

# Chapter 4: Unconstrained Nonlinear Optimization

Edoardo Amaldi

DEIB – Politecnico di Milano  
edoardo.amaldi@polimi.it



Academic year 2023-24

## 4.1 Examples

### 1) Statistical estimation

Random variable  $X$  with density  $f(x, \underline{\theta})$ , where  $\underline{\theta} \in \mathbb{R}^m$  is parameter vector, and independent observations  $x_1, \dots, x_n$ .

Maximum likelihood: Estimates  $\hat{\underline{\theta}}$  of  $\underline{\theta}$  are derived by maximizing

$$L(\underline{\theta}) = f(x_1, \underline{\theta}) f(x_2, \underline{\theta}) \dots f(x_n, \underline{\theta})$$

*product of non linear functions  
(well, non linear actually)  
and unconstrained*

Assumption:  $\exists \underline{\theta}$  for which all factors are positive.

Since  $\ln(\cdot)$  is monotonically increasing,  $\hat{\underline{\theta}}$  also maximizes

$$\ln(L(\underline{\theta})) = \sum_{j=1}^n \ln(f(x_j, \underline{\theta}))$$

If  $f$  is differentiable w.r.t.  $\underline{\theta}$  at  $\hat{\underline{\theta}}$ , necessary optimality conditions:

$$\sum_{j=1}^n \frac{\nabla_{\underline{\theta}} f(x_j, \hat{\underline{\theta}})}{f(x_j, \hat{\underline{\theta}})} = \underline{0}$$

For Gaussian density

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp -\frac{(x - \mu)^2}{2\sigma^2}$$

and  $\underline{\theta} = (\mu, \sigma)$ , we obtain

$$\ln(L(\underline{\theta})) = \ln\left(\frac{1}{\sigma\sqrt{2\pi}}\right)^n \prod_{j=1}^n \exp -\frac{(x_j - \mu)^2}{2\sigma^2} = -\frac{n}{2} \ln(2\pi) - n \ln(\sigma) +$$
$$-\frac{1}{2\sigma^2} \sum_{j=1}^n (x_j - \mu)^2$$

Minimum is achieved in a stationary point:

and

$$\left. \begin{aligned} \frac{\partial[\ln(L(\underline{\theta}))]}{\partial \mu} &= \frac{1}{\sigma^2} \sum_{j=1}^n (x_j - \mu) = 0 \\ \frac{\partial[\ln(L(\underline{\theta}))]}{\partial \sigma} &= -\frac{n}{\sigma} + \frac{1}{\sigma^3} \sum_{j=1}^n (x_j - \mu)^2 = 0 \end{aligned} \right\}$$

Thus

$$\Rightarrow \begin{aligned} \hat{\mu} &= \frac{1}{n} \sum_{j=1}^n x_j \\ \hat{\sigma} &= \sqrt{\frac{1}{n} \sum_{j=1}^n (x_j - \hat{\mu})^2} \end{aligned}$$

## 2) Training multilayer neural networks

Supervised learning:

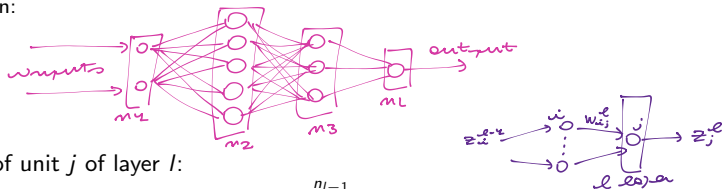
Given a training set  $T = \{(\underline{x}^1, \underline{y}^1), \dots, (\underline{x}^p, \underline{y}^p)\}$  where  $\underline{y}^k \in [0, 1]^{n_{out}}$  desired output for  $\underline{x}^k \in \mathbb{R}^{n_{in}}$ , construct a model that maps  $\underline{x}^k$ 's into  $\underline{y}^k$ 's as well as possible.

Multilayer networks:

$L$  layers with  $n_l$  units in layer  $l$ ,  $n_1 = n_{in}$  and  $n_L = n_{out}$ .

First layer of inputs  $x_1, \dots, x_{n_1}$ , other layers with activation units.

Illustration:



Output of unit  $j$  of layer  $l$ :

$$z_j^l = \phi\left(\sum_{i=1}^{n_{l-1}} w_{ij}^l z_i^{l-1} - w_{0j}^l\right)$$

where weights  $w_{ij}$  to be determined and  $\phi: \mathbb{R} \rightarrow \mathbb{R}$  is sigmoid  $\phi(t) = \frac{1}{1+e^{-t}}$ .

A multilayer network defines a mapping  $h(\underline{w}, \cdot)$  from  $\mathbb{R}^{n_1}$  to  $\mathbb{R}^{n_L}$  parametrized by  $\underline{w} = \{w_{ij}^l : l = 1, \dots, L; i = 1, \dots, n_{l-1}; j = 1, \dots, n_l\}$ .

Training problem: Given  $T = \{(\underline{x}^1, \underline{y}^1), \dots, (\underline{x}^p, \underline{y}^p)\}$ , determine values of  $\underline{w}$  which approximate as well as possible the mapping underlying  $T$ .

In general one minimizes

$$\frac{1}{2} \sum_{k=1}^p (\|\underline{y}^k - h(\underline{w}, \underline{x}^k)\|^2)$$

*well, non linear  
and non convex  
minimization*

challenging (non convex)



Example 1.5.3 of D. Bertsekas, Nonlinear Programming, Athena Scientific 1999.

## 4.2 Optimality conditions

Generic problem:

$$\min_{\underline{x} \in S} f(\underline{x})$$

where  $S \subseteq \mathbb{R}^n$ ,  $f : S \rightarrow \mathbb{R}$  and  $f \in \mathcal{C}^1$  or  $\mathcal{C}^2$ .

*simple constraints, i.e. the non-restrictive*

Unconstrained case:  $S = \mathbb{R}^n$

$\Rightarrow$  we don't care about " $\underline{d}$ ", we consider  $\underline{d}$  just as the direction he provides

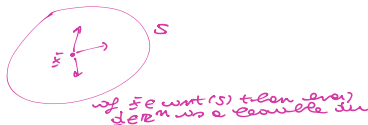
**Definition:**  $\underline{d} \in \mathbb{R}^n$  is a feasible direction at  $\underline{\bar{x}}$  if

*smallest step - another we could take*

$$\exists \bar{\alpha} > 0 \text{ such that } \underline{\bar{x}} + \alpha \underline{d} \in S \quad \forall \alpha \in [0, \bar{\alpha}] \quad (1)$$

Illustrations:

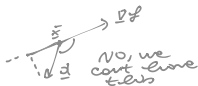
*we wif we can make strictly positive steps in the direction even less  $\underline{d}$  and still remaining in  $S$*



## First order necessary local optimality conditions:

If  $f \in C^1$  on  $S$  and  $\bar{x}$  is a local minimum of  $f$  over  $S$ , then for any feasible direction  $d \in \mathbb{R}^n$  at  $\bar{x}$

$$\nabla^t f(\bar{x})d \geq 0,$$



namely all feasible directions are ascent directions.

*we worsening direction*

Proof:

Consider  $\varphi: [0, \bar{\alpha}] \rightarrow \mathbb{R}$  st  $\varphi(\alpha) = f(\bar{x} + \alpha d)$ .

Since  $\bar{x}$  is a local min (l.e) on  $S$  of  $f$  on  $S$ , then  $\alpha = 0$  is a local min of  $\varphi$  on  $[0, \bar{\alpha}]$  (as that we remain in  $S$ ).

Being  $f \in C^1$ , also  $\varphi \in C^1$ , and the Taylor series of  $\varphi$  at  $\alpha = 0$  is

$$\varphi(\alpha) = \varphi(0) + \alpha \varphi'(0) + o(\alpha) \quad \sim \text{note: } o(\alpha) = o(\alpha) \text{ wif } \alpha \rightarrow 0 \text{ faster than } \alpha, \text{ when } \alpha \rightarrow 0 \text{ w itself}$$

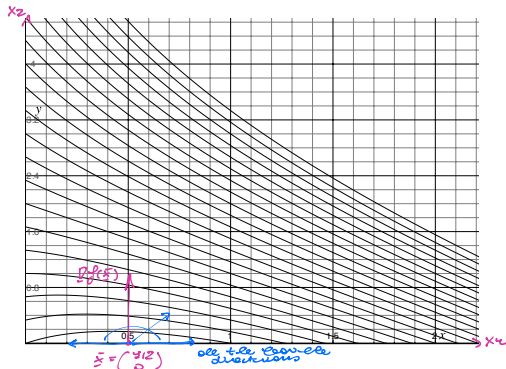
Suppose that  $\varphi'(0) > 0$ . Then wif  $\alpha \rightarrow 0^+$  we can neglect the remainder  $o(\alpha)$  and we have

$$\varphi(\alpha) - \varphi(0) = \alpha \varphi'(0) > 0 \Rightarrow 0 \text{ would not be a local min of } \varphi, \text{ absurd}$$

So we must have  $\varphi'(0) \geq 0$ , and since  $\varphi'(\alpha) = \nabla f(\bar{x} + \alpha d) \cdot d$  we have the condition of the theorem.

Example:

$$\min_{x_1, x_2 \geq 0} f(x_1, x_2) = x_1^2 - x_1 + x_2 + x_1 x_2$$



$\Rightarrow \bar{x}$  is a (global) min because  
 $Df(\bar{x}) \cdot d \geq 0$  for all possible directions  
 $d$  at  $\bar{x}$ , even if  $Df(\bar{x}) \neq 0$

— the first order "classical" conditions (of derivative = 0) are not enough here, we have to introduce possible directions



## Second order/necessary/local optimality conditions:

If  $f \in \mathcal{C}^2$  on  $S$  and  $\bar{x}$  is a local minimum of  $f$  over  $S$  then

i)  $\nabla^t f(\bar{x})\underline{d} \geq 0$

$\forall \underline{d} \in \mathbb{R}^n$  feasible direction at  $\bar{x}$ ,

*the 1st order condition we just saw*

ii) if  $\nabla^t f(\bar{x})\underline{d} = 0$

*the vectors are orthogonal*

then  $\underline{d}^t \nabla^2 f(\bar{x})\underline{d} \geq 0$ .  $\Leftrightarrow \phi''(0) \geq 0$

*d.H.d. is  $\geq 0$ , so that  $\underline{d}$  at which we have  $\perp$ .*

Proof:

Similarly for (i). Suppose  $\nabla^t f(\bar{x})\underline{d} = 0$ , then *we can expand the Taylor series and we get*

$$\phi(\alpha) = \phi(0) + \underbrace{\alpha \phi'(0)}_0 + \frac{1}{2} \alpha^2 \phi''(0) + o(\alpha^2).$$

*if  $\phi''(0) = \underline{d}^t \nabla^2 f(\bar{x}) \underline{d} < 0$ , we can reason as before and we would get*

$$\phi(\alpha) - \phi(0) = \frac{1}{2} \alpha^2 \phi''(0) < 0 \Rightarrow$$

*for sufficiently small  $\alpha$  we would get a better loc min of  $\phi$  than  $0$ , absurd*

*So we must have  $\phi''(0) \geq 0$  which is the thesis*

**Corollary:** (Unconstrained case)

If  $f \in \mathcal{C}^2$  on  $S$  and  $\bar{x} \in \text{int}(S)$  is a local minimum of  $f$  over  $S$ , then

- 1  $\nabla f(\bar{x}) = 0$  (stationarity condition)
- 2  $\nabla^2 f(\bar{x})$  is positive semidefinite.

Proof: Since  $\bar{x} \in \text{int}(S)$ , all the vectors  $d \in \mathbb{R}^n$  are feasible directions at  $\bar{x}$ .

So from the previous condition (i) we have  $\nabla f(\bar{x}) \cdot d \geq 0 \forall d$  and  $-d$ , so we get (4) here.

We use (2) as a consequence of condition (ii) since  $d \cdot (\nabla^2 f(\bar{x})) d \geq 0 \forall d \in \mathbb{R}^n$ , the Hessian matrix  $\nabla^2 f(\bar{x})$  is positive semidefinite.

*we check not only one of the previous conditions*

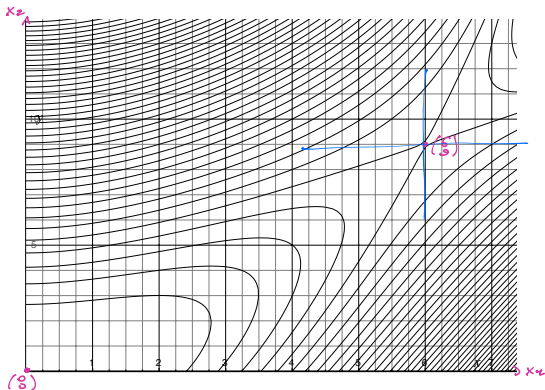
Types of candidate points: local minima, local maxima and saddle points.

Above optimality conditions are not sufficient

*eg just take  $f(x) = x^3$ , for which  $x=0$  is not a local minimum, even if it satisfies the conditions*

Example:

$$\min_{x_1, x_2 \geq 0} f(x_1, x_2) = x_1^3 - x_1^2 x_2 + 2x_2^2$$



Candidate points are  $(0,0)$  and  $(6,9)$ .  
-  $(6,9)$  is not a local min, even though for  $x_1=6$  and  $x_2=9$  is a local min w.r.t  $x_2$  and for  $x_2=9$  and  $x_1=6$  is a local min w.r.t  $x_1$  ( $\Rightarrow$  is a saddle)

## Sufficient | local optimality conditions:

extension (on a general set  $S$ ) of the classical opt conditions

If  $f \in C^2$  on  $S$  and  $\bar{x} \in \text{int}(S)$  such that  $\nabla f(\bar{x}) = \underline{0}$  and  $\nabla^2 f(\bar{x})$  is positive definite, then  $\bar{x}$  is a strict local minimum of  $f$  over  $S$ , namely

$$f(x) > f(\bar{x}) \quad \forall x \in \mathcal{N}_\epsilon(\bar{x}) \cap S.$$

Proof:

Let  $\underline{d} \in \mathcal{B}_\epsilon(0)$  be any feasible direction such that  $\bar{x} + \underline{d} \in S \cap \mathcal{B}_\epsilon(\bar{x})$ .



Then

$$f(\bar{x} + \underline{d}) = f(\bar{x}) + \underbrace{\nabla^t f(\bar{x})}_{0} \underline{d} + \frac{1}{2} \underline{d}^t \nabla^2 f(\bar{x}) \underline{d} + o(\|\underline{d}\|^2)$$

Since  $\nabla^2 f(\bar{x})$  is pos def, then  $\exists \alpha > 0$  st  $\underline{d} \cdot (\nabla^2 f(\bar{x}) \underline{d}) \geq \alpha \cdot \|\underline{d}\|^2$  (where  $\alpha$  is related to the smallest evalue of the Hessian).

Thus for  $\|\underline{d}\|$  sufficiently small, we have that

$$f(\bar{x} + \underline{d}) - f(\bar{x}) \geq \frac{\alpha}{2} \|\underline{d}\|^2 > 0$$

$\Rightarrow f(\bar{x}) < f(\bar{x} + \underline{d}) \Rightarrow \bar{x}$  is a strict loc min along the direction  $\underline{d}$

Since this holds  $\forall \underline{d} \in \mathbb{R}^n$  such that  $\bar{x} + \underline{d} \in S \cap \mathcal{B}_\epsilon(\bar{x})$ ,  $f$  is locally strictly convex.

# Convex problems

$$\min_{x \in C \subseteq \mathbb{R}^n} f(x) \quad \text{where } C \subseteq \mathbb{R}^n \text{ convex and } f: C \rightarrow \mathbb{R} \text{ convex}$$

Every local minimum is a global minimum.

*a stronger result for convex problems*

## Necessary and sufficient (NS) conditions:

Let  $f$  be convex and  $C^1$  on  $C \subseteq \mathbb{R}^n$  convex.  $x^*$  is a global minimum of  $f$  on  $C$  if and only if

$$\nabla^t f(x^*)(y - x^*) \geq 0 \quad \forall y \in C.$$

Proof:

nec. cond: w/  $y \in C$  and  $x^*$  is a lsc min (or global min) then we have that (all feasible dirs are ascending, i.e.)

$$\nabla f(x^*) \cdot d \geq 0 \quad \forall d \text{ feasible dir at } x^* \\ \text{we } \forall d = y - x^* \text{ with } y \in C$$

Suf. cond: here we use the characterization of convexity:

$$\begin{aligned} f \text{ is convex} & \Leftrightarrow f(y) \geq f(x^*) + \underbrace{\nabla f(x^*) \cdot (y - x^*)}_{\substack{\text{ZO from} \\ \text{nec cond}}} \quad \forall y \in C \\ & \Rightarrow f(y) \geq f(x^*) \quad \forall y \in C \\ & \Rightarrow x^* \text{ is a global opt} \end{aligned}$$

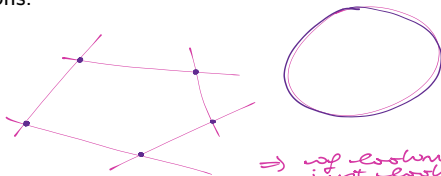
Recall: Given any  $C \subseteq \mathbb{R}^n$  convex,  $\underline{x} \in C$  is an *extreme point* of  $C$  if it cannot be expressed as a convex combination of two different points of  $C$ .

**Property:** (maximization of convex functions)

*eg we may want to maximize a linear function (which may be convex & concave)*

Let  $f$  be convex defined on  $C$  convex, bounded and closed. If  $f$  has a (finite) maximum over  $C$ , then  $\exists$  an optimal extreme point of  $C$ .

Illustrations:



$\Rightarrow$  if looking for max points just look at extreme points  
It is more or less the main result of linear programming (LP)

Special case: Linear programming

## 4.3 Iterative methods and convergence

relevant since all methods for constrained or not solvers are iterative in nature

Generic Nonlinear Optimization (N.O.) problem:

$$\begin{aligned} \min \quad & f(\underline{x}) \\ \text{s.t.} \quad & g_i(\underline{x}) \leq 0 \quad 1 \leq i \leq m \\ & \underline{x} \in S \subseteq \mathbb{R}^n \end{aligned}$$

new notation (w.l.o.g.) for the constraints

If  $X = \{\underline{x} \in S : g_i(\underline{x}) \leq 0, 1 \leq i \leq m\} \subset \mathbb{R}^n$  then constrained problem.

Difficulty depends on  $f$  and  $X$ . Usually  $f$  and  $g_i$  are at least continuously differentiable.

In some cases (e.g., LP and combinatorial optimization) an optimal solution can be found in a finite number of elementary operations.

Efficiency depends on how this number grows with the instance size (polynomial vs exponential).

(we # variables, # constraints, ecc)

## Most N.O. methods are iterative

- start from  $\underline{x}_0 \in X$
- generate a sequence  $\{\underline{x}_k\}_{k \geq 0}$  that “converges” to a point of  $\Omega = \{\text{“desired solutions”}\}$ .

## Different meanings of “converge” and “desired solutions”:

- $\{\underline{x}_k\}_{k \geq 0}$  converges to a point of  $\Omega$   
or  $\exists$  a limit point of  $\{\underline{x}_k\}_{k \geq 0}$  which belongs to  $\Omega$   
*a good estimate after  
sufficient number of  
iterations*
  - $\Omega =$  set of global optima  
or  $\Omega =$  set of candidate points satisfying 1st/2nd order necessary optimality conditions  
*where have we really  
converged in the  
classical sense*
- eg. wif  $X = \mathbb{R}^m$  we can have  
 $\Omega = \{x \in \mathbb{R}^m : \nabla f(x) = 0\}$*

Often but not always descent methods:  $f(\underline{x}_{k+1}) < f(\underline{x}_k)$  for each  $k$



Interested in robust and efficient methods.

1) **Robustness** associated to global convergence

*ie morally, it does not matter the so close from where we start*

**Definition:** An algorithm is **globally (locally) convergent** if  $\{\underline{x}_k\}_{k \geq 0}$  satisfies one of previous properties for any  $\underline{x}_0 \in X$  (only for  $\underline{x}_0$  in a neighborhood of an  $\underline{x}^* \in \Omega$ ).

*a candidate point*

2) **Efficiency** characterized by convergence speed

Assume that  $\lim_{k \rightarrow \infty} \underline{x}_k = \underline{x}^*$  where  $\underline{x}^* \in \Omega$

**Definitions:**  $\{\underline{x}_k\}_{k \geq 0}$  converges to  $\underline{x}^*$  with order  $p \geq 1$  if  $\exists r > 0$  and  $k_0 \in \mathbb{N}$  such that

$$\|\underline{x}_{k+1} - \underline{x}^*\| \leq r \|\underline{x}_k - \underline{x}^*\|^p \quad \forall k \geq k_0.$$

Largest  $p$  is the **order of convergence** and smallest  $r > 0$  is the **rate**.

If  $p = 1$  and  $r < 1$  **linear** convergence, if  $p = 1$  and  $r \geq 1$  **sublinear** convergence.

*this one is slower than*

N.B.: If  $p = 1$  the distance w.r.t.  $\underline{x}^*$  decreases at each iteration by a factor  $r$ .

Example: (E1) Consider  $1 + \frac{c}{k} \xrightarrow{k \rightarrow \infty} 1$

- verify that this is a linear convergence  
 - and that the ratio is  $\frac{c}{k}$ .

moreover, we need to study the ratio of LHS/RHS:

$$\frac{\|c_n(k+1)\|}{\|c_n(k)\|} = \frac{|ck+4-c|}{|ck-c|} = \frac{\frac{c}{k+4}}{\left(\frac{c}{k}\right)^p} = \frac{k^p}{k+4} \xrightarrow{k \rightarrow \infty} k^{p-1}$$

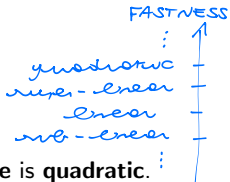
$\Rightarrow$  one can test  $p$   
 we can test  $p > 1$

(E2) Study similarly  $ck = c + \frac{c}{2k}$  (let:  $r = \frac{1}{2}$ )

**Definition:** The convergence is **superlinear** if there exists  $\{r_k\}_{k \geq 0}$  with  $\lim_{k \rightarrow \infty} r_k = 0$  such that

$$\|x_{k+1} - x^*\| \leq r_k \|x_k - x^*\| \quad \forall k \geq k_0.$$

Example:  $1 + \frac{1}{k^k}$



**Definition:** If  $p = 2$  (and  $r$  not necessarily  $< 1$ ), the convergence is **quadratic**.

Example:  $1 + \frac{1}{2^{2^k}}$

## 4.4 Line search methods

Unconstrained optimization problem:

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

with  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  of class  $\mathcal{C}^1$  or  $\mathcal{C}^2$  and bounded below.

*otherwise the problem is unbounded (and is not interesting)*

Iterative methods: start from  $\underline{x}_0 \in \mathbb{R}^n$  and generate  $\{\underline{x}_k\}_{k \geq 0}$  “converging” to an  $\bar{\underline{x}} \in \Omega$ .

See Chap. 3 of J. Nocedal, S. Wright, Numerical Optimization, Springer 1999.

## 1) General scheme

Select  $\underline{x}_0$  and  $\varepsilon > 0$ , set  $k := 0$

**Repeat**

Choose search direction  $\underline{d}_k \in \mathbb{R}^n$

Determine step length  $\alpha_k > 0$  along  $\underline{d}_k$  s.t.

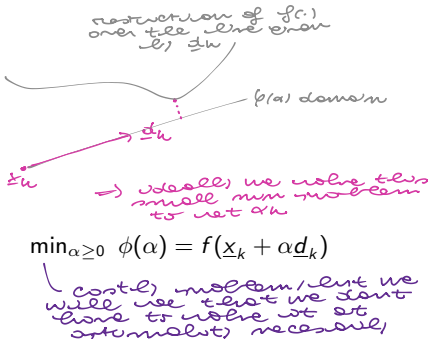
Set  $\underline{x}_{k+1} := \underline{x}_k + \alpha_k \underline{d}_k$  and  $k := k + 1$

**Until** termination criterion is satisfied

Termination criterion:  $\|\nabla f(\underline{x}_k)\| < \varepsilon$  or  $|f(\underline{x}_k) - f(\underline{x}_{k+1})| < \varepsilon$  or  $\|\underline{x}_{k+1} - \underline{x}_k\| < \varepsilon$

Often approximate  $\alpha_k$  (also  $f(\underline{x}_{k+1}) < f(\underline{x}_k) \forall k \geq 0$ ).

Flexibility in choice of  $\underline{d}_k$  and  $\alpha_k$ , efficiency depends on both.



## 2) Search directions

In many *line search methods*, i.e., iterative methods based on search directions,

$$\underline{d}_k = -D_k \nabla f(\underline{x}_k)$$

with positive definite  $n \times n$  matrix  $D_k$ .

$\underline{d}_k$  is a descent direction because of the matrix  $D_k$  being pos def

Now also then we can compute

$$\begin{aligned} \nabla f(\underline{x}_k) \cdot \underline{d}_k &= \nabla f(\underline{x}_k) \cdot (-D_k \nabla f(\underline{x}_k)) = \\ &= -(\underbrace{\underline{z} \cdot A \underline{z}}_{>0}) < 0 \end{aligned}$$

## Example 1: Gradient method

Given  $f \in \mathcal{C}^1$ , consider linear approximation of  $f(\underline{x}_k + \underline{d})$  at  $\underline{x}_k$

$$l_k(\underline{d}) := f(\underline{x}_k) + \nabla^t f(\underline{x}_k) \underline{d}$$

and choose  $\underline{d}_k \in \mathbb{R}^n$  minimizing  $l_k(\underline{d})$  over sphere of radius  $\|\nabla f(\underline{x}_k)\|$ :

*then about wh  
we are free*

$$\begin{aligned} \min \quad & \nabla^t f(\underline{x}_k) \underline{d} \\ \text{s.t.} \quad & \|\underline{d}\| = \|\nabla f(\underline{x}_k)\|. \end{aligned} \tag{1}$$

*region where we  
look for  $\underline{d}$*

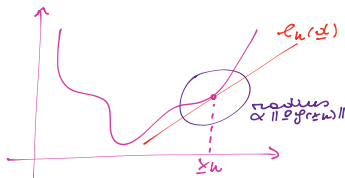
*since  $\nabla f(\underline{x}_k) \cdot \underline{d} = \|\nabla f(\underline{x}_k)\| \cdot \|\underline{d}\| \cdot \cos \theta$   
we minimize (4) when  $\cos(\theta) = -1$   
that is  $\theta = \pi$*

**Steepest descent direction:**

$$\underline{d}_k = -\nabla f(\underline{x}_k)$$

where  $\underline{D}_k = I_n$ .

Clearly  $\underline{d}_k$  is a descent direction if  $\nabla f(\underline{x}_k) \neq \underline{0}$ .



## Example 2: Newton method

Given  $f \in \mathcal{C}^2$  and  $H(\underline{x}_k) = \nabla^2 f(\underline{x}_k)$ .

Consider quadratic approximation of  $f(\underline{x}_k + \underline{d})$  at  $\underline{x}_k$

$$q_k(\underline{d}) := f(\underline{x}_k) + \nabla^t f(\underline{x}_k) \underline{d} + \frac{1}{2} \underline{d}^t H(\underline{x}_k) \underline{d}$$

and choose  $\underline{d}_k \in \mathbb{R}^n$  and  $\alpha_k$  leading to a stationary point of  $q_k(\underline{d})$ .

Since  $\nabla_{\underline{d}} q_k(\underline{d}) = \underline{0}$  implies  $\nabla^t f(\underline{x}_k) + \underline{d}^t H(\underline{x}_k) = \underline{0}$ , if  $H^{-1}(\underline{x}_k)$  exists then

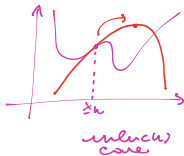
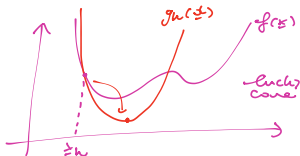
**Newton direction:**

$$\underline{d}_k = -H^{-1}(\underline{x}_k) \nabla f(\underline{x}_k),$$

where  $\underline{D}_k = H^{-1}(\underline{x}_k)$ .

If  $H(\underline{x}_k)$  is p.d. and  $\nabla f(\underline{x}_k) \neq \underline{0}$ ,  $\underline{d}_k$  is a descent direction.

