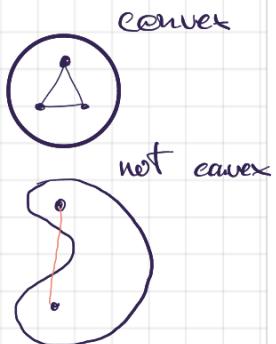


## Problem definition - Convexity

Def A set  $C \subset \mathbb{R}^n$  is convex if  $\forall x, y$  the segment joining the two points is contained in  $C$

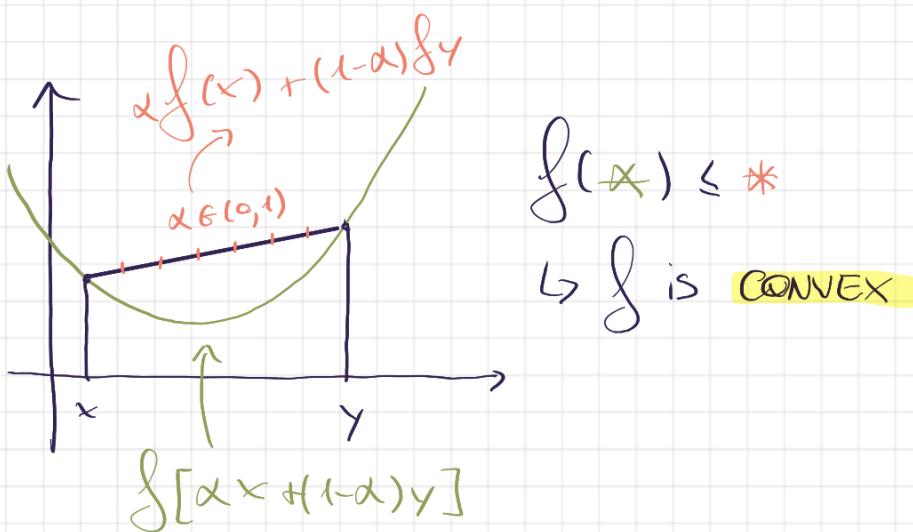


$$\alpha x + (1-\alpha)y \in C \quad \forall \alpha \in [0,1]$$

Def Given a convex set  $C \subset \mathbb{R}^n$ , a function  $f: C \rightarrow \mathbb{R}$  is convex if

$$[f(\alpha x + (1-\alpha)y) \leq \alpha f(x) + (1-\alpha)f(y)]$$

$$\hookrightarrow \forall x, y \in C, \forall \alpha \in [0,1]$$



$f(x) \leq *$

$\hookrightarrow f$  is **convex**

- PROPERTY - - - - -

If "f" is convex and the set of feasible solutions (set of points satisfying the constraints) is a convex set, a local minimum is also a global minimum

Sel4

TODO

# - Gradient methods - - - - -

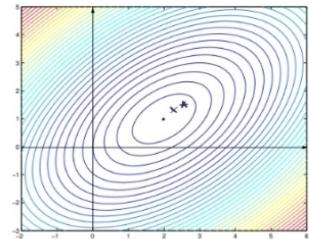
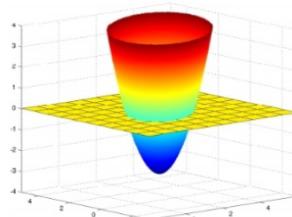
↳ if  $A$  is a symmetric positive definite (SPD) matrix, then the solution of  $Ax = b$  is the UNIQUE MINIMUM of:

$$J: \mathbb{R}^n \rightarrow \mathbb{R} \quad J(x) = \frac{1}{2} x^T A x - x^T b \quad (\approx \|Ax - b\|)$$

! Since  $A$  is SPD and  $A$  is the Hessian of  $J(x)$ , then  $J(x)$  is strictly convex

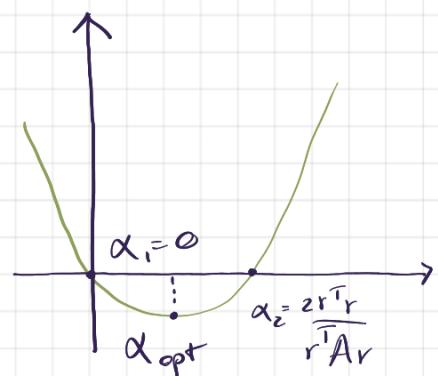
## STEEPEST DESCENT

↳ IDEA: build an iterative method by moving, at each step, along the direction  $-\nabla J(x_k)$



$$\begin{aligned} &\text{↳ direction}(x) = -\nabla J(x) = b - Ax = r(x) \\ &\text{↳ step}(x) = \alpha_{\text{opt}} = \frac{r^T r}{r^T A r} \end{aligned}$$

↳ once we have a direction, knowing  $J(x)$  is a paraboloid, on the plane given by  $-\nabla J(x)$  will result a parabola: get the minimum for optimal step



## ↳ IMPLEMENTATION:

matrix mult is the main cost  $\rightarrow$  if A sparse  
 $(n^2)$  it's better

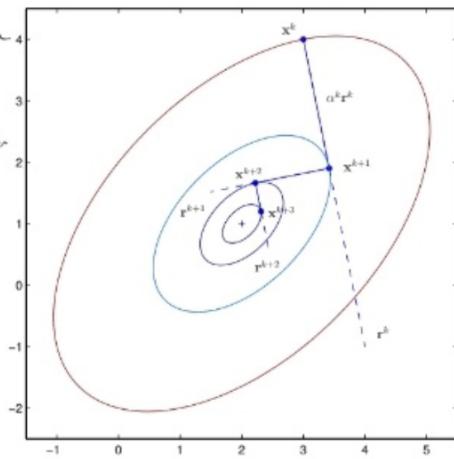
Given  $x_0$ , for  $k \geq 0$

- Compute  $r_k = b - Ax_k \rightarrow$  direction
- Compute  $\alpha_k = \frac{r_k^T r_k}{r_k^T A r_k} \rightarrow$  step
- Update  $x_{k+1} = x_k + \alpha_k r_k \rightarrow$  update

Note that 2 matrix-vector products are performed at each iteration. We can do better:

Given  $x_0, r_0 = b - Ax_0$ , for  $k \geq 0$

- Compute  $z_k = Ar_k$
- Compute  $\alpha_k = \frac{r_k^T r_k}{r_k^T z_k}$
- Update  $x_{k+1} = x_k + \alpha_k r_k$
- Update  $r_{k+1} = r_k - \alpha_k z_k$



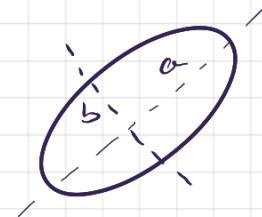
↳ BETTER COMPUTATIONAL COST (only 1 matrix-multiplication per iteration)



34/120

## ↳ CONVERGENCE PROPERTIES:

$$\kappa_2(A) = \frac{\max(\lambda(A))}{\min(\lambda(A))} \stackrel{\text{SPD matrix}}{=} \frac{a}{b}$$



↳ bigger  $\kappa_2(A)$  = "more bounces" = harder convergence

~~TO DO:~~

Energy norm?

## CONJUGATE GRADIENT METHODS

Def A vector  $p_n$  is a **descent vector**  
 if  $\nabla J(x_n)^T p_n < 0$   
 for a short enough step

↳ **IDEA**: There may be better  
 descent vectors to improve convergence

### • In Steepest Descent

↳ all new directions were parallel to the ones before defined in the linear space of the new gradient

$$x_{k+1} = \arg \min J(x_k + v)$$

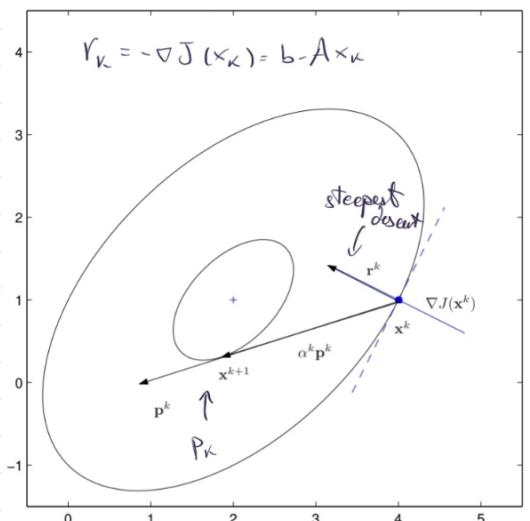
$$v \in L(r_k)$$

$$d_k - r_k \rightarrow d_{k+1}^T d_k = 0$$

! If vectors  $d_0, \dots, d_k$  form a  $A$ -conjugate set, they are also independent

↳ in order to compute  $[d_{k+1} = -\nabla J(x_{k+1}) + \beta_k d_k]$

it needs to compute  $\beta_k$



### • in Conjugate Descent

↳ all new directions belong to the linear space of every direction took before

$$x_{k+1} = \arg \min J(x_k + v)$$

$$v \in L(d_0, \dots, d_k)$$

↳  $d_{k+1}$  and  $d_k$  are conjugate to  $A$

$$d_k = r_k + \beta_k d_{k-1} \rightarrow d_{k+1}^T A d_k = 0$$

↳ it can be proven that  $\beta_k$  can be computed in such way that  $d_{k+1}$  is A-conjugate to all previous directions:

$$\beta_k = -\frac{(d_k)^T A r_{k+1}}{d_k^T A d_k}$$

## ↳ IMPLEMENTATION

Given  $x^0$

Compute  $r^0 = b - Ax^0$

Set  $d^0 = r^0$

For  $k \geq 0$

$$z_k = Ad_k$$

$$\alpha_k = \frac{r_k^T d_k}{d_k^T z_k} \quad \text{only computed once}$$

$$x_{k+1} = x_k + \alpha_k d_k$$

$$r_{k+1} = r_k - \alpha_k z_k$$

$$\beta_{k+1} = -\frac{(r_{k+1})^T r_{k+1}}{r_k^T r_k}$$

$$d_{k+1} = r_{k+1} + \beta_{k+1} d_k$$

→ one line in MATLAB:

`x = pcg (A, b, tol, max_iter, pre-conditioner)`

## ↳ CONVERGENCE PROPERTY

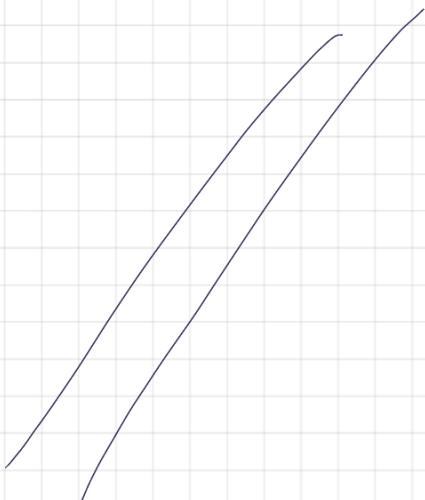
- Finite Termination property: The method WILL stop after no more than "n" steps (as in  $A \in \mathbb{R}^{n,n}$ )

↳ if "n" is very large it is not useful

↳ very depending on  $K_2(A) \rightarrow$  PRECONDITIONING

## PRECONDITIONING

- ↳ look for  $P \in \mathbb{R}^{n,n}$  invertible such that  $\kappa_2(P^{-1}A) \ll \kappa_2(A)$   
in order to:  $Ax = b \rightarrow P^{-1}Ax = P^{-1}b$
- ↳ best conditioner  $P = A$  BUT COST TOO HIGH
- ↳ IDEA: build  $P = \text{approx. of } A$
- ↳ if  $A$  is SPD (symmetric positive definite)  
one can compute an INCOMPLETE CHOLESKY FACTORIZATION



## - Unconstrained optimization (general case) - - -

↳  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  smooth (first order derivable) nonlinear function

we want  $x^*$  such as  $f(x^*) = \min f(x)$

Def FIRST ORDER NECESSARY CONDITIONS

↳ if  $x^*$  is minimum, then  $\nabla f(x^*) = 0$

Def SECOND ORDER NECESSARY CONDITIONS

↳ if  $x^*$  is minimum, then  $\nabla^2 f(x^*)$  is positive semidefinite

? //

↳  $x^*$  minimum  $\Rightarrow \begin{cases} \nabla f(x^*) = 0 \\ \nabla^2 f(x^*) \text{ semidefinite positive} \end{cases}$

Def SECOND ORDER SUFFICIENT CONDITIONS

↳  $\begin{cases} \nabla f(x^*) = 0 \\ \nabla^2 f(x^*) \text{ positive definite} \end{cases} \Rightarrow x^* \text{ strict local minimum}$

! ----- !  
! If  $f$  is convex any local minimum is a  
! global minimum !

