \in \mathcal{D}_x} \sup_{\lambda \in \mathcal{D}_\lambda} L(x, \lambda)$$

## Proof

- one always has  $\sup_{\lambda} \inf_x L(x, \lambda) \leq \inf_x \sup_{\lambda} L(x, \lambda)$   
the difference is called the *duality gap*, generally  $> 0$
- from the def. of a saddle point, one has

$$\sup_{\lambda} L(x^*, \lambda) = L(x^*, \lambda^*) = \inf_x L(x, \lambda^*)$$

then

$$\begin{aligned} \inf_x [\sup_{\lambda} L(x, \lambda)] &\leq \sup_{\lambda} L(x^*, \lambda) \\ \inf_x L(x, \lambda^*) &\leq \sup_{\lambda} [\inf_x L(x, \lambda)] \end{aligned}$$

- Morality: one can look for  $\lambda^*$  first, and then for  $x^*$ ...



## Lemma

If  $(x^*, \lambda^*)$  is a saddle point of  $L$  in  $\mathcal{D}_x \times \mathcal{D}_\lambda$ , then

$$\sup_{\lambda \in \mathcal{D}_\lambda} \inf_{x \in \mathcal{D}_x} L(x, \lambda) = L(x^*, \lambda^*) = \inf_{x \in \mathcal{D}_x} \sup_{\lambda \in \mathcal{D}_\lambda} L(x, \lambda)$$

## Proof

- one always has  $\sup_{\lambda} \inf_x L(x, \lambda) \leq \inf_x \sup_{\lambda} L(x, \lambda)$   
the difference is called the *duality gap*, generally  $> 0$
- from the def. of a saddle point, one has

$$\sup_{\lambda} L(x^*, \lambda) = L(x^*, \lambda^*) = \inf_x L(x, \lambda^*)$$

then

$$\begin{aligned} \inf_x [\sup_{\lambda} L(x, \lambda)] &\leq \sup_{\lambda} L(x^*, \lambda) \\ \inf_x L(x, \lambda^*) &\leq \sup_{\lambda} [\inf_x L(x, \lambda)] \end{aligned}$$

- Morality: one can look for  $\lambda^*$  first, and then for  $x^*$ ...

## Lemma

If  $(x^*, \lambda^*)$  is a saddle point of  $L$  in  $\mathcal{D}_x \times \mathcal{D}_\lambda$ , then

$$\sup_{\lambda \in \mathcal{D}_\lambda} \inf_{x \in \mathcal{D}_x} L(x, \lambda) = L(x^*, \lambda^*) = \inf_{x \in \mathcal{D}_x} \sup_{\lambda \in \mathcal{D}_\lambda} L(x, \lambda)$$

## Proof

- one always has  $\sup_{\lambda} \inf_x L(x, \lambda) \leq \inf_x \sup_{\lambda} L(x, \lambda)$   
the difference is called the *duality gap*, generally  $> 0$
- from the def. of a saddle point, one has

$$\sup_{\lambda} L(x^*, \lambda) = L(x^*, \lambda^*) = \inf_x L(x, \lambda^*)$$

then

$$\begin{aligned} \inf_x [\sup_{\lambda} L(x, \lambda)] &\leq \sup_{\lambda} L(x^*, \lambda) \\ \inf_x L(x, \lambda^*) &\leq \sup_{\lambda} [\inf_x L(x, \lambda)] \end{aligned}$$

- Morality: one can look for  $\lambda^*$  first, and then for  $x^*$ ...

# Saddle points of the Lagrangian

## Theorem

If  $(x^*, \lambda^*)$  is a saddle point of the Lagrangian  $L$  in  $\mathbb{R}^d \times \mathbb{R}_+^m$ , then  $x^*$  is a solution of the primal problem (P)

$$(P) \quad \min_x f(x) \quad \text{s.t.} \quad \theta_i(x) \leq 0, \quad 1 \leq i \leq m$$

## Proof

- From  $L(x^*, \lambda) \leq L(x^*, \lambda^*)$ ,  $\forall \lambda \in \mathcal{D}_\lambda = \mathbb{R}_+^m$

$$f(x^*) + \sum_i \lambda_i \theta_i(x^*) \leq f(x^*) + \sum_i \lambda_i^* \theta_i(x^*)$$

$$\sum_i (\lambda_i - \lambda_i^*) \theta_i(x^*) \leq 0$$

whence  $\theta_i(x^*) \leq 0$  by  $\lambda_i \rightarrow +\infty$ :  $x^*$  satisfies constraints