If  $H(\underline{x}_k)$  is not p.d.,  $\underline{d}_k$  may not be defined ( $\nexists H^{-1}(\underline{x}_k)$ ) or may be an ascent direction.



$$\begin{aligned} \nabla f(\underline{x}_k) \cdot \underline{d}_k &= \nabla f \cdot (-H^{-1} \nabla f) = \\ &= -(\nabla \cdot H^{-1} \nabla) < 0 \\ &\quad \text{no def} \end{aligned}$$

### 3) Step length

To guarantee global convergence, an approximate solution  $\alpha_k$  of line search:

$$\min_{\alpha \geq 0} \phi(\alpha) = f(\underline{x}_k + \alpha \underline{d}_k).$$

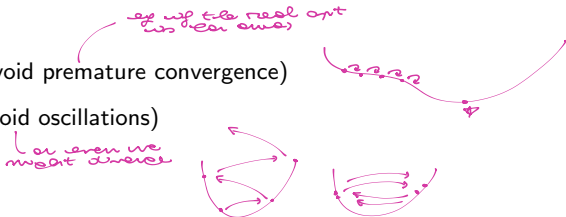
is sufficient.

Different methods to generate  $\alpha_k$  and stop when appropriate conditions are satisfied (simple, after a few iterations).

$f(\underline{x}_k + \alpha_k \underline{d}_k) < f(\underline{x}_k)$  does not suffice.

Basic principles:

- $\alpha$  must not be too small (to avoid premature convergence)
- $\alpha$  must not be too large (to avoid oscillations)





## Wolfe conditions:

Sufficient reduction <sup>(decrease)</sup>

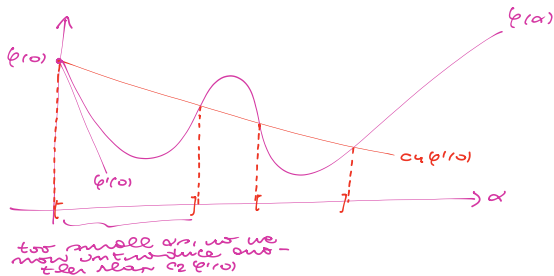
$$\phi(\alpha) \leq \phi(0) + c_1 \alpha \phi'(0) \quad \text{con } c_1 \in [0, 1]$$

which is equivalent to

$$f(\underline{x}_k + \alpha \underline{d}_k) \leq \underbrace{f(\underline{x}_k)}_{\phi(0)} + c_1 \alpha \underbrace{\nabla^t f(\underline{x}_k) \underline{d}_k}_{\phi'(0)} \quad (\text{Armijo criterion})$$

$\phi'(0) < 0$  since  $\underline{d}_k$  is a descent direction,  $c_1 \leq 1/2$  so that it is satisfied by the minimum of a quadratic convex  $\phi(\alpha)$  (exercise set n.6).

Illustration:



To avoid too small steps also condition:

$$\Rightarrow c_4 \leq c_2 \leq c_4$$

$$\phi'(\alpha) \geq c_2 \phi'(0) \quad \text{con } c_2 \in (c_1, 1)$$

which is equivalent to

$$\nabla^t f(\underline{x}_k + \alpha \underline{d}_k) \underline{d}_k \geq c_2 \nabla^t f(\underline{x}_k) \underline{d}_k.$$

In general  $c_2 = 0.9$  for (quasi)-Newton and  $c_2 = 0.1$  for non-linear conjugate gradient.

Weak Wolfe conditions:

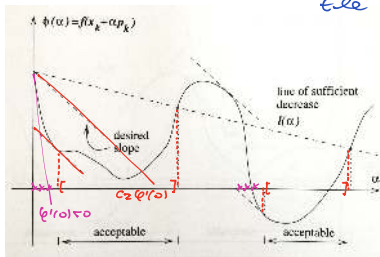
$$\phi(\alpha) \leq \phi(0) + c_1 \alpha \phi'(0) \quad (2)$$

$$\phi'(\alpha) \geq c_2 \phi'(0) \quad (3)$$

with  $0 < c_1 < c_2 < 1$ .

there are also other conditions but these are the most common ones

we can also cut too large values  $\alpha$ , using the old value  $\alpha_{j-1}$



we run through the  $\alpha$  axis, check the value there of  $\phi(\alpha)$ , and leaving on the conditions we decide to have (accept) or not test a value

See Chap. 3 of J. Nocedal, S. Wright, Numerical Optimization, Springer 1999.

## Strong Wolfe conditions:

$$\phi(\alpha) \leq \phi(0) + c_1 \alpha \phi'(0) \quad (4)$$

$$|\phi'(\alpha)| \leq c_2 |\phi'(0)| \quad (5)$$

with  $0 < c_1 < c_2 < 1$ .

Exclude values of  $\alpha$  with  $\phi'(\alpha)$  too positive, far from stationary points of  $\phi$ .

Conditions are invariant w.r.t. affine transformation of the variables.

## **Proposition:**

If  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is  $C^1$  and  $\underline{d}_k$  descent direction at  $\underline{x}_k$  such that  $f$  is bounded below along  $\{\underline{x}_k + \alpha \underline{d}_k : \alpha > 0\}$ . Then if  $0 < c_1 < c_2 < 1$  there exist intervals of step lengths satisfying the Wolfe conditions (weak and strong).

Simple consequence of the mean value theorem.

*still we need to see how can we leave, and which values we will focus on the weak Wolfe case.*

## Methods for 1-D search

Many methods (with/without derivatives) to determine an approximate solution  $\alpha_k$  of

$$\min_{\alpha \geq 0} \phi(\alpha) = f(\underline{x}_k + \alpha \underline{d}_k)$$

satisfying appropriate conditions (e.g. Wolfe) which guarantee global convergence.

In general, two phases:

- determine  $[\alpha_{min}, \alpha_{max}]$  containing “acceptable” step lengths (“bracketing phase”),
- select a good value  $\alpha$  within  $[\alpha_{min}, \alpha_{max}]$  via bisection or interpolation.

## Bisection

$\phi \in \mathcal{C}^1$ ,  $\phi'(0) < 0$  since  $\underline{d}_k$  descent direction and  $\exists \bar{\alpha}$  such that  $\phi'(\alpha) > 0$  for  $\alpha \geq \bar{\alpha}$ .

Start from  $[\alpha_{min}, \alpha_{max}]$  with  $\phi'(\alpha_{min}) < 0$  and  $\phi'(\alpha_{max}) > 0$  and iteratively reduce it.

Iteration: set  $\tilde{\alpha} = \frac{1}{2}(\alpha_{min} + \alpha_{max})$

**if**  $\phi'(\tilde{\alpha}) > 0$  **then**  $\alpha_{max} := \tilde{\alpha}$

**if**  $\phi'(\tilde{\alpha}) < 0$  **then**  $\alpha_{min} := \tilde{\alpha}$

Linear convergence with rate 1/2

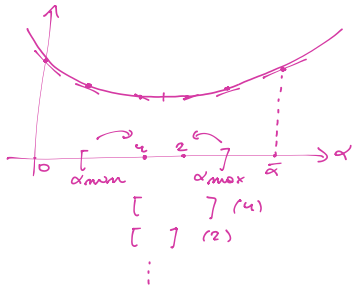
To find initial  $[\alpha_{min}, \alpha_{max}]$ :

1)  $\alpha_{min} := 0$  e  $s := s_0$

2) compute  $\phi'(s)$

**if**  $\phi'(s) < 0$  **then**  $\alpha_{min} := s$ ,  $s := 2s$ , **goto** 2)

**if**  $\phi'(s) > 0$  **then**  $\alpha_{max} := s$ , **stop**



Adaptation to determine  $\alpha_k$  satisfying Wolfe conditions.

Procedure:

i) select  $\alpha > 0$  and set  $\alpha_{min} = \alpha_{max} = 0$

ii) **if**  $\alpha$  satisfies Wolfe (2) **then** goto iii)

**else**  $\alpha_{max} := \alpha$ ,  $\alpha := \frac{\alpha_{min} + \alpha_{max}}{2}$ , **goto** ii)

iii) **if**  $\alpha$  satisfies Wolfe (3) **then**  $\alpha_k = \alpha$ , **stop**

**else**  $\alpha_{min} := \alpha$

$$\alpha := \begin{cases} 2\alpha_{min} & \text{if } \alpha_{max} = 0 \\ \frac{\alpha_{min} + \alpha_{max}}{2} & \text{if } \alpha_{max} > 0 \end{cases}$$

**goto** ii)

**Proposition:** If  $f \in \mathcal{C}^1$  is bounded below along ray  $\{\underline{x}_k + \alpha \underline{d}_k : \alpha \geq 0\}$ , the procedure stops after a finite number of iterations and yields  $\alpha_k$  satisfying Wolfe conditions.

## 4) Global convergence of line search methods

Suitable assumptions on  $\alpha_k$  and  $\underline{d}_k$  can guarantee global convergence.

*steepest descent  
(with -) direction*

**Key aspect:** angle  $\theta_k$  between  $\underline{d}_k$  and  $-\nabla f(\underline{x}_k)$

$$\cos(\theta_k) = -\frac{\nabla^t f(\underline{x}_k) \underline{d}_k}{\|\nabla f(\underline{x}_k)\| \|\underline{d}_k\|}$$

*since  $\underline{a} \cdot \underline{b} = \|\underline{a}\| \|\underline{b}\| \cos(\theta)$*

General result showing how far  $\underline{d}_k$  can deviate from  $-\nabla f(\underline{x}_k)$  and still give rise to globally convergent iterations.

For a proof assuming weak Wolfe conditions, see J. Nocedal, S. Wright, Numerical Optimization, Springer 1999, p. 43-44.

## Theorem: (Zoutendijk)

Consider any line search method iteration with descent  $d_k$  and  $\alpha_k$  satisfying Wolfe conditions. Suppose  $f$  is bounded below on  $\mathbb{R}^n$ ,  $f \in C^1$  on open set  $N$  containing  $L_0 = \{x \in \mathbb{R}^n : f(x) \leq f(x_0)\}$  and  $\nabla f(x)$  is Lipschitz continuous on  $N$ , i.e.,  $\exists L > 0$  such that

$$\|\nabla f(x) - \nabla f(\bar{x})\| \leq L\|x - \bar{x}\| \quad \forall x, \bar{x} \in N.$$

moreover, bounded curvature

Then

$$\sum_{k \geq 0} \cos^2(\theta_k) \|\nabla f(x_k)\|^2 < +\infty. \quad (6)$$

we sum over all iterations

if the series does not diverge, it means that the argument  $\rightarrow 0$  for  $k \rightarrow \infty$ , that is:

$$\cos(\theta_k)^2 \|\nabla f(x_k)\|^2 \xrightarrow{k \rightarrow \infty} 0$$

$$\text{if } \theta_k \rightarrow 0 \Rightarrow \cos(\theta_k) \rightarrow 1$$

this holds if the directions  $d_k$  are not too close to orthogonal (not) with  $-\nabla f(x_k)$

ie  $d_k$  are still "close" to the  $-\nabla f(x_k)$  direction

$\Rightarrow$  the gradient method is globally convergent, ie  $\lim_{k \rightarrow \infty} \|\nabla f(x_k)\| = 0 \quad \forall x_0$

Consequence: The gradient method ( $\cos \theta_k = 1$ ) satisfying Wolfe conditions is globally convergent.



If  $D_k$  symmetric and p.d.  $\forall k \geq 0$  and  $\exists$  constant  $M$  such that

$$\|D_k\| \|D_k^{-1}\| \leq M \quad \forall k \geq 0$$

(bounded condition number), it can be verified that

$$\cos \theta_k \geq 1/M.$$

*the  $\theta$  of before w/p  
these are related  
to the condition  
number of  $D_k$*

In such cases Newton and quasi-Newton methods are globally convergent.

## 4.5 Gradient method

(+ computationally - best  
+ globally convergent  
- slow convergence)

Given  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  with  $f \in C^1$ , look for a stationary point.

**Gradient method** with exact 1-D search:

Choose  $\underline{x}_0$ , set  $k := 0$

Iteration  $k$ :

$$\underline{d}_k := -\nabla f(\underline{x}_k)$$

Determine  $\alpha_k > 0$  such that  $\min_{\alpha \geq 0} \phi(\alpha) = f(\underline{x}_k + \alpha \underline{d}_k)$

$$\underline{x}_{k+1} := \underline{x}_k + \alpha_k \underline{d}_k$$

$$k := k + 1$$

exact since we set/choose  
the  $\alpha$  by solving it  
analytically that we  
reach no local

$f(\cdot)$  restricted on values  
on the direction  $\underline{d}_k$

Termination criteria:  $\|\nabla f(\underline{x}_k)\| < \varepsilon$  or  $|f(\underline{x}_k) - f(\underline{x}_{k+1})| < \varepsilon$  or  $\|\underline{x}_{k+1} - \underline{x}_k\| < \varepsilon$ .

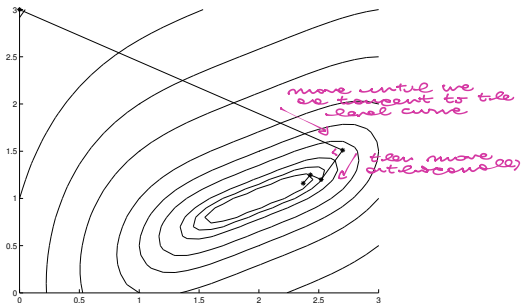
**Property:** If 1-D search is exact, the successive directions are orthogonal.

Since we have that

$$\phi'(\alpha) = \nabla f(\underline{x}_k + \alpha \underline{d}_k) \cdot \underline{d}_k = \nabla f(\underline{x}_k + \alpha \underline{d}_k) \cdot \underline{d}_k = 0$$

$$\Rightarrow \underline{d}_{k+1} \cdot \underline{d}_k = -\nabla f(\underline{x}_k + \alpha \underline{d}_k) \cdot \underline{d}_k = 0$$

then we  
get news!  
e.g.



Example: zig-zag trajectory, very slow convergence

We first consider the case of quadratic strictly convex functions.

Any  $C^2$  function can be well approximated around any local/global minimum by such a function.

## Quadratic strictly convex functions:

$$f(\underline{x}) = \frac{1}{2} \underline{x}^t Q \underline{x} - \underline{b}^t \underline{x} \quad \text{with } Q \text{ symmetric and p.d.}$$

*the Hessian matrix*  
*it's line of strictly convex*

Global minimum is unique solution of  $Q\underline{x} = \underline{b}$  ( $\nabla f(\underline{x}) = \underline{0}$ ) and  $\alpha_k$  can be determined explicitly:

$$\phi(\alpha) = f(\underbrace{\underline{x}_k - \alpha \nabla f(\underline{x}_k)}_{\underline{z}_k}) = \frac{1}{2} (\underline{x}_k - \alpha \nabla f(\underline{x}_k))^t Q (\underline{x}_k - \alpha \nabla f(\underline{x}_k)) - \underline{b}^t (\underline{x}_k - \alpha \nabla f(\underline{x}_k))$$

$$\Rightarrow \phi'(\alpha) = \frac{\partial}{\partial \alpha} \phi(\alpha) = (-\nabla f(\underline{z}_k))^t Q (\underline{z}_k - \alpha \nabla f(\underline{z}_k)) - \underline{b}^t (-\nabla f(\underline{z}_k)) = 0$$

since  $\nabla f(\underline{z}_k)^t = \underline{z}_k^t Q - \underline{b}^t$  (if derivative of wrt  $\underline{z}$ ) we  
and  $\underline{b}^t = \underline{z}_k^t Q - \nabla f(\underline{z}_k)^t$  and we plug it in to get

$$\phi'(\alpha) = -\nabla f(\underline{z}_k)^t Q \underline{z}_k + \alpha \nabla f(\underline{z}_k)^t Q \nabla f(\underline{z}_k) + (-\nabla f(\underline{z}_k)^t + \underline{z}_k^t Q \nabla f(\underline{z}_k)) = 0$$

we solve for  $\alpha$  we get

$$\alpha_k = \frac{\nabla f(\underline{z}_k)^t \nabla f(\underline{z}_k)}{\nabla f(\underline{z}_k)^t Q \nabla f(\underline{z}_k)} = \frac{\underline{z}_k^t Q \underline{z}_k}{\underline{z}_k^t Q \underline{z}_k}$$

## Convergence analysis

Often consider convergence rate of  $f(\underline{x}_k) \rightarrow f(\underline{x}^*)$  instead of  $\|\underline{x}_k - \underline{x}^*\| \rightarrow 0$  when  $k \rightarrow \infty$ .

*is actually its independent  
the choice of what to monitor*

**Proposition:** If  $H(\underline{x}^*)$  is p.d.,  $\underline{x}_k$  converges (super)linearly at  $\underline{x}^*$  w.r.t.  $|f(\underline{x}_k) - f(\underline{x}^*)|$  if and only if it converges in the same way w.r.t.  $\|\underline{x}_k - \underline{x}^*\|$ .

Indeed

$$f(\underline{x}) \approx f(\underline{x}^*) + \frac{1}{2}(\underline{x} - \underline{x}^*)^t H(\underline{x}^*)(\underline{x} - \underline{x}^*)$$

and  $\exists$  a neighborhood  $N(\underline{x}^*)$  such that

*from these bounds we set  
the convergence relation  
of the theorem above*

$$\lambda'_1 \|\underline{x} - \underline{x}^*\|^2 \leq |f(\underline{x}) - f(\underline{x}^*)| \leq \lambda'_n \|\underline{x} - \underline{x}^*\|^2 \quad \forall \underline{x} \in N(\underline{x}^*)$$

with  $\lambda'_1 = \lambda_1 - \varepsilon > 0$  and  $\lambda'_n = \lambda_n + \varepsilon$ , where  $\varepsilon > 0$  and  $0 < \lambda_1 \leq \dots \leq \lambda_n$  are the eigenvalues of  $H(\underline{x}^*)$ . *assumed to be positive def (which was a strong assumption)*

N.B.: This equivalence does not hold in general (e.g., functions non everywhere  $\mathcal{C}^1$ )

Quadratic strictly convex functions:

$$f(\underline{x}) = \frac{1}{2} \underline{x}^t Q \underline{x} - \underline{b}^t \underline{x} \text{ and weighted norm } \|\underline{x}\|_Q^2 := \underline{x}^t Q \underline{x}.$$

Since  $Q \succeq \rho I = \underline{b}$  (?) we have that

$$\frac{1}{2} \|\underline{x} - \underline{x}^*\|_Q^2 = \frac{1}{2} [(\underline{x} - \underline{x}^*)^t Q (\underline{x} - \underline{x}^*)] = \dots = f(\underline{x}) - f(\underline{x}^*)$$

while in practice we rarely  
use the exact version of  
the method

**Theorem:** If gradient method with exact 1-D search is applied to any quadratic strictly convex  $f \in \mathcal{C}^2$ , for any  $\underline{x}_0$  we have  $\lim_{k \rightarrow \infty} \underline{x}_k = \underline{x}^*$  and

$$\|\underline{x}_{k+1} - \underline{x}^*\|_Q^2 \leq \left( \frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1} \right)^2 \|\underline{x}_k - \underline{x}^*\|_Q^2, \quad (1)$$

*the method is (globally) convergent*

where  $0 < \lambda_1 \leq \dots \leq \lambda_n$  are the eigenvalues of  $Q$ .

*convergence is linear  
(due to the square being  
over) and the rate  
it depends on  $Q$*

Proof sketch:

Zoutendijk's theorem implies global convergence.

Since exact 1-D search, easy to verify that

$$\|\underline{x}_{k+1} - \underline{x}^*\|_Q^2 = \left( 1 - \frac{\underline{g}_k^t \underline{g}_k}{(\underline{g}_k^t Q \underline{g}_k)(\underline{g}_k^t Q^{-1} \underline{g}_k)} \right) \|\underline{x}_k - \underline{x}^*\|_Q^2,$$

where  $\underline{g}_k = Q \underline{x}_k - \underline{b} = \nabla f(\underline{x}_k)$

Then just apply Kantorovich inequality:

If  $Q$  p.d. (with  $\lambda_1$  and  $\lambda_n$  smallest and largest eigenvalues), for each  $\underline{x} \neq \underline{0}$  we have

$$\frac{(\underline{x}^t \underline{x})^2}{(\underline{x}^t Q \underline{x})(\underline{x}^t Q^{-1} \underline{x})} \geq \frac{4\lambda_n \lambda_1}{(\lambda_n + \lambda_1)^2}.$$

□

If  $\lambda_1 = \lambda_n$  ( $Q = \gamma I$ ), method "converges" in one iteration.

Upper bound (1) is reached for some choices of  $\underline{x}_0$  (Aikake).

Linear convergence whose rate depends on condition number  $\kappa = \frac{\lambda_n}{\lambda_1}$  of  $Q$ :

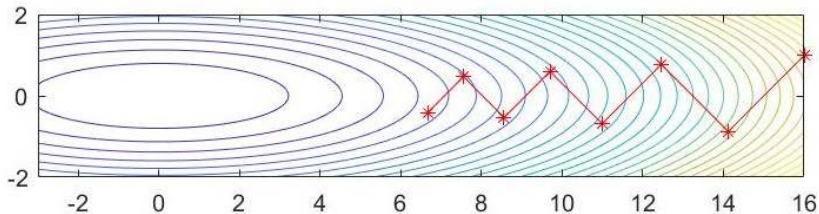
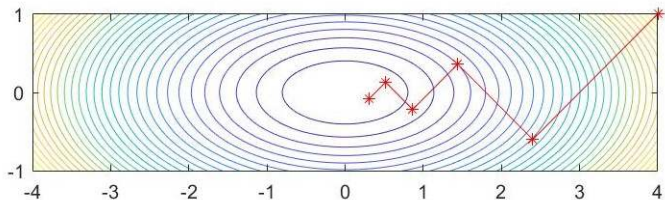
$$r = \left( \frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1} \right) = \left( \frac{\kappa - 1}{\kappa + 1} \right)$$

the closer  $\kappa$  to 1 the smaller  $r$ ; if the spektrum of  $Q$  is very wide then  $\kappa \gg 1$  and  $r \approx 1$ .

*which was slow convergence*

Example:

$$\min f(x_1, x_2) = \frac{1}{2}x_1^2 + \frac{a}{2}x_2^2 \quad \text{with } a \geq 1 \text{ and eigenvalues } \frac{1}{2} \text{ and } \frac{a}{2}$$



Some points of  $\{\underline{x}_k\}$  for  $a = 4$  (top) and  $a = 16$  (bottom), starting from  $\underline{x}_0 = \begin{pmatrix} a \\ 1 \end{pmatrix}$ .



Arbitrary nonlinear functions:

**Theorem:** If  $f \in \mathcal{C}^2$  and gradient method with exact 1-D search converges to  $\underline{x}^*$  with  $H(\underline{x}^*)$  p.d., then

*we have a similar result but involving the f.c. terms*

$$f(\underline{x}_{k+1}) - f(\underline{x}^*) \leq \left( \frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1} \right)^2 [f(\underline{x}_k) - f(\underline{x}^*)]$$

where  $0 < \lambda_1 \leq \dots \leq \lambda_n$  are eigenvalues of  $H(\underline{x}^*)$ .

We cannot expect better convergence with inexact (approximate) 1-D search.

$\alpha_k$  minimizing  $\phi(\alpha)$  may not be the best choice, we could try to "extract" 2nd order information about  $f(\underline{x})$ . *well we assumed  $\mathcal{C}^2$  so we are good*

Example: for  $f(\underline{x})$  quadratic strictly convex,  $\alpha_k = 1/\lambda_{k+1}$  lead to  $\underline{x}^*$  in at most  $n$  iterations!

*size of the space ( $\mathbb{R}^n$ ) on which we are optimizing*

## 4.6 Newton method

Let  $f \in \mathcal{C}^2$  and  $H(\underline{x}) = \nabla^2 f(\underline{x})$ .

Consider quadratic approximation of  $f(\underline{x})$  at  $\underline{x}_k$ :

$$q_k(\underline{x}) := f(\underline{x}_k) + \nabla^t f(\underline{x}_k)(\underline{x} - \underline{x}_k) + \frac{1}{2}(\underline{x} - \underline{x}_k)^t H(\underline{x}_k)(\underline{x} - \underline{x}_k)$$

and choose as  $\underline{x}_{k+1}$  a stationary point of  $q_k(\underline{x})$ , namely

$$\nabla f(\underline{x}_k) + H(\underline{x}_k)(\underline{x}_{k+1} - \underline{x}_k) = \underline{0}.$$

If  $H(\underline{x}_k)$  is not singular,  $H^{-1}(\underline{x}_k)$  exists and

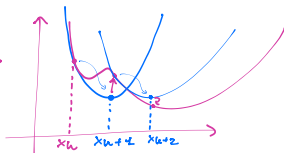
$$\underline{x}_{k+1} := \underline{x}_k - \underbrace{H^{-1}(\underline{x}_k)}_{\underline{d}_k} \nabla f(\underline{x}_k).$$

If  $H(\underline{x}_k)$  is p.d.,  $f \in \mathcal{C}^2$  implies that  $H^{-1}(\underline{x}_k)$  p.d. over  $N(\underline{x}_k)$  and iteration is well defined, otherwise  $\underline{d}_k$  may not be a descent direction.

In the "pure" Newton method,  $\alpha_k = 1$  for each  $k$ .

For  $f$  quadratic and strictly convex, global minimum in a single iteration.

*but this is of course a "local" case, but with H we could have performance issues*



*due to the convex level set inversion here*

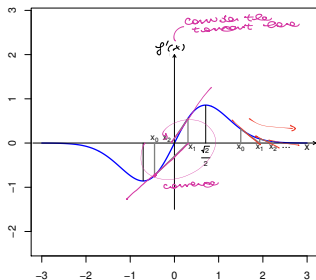
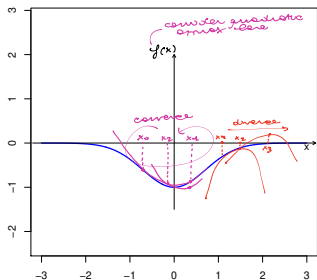
- expensive computation
- (scale) dependent
- + fast convergence

*the extreme opposite of the previous method*

**Property:** Newton method is invariant w.r.t. affine and non singular coordinate changes (see exercise set 6).

**Observation:** Newton method is not globally convergent, but very fast local convergence if  $x_0$  is sufficiently close to a desired solution.

Example:  $\min_{x \in \mathbb{R}} f(x) = -\exp(-x^2)$  with global minimum  $x^* = 0$  and  $f'(x) = 2x \exp(-x^2)$



If  $-0.2 \leq x_0 \leq 0.2$ ,  $\{x_k\}_{k \in \mathbb{N}}$  converges at  $x^* = 0$ . If  $x_0 > 1$ ,  $\{x_k\}_{k \in \mathbb{N}}$  diverges.

**Alternative interpretation** of Newton method (1-D case):

$f(x) \in \mathcal{C}^2$  and look for  $x^*$  such that  $f'(x) = 0$ .

Method of tangents (Newton-Raphson) to determine the zeros of a 1-D function:

At iteration  $k$ ,  $f'(x)$  is approximated with the tangent at  $x_k$

$$z = f'(x_k) + f''(x_k)(x - x_k)$$

$x_{k+1}$  corresponds to the intersection with the  $x$ -axis:  $x_{k+1} = x_k - \frac{f'(x_k)}{f''(x_k)}$

$n$ -D case: Determine a stationary point of  $f(\underline{x})$  by solving non linear system  $\nabla f(\underline{x}) = \underline{0}$  with "Newton-Raphson" method.

**Theorem:** (proof see Nocedal and Wright, 1999 edition, p. 52-53)

Suppose  $f \in \mathcal{C}^2$  and  $\underline{x}^*$  <sup>stationary point</sup> such that  $\nabla f(\underline{x}^*) = \underline{0}$  and  $H(\underline{x}^*)$  <sup>strict local min</sup> p.d. and  $\exists L > 0$  such that

$$\|H(\underline{x}) - H(\underline{y})\| \leq L\|\underline{x} - \underline{y}\| \quad \forall \underline{x}, \underline{y} \in N(\underline{x}^*) \quad \text{(descent) neighborhood of the convex}$$

then, for  $\underline{x}_0$  sufficiently close to local minimum  $\underline{x}^*$ ,

- i)  $\{\underline{x}_k\} \rightarrow \underline{x}^*$  with a quadratic convergence order,
- ii)  $\{\|\nabla f(\underline{x}_k)\|\} \rightarrow 0$  quadratically when  $k \rightarrow \infty$ .

new) new) fast convergence, but  $\underline{x}_0$  may be required to be slow a lot close to  $\underline{x}^*$

Disadvantages:

$\Rightarrow$  idea: end a trust off method

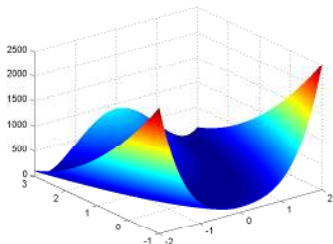
- If  $H(\underline{x}_k)$  is singular the step is not defined.
- If  $H^{-1}(\underline{x}_k)$  is not p.d., Newton direction may not be descent direction.
- Even for a descent direction  $\alpha_k = 1$  <sup>"take the pure newton step"</sup> may increase the value of  $f$ .
- Computation of  $H^{-1}(\underline{x}_k)$  at each iteration (  $O(n^3)$  complexity ).
- Only locally convergent: if  $\underline{x}_0$  is not close enough to  $\underline{x}^*$ ,  $\{\underline{x}_k\}_{k \geq 0}$  may not converge.
- Since  $\{\underline{x}_k\}_{k \geq 0}$  converges from any  $\underline{x}_0$  sufficiently close to any stationary point with non singular  $\nabla^2 f(\underline{x})$ , it may converge to local maxima.