! ----- !  
! If  $f$  is also differentiable, then any stationary  
! point is a global minimum of  $f$ . !

# GRADIENT METHODS

8/55

## STEEPEST DESCENT METHODS

↪ we want to find direction  $p_k$  and step length  $\alpha$

$$x_{k+1} = x_k + \alpha p_k \text{ and let } \phi(\alpha) = f(x_k + \alpha p_k)$$

↪ once you decide  $p_k$  you have only one dimension to work in

NEWTON METHOD

CHOOSING STEP LENGTH  $\alpha$  - BACKTRACKING ARMJO

NON-LINEAR CONJUGATE GRADIENT METHOD

FR / PR

## VARIANTS OF NEWTON METHOD

↳ the computation of the Newton step may be critical in several respects:

- ① assembling  $\nabla^2 f(x_k)$  is costly
- ② solving the linear system ( $m(p) \rightarrow \nabla m(p) = 0 \rightarrow p$ ) is costly
- ③  $\nabla^2 f(x_k)$  may not be SPD ( $p \rightarrow$  not a descent direction)

### INEXACT NEWTON METHOD

## FINITE DIFFERENCE APPROXIMATION

↳ problems:

- computing derivatives ( $\nabla f$ ,  $\nabla^2 f$ ) is very costly
- if you don't know  $f$  analytically you can't compute derivatives

↳ need approximation rules for  $f$  derivatives  
which only calls for function evaluation  
at some given points.

### • COMPUTING DERIVATIVE $f'(\bar{x})$

↳ approximate  $f$  around  $\bar{x}$  with 1st order Taylor

$$f(x) = f(\bar{x}) + f'(\bar{x})(x - \bar{x}) + \frac{1}{2} f''(\xi)(x - \bar{x})$$

↑ Lagrange remainder

↳ choose  $h > 0$  and take  $x = \bar{x} + h$

$$f(\bar{x} + h) = f(\bar{x}) + f'(\bar{x})h + f''(\xi)h^2$$

$$\text{so } f'(x) = \frac{f(\bar{x} + h) - f(\bar{x})}{h} - \frac{1}{2} f''(\xi)h$$

then 
$$\left[ f'(x) \approx \frac{f(\bar{x} + h) - f(\bar{x})}{h} \right] \rightarrow$$
 you can approximate it even if you don't know  $f$  analytically

with error  $\frac{1}{2} f''(\xi)h \rightarrow O(h)$

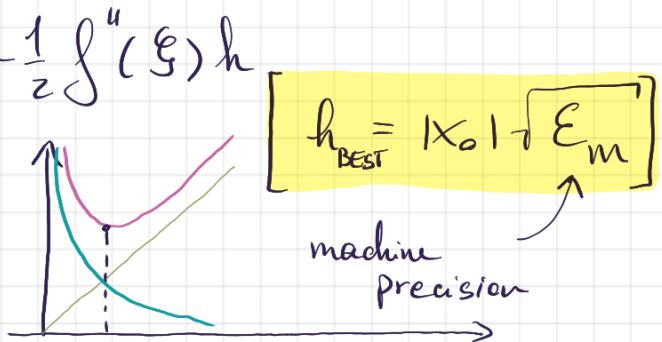
! Ideally  $h$  as small as possible

↳ problem with small  $h$ :

- smaller analytical error  $\sim \frac{1}{2} f''(\xi) h$
- bigger computer error

$$\text{tot\_error} = C_1 h + C_2 \frac{1}{h}$$

analytical  $E \nearrow$  numerical error



↳ one can either go for

$$x = \bar{x} + h \quad (\text{forward finite differences})$$

$$x = \bar{x} - h \quad (\text{backward finite differences})$$

↳ or getting a (better) centred finite difference:

↳ take both FW and BW finite differences of second order:

$$f(x_0 + h) = f(x_0) + f'(x_0)h + \frac{1}{2} f''(x_0)h^2 + \frac{1}{6} f'''(\xi_+)h^3$$

$$f(x_0 - h) = f(x_0) - f'(x_0)h + \frac{1}{2} f''(x_0)h^2 - \frac{1}{6} f'''(\xi_-)h^3$$

↳ subtracting:

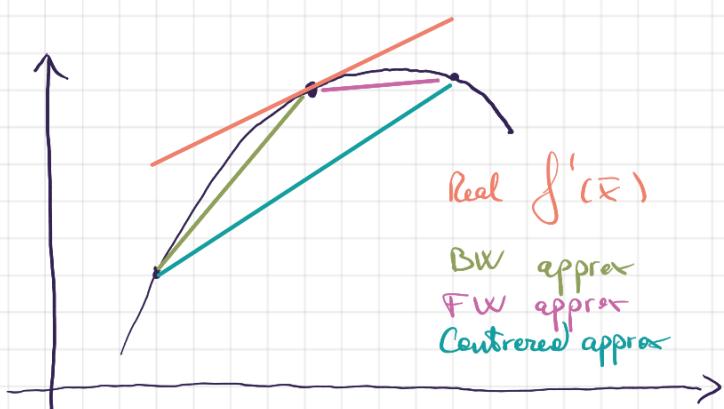
$$f'(x_0) = \frac{f(x_0 + h) - f(x_0 - h)}{2h} + \frac{1}{6}(f'''(\xi_+) + f'''(\xi_-))h^2$$

error  $\mathcal{O}(h^2)$

having  $\xi_+ \in [\bar{x}, \bar{x} + h]$

$\xi_- \in [\bar{x} - h, \bar{x}]$

→ Cost: higher numbers of function evaluation  
(twice more)



- COMPUTING SECOND DERIVATIVE  $f''(\bar{x})$

↳ approximate 3<sup>rd</sup> order Taylor expansion

$$f(x_0 + h) = f(x_0) + f'(x_0)h + \frac{1}{2}f''(x_0)h^2 + \frac{1}{6}f'''(x_0)h^3 + \frac{1}{24}f^{(4)}(\xi_+)h^4$$

$$f(x_0 - h) = f(x_0) - f'(x_0)h + \frac{1}{2}f''(x_0)h^2 - \frac{1}{6}f'''(x_0)h^3 + \frac{1}{24}f^{(4)}(\xi_-)h^4$$

↳ summing:

$$\left[ f''(x_0) = \frac{f(x_0 + h) - 2f(x_0) + f(x_0 - h)}{h^2} \right] + \frac{1}{24}(f'''(\xi_+) + f'''(\xi_-))h^2$$

error  $O(h^2)$

- COMPUTING GRADIENT  $\nabla f(\bar{x})$

↳ given  $f: \mathbb{R}^n \rightarrow \mathbb{R}$   $f(x)$   $x = x_1, \dots, x_n$

↳ Knowing  $\nabla f(x) = \left( \frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right)$

↳ you have to compute 'n' derivatives:

$$\frac{\partial f}{\partial x_i}(x) = \frac{\partial f}{\partial x_i}(x_1, \dots, x_n) \simeq \frac{f(x_1, \dots, x_i + h, \dots, x_n) - f(x_1, \dots, x_n)}{h}$$

↳ now, one may generalize:

$$\begin{aligned} (x_1, \dots, x_i + h, \dots, x_n) &= (x_1, \dots, x_n) + h(0, \dots, \underbrace{1}_{\text{pos. } i}, \dots, 0) \\ &= (x_1, \dots, x_n) + h\mathbf{e}_i = x + h\mathbf{e}_i \end{aligned}$$

where  $\mathbf{e}_i$  is the i-th vector of the canonical basis in  $\mathbb{R}^n$

↳ so in compact form:  $\left[ \frac{\partial f}{\partial x_i}(x) \simeq \frac{f(x + h\mathbf{e}_i) - f(x)}{h} \right]$

(one may use also backward or centered variants)

↳ Costs: FW/BW; n+1 evaluations

Centered; 2n evaluations

↳ choice of  $h = \sqrt{\epsilon_m}$

## • COMPUTING JACOBIAN

↳ given  $F: \mathbb{R}^n \rightarrow \mathbb{R}^n$

$$F'(x) \in \mathbb{R}^{n \times n}$$

$$\left[ F'(x) \right]_{ij} = \frac{\partial F_i}{\partial x_j}(x) \rightarrow F'(x) = \begin{pmatrix} \frac{\partial F_1(x)}{\partial x_1} \\ \frac{\partial F_2(x)}{\partial x_2} \\ \vdots \\ \frac{\partial F_n(x)}{\partial x_n} \end{pmatrix}$$

↳ naive idea: repeat construction of  $\nabla f$  to  $F_i$   
(quite costly)

↳ we may work on columns:

$$F'(x) = \left( \frac{\partial F(x)}{\partial x_1} \mid \cdots \mid \frac{\partial F(x)}{\partial x_n} \right)$$

$$\hookrightarrow \text{now } \frac{\partial F}{\partial x_i}(x) \simeq \frac{F(x + h e_i) - F(x)}{h}$$

↳ cost:  $n+1$  evaluations

(only one per column because we  
always have the same of one one column)

↳ the cost is still high for a high ' $n$ '

↳ we may optimize it:  

- matrix free implementation
- sparse Jacobians

## • COMPUTING JACOBIAN: MATRIX-FREE IMPLEMENTATION

↳ if  $n$  is large: costly to build and to store

↳ generally you don't need all  $F'(x)$

↳ you need to know  $F'(x) \cdot d$  so you can replace

given  $d$ , gives the  $F'(x)$  with

a function that, same result as  $F'(x) \cdot d$

↳ for simplicity  $\nabla f \cdot d = \nabla f$

↳ recall the definition of the directional derivative  
of  $f$  in direction ' $d$ '

$$D(f(x); d) := \lim_{h \rightarrow 0} \frac{f(x + hd) - f(x)}{h} = \nabla f(x)^T d \rightarrow \nabla f(x) \text{ is column vector}$$

↳ so we may approximate  $F'(x) \cdot d \approx \frac{F(x+hd) - F(x)}{h}$

↳ computational gain: ① No need for matrix allocation

② Only two function evaluations

! The ② is only apparent since you may need  
to frequently change ' $d$ ' (anyway  $k \ll n$ )

## • COMPUTE JACOBIANS: SPARSE JACOBIANS

↳ consider  $F(x)$  like :

$$F(x) = \begin{pmatrix} F_1(x_1, x_2) \\ F_2(x_1, x_2, x_3) \\ F_3(x_2, x_3, x_4) \\ \vdots \\ F_n(x_{n-1}, x_n) \end{pmatrix}, \text{ then } F'(x) = \begin{pmatrix} \times & \times & & & \\ \times & \times & \times & & \\ & \times & \times & \times & \\ & & \times & \times & \times \\ & & & \times & \times \end{pmatrix}$$

tridiagonal

↳ if we approximate  $F'(x)$   $\frac{\partial F}{\partial x_i}(x) \simeq \frac{F(x + he_i) - F(x)}{h}$   
column-wise:

↳ since  $F_3$  (and so on) does not depend on  $x$ ,  
then  $F_3(x + he_i) = F_3(x)$ :

$$\frac{F(x + he_1) - F(x)}{h} = \begin{pmatrix} \frac{F_1(x + he_1) - F_1(x)}{h} \\ \frac{F_2(x + he_1) - F_2(x)}{h} \\ \frac{F_3(x + he_1) - F_3(x)}{h} \\ \vdots \end{pmatrix} \simeq \begin{pmatrix} \frac{\partial F_1}{\partial x_1}(x) \\ \frac{\partial F_2}{\partial x_1}(x) \\ 0 \\ \vdots \end{pmatrix} \rightarrow \text{so it's useless to compute zero components}$$

↳ we can compute larger perturbations for another column that has different (mutually exclusive) dependencies:

↳ for this particular example we need to use only 3 perturbations:

$$P_1 = h (e_1 + e_n + e_2 + \dots)$$

$$P_2 = h (e_2 + e_5 + e_8 + \dots)$$

$$P_3 = h (e_3 + e_6 + e_9 + \dots)$$

$$\frac{F(x + h(e_1 + e_4)) - F(x)}{h} = \begin{pmatrix} \frac{F_1(x + he_1) - F_1(x)}{h} \\ \frac{F_2(x + he_1) - F_2(x)}{h} \\ \frac{F_3(x + he_4) - F_3(x)}{h} \\ \frac{F_4(x + he_4) - F_4(x)}{h} \\ \frac{F_5(x + he_4) - F_5(x)}{h} \\ 0 \end{pmatrix} \simeq \begin{pmatrix} \frac{\partial F_1}{\partial x_1}(x) \\ \frac{\partial F_2}{\partial x_1}(x) \\ \frac{\partial F_3}{\partial x_4}(x) \\ \frac{\partial F_4}{\partial x_4}(x) \\ \frac{\partial F_5}{\partial x_4}(x) \\ 0 \end{pmatrix}$$

$$F'(x) = \begin{pmatrix} \times & \times & & & & \\ \times & \times & \times & & & \\ & \times & \times & \times & & \\ & & \times & \times & \times & \\ & & & \times & \times & \\ & & & & \times & \times \end{pmatrix}$$

→ because we have a tridiagonal  $F'(x)$

↳ This can be generalized to arbitrary functions

$F : \mathbb{R}^n \rightarrow \mathbb{R}^m$  by means of graphs and graph colouring

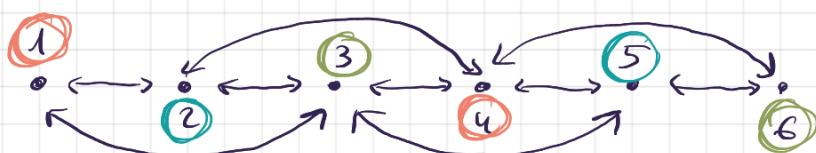
↳ build column incidence graphs:

↳ graph with "n" nodes (number of columns)

↳ connect nodes "i" and "j" if the two columns have at least one dependency in common

↳ assign color to the nodes: nodes not connected by arcs have the same color

example:  $F_i(x_{i-1}, x_i, x_{i+1}) \quad n=6$



↳ the number of colors correspond to the number of perturbations needed

↳ if  $i_1, i_2$  have the same color

We use perturbation  $p = h(e_{i_1} + e_{i_2})$

- COMPUTING HESSIAN  $\nabla^2 f(x)$

↳ note that Hessian of  $f(x)$  is the Jacobian of  $\nabla f(x)$

Taylor's expansion applied to  $\nabla f$ :

so:

$$\nabla f(x + p) = \nabla f(x) + \nabla^2 f(x)p + \mathcal{O}(\|p\|^2)$$

Taking  $p = he_i$ , we get

$$\nabla^2 f(x)e_i = \frac{(\nabla f)}{\partial x_i}(x) \simeq \frac{\nabla f(x + he_i) - \nabla f(x)}{h}$$

with error  $\mathcal{O}(h)$

↳ if  $\nabla f$  is not available you can approximate through centered finite differences:

$$\frac{\partial f}{\partial x_i \partial x_j}(x) = \frac{f(x + he_i + he_j) - f(x + he_i) - f(x + he_j) + f(x)}{h^2} + \mathcal{O}(h)$$



Use this approximations in order to solve optimization problems with gradient based methods

## SUMMARY - FINITE DIFFERENCE APPROXIMATIONS

-  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  ————— FORMULA ————— ERROR ————— TYPE

$$f'(x_0) \left\{ \begin{array}{l} f'(x_0) = \frac{f(x_0 + h) - f(x_0)}{h} - \frac{1}{2} f''(\xi_+) h, \\ f'(x_0) = \frac{f(x_0) - f(x_0 - h)}{h} + \frac{1}{2} f''(\xi_-) h \end{array} \right.$$

$$f'(x_0) = \frac{f(x_0 + h) - f(x_0 - h)}{2h} + \frac{1}{6} (f'''(\xi_+) + f'''(\xi_-)) h^2$$

$\mathcal{O}(h)$  forward

$\mathcal{O}(h)$  backward

$\mathcal{O}(h^2)$  centered

$$f''(x_0) = \frac{f(x_0 + h) - 2f(x_0) + f(x_0 - h)}{h^2} + \frac{1}{24} (f''''(\xi_+) + f''''(\xi_-)) h^2$$

$\mathcal{O}(h^2)$  centered

$$\nabla f(x) \quad \frac{\partial f}{\partial x_i}(x) \simeq \frac{f(x + he_i) - f(x)}{h} \quad \begin{cases} \text{FW/BW cost} = n+1 \text{ eval} \\ \text{CENTERED cost} = 2n \text{ eval} \end{cases}$$

$$\nabla^2 f(x) \rightarrow \text{Jacobian of grad} \quad \nabla^2 f(x)e_i = \frac{(\nabla f)}{\partial x_i}(x) \simeq \frac{\nabla f(x + he_i) - \nabla f(x)}{h} \quad (\text{grad cost}) + n+1 \text{ eval}$$

$\hookrightarrow$  Second derivatives  $\frac{\partial f}{\partial x_i \partial x_j}(x) = \frac{f(x + he_i + he_j) - f(x + he_i) - f(x + he_j) + f(x)}{h^2}$

$3n^2 + 1$  eval

--  $F: \mathbb{R}^n \rightarrow \mathbb{R}^m$  —————

$$F'(x) = \begin{pmatrix} \frac{\partial F(x)}{\partial x_1} & | & \dots & | & \frac{\partial F(x)}{\partial x_n} \end{pmatrix} \text{ having } \frac{\partial F}{\partial x_i}(x) \simeq \frac{F(x + he_i) - F(x)}{h}$$

$\hookrightarrow$  cost =  $n+1$  evaluations ( $n$  for  $F(x)$  and  $n$  for  $F(x + he_i)$ )

$\hookrightarrow$  optimizations: • MATRIX-FREE IMPLEMENTATIONS

- SPARSE JACOBIANS

## NON GRADIENT METHODS

### NEDLER-MEAD METHOD

↳ direct search method ← **symplex method**

Def A **SYMPLEX**  $S$  in  $\mathbb{R}^n$  is the convex hull of  $n+1$  points  $x_i \in \mathbb{R}^n$ ,

$$S = \left\{ y \in \mathbb{R}^n : y = \sum_{i=0}^{n+1} \lambda_i x_i, \lambda_i \geq 0, \sum \lambda_i = 1 \right\}$$

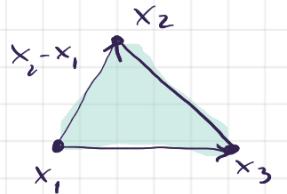
↳ a symplex  $S$  is **NON-SINGULAR** if the " $n$ " vectors

$x_2 - x_1$   
 $x_3 - x_1$   
⋮  
 $x_{n+1} - x_1$

} are linearly independent

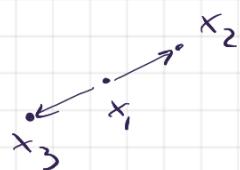
### EXAMPLE

$$\mathbb{R}^2 \rightarrow n=2$$



$$y = \lambda_1 x_1 + \lambda_2 x_2 + \lambda_3 x_3$$

non-singular



singular because  $x_2 - x_1 \parallel x_3 - x_1$

- HOW DOES NM-METHOD WORK ?

↳ based on 4 parameters

(1)	$\rho$ reflection	$\rho > 0$	(typically $\rho = 1$ )
(2)	$\chi$ expansion	$\chi > 1, \chi > \rho$	(typically $\chi = 2$ )
(3)	$\gamma$ contraction	$0 < \gamma < 1$	(typically $\gamma = \frac{1}{2}$ )
(4)	$\sigma$ shrinking	$0 < \sigma < 1$	(typically $\sigma = \frac{1}{2}$ )

IDEA: start with a simplex and, at every step  
 replace our points with a suitable new point  
 that improve your situation  
 ↳ the point chosen is the one with worse value

- DESCRIPTION OF A SINGLE ITERATION

$S_k$ : given non singular simplex

$x_1^{(k)}, \dots, x_{n+1}^{(k)}$ : points of  $S_k$

↳ assume that points  $x_i^{(k)}$  are ordered in such a way  
 that  $f(x_1^{(k)}) \leq \dots \leq f(x_{n+1}^{(k)})$

↳ for simplicity  $f_1^{(k)} = f(x_1^{(k)})$

# NUMERICAL SOLUTION OF NL EQUATIONS

↳ having  $f: \mathbb{R}^n \rightarrow \mathbb{R}$   
 find  $x^* \in \mathbb{R}$  st  $f(x^*) = 0$  → NON-LINEAR EQUATIONS

## NEWTON METHOD

### ALGORITHM

↳ at  $x_k$  linearize  $f(x_k)$   
 (1<sup>st</sup> order Taylor polynomial)

↳ get the root of  
 the linearized function  
 and use it as  $x_{k+1}$

↳ repeat until you find  $x^*$

↳ mathematically:  $f(x) \approx f(x_k) + f'(x_k)(x - x_k) = 0$

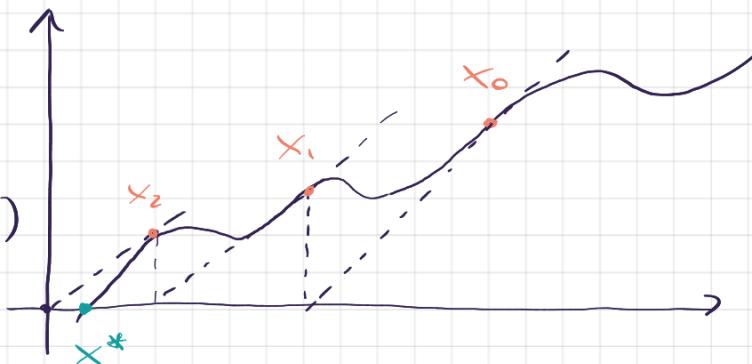
$$\Rightarrow \text{so } x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)} ; f'(x_k) \neq 0$$

### BEHAVIOR

↳ good convergence (quadratic with some assumptions)

↳ bad for robustness (the linearized function  
 may bring you in the opposite direction)

↳ quite problematic the choice of the starting point



!.

Typically, to improve robustness is to endow it with  
 a line search strategy applied to a merit function  
 (ie a function whose unconstrained minimum  
 corresponds to a solution of the nonlinear eq.)