# Saddle points of the Lagrangian

## Theorem

If  $(x^*, \lambda^*)$  is a saddle point of the Lagrangian  $L$  in  $\mathbb{R}^d \times \mathbb{R}_+^m$ , then  $x^*$  is a solution of the primal problem (P)

$$(P) \quad \min_x f(x) \quad \text{s.t.} \quad \theta_i(x) \leq 0, \quad 1 \leq i \leq m$$

## Proof

- From  $L(x^*, \lambda) \leq L(x^*, \lambda^*)$ ,  $\forall \lambda \in \mathcal{D}_\lambda = \mathbb{R}_+^m$

$$f(x^*) + \sum_i \lambda_i \theta_i(x^*) \leq f(x^*) + \sum_i \lambda_i^* \theta_i(x^*)$$

$$\sum_i (\lambda_i - \lambda_i^*) \theta_i(x^*) \leq 0$$

whence  $\theta_i(x^*) \leq 0$  by  $\lambda_i \rightarrow +\infty$ :  $x^*$  satisfies constraints

- Moreover,  $\sum_i -\lambda_i^* \theta_i(x^*) \leq 0$ , by  $\lambda_i = 0$ ,  
and so  $\sum_i \lambda_i^* \theta_i(x^*) = 0$  (complementarity condition)
- From  $L(x^*, \lambda^*) \leq L(x, \lambda^*)$ ,  $\forall x \in \mathbb{R}^d$

$$f(x^*) + \sum_i \lambda_i^* \theta_i(x^*) \leq f(x) + \sum_i \lambda_i^* \theta_i(x)$$

so for all admissible  $x$ , i.e. such that  $\theta_i(x) \leq 0$ ,  $1 \leq i \leq m$

$$f(x^*) \leq f(x)$$

### Summary :

saddle points of the Lagrangian, when they exist, give solutions to the optimization problem.

*But they don't always exist...*

- Moreover,  $\sum_i -\lambda_i^* \theta_i(x^*) \leq 0$ , by  $\lambda_i = 0$ ,  
and so  $\sum_i \lambda_i^* \theta_i(x^*) = 0$  (complementarity condition)
- From  $L(x^*, \lambda^*) \leq L(x, \lambda^*)$ ,  $\forall x \in \mathbb{R}^d$

$$f(x^*) + \sum_i \lambda_i^* \theta_i(x^*) \leq f(x) + \sum_i \lambda_i^* \theta_i(x)$$

so for all admissible  $x$ , i.e. such that  $\theta_i(x) \leq 0$ ,  $1 \leq i \leq m$

$$f(x^*) \leq f(x)$$

### Summary :

saddle points of the Lagrangian, when they exist, give solutions to the optimization problem.

*But they don't always exist...*

- Moreover,  $\sum_i -\lambda_i^* \theta_i(x^*) \leq 0$ , by  $\lambda_i = 0$ ,  
and so  $\sum_i \lambda_i^* \theta_i(x^*) = 0$  (complementarity condition)
- From  $L(x^*, \lambda^*) \leq L(x, \lambda^*)$ ,  $\forall x \in \mathbb{R}^d$

$$f(x^*) + \sum_i \lambda_i^* \theta_i(x^*) \leq f(x) + \sum_i \lambda_i^* \theta_i(x)$$

so for all admissible  $x$ , i.e. such that  $\theta_i(x) \leq 0$ ,  $1 \leq i \leq m$

$$f(x^*) \leq f(x)$$

### Summary :

saddle points of the Lagrangian, when they exist, give solutions to the optimization problem.

*But they don't always exist...*

# Existence of saddle points

## Theorem

If  $f$  and the constraints  $\theta_i$  are **convex** functions of  $x$  in  $\mathbb{R}^d$ , and if  $x^* = \arg \min_x f(x)$  in  $\{x : \theta_i(x) \leq 0, 1 \leq i \leq m\}$  is regular then  $x^*$  corresponds to a saddle point  $(x^*, \lambda^*)$  of the Lagrangian

**Proof:** from Kuhn-Tucker, derive the saddle point property

- $$L(x^*, \lambda) = f(x^*) + \sum_i \lambda_i \theta_i(x^*)$$

$$\leq f(x^*) = f(x^*) + \sum_i \lambda_i^* \theta_i(x^*) = L(x^*, \lambda^*)$$
 using admissibility of  $x^*$ , positivity of  $\lambda_i$  and complementarity
- $$L(x, \lambda^*) = f(x) + \sum_i \lambda_i^* \theta_i(x)$$
 is a convex function of  $x$   
 From the stationarity of  $L$ , one has
 
$$\nabla_x L(x^*, \lambda^*) = \nabla f(x^*) + \sum_i \lambda_i^* \nabla \theta_i(x^*) = 0$$
 sufficient to show that  $x^*$  is a minimum of the convex function  $L(x, \lambda^*)$