For a comparison between gradient and Newton methods, see Nocedal and Wright, Numerical Optimization, Edition 1999, p. 199.

Rosenbrock function

$$f(\underline{x}) = 100(x_2 - x_1^2)^2 + (1 - x_1)^2.$$

quadratic and nonconvex.



**Fourth computer laboratory:** explore the considerable difference in convergence speed between various line search methods.

## Modifications and extensions

1) If  $\alpha_k = 1$  does not satisfy Wolfe (alternative) conditions then inexact 1-D search.

*on update are useful for local convergence*

2) To guarantee global convergence *(gradient)* *I* or *(Newton)* *H<sup>-1</sup>*

$$\underline{d}_k = -D_k \nabla f(\underline{x}_k)$$

with  $D_k \neq [\nabla^2 f(\underline{x}_k)]^{-1}$ . If  $D_k$  is symmetric and p.d.,  $\underline{d}_k$  is a descent direction.

Trade-off between steepest descent and Newton directions:

$$D_k := (\varepsilon_k I + \nabla^2 f(\underline{x}_k))^{-1}$$

where  $\varepsilon_k > 0$  are smallest values such that eigenvalues of  $(\varepsilon_k I + \nabla^2 f(\underline{x}_k))$  are  $\geq \delta > 0$ .

Such  $\varepsilon_k$  making  $D_k$  p.d. always exist.

*but we still need to invert the matrix*

Coincides with “pure” Newton method when getting closer to a local minimum.

### 3) Trust region methods

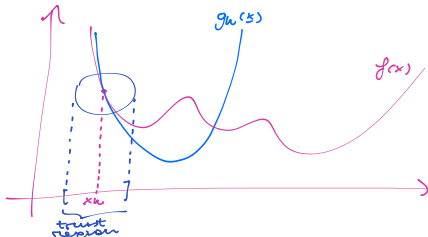
another way to elaborate the Newton method

Idea: simultaneously determine  $\underline{d}_k$  and  $\alpha_k$  by minimizing local quadratic approximation  $q_k(\underline{x})$  at  $\underline{x}_k$  over a trust region on which  $q_k(\underline{x})$  provides a good approximation of  $f(\underline{x})$ .

Example:  $B_k = \{\underline{x} \in \mathbb{R}^n : \|\underline{x} - \underline{x}_k\| \leq \Delta_k\}$

it becomes constrained opt. problem

Illustration:



min  $q_k(x)$   
st  $x \in B_k$

In general, *trust region subproblem*  $\min_{\underline{x} \in B_k} q_k(\underline{x})$  can be solved in closed form or it has low computational requirements.

The trust region size (e.g.  $\Delta_k$ ) is updated adaptively during the iterations based on an estimate of the quality (e.g.  $\max |f(\underline{x}) - q_k(\underline{x})|$ ) of the quadratic approximation over it.

$q_k(\cdot)$  and  $f(\cdot)$  similar  $\Rightarrow$  longer  $\Delta_k$   
" " " " different  $\Rightarrow$  smaller  $\Delta_k$

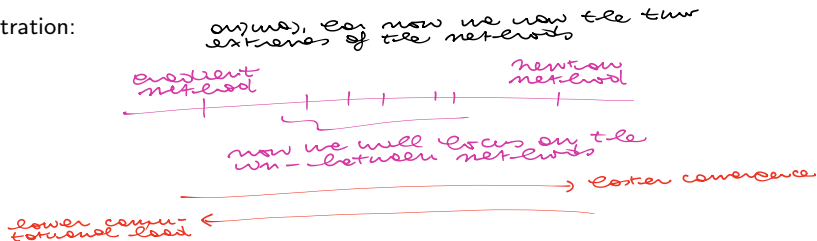


### 3) Trust region methods

**Idea:** simultaneously determine  $\underline{d}_k$  and  $\alpha_k$  by minimizing local quadratic approximation  $q_k(\underline{x})$  at  $\underline{x}_k$  over a *trust region* on which  $q_k(\underline{x})$  provides a good approximation of  $f(\underline{x})$ .

Example:  $B_k = \{\underline{x} \in \mathbb{R}^n : \|\underline{x} - \underline{x}_k\| \leq \Delta_k\}$

Illustration:



In general, *trust region subproblem*  $\min_{\underline{x} \in B_k} q_k(\underline{x})$  can be solved in closed form or it has low computational requirements.

The trust region size (e.g.  $\Delta_k$ ) is updated adaptively during the iterations based on an estimate of the quality (e.g.  $\max |f(\underline{x}) - q_k(\underline{x})|$ ) of the quadratic approximation over it.



## 4.7 Conjugate direction methods

Aim: faster convergence than gradient method and lower computational load than Newton method.

First consider quadratic strictly convex functions

*even when it's called slow  
non-linear functions will  
be of that shape*

$$\min_{x \in \mathbb{R}^n} q(x) = \frac{1}{2} x^t Q x - b^t x$$

with  $Q$   $n \times n$  symmetric and p.d.

**Definition:** Given  $n \times n$  and symmetric  $Q$ , two nonzero  $d_1, d_2 \in \mathbb{R}^n$  are  $Q$ -conjugate if  $d_1^t Q d_2 = 0$ . *extension of orthogonality* *why are they useful?*

Example:

$$f(x_1, x_2) = 12x_2 + 4x_1^2 + 4x_2^2 - 4x_1x_2$$
$$\Rightarrow Q = \begin{pmatrix} 8 & -4 \\ -4 & 8 \end{pmatrix} \quad \text{if } d_1 = \begin{pmatrix} a \\ 0 \end{pmatrix} \text{ then } d_2 = \begin{pmatrix} 0 \\ b \end{pmatrix} \text{ must satisfy}$$
$$d_1^t Q d_2 = 8a - 4b = 0 \Rightarrow b = 2a$$

*and choosing  $a$  we can get a conjugate direction*

**Proposition:** If  $Q$  p.d. and nonzero  $d_0, \dots, d_k$  are mutually  $Q$ -conjugate, then  $d_0, \dots, d_k$  are linearly independent.

Proof: *if  $\exists \lambda_0, \dots, \lambda_k$  st  $\sum d_i \lambda_i = 0$  then we can pre/post multiply either sides by  $Q d_i^T$  we get*

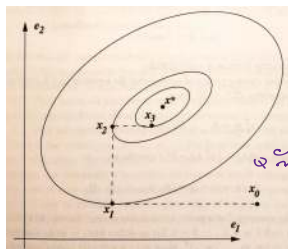
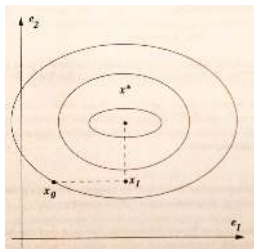
$$0 + \lambda_i d_i (Q d_i^T) = 0 (Q d_i^T) = 0$$

$$\Rightarrow \lambda_i \underbrace{(d_i^T Q d_i)}_{\substack{\neq 0 \text{ since } Q \\ \text{not self}}} = 0 \quad (\Leftrightarrow) \quad \lambda_i = 0 \quad \forall i$$

## Geometric/algebraic interpretation

If  $Q$  is diagonal,  $q(x)$  can be minimized via 1-D search along coordinate directions.

*easy to search*



*core of  $Q$  not diagonal*

Nocedal and Wright, Numerical Optimization, Edition 1999, p. 104-105.



## Theorem: (Conjugate directions)

Let  $\{\underline{d}_i\}_{i=0}^{n-1}$  be  $n$  nonzero mutually  $Q$ -conjugate directions.

For any  $\underline{x}_0 \in \mathbb{R}^n$ ,  $\{\underline{x}_k\}_{k \geq 0}$  generated according to *the usual line search scheme*

$$\underline{x}_{k+1} = \underline{x}_k + \alpha_k \underline{d}_k \quad \text{--- Ent with the line search scheme} \quad (1)$$

with

$$\alpha_k = -\frac{\underline{g}_k^t \underline{d}_k}{\underline{d}_k^t Q \underline{d}_k} \quad \text{and} \quad \underline{g}_k := \nabla q(\underline{x}_k) = Q \underline{x}_k - \underline{b}$$

*since  $q(\cdot)$  is a quadratic form*

terminates to the (unique) global optimal solution  $\underline{x}^*$  of  $q(\underline{x})$  in at most  $n$  iterations, that is

$$\underline{x}_n = \underline{x}_0 + \sum_{k=0}^{n-1} \alpha_k \underline{d}_k = \underline{x}^*.$$

Proof:

Since  $\underline{d}_k$ 's are linearly independent,  $\exists \alpha_k$ 's such that

$$\underline{x}^* - \underline{x}_0 = \alpha_0 \underline{d}_0 + \dots + \alpha_{n-1} \underline{d}_{n-1}.$$

Multiplying by  $\underline{d}_k^T \nabla$  we get  
 $\underline{d}_k^T \nabla (\underline{x}^* - \underline{x}_0) =$  (all terms except the  $\alpha_k$  term)

$$\Rightarrow \alpha_k = \frac{\underline{d}_k^T \nabla (\underline{x}^* - \underline{x}_0)}{\underline{d}_k^T \nabla \underline{d}_k}$$

Now, following the iterative process (4) of the line search scheme we get

$$\underline{z}_k - \underline{x}_0 = \alpha_0 \underline{d}_0 + \dots + \alpha_{k-1} \underline{d}_{k-1}$$

and since the mutual  $\nabla$ -orthogonality of the  $\underline{d}_k$ 's implies that  $\underline{d}_k^T \nabla (\dots) = \underline{d}_k^T \nabla (\dots)$

$$\Rightarrow \underline{d}_k^T \nabla (\underline{z}_k - \underline{x}_0) = 0$$

So now we can rewrite

$$\alpha_k = \frac{\underline{d}_k^T \nabla (\underline{x}^* - \underline{z}_k + \underline{z}_k - \underline{x}_0)}{\underline{d}_k^T \nabla \underline{d}_k} = \frac{\underline{d}_k^T \nabla (\underline{x}^* - \underline{z}_k)}{\underline{d}_k^T \nabla \underline{d}_k} =$$

$$= \frac{\underline{d}_k^T \underline{g}_k}{\underline{d}_k^T \nabla \underline{d}_k}$$

since  $\underline{g}_k = \nabla g(\underline{z}_k) = \nabla \underline{c} \underline{z}_k - \underline{b}$  and the fact that  $\nabla \underline{c} \underline{x}^* = \underline{b}$

## Property: (Expanding subspace)

Let  $\underline{d}_0, \dots, \underline{d}_{n-1}$  be nonzero mutually  $Q$ -conjugate vectors. Then, for any  $\underline{x}_0 \in \mathbb{R}^n$ ,

$\{\underline{x}_k\}_{k \geq 0}$  generated according to

*this sequence*

$$\underline{x}_{k+1} = \underline{x}_k + \alpha_k \underline{d}_k \quad \text{with} \quad \alpha_k = -\frac{\underline{g}_k^t \underline{d}_k}{\underline{d}_k^t Q \underline{d}_k}$$

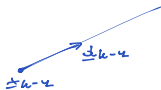
is such that

$$\underline{x}_k = \underline{x}_0 + \sum_{j=0}^{k-1} \alpha_j \underline{d}_j$$

*all what we just derived before*

minimizes  $q(\underline{x}) = \frac{1}{2} \underline{x}^t Q \underline{x} - \underline{b}^t \underline{x}$  not only on the line

$$\{ \underline{x} \in \mathbb{R}^n : \underline{x} = \underline{x}_{k-1} + \alpha \underline{d}_{k-1}, \alpha \in \mathbb{R} \}$$



but also on the affine subspace  $V_k = \{ \underline{x} \in \mathbb{R}^n : \underline{x} = \underline{x}_0 + \text{span}\{\underline{d}_0, \dots, \underline{d}_{k-1}\} \}$ .

In particular,  $\underline{x}_n$  is the global optimum of  $q(\underline{x})$  on  $\mathbb{R}^n$ .

*since all the d's are a basis of the full R^n*



Proof:

*Since  $q(\underline{x})$  is strictly convex,  $\exists \Phi$  is a global min of  $q(\underline{x})$  over  $\mathbb{C} \subseteq \mathbb{R}^n$   $\Leftrightarrow$*

*more @ it works since we optimize one direction at a time (moving the exact that the) are  $\perp$*



## Property: (Expanding subspace)

Let  $\underline{d}_0, \dots, \underline{d}_{n-1}$  be nonzero mutually  $Q$ -conjugate vectors. Then, for any  $\underline{x}_0 \in \mathbb{R}^n$ ,

$\{\underline{x}_k\}_{k \geq 0}$  generated according to

$$\underline{x}_{k+1} = \underline{x}_k + \alpha_k \underline{d}_k \quad \text{with} \quad \alpha_k = -\frac{\underline{g}_k^t \underline{d}_k}{\underline{d}_k^t Q \underline{d}_k}$$

is such that

$$\underline{x}_k = \underline{x}_0 + \sum_{j=0}^{k-1} \alpha_j \underline{d}_j$$

minimizes  $q(\underline{x}) = \frac{1}{2} \underline{x}^t Q \underline{x} - \underline{b}^t \underline{x}$  not only on the line

$$\{ \underline{x} \in \mathbb{R}^n : \underline{x} = \underline{x}_{k-1} + \alpha \underline{d}_{k-1}, \alpha \in \mathbb{R} \}$$

but also on the affine subspace  $V_k = \{ \underline{x} \in \mathbb{R}^n : \underline{x} = \underline{x}_0 + \text{span}\{\underline{d}_0, \dots, \underline{d}_{k-1}\} \}$ .

In particular,  $\underline{x}_n$  is the global optimum of  $q(\underline{x})$  on  $\mathbb{R}^n$ .

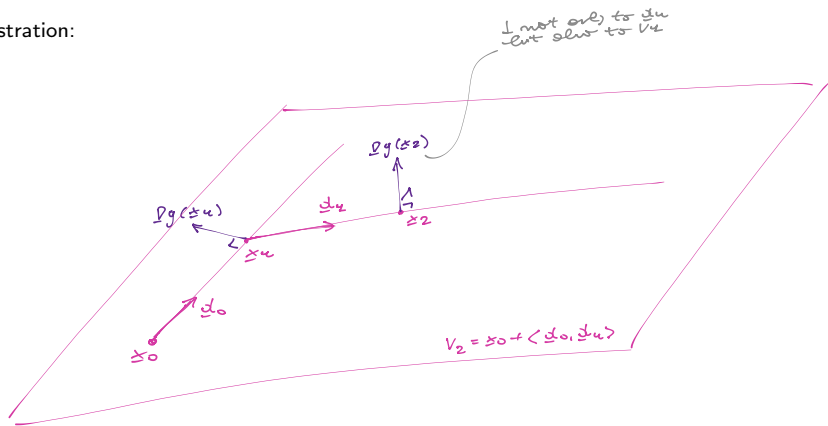
Proof:

Since  $q(\underline{x})$  is strictly convex,

$\underline{x}^*$  is a global min of  $q(\underline{x})$  over  $C \subseteq \mathbb{R}^n$   $\Leftrightarrow \nabla q(\underline{x}^*)^t (\underline{x} - \underline{x}^*) \geq 0$  for all  $\underline{x} \in C$

So we need to check that  $\nabla q(\underline{x}_k)^t (\underline{x} - \underline{x}_k) \geq 0$  for all  $\underline{x} \in V_k$  for  $k=1, \dots, n$

Illustration:



**Consequence:** In conjugate direction method the gradients  $\underline{g}_k$  satisfy  $\underline{g}_k^t \underline{d}_i = 0$  for all  $i$  with  $1 \leq i < k$ .

## 4.7.1 Conjugate gradient method for quadratic convex functions

Initialization: Arbitrary  $\underline{x}_0$ ,  $\underline{g}_0 = \nabla q(\underline{x}_0) = Q\underline{x}_0 - \underline{b}$ ,  $\underline{d}_0 := -\underline{g}_0$  and  $k = 0$

Iteration:  $\underline{x}_{k+1} := \underline{x}_k + \alpha_k \underline{d}_k$  with  $\alpha_k = -\frac{\underline{g}_k^t \underline{d}_k}{\underline{d}_k^t Q \underline{d}_k}$  (exact 1-D search)

$$\underline{d}_{k+1} := -\underline{g}_{k+1} + \beta_k \underline{d}_k \quad \text{with} \quad \beta_k = \frac{\underline{g}_{k+1}^t Q \underline{d}_k}{\underline{d}_k^t Q \underline{d}_k}$$

*steepest descent dir*

*but there was null matrix & here*

*this solves the problem of how to derive the Q-conj directions*

Observations:

- $\alpha_k = -\frac{\underline{g}_k^t \underline{d}_k}{\underline{d}_k^t Q \underline{d}_k}$  minimizes  $q(\underline{x})$  along line through  $\underline{x}_k$  generated by  $\underline{d}_k$

$$\frac{dq(\underline{x}_k + \alpha \underline{d}_k)}{d\alpha} = \underline{d}_k^t Q(\underline{x}_k + \alpha \underline{d}_k) - \underline{b}^t \underline{d}_k \stackrel{!}{=} 0 \Rightarrow (\dots) \text{ we set result, that derivation of } \alpha_k$$

- Limited computational requirements, no matrix inversions are needed.

To show that global optimal solution is found after at most  $n$  iterations, just verify that directions are mutually Q-conjugated.

*it's just algebra*

## Proposition:

At each iteration  $k$  in which the optimum solution of  $q(x)$  has not yet been found ( $\underline{g}_i \neq \underline{0}$  for  $i = 0, \dots, k$ )

i)  $\underline{d}_0, \dots, \underline{d}_{k+1}$  generated are mutually Q-conjugate

(ii)  $\alpha_k = \frac{\underline{g}_k^t \underline{g}_k}{\underline{d}_k^t Q \underline{d}_k} \neq 0$  *not really relevant, it is just needed for the next step*

iii)  $\beta_k = \frac{\underline{g}_{k+1}^t (\underline{g}_{k+1} - \underline{g}_k)}{\underline{g}_k^t \underline{g}_k} = \frac{\underline{g}_{k+1}^t \underline{g}_{k+1}}{\underline{g}_k^t \underline{g}_k}$  *we then learn in the new step, for (strictly) convex, since it does not need the matrix Q (which is constant only for strictly convex functions)*

$$= \frac{\|\underline{g}_{k+1}\|^2}{\|\underline{g}_k\|^2}$$

Advantages: No need for matrix inversions, limited computational requirements.

Disadvantages:

- Exact or at least accurate 1-D search otherwise the directions may lose Q-conjugacy.
- The method is not invariant w.r.t. affine transformations of the coordinates.

*which is a really bad news*

**Fourth computer laboratory:** compare the convergence speed of gradient, conjugate gradient and Newton methods.

## 4.7.2 Conjugate direction methods

vs conjugate gradient (of course)

the real method to use in the case of large scale (even  $n$ ) problems

For arbitrary functions with large  $n$ , approximate  $\alpha_k$  and  $\beta_k$  must not depend on Hessian.

Arbitrary  $\underline{x}_0$  and  $\underline{d}_0 = -\nabla f(\underline{x}_0)$

search for the  $\mathcal{Q}$ -conj direction

$\underline{x}_{k+1} := \underline{x}_k + \alpha_k \underline{d}_k$  with inexact 1-D search and  $\underline{d}_{k+1} = -\nabla f(\underline{x}_{k+1}) + \beta_k \underline{d}_k$ .

Most popular formulae for  $\beta_k$ :

$$\beta_k^{FR} = \frac{\|\nabla f(\underline{x}_{k+1})\|^2}{\|\nabla f(\underline{x}_k)\|^2}$$

Fletcher-Reeves

RHS symmetric expr of the previous one

$$\beta_k^{PR} = \frac{\nabla^t f(\underline{x}_{k+1})(\nabla f(\underline{x}_{k+1}) - \nabla f(\underline{x}_k))}{\|\nabla f(\underline{x}_k)\|^2}$$

Polak-Ribière

LHS expr of the previous one

Observation:  $\underline{d}_k$  is a descent direction if exact 1-D search

$$\nabla^t f(\underline{x}_k) \underline{d}_k = -\|\nabla^t f(\underline{x}_k)\|^2 + \beta_{k-1} \nabla^t f(\underline{x}_k) \underline{d}_{k-1} = -\|\nabla^t f(\underline{x}_k)\|^2 < 0.$$

o.k. from 1 (above)

For quadratic functions the method coincides with CG method.

while we denote as CD this extension to arbitrary functions

For nonquadratic functions, Polak-Ribière version turns out to be more efficient than Fletcher-Reeves one.

### Observations

*no computational cost  
is very low*

- At each iteration it suffices to store  $\underline{x}_k$ ,  $\nabla f(\underline{x}_k)$ ,  $\nabla f(\underline{x}_{k+1})$  and  $\underline{d}_k$ .
- Version with “restart” in which  $\beta_k = 0$  every  $m$  iterations ( $m \ll n$ ) is globally convergent.

*where ever, more time we take a  
steepest descent step*

When  $\beta_k = 0$ ,  $\underline{d}_{k+1} = -\nabla f(\underline{x}_{k+1})$  and all previous information is lost.

For large  $n$ , we hope to find a solution way before  $n$  iterations!

*no this restart version is just  
a theoretical comment, useless  
in practice*

## 4.7.3 Convergence

### Convergence for quadratic functions

Let  $q(\underline{x}) = \frac{1}{2}\underline{x}^t Q \underline{x} - \underline{b}^t \underline{x}$  be quadratic strictly convex with  $\lambda_1 \leq \dots \leq \lambda_n$  the eigenvalues of  $Q$ , then

$$\| \underline{x}_{k+1} - \underline{x}^* \|_Q^2 \leq \left( \frac{\lambda_{n-k} - \lambda_1}{\lambda_{n-k} + \lambda_1} \right)^2 \| \underline{x}_0 - \underline{x}^* \|_Q^2$$

*- still Q norm*  
*- still uses*  
*BUT*

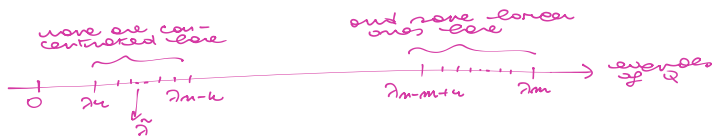
where  $\| \underline{x} - \underline{x}^* \|_Q^2 = (\underline{x}^t - \underline{x}^{*t}) Q (\underline{x} - \underline{x}^*) = 2(q(\underline{x}) - q(\underline{x}^*))$ .

*- error is no error*  
*- neglect more of the*  
*lowest eigenvalues*

If  $m$  large eigenvalues and other  $n - m$  “concentrated” around a  $\tilde{\lambda}$ , after  $m + 1$  iterations

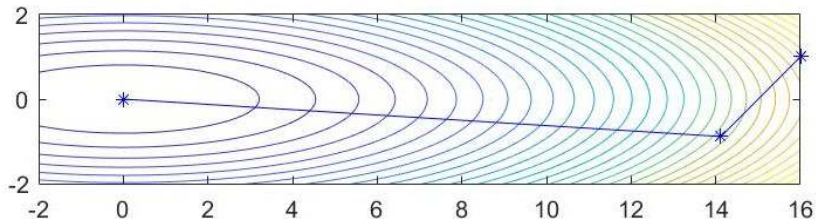
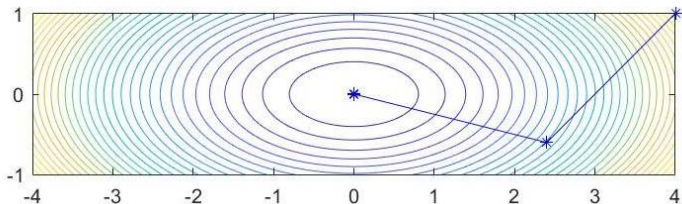
$\| \underline{x}_{m+1} - \underline{x}^* \|_Q \approx \varepsilon \| \underline{x}_0 - \underline{x}^* \|_Q$  with  $\varepsilon = (\lambda_{n-m} - \lambda_1) / 2\tilde{\lambda}$ , that is, we have an accurate estimate of the solution after  $m + 1$  iterations.

Illustration:



Example:

$$\min f(x_1, x_2) = \frac{1}{2}x_1^2 + \frac{a}{2}x_2^2 \quad \text{with } a \geq 1 \text{ and hence eigenvalues } \frac{1}{2} \text{ and } \frac{a}{2}$$



The sequence  $\{x_k\}$  for  $a = 4$  (top) and  $a = 16$  (bottom), starting from  $x_0 = \begin{pmatrix} a \\ 1 \end{pmatrix}$ .



## Convergence for arbitrary functions

1) If  $f \in \mathcal{C}^2$  and  $\{\underline{x}_k\}_{k \geq 0}$  generated by the F-R method with exact 1-D search converges to  $\underline{x}^*$  with p.d.  $H(\underline{x}^*)$ , then

*we  $\pm$   $\nabla$  has a strict local min*

$$\lim_{k \rightarrow \infty} \frac{\|\underline{x}_{k+n} - \underline{x}^*\|}{\|\underline{x}_k - \underline{x}^*\|} = 0,$$

$$\frac{\|\text{err}(k+n)\|}{\|\text{err}(k)\|}$$

*which resembles the ratio  $r \rightarrow 0$*

namely convergence is superlinear within  $n$  iterations.

*no new error*

Similar result also for inexact 1-D search.

2) Global convergence of F-R method even without "restart" (for P-R?).

Zoutendijk's theorem implies:

For F-R method with inexact 1-D search satisfying strong Wolfe conditions with  $0 < c_1 < c_2 < 1/2$ , we have

$$\liminf_{k \rightarrow \infty} \|\nabla f(\underline{x}_k)\| = 0.$$

A sub-sequence has  $\|\nabla f(\underline{x}_k)\|$  that converges to 0.

*due to the inf presence*

*to solve the issue of not being CG invariant to affine transformations*

## 4.7.4 Preconditioned conjugate gradient method

The conjugate gradient method (CG) can be accelerated by a variable change  $\underline{x} = \underline{S}\underline{y}$ , where  $\underline{S}$  is  $n \times n$  symmetric and non singular.

By applying CG to

$$h(\underline{y}) = q(\underline{S}\underline{y}) = \frac{1}{2}\underline{y}^t \underline{S} \underline{Q} \underline{S} \underline{y} - \underline{b}^t \underline{S} \underline{y}$$

*issue: find S at S^T S becomes well conditioned*

we obtain

$$\underline{y}_{k+1} = \underline{y}_k + \alpha_k \tilde{\underline{d}}_k$$

*direction in the space*

with  $\alpha_k$  determined by 1-D search,  $\tilde{\underline{d}}_0 = -\nabla h(\underline{y}_0)$  and  $\tilde{\underline{d}}_k = -\nabla h(\underline{y}_k) + \beta_{k-1} \tilde{\underline{d}}_{k-1}$  for  $k = 1, \dots, n-1$  where

$$\beta_{k-1} = \frac{\nabla^t h(\underline{y}_k) \nabla h(\underline{y}_k)}{\nabla^t h(\underline{y}_{k-1}) \nabla h(\underline{y}_{k-1})}$$

Setting  $\underline{x}_k = \underline{S}\underline{y}_k$ ,  $\nabla h(\underline{y}_k) = \underline{S}\underline{g}_k$ ,  $\underline{d}_k = \tilde{\underline{d}}_k$ , we obtain the equivalent **preconditioned conjugate gradient method**:

$$\underline{x}_{k+1} = \underline{x}_k + \alpha_k \underline{d}_k$$

with  $\alpha_k$  determined by 1-D search,  $\underline{d}_0 = -\underline{S}\underline{g}_0$  and

$$\underline{d}_k = -\underline{S}\underline{g}_k + \beta_{k-1} \underline{d}_{k-1} \quad \text{for } k = 1, \dots, n-1$$

where

$$\beta_{k-1} = \frac{\underline{\mathbf{g}}_k^t \mathbf{S}^2 \underline{\mathbf{g}}_k}{\underline{\mathbf{g}}_{k-1}^t \mathbf{S}^2 \underline{\mathbf{g}}_{k-1}}.$$

Clearly when  $S = I$  it coincides with the standard CG method.