# NUMERICAL SOLUTIONS FOR NON-LINEAR SYSTEMS OF EQUATIONS

$$F \in \mathbb{R}^n \rightarrow \mathbb{R}^m$$

→ we want to find  $F(x) = 0$

$$F(x) = \begin{pmatrix} F_1(x_1, \dots, x_n) \\ \vdots \\ F_m(x_1, \dots, x_n) \end{pmatrix}$$

$$\text{which means } \begin{pmatrix} F_1(x) = 0 \\ F_2(x) = 0 \\ \vdots \\ F_m(x) = 0 \end{pmatrix}$$

↳ important for Constrained Optimization

## NEWTON METHOD

↳ it works in a similar fashion as the one for single NC equations:

- take  $x^{(k)}$
- linearize  $F$  at  $x^{(k)}$
- replace  $F(x) = 0$  by its lin  $F^* = 0$
- the solution is  $x^{(k+1)}$
- repeat

↳ even in this case convergence is good if  $F'(x)$  is non-singular

↳ the robustness is NOT satisfying, the method is typically enhanced by the use of line search applied on to the merit function:

$$f(x) = \frac{1}{2} \|F(x)\|^2$$

Jacobian  
↓

$$F(x) \approx F(x^{(k)}) + F'(x^{(k)})(x - x^{(k)})$$

$$F(x^{(k)}) + F'(x^{(k)})(x - x^{(k)}) = 0$$

$$F'(x^{(k)})x = -F(x^{(k)}) + F'(x^{(k)}) \cdot x^{(k)}$$

↑ the only unknown in a linear system

$$A \cdot x = b$$

$F'(x^{(k)})$  should be non-singular

Computationally better

$$F(x^{(k)}) + F'(x^{(k)}) \underbrace{(x - x^{(k)})}_s = 0$$

$$F'(x^{(k)}) \cdot s = -F(x^{(k)})$$

$$\hookrightarrow \text{then } x^{(k+1)} = x^{(k)} + s$$

## NON LINEAR LEAST SQUARES PROBLEM

↳ already discussed  $F: \mathbb{R}^n \rightarrow \mathbb{R}^n$

finding  $x^*$  st.  $F(x^*) = 0 \in \mathbb{R}^n$

↳ now we consider  $F: \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  $m > n$

$F(x) = 0 \in \mathbb{R}^m$  represents the overdetermined system:

$$\left\{ \begin{array}{l} F_1(x_1, \dots, x_n) = 0 \\ \vdots \\ F_m(x_1, \dots, x_n) = 0 \end{array} \right. \quad \begin{array}{l} \text{can fail at having a solution} \\ \rightarrow (\text{no point belongs simultaneously} \\ \text{to all } F_i) \end{array}$$

↳ change the problem ( $F(x) = 0$ )  $\rightarrow$  ("all  $F_i$  as small as pos")

so now we want  $\left[ \min \frac{1}{2} \sum_i F_i^2(x^*) \right]$

↳ we call  $f(x) = \frac{1}{2} \sum F_i^2(x) = \frac{1}{2} \|F(x)\|^2$

the **MERIT FUNCTION**



In the linear case  $F(x) = Ax = y$

so the problem translates in finding  $\min_{\text{const}} \frac{1}{2} \|Ax - y\|_2^2$

$$\min \frac{1}{2} \|Ax - y\|_2^2 = \min \left[ \frac{1}{2} x^T A^T A x - x^T A^T y + \frac{1}{2} y^T y \right]$$

↳ which you can find is a convex function

(since  $A^T A$  semi pos def) so it's solvable

## GAUSS-NEWTON METHOD

↪ GN is an approximation of Newton method based on the structure of the non linear least square problem:

$$\left\{ \begin{array}{l} \nabla^2 f(x^{(k)}) p^{(k)} + \nabla f(x^{(k)}) = 0 \\ x^{(k+1)} = x^{(k)} + p^{(k)} \end{array} \right. \quad \begin{array}{l} \rightarrow \text{need } \nabla f \text{ and } \nabla^2 f \\ \text{for Newton method} \end{array}$$

$$\left[ f(x) = \frac{1}{2} \| F(x) \|^2 = \frac{1}{2} F^T(x) F(x) = \frac{1}{2} \sum F_i^2(x) \right] \rightarrow \begin{array}{l} \text{merit function} \\ \text{for NLSQ problem} \end{array}$$

•  $\nabla f$ ?  $\rightarrow \frac{\partial f}{\partial x_i} = F_1 \frac{\partial F_1}{\partial x_i} + \dots + F_m \frac{\partial F_m}{\partial x_i}$

$\hookrightarrow \left[ \nabla f(x) = [F'(x)]^T F(x) \right]$

$F'(x) = x \begin{bmatrix} & & & \\ & & & \\ & & & \\ & & & \end{bmatrix}$

•  $\nabla^2 f$ ? (I neglect dependency of  $x$  to ease notation)

$$\begin{aligned} \hookrightarrow \frac{\partial f}{\partial x_i \partial x_j} &= \frac{\partial F_1}{\partial x_j} \frac{\partial F_1}{\partial x_i} + F_1 \frac{\partial^2 F_1}{\partial x_i \partial x_j} + \dots + \frac{\partial F_m}{\partial x_j} \frac{\partial F_m}{\partial x_i} + F_m \frac{\partial^2 F_m}{\partial x_i \partial x_j} \\ &= \sum_{k=1}^m \left( \underbrace{\frac{\partial F_k}{\partial x_i}}_{[F'(x)]_{ki}} \cdot \underbrace{\frac{\partial F_k}{\partial x_j}}_{[F'(x)]_{kj}} + F_k \underbrace{\frac{\partial^2 F_k}{\partial x_i \partial x_j}}_{F_k \cdot \nabla^2 F_k(x)} \right) \\ &\qquad \qquad \qquad \underbrace{[F'(x)]_{ki} [F'(x)]_{kj}}_{(F'(x))^T F'(x)} \end{aligned}$$

$$= [F'(x)]^T F'(x) + \underbrace{\sum F_k (\nabla^2 F_k(x))}_{\hookrightarrow \text{it is very costly because}}$$

this has to compute a Hess for every  $F_i$

↳ so in the GN method you simply avoid it :  $\boxed{\nabla^2 f(x) \approx [F'(x)]^T F'(x)}$

↳ so the step vector will be :

$$N: P_k^{(N)} \text{ s.t. } \nabla^2 f(x_k) P_k^{(N)} = -\nabla f(x_k)$$

$$\boxed{GN: P_k^{(GN)} \text{ s.t. } [F'(x_k)]^T F'(x_k) P_k^{(GN)} = -F'(x_k) F(x_k)}$$

explanation of why  
GN is linearization of  
NLSQ problem

25/10

## --- Constrained Optimization ---

↳ given  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  (continuously differentiable)

$\min_{x \in X} f(x)$  with  $X \subset \mathbb{R}^n$  non-empty, convex, closed.

↳ if a vector  $x \in X$  then is said to be a  
**FEASIBLE POINT**

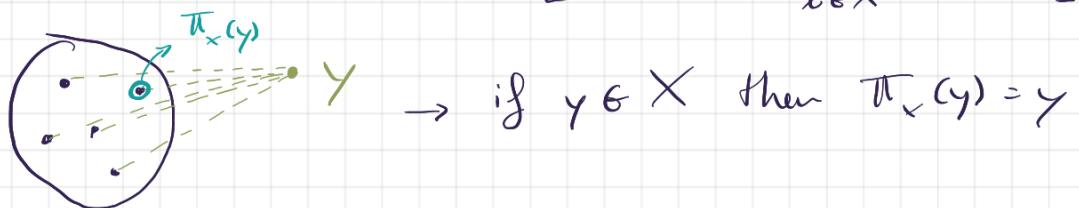
## PROJECTED GRADIENT METHOD

Def PROJECTION ONTO A SET

↳ given  $X \subset \mathbb{R}^n$

↳ given  $y \in \mathbb{R}^n$

↳ the projection of  $y$  on  $X$  is:  $\left[ \Pi_X(y) = \arg \min_{x \in X} \|y - x\|^2 \right]$



↳ To compute the projection you need to solve a minimization problem

↳ there are some cases in which is easier to compute projections:

CASE 1  $X$  is a sphere in  $\mathbb{R}^n$

$$X = \left\{ x \in \mathbb{R}^n \mid \|x\|_2 \leq R \right\}$$

$$\text{so } \Pi_X(y) = \begin{cases} y & \text{if } \|y\| \leq R \\ \frac{R}{\|y\|}y & \text{if } \|y\| > R \end{cases}$$

CASE 2  $X$  is a n-dim "box"

$$X = \left\{ x \in \mathbb{R}^n \mid L < x < U \text{ (component wise)} \right\}$$

$$\text{so } [\Pi_X(y)]_i = \begin{cases} y_i & L_i < y_i < U_i \\ L_i & y_i < L_i \\ U_i & y_i > U_i \end{cases}$$

↳ PROJECTED GRADIENT METHOD IDEA

↳ given  $x_k$  compute  $x_{k+1}$  with steepest descent

↳ if  $x_{k+1} \in X$  is ok

↳ if  $x_{k+1} \notin X$  then use projection  $\Pi_X(x_{k+1})$

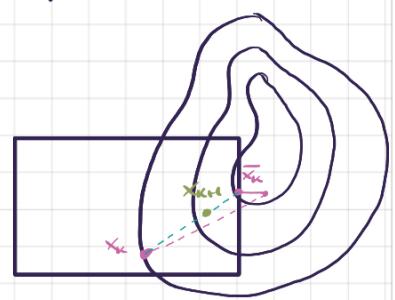
$$\text{↳ so } x_{k+1} = \Pi_X(x_k - \alpha_k \nabla f(x_k)) \quad \alpha_k > 0$$

↳ It is actually preferred to be used in two steps :

$$\textcircled{1} \quad \bar{x}_k = \Pi_X(x_k - \gamma_k \nabla f(x_k)) \quad \gamma > 0$$

\textcircled{2} move along the direction  $\bar{x}_k - x_k$  with suitable steplength  $\alpha_k > 0$

$$x_{k+1} = x_k + \alpha_k (\bar{x}_k - x_k)$$



dir  $\bar{x}_k - x_k \rightarrow$  get  $x_{k+1}$  there

↳ how do we choose  $\gamma_k$  and  $\alpha_k$ ?

① "Limited minimization rule"

↳  $\gamma_k = \gamma$  constant

↳  $\alpha_k = [0, 1]$  chosen to satisfy

$$f(x_k + \alpha_k (\bar{x}_k - x_k)) = \min_{\alpha \in [0, 1]} f(x_k + \alpha (\bar{x}_k - x_k))$$

② Armijo rule

↳  $\gamma_k = \gamma$  constant

↳  $\alpha_k$  chosen to satisfy Armijo condition

## LAGRANGE MULTIPLIER THEORY

↳ consider constrained opt. problems as

$$\min f(x) \quad f: \mathbb{R}^n \rightarrow \mathbb{R}$$

$$\text{st } h_i(x) = 0 \quad i=1 \dots n \quad h_i: \mathbb{R}^n \rightarrow \mathbb{R} \quad (\text{equality})$$

$$g_j(x) \leq 0 \quad j=1 \dots p \quad g_j: \mathbb{R}^n \rightarrow \mathbb{R} \quad (\text{inequality})$$

↳ written in compact form as:

$$h(x) = \begin{pmatrix} h_1(x) \\ \vdots \\ h_m(x) \end{pmatrix} \quad g(x) = \begin{pmatrix} g_1(x) \\ \vdots \\ g_p(x) \end{pmatrix}$$

$$h(x): \mathbb{R}^n \rightarrow \mathbb{R}^m \quad g(x): \mathbb{R}^n \rightarrow \mathbb{R}^m$$

↳ so the problem translates:

$$\min f(x)$$

$$\text{st } h(x) = 0 \quad \text{componentwise}$$

$$g(x) \leq 0 \quad \text{componentwise}$$

## • Case 1: Equality constraints only

$$(p) \quad \begin{aligned} & \min f(x) \\ & \text{s.t. } h_i(x) = 0 \end{aligned}$$

$$\begin{aligned} f: \mathbb{R}^n &\rightarrow \mathbb{R} \\ h_i: \mathbb{R}^n &\rightarrow \mathbb{R}^m \end{aligned}$$