# Existence of saddle points

## Theorem

If  $f$  and the constraints  $\theta_i$  are **convex** functions of  $x$  in  $\mathbb{R}^d$ , and if  $x^* = \arg \min_x f(x)$  in  $\{x : \theta_i(x) \leq 0, 1 \leq i \leq m\}$  is regular then  $x^*$  corresponds to a saddle point  $(x^*, \lambda^*)$  of the Lagrangian

**Proof:** from Kuhn-Tucker, derive the saddle point property

- $$L(x^*, \lambda) = f(x^*) + \sum_i \lambda_i \theta_i(x^*)$$

$$\leq f(x^*) = f(x^*) + \sum_i \lambda_i \theta_i(x^*) = L(x^*, \lambda^*)$$
 using admissibility of  $x^*$ , positivity of  $\lambda_i$  and complementarity
- $$L(x, \lambda^*) = f(x) + \sum_i \lambda_i^* \theta_i(x)$$
 is a convex function of  $x$   
 From the stationarity of  $L$ , one has
 
$$\nabla_x L(x^*, \lambda^*) = \nabla f(x^*) + \sum_i \lambda_i^* \nabla \theta_i(x^*) = 0$$
 sufficient to show that  $x^*$  is a minimum of the convex function  $L(x, \lambda^*)$

# Existence of saddle points

## Theorem

If  $f$  and the constraints  $\theta_i$  are **convex** functions of  $x$  in  $\mathbb{R}^d$ , and if  $x^* = \arg \min_x f(x)$  in  $\{x : \theta_i(x) \leq 0, 1 \leq i \leq m\}$  is regular then  $x^*$  corresponds to a saddle point  $(x^*, \lambda^*)$  of the Lagrangian

**Proof:** from Kuhn-Tucker, derive the saddle point property

- $$L(x^*, \lambda) = f(x^*) + \sum_i \lambda_i \theta_i(x^*)$$

$$\leq f(x^*) = f(x^*) + \sum_i \lambda_i \theta_i(x^*) = L(x^*, \lambda^*)$$
 using admissibility of  $x^*$ , positivity of  $\lambda_i$  and complementarity
- $$L(x, \lambda^*) = f(x) + \sum_i \lambda_i^* \theta_i(x)$$
 is a convex function of  $x$   
 From the stationarity of  $L$ , one has
 
$$\nabla_x L(x^*, \lambda^*) = \nabla f(x^*) + \sum_i \lambda_i^* \nabla \theta_i(x^*) = 0$$
 sufficient to show that  $x^*$  is a minimum of the convex function  $L(x, \lambda^*)$

# Dual problem

**Summary:** Provided the Lagrangian has saddle points

- Solutions to  $(P) \min_x f(x) \text{ s.t. } \theta_i(x) \leq 0, 1 \leq i \leq m$  are the 1st argument of a saddle point  $(x^*, \lambda^*)$  of the Lagrangian  $L(x, \lambda)$
- If  $\lambda^*$  were known, amounts to solving an **unconstrained** problem

$$x^* = \arg \min_x L(x, \lambda^*)$$

- How to find such a  $\lambda^*$  ?

One has  $L(x^*, \lambda^*) = \max_{\lambda \in \mathbb{R}_+^m} \min_x L(x, \lambda)$ ,

so  $\lambda^*$  should be a solution of the **dual problem**

$$(D) \max_{\lambda} g(\lambda), \text{ s.t. } \lambda \in \mathbb{R}_+^m, \text{ where } g(\lambda) = \min_x L(x, \lambda)$$

# Dual problem

**Summary:** Provided the Lagrangian has saddle points

- Solutions to  $(P) \min_x f(x)$  s.t.  $\theta_i(x) \leq 0, 1 \leq i \leq m$  are the 1st argument of a saddle point  $(x^*, \lambda^*)$  of the Lagrangian  $L(x, \lambda)$
- If  $\lambda^*$  were known, amounts to solving an **unconstrained** problem

$$x^* = \arg \min_x L(x, \lambda^*)$$

- How to find such a  $\lambda^*$  ?

One has  $L(x^*, \lambda^*) = \max_{\lambda \in \mathbb{R}_+^m} \min_x L(x, \lambda)$ ,

so  $\lambda^*$  should be a solution of the **dual problem**

$$(D) \max_{\lambda} g(\lambda), \text{ s.t. } \lambda \in \mathbb{R}_+^m, \text{ where } g(\lambda) = \min_x L(x, \lambda)$$

# Dual problem

**Summary:** Provided the Lagrangian has saddle points

- Solutions to  $(P) \min_x f(x)$  s.t.  $\theta_i(x) \leq 0, 1 \leq i \leq m$  are the 1st argument of a saddle point  $(x^*, \lambda^*)$  of the Lagrangian  $L(x, \lambda)$
- If  $\lambda^*$  were known, amounts to solving an **unconstrained** problem

$$x^* = \arg \min_x L(x, \lambda^*)$$

- How to find such a  $\lambda^*$  ?