Since  $\nabla^2 h(\underline{y}) = SQS$ ,  $\tilde{\underline{d}}_0, \dots, \tilde{\underline{d}}_{n-1}$  are  $(SQS)$ -conjugate. Moreover  $\underline{d}_k = S\tilde{\underline{d}}_k$  implies that  $\underline{d}_0, \dots, \underline{d}_{n-1}$  are  $Q$ -conjugate.

To achieve faster convergence, we look for  $S$  such that  $SQS$  has a smaller condition number than  $Q$  or eigenvalues that are distributed into “groups”.

Recall: a good approximate solution can be found in a number of iterations not much larger than the number of groups.

## 4.8 Quasi-Newton methods

Instead of using/inverting  $\nabla^2 f(\underline{x}_k)$ , second order derivative information is extracted from variations in  $\nabla f(\underline{x})$ .

*avoid and iterative*

Generate  $\{H_k\}$  of symmetric p.d. approximations of  $[\nabla^2 f(\underline{x}_k)]^{-1}$  and take

$$\underline{x}_{k+1} = \underline{x}_k + \alpha_k \underline{d}_k \quad \text{with} \quad \underline{d}_k = -H_k \nabla f(\underline{x}_k),$$

where  $\alpha_k > 0$  minimizes  $f(\underline{x})$  along  $\underline{d}_k$  or satisfies some inexact 1-D search conditions.

**Advantages** w.r.t. Newton method:

- since  $H_k$ 's are symmetric and p.d., always well defined and descent direction,
- only involves first order derivatives,
- $H_k$  is constructed iteratively, each iteration is  $O(n^2)$ . *which we could write  $O(n^3)$  of real-Newton method (see matrix inversion tree)*

**Disadvantages** w.r.t. conjugate direction methods: requires storing/handling matrices.

*rather than the conjugate vectors of CG*

**Idea:** Second order derivative information is extracted from  $\nabla f(\underline{x}_k)$  and  $\nabla f(\underline{x}_{k+1})$ .

Quadratic approximation of  $f(\underline{x})$  around  $\underline{x}_k$ :

$$f(\underline{x}_k + \underline{\delta}) \approx f(\underline{x}_k) + \underline{\delta}^t \nabla f(\underline{x}_k) + \frac{1}{2} \underline{\delta}^t \nabla^2 f(\underline{x}_k) \underline{\delta}.$$

Differentiating <sup>(wrt  $\underline{\delta}$ )</sup> we obtain

$$\nabla f(\underline{x}_k + \underline{\delta}) \approx \nabla f(\underline{x}_k) + \nabla^2 f(\underline{x}_k) \underline{\delta}.$$

Substituting  $\underline{\delta}$  with  $\underline{\delta}_k := \underline{x}_{k+1} - \underline{x}_k$  we can move the gradients to the LHS case, and so we also define  $\underline{g}_k = \nabla f(\underline{x}_{k+1}) - \nabla f(\underline{x}_k)$ .

So we have

$$\underline{g}_k \approx \nabla^2 f(\underline{x}_k) \underline{\delta}_k \Leftrightarrow [\nabla^2 f(\underline{x}_k)]^{-1} \underline{g}_k \approx \underline{\delta}_k$$

and this was the rule we was solving for  $H_{k+1}$  matrix ( $H_{k+1}$  since we need  $k+1$  iter to compute  $\underline{\delta}_k$  and  $\underline{g}_k$ )

Since  $\underline{\delta}_k$  and  $\underline{g}_k$  can only be determined after 1-D search, we select  $H_{k+1}$  symmetric and p.d. such that

$$\boxed{H_{k+1} \underline{g}_k = \underline{\delta}_k} \quad \text{(secant condition).} \quad (1)$$

$H_{k+1}$  is not univocally defined:  $n$  equations and  $n(n+1)/2$  degrees of freedom.

(variables)

unusable part we make for 

How do we update  $H_k$ ?

Simple way is by successive updates:

$$H_{k+1} = H_k + a_k \underline{u} \underline{u}^t \quad (2)$$

*a scalar expr on a symmetric rank 1 matrix*

where  $\underline{u} \underline{u}^t$  symmetric matrix of rank 1 and  $a_k$  proportionality coefficient.

To satisfy (1) we must have

$$H_k \underline{\gamma}_k + a_k \underline{u} \underline{u}^t \underline{\gamma}_k = \underline{\delta}_k$$

*scalar*  $\underline{u} \underline{u}^t (\underline{u}^t \underline{\delta}_k) = \underline{\delta}_k - H_k \underline{\delta}_k$   
*scalar*  
 $\Rightarrow \underline{u} \parallel \underline{\delta}_k - H_k \underline{\delta}_k$

and hence  $\underline{u}$  and  $(\underline{\delta}_k - H_k \underline{\gamma}_k)$  must be collinear.

Since  $a_k$  accounts for proportionality, we can set  $\underline{u} = \underline{\delta}_k - H_k \underline{\gamma}_k$  and hence  $a_k \underline{u}^t \underline{\gamma}_k = 1$ .

$$\Rightarrow a_k = \frac{1}{\underline{u}^t \underline{\delta}_k}$$

Rank one update formula:

$$H_{k+1} = H_k + \frac{(\underline{\delta}_k - H_k \underline{\gamma}_k)(\underline{\delta}_k - H_k \underline{\gamma}_k)^t}{(\underline{\delta}_k - H_k \underline{\gamma}_k)^t \underline{\gamma}_k} \quad (3)$$

*scalar*  $\frac{1}{\underline{u}^t \underline{\delta}_k} = a_k$

## Properties

1 For quadratic strictly convex functions,  $H_n = Q^{-1}$  in at most  $n$  iterations, even with inexact 1-D search.

*eventually  $H_k$  becomes  $Q^{-1}$*

2 No guarantee that  $H_k$  is p.d.!

*really, but, we now improve  $\hookrightarrow$*

## Rank two updates

$$H_{k+1} = H_k + a_k \underline{u} \underline{u}^t + b_k \underline{v} \underline{v}^t \quad (4)$$

are more interesting.

To satisfy (1) we have

$$H_k \underline{\gamma}_k + a_k \underbrace{\underline{u} \underline{u}^t}_{\stackrel{!}{=} \delta_k} \underline{\gamma}_k + b_k \underbrace{\underline{v} \underline{v}^t}_{\stackrel{!}{=} H_k \underline{\gamma}_k} \underline{\gamma}_k = \delta_k$$

view the condition as a vector equation with free choices

where  $\underline{u}, \underline{v}$  are not determined univocally.

Setting  $\underline{u} = \delta_k$  and  $\underline{v} = H_k \underline{\gamma}_k$ , we obtain  $a_k \underline{u}^t \underline{\gamma}_k = 1$  and  $b_k \underline{v}^t \underline{\gamma}_k = -1$

and hence the rank two update formula:

$$H_{k+1} = H_k + \frac{\delta_k \delta_k^t}{\delta_k^t \underline{\gamma}_k} - \frac{H_k \underline{\gamma}_k \underline{\gamma}_k^t H_k}{\underline{\gamma}_k^t H_k \underline{\gamma}_k}$$

Davidon-Fletcher-Powell (DFP) (5)

which has more exact properties

**Proposition:** If

$$\delta_{k\underline{k}}^t \gamma_k > 0 \quad \forall k \quad (\text{curvature condition}),$$

as it involves  $\delta_k$  we use the notation of the gradient

the DFP method preserves the positive definiteness of  $H_k$ , i.e., if  $H_0$  is p.d. then  $H_k$  is p.d. for all  $k \geq 1$ .

**Proof:** By induction, suppose that  $H_n$  is p.d. then  $\exists^T H_n u + u \geq 0 \quad \forall u \neq 0$  (we need to show this)

$\forall u$   $H_n$  is p.d. it admits a Cholesky factorization  $H_n = L_n L_n^T$

rearrange the  $u$  vector (convenience) and let  $z = L_n^T u$  and  $b = L_n^T u$

we get (rearrange  $z$  and the  $H_n u$  update rule):

$$\exists^T \left( H - \frac{H z z^T H}{z^T H z} \right) z = z^T z - \frac{(z^T b)^2}{b^T b} \geq 0 \quad (\#)$$

rearrange the term (B)

because of Cauchy-Schwarz  $|z^T b| \leq \|z\| \|b\|$

Since  $z \neq 0$ , the equality holds if  $z$  and  $b$  are collinear, we will see  $z$  and  $z$  are.

Since  $\delta^T z > 0$  (curvature cond) we have that (rearrange the term (B)):

$$z \left( \frac{\delta \delta^T}{\delta^T z} \right) z \geq 0$$

we exist if  $z$  and  $z$  are collinear and we have (#) we get the sum of two  $\geq 0$  terms



**Fact:** The curvature condition  $\underline{\delta}_k^t \underline{\gamma}_k > 0$  holds for every  $k \geq 0$  provided that the 1-D search satisfies (weak or strong) Wolfe conditions.

Proof\*:

$$\begin{aligned} &= \nabla g(\underline{x}_k + \underline{u}) - \nabla g(\underline{x}_k) = \begin{cases} g(\underline{x}) = \underline{x}^T Q \underline{x} - \underline{b}^T \underline{x} \\ \nabla g(\underline{x}) = 2Q\underline{x} - \underline{b}^T \end{cases} \\ &= (2Q\underline{x} - \underline{b}) \Big|_{\underline{x} = \underline{x}_k + \underline{u}} - (2Q\underline{x} - \underline{b}) \Big|_{\underline{x} = \underline{x}_k} = 2Q(\underline{x}_k + \underline{u} - \underline{x}_k) = 2Q\underline{u} \end{aligned}$$

For quadratic strictly convex functions,  $\underline{\gamma}_k = Q\underline{\delta}_k$  implies  $\underline{\delta}_k^t Q \underline{\delta}_k = \underline{\delta}_k^t \underline{\gamma}_k > 0$  because  $Q$  is p.d.

For arbitrary functions:

Weak Wolfe conditions:

$$f(\underline{x}_k + \alpha_k \underline{d}_k) \leq f(\underline{x}_k) + c_1 \alpha_k \nabla^t f(\underline{x}_k) \underline{d}_k \quad (\text{Armijo criterion}) \quad (6)$$

$$\nabla^t f(\underline{x}_k + \alpha_k \underline{d}_k) \underline{d}_k \geq c_2 \nabla^t f(\underline{x}_k) \underline{d}_k \quad (7)$$

with  $0 < c_1 < c_2 < 1$ .

Since  $\underline{\delta}_k = \alpha_k \underline{d}_k$ , (7) implies

$$\nabla^t f(\underline{x}_{k+1}) \underline{\delta}_k \geq c_2 \nabla^t f(\underline{x}_k) \underline{\delta}_k,$$

which in turn implies

$$\underline{\gamma}_k^t \underline{\delta}_k \geq (c_2 - 1) \alpha_k \nabla^t f(\underline{x}_k) \underline{d}_k$$

with  $(c_2 - 1) < 0$ ,  $\alpha_k > 0$ , and  $\nabla^t f(\underline{x}_k) \underline{d}_k < 0$  because  $\underline{d}_k$  is a descent direction. □

## Properties

For quadratic strictly convex functions, DFP method with exact 1-D search:

- 1 terminates in at most  $n$  iterations with  $H_n = Q^{-1}$ ,
- 2 generates  $Q$ -conjugate directions (from  $H_0 = I$  it generates CG directions),
- 3 secant condition is hereditary, i.e.,  $H_i \underline{\gamma}_j = \underline{\delta}_j$  for  $j = 0, \dots, i - 1$ .

*we not only  $H_{k+1} \underline{\delta}_k = \underline{\delta}_k$   
but this is preserved also  
on rest  $\underline{\delta}_k$  and  $\underline{\delta}_h$  vectors*

For arbitrary functions:

- 4 if  $\underline{\delta}_k^t \underline{\gamma}_k > 0$  (curvature condition), all  $H_k$  are p.d. if  $H_0$  is p.d. (hence descent method),
- 5 each iteration is  $O(n^2)$ ,
- 6 superlinear convergence rate (in general only local),
- 7 if  $f(x)$  convex, DFP method with exact 1-D search is globally convergent.

*much better  
than linear*

# BFGS method

is that the idea is to make the approximation matrix easier invertible (analytically)

We can construct an approximation of  $\nabla^2 f(\underline{x}_k)$  rather than of  $[\nabla^2 f(\underline{x}_k)]^{-1}$ .

Since we aim at  $B_k \approx \nabla^2 f(\underline{x}_k)$ ,  $B_k$  must satisfy  $B_{k+1}\underline{\delta}_k = \underline{\gamma}_k$ .

Taking  $B_{k+1} = B_k + a_k \underline{u} \underline{u}^t + b_k \underline{v} \underline{v}^t$ , with similar manipulations, we have:

seem using a rank 2 update

$$B_{k+1} = B_k + \frac{\underline{\gamma}_k \underline{\gamma}_k^t}{\underline{\gamma}_k^t \underline{\delta}_k} - \frac{B_k \underline{\delta}_k \underline{\delta}_k^t B_k}{\underline{\delta}_k^t B_k \underline{\delta}_k}$$

you come of here with  $\underline{\delta}$  and  $\underline{\gamma}$  exchanged

(8)

which should be inverted to obtain  $H_{k+1} = B_{k+1}^{-1}$

By applying twice Sherman–Morrison identity

$$(A + \underline{a} \underline{b}^t)^{-1} = A^{-1} - \frac{A^{-1} \underline{a} \underline{b}^t A^{-1}}{1 + \underline{b}^t A^{-1} \underline{a}}, \quad A \in \mathbb{R}^{n \times n} \text{ non singular, } \underline{a}, \underline{b} \in \mathbb{R}^n, \text{ denominator } \neq 0,$$

we obtain the Broyden Fletcher Goldfarb and Shanno (BFGS) update formula:

$$H_{k+1} = H_k + \left( 1 + \frac{\underline{\gamma}_k^t H_k \underline{\gamma}_k}{\underline{\delta}_k^t \underline{\gamma}_k} \right) \frac{\underline{\delta}_k \underline{\delta}_k^t}{\underline{\delta}_k^t \underline{\gamma}_k} - \frac{H_k \underline{\gamma}_k \underline{\delta}_k^t + \underline{\delta}_k \underline{\gamma}_k^t H_k}{\underline{\delta}_k^t \underline{\gamma}_k}$$
(9)

Indeed  $B_{k+1} H_{k+1} = I$  if  $B_k H_k = I$ .

The BFGS method has same properties 1 to 5 as DFP method.

In practice, it is more robust w.r.t. to rounding errors and inexact 1-D search.

*no BFGS actually is better than DFP*

BFGS and DFP are two extreme cases of unique Broyden family of update formulae:

$$H_{k+1} = (1 - \phi)H_{k+1}^{\text{DFP}} + \phi H_{k+1}^{\text{BFGS}}$$

*convex comb of the DFP and BFGS matrices*

with  $0 \leq \phi \leq 1$ .

**Properties:** (Broyden family)

- $H_{k+1}$  satisfies secant condition and is p.d. if  $\delta_k^t \gamma_k > 0$ .
  - Methods invariant w.r.t. affine variable transformations.
  - If  $f(x)$  quadratic strictly convex, methods with exact 1-D search find  $x^*$  in at most  $n$  iterations ( $H_n = Q^{-1}$ ) and the generated directions are  $Q$ -conjugate.
  - Quasi-Newton methods are much less "sensitive" to inexact 1-D search than CD ones.
- we get the best of CG and the good properties of the newton method*
- and hence this true (see exercise) (sensitive vs a new) good property*

# Convergence of quasi-Newton methods

Complex analysis because approximation of Hessian (inverse) is updated at each iteration.

Convergence speed for  $\{B_k\}$  or  $\{H_k\}$  with inexact 1-D search (Wolfe cond.) where  $\alpha_k = 1$  is tried first:

**Theorem:** (Dennis and Moré)

Consider  $f \in \mathcal{C}^3$  and quasi-Newton method with  $B_k$  p.d. and  $\alpha_k = 1$  for each  $k$ .  
If  $\lim_{k \rightarrow \infty} \underline{x}_k = \underline{x}^*$  with  $\nabla f(\underline{x}^*) = \underline{0}$  and  $\nabla^2 f(\underline{x}^*)$  is p.d.,  $\{\underline{x}_k\}$  converges superlinearly if and only if

*this as in Moré's exacted lemma (see the Po)*  
*stationary point*      *local opt*

$$\lim_{k \rightarrow \infty} \frac{\|(B_k - \nabla^2 f(\underline{x}^*))\underline{d}_k\|}{\|\underline{d}_k\|} = 0. \quad (10)$$

If quasi-Newton  $\underline{d}_k$  approximates Newton direction well enough,  $\alpha_k = 1$  satisfies Wolfe cond. when  $\underline{x}_k \rightarrow \underline{x}^*$ .

*(Po)*  
*we don't ask  $B_k \rightarrow \nabla^2 f(\underline{x}^*)$  but just that the  $B_k$  accuracy (in approximation of the Hessian) increases along iterations  $k$*

Observation: No need that  $B_k \rightarrow \nabla^2 f(\underline{x}^*)$ , it suffices that  $B_k$ 's become increasingly accurate approximations of  $\nabla^2 f(\underline{x}^*)$  along  $\underline{d}_k$ !

The necessary and sufficient condition (10) is satisfied by quasi-Newton methods such as BFGS and DFP.

Comparing the convergence rates of gradient, Newton and BFGS methods:

for Rosenbrock's function, see p. 199 (Chap. 8) of J. Nocedal, S. Wright, Numerical Optimization, Springer, 1999.

Global convergence:

Under some assumptions, can guarantee global convergence for arbitrary functions with inexact 1-D search.

In general "classical" globalization techniques (restart or trust region) are not adopted because no examples of non convergence are known.

Widely used: quasi-Newton methods with BFGS and DFP updates and 1-D search procedures satisfying Wolfe conditions.

# Chapter 5: Constrained nonlinear optimization

Edoardo Amaldi

DEIB – Politecnico di Milano  
edoardo.amaldi@polimi.it

Course material on WeBeep 2022-23 - Optimization



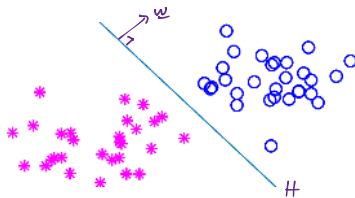
Academic year 2023-24

## 5.1 Example: Design linear classifiers and train SVMs

Support Vector Machines (SVMs) for binary classification.

Training set  $T = \{(\underline{x}^i, y^i) : \underline{x}^i \in \mathbb{R}^n, y^i \in \{-1, 1\}, i = 1, \dots, p\}$ .

Linear classifier: Suppose  $T$  is linearly separable



hyperplane  $H(\underline{w}, b) = \{\underline{x} \in \mathbb{R}^n : \underline{w}^T \underline{x} = b\}$  separates the points of the two classes if

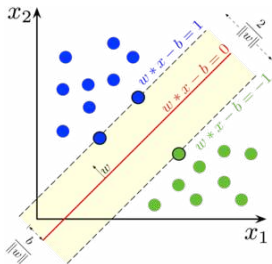
$$\begin{aligned} \underline{w}^T \underline{x}_i - b &\geq 1 & \text{for } i: \quad y_i = +1 \\ \underline{w}^T \underline{x}_i - b &\leq -1 & \text{for } i: \quad y_i = -1 \end{aligned}$$

$H$  not unique.

*the "best" one would be the one with the largest margin*



If  $T$  is linearly separable,  $H$  with **largest margin** (min distance from  $H$  to any  $x^i$ ) is the most robust w.r.t. noise.



Since width =  $\frac{2}{\|w\|}$ , hard-margin linear SVM training:

*we assume there is a linear separation*

*and we look for a linear separation*

$$\max \text{width} = \max \frac{2}{\|w\|_2} = \min \frac{\|w\|_2}{2}$$

$$\text{st } \sum_i (w^T x_i - b) - 1 \geq 0 \quad \forall i \quad (\text{we all points are correctly classified})$$

$$\begin{matrix} w \in \mathbb{R}^m \\ b \in \mathbb{R} \end{matrix}$$

*quadratic obj function linear constraints*

Remark:  $H$  with maximum margin is completely determined by the *support vectors* (closest  $x^i$ 's to  $H$ ).

*(as we also reduce the # of the constraints, never for now were the we are constraining each point)*

Decision function:  $h(\underline{w}, b, \underline{x}) = \underline{w}^t \underline{x} - b$ .

Extensions:

- 1) Soft margin for nonlinearly separable  $T$  (not convex)
- 2) Nonlinear classifiers by applying kernels.

See Computer Lab 5.

For other applications see e.g. Chap. 6-8 of S. Boyd and L. Vandenberghe, *Convex Optimization*, Cambridge Press, 2004.

## 5.2 Necessary optimality conditions

*generic problem structure:*

Consider

$$\begin{array}{ll} \min & f(x) \\ \text{s.t.} & g_i(\underline{x}) \leq \underline{0} \\ & \underline{x} \in \mathbb{R}^n \end{array} \quad \begin{array}{l} i \in I = \{1, \dots, m\} \\ | \\ \text{also since we can always} \\ \text{change us in the form} \end{array} \quad (1)$$

where  $f, g_i \in C^1$ .

Assumption: Feasible region  $S = \{ \underline{x} \in \mathbb{R}^n : g_i(\underline{x}) \leq 0, \forall i \in I \} \neq \emptyset$  but its interior can be empty.

Definitions: For each  $\bar{x} \in S$

*even we a feasible point*

- $\mathcal{D}(\bar{x}) = \{ \underline{d} \in \mathbb{R}^n : \exists \bar{\alpha} > 0 \text{ such that } \bar{x} + \alpha \underline{d} \in S, \forall \alpha \in [0, \bar{\alpha}] \}$

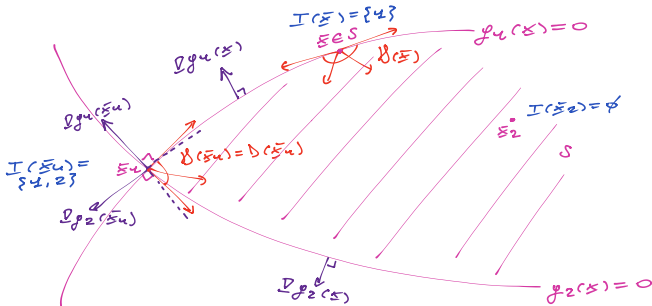
cone of the feasible directions.

- $I(\bar{x}) = \{ i \in I : g_i(\bar{x}) = 0 \} \subseteq I$  set of indices of the active constraints.

*ie the constraints which are activated at equilibrium*

- $D(\bar{x}) = \{ \underline{d} \in \mathbb{R}^n : \nabla^t g_i(\bar{x}) \underline{d} \leq 0, \forall i \in I(\bar{x}) \}$

cone of the directions constrained by the gradients of the active constraints.



Definitions: For each  $\bar{x} \in S$

*even we a feasible point*

- $\mathcal{D}(\bar{x}) = \{ \underline{d} \in \mathbb{R}^n : \exists \bar{\alpha} > 0 \text{ such that } \bar{x} + \alpha \underline{d} \in S, \forall \alpha \in [0, \bar{\alpha}] \}$

cone of the feasible directions.

- $I(\bar{x}) = \{ i \in I : g_i(\bar{x}) = 0 \} \subseteq I$  set of indices of the active constraints.

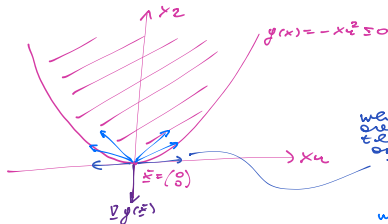
*ie the constraints which are not fixed or redundant*

- $D(\bar{x}) = \{ \underline{d} \in \mathbb{R}^n : \nabla^t g_i(\bar{x}) \underline{d} \leq 0, \forall i \in I(\bar{x}) \}$

cone of the directions constrained by the gradients of the active constraints.

*(another example)*

*because those directions*



*while these directions are good for  $D(\bar{x})$  since their projection is 0 on the direction of the gradient, we  $\nabla g(\bar{x}) = (-1, -1) = 0$*

*while the above ones are even closer  $\nabla g(\bar{x}) = (-1, -1) = 0$*

$S(\bar{x})$ : *shaded region with diagonal lines*

$D(\bar{x})$ : *shaded region with diagonal lines*

*x2-axis excluded*

*x1-axis included*

**Property:**  $\overline{\mathcal{D}(\bar{x})} \subseteq D(\bar{x})$  for all  $\bar{x} \in S$ .

Proof:

Given any  $\underline{d} \in \mathcal{D}(\bar{x})$ , for sufficiently small  $\alpha$  we have

$$0 = f_i(\bar{x} + \alpha \underline{d}) = \cancel{f_i(\bar{x})} + \alpha \nabla f_i(\bar{x})^T \underline{d} + o(\alpha)$$

$$\Rightarrow \underbrace{\nabla f_i(\bar{x})^T \underline{d}}_{\alpha > 0} \leq 0$$

$$\forall i \in I(\bar{x}) \\ \text{with } \alpha f_i(\bar{x}) = 0$$

$$\Rightarrow \nabla f_i(\bar{x})^T \underline{d} \leq 0 \quad \forall i \in I(\bar{x})$$

that is  $\underline{d} \in \mathcal{K}(\bar{x})$  and hence  $\mathcal{K}(\bar{x}) \subseteq D(\bar{x})$ .

But  $\mathcal{K}(\bar{x})$  is a closed set, then we have also that  $\overline{\mathcal{K}(\bar{x})} \subseteq D(\bar{x})$

Not all  $\underline{d} \in D(\bar{x})$  are feasible directions. *as we saw in the example of  $-x_1^2$*



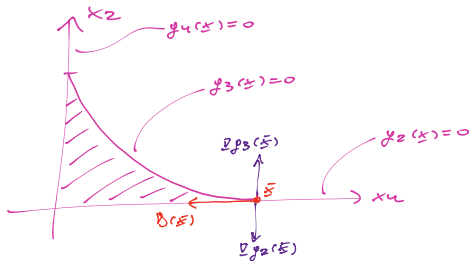
Example:

$$g_1(\underline{x}) = -x_1 \leq 0$$

$$g_2(\underline{x}) = -x_2 \leq 0$$

$$g_3(\underline{x}) = -(1-x_1)^3 + x_2 \leq 0$$

$\mathcal{D}(\bar{x})$ :   
 $D(\bar{x})$ : 



At  $\bar{x} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$  we have  $\mathcal{D}(\bar{x}) = \{(\alpha, 0) : \alpha < 0\}$

Since  $\nabla g_2(\bar{x}) = \begin{pmatrix} 0 \\ -1 \end{pmatrix}$  and  $\nabla g_3(\bar{x}) = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$  we have  
 $\mathcal{D}(\bar{x}) \subset D(\bar{x}) = \left\{ \begin{pmatrix} \alpha \\ \beta \end{pmatrix} : \alpha \in \mathbb{R} \right\}$

**Theorem:** (Extension of first order necessary optimality conditions)

If  $f \in C^1$  on  $S$  and  $\bar{x} \in S$  is a local minimum of  $f$  on  $S$ , then

$$\nabla f^t(\bar{x})\underline{d} \geq 0 \quad \forall \underline{d} \in \overline{\mathcal{D}}(\bar{x}),$$

that is, all feasible directions are ascent directions.

*where now we characterize more precisely the "feasible directions" in terms of  $\mathcal{D}(f)$*

Proof:

The result holds  $\forall \underline{d} \in \overline{\mathcal{D}}(\bar{x})$ .

For every  $\underline{d} \in \overline{\mathcal{D}}(\bar{x})$ ,  $\exists$  a sequence  $\{\underline{d}^k\}$  with  $\underline{d}^k \in \mathcal{D}(\bar{x})$  such that  $\lim_{k \rightarrow \infty} \underline{d}^k = \underline{d}$ .

Since  $\nabla f^t(\bar{x})\underline{d}^k \geq 0, \forall k$ , then  $\lim_{k \rightarrow \infty} \nabla f^t(\bar{x})\underline{d}^k = \nabla f^t(\bar{x})\underline{d} \geq 0$ .

But  $\overline{\mathcal{D}}(\bar{x})$  is difficult to characterize.

*where  $\mathcal{D}(f)$  may be easy to characterize, no*



Since  $D(\bar{x})$  is well characterized, we introduce further conditions.