### LAGRANGE MULTIPLIER THEOREM

↳ let  $x^*$  be the local minimum for (p)  
and assume that  $\{h_i(x^*)\}_{i=1, m}$  is a set  
of linearly independent vectors. ( $\leftarrow$  LIP)

↳ there is an unique vector  $\lambda \in \mathbb{R}^m$  st.  $\left[ \nabla f(x^*) + \sum \lambda_i \nabla h_i(x^*) = 0 \right]$   
↳ the elements of  $\lambda$  are called LAGRANGE MULTIPLIERS

### LAGRANGIAN FUNCTION

↳ introduce  $L(x, \lambda) = f(x) + \sum \lambda_i h_i(x) = f(x) + \lambda^T h(x)$   
called LAGRANGIAN FUNCTION associated to (p)

-----  
? What is the connection between LAGRANGE MULTIPLIERS  
and LAGRANGIAN FUNCTION ?  
-----

$$L(x, \lambda) = f(x) + \sum \lambda_i h_i(x) \rightarrow \nabla_x L = \underbrace{\nabla f(x^*) + \sum \lambda_i \nabla h_i(x^*)}_{\text{Lagrange mult theorem}}$$

$$\text{but also } \nabla_{\lambda} L(x, \lambda) = h(x)$$

→ in a nutshell  $\nabla_{x, \lambda} L(x^*, \lambda^*) = 0$  is the  
solution to our (p)

## • Case 2: Adding inequality constraints

$$\begin{array}{ll} \min f(x) & f: \mathbb{R}^n \rightarrow \mathbb{R} \\ \text{st } h(x) = 0 & h: \mathbb{R}^n \rightarrow \mathbb{R}^m \\ g(x) \leq 0 & g: \mathbb{R}^n \rightarrow \mathbb{R}^p \end{array}$$

### Def ACTIVE SET

↳ the active set for the inequality constraint  $g_i(x) \leq 0$  at a point  $x \in \mathbb{R}^n$  is defined as:

$$A(x) = \left\{ j \in \{1, \dots, p\} \text{ st } g_j(x) = 0 \right\}$$

↳ basically an inequality constraint  $g_i(x) \leq 0$  is active at  $x$  if  $g_i(x) = 0$

### EXAMPLE

$$\begin{array}{l} \min x^2 + 1 \\ \text{st } 1-x \leq 0 \\ \quad x-2 \leq 0 \end{array} \quad \left| \begin{array}{l} \rightarrow 1 \leq x \leq 2 \\ \rightarrow X = [1, 2] \end{array} \right. \quad (\text{feasible set})$$

$$\begin{aligned} \text{given } x \in X \rightarrow A(x) &= \left\{ j \in \{1, 2\} \mid g_j(x) = 0 \right\} \\ &= \begin{cases} \{1\} & \text{if } x=1 \\ \{2\} & \text{if } x=2 \end{cases} \end{aligned}$$

## KARUSH - KUHN - TUCKER THEOREM [ KKT CONDITIONS ]

↳ let us introduce the Lagrangian function for (P)

$$\mathcal{L}(x, \lambda, \mu) = f(x) + \sum_{i=1}^m \lambda_i h_i(x) + \sum_{j=1}^p \mu_j g_j(x)$$

↳ let  $x^*$  be a LOCAL MINIMUM point for (P)

and assume  $x^*$  is REGULAR (ie. the set of vectors  $\nabla h_i(x^*) \quad \forall i=1\dots m \quad \& \quad \nabla g_j(x^*) \quad \forall j \in A(x^*)$ ) is a set of linearly independent vectors.

↳ then there exists a UNIQUE  $x^* \in \mathbb{R}^m$  and a UNIQUE  $\mu^* \in \mathbb{R}^p$

such that : 
$$\begin{cases} \nabla_x \mathcal{L}(x^*, \lambda^*, \mu^*) = 0 \\ \nabla_\lambda \mathcal{L}(x^*, \lambda^*, \mu^*) = 0 \\ g(x^*) \leq 0 \\ \mu_j^* \geq 0 \quad \forall j = 1 \dots p \\ \mu_j^* = 0 \quad \forall j \notin A(x^*) \end{cases}$$
 which means  $\mu_j^* g_j(x^*) = 0$

• REMARK :  $\mu_j^* = 0$  if  $g_j(x^*) < 0$   
 $\mu_j^* \geq 0$  if  $g_j(x^*) = 0$   $\rightarrow \mu_j^* g_j(x^*) = 0$

↳ thus KKT CONDITIONS

are EQUIVALENTLY WRITTEN :

$$\begin{cases} \nabla_x \mathcal{L}(x^*, \lambda^*, \mu^*) = 0 \\ \nabla_\lambda \mathcal{L}(x^*, \lambda^*, \mu^*) = 0 \\ g(x^*) \leq 0 \\ \mu^* \geq 0 \\ \mu_j^* g_j(x^*) = 0 \end{cases}$$

(COMPLEMENTARITY CONDITIONS)

↳ since both  $g(x^*)$  and  $\mu^*$  should not change sign

$\mu_j^* g_j(x^*) = 0$  can be replaced by  $(\mu^*)^\top g(x^*) = 0$

(in general  $v^\top v = 0$  doesn't mean  $v_i v_i = 0 \quad \forall i$ )

**REMARK:** If  $j \notin A(x^*) \rightarrow \mu_j^* = 0$   
but if  $j \in A(x^*) \rightarrow \mu_j^* > 0$  (so both  $\langle \cdot \rangle_{\geq_0}^{=0}$  can be)

↳ if  $j \notin A(x^*)$   $\begin{cases} \mu_j^* = 0 & \text{NOT STRICT COMPLEMENTARITY} \\ \mu_j^* > 0 & \text{STRICT COMPLEMENTARITY} \end{cases}$

↳ strict complementarity  $\rightarrow$  stronger convergence

# QUADRATIC PROGRAMMING

EQUALITY CONSTRAINTS ONLY

$$(QP) \quad \min \frac{1}{2} x^T Q x + c^T x \\ \text{st} \quad Ax = b$$

$$\left| \begin{array}{l} Q \in \mathbb{R}^{n \times n} \\ \text{symmetric positive semi-definite} \\ c \in \mathbb{R}^n \quad x \in \mathbb{R}^n \\ A \in \mathbb{R}^{m \times n} \\ b \in \mathbb{R}^m \end{array} \right\} m = \text{number of constraints}$$

↪ apply KKT:

$$\begin{aligned} \mathcal{L}(x, \lambda) &= \frac{1}{2} x^T Q x + c^T x + \sum \lambda_i \cdot [Ax - b] \\ &= \frac{1}{2} x^T Q x + c^T x + \lambda^T (Ax - b) \end{aligned}$$

$$\begin{cases} \nabla_x \mathcal{L}(x, \lambda) = Qx + c + A^T \lambda \\ \nabla_\lambda \mathcal{L}(x, \lambda) = Ax - b \end{cases}$$

$$\begin{aligned} \hookrightarrow \text{so at minimum we have: } & \begin{cases} Qx + c + A^T \lambda = 0 \\ Ax - b = 0 \end{cases} \\ (\text{KKT conditions for QP}) \end{aligned}$$

↪ build  $K \omega^* = d$  such that

$$K = \begin{bmatrix} Q & A^T \\ A & \emptyset \end{bmatrix} \quad \omega^* = \begin{pmatrix} x^* \\ \lambda^* \end{pmatrix} \quad d = \begin{pmatrix} -c \\ b \end{pmatrix}$$

↳ remember that  $h(x) = 0 \rightarrow Ax = b$   
 $\nabla h(x) \rightarrow A$

and we need  $h(x)$  linearly independent  
which means  $A$  is full rank

### THEOREM

↳ Consider  $K = \begin{bmatrix} Q & A^T \\ A & \emptyset \end{bmatrix} \rightarrow Q$  is symmetric positive semidef

- If  $A$  full (row) rank and given  $Z$  such that  
columns ( $Z$ ): basis for  $\text{Ker}(A)$

↳ Then  $Z^T Q Z$  is SPD ( $Q$  is positive definite on  $\text{Ker}(A)$ )

↳ Then  $K$  is non-singular, hence KKT have unique solution  $(x^*, \lambda^*)$   
( $x^*$  local minimizer which is also a global one)

### PROPERTIES OF $K$

- $K$  is Symmetric
- Def: the INERTIA of a symmetric matrix is a triplet of numbers denoting the # of positive, negative and null eigenvalues  $(P, N, \emptyset)$

↳  $\text{INERTIA}(K) = (n, m, \emptyset)$  if  $Z^T Q Z$  is pos. def.  
dimensionality  $\uparrow$   $\nwarrow$  # constraints

↳  $K$  is INDEFINITE

## HOW TO SOLVE KKT CONDITIONS FOR QP

### ① Full system factorization

$$K = LDL^T \quad / \quad LU \text{ decomposition}$$

↳ GOOD APPROACH IF:  $K$  not large and not sparse

### ② Schur complement approach

$$\begin{cases} Qx + A^T x = -c \\ Ax = b \end{cases} \rightarrow \text{if } Q \text{ is non singular you can eliminate } x \text{ from the system}$$

$$x = -Q^{-1}A^T\lambda - Q^{-1}c$$

$$A(-Q^{-1}A^T\lambda - Q^{-1}c) = b$$

$$\underbrace{AQ^{-1}A^T\lambda}_{\text{has dim } m} = -b - AQ^{-1}c \rightarrow$$

(from  $m+n$  to  $m$ )

$$\left\{ \begin{array}{l} \hat{Q} = AQ^{-1}A^T \text{ is SPD} \\ \text{under assumptions:} \\ \cdot A \text{ full rank} \\ \cdot Q \text{ non singular} \end{array} \right.$$

↳ GOOD APPROACH IF:  $m$  is small and

$Q$  easy to invert

### ③ Null-Space method

↳ suitable also for  $A$  singular

↳ assume we know one of  $\hat{x}$  st  $A\hat{x} = b$

(we must have many, as  $m < n$ )

↳ assume also that full rank  $Z \in \mathbb{R}^{n \times (n-m)}$

is given such that  $AZ = \emptyset$

(note that  $n-m = \dim(\ker(A))$ )

↳  $\left\{ \text{column } z_i \text{ of } Z : Az_i = \emptyset \right\}$

$z_i \in \ker(A)$

$z_i$  are linearly independent

$\Rightarrow Z$ : columns form  
a basis for  $\ker(A)$

↳ we can write  $[x = Zv + \hat{x}]$

where  $Z$ : generic element in  $\ker(A)$

$\hat{x}$ : particular solution

$v$ : arbitrary vector in  $\mathbb{R}^{n-m}$

↳ replace it in 1st equation

$$Qx - A^T \lambda = -c \rightarrow QZv + Q\hat{x} - A^T \lambda = -c$$

↳ premultiply  $Z^T$

because  $AZ = Z^T A^T = 0$

$$Z^T Q Z v + Z^T Q \hat{x} + Z^T A^T \lambda = -Z^T c$$

$$\boxed{Z^T Q Z v = -Z^T (c + Q \hat{x})} \rightarrow \text{being a } \mathbb{R}^{(n-m) \times (n-m)} \text{ system}$$

↳ LINSYS with  $Z^T Q Z$  symmetric positive definite  
with dimension  $n-m$

↳ you get  $v \rightarrow x = Zv + \hat{x}$

**REMARK:** If you also need  $\lambda$ :

$$Qx + A^T \lambda = -c$$

$$A^T \lambda = -c - Qx$$

$$AA^T \lambda = -Ae - AQx \quad \text{and solve for } \lambda$$

**REMARK:** How to obtain  $Z$ ?

↳ if  $A$  is full rank

$$A = \left[ \begin{array}{c|c} A_1 & A_2 \\ \hline \underbrace{\hspace{1cm}}_m & \underbrace{\hspace{1cm}}_{n-m} \end{array} \right] \left\{ \begin{array}{l} m \\ A_1 \text{ square} \\ A_2 \end{array} \right\} \rightarrow \begin{array}{l} \mathbb{R}^{m \times m} \\ \mathbb{R}^{n-m \times m} \end{array}$$