One has  $L(x^*, \lambda^*) = \max_{\lambda \in \mathbb{R}_+^m} \min_x L(x, \lambda)$ ,

so  $\lambda^*$  should be a solution of the **dual problem**

$$(D) \max_{\lambda} g(\lambda), \text{ s.t. } \lambda \in \mathbb{R}_+^m, \text{ where } g(\lambda) = \min_x L(x, \lambda)$$



Under some conditions, it is equivalent to solve the (P) or (D) :

### Theorem

- If the  $\theta_i$  are continuous over  $\mathbb{R}^d$ , and  $\forall \lambda \in \mathbb{R}_+^m$ ,  $x^*(\lambda) = \arg \min_x L(x, \lambda)$  is unique, and  $x^*(\lambda)$  is a continuous function of  $\lambda$  then  $\lambda^*$  solves (D)  $\Rightarrow x^*(\lambda^*)$  solves (P)
- If (P) has at least one solution  $x^*$ ,  $f$  and the  $\theta_i$  are convex and  $x^*$  is regular, then (D) has at least a solution  $\lambda^*$ .

### Remark

(D) is still an optimization problem under constraints...

... but constraints  $\lambda \in \mathbb{R}_+^m$  are much simpler to handle !





































# Outline

- 1 Optimization without constraints
  - Optimization scheme
  - Linear search methods
  - Gradient descent
  - Conjugate gradient
  - Newton method
  - Quasi-Newton methods
- 2 Optimization under constraints
  - Lagrange
  - Equality constraints
  - Inequality constraints
  - Dual problem - Resolution by duality
  - Numerical methods
    - Penalty functions
    - Projected gradient: equality constraints
    - Projected gradient: inequality constraints
- 3 Conclusion



# Conclusion

When facing a constrained optimization problem...

...one reflex : **BUILD THE LAGRANGIAN !**

$$L(x, \lambda) = f(x) + \sum_j \lambda_j \theta_j(x)$$

- solve  $\nabla L(x, \lambda) = 0$  to find a candidate optimum  $(x^*, \lambda^*)$
- check the positivity of its Hessian  $\nabla_x^2 L(x^*, \lambda^*)$  to check if  $x^*$  is a min, a max or a saddle point of  $f$ .

# Conclusion

When facing a constrained optimization problem...

...one reflex : **BUILD THE LAGRANGIAN !**

$$L(x, \lambda) = f(x) + \sum_j \lambda_j \theta_j(x)$$

- solve  $\nabla L(x, \lambda) = 0$  to find a candidate optimum  $(x^*, \lambda^*)$
- check the positivity of its Hessian  $\nabla_x^2 L(x^*, \lambda^*)$  to check if  $x^*$  is a min, a max or a saddle point of  $f$ .

# Conclusion

When facing a constrained optimization problem...

...one reflex : **BUILD THE LAGRANGIAN !**

$$L(x, \lambda) = f(x) + \sum_j \lambda_j \theta_j(x)$$

- solve  $\nabla L(x, \lambda) = 0$  to find a candidate optimum  $(x^*, \lambda^*)$
- check the positivity of its Hessian  $\nabla_x^2 L(x^*, \lambda^*)$  to check if  $x^*$  is a min, a max or a saddle point of  $f$ .

# Conclusion

When facing a constrained optimization problem...

...one reflex : **BUILD THE LAGRANGIAN !**

$$L(x, \lambda) = f(x) + \sum_j \lambda_j \theta_j(x)$$

- solve  $\nabla L(x, \lambda) = 0$  to find a candidate optimum  $(x^*, \lambda^*)$
- check the positivity of its Hessian  $\nabla_x^2 L(x^*, \lambda^*)$  to check if  $x^*$  is a min, a max or a saddle point of  $f$ .

# Conclusion

When facing a constrained optimization problem...

...one reflex : **BUILD THE LAGRANGIAN !**

$$L(x, \lambda) = f(x) + \sum_j \lambda_j \theta_j(x)$$

- solve  $\nabla L(x, \lambda) = 0$  to find a candidate optimum  $(x^*, \lambda^*)$
- check the positivity of its Hessian  $\nabla_x^2 L(x^*, \lambda^*)$  to check if  $x^*$  is a min, a max or a saddle point of  $f$ .