**Definition:** (Constraint Qualification CQ – Zangwill)

The constraint qualification assumption holds at  $\bar{x} \in S$  if  $\overline{\mathcal{D}}(\bar{x}) = D(\bar{x})$

## Theorem: (Karush-Kuhn-Tucker necessary optimality conditions)

Suppose  $f, g_i \in C^1$  and CQ assumption holds at  $\bar{x} \in \{\underline{x} \in \mathbb{R}^n : g_i(\underline{x}) \leq 0, \forall i \in I\}$ . *feasible region*

If  $\bar{x}$  is a local minimum of  $f$  over  $S$  then  $\exists u_1, \dots, u_m \geq 0$  (KKT-multipliers) such that:

$$\nabla f(\bar{x}) + \sum_{i \in I(\bar{x})} u_i \nabla g_i(\bar{x}) = \mathbf{0} \equiv \begin{cases} \nabla f(\bar{x}) + \sum_{i=1}^m u_i \nabla g_i(\bar{x}) = \mathbf{0} \\ u_i g_i(\bar{x}) = 0 \quad \forall i \in I \end{cases}$$

$$\Leftrightarrow \nabla f(\bar{x}) = \sum_{i \in I(\bar{x})} (-u_i) \nabla g_i(\bar{x})$$

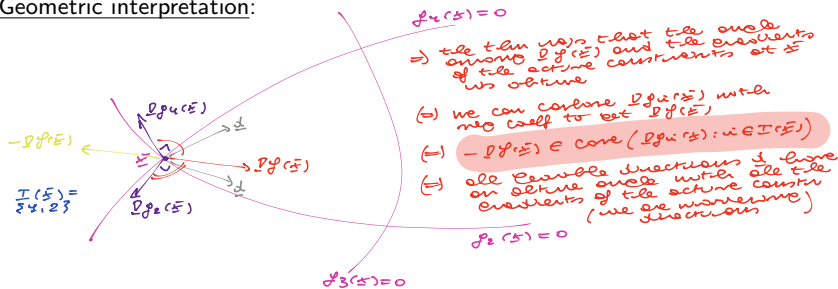
$\bar{x}$  must also satisfy all the constraints  $g_i(\underline{x}) \leq 0, \forall i \in I$ .

*ie  $\bar{x}$  must be in feasible region*

*a rewriting where all the indices  $i \in I$  (all the + that new constraint)*

*not linear constraints - we're taking uni and  $\bar{x}$  are unknowns -  $\rightarrow$  we will study cones*

### Geometric interpretation:



**Theorem:** (Karush-Kuhn-Tucker necessary optimality conditions)

Suppose  $f, g_i \in C^1$  and CQ assumption holds at  $\bar{x} \in \{x \in \mathbb{R}^n : g_i(x) \leq 0, \forall i \in I\}$ .

If  $\bar{x}$  is a local minimum of  $f$  over  $S$  then  $\exists u_1, \dots, u_m \geq 0$  (KKT-multipliers) such that:

$$\nabla f(\bar{x}) + \sum_{i \in I(\bar{x})} u_i \nabla g_i(\bar{x}) = \underline{0} \equiv \begin{cases} \nabla f(\bar{x}) + \sum_{i=1}^m u_i \nabla g_i(\bar{x}) = \underline{0} \\ u_i g_i(\bar{x}) = 0 \quad \forall i \in I \end{cases}$$

$\bar{x}$  must also satisfy all the constraints  $g_i(x) \leq 0, \forall i \in I$ .

Geometric interpretation:

$u_i > 0 \Rightarrow$  the corresponding constraint must be satisfied at equality

to check or plus we will develop the case

Proof:

Assuming CQ holds at  $\bar{x}$ , we have  $\overline{\mathcal{D}}(\bar{x}) = D(\bar{x})$ .

*rec cond*

NC for  $\bar{x}$  to be a local minimum of  $f$  over  $S$  is

*that we can  
 $\exists \underline{d} \in D(\bar{x})$*

*(thanks to the  
CQ assumption)*

$$\nabla^t f(\bar{x}) \underline{d} \geq 0, \quad \forall \underline{d} \text{ such that } \nabla^t g_i(\bar{x}) \underline{d} \leq 0 \quad \forall i \in I(\bar{x}). \quad (2)$$

Farkas Lemma:

$$\left\{ \begin{array}{l} A\underline{u} = \underline{b} \\ \underline{u} \geq 0 \end{array} \right\} \text{ has a solution} \Leftrightarrow \left\{ \begin{array}{l} \underline{b}^t \underline{d} \geq 0 \\ \forall \underline{d} \text{ such that } \underline{d}^t A \geq 0 \end{array} \right.$$

*thinking about the structure of the rhs  
before we can take*

$$\underbrace{\underline{b}}_{\underline{D}f(\bar{x})} + \sum_{i \in I(\bar{x})} \underbrace{\mu_i \underline{D}g_i(\bar{x})}_{A} = 0$$

$$A = \left( \begin{array}{ccc} | & & | \\ \underline{D}g_{i_1}(\bar{x}) & \dots & \underline{D}g_{i_\ell}(\bar{x}) \\ | & & | \end{array} \right)$$

*we can also  
- m rows  
- |I|=l cols*

Then (2) is equivalent to  $\underline{b}^t \underline{d} \geq 0 \quad \forall \underline{d} \in \mathbb{R}^n : \underline{d}^t A \geq 0$  (2\*)

But according to Farkas Lemma we have

$$(2^*) \Leftrightarrow \exists \underline{u} \geq 0 : A\underline{u} = \underline{b}$$

$$\Leftrightarrow \exists \mu_i \geq 0 \forall i \in I(\bar{x}) : \underline{D}f(\bar{x}) = \sum_{i \in I(\bar{x})} (-\mu_i) \underline{D}g_i(\bar{x})$$

### Example 1:

$$\begin{aligned} \min \quad & f(\underline{x}) = x_1 + x_2 \\ \text{s.t.} \quad & g_1(\underline{x}) = x_1^2 + x_2^2 \leq 2 \\ & g_2(\underline{x}) = -x_2 \leq 0 \end{aligned}$$

*actually*  
 $g_4(\underline{x}) = x_1^2 + x_2^2 - 2 \geq 0$   
 so we must have  
 zero on the RHS

KKT conditions: the candidate points  $\bar{x}$  must satisfy the following conditions:

$$\nabla f(\bar{x}) + \sum_{i \in I(\bar{x})} (-\mu_i) \nabla g_i(\bar{x}) = 0$$

*well we use the other convention (lower cost) exercise*

$$\begin{pmatrix} 1 \\ 1 \end{pmatrix} + \left[ \mu_1 \begin{pmatrix} 2\bar{x}_1 \\ 2\bar{x}_2 \end{pmatrix} + \mu_2 \begin{pmatrix} 0 \\ -1 \end{pmatrix} \right] = 0$$

$$\mu_1 g_1(\bar{x}) = \mu_1 (\bar{x}_1^2 + \bar{x}_2^2 - 2) = 0$$

$$\mu_2 g_2(\bar{x}) = \mu_2 (-\bar{x}_2) = 0$$

$$\mu_1 \geq 0$$

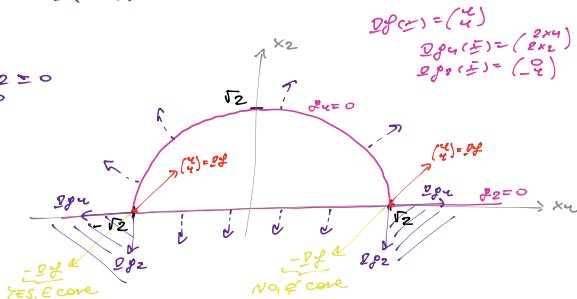
$$\mu_2 \geq 0$$

$$\bar{x}_1^2 + \bar{x}_2^2 - 2 \geq 0$$

$$-\bar{x}_2 \geq 0$$

KKT cond

feasibility cond of  $\bar{x}$



$$\begin{pmatrix} 1 \\ 1 \end{pmatrix} + u_1 \begin{pmatrix} 2x_1 \\ 2x_2 \end{pmatrix} + u_2 \begin{pmatrix} 0 \\ -1 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$$u_1(x_1^2 + x_2^2 - 2) = 0$$

$$u_2(-x_2) = 0$$

$$x_1^2 + x_2^2 \leq 2$$

$$-x_2 \leq 0$$

$$u_1 \geq 0, u_2 \geq 0$$

Four cases:

(C4)  $\begin{matrix} u_1 = 0 \\ u_2 = 0 \end{matrix} \Rightarrow \begin{pmatrix} 1 \\ 1 \end{pmatrix} + \sum 0 \cdot \parallel = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \neq \begin{pmatrix} 0 \\ 0 \end{pmatrix}$  impossible

(C2)  $\begin{matrix} u_1 = 0 \\ u_2 > 0 \end{matrix} \Rightarrow \begin{pmatrix} 1 \\ 1 \end{pmatrix} + 0 \cdot \parallel + u_2 \begin{pmatrix} 0 \\ -1 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$  impossible

(C3)  $\begin{matrix} u_1 > 0 \\ u_2 = 0 \end{matrix} \Rightarrow \begin{pmatrix} 1 \\ 1 \end{pmatrix} + u_1 \begin{pmatrix} 2x_1 \\ 2x_2 \end{pmatrix} + 0 \cdot \parallel = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$

$$1 + 2u_1 x_2 = 0 \quad \text{impossible}$$

$$u_1 = -\frac{1}{2x_2} \leq 0$$

↑  
this is  $\leq 0$   
i.e. contradiction

we lie on  
the circle

(C6)  $\begin{matrix} u_1 > 0 \\ u_2 > 0 \\ x_2 = 0 \end{matrix} \Rightarrow \begin{cases} x_1^2 + x_2^2 - 2 = 0 \\ x_2 = 0 \end{cases} \Rightarrow$  two sets:  
 $A = \begin{pmatrix} -\sqrt{2} \\ 0 \end{pmatrix}, B = \begin{pmatrix} \sqrt{2} \\ 0 \end{pmatrix}$

We now check to have  $u_i \geq 0$  w.r. to  
the real and complex parts  
we verify that constraints are satisfied

$$\begin{pmatrix} 1 \\ 1 \end{pmatrix} + u_1 \begin{pmatrix} 2x_1 \\ 2x_2 \end{pmatrix} + u_2 \begin{pmatrix} 0 \\ -1 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$$u_1(x_1^2 + x_2^2 - 2) = 0$$

$$u_2(-x_2) = 0$$

$$x_1^2 + x_2^2 \leq 2$$

$$-x_2 \leq 0$$

$$u_1 \geq 0, u_2 \geq 0$$

Four cases:

*we lie on the circle*

(CG)  $\begin{cases} u_1 > 0 \\ u_2 > 0 \\ x_2 = 0 \end{cases} \Rightarrow \begin{cases} x_1^2 + x_2^2 - 2 = 0 \\ x_2 = 0 \end{cases} \Rightarrow$  two sets:  $B = \begin{pmatrix} -\sqrt{2} \\ 0 \end{pmatrix}, A = \begin{pmatrix} \sqrt{2} \\ 0 \end{pmatrix}$

*We now check to have  $u_i \geq 0$  w.r. to the real) are candidate points*

*we verify KKT cond. (eval.)*

$$A: \begin{pmatrix} 4 \\ 4 \end{pmatrix} + u_1 \begin{pmatrix} 2\sqrt{2} \\ 0 \end{pmatrix} + u_2 \begin{pmatrix} 0 \\ -4 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \Rightarrow \begin{cases} u_2 = 4 \text{ good} \\ u_1 = -\frac{4}{2\sqrt{2}} < 0 \text{ bad} \end{cases}$$

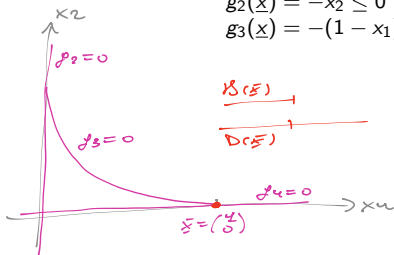
$$B: \begin{pmatrix} 4 \\ 4 \end{pmatrix} + u_1 \begin{pmatrix} -2\sqrt{2} \\ 0 \end{pmatrix} + u_2 \begin{pmatrix} 0 \\ -4 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \Rightarrow \begin{cases} u_2 = 4 \text{ good} \\ u_1 = \frac{4}{2\sqrt{2}} > 0 \text{ good} \end{cases}$$

$\Rightarrow$  only B was a candidate point (and indeed w.r. to the local & abs minimum)

If CQ assumption does not hold at  $\bar{x}$ , KKT conditions need not be necessary for local optimality.

Example 2:

$$\begin{aligned} \min \quad & f(x) = -x_1 \\ \text{s.t.} \quad & g_1(x) = -x_1 \leq 0 \\ & g_2(x) = -x_2 \leq 0 \\ & g_3(x) = -(1-x_1)^3 + x_2 \leq 0 \end{aligned}$$



We know that  $Df(\bar{x}) \supset \bar{J}(\bar{x})$  for  $\bar{x} = (\frac{1}{2}, 0)$

now  $\bar{x}$  is still the min for  $f(x)$ , but no multipliers exist to meet KKT conditions.

Since  $(u_1, u_2) = (0, 0)$  we have

$$(-\frac{1}{2}) + u_1(-\frac{1}{2}) + u_3(\frac{3}{4}) \neq (0) \quad \forall u_2, u_3 \geq 0$$

since constraint 2 is active not



*how do we check to have CQ assumption?*

**Proposition:** (Sufficient conditions for Constraint Qualification)

1) If

- all  $g_i$  are linear functions (Karlin)

or

- all  $g_i$  are convex and  $\exists \underline{a}$  such that  $g_i(\underline{a}) < 0, \forall i \in I$ , (Slater)  
*we  $\rightarrow$  an interior feasible point*

CQ assumption holds at every  $\underline{x} \in S$ .

2) If  $\nabla g_i(\bar{x}), i \in I(\bar{x})$ , are linearly independent, CQ assumption holds at  $\bar{x} \in S$ .

*the gradients of the active constraints*

N.B.: When the gradients of the active constraints are linearly independent, KKT multiplier vector is unique.

**Theorem:** (Necessary and sufficient conditions – convex problems)

If  $f \in C^1$ ,  $g_i \in C^1 \forall i \in I$  are convex, and  $\exists a$  such that  $g_i(a) < 0, \forall i \in I$ , then

$x^* \in S$  is a global minimum if and only if  $\exists u_1, \dots, u_m \geq 0$  such that

$$\begin{cases} \nabla f(x^*) + \sum_{i=1}^m u_i \nabla g_i(x^*) = 0 \\ u_i g_i(x^*) = 0 \quad \forall i \in I. \end{cases}$$

*noting the KKT conditions*  
 + (convexity) & ( $u_i \geq 0$ )

Proof:

( $\Rightarrow$ ) If  $x^* \in S$  is a (loc) min, then  $\exists \nabla$  necessary KKT conditions through Slater condition.

( $\Leftarrow$ ) If we have the KKT cond then  $Df(x^*)^T \exists \lambda \geq 0 \quad \forall \lambda \in \mathcal{N}(x^*) = \mathcal{N}(S^*)$  (from CQ and Farkas lemma)  
 Therefore  $x^* \in S$  is the (global) min

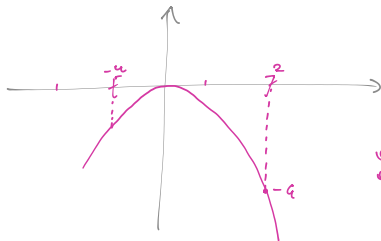
For Linear Programs, it amounts to the complementary slackness theorem.

Remark: Result holds under milder convexity conditions ( $f$  pseudoconvex and the  $g_i$ 's quasiconvex).

If  $f$  is not convex, KKT conditions are not sufficient.

Example 3:

$$\begin{aligned} \min f(x) &= -x^2 \\ g_1(x) &= -2 + x \leq 0 \\ g_2(x) &= -x - 1 \leq 0 \end{aligned} \quad \left. \vphantom{\begin{aligned} \min f(x) &= -x^2 \\ g_1(x) &= -2 + x \leq 0 \\ g_2(x) &= -x - 1 \leq 0 \end{aligned}} \right\} x \in [-4, 2]$$



At  $x = -4$  the KKT cond are satisfied but the candidate point is a local max and not a global min

In convexity is really important to get well cond. for global optimum!

# General case *with slow equality constraints*

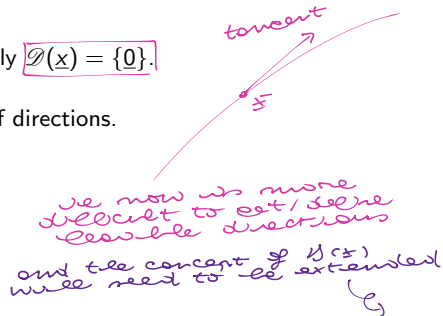
Consider

$$\begin{aligned} \min \quad & f(x) \\ \text{s.t.} \quad & g_i(x) \leq 0 \quad i \in I = \{1, \dots, m\} \\ & h_l(x) = 0 \quad l \in L = \{1, \dots, p\} \\ & x \in X \subseteq \mathbb{R}^n \end{aligned}$$

where  $f, g_i, h_l \in C^1$ .

When nonlinear equality constraints, usually  $\mathcal{D}(x) = \{0\}$ .

Extend previous results by defining cone of directions.



**Definition:** Closed cone of the tangents at  $\bar{x}$

$$\mathcal{T}(\bar{x}) = \left\{ \underline{d} \in \mathbb{R}^n : \underline{d} = \lambda \lim_{k \rightarrow \infty} \frac{\underline{x}^k - \bar{x}}{\|\underline{x}^k - \bar{x}\|}, \lambda \geq 0, \left\{ \underline{x}^k \right\}_{k \rightarrow \infty} \rightarrow \bar{x} \text{ with } \underline{x}^k \neq \bar{x} \right\}$$

with  $\exists \epsilon \exists CS$  are possible sets

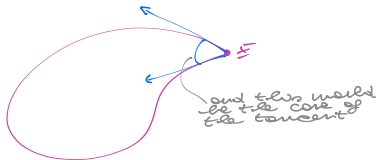
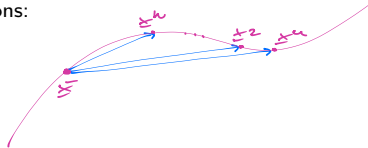
limit direction

(if  $\underline{d} \in \mathcal{B}(\bar{x})$  and  $\exists$  a sequence  $(\underline{x}^k) \rightarrow \bar{x}$ , then directions of the cone?  $\underline{x}^k - \bar{x}$  converge to  $\underline{d}$ )

if the corresponding clouds that we get with the sequence

we can also cone these through the interior

Illustrations:



**Definition:** (Constraint Qualification CQ - Abadie)

The CQ assumption holds at  $\bar{x} \in S$  if  $\mathcal{T}(\bar{x}) = D(\bar{x}) \cap H(\bar{x})$  where

$$\mathcal{B}(\bar{x}) = \left\{ \underline{d} \in \mathbb{R}^m : \nabla g_i(\bar{x}) \cdot \underline{d} > 0 \quad \forall i \in I(\bar{x}) \right\} \quad \left\{ \begin{array}{l} \text{for all the active} \\ \text{(w/eq)} \text{ constraints} \end{array} \right.$$

$$\mathcal{H}(\bar{x}) = \left\{ \underline{d} \in \mathbb{R}^m : \nabla g_i(\bar{x}) \cdot \underline{d} = 0 \quad \forall i \in L \right\} \quad \left\{ \begin{array}{l} \text{for all the equality constraints} \\ \text{(which become at equality)} \\ \text{or active (i.e. derivation)} \end{array} \right.$$

**Theorem:** (General KKT necessary optimality conditions)

Suppose  $f \in C^1$ ,  $g_i \in C^1 \forall i$ ,  $h_l \in C^1 \forall l$  and CQ assumption holds at  $\bar{x} \in S$ . *by the one of above*

If  $\bar{x}$  is a local minimum of  $f$  over  $S$  then  $\exists u_i \geq 0, \forall i \in I(\bar{x})$  and  $v_l \in \mathbb{R}, \forall l \in L$  such that

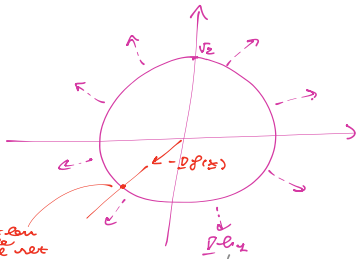
$$\nabla f(\bar{x}) + \sum_{i \in I(\bar{x})} u_i \nabla g_i(\bar{x}) + \sum_{l \in L} v_l \nabla h_l(\bar{x}) = \underline{0}.$$

N.B.: If only equalities, KKT conditions coincide with classical Lagrange optimality conditions.

$\Rightarrow$  we can express  $-\nabla f(\bar{x})$  as a linear combination of all the active constraints

**Example 1:**

$$\begin{aligned} \min \quad & f(\underline{x}) = x_1 + x_2 \\ \text{s.t.} \quad & x_1^2 + x_2^2 = 2 \end{aligned}$$



$$-\nabla f(\bar{x}) = \begin{pmatrix} -4 \\ -4 \end{pmatrix}$$

$$\Theta_4(\bar{x}) = x_1^2 + x_2^2 - 2$$

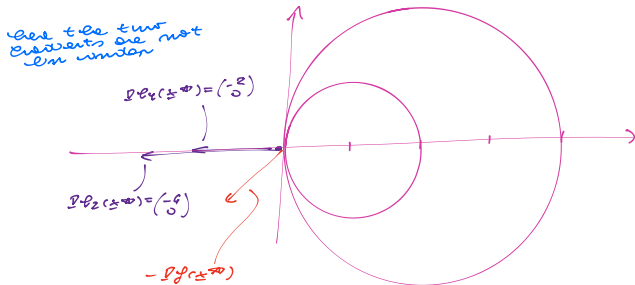
$$\nabla \Theta_4(\bar{x}) = \begin{pmatrix} 2x_1 \\ 2x_2 \end{pmatrix}$$

$\begin{pmatrix} -4 \\ -4 \end{pmatrix}$  is in the optimal set

can show we're on the tangent

## Example 2:

$$\begin{aligned} \min \quad & f(x) = x_1 + x_2 \\ \text{s.t.} \quad & (x_1 - 1)^2 + x_2^2 - 1 = 0 \quad R=4 \quad C=(4,0) \\ & (x_1 - 2)^2 + x_2^2 - 4 = 0 \quad R=2 \quad C=(2,0) \end{aligned}$$



we use that we can't express w/o  $\geq$  even comb. of the two constraints alone

but  $x^* = (0, 0)$  is the optimal set (same the only feasible point)

$\Rightarrow \exists r_1, r_2$  at the left hand side are not valid (since  $\geq$  here is an inequality point)

**Proposition:** (Sufficient conditions for CQ)

- If  $g_i$  convex,  $h_l$  linear and  $\exists \underline{a} \in X$  such that  $g_i(\underline{a}) < 0, \forall i \in I$  and  $h_l(\underline{a}) = 0 \forall l \in L$ , then CQ assumption holds at every  $\underline{x} \in S$ .
- If  $\nabla g_i(\bar{x}), \forall i \in I(\bar{x})$ , and  $\nabla h_l(\bar{x}), \forall l \in L$ , are linearly independent then CQ assumption holds at  $\bar{x} \in S$ .

*↪ extension  
considering all eq  
constraint & active  
constraint*



## 5.3 Sufficient optimality conditions

Generic NLP

$$(P) \quad \begin{cases} \min & f(\underline{x}) \\ \text{s.t.} & g_i(\underline{x}) \leq 0 \quad \forall i \in I = \{1, \dots, m\} \\ & \underline{x} \in X \subseteq \mathbb{R}^n \end{cases}$$

*not convex*

where  $X$  is an arbitrary subset (even discrete).

## Definitions

- The **Lagrange function** associated with (P) is

$$L(\underline{x}, \underline{u}) = f(\underline{x}) + \sum_{i \in I} u_i g_i(\underline{x}) \quad \forall \underline{x} \in X \text{ and } \underline{u} \geq \underline{0}$$

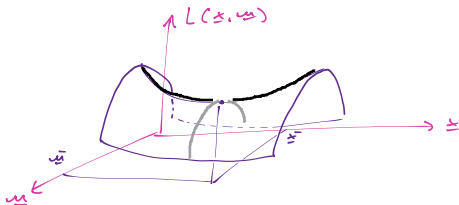
N.B.:  $\underline{u} \geq \underline{0}$  since  $g_i(\underline{x}) \leq 0$ .

- $(\bar{\underline{x}}, \bar{\underline{u}})$  with  $\bar{\underline{x}} \in X$  and  $\bar{\underline{u}} \geq \underline{0}$  is a **saddle point** of  $L(\underline{x}, \underline{u})$

$$\text{if } L(\bar{\underline{x}}, \bar{\underline{u}}) \leq L(\underline{x}, \bar{\underline{u}}) \quad \forall \underline{x} \in X \quad \text{and} \quad L(\bar{\underline{x}}, \bar{\underline{u}}) \geq L(\bar{\underline{x}}, \underline{u}) \quad \forall \underline{u} \geq \underline{0},$$

that is,  $\bar{\underline{x}}$  minimizes  $L(\underline{x}, \bar{\underline{u}})$  over  $X$  and  $\bar{\underline{u}}$  maximizes  $L(\bar{\underline{x}}, \underline{u})$  over  $\mathbb{R}^m$ .

Illustration:



*→ (we had) w/ constr  
is not satisfied  
so (we avoid) w/ constr  
is satisfied*

*net constr  
remains free*

*$\bar{x}$  is the min  
at  $u$  fixed*

*$\bar{u}$  is the max  
at  $x$  fixed*

**Proposition:** (Characterization of saddle points)

$(\bar{x}, \bar{u})$  with  $\bar{x} \in X$  and  $\bar{u} \geq 0$  is a saddle point of  $L(x, u)$  if and only if

i)  $L(\bar{x}, \bar{u}) = \min_{x \in X} L(x, \bar{u})$

*this is trivial, was part of the definition*

ii)  $g_i(\bar{x}) \leq 0 \quad \forall i \in I$

*$\bar{x}$  is a feasible pt (since we only w/ to network) all the constraints are satisfied*

iii)  $\bar{u}_i g_i(\bar{x}) = 0 \quad \forall i \in I$

*complementary conditions:  
(if one of them is  $> 0$ )  $\Rightarrow$  (the other has to be  $= 0$ )*

Proof\*:

*Consequence of the definition of the Lagrangian function and a lot of relations with duality*



*a new strong result*  
**Theorem:** (Sufficient optimality condition)

If  $(\bar{x}, \bar{u})$  is a saddle point of  $L(x, u)$ , then  $\bar{x}$  is a global minimum of problem (P).

**Proof:** *condition (u) of the characterization*  $\Rightarrow L(\bar{x}, \bar{u}) \leq L(x, \bar{u}) \quad \forall x \in X$

*2nd order condition of the characterization*  $\Rightarrow f(\bar{x}) + \sum_{u \in I} \bar{u}_i g_i(\bar{x}) \leq f(x) + \sum_{u \in I} \bar{u}_i g_i(x) \quad \forall x \in X$   
*this is = 0 by character (3)*

$\Rightarrow f(\bar{x}) \leq f(x) + \sum_{u \in I} \bar{u}_i g_i(x) \quad \forall x \in X$   
*this will be 0 by character (2)* } *and in this case we make the other side of (P)*

*that's because since  $\bar{u} \geq 0$  we have*  
 $f(\bar{x}) \leq f(x) \quad \forall x \in X \text{ at } g_i(x) \leq 0 \quad \forall i \in I$   
*we conclude for the problem (P)*

**Observations:**

- Result applies to any mathematical program (convex or not, with  $f$  and  $g_i$  differentiable or not,  $X$  continuous or discrete, ...).
- For some problems a saddle point may not exist, in general for nonconvex problems.

Example:

$$\begin{aligned} \min \quad & f(x) = -x^2 \\ \text{s.t.} \quad & 2x - 1 \leq 0 \\ & 0 \leq x \leq 1 \end{aligned}$$

where  $g(x) = 2x - 1$  and  $X = \{x : 0 \leq x \leq 1\}$

$$L(x, u) = f(x) + \sum_{i \in I} u_i g_i(x) = \begin{cases} |I|=1 & (-x^2) + u(2x-1) \end{cases}$$

which is concave w.r.t  $x \Rightarrow \max_{x \in [0, 1]} L(x, u)$   
(for a fixed  $u \geq 0$ )  
is reached at an extreme point, i.e.  $x=0$  or  $x=1$

while we can leave, and  $x^* = \frac{u}{2}$  but  $\nexists$  saddle point of the type  $(\frac{u}{2}, u)$

$\Rightarrow$  we indeed have only a  
SUF cond, not NEC, but:

**Theorem:** (saddle point for convex problems)

Suppose  $f$  and  $g_i, \forall i \in I$  are convex,  $X \subseteq \mathbb{R}^n$  is convex and  $\exists \underline{a} \in X$  such that  $g(\underline{a}) < 0$ .  
*For an interior saddle point*

If (P) has an optimal solution  $\bar{x}, \exists \bar{u} \geq 0$  such that  $(\bar{x}, \bar{u})$  is a saddle point of  $L(x, u)$ .

we need NEC & SUF  
for convex problems  
(else we'd have a concave cone)

## Connection with KKT conditions for convex problems

If  $f$  and  $g_i \in C^1$  are convex,  $X = \mathbb{R}^n$  and  $\exists \underline{a} \in X$  such that  $g(\underline{a}) < 0$ , then  $\bar{x}$  is an optimal solution if and only if  $\bar{x}$  satisfies the KKT conditions.