↳  $A_1$  square NON singular

(if the first  $m$  columns are NOT linearly independent we need to reorder)

$$\begin{aligned} \Rightarrow Z &= \left[ \begin{array}{c|c} -A_1^{-1}A_2 & \\ \hline \cdots & \\ I & \end{array} \right] \left\{ \begin{array}{l} m \\ n-m \end{array} \right\} \rightarrow AZ = [A_1 \mid A_2] \left[ \begin{array}{c|c} -A_1^{-1}A_2 & \\ \hline \cdots & \\ I & \end{array} \right] \\ &= -A_1^{-1}A_1A_2 + A_2 \\ &= \emptyset \end{aligned}$$

**REMARK:** How to obtain  $\hat{x}$ ?

$$A\hat{x} = b \rightarrow \left[ \begin{array}{c|c} A_1 & A_2 \end{array} \right] \left[ \begin{array}{c} \hat{x}_1 \\ \hline \hat{x}_2 \end{array} \right] = b \rightarrow A_1\hat{x}_1 + A_2\hat{x}_2 = b$$

$$\text{take } \hat{x}_2 = 0, \text{ then } \hat{x}_1 = A_1^{-1}b \rightarrow \hat{x} = \left[ \begin{array}{c} A_1^{-1}b \\ \hline 0 \end{array} \right]$$

# Duality theory

↳ problems with inequality constraints only

$$\left[ \begin{array}{ll} (\text{P}) & \min f(x) \quad f: \mathbb{R}^n \rightarrow \mathbb{R} \\ \text{st} & g(x) \leq 0 \quad g: \mathbb{R}^n \rightarrow \mathbb{R}^m \end{array} \right]$$

↳ Lagrangian function:  $\mathcal{L}(x, \lambda) := f(x) + \sum \lambda_i g_i(x)$   
 $\mathcal{L}: \mathbb{R}^{n+m} \rightarrow \mathbb{R}$

↳ Define  $q(\lambda) := \inf_x \mathcal{L}(x, \lambda)$



What is the infimum of a set?

↳ given  $A \subset \mathbb{R}$ ,  $\inf A$  is a number  $c$  such that:

$$\textcircled{1} \quad c \leq x \quad \forall x \in A$$

$$\textcircled{2} \quad \nexists b \text{ st } b \leq x, \forall x \in A, c < b$$

$$\text{Example 1: } A = \left\{ \frac{1}{n} \right\}_{n \in \mathbb{N}} = \left\{ 1, \frac{1}{2}, \frac{1}{3} \right\} \rightarrow \inf A = 0$$

↳ MINIMUM vs INFIMUM?  $\rightarrow \text{MIN} \in A$  always

$\text{INF} \in A$  not always

↳ we introduce the set

$$\mathcal{D} = \left\{ \lambda \text{ st } q(\lambda) > -\infty \right\} \text{ dual problem for (P)}$$

$$(\text{D}) \quad \max q(\lambda) \iff (\text{P}) \quad \min f(x)$$

$$\text{st } \lambda \geq 0 \quad \text{st } g(x) \leq 0$$

## EXAMPLE

$$\min \frac{1}{2} (x_1^2 + x_2^2)$$

$$\text{st } 1-x_1 \leq 0$$

↳ we can see  $x^* = (1, 0)$

$$\mathcal{L}(x, \lambda) = \frac{1}{2} (x_1^2 + x_2^2) + \lambda (1 - x_1)$$

↳ for a fixed  $\lambda$ , what is  $\inf \mathcal{L}(x, \lambda)$ ?

↳ let us compute  $\nabla_x \mathcal{L}(x, \lambda)$  to minimize  $\mathcal{L}(x, \lambda)$

for fixed  $\lambda$ :

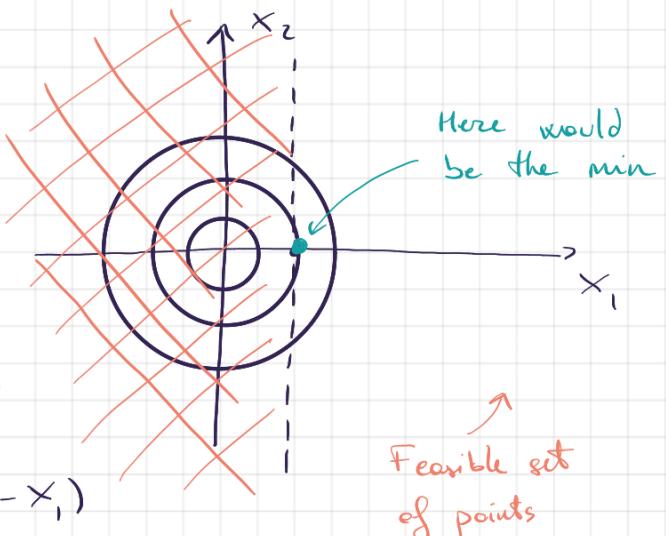
$$\nabla_x \mathcal{L}(x, \lambda) = \begin{pmatrix} x_1 - \lambda \\ x_2 \end{pmatrix} = 0 \rightarrow \begin{cases} x_1 = \lambda \\ x_2 = 0 \end{cases}$$

$$\mathcal{L}(x, \lambda) \Big|_{x=(\lambda, 0)} = \frac{1}{2} \lambda^2 + \lambda (1 - \lambda) = -\frac{1}{2} \lambda^2 + \lambda$$

↳ we found the dual problem:

$$(D) \max -\frac{1}{2} \lambda^2 + \lambda$$

st  $\lambda \geq 0$



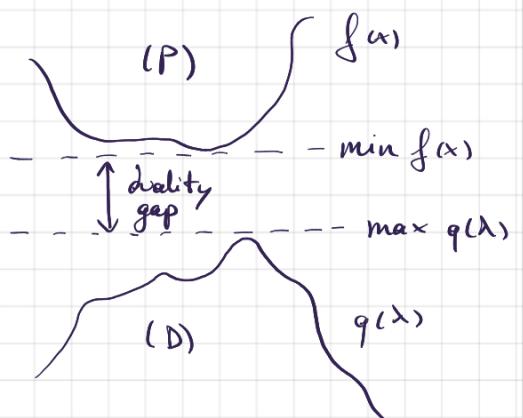
## THEOREM

The function  $q(\lambda)$  is concave and  $\mathcal{D}$  is convex

## WEAK DUALITY THEOREM

- ↳ Relationship between  $(P)$  and  $(D)$
- ↳ for any feasible  $x$  for  $(P)$  and any feasible  $\lambda$  for  $(D)$

$$q(\lambda) \leq f(x)$$



## THEOREM $\rightarrow (P) - (D)$ pair KKT conditions?

- ↳ let  $x^*$  be a solution to  $(P)$  with  $f, g$  convex functions smooth at  $x^*$
- ↳ if  $(x^*, \lambda^*)$  satisfies KKT conditions for  $(P)$  then  $\lambda^*$  is a solution for  $(D)$
- ↳ the Lagrange multiplier of  $(P)$  is the solution to  $(D)$

## BACK TO EXAMPLE

$$\begin{array}{ll} \min & \frac{1}{2} (x_1^2 + x_2^2) \\ \text{st} & 1-x_1 \leq 0 \end{array} \quad \Rightarrow \quad \mathcal{L}(x, \lambda) = \frac{1}{2} (x_1^2 + x_2^2) + \lambda(1-x_1)$$

KKT CONDITIONS

$$\left\{ \begin{array}{l} x_1 - \lambda = 0 \\ x_2 = 0 \\ 1 - x_1 \leq 0 \\ \lambda \geq 0 \\ \lambda(1-x_1) = 0 \end{array} \right\} \quad \begin{array}{l} \nabla_x \mathcal{L}(x, \lambda) = 0 \\ g(x) \leq 0 \\ \text{non negativity of } \mathcal{L} \text{ multipliers} \\ \text{complementarity condition} \end{array}$$

doesn't satisfy  $1-x_1 \leq 0$

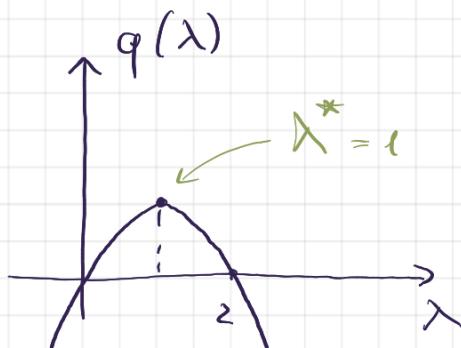
↳ we find  $\begin{pmatrix} x_1 = \lambda \\ x_2 = 0 \end{pmatrix}, \lambda(1-\lambda) = 0 \quad \begin{cases} \lambda = 0 \\ \lambda = 1 \end{cases}$

$$\text{so } \omega^* = \underbrace{\begin{pmatrix} 1 & 0 \\ x^* & \lambda^* \end{pmatrix}}_{\omega^*}'$$

↳ consider now the dual

$$\max \lambda(1-\lambda)$$

$$\text{st } \lambda \geq 0$$



## DT - LINEAR PROGRAMMING

$$\left[ \begin{array}{lll} (\text{LP}) & \min & c^T x \\ \text{st} & b - Ax \leq 0 & x \in \mathbb{R}^n \quad c \in \mathbb{R}^n \\ & & b \in \mathbb{R}^m \quad A \in \mathbb{R}^{m \times n} \end{array} \right]$$

↳ dual objective :  $q(\lambda) = \inf_x \mathcal{L}(x, \lambda)$

$$\begin{aligned} \mathcal{L}(x, \lambda) &= c^T x + \lambda^T (b - Ax) \\ &= (c - A^T \lambda)^T x - \lambda^T b \end{aligned}$$

- CASE 1:  $c - A^T \lambda \neq 0$

then  $\inf_x \mathcal{L}(x, \lambda) = \inf_x \underbrace{(c - A^T \lambda)^T x}_{\text{unbounded from below (it's a simple straight line)}} + \lambda^T b = -\infty$

- CASE 2:  $c - A^T \lambda = 0$

then  $\inf_x \mathcal{L}(x, \lambda) = \inf_x \lambda^T b = \lambda^T b$

$$q(\lambda) = \lambda^T b$$

# INTERIOR POINT METHODS

## LINEAR PROGRAMMING

↳ consider the following

$$(LP) \quad \begin{aligned} & \min c^T x \\ \text{(1)} \quad & \text{st } Ax = b \\ & x \geq 0 \end{aligned} \quad \begin{aligned} & x \in \mathbb{R}^n, c \in \mathbb{R}^n \\ & b \in \mathbb{R}^m, A \in \mathbb{R}^{m \times n} \text{ full rank} \end{aligned}$$

! This way of writing the equivalence is the "canonical form" and it's equivalent to  $b - Ax \leq 0$  written before BUT with different meanings for  $c, x, A, b$ :

↳ start from (1)  $\rightarrow b - Ax \leq 0$  form

↳ introduce  $s = Ax - b$  (slack variable)

so rewrite (1) as

$$\begin{cases} \min c^T x \\ \text{st } s \geq 0 \\ s = Ax - b \end{cases}$$

↳ introduce  $\begin{cases} x^+ = \max(x, 0) \\ x^- = \max(-x, 0) \end{cases} \rightarrow x = \begin{pmatrix} x^+ \\ -x^- \\ 0 \end{pmatrix} \quad x^+ = \begin{pmatrix} 7 \\ 0 \\ 5 \end{pmatrix} \quad x^- = \begin{pmatrix} 0 \\ 3 \\ 0 \end{pmatrix}$

so  $x = x^+ - x^-$

↳ so write (1) as

$$\min c^T(x^+ - x^-) = c^T x^+ - c^T x^- + 0s$$

$$\text{st } s = Ax^+ - Ax^- - b$$

$$s \geq 0 \quad x^+ \geq 0 \quad x^- \geq 0$$

↳ if we introduce  $\hat{c} = \begin{pmatrix} c \\ -c \\ 0 \end{pmatrix} \quad \hat{x} = \begin{pmatrix} x^+ \\ x^- \\ 0 \end{pmatrix} \quad \hat{A} = [A \quad -A \quad -I]$

↳ we get

$$\min \hat{c}^T \hat{x}$$

$$(2) \quad \hat{A} \hat{x} = b$$

$$\hat{x} \geq 0$$

↳ let's write the dual for (LP) ②

$$(P-LP) \quad \min c^T x$$

$$\text{st} \quad Ax = b$$

$$x \geq 0$$

A full rank

$\lambda \rightarrow$  equality constraints  $\lambda \in \mathbb{R}^m$

$s \rightarrow$  inequality constraints  $s \in \mathbb{R}^n$

because we need  $\leq$  ineq. const.

$$\begin{aligned} \mathcal{L}(x, \lambda, s) &= c^T x + s^T (-x) + \lambda^T (b - Ax) \\ &= (c - s - A^T \lambda)^T x + \lambda^T b \end{aligned}$$

$$\hookrightarrow q(\lambda, s) = \min_x \mathcal{L}(x, \lambda, s) \rightarrow \text{for fixed } \lambda, s$$

$$\nabla_x \mathcal{L}(x, \lambda, s) = c - s - A^T \lambda = 0$$

↪ replace  $c - s - A^T \lambda = 0$  in  $\mathcal{L}(x, \lambda, s)$   
and get  $q(\lambda, s) = \lambda^T b$

$$(D-LP) \quad \max \lambda^T b$$

$$\text{st} \quad c - s - A^T \lambda = 0$$

$$s \geq 0, \lambda \geq 0$$

-----  
 $x$ : Primal var.  
 $\lambda, s$ : Dual var  
-----

! Remember Weak Duality holds:

$$\lambda^T b \leq c^T x \quad \forall x \mid x \text{ is primal feasible}$$

$$\forall \lambda \mid \lambda \text{ is dual feasible}$$

$$X = \{x \in \mathbb{R}^n \mid Ax = b, x \geq 0\} \rightarrow \text{primal feasible}$$

$$Y = \{(\lambda, s) \in \mathbb{R}^{m+n} \mid A^T \lambda + s = c, s \geq 0\} \rightarrow \text{dual feasible}$$

## STRONG DUALITY THEOREM (FOR LP)

↳ if either  $X \neq \emptyset$  or  $Y \neq \emptyset$   
 and either P-LP or D-LP is solvable  
 then:  $\min_x C^T x = \max_{\lambda} b^T \lambda$

↳ thus the DUALITY GAP is vanishing

## KKT CONDITIONS FOR LP

$$\begin{aligned} (\text{P-LP}) \quad & \min c^T x \\ \text{st} \quad & Ax = b \\ & x \geq 0 \end{aligned}$$

$$\mathcal{L}(x, s, \lambda) = c^T x - s^T x + \lambda^T (b - Ax)$$

$$\text{KKT: } \left\{ \begin{array}{l} \nabla_x \mathcal{L} = 0 \\ \nabla_{\lambda} \mathcal{L} = 0 \\ x \geq 0 \\ s \geq 0 \\ x^T s = 0 \quad (\text{complementarity condition}) \end{array} \right.$$

$$\text{KKT: } \left\{ \begin{array}{l} c - s - A^T \lambda = 0 \\ b - Ax = 0 \\ x \geq 0 \\ s \geq 0 \\ x^T s = 0 \end{array} \right.$$

$$\begin{aligned} (\text{D-LP}) \quad & \max b^T \lambda \\ \text{st} \quad & c - s - A^T \lambda = 0 \\ & s \geq 0 \end{aligned}$$

<< can be re-written as: >>

$$\begin{aligned} & \max b^T \lambda \\ \text{st} \quad & c - A^T \lambda \geq 0 \end{aligned}$$

$$\mathcal{L}(\lambda, x) \xrightarrow{\text{multiplier for ineq. constr.}} -b^T \lambda + x^T (A^T \lambda - c)$$

$$\text{KKT: } \left\{ \begin{array}{l} Ax - b = 0 \\ c - A^T \lambda \geq 0 \\ x \geq 0 \\ x^T (c - A^T \lambda) = 0 \end{array} \right.$$

$$\text{KKT: } \left\{ \begin{array}{l} Ax - b = 0 \\ s \geq 0 \\ x \geq 0 \\ x^T s = 0 \\ s = c - A^T \lambda \end{array} \right. \leftarrow \text{reintroduced}$$

! We can observe that we have the same KKT - conditions for both problems

↳ indeed they are called **KKT conditions for the primal-dual pair**

### INTERIOR POINT METHOD IDEA (FOR LP)

↳ it tries to find a vector  $(x, \lambda, s)$  satisfying:

$$\text{KKT} \left\{ \begin{array}{l} c - s - A^T \lambda = 0 \\ b - Ax = 0 \\ x \geq 0 \\ s \geq 0 \\ x^T s = 0 \end{array} \right. \rightarrow \begin{array}{l} \text{constrained} \\ \text{non linear} \\ \text{equations} \end{array}$$

↳ thus we can define:

$$F(x, \lambda, s) = \begin{pmatrix} c - s - A^T \lambda \\ Ax - b \\ x^T s \end{pmatrix} \rightarrow \begin{array}{l} \mathbb{R}^n \\ \mathbb{R}^m \\ \mathbb{R} \end{array}$$

↳ PROBLEM:  $F: \mathbb{R}^{2n+m} \rightarrow \mathbb{R}^{n+m+1}$  not squared

↳ TRICK: replace  $x^T s = \sum x_i s_i = x_1 s_1 + \dots + x_n s_n = 0$   
with  $XSe$

$$\text{where } X = \begin{pmatrix} x_1 & x_2 & \dots & x_n \end{pmatrix}, S = \begin{pmatrix} s_1 & s_2 & \dots & s_n \end{pmatrix}, e = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}$$

$$\text{so } XSe = \begin{pmatrix} x_1 s_1 \\ x_2 s_2 \\ \vdots \\ x_n s_n \end{pmatrix} \in \mathbb{R}^n$$

↳ we get  $F(x, \lambda, s) = \begin{bmatrix} A^T \lambda + s - c \\ Ax - b \\ x^T s_e \end{bmatrix}$

having  $F: \mathbb{R}^{2n+m} \rightarrow \mathbb{R}^{2n+m}$

↳ so we need to find  $(x, \lambda, s)$  st  $\begin{cases} F(x, \lambda, s) = 0 \\ x, s \geq 0 \end{cases}$

### INTERIOR POINT METHOD ALGORITHM FOR LP

- ↳ iteratively choose  $(x_k, \lambda_k, s_k)$  satisfying  $x_k > 0, s_k > 0$
- ↳ use perturbed Newton method
- ↳ ALGORITHM:

- Given  $(x_0, \lambda_0, s_0)$
- For  $k \geq 0$ 
  - Compute  $(\Delta x_k, \Delta \lambda_k, \Delta s_k)$  with Newton method

$$F'(x_k, \lambda_k, s_k) \begin{pmatrix} \Delta x_k \\ \Delta \lambda_k \\ \Delta s_k \end{pmatrix} = -F(x_k, \lambda_k, s_k)$$

$$\begin{bmatrix} 0 & A^T & I_n \\ A & 0 & 0 \\ 0 & 0 & X \end{bmatrix} \quad \text{sparse!}$$

- Update  $\begin{bmatrix} x_{k+1} \\ \lambda_{k+1} \\ s_{k+1} \end{bmatrix} = \begin{bmatrix} x_k \\ \lambda_k \\ s_k \end{bmatrix} + \alpha_k \begin{bmatrix} x_k \\ \lambda_k \\ s_k \end{bmatrix}$

$\alpha_k$  given by line search in order to satisfy  $x_{k+1} > 0$  and  $s_{k+1} > 0$

## IMPROVEMENTS FOR IPM FOR LP

(i) Consider a less aggressive choice

$$F(x, \lambda, s) = 0 \xrightarrow{\text{X}} F(x, \lambda, s) \leq \text{tol}$$

↳ in particular we focus on  $XSe$   
at each iteration  $XSe$  is "sufficiently small"  
and approaches zero towards the solution

↳ we need a measure for the distance  
from the solution:

$$\mu = \frac{\sum x_i s_i}{n} \xrightarrow{k \rightarrow \infty} 0 \quad \rightarrow n\mu = x^T s =: \text{duality gap}$$

↳ we can prove  $x^T s$  is the duality gap actually:

$$\begin{aligned} \text{Duality Gap} &=: c^T x - b^T \lambda = \\ &= (A^T \lambda + s)^T x - b^T \lambda \rightarrow \text{use } c = A^T \lambda + s \\ &= \dots = 0 \\ &= \lambda^T (A x - b) + x^T s = x^T s \end{aligned}$$

↳ thus by the strong duality theorem we  
know that  $x^T s$  vanishes at the solution

↳ so we can impose

$$XSe = 0 \xrightarrow{\text{X}} XSe - \theta \mu e = 0$$

↳  $\mu$  will be updated at each step

↳  $\theta \in (0, 1)$  called centrality parameter

(2)

A different role of choosing  $\alpha_k$

↳ for simplicity we ignore  $\Theta$  for now

Def

### ANALYTIC CENTER

↳ for any fixed  $\mu > 0$ , solution to

$F(x, \lambda, s; \mu) = 0$  is called analytic center  
or  $\mu$ -center and it's denoted by  
 $(x(\mu), \lambda(\mu), s(\mu))$

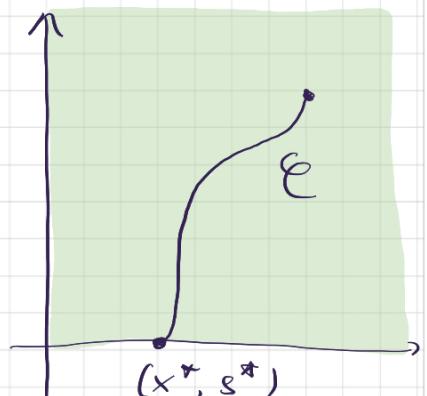
Def

### CENTRAL PATH

↳ the set

$$\mathcal{C} = \{(x(\mu), \lambda(\mu), s(\mu)) : \mu > 0\}$$

is called primal-dual central path



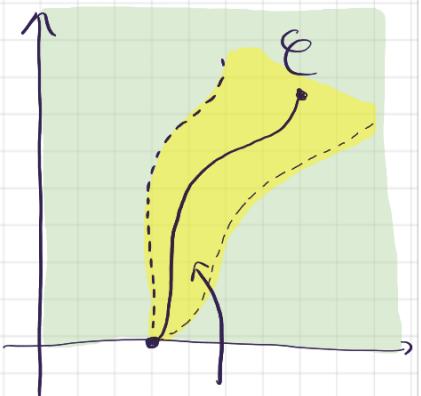
for  $\mu = 0$

↳ we now want to follow the central path but NOT exactly,  
just stay close

↳ we want to be more strict as we get close to  $\mu = 0$

$$XSe = \mu e \rightarrow XSe \approx \mu e$$

↳ we introduce the neighborhood of  $\mathcal{C}$



$$N_2(0) = \{(x, \lambda, s) \text{ strictly feasible} : \|XSe - \mu e\|_2 \leq \theta \mu\}$$

$$(x, \lambda, s) \in \overset{\circ}{F} = \left\{ (x, \lambda, s) \in \mathbb{R}^{n+m} : \begin{array}{l} Ax = b \\ c = s + A^T \lambda \\ x, s \geq 0 \end{array} \right\}$$

strictly

## SHORT STEP PATH FOLLOWING FEASIBLE (PM)

$(\Delta x_k, \Delta \lambda_k, \Delta s_k)$  follows  $\mathcal{C}$   
is actually short

all iterates  
are feasible

- start with  $(x_0, \lambda_0, s_0) \in \mathcal{N}_2(\theta)$
- for  $k \geq 0$ 
  - Compute Newton step ensuring  $(x_{k+1}, \lambda_{k+1}, s_{k+1})$  stay feasible:

$$F'(x_k, \lambda_k, s_k; \mu_k) \begin{pmatrix} \Delta x_k \\ \Delta \lambda_k \\ \Delta s_k \end{pmatrix} = -F(x_k, \lambda_k, s_k; \mu_k)$$

$$\begin{bmatrix} 0 & A^T & I_n \\ A & 0 & 0 \\ S & 0 & X \end{bmatrix} \begin{pmatrix} \Delta x_k \\ \Delta \lambda_k \\ \Delta s_k \end{pmatrix} = \begin{bmatrix} c - s_k - A^T \lambda_k > 0 \\ b - Ax_k > 0 \\ -XSe + \theta \mu_k e \end{bmatrix}$$

since is feasible  
it satisfies  
first two  
conditions



How do we know we are and stay feasible?