Proof:  $\bar{x}$  is an opt. pt  $\Leftrightarrow \exists \bar{u} \geq 0 : (\bar{x}, \bar{u})$  is a saddle point for  $L(\bar{x}, \bar{u})$

( $\Leftarrow$ ) It is the SFC condition for global optimality  
( $\Rightarrow$ ) It is the previous theorem (about the always existence of a saddle point)

$$\text{Since } L(\bar{x}, \bar{u}) = f(\bar{x}) + \sum_{i \in I} \bar{u}_i g_i(\bar{x}) = f(\bar{x}) + \bar{u} \cdot \bar{g}(\bar{x})$$

is convex, then we can look (due to condition (4)) for the minimum point (i.e. improving  $\rightarrow$  is a stationary point)

$$\nabla_{\bar{x}} L(\bar{x}, \bar{u}) = 0$$

which coincides with the KKT conditions

$$\nabla f(\bar{x}) + \sum_{i \in I} \bar{u}_i \nabla g_i(\bar{x}) = 0$$

and the CV assumption holds at every feasible pt

- N.B.: 1) Without convexity assumption a stationary point  $\bar{x}$  may not minimize  $L(\bar{x}, \bar{u})$ .  
2) KKT multipliers are then identical to Lagrange multipliers at the saddle point.

## 5.4 Lagrangian duality

Generic NLP:

$$(P) \quad \begin{cases} \min & f(\underline{x}) \\ \text{s.t.} & g_i(\underline{x}) \leq 0 \quad \forall i \in I = \{1, \dots, m\} \\ & \underline{x} \in X \subseteq \mathbb{R}^n \end{cases}$$

To any minimization NLP we can associate a maximization NLP such that, under some assumptions, the objective function values of respective optimal solutions coincide.

Tackle the primal problem ( $P$ ) indirectly, by solving the dual (second) problem.

To try to solve ( $P$ ), we can look for a saddle point of the Lagrange function.

*this can be  
convenient in  
some cases*



## Dual function:

$$w(\underline{u}) = \min_{\underline{x} \in X} L(\underline{x}, \underline{u}) \quad \forall \underline{u} \geq \underline{0}$$

For every  $\underline{u} \geq \underline{0}$  we fix, we set  $w(\underline{u}) = \min_{\underline{x} \in X} L(\underline{x}, \underline{u})$  problem over the  $\underline{x} \in X$

Well-defined if, for instance,  $f$  and the  $g_i$ 's are continuous and  $X$  is compact.

Search for a saddle point (if  $\exists$ ):

$\Rightarrow$  objective function will look for a saddle point:

Dual problem:  $(D) \begin{cases} \max w(\underline{u}) \\ \underline{u} \geq \underline{0} \end{cases}$

$$\max_{\underline{u} \geq \underline{0}} \left( \min_{\underline{x} \in X} L(\underline{x}, \underline{u}) \right)$$

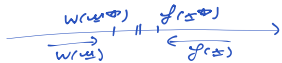
N.B.:  $w(\underline{u})$  and  $(D)$  are defined even if no saddle point exists.

and we relax since we have removed the algebraic constraints. But as of course the function  $w(\underline{u})$  may not have good properties, as a trade-off

Observations:

- 1) Different Lagrangian duals of  $(P)$  depending on which  $g_i(\underline{x}) < 0$  are dualized. Choice affects optimal value of  $(D)$  and complexity to evaluate  $w(\underline{u})$ . we could decide to hear out none of the constraints
- 2) Lagrangian dual is useful to solve large-scale LPs and (non)convex/discrete optimization problems.

**Theorem:** (Weak duality)



For every feasible  $\underline{x}$  of (P) and feasible  $\underline{u} \geq \underline{0}$  of (D), we have  $w(\underline{u}) \leq f(\underline{x})$ .

Proof: B) definition of Lagrangian function:  
 $w(\underline{u}) = \min_{\underline{x} \in X} (f(\underline{x}) + \underline{u}^T g(\underline{x})) \leq f(\underline{x}) + \underline{u}^T g(\underline{x}) \quad \forall \underline{x} \in X, \underline{u} \geq \underline{0}$   
*(no  $\underline{x} \in X$  is feasible)*

So wif  $g(\underline{x}) \leq \underline{0}$  then we have  $w(\underline{u}) \leq f(\underline{x})$ .  
*(due to  $\underline{u} \geq \underline{0}$ )*

In particular, for every  $\underline{u} \geq \underline{0}$  we have  $w(\underline{u}) \leq f(\underline{x}^*)$  for an optimal  $\underline{x}^*$  of (P).

$\Rightarrow$  the term gives a lower bound on the least opt val of function ( $f(\underline{x}^*)$ ) value

**Consequence:**

If a feasible solution  $\underline{\bar{x}}$  of (P) and  $\underline{\bar{u}} \geq \underline{0}$  satisfy  $w(\underline{\bar{u}}) = f(\underline{\bar{x}})$ ,  $\underline{\bar{x}}$  is optimal for (P) and  $\underline{\bar{u}}$  is optimal for (D).

For Linear Programs the objective function values of optimal solutions of (P) and (D) coincide, for NLPs this is not always the case.

## Theorem: (Strong duality)

i) If  $(P)$  has a saddle point  $(\bar{x}, \bar{u})$ , then   
 - not only  $\bar{x}$  is optimal for  $(P)$    
 - but also  $\bar{u}$  is optimal for  $(D)$

$$\begin{cases} \max_{\underline{u} \geq 0} w(\underline{u}) \\ \end{cases} = \boxed{w(\bar{u}) = f(\bar{x})} = \min \{ f(\underline{x}) : \underline{g}(\underline{x}) \leq 0, \underline{x} \in X \}.$$

ii) If  $\exists$  a feasible  $\bar{x}$  of  $(P)$  and  $\bar{u} \geq 0$  such that  $w(\bar{u}) = f(\bar{x})$ , then  $(\bar{x}, \bar{u})$  is a saddle point of  $L(\underline{x}, \underline{u})$ .

and also the converse result holds:   
 $w(\bar{u}) = f(\bar{x}) \Rightarrow (\bar{x}, \bar{u})$  is a saddle point

Proof:

(i) Since  $(\bar{x}, \bar{u})$  is a saddle point we have that   
 $w(\bar{u}) = L(\bar{x}, \bar{u}) = \min_{\underline{x} \in X} L(\underline{x}, \bar{u})$  — Conv. case (4)

Moreover, by definition of Lagrangian function   
 $L(\bar{x}, \bar{u}) = f(\bar{x}) + \bar{u}^T g(\bar{x}) = f(\bar{x}) = \begin{cases} \min_{\underline{x} \in X} f(\underline{x}) \\ \text{st } g(\underline{x}) \leq 0 \end{cases}$    
 = 0 Conv. case (3)   
 SDC (weak opt. condition)

$$\Rightarrow w(\bar{u}) = f(\bar{x})$$

Because of weak duality we have  $w(\underline{u}) \leq f(\bar{x}) \forall \underline{u} \geq 0$

$$\text{where } w(\bar{u}) = \begin{cases} \max_{\underline{u} \geq 0} w(\underline{u}) \\ \end{cases}$$

(iii) Let  $\bar{x}$  be a feasible set of  $(P)$  and  $\bar{u} \geq 0$  st  $w(\bar{x}) = f(\bar{x})$   
 (ie the statement)

By definition of  $w(\cdot)$  we have:

$$w(\bar{x}) \geq f(\bar{x}) + \bar{u}^T g(\bar{x}) \quad \forall \bar{x} \in X$$

Then  $\bar{x} = \bar{x}$  we have

$$f(\bar{x}) = w(\bar{x}) \geq f(\bar{x}) + \bar{u}^T g(\bar{x})$$

$$\Rightarrow w(\bar{x}) - f(\bar{x}) = 0 \quad \square \quad \bar{u}^T g(\bar{x})$$

$$\Rightarrow \bar{u}_i g_i(\bar{x}) = 0 \quad \forall i \in I \quad \left\{ \begin{array}{l} \text{clear (3)} \\ \text{so } \bar{x} \text{ is feasible} \end{array} \right.$$

Since  $\bar{u} \geq 0$ ,  $g_i(\bar{x}) \leq 0 \quad \forall i \in I$  and  $w(\bar{x}) = L(\bar{x}, \bar{u}) = \min_{\delta \in X} L(\bar{x}, \bar{u})$   
 clear (2) clear (4)

$\Rightarrow$  we have all the characterization points and  $w(\bar{x}, \bar{u})$  is a saddle point

## Consequence:

If  $f, g_i$ 's and  $X \subseteq \mathbb{R}^n$  are convex,  $\exists \underline{a}$  such that  $\underline{g}(\underline{a}) < \underline{0}$  and  $(P)$  has a finite optimal solution,  $\exists \underline{a}$  saddle point  $(\bar{x}, \bar{u})$  and i) holds:

$$\left\{ \begin{array}{l} \max w(\underline{u}) \\ \underline{u} \geq \underline{0} \end{array} \right. = \min \{ f(\underline{x}) : \underline{g}(\underline{x}) \leq \underline{0}, \underline{x} \in X \}.$$

N.B.: *Strong duality*, the optimal values of the two objective functions coincide.

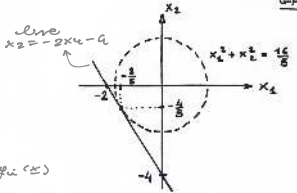
In general, we can have a duality gap ( $<$  instead of  $=$ ).

$\uparrow$   
 we can non-convex problems

$$1) \begin{cases} \text{MIN} & x_1^2 + x_2^2 \\ \text{s.t.} & 2x_1 + x_2 \leq -4 \end{cases}$$

$f$  &  $g$  CONVEX }  $\Rightarrow \exists$  SADDLE-PT  
 $\Xi = \begin{pmatrix} -\frac{4}{3} \\ \frac{8}{3} \end{pmatrix}$  }  $\Rightarrow$  NO DUALITY GAP

$g(x) = 2x_1 + x_2 - a \leq 0$



$x^* = \left(-\frac{4}{3}, \frac{8}{3}\right)$

$f(x^*) = \frac{16}{9}$

$f(x) + \sum_{i \in I} \mu_i g_i(x)$

$L(x, \mu) = x_1^2 + x_2^2 + 2\mu x_1 + \mu x_2 + 4\mu$  CONVEX

Case we minimize wrt  $x$

$$\begin{cases} \frac{\partial L}{\partial x_1} = 2x_1 + 2\mu = 0 \\ \frac{\partial L}{\partial x_2} = 2x_2 + \mu = 0 \end{cases}$$

$\Rightarrow \begin{cases} x_1 = -\mu \\ x_2 = -\frac{\mu}{2} \end{cases}$

we plug them into  $L(x, \mu)$

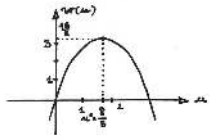
$\Rightarrow w(\mu) = -\frac{5\mu^2}{4} + 4\mu$   
 CONCAVE

(D) MAX  $w(\mu)$   
 $\mu \geq 0$

$\frac{\partial w(\mu)}{\partial \mu} = -\frac{5\mu}{2} + 4 \Rightarrow \mu^* = \frac{8}{5} \geq 0$

Then  $w(\mu^*) = -\frac{5}{4} \left(\frac{8}{5}\right)^2 + \frac{32}{5} = \frac{16}{5}$   
 $= f(x^*)$

NO DUALITY GAP



2) 
$$\begin{cases} \text{MIN} & -2x_1 + x_2 \\ \text{s.c.} & x_1 + x_2 - 3 = 0 \\ & (x_1, x_2) \in X = \left\{ \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 4 \end{pmatrix}, \begin{pmatrix} 4 \\ 0 \end{pmatrix}, \begin{pmatrix} 4 \\ 4 \end{pmatrix}, \begin{pmatrix} 1 \\ 2 \end{pmatrix}, \begin{pmatrix} 2 \\ 1 \end{pmatrix} \right\} \end{cases}$$

INTEGER PROGRAM

(P) :  $x^* = \begin{pmatrix} 2 \\ 1 \end{pmatrix}, f(x^*) = -3$

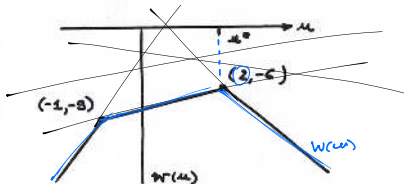
$L(x, \mu) = -2x_1 + x_2 + \mu(x_1 + x_2 - 3)$        $\mu$  UNRESTRICTED

$$W(\mu) = \min_{x \in X} L(x, \mu) = \begin{cases} -4 + 5\mu & \text{FOR } \mu \leq -1 \\ -8 + \mu & \text{FOR } -1 \leq \mu \leq 2 \\ -3\mu & \text{FOR } \mu \geq 2 \end{cases}$$

For each  $x \in X$  we get a different linear function with  $\mu$ , we can take the lower envelope

$\mu^* = 2, W(\mu^*) = -6 < f(x^*)$

$\exists$  DUALITY GAP



Since under certain conditions we can solve (P) indirectly by solving (D)

**Property 1:** The dual function  $w(\underline{u})$  is concave. *regardless of our assumption*

Proof\*:

*which is a good news since we are maximizing wt so in general solving the dual (D) is a lot easier than (P)*

**Observations:**

• If  $X \subseteq \mathbb{Z}^n$ ,  $w(\underline{u})$  is not everywhere continuously differentiable. Concave piecewise linear function, lower envelope of a (in)finite family of hyperplanes in  $\mathbb{R}^{n+1}$ .

• In general (D) is easier than (P).

• Since  $w(\underline{u})$  is concave local optima are global optima, but need for ad hoc solution method: subgradient method.

*test we solved with abs in the case of non-convex piecewise linear functions (P) vs (D) - lower envelope*

*the set of all the  $\underline{x}$ 's for which the min is reached, where  $w(\underline{u}) = \min_{\underline{x} \in X} L(\underline{x}, \underline{u})$*

**Property 2:** For  $\underline{u} \in \mathbb{R}_+^m$  let  $X(\underline{u}) = \{ \underline{x} \in X : f(\underline{x}) + \underline{u}^t \underline{g}(\underline{x}) = w(\underline{u}) \}$  then

$\underline{g}(\underline{x})$  is a subgradient of  $w(\underline{u})$  at  $\underline{u}$  for each  $\underline{x} \in X(\underline{u})$ .

*no wt's extremely easy (another good news) getting subgradients*

Proof\*:

**Observations:**

- Every subgradient of  $w(\underline{u})$  at  $\underline{u}$  can be expressed as a convex combination of the subgradients  $\underline{g}(\underline{x})$  with  $\underline{x} \in X(\underline{u})$ .
- If  $w$  is continuously differentiable at  $\underline{u}$ ,  $X(\underline{u})$  contains a single element  $\underline{\tilde{x}}$  and  $\underline{g}(\underline{\tilde{x}})$  is the gradient of  $w(\underline{u})$  at  $\underline{u}$ .



# Summary

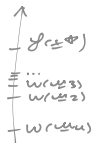
- In general  $(D)$  is easier than  $(P)$  – even if no saddle point exists.
- If a saddle point exists: we can solve  $(D)$  and derive optimal  $\underline{x}^*$  of  $(P)$  by minimizing  $L(\underline{x}, \underline{u}^*)$  over  $X$ , ensuring  $g_i(\underline{x}^*) \leq 0$  and  $u_i^* g_i(\underline{x}^*) = 0 \quad \forall i \in I$ .
- If no saddle point exists: optimal  $\underline{u}^*$  of  $(D)$  gives a lower bound  $w(\underline{u}^*)$  for  $f(\underline{x}^*)$ .

Find  $\underline{u}^* \geq \underline{0}$  maximizing  $w(\underline{u})$  by using the subgradient method that generates  $\{\underline{u}^k\} \rightarrow \underline{u}^*$  when  $k \rightarrow \infty$ .

For each  $\underline{u}^k$ , we have a lower bound  $w(\underline{u}^k)$  for  $f(\underline{x}^*)$  and we determine  $\underline{x}^k$  that minimizes  $L(\underline{x}, \underline{u}^k)$  over  $X$ .

*a sequence of  
problems which are  
easier to solve*

*a sequence of  
lower bounds*



## 5.5 Second order optimality conditions

Nonlinear program:

$$(P) \quad \begin{array}{ll} \min & f(\underline{x}) \\ \text{s.t.} & g_i(\underline{x}) \leq 0 \quad i \in I = \{1, \dots, m\} \\ & h_l(\underline{x}) = 0 \quad l \in L = \{1, \dots, k\} \\ & \underline{x} \in X \subseteq \mathbb{R}^n \end{array}$$

with  $f$ ,  $g_i$ 's and  $h_l$ 's of class  $C^2$  and  $X$  open subset of  $\mathbb{R}^n$ .

Lagrange function:

$$L(\underline{x}, \underline{u}, \underline{v}) = f(\underline{x}) + \sum_{i=1}^m u_i g_i(\underline{x}) + \sum_{l=1}^k v_l h_l(\underline{x}) = \boxed{f(\underline{x}) + \underline{u}^t \underline{g}(\underline{x}) + \underline{v}^t \underline{h}(\underline{x})}$$

with  $\underline{u} \geq 0$  and  $\underline{v} \in \mathbb{R}^k$ .

Hessian submatrix w.r.t. the variables  $x_j$ :

$$\underbrace{\frac{\partial^2}{\partial x^2} L(\underline{x}, \underline{u}, \underline{v})}_{\text{w.r.t } \underline{x}} = \left( \nabla^2 f(\underline{x}) \right) + \left( \sum_{i=1}^m u_i \nabla^2 g_i(\underline{x}) \right) + \left( \sum_{l=1}^k v_l \nabla^2 h_l(\underline{x}) \right)$$

## Second order KKT necessary conditions:

If  $\bar{x}$  is a local minimum of  $(P)$  and  $\nabla g_i(\bar{x})$ , with  $i \in I(\bar{x})$ , and  $\nabla h_l(\bar{x})$ , with  $l \in L$ , are linearly independent, then  $\bar{x}$  and some  $(\bar{u}, \bar{v})$  satisfy the KKT conditions:

*Can take CQ on*

*taking the  $\lambda$  of the second order condition and setting it = 0*

$$\nabla_x L(\underline{x}, \underline{u}, \underline{v}) = \nabla f(\underline{x}) + \sum_{i=1}^m u_i \nabla g_i(\underline{x}) + \sum_{l=1}^k v_l \nabla h_l(\underline{x}) = \underline{0}$$

$$g_i(\underline{x}) \leq 0$$

$$i \in I = \{1, \dots, m\}$$

$$h_l(\underline{x}) = 0$$

$$l \in L = \{1, \dots, k\}$$

$$u_i g_i(\underline{x}) = 0$$

$$i \in I$$

$$\underline{u} \geq \underline{0}, \underline{v} \in \mathbb{R}^k.$$

*the 2nd order part*

Moreover, every  $\underline{d} \in \mathbb{R}^n$  such that

$$\nabla^t g_i(\bar{x}) \underline{d} \leq 0 \quad i \in I(\bar{x})$$

$$\nabla^t h_l(\bar{x}) \underline{d} = 0 \quad l \in L$$

must satisfy

$$\underline{d}^t \nabla_{xx}^2 L(\bar{x}, \bar{u}, \bar{v}) \underline{d} \geq 0$$

*Can show  $\underline{d}$ , the second order condition  $\nabla_{xx}^2 L(\cdot)$  is positive-def*

## Second order KKT sufficient conditions:

Let  $\bar{x}$  satisfies with  $(\bar{u}, \bar{v})$  the previous KKT conditions.

If

$$\underline{d}^t \nabla_{xx}^2 L(\bar{x}, \bar{u}, \bar{v}) \underline{d} > 0$$

for each  $\underline{d} \neq \underline{0}$  such that

$$\nabla^t g_i(\bar{x}) \underline{d} = 0 \quad i \in I^+$$

$$\nabla^t g_i(\bar{x}) \underline{d} \leq 0 \quad i \in I^0$$

$$\nabla^t h_l(\bar{x}) \underline{d} = 0 \quad l = 1, \dots, k$$

*makes sense since if  $u_i > 0$  then  $g_i(\bar{x}) = 0$  by the complementarity slackness conditions*

where  $I^+ = \{i \in I : u_i > 0\}$  and  $I^0 = \{i \in I : u_i = 0\}$ ,

then  $\bar{x}$  is a strict local minimum of  $(P)$ .

*local or these conditions are, to or, not, also the non-convex ones*

See Chap. 12 of J. Nocedal and S. Wright, Numerical Optimization, Springer 1999

## 5.6 Quadratic programming (QP, like LP was linear) maximize

Optimize a quadratic function subject to linear constraints:

$$(P) \quad \begin{aligned} \min \quad & \frac{1}{2} \underline{x}^t Q \underline{x} + \underline{c}^t \underline{x} \\ \text{s.t.} \quad & \underline{a}_i^t \underline{x} \leq b_i \quad i \in I \\ & \underline{a}_i^t \underline{x} \equiv b_i \quad i \in E \\ & \underline{x} \in \mathbb{R}^n, \end{aligned}$$

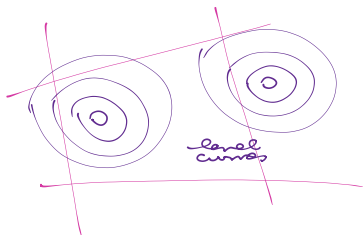
where  $Q \in \mathbb{R}^{n \times n}$ .

Without loss of generality:  $Q$  is symmetric (same function value with  $\bar{Q}$  not symmetric and  $Q = \frac{1}{2}(\bar{Q} + \bar{Q}^t)$ ).

Difficulty depends on  $Q$ : if  $Q$  positive (semi)definite,  $(P)$  convex, otherwise can have a large number of local optima.

Example:  $\min\{\overbrace{-\underline{x}^t \underline{x}}^{\|\underline{x}\| \leq 1} : \underbrace{-1 \leq x_i \leq 1, i = 1, \dots, n}_{\text{exercise with } \mathbb{R}^n}\}$  where all  $2^n$  vertices of  $\{-1, 1\}^n$  are local minima.

## Illustrations of convex Quadratic Programs (QPs):



optimum points  
could be on the  
boundary or else  
in the inside  
  
⇒ more complex  
to solve QPs

QPs are the simplest NLP problems besides Linear Programs. Efficient QP algorithms are available.

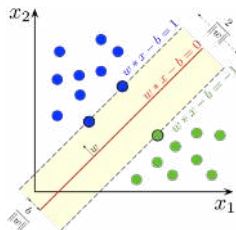
Many direct applications (for portfolio optimization see exercise 9.1).

## Example: Training linear Support Vector Machines (SVMs)

Training set  $T = \{(\underline{x}^i, y^i) : \underline{x}^i \in \mathbb{R}^n, y^i \in \{-1, 1\}, i = 1, \dots, p\}$ .

Linear decision function:  $f(\underline{w}, b, \underline{x}) = \underline{w}^t \underline{x} - b$ .

Separating hyperplane with largest margin (width  $\frac{2}{\|\underline{w}\|}$ ) guarantees best generalization.



Hard-margin linear SVM training:

$$\begin{aligned} \min_{\substack{\underline{w} \in \mathbb{R}^n \\ b \in \mathbb{R}}} & \quad \frac{1}{2} \|\underline{w}\|^2 \quad \Rightarrow \text{quadratic problem} \\ \text{s.t.} & \quad y^i (\underline{w}^t \underline{x}^i - b) - 1 \geq 0 \quad i = 1, \dots, p. \end{aligned}$$

strictly convex function but possibly huge number of linear constraints.

*is the solving rules will be to move to the dual (to have one or more constraints)*

Reformulated as QP with a single constraint using duality:

$$L(\underline{w}, b, \underline{u}) = \frac{1}{2} \|\underline{w}\|^2 - \sum_{i=1}^p u_i (y^i (\underline{w}^t \underline{x}^i - b) - 1)$$

...  $f(\underline{x}) = \sum_{w \in \mathbb{I}} u_i g_i(\underline{x})$  where the opt moves to the min w.r.t  $\underline{x} = (w, b)$

Dual: we look for a feasible point, so we want to see what

$$\max_{\underline{u} \geq 0} \left( \min_{w, b} L(w, b, \underline{u}) \right)$$

$$\Leftrightarrow \max_{\underline{u} \geq 0} L(w, b, \underline{u})$$

$$\text{st } \left. \begin{aligned} \frac{\partial}{\partial w} L(w, b, \underline{u}) &= 0 \\ \frac{\partial}{\partial b} L(w, b, \underline{u}) &= 0 \end{aligned} \right\}$$

we set this gradient = 0 w.r.t to get a min

$$w = \sum_{i=1}^p u_i y_i x_i$$

$$\sum_{i=1}^p u_i y_i = 0$$

which is a hyperplane, and being a convex exercise there is no perfect correspondence of using primal or dual



## 5.6.1 QP with only equality constraints

Consider


$$\min \left\{ \frac{1}{2} \underline{x}^t Q \underline{x} + \underline{c}^t \underline{x} : A \underline{x} = \underline{b} \right\} \quad (1)$$

$\rightarrow \sum_{i \in T} u_i^T x = b_i$

where  $A \in \mathbb{R}^{m \times n}$ .

*this defines a subspace of  $\mathbb{R}^m$  for solving these problems (by constraining) via eqs)*

Since only linear equations, CQ assumption is satisfied at every feasible point and simple KKT conditions:



$$\underbrace{Q \underline{x} + \underline{c}}_{\text{gradient of the obj. function}} + \sum_{i=1}^m \underbrace{u_i \underline{a}_i}_{\text{gradient of the active constraint (see of the row case)}} = \underline{0}$$

$$A \underline{x} = \underline{b}$$

N.B.: Complementary slackness constraints are automatically satisfied.

More or less direct solution of the linear system:

$$\begin{pmatrix} Q & A^t \\ A & 0 \end{pmatrix} \begin{pmatrix} \underline{x} \\ \underline{u} \end{pmatrix} = \begin{pmatrix} -\underline{c} \\ \underline{b} \end{pmatrix}.$$

If  $A$  of full rank and  $Q$  is p.d. on subspace  $\{\underline{x} \in \mathbb{R}^n : A \underline{x} = \underline{0}\}$ , matrix is non singular.

## Null-space method

Determine  $Z \in \mathbb{R}^{n \times (n-m)}$  whose columns span the null space  $\{\underline{x} \in \mathbb{R}^n : A\underline{x} = \underline{0}\}$  of  $A$ .  
*ker(A)*

$Z$  can be computed by (sub)matrix factorization of  $A$  (if  $A$  sparse by  $LU$  factorization).

Given feasible  $\underline{x}_0$ , any other feasible solution

$$\underline{x} = \underline{x}_0 + Z\underline{w}$$

*we we perform a  
change of variables*

for an appropriate  $\underline{w} \in \mathbb{R}^{n-m}$ .

(1) is equivalent to unconstrained QP:  $\frac{1}{2}(\underline{x}_0 + Z\underline{w})^T Q (\underline{x}_0 + Z\underline{w}) + \underline{c}^T (\underline{x}_0 + Z\underline{w})$

$$\min_{\underline{w} \in \mathbb{R}^{n-m}} \left[ \frac{1}{2} \underline{w}^T (Z^T Q Z) \underline{w} + (\underline{Q} \underline{x}_0 + \underline{c})^T Z \underline{w} \right]$$

*if  $Z^T Q Z$  is not def  $\Rightarrow \exists!$  optimal w if set nonempty*  
 $\nabla(\text{test obj}) = 0$  test w  
 $(Z^T Q Z) \underline{w} = -Z^T (\underline{Q} \underline{x}_0 + \underline{c})$

Also other methods but null-space ones are widely used.

## 5.6.2 QP with equality and inequality constraints

### Active-set methods

$$(P) \quad \begin{array}{ll} \min & q(\underline{x}) = \frac{1}{2} \underline{x}^t Q \underline{x} + \underline{c}^t \underline{x} \\ \text{s.t.} & \underline{a}_i^t \underline{x} \leq b_i \quad i \in I \\ & \underline{a}_i^t \underline{x} = b_i \quad i \in E \\ & \underline{x} \in \mathbb{R}^n \end{array}$$

where  $Q \in \mathbb{R}^{n \times n}$ .

*the set of active constraints around  $\underline{x}^*$*

Idea: Determine  $I(\underline{x}^*) = \{i \in I : \underline{a}_i^t \underline{x}^* = b_i\}$  where  $\underline{x}^*$  is an optimal solution, by solving a sequence of QPs with only equality constraints.

idea: if we know which of the inequality constraints are active at the opt set, then we would come back to the previous case of only eq constraints

# Active-set method for convex QPs

*eg. to travel the exact part of simplex methods on LPs*

Initialization: Find initial feasible  $\underline{x}_0$  and

choose  $W_0 \subseteq \underbrace{\{i \in I : \underline{a}_i^t \underline{x}_0 = b_i\}}_{I(\underline{x}_0)} \cup E$  of the active constraints at  $\underline{x}_0$ , with  $E \subseteq W_0$ .

$\Rightarrow$   $W_0$  is like the starting working set on which we optimize of course

Iteration  $k$ :

Given current feasible  $\underline{x}_k$ , determine  $\underline{d}_k$  by solving the subproblem:

$$\min \{ q(\underline{x}_k + \underline{d}) : \underline{a}_i^t(\underline{x}_k + \underline{d}) = b_i, i \in W_k \}, \quad (2)$$

*only consider constraints according to  $W_k$*

*$\underline{x}_k + \underline{d}$  must satisfy all the  $W_k$  constraints (and only them)*

where  $W_k$  is current working set, with  $W_k \subseteq \{i \in I : \underline{a}_i^t \underline{x}_k = b_i\} \cup E$ .

*a choice (of indexes) among the active constraints*

*plus the full index set of eq. constraints*

(2) is equivalent to:

*$\underline{a}_i^t \underline{x}_k = b_i$  already construction*

$$\min \{ q(\underline{x}_k + \underline{d}) : \underline{a}_i^t \underline{d} = 0, i \in W_k \}. \quad (3)$$

*no we are just left with these*

N.B.: If  $Z^t Q Z$  is p.d. (always true if  $Q$  is p.d.), (3) has a unique solution  $\underline{d}_k$ .

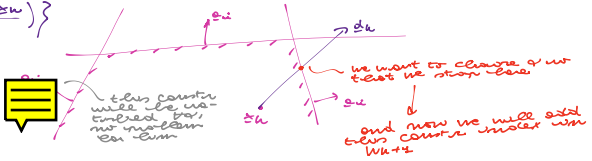
Based on solution  $\underline{d}_k$  of (3), we determine  $\alpha_k$ ,  $\underline{x}_{k+1} = \underline{x}_k + \alpha_k \underline{d}_k$  and  $W_{k+1}$ .

the ones which can now be left out

- If  $\underline{d}_k \neq \underline{0}$ , we determine the largest  $\alpha$  satisfying all constraints not in  $W_k$ :

$$\alpha_k = \min \left\{ 1, \min_{\substack{i \notin W_k: \\ a_i^T \underline{d}_k > 0}} \left( \frac{b_i - a_i^T \underline{x}_k}{a_i^T \underline{d}_k} \right) \right\}$$

*we ignore the ones that are 1 (we already handled)*



and set  $\underline{x}_{k+1} = \underline{x}_k + \alpha_k \underline{d}_k$ .

$W_{k+1} = W_k \cup \{i'\}$  where  $i'$  is index of one constraint becoming active at  $\underline{x}_{k+1}$ .

- If  $\underline{d}_k = \underline{0}$ ,  $\underline{x}_k$  is a minimum over subspace defined by  $W_k$  and we set  $\underline{x}_{k+1} = \underline{x}_k$ .  
*but we have to check if it's optimal, we are not done yet*

KKT conditions of (3) imply there are multipliers  $u_i^k$  such that:

$$\left( Q \underline{x}_k + \underline{c} \right) + \sum_{i \in W_k} u_i^k \underline{a}_i = \underline{0}. \tag{4}$$

If  $u_i^k \geq 0$  for every  $i \in W_k \cap I$  Then  $\underline{x}_k$  is a local optimum of original QP

Else  $W_{k+1} = W_k \setminus \{i'\}$  where  $i'$  is the index with the most negative  $u_{i'}^k$ .

*we set  $u_{i'}^k < 0$*

- If  $\underline{d}_k \neq \underline{0}$ , we determine the largest  $\alpha$  satisfying all constraints not in  $W_k$ :



and set  $\underline{x}_{k+1} = \underline{x}_k + \alpha_k \underline{d}_k$ .

$W_{k+1} = W_k \cup \{i'\}$  where  $i'$  is index of one constraint becoming active at  $\underline{x}_{k+1}$ .

- If  $\underline{d}_k = \underline{0}$ ,  $\underline{x}_k$  is a minimum over subspace defined by  $W_k$  and we set  $\underline{x}_{k+1} = \underline{x}_k$ .

KKT conditions of (3) imply there are multipliers  $u_i^k$  such that:

$$Q\underline{x}_k + \underline{c} + \sum_{i \in W_k} u_i^k \underline{a}_i = \underline{0}. \quad (4)$$

If  $u_i^k \geq 0$  for every  $i \in W_k \cap I$  Then  $\underline{x}_k$  is a local optimum of original QP

Else  $W_{k+1} = W_k \setminus \{i'\}$  where  $i'$  is the index with the most negative  $u_{i'}^k$ .

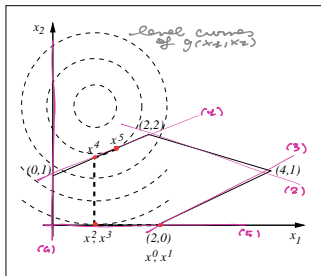
**Proposition:** If  $Q$  is p.d. ( $q$  is strictly convex), the method (with anti-cycling rule) finds an optimal solution within a finite number of iterations.

*works because of:*  
 Note: finite number of working sets. and  
*of changes on the Ws*

Example:

$$\begin{aligned} \min \quad & q(x_1, x_2) = (x_1 - 1)^2 + (x_2 - 2.5)^2 \\ \text{s.t.} \quad & -x_1 + 2x_2 - 2 \leq 0 & (1) \\ & x_1 + 2x_2 - 6 \leq 0 & (2) \\ & x_1 - 2x_2 - 2 \leq 0 & (3) \\ & -x_1 \leq 0 & (4) \\ & -x_2 \leq 0 & (5) \end{aligned}$$

Figure:



From J. Nocedal, S. Wright, Numerical Optimization, First Edition, Springer 1999, p. 462-463.

Iteration 0:

$$\underline{x}_0 = \begin{pmatrix} 2 \\ 0 \end{pmatrix} \text{ and we take } W_0 = \{3, 5\}.$$

Since  $\underline{x}_0$  is a vertex of the feasible solution polyhedron,

$\underline{x}_0$  minimizes  $q(\underline{x})$  w.r.t.  $W_0$  and

$$\underline{d}_0 = \underline{0} \text{ is optimal solution of } \min\{ q(\underline{x}_0 + \underline{d}) : \underline{a}_i^t \underline{d} = 0, i \in W_0 \}.$$

*since we took  $\underline{x}_0$  on the intersection of the two constraints, and not just a "random inside" feasible set*

$$\text{Thus } \underline{x}_1 = \underline{x}_0 + \alpha_0 \underline{d}_0 = \underline{x}_0.$$

KKT conditions:

$$\nabla q(\underline{x}_0) = \begin{pmatrix} 2 \\ -5 \end{pmatrix} = u_3 \begin{pmatrix} -1 \\ 2 \end{pmatrix} + u_5 \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$

we obtain the multipliers  $\begin{pmatrix} u_3 \\ u_5 \end{pmatrix} = \begin{pmatrix} -2 \\ -1 \end{pmatrix}$  for the active constraints.

Since  $u_3 < u_5 < 0$ , we set  $W_1 = W_0 \setminus \{3\} = \{5\}$ .

*now we are just optimizing moving on the  $x_2$  axis*



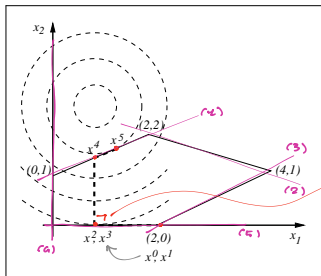
**Proposition:** If  $Q$  is p.d. ( $q$  is strictly convex), the method (with anti-cycling rule) finds an optimal solution within a finite number of iterations.

*works because of:*  
 Note: finite number of working sets. and  
*of changes on the Wk*

Example:

$$\begin{aligned} \min \quad & q(x_1, x_2) = (x_1 - 1)^2 + (x_2 - 2.5)^2 \\ \text{s.t.} \quad & -x_1 + 2x_2 - 2 \leq 0 \quad (1) \\ & x_1 + 2x_2 - 6 \leq 0 \quad (2) \\ & x_1 - 2x_2 - 2 \leq 0 \quad (3) \\ & -x_1 \leq 0 \quad (4) \\ & -x_2 \leq 0 \quad (5) \end{aligned}$$

Figure:



*x2 becomes optimal  
 as it's the furthest  
 away from the center  
 of the level curves*

From J. Nocedal, S. Wright, Numerical Optimization, First Edition, Springer 1999, p. 462-463.

### Iteration 1:

Optimal solution of  $\min\{ q(\underline{x}_1 + \underline{d}) : \underline{a}_i^t \underline{d} = 0, i \in W_1 \}$  is  $\underline{d}_1 = \begin{pmatrix} -1 \\ 0 \end{pmatrix}$ .

Since  $\underline{d}_1$  does not violate any constraint with indices not in  $W_1$ ,  $\alpha_1 = 1$  and  $\underline{x}_2 = \underline{x}_1 + \alpha_1 \underline{d}_1 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$ .

Since at  $\underline{x}_2$  no other constraints are active, we set  $W_2 = W_1 = \{5\}$ .

### Iteration 2:

Optimal solution of  $\min\{ q(\underline{x}_2 + \underline{d}) : \underline{a}_i^t \underline{d} = 0, i \in W_2 \}$  is  $\underline{d}_2 = \underline{0}$ .

From KKT conditions

$$\nabla q(\underline{x}_2) = \begin{pmatrix} 0 \\ -5 \end{pmatrix} = u_5 \begin{pmatrix} 0 \\ 1 \end{pmatrix},$$

we obtain  $u_5 = -5$ .

Thus  $\underline{x}_3 = \underline{x}_2$  and we set  $W_3 = W_2 \setminus \{5\} = \emptyset$ .

*as we get to optimize (less) and we moved more vertically towards the level curves center*

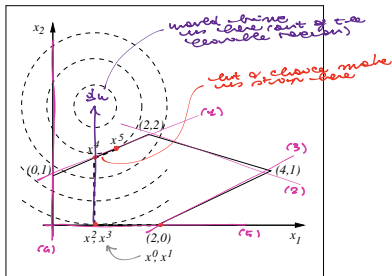
**Proposition:** If  $Q$  is p.d. ( $q$  is strictly convex), the method (with anti-cycling rule) finds an optimal solution within a finite number of iterations.

*works because of:*  
 Note: finite number of working sets. and  
*of choices for the Wk*

Example:

$$\begin{aligned} \min \quad & q(x_1, x_2) = (x_1 - 1)^2 + (x_2 - 2.5)^2 \\ \text{s.t.} \quad & -x_1 + 2x_2 - 2 \leq 0 \quad (1) \\ & x_1 + 2x_2 - 6 \leq 0 \quad (2) \\ & x_1 - 2x_2 - 2 \leq 0 \quad (3) \\ & -x_1 \leq 0 \quad (4) \\ & -x_2 \leq 0 \quad (5) \end{aligned}$$

Figure:



From J. Nocedal, S. Wright, Numerical Optimization, First Edition, Springer 1999, p. 462-463.

### Iteration 3:

Optimal solution of  $\min\{ q(\underline{x}_3 + \underline{d}) : \underline{a}_i^t \underline{d} = 0, i \in W_3 \}$  is  $\underline{d}_3 = \begin{pmatrix} 0 \\ 2.5 \end{pmatrix}$ .

Since  $\underline{d}_3$  violates constraints (1) and (2) which are not in  $W_1$ ,  $\alpha_3 = 0.6$  and  $\underline{x}_4 = \underline{x}_3 + \alpha_3 \underline{d}_3 = \begin{pmatrix} 1 \\ 1.5 \end{pmatrix}$ .

Since at  $\underline{x}_4$  only constraint (1) becomes active, we set  $W_4 = \{1\}$ .

### Iteration 4:

Optimal solution of  $\min\{ q(\underline{x}_4 + \underline{d}) : \underline{a}_i^t \underline{d} = 0, i \in W_4 \}$  is  $\underline{d}_4 = \begin{pmatrix} 0.4 \\ 0.2 \end{pmatrix}$ .

Since  $\underline{x}_4 + \underline{d}_4 = \begin{pmatrix} 1.4 \\ 1.7 \end{pmatrix}$  satisfies all the constraints with indices not in  $W_1$ , we take  $\alpha_4 = 1$ , set  $\underline{x}_5 = \underline{x}_4 + \underline{d}_4$  and  $W_5 = W_4 = \{1\}$ .

Iteration 5:

Optimal solution of  $\min\{ q(\underline{x}_5 + \underline{d}) : \underline{a}_i^t \underline{d} = 0, i \in W_5 \}$  is  $\underline{d}_5 = 0$ .

Solving the KKT conditions (4) we obtain  $u_1 = 1.25 \geq 0$ .

Thus  $\underline{x}_5 = \begin{pmatrix} 1.4 \\ 1.7 \end{pmatrix}$  is optimal for the original problem.

## 5.6.3 Non convex QP and solvers

If  $Q$  has some negative eigenvalues, the active-set method for convex QP can be adapted by modifying  $\underline{d}_k$  and  $\alpha_k$  in certain situations.

See J. Nocedal, S. Wright, Numerical Optimization, First edition, Springer 1999, p. 468-474.

Since  $W_k$  may change by just one index at every iteration, efficient QP solvers proceed by successive updates of the factors computed at the previous iterations.

Available active-set-based solvers: LINDO, QPOPT, NAG Library, Matlab,...

## 5.7 Penalty method and augmented Lagrangian method

Generic NLP:

*and changed the function name*

*now so far we've (relative to 5.0)*

$$\begin{aligned} \min \quad & f(\underline{x}) \\ \text{s.t.} \quad & c_i(\underline{x}) \geq 0 \quad i \in I \\ & c_i(\underline{x}) = 0 \quad i \in E \\ & \underline{x} \in \mathbb{R}^n \end{aligned} \tag{1}$$

where  $f$  and  $c_i$ 's are of class  $\mathcal{C}^1$  or  $\mathcal{C}^2$ .

Notation, examples and proofs: see Chapter 17 of J. Nocedal, S. Wright, Numerical Optimization, Springer, 1999, p. 491-500.

## 5.7.1 Quadratic penalty method

*to reduce to unconstrained opt* *with the obj function*

**Idea:** Delete constraints, penalize their violation and solve a sequence of unconstrained optimization problems.

Description for

$$\begin{aligned} \min \quad & f(\underline{x}) \\ \text{s.t.} \quad & c_i(\underline{x}) = 0 \quad i \in E = \{1, \dots, m\} \\ & \underline{x} \in \mathbb{R}^n. \end{aligned} \tag{2}$$

**Definition:** The quadratic penalty function problem associated to (2) is

$$\min_{\underline{x} \in \mathbb{R}^n} Q(\underline{x}, \mu) = f(\underline{x}) + \frac{1}{2\mu} \sum_{i \in E} c_i^2(\underline{x}) \tag{3}$$

with penalty parameter  $\mu > 0$ .

*to prevent the penalization*

*the more penalties more the closer violations (infeasibilities since with the constraint  $c_i(\underline{x}) = 0$ , were eq constraint)*

We consider  $\{\mu_k\}_{k \geq 1}$  with  $\lim_{k \rightarrow \infty} \mu_k = 0$  and, for each  $k$ , we determine an approximate solution  $\underline{x}_k$  of (3) using an unconstrained optimization method.

*we trust the penalization goes to the (is inevitable) we "force" the constraints*



Example:

$$\begin{aligned} \min \quad & x_1 + x_2 \\ \text{s.t.} \quad & x_1^2 + x_2^2 - 2 = 0 \end{aligned}$$

with optimal solution  $(-1, -1)^t$ .

Quadratic penalty problem:

$$\begin{aligned} \min_{x \in \mathbb{R}^2} Q(x; \mu) &= f(x) + \frac{\mu}{2\mu} \sum_{w \in E} c_w(x)^2 = \\ &= x_1 + x_2 + \frac{\mu}{2\mu} (x_1^2 + x_2^2 - 2) \end{aligned}$$

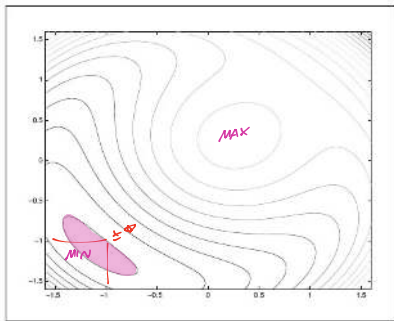


Figure 17.1 Contours of  $Q(x; \mu)$  from (17.4) for  $\mu = 1$

for  $\mu = 1$  the minimum of  $Q(x; \mu)$  is close to  $x^* = (-1, -1)$

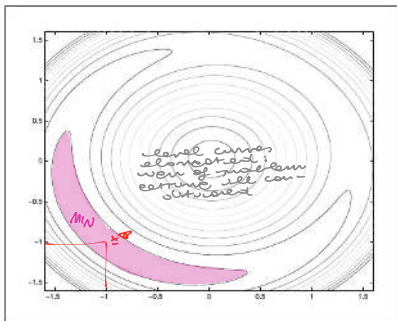


Figure 17.2 Contours of  $Q(x; \mu)$  from (17.4) for  $\mu = 10$

when  $\mu$  is large enough, all contours are attracted

## General scheme

*we want to be more and more accurate*

0) Select  $\varepsilon > 0$ ,  $\mu_0 > 0$ , sequence of tolerances  $\{\tau_k\}_{k \geq 0}$  with  $\tau_k > 0$  and  $\lim_{k \rightarrow \infty} \tau_k = 0$ .

Choose initial  $\underline{x}_0^s$  and set  $k = 0$ .  
*near starting set*

1) Determine an approximate minimizer  $\underline{x}_k$  of  $Q(\underline{x}, \mu_k)$  starting from  $\underline{x}_k^s$  and terminate when  $\|\nabla Q(\underline{x}, \mu_k)\| \leq \tau_k$ .  
*ie when the required accuracy is reached*

terminate this step (2)

2) If termination condition is satisfied (e.g.  $|f(\underline{x}_{k-1}) - f(\underline{x}_k)| < \varepsilon$ )

Then return solution  $\underline{x}_k$

Else choose  $\mu_{k+1} \in (0, \mu_k)$  and starting  $\underline{x}_{k+1}^s$ , set  $k = k + 1$  and Goto 1)

*decrease  $\mu_k$  ie increase regularization*  
 $\dots \rightarrow \frac{\mu}{2\mu} \leq \dots \rightarrow$  *smaller  $\mu$  lower regularization*

Choices:

- For convergence results, it suffices that  $\lim_{k \rightarrow \infty} \tau_k = 0$ .
- $\{\mu_k\}_{k \geq 0}$  generated adaptively starting from  $\mu_0$ : if minimization of  $Q(\underline{x}, \mu_k)$  is "difficult" set e.g.  $\mu_{k+1} = 0.7\mu_k$ , otherwise  $\mu_{k+1} = 0.1\mu_k$ .  
*smaller/lower decrease if problem is difficult*
- Judicious choice of the starting  $\underline{x}_k^s$  when solving unconstrained penalty problem at each iteration:  $\underline{x}_{k+1}^s := \underline{x}_k$   
*the set that we est using the lra unconstrained problem*

# Convergence

unresolutive, we are online to solve at optimum; each sub-problem ( $\tau_k = 0$   $\forall k$ )

interesting but unresolutive

**Theorem 1:** Suppose each  $x_k$  is a global minimizer of  $Q(x, \mu_k)$  and  $\lim_{k \rightarrow \infty} \mu_k = 0$ , then every limit point  $x^*$  of  $\{x_k\}_{k \geq 0}$  generated with above scheme ( $\tau_k = 0, \forall k \geq 0$ ) is a global minimum of problem (2).

two cases to two - notice w/  $\|\nabla Q(x_k, \mu_k)\| = 0$

Proof:

Let  $\bar{x}$  be an optimal solution of (2).

Since  $\bar{x}$  is a global min of  $Q(x, \mu)$  and  $\bar{x}$  is feasible for (2, the overall problem) then

$$Q(\bar{x}, \mu_k) \leq Q(x_k, \mu_k) \quad \forall k$$

global min
just a feasible set (for Q)

since  $\bar{x}$  is feasible and opt for (2)

morely:

$$f(\bar{x}) + \frac{\mu}{24\mu} \sum_{\omega=4}^m c_{\omega}(\bar{x})^2 \leq f(x_k) + \frac{\mu}{24\mu} \sum_{\omega=4}^m c_{\omega}(x_k)^2 = f(x_k) \quad (6)$$

$$\Leftrightarrow \sum_{\omega=4}^m c_{\omega}(\bar{x})^2 \leq 24\mu [f(x_k) - f(\bar{x})] \quad \forall k \quad (7)$$

Consider an sub-sequence of  $\{x_k\}_{k \geq 0}$  with  $k \in K: \lim_{k \in K} x_k = \bar{x}$

B) Letting  $k \rightarrow \infty$  with  $k \in K$  we obtain in (7), the following limit:

$$\sum_{\omega=4}^m c_{\omega}(\bar{x})^2 = \lim_{k \in K} \sum_{\omega=4}^m c_{\omega}(x_k)^2 \leq 0 \Rightarrow c_{\omega}(\bar{x}) = 0 \quad \forall \omega \in F$$

(we see limit point is feasible for (2))

now we show that  $\pm \Phi$  is not one) (feasible but slow the optimal set.

B) let's assume  $k \rightarrow +\infty$  with  $k \in \mathcal{K}$  with (a) we obtain

$$\underline{f(\pm \Phi)} \leq f(\pm \Phi) + \lim_{k \in \mathcal{K}} \frac{\epsilon}{2\gamma_k} \sum_{w=1}^m c_w(\pm u)^2 \leq \underline{f(\bar{x})}$$

$\Rightarrow \pm \Phi$  is an opt set of (2)

no need to solve at opt the unconstrained problems

more restrictive

we demand more and more precision

**Theorem 2:** If

- tolerances  $\tau_k > 0$  satisfy  $\lim_{k \rightarrow \infty} \tau_k = 0$
- $\lim_{k \rightarrow \infty} \mu_k = 0$ ,

then every limit point  $\underline{x}^*$  of  $\{\underline{x}_k\}_{k>0}$  at which all  $\nabla c_i(\underline{x}^*)$ , with  $i \in E$ , are linearly independent is a KKT point of problem (2).

For such points, the subsequence defined by  $\mathcal{K}$  with  $\lim_{k \in \mathcal{K}} \underline{x}_k = \underline{x}^*$  satisfies

$$\lim_{k \in \mathcal{K}} -\frac{c_i(\underline{x}_k)}{\mu_k} = u_i^* \quad \forall i \in E,$$

where  $\underline{u}^*$  satisfies with  $\underline{x}^*$  the KKT conditions for problem (2).

the corresponding optimal kkt multipliers

if  $> 0$  and  $= 0$  we are not error  
 as we will try to drive this to zero

$\Rightarrow -c_w(\pm u) \approx \gamma_k u_i^2$   
 and as we can tune the  $\gamma_k$  arbitrarily, more precise, during the method almost lim, according to  $u_i^*$ . This will be the idea of the outer-iterated Lagrangian method

Observation: (4) implies that

- i) The minimizer  $\underline{x}_k$  of  $Q(\underline{x}, \mu_k)$  does not satisfy  $c_i(\underline{x}) = 0$  exactly for all  $i \in E$  ( $c_i(\underline{x}_k) = -\mu_k u_i^*$ ). To obtain a feasible solution, we must  $\mu_k \rightarrow 0$ .
- ii) In some circumstances  $-\frac{c_i(\underline{x}_k)}{\mu_k}$  may be used as estimates of  $u_i^*$ .

Recall: Lagrange function for problem (2) is

$$L(\underline{x}, \underline{u}) = f(\underline{x}) - \sum_{i=1}^m u_i c_i(\underline{x}) \quad (5)$$

*when unconstrained due to equality constraint*

and KKT conditions require that, apart from  $c_i(\underline{x}) = 0$

*no complementary slackness  
now since we have equality constraint*

$$\nabla_{\underline{x}} L(\underline{x}, \underline{u}) = \nabla f(\underline{x}) - \sum_{i=1}^m u_i \nabla c_i(\underline{x}) = \underline{0}. \quad (6)$$

By comparing

$$Q(\underline{x}, \mu) = f(\underline{x}) + \frac{1}{2\mu} \sum_{i=1}^m c_i(\underline{x})^2$$
$$\nabla_{\underline{x}} Q(\underline{x}, \mu) = \nabla f(\underline{x}) + \frac{1}{\mu} \sum_{i=1}^m c_i(\underline{x}) \nabla c_i(\underline{x}) = \underline{0} \quad (7)$$

and (6), it appears that  $-\frac{c_i(\underline{x})}{\mu}$  has been substituted with  $u_i$ .

It can be proved that if  $\tau_k \rightarrow 0$  then  $\underline{x}_k \rightarrow \underline{x}^*$  and  $-\frac{c_i(\underline{x}_k)}{\mu_k} \rightarrow u_i^*$   $i = 1, 2, \dots, m$ .

*that's the main problem of this method*

Observation: When  $\mu_k \rightarrow 0$  the quadratic penalty problem (3) becomes ill conditioned.

$$\nabla_{\underline{x}\underline{x}}^2 Q(\underline{x}, \mu_k) = \nabla^2 f(\underline{x}) + \frac{1}{\mu_k} A^t(\underline{x}) A(\underline{x}) + \frac{1}{\mu_k} \sum_{i=1}^m c_i(\underline{x}) \nabla^2 c_i(\underline{x}) \quad (8)$$

*min Q(x, y\_k)  
x ∈ ℝ<sup>m</sup>*

*we see this looking at the lesson*

where  $A^t(\underline{x}) = [\nabla c_1(\underline{x}); \dots; \nabla c_m(\underline{x})]$  and  $A \in \mathbb{R}^{m \times n}$  of full rank  $m \leq n$ , usually  $m < n$ .

When  $\underline{x}$  is close to minimizer of  $Q(\underline{x}, \mu_k)$  and assumptions of Theorem 2 are satisfied, (4) implies that

$$\nabla_{\underline{x}\underline{x}}^2 Q(\underline{x}, \mu_k) \approx \nabla_{\underline{x}\underline{x}}^2 L(\underline{x}, \underline{u}^*) + \frac{1}{\mu_k} A^t(\underline{x}) A(\underline{x}). \quad (9)$$

*optimal wrt constraint*

Since  $\nabla_{\underline{x}\underline{x}}^2 L(\underline{x}, \underline{u}^*)$  does not depend on  $\mu_k$  and  $\frac{1}{\mu_k} A^t(\underline{x}) A(\underline{x})$  has  $n - m$  eigenvalues of value 0 and  $m$  eigenvalues of value  $O(1/\mu_k)$ , numerical issues arise when  $\mu_k \rightarrow 0$ .

For convenience, this approach can also be extended to inequality constraints

## Problems with both equality and inequality constraints:

### Quadratic penalty problem

$$\min_{\underline{x} \in \mathbb{R}^n} Q(\underline{x}, \mu) = f(\underline{x}) + \frac{1}{2\mu} \sum_{i \in E} c_i^2(\underline{x}) + \frac{1}{2\mu} \sum_{i \in I} ([c_i(\underline{x})]^-)^2 \quad (10)$$

$= \begin{cases} c_i(\underline{x}) & \text{if } c_i(\underline{x}) < 0 \\ 0 & \text{otherwise} \end{cases}$

where  $[y]^-$  denotes  $\max(-y, 0)$ .

also here we may set differentiable, not always necessary

Other penalty functions are available.

If only equality constraints, the exact penalty problem is

$$\min_{\underline{x} \in \mathbb{R}^n} Q(\underline{x}, \mu) = f(\underline{x}) + \frac{1}{2\mu} \sum_{i \in E} |c_i(\underline{x})|. \quad (11)$$

N.B.: Q is not everywhere differentiable.

due to the 1.1

## 5.7.2 Augmented Lagrangian method

Idea: Reduce ill-conditioning issues of the unconstrained subproblems (in quadratic penalty method) by introducing explicit estimates of the Lagrange multipliers.

Description for

$$\begin{aligned}
 \min \quad & f(\underline{x}) \\
 \text{s.t.} \quad & c_i(\underline{x}) = 0 \quad i \in E = \{1, \dots, m\} \\
 & \underline{x} \in \mathbb{R}^n.
 \end{aligned} \tag{12}$$

*same of above*

**Definition:** The augmented Lagrange function associated to problem (12) is

*we add a quadratic penalty term  
 but to the augmented function  
 (rather than to the obj. function)*

$$L_A(\underline{x}, \underline{u}, \gamma) = \underbrace{\left[ f(\underline{x}) + \sum_{i=1}^m u_i c_i(\underline{x}) \right]}_{L(\underline{x}, \underline{u})} + \frac{\gamma}{2\gamma} \left[ \sum_{i=1}^m c_i(\underline{x})^2 \right]$$

*see mult*      *penalty term*

*we are indifferent  
 as we have equality constraint*

Since the KKT could require that  $\nabla_x L(\underline{x}^*, \underline{u}^*) = 0$  and  $c_i(\underline{x}^*) = 0$  but then at optimality

$$L_A(\underline{x}, \underline{u}, \gamma) = L(\underline{x}, \underline{u}) \Rightarrow \text{no need for } \gamma > 0$$



## Similar approach:

At each iteration:  $\mu_k > 0$  and determine an approximate minimizer  $\underline{x}_k$  of  $L_A(\underline{x}, \underline{u}^k, \mu_k)$  via an unconstrained optimization method, where  $\underline{u}^k$  is an updated estimate.

*we need to understand how to update  $\underline{u}$*

Differentiating w.r.t.  $\underline{x}$ , we obtain

$$\nabla_{\underline{x}} L_A(\underline{x}, \underline{u}, \mu) = \nabla f(\underline{x}) - \sum_{i=1}^m \left( u_i - \frac{c_i(\underline{x})}{\mu} \right) \nabla c_i(\underline{x}).$$

$$L_A(c, \underline{x}, \underline{u}, \mu) = \underbrace{f(\underline{x})}_{\text{obj}} + \sum_{i=1}^m \underbrace{u_i c_i(\underline{x})}_{\text{constraints}} + \frac{1}{2\mu} \left[ \sum_{i=1}^m \underbrace{u_i^2}_{\text{weights}} \right]$$

Considerations similar to those in proof of Theorem 2 allow to establish that

$$u_i^* \approx u_i^k - \frac{c_i(\underline{x}_k)}{\mu_k} \quad i \in E, \quad (13)$$

*relation that should be understood (to call)*

which is equivalent to

$$c_i(\underline{x}_k) \approx \mu_k (u_i^k - u_i^*) \quad i \in E. \quad (14)$$

*ca tries to get towards 0 (due to less weight)*

*we can keep  $\mu_k$  positive (to avoid ill-posed problems)*

*and instead we drive these  $\rightarrow 0$*

## General scheme

- 0) Choose  $\varepsilon > 0$ ,  $\mu_0 > 0$ , tolerances  $\{\tau_k\}_{k \geq 0}$  with  $\tau_k > 0$  and  $\lim_{k \rightarrow \infty} \tau_k = 0$ ,  $\underline{x}_0^s$  and initial  $\underline{u}^0$ , set  $k := 0$ .
- 1) Determine an approximate minimizer  $\underline{x}_k$  of  $L_A(\underline{x}, \underline{u}^k, \mu_k)$  starting from  $\underline{x}_k^s$  and terminate when  $\|\nabla_{\underline{x}} L_A(\underline{x}, \underline{u}^k, \mu_k)\| \leq \tau_k$ .
- 2) If overall termination condition is satisfied (e.g.,  $|f(\underline{x}_{k-1}) - f(\underline{x}_k)| < \varepsilon$ )

Then Stop

Else set 
$$u_i^{k+1} = u_i^k - \frac{c_i(\underline{x}_k)}{\mu_k} \quad \text{for } i \in E \quad (15)$$

choose  $\mu_{k+1} \in (0, \mu_k)$  and next starting solution  $\underline{x}_{k+1}^s$

set  $k := k + 1$  and Goto 1)

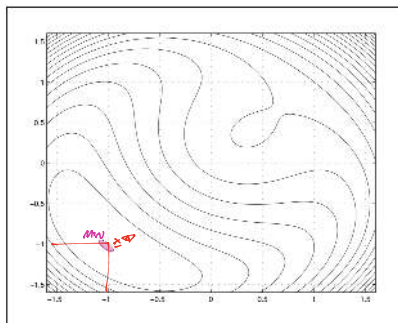
Including in  $L_A$  an additional term related to the Lagrange multipliers leads to substantial improvements w.r.t. the quadratic penalty method.

Example:

$$\begin{aligned} \min \quad & x_1 + x_2 \\ \text{s.t.} \quad & x_1^2 + x_2^2 - 2 = 0 \end{aligned}$$

with optimal solution  $\underline{x}^* = (-1, -1)^t$ , optimal multiplier  $u^* = -0.5$  and unconstrained optimization subproblem:

From J. Nocedal, S. Wright, Numerical Optimization, Springer, 1999, p. 513-514.



*outer level curves*

Figure 17.5 Contours of  $\mathcal{L}_A(x, \lambda; \mu)$  from (17.40) for  $\lambda = -0.4$  and  $\mu = 1$ .

Including in  $L_A$  an additional term related to the Lagrange multipliers leads to substantial improvements w.r.t. the quadratic penalty method.

Example:

$$\begin{aligned} \min \quad & x_1 + x_2 \\ \text{s.t.} \quad & x_1^2 + x_2^2 - 2 = 0 \end{aligned}$$

with optimal solution  $\underline{x}^* = (-1, -1)^t$ , optimal multiplier  $u^* = -0.5$  and unconstrained optimization subproblem:

From J. Nocedal, S. Wright, Numerical Optimization, Springer, 1999, p. 513-514.

$$\min_{\underline{x} \in \mathbb{R}^2} L_A(\underline{x}, u, \mu) = [ (x_1 + x_2) - u(x_1^2 + x_2^2 - 2) ] + \left[ \frac{\mu}{24} (x_1^2 + x_2^2 - 2)^2 \right]$$

Suppose that  $\mu_k = \mu$  and estimate  $u^k = -0.6$ .  
Contours of  $L_A(\pm 1, -0.6, \mu)$  are smoother than those of  $Q(\pm 1, \mu)$   
but the minimum

$$\underline{x}^k = \begin{pmatrix} -2.02 \\ -4.02 \end{pmatrix} \text{ of } L_A(-)$$

is much closer to  $\underline{x}^*$  than the

$$\begin{pmatrix} -2.4 \\ -2.4 \end{pmatrix} \text{ of } Q(-)$$

$\Rightarrow$  reference line improved

### Theorem 3:

Let  $\underline{x}^*$  be a local minimum of (12) at which the  $\nabla c_i(\underline{x}^*)$ ,  $i \in E$ , are linearly independent and 2nd order sufficient optimality conditions are satisfied for  $\underline{u} = \underline{u}^*$ .

Then  $\exists \bar{\mu} > 0$  such that for all  $\mu \in (0, \bar{\mu}]$ ,  $\underline{x}^*$  is a strict local minimum of  $L_A(\underline{x}, \underline{u}^*, \mu)$ .

N.B.: In general  $\underline{u}^*$  is unknown.

The next result

*the optimal multiplier*

- concerns the more realistic case in which  $\underline{u} \neq \underline{u}^*$ ,
- provides conditions under which  $\exists$  a minimizer of  $L_A$  close to  $\underline{x}^*$  and error bounds on  $\underline{x}_k$  and on  $\underline{u}^{k+1}$ .

#### Theorem 4:

Suppose the assumptions of Theorem 3 are satisfied at  $\underline{x}^*$  and  $\underline{u}^*$ , and let  $\bar{\mu} > 0$  be the corresponding threshold.

Then  $\exists$  scalars  $\delta > 0$ ,  $\varepsilon > 0$ , and  $M$  such that

i) For all  $\underline{u}^k$  and  $\mu_k$  satisfying

$$\|\underline{u}^k - \underline{u}^*\| \leq \delta/\mu_k, \quad \mu_k \leq \bar{\mu}, \quad (16)$$

the problem

*$\Rightarrow$  local, the problem is well defined*

$$\min_{\underline{x} \in \mathbb{R}^n : \|\underline{x} - \underline{x}^*\| \leq \varepsilon} L_A(\underline{x}, \underline{u}^k, \mu_k)$$

has a unique solution  $\underline{x}_k$ . Moreover, we have  $\|\underline{x}_k - \underline{x}^*\| \leq M\mu_k\|\underline{u}^k - \underline{u}^*\|$ .

ii) For all  $\underline{u}^k$  and  $\mu_k$  satisfying (16), we have

*moreover we have an error bound (which depends on the quality of the initial guess)*

$$\|\underline{u}^{k+1} - \underline{u}^*\| \leq M\mu_k\|\underline{u}^k - \underline{u}^*\|,$$

where  $\underline{u}^{k+1}$  is given by the formula (15). *good news about convergence speed*

iii) For all  $\underline{u}^k$  and  $\mu_k$  satisfying (16), the matrix  $\nabla_{\underline{x}\underline{x}}^2 L_A(\underline{x}_k, \underline{u}^k, \mu_k)$  is positive definite and the  $\nabla c_i(\underline{x}_k)$ , with  $i \in E$ , are linearly independent. *at  $\underline{x}_k$*

*we need, it's all well defined*

Problems with also inequality constraints:

We can introduce slack variables and substitute  $c_i(\underline{x}) \geq 0$ ,  $i \in I$ , with

$$c_i(\underline{x}) - s_i = 0, \quad s_i \geq 0, \quad i \in I.$$

In LANCELOT solver, the bounds on the variables are explicitly taken into account in the subproblem

$$\min_{l_{inf} \leq \underline{x} \leq l_{sup}} L_A(\underline{x}, \underline{u}^k, \mu_k).$$

## 5.8 Barrier method

Description for:

$$\begin{aligned}
 \min \quad & f(\underline{x}) \\
 \text{s.t.} \quad & c_i(\underline{x}) \geq 0 \quad i \in I = \{1, \dots, m\} \\
 & \underline{x} \in \mathbb{R}^n.
 \end{aligned} \tag{1}$$

Notation and examples: Chapter 17 of J. Nocedal, S. Wright, Numerical Optimization, Springer, 1999, p. 498-508.

**Definition:** Let

$$X^\circ = \text{int}(\{\underline{x} \in \mathbb{R}^n : c_i(\underline{x}) \geq 0, i \in I\}) \neq \emptyset,$$

a function defined on  $\mathbb{R}^n$  is a **barrier function** if it is continuous over  $X^\circ$ , tends to  $\infty$  when approaching  $\partial X$  and has value  $\infty$  on  $\mathbb{R}^n \setminus X^\circ$ .

Example: **Logarithmic barrier function** for  $c_i(\underline{x}) \geq 0$ :



*still a "penalty" to enter, where we try to minimize the distance to  $\partial X$*



*it does not have to go to zero in  $x^\circ$ , it just needs to be continuous*



Idea: Add to objective function the barrier terms associated to the constraints and solve a sequence of

**Definition:** The logarithmic barrier problem associated to problem (1) is

$$\min_{\underline{x} \in \mathbb{R}^n} P(\underline{x}, \mu) = f(\underline{x}) - \mu \sum_{i \in I} \ln c_i(\underline{x}), \quad (2)$$

with barrier parameter  $\mu > 0$ .

N.B.: When  $\mu \rightarrow 0$  the barrier term becomes negligible.

*since  $a \cdot b \rightarrow 0$   
 $a \rightarrow 0$   
 $b \rightarrow 0$*

*this constraint leads to  $le \geq 0$   
 $\rightarrow -\infty$  if we move towards not respecting it  
 while the (-) we recover to end hence the revolution  
 $f(\underline{x}) + \mu \sum [-\ln c_i(\underline{x})]$*

We consider  $\{\mu_k\}$  with  $\lim_{k \rightarrow \infty} \mu_k = 0$ , start from  $\underline{x}_0 \in X^\circ$  and, for each  $k$ , determine an approximate minimizer  $\underline{x}_k$  of  $P(\underline{x}, \mu_k)$  with an unconstrained optimization method.

Example 1:

$$\begin{array}{ll} \min & x \\ \text{s.t.} & x \geq 0 \\ & 1 - x \geq 0 \end{array}$$

with optimal solution  $x^* = 0$  and logarithmic barrier problem:

$$\min_{x \in \mathbb{R}} P(x, \mu) = x - \mu \ln x - \mu \ln(1 - x).$$

Compare  $P(x, \mu)$  for values of  $\mu$  from 1 to 0.01.

See J. Nocedal, S. Wright, Numerical Optimization, Springer, 1999, p. 499-500.

Example 2:

$$\begin{aligned} \min \quad & (x_1 + 0.5)^2 + (x_2 - 0.5)^2 \\ \text{s.t.} \quad & x_1 \in [0, 1] \\ & x_2 \in [0, 1] \end{aligned}$$

with optimal solution  $\underline{x}^* = (0, 0.5)^t$  and logarithmic barrier problem:

Compare contours of  $P(\underline{x}, \mu)$  for values of  $\mu$  from 1 to 0.01.

For  $\mu = 0.01$ , around  $\underline{x}^*$  (more elongated and less elliptical) indicate possible numerical problems.

See J. Nocedal, S. Wright, Numerical Optimization, Springer, 1999, p. 500-502.

## General scheme

- 0) Choose  $\varepsilon > 0$ ,  $\mu_0 > 0$ , tolerances  $\{\tau_k\}_{k \geq 0}$  with  $\tau_k > 0$  and  $\lim_{k \rightarrow \infty} \tau_k = 0$ , initial point  $\underline{x}_0^s$ . Set  $k := 0$ .
- 1) Determine an approximate minimizer  $\underline{x}_k$  of  $P(\underline{x}, \mu_k)$  starting from  $\underline{x}_k^s$  and terminate when  $\|\nabla P(\underline{x}, \mu_k)\| \leq \tau_k$ .
- 2) If overall termination condition is satisfied (e.g.  $|f(\underline{x}_{k-1}) - f(\underline{x}_k)| < \varepsilon$ )  
Then Stop  
Else select  $\mu_{k+1} \in (0, \mu_k)$  and  $\underline{x}_{k+1}^s$ , set  $k := k + 1$  and Goto 1)

Since  $\underline{x}_0 \in X^\circ$  the sequence  $\{\underline{x}_k\}$  remains in  $X^\circ$ , the algorithm is an

interior point method.

*we will never get a point on the  $\partial X$ , on extreme point*

Important connection between a minimum of  $P(\underline{x}, \mu)$ , denoted  $\underline{x}(\mu)$ , and a point  $(\underline{x}^*, \underline{u}^*)$  satisfying the KKT conditions of problem (1), namely

$$\nabla_{\underline{x}} L(\underline{x}, \underline{u}) = \nabla f(\underline{x}) - \sum_{i=1}^m u_i \nabla c_i(\underline{x}) = \underline{0} \quad (3)$$

$$c_i(\underline{x}) \geq 0 \quad \forall i \in I \quad (4)$$

$$u_i c_i(\underline{x}) = 0 \quad \forall i \in I \quad (5)$$

$$u_i \geq 0 \quad \forall i \in I. \quad (6)$$

In a minimizer  $\underline{x}(\mu)$  of  $P(\underline{x}, \mu)$ , we have

$$\nabla_{\underline{x}} P(\underline{x}, \mu) = \nabla f(\underline{x}) - \sum_{i=1}^m \underbrace{\frac{\mu}{c_i(\underline{x})}}_{\text{KKT multipliers}} \nabla c_i(\underline{x}) = \underline{0}. \quad (7)$$

By defining the estimates of the multipliers

$$u_i(\mu) := \frac{\mu}{c_i(\underline{x}(\mu))} \quad \text{with } i = 1, \dots, m, \quad (8)$$

(7) can be rewritten as

$$\nabla f(\underline{x}) - \sum_{i=1}^m u_i(\mu) \nabla c_i(\underline{x}) = \underline{0} \quad (9)$$

which is equivalent to (3).

Observation: For  $\mu > 0$  the KKT conditions (3)-(6) hold except (5) because

$$u_i(\mu)c_i(\underline{x}(\mu)) = \mu \quad \text{for } i = 1, \dots, m.$$

When  $\mu \rightarrow 0$ , a minimizer  $\underline{x}(\mu)$  of  $P(\underline{x}, \mu)$  and the associated estimate

$$u_i(\mu) := \frac{\mu}{c_i(\underline{x}(\mu))} \quad \text{with } i = 1, \dots, m$$

tend to progressively satisfy the KKT conditions of problem (1).

Thus we generate points on the so-called **central path**

$$\{(\underline{x}(\mu), \underline{u}(\mu)) : \mu > 0\}$$

defined by (8).

## Theorem:

Suppose that  $X^\circ \neq \emptyset$  and  $\underline{x}^*$  is a local minimum of (1) at which the KKT conditions are satisfied for some  $\underline{u}^*$ .

Moreover, suppose that

- gradients of the active constraints at  $\underline{x}^*$  are linearly independent,
- strict complementarity conditions are satisfied at  $\underline{x}^*$  ( $\forall i \in I$  exactly one of  $c_i(\underline{x}^*)$  or  $u_i^*$  is equal to 0),
- 2nd order sufficient conditions are satisfied at  $(\underline{x}^*, \underline{u}^*)$ .

Then

- $\exists$  unique continuously differentiable vector function  $\underline{x}(\mu)$  s.t.  $\lim_{\mu \rightarrow 0_+} \underline{x}(\mu) = \underline{x}^*$ .  
For all sufficiently small  $\mu$ ,  $\underline{x}(\mu)$  is a local minimum of  $P(\underline{x}, \mu)$  in some neighborhood of  $\underline{x}^*$ .
- For  $\underline{x}(\mu)$  in (i), the Lagrange multiplier estimates  $\underline{u}(\mu)$  defined by
$$u_i(\mu) = \mu / c_i(\underline{x}(\mu)) \quad i = 1, \dots, m,$$
converge to  $\underline{u}^*$  when  $\mu \rightarrow 0_+$ .
- $\nabla_{\underline{x}\underline{x}}^2 P(\underline{x}, \mu)$  is positive definite for all sufficiently small  $\mu$ .

If also equality constraints, one may include quadratic penalty terms (combined log-barrier/quadratic penalty function problem).

**Sixth computer lab:** application of the logarithmic barrier method to LP.

An **interior point method** for LP

$$\begin{aligned} \min \quad & \underline{c}^t \underline{x} \\ \text{s.t.} \quad & A\underline{x} = \underline{b} \\ & \underline{x} \geq \underline{0} \end{aligned} \tag{10}$$

is obtained by applying the barrier method to constraints (11) and by adapting the Newton method to account for (10).

Unlike for Simplex method, such interior point method for LP can be proved to be a polynomial time algorithm.



## 5.9 Introduction to sequential quadratic programming

Generic NLP:


$$(P) \quad \begin{array}{ll} \min & f(\underline{x}) \\ \text{s.t.} & g_i(\underline{x}) \leq 0 \quad i \in I = \{1, \dots, m\} \\ & h_l(\underline{x}) = 0 \quad l \in E = \{1, \dots, p\} \\ & \underline{x} \in \mathbb{R}^n \end{array}$$

where  $f$ ,  $g_i$ 's and  $h_l$ 's are of class  $C^2$ .

Idea: Extend the Newton method to nonlinearly constrained problems.

Given a current  $\underline{x}_k$ , we could try to determine an improving direction  $\underline{d}_k$  by solving the quadratic approximation of (P):

$$(QA_k) \quad \begin{array}{ll} \min & \frac{1}{2} \underline{d}^t \nabla^2 f(\underline{x}_k) \underline{d} + \nabla^t f(\underline{x}_k) \underline{d} + f(\underline{x}_k) \\ \text{s.t.} & \frac{1}{2} \underline{d}^t \nabla^2 g_i(\underline{x}_k) \underline{d} + \nabla^t g_i(\underline{x}_k) \underline{d} + g_i(\underline{x}_k) \leq 0 \quad i \in I = \{1, \dots, m\} \\ & \frac{1}{2} \underline{d}^t \nabla^2 h_l(\underline{x}_k) \underline{d} + \nabla^t h_l(\underline{x}_k) \underline{d} + h_l(\underline{x}_k) = 0 \quad l \in E = \{1, \dots, p\} \end{array}$$

*quadratic approx of  $f(\cdot)$  at  $\underline{x}_k$*  

*quadratic approx of the  $g_i(\cdot)$  and  $h_l(\cdot)$*

but difficult because of the quadratic constraints.

### Observation:

If  $(\underline{d}^*, \underline{\eta}^*, \underline{\rho}^*)$  is a stationary point of the Lagrange function associated to  $(QA_k)$  it is also a stationary point of the Lagrange function associated to the Quadratic Program:

$$\begin{aligned} \min \quad & \frac{1}{2} \underline{d}^t \nabla_{\underline{x}\underline{x}}^2 L(\underline{x}_k, \underline{\eta}^*, \underline{\rho}^*) \underline{d} + \nabla^t f(\underline{x}_k) \underline{d} + f(\underline{x}_k) \\ (QPA_k) \quad \text{s.t.} \quad & \nabla^t g_i(\underline{x}_k) \underline{d} + g_i(\underline{x}_k) \leq 0 & i \in I = \{1, \dots, m\} \\ & \nabla^t h_l(\underline{x}_k) \underline{d} + h_l(\underline{x}_k) = 0 & l \in E = \{1, \dots, p\} \end{aligned}$$

All constraints are linear (approximations).

To obtain a good approximation of  $(P)$  via Quadratic Programs, the objective function must include not only a quadratic model of  $f$  but also 2nd order information of the  $g_i$ 's.

## General scheme

Let  $\underline{x}_0$ ,  $\underline{u}_0$  and  $\underline{v}_0$  be estimates of a solution of  $(P)$  and of the corresponding multipliers.

Iteration  $k$ :

Given  $(\underline{x}_k, \underline{u}_k, \underline{v}_k)$  determine  $\underline{d}_k$  and the corresponding multipliers  $(\underline{\eta}_k, \underline{\rho}_k)$  of the Quadratic Program:

$$\begin{aligned} \min \quad & \frac{1}{2} \underline{d}^t \nabla_{\underline{x}\underline{x}}^2 L(\underline{x}_k, \underline{u}_k, \underline{v}_k) \underline{d} + \nabla^t f(\underline{x}_k) \underline{d} \\ (QP_k) \quad \text{s.t.} \quad & \nabla^t g_i(\underline{x}_k) \underline{d} + g_i(\underline{x}_k) \leq 0 & i \in I = \{1, \dots, m\} \\ & \nabla^t h_l(\underline{x}_k) \underline{d} + h_l(\underline{x}_k) = 0 & l \in E = \{1, \dots, p\} \end{aligned}$$

Set  $\underline{x}_{k+1} := \underline{x}_k + \underline{d}_k$ ,  $\underline{u}_{k+1} := \underline{\eta}_k$  and  $\underline{v}_{k+1} := \underline{\rho}_k$

Although  $(QP_k)$  derives from  $(QPA_k)$  by substituting the optimal multipliers with the current estimates, it can be proved that:

- feasible region of the subproblem  $(QP_k)$  is a *linear approximation* of that of the original problem,
- Lagrange function  $L_Q(\underline{d}, \underline{\eta}, \underline{\rho})$  of  $(QP_k)$  is a *quadratic approximation* of that of  $(P)$ .

An iteration of the Sequential Quadratic Programming method (SQP) is equivalent to:

- carry out one iteration of the Newton method for the Lagrange function,
- enforce feasibility with respect to the linearization of the feasible region.

The SQP method is well defined:

**Proposition:**

$(\underline{x}^*, \underline{u}^*, \underline{v}^*)$  is a KKT point of  $(P)$  if and only if  $(\underline{d}^*, \underline{\eta}^*, \underline{\rho}^*) = (\underline{0}, \underline{u}^*, \underline{v}^*)$  is a KKT point of  $(QP_k)$ .

Convergence properties similar to those for Newton method:

Quadratic local convergence if

- (i) Hessian matrices of the objective function and constraints are Lipschitz continuous,
- (ii) constraint qualification assumption is satisfied,
- (iii) 2nd order sufficient optimality conditions and strict complementarity conditions are satisfied.

To guarantee global convergence:

- 1-D search that minimizes an appropriate merit function such as

$$M(\underline{x}; \mu) = f(\underline{x}) + \frac{1}{2\mu} \left( \sum_{i=1}^m \max\{0, g_i(\underline{x})\} + \sum_{l=1}^p |h_l(\underline{x})| \right)$$

- or trust region based approach.

Quasi-Newton versions (without 2nd order derivatives) have also been investigated.

Several SQP codes are available (SQP, NPSOL, SNOPT, Matlab,...).

# Subgradient method

Consider  $\min_{\underline{x} \in \mathbb{R}^n} f(\underline{x})$  with  $f$  convex.

Start from an arbitrary  $\underline{x}_0$ .

At  $k$ -th iteration: consider  $\underline{\gamma}_k \in \partial f(\underline{x}_k)$  and set

$$\underline{x}_{k+1} := \underline{x}_k - \alpha_k \underline{\gamma}_k$$

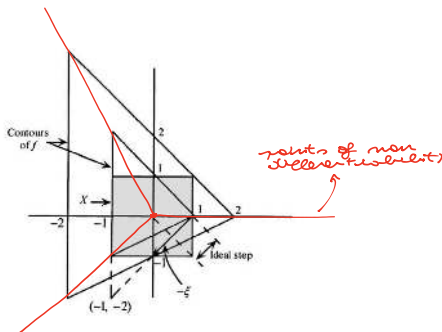
with  $\alpha_k > 0$

*about the  
choice of  $\alpha_k$  ↪*

**Observation:** No 1-D search (optimization) because for nondifferentiable functions a subgradient  $\underline{\gamma} \in \partial f(\underline{x})$  is not necessarily a descent direction!

Example:  $\min_{-1 \leq x_1, x_2 \leq 1} f(x_1, x_2)$  with  $f(x_1, x_2) = \max\{-x_1, x_1 + x_2, x_1 - 2x_2\}$

Level curves in black, points of nondifferentiability  $(t, 0)$ ,  $(-t, 2t)$  and  $(-t, -t)$  for  $t \geq 0$ , global minimum  $\underline{x}^* = (0, 0)$ .



At  $\underline{x}_k = (1, 0)^t$  consider  $\underline{\gamma}_k = (1, 1) \in \partial f(\underline{x}_k)$ ,  $f(\underline{x})$  increases along  $\{\underline{x} \in \mathbb{R}^2 : \underline{x} = \underline{x}_k - \alpha_k \underline{\gamma}_k, \alpha_k \geq 0\}$  but if  $\alpha_k$  is sufficiently small then  $\underline{x}_{k+1} = \underline{x}_k - \alpha_k \underline{\gamma}_k$  is closer to  $\underline{x}^*$ .  
*good idea to choose*

From Chapter 8, Bazaraa et al., Nonlinear Programming, Wiley, 2006, p. 436-437

## Theorem:

If  $f$  is convex,  $\lim_{\|\underline{x}\| \rightarrow \infty} f(\underline{x}) = +\infty$ ,  $\lim_{k \rightarrow \infty} \alpha_k = 0$  (small steps) and  $\sum_{k=0}^{\infty} \alpha_k = \infty$  (but not too small), the subgradient method terminates after a finite number of iterations with an optimal solution  $\underline{x}^*$  or infinite sequence  $\{\underline{x}_k\}$  admits a subsequence converging to  $\underline{x}^*$ .

## Stepsize:

In practice  $\{\alpha_k\}$  as above (e.g.,  $\alpha_k = 1/k$ ) are too slow.

An option:  $\alpha_k = \alpha_0 \rho^k$  for a given  $\rho < 1$ . A more popular one (min problems):

$$\alpha_k = \varepsilon_k \frac{f(\underline{x}_k) - \hat{f}}{\|\underline{\gamma}_k\|^2},$$

where  $0 < \varepsilon_k < 2$  and  $\hat{f}$  is either the optimal value  $f(\underline{x}^*)$  or an estimate.

Stopping criterion: prescribed maximum number of iterations  
(even if  $\underline{0} \in \partial f(\underline{x}_k)$  it may non be considered at  $\underline{x}_k$ ).

Need to store the best solution  $\underline{x}_k$  found. *so we may take a "wrong" subgradient*

Simple extension for bounds (projections).