↳ it can be proven that by choosing

$$\theta = 0, 1 \quad \theta = 1 - \frac{\sigma_f}{\sqrt{n}}$$

all iterates belong to  $\mathcal{N}_2(\theta)$

↳ we use  $\theta$  to reduce  $\mu$ :  $\mu_{k+1} = \theta \mu_k$



DRAWBACK: Since  $\mathcal{N}_2(\theta)$  is very narrow we make little progress at each iteration

## LONG STEP PATH FOLLOWING FEASIBLE IPM

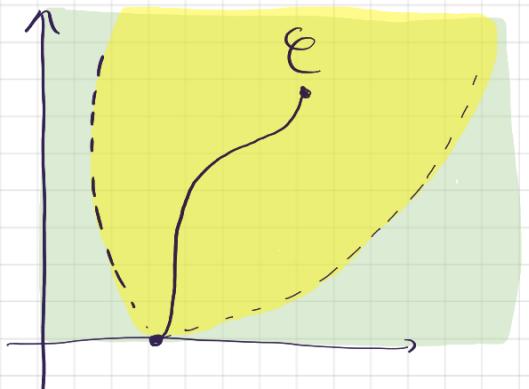
↳ we can enlarge  $N_2(\theta)$  to  $N_\infty(\gamma)$

$$N_\infty(\gamma) = \left\{ (x, \lambda, s) \in \mathcal{F}^0 : x_i s_i \geq \gamma \mu \quad \forall i \right\}$$

$\gamma \approx 10^{-3}$

↳ basically we ask that no pair  $x_i, s_i$  that approaches zero faster than the others

↳ this way we can take generally longer steps



## THEOREM

↳ given  $(x_0, \lambda_0, s_0) \in N_\infty(\gamma)$

given  $\epsilon > 0$

↳ asymptotically you get

$\mu_K \leq \epsilon \mu_0$  within  $\mathcal{O}(n \log(\frac{1}{\epsilon}))$  iterations

## PREDICTOR-CORRECTOR IPM IMPLEMENTATION

↳ consider the original KKT (P)-(D) pair

$$\begin{cases} c - s - A^T \Delta = 0 \\ Ax - b = 0 \\ XS e = 0 \\ X, S \geq 0 \end{cases} \quad \text{without } -\mu e$$

↳ also consider the Newton step:

↳ called **affine scaling system**

$$\begin{bmatrix} 0 & A^T & I_n \\ A & 0 & 0 \\ S & 0 & X \end{bmatrix} \begin{bmatrix} \Delta X_k^{AFF} \\ \Delta \Delta_k^{AFF} \\ \Delta S_k^{AFF} \end{bmatrix} = \begin{bmatrix} c - S_k - A^T \Delta_k \\ b - Ax_k \\ -XS_e \end{bmatrix}$$

↳ let's consider what happens for  $X_k + \Delta X_k^{AFF}$  and  $S_k + \Delta S_k^{AFF}$   
 (we are seeking  $X, S$  satisfying  $XS_e = 0$  having  $X, S \geq 0$ )

↳  $XS_e = 0$  at solution,

$$\text{so in general } X_k S_k \neq 0 \quad = \text{diag}(\Delta X_k) \quad = \text{diag}(\Delta S_k)$$

$$X_{k+1} = X_k + \Delta X_k^{AFF}; \quad S_{k+1} = S_k + \Delta S_k^{AFF}$$

$$(X S_e)_{k+1} = (X_k + \Delta X_k^{AFF})(S_k + \Delta S_k^{AFF})$$

$$= X_k S_k + X_k \Delta S_k^{AFF} + S_k \Delta X_k^{AFF} + \Delta X_k^{AFF} \Delta S_k^{AFF}$$

$\Rightarrow 0 \rightarrow$  according to  $\star$  (3rd block eq)

$$X_k \Delta S_k^{AFF} + S_k \Delta X_k^{AFF} = -X_k S_k e$$

$$= \Delta X_k^{AFF} \Delta S_k^{AFF}$$

! The situation improves if  $\Delta X_k^{AFF} \Delta S_k^{AFF} \approx 0$

↳ Predictor-Corrector IDEA:

- Compute the Newton step ( $\Delta x_k^A \Delta \lambda_k^A \Delta s_k^A$ ) → PREDICTOR

$$\begin{bmatrix} 0 & A^T & I_n \\ A & 0 & 0 \\ S & 0 & X \end{bmatrix} \begin{bmatrix} \Delta x_k^{AFF} \\ \Delta \lambda_k^{AFF} \\ \Delta s_k^{AFF} \end{bmatrix} = \begin{bmatrix} C - S_k - A^T \lambda_k \\ b - Ax_k \\ -XSe \end{bmatrix}$$

- Compute a new direction ( $\Delta x_k^C \Delta \lambda_k^e \Delta s_k^c$ )  
that pushes towards  $\Delta x_k^A \Delta s_k^A = 0$

i.e. solve

$$\begin{bmatrix} 0 & A^T & I_n \\ A & 0 & 0 \\ S & 0 & X \end{bmatrix} \begin{bmatrix} \Delta x_k^{Corr} \\ \Delta \lambda_k^{Corr} \\ \Delta s_k^{Corr} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ -\Delta x_k^A \Delta s_k^A \end{bmatrix}$$

I assume  
To have a  
feasible iterate

Implementation

18/12

70%

I love you  
pusone

I love you  
more pusone

PUS1

- Summarizing Predictor-Corrector IPM

- ① Solve the affine scaling system
- ② Compute  $\alpha$  st  $x_{k+1} = x_k + \alpha \Delta x_k^A \geq 0$   
 $s_{k+1} = s_k + \alpha \Delta s_k^A \geq 0$

- ③ Build  $\mu_k^{AB}$ ,  $\Omega_k$
- ④ Solve for  $(\Delta x_k^{\text{new}}, \Delta s_k^{\text{new}}, \Delta S_k^{\text{new}})$  the new system  
merging predictor & corrector + perturbation
- ⑤ Compute  $\alpha$  st  $x_{k+1} = x_k + \alpha \Delta x_k^{\text{new}}$   
 $s_{k+1} = s_k + \alpha \Delta s_k^{\text{new}}$
- ⑥ Compute  $(x_{k+1}, \Delta x_{k+1}, s_{k+1})$



## INTERIOR POINT METHODS FOR QP

↳ also having inequality constraints

$$(QP) \quad \begin{aligned} \min \quad & \frac{1}{2} x^T Q x + c^T x \\ \text{st.} \quad & A x \geq b \end{aligned} \quad \rightarrow \quad \begin{array}{l} (\text{KKT}) \\ \left\{ \begin{array}{l} Q x + c - A^T \lambda = 0 \\ \lambda \geq 0 \\ b - A x \leq 0 \\ A^T (b - A x) = 0 \end{array} \right. \end{array}$$

$$\begin{array}{l} (\text{KKT}) \\ \left\{ \begin{array}{l} Q x + c - A^T \lambda = 0 \\ \lambda \geq 0 \\ y \geq 0 \\ A^T y = 0 \rightarrow \Delta Y_e = 0 \\ y = A x - b \end{array} \right. \end{array}$$

$\xrightarrow{\text{diag } (\lambda)}$        $\xrightarrow{\text{diag } (y)}$

$$\begin{array}{l} \text{↳ now solve } \left\{ \begin{array}{l} F(x, y, \lambda) = 0 \\ \Delta Y_e \geq 0 \end{array} \right. \quad \text{being } F(x, y, \lambda) = \begin{bmatrix} Q x + c - A^T y \\ A x - b - y \\ \Delta Y_e \end{bmatrix} \end{array}$$

$$\begin{array}{l} \text{↳ better having perturbed } F(x, y, \lambda; \mu, \sigma) = \begin{bmatrix} Q x + c - A^T y \\ A x - b - y \\ \Delta Y_e - \sigma \mu e \end{bmatrix} \end{array}$$



## PENALTY METHODS

↪ consider a purely non-linear program

$$\begin{array}{ll} \min f(x) & f: \mathbb{R}^n \rightarrow \mathbb{R}^n \\ \text{st } h(x) = 0 & h: \mathbb{R}^m \rightarrow \mathbb{R}^m \\ g(x) \leq 0 & g: \mathbb{R}^n \rightarrow \mathbb{R}^p \end{array}$$

**IDEA:** plug the constraints in the cost function  
 ↪ we have an unconstrained problem

### QUADRATIC PENALTY METHODS

↪ first analyze equality constraints

$$(P_1) \min f(x) \quad \leftarrow Q(x; \mu) := f(x) + \frac{\mu}{2} \|h(x)\|_2^2$$

↙ Parameter

$$\text{st } h(x) = 0$$

! The parameter may make the problem unbounded from below, so we have to be careful in choosing it

↪ what about inequality constraints?

↪ we introduce  $g^+(x) = \max(g(x), 0)$

$$g = \begin{pmatrix} 7 \\ -5 \\ x-1 \end{pmatrix} \rightarrow g^+ = \begin{pmatrix} 7 \\ 0 \\ \max(x-1, 0) \end{pmatrix}$$

and replace  $g(x) \leq 0$  with  $g^+(x) = 0$

$$Q(x; \mu) = f(x) + \frac{\mu}{2} \|h(x)\|^2 + \frac{\mu}{2} \|g^+(x)\|^2$$

## QUADRATIC PENALTY METHOD ALGORITHM

- Given  $x_0^{(0)}, \mu_0 > 0$
  - For  $K > 0 \leftarrow$  outer iter
    - Build  $Q(x; \mu_K)$  and find an APPROXIMATE sol  
 $\min Q(x; \mu_K)$  starting from  $x_K^{(0)}$
    - Update  $\mu_{K+1}$ , choose  $x_{K+1}^{(0)}$  for next step
- Until convergence

NOTE :  $\mu_{K+1} = \rho \mu_K$  ,  $\rho = \begin{cases} 10 & \rightarrow \text{if inner iterations small} \\ 1,5 & \rightarrow \text{if inner iterations large} \end{cases}$

## PROBLEMS WITH Q.PENALTY METHODS

↳ in general  $x_K \approx \operatorname{argmin} Q(x; \mu_K)$

but  $x_K \neq \operatorname{argmin} f(x)$   
 st  $h(x) = 0$

↳ you reach  $x_K \approx \operatorname{argmin} f(x)$  only at CONVERGENCE

↳ furthermore  $h(x_K) \neq 0$

↳ it can be PROVEN that :

$$h_i^*(x) \approx -\frac{\lambda_i^*}{\mu_K} \rightarrow i\text{-th Lagrange Multiplier associated to } h_i(x) = 0 \text{ in KKT theorem}$$

↳ penalty parameter at iteration K

↳ remember  $\mu_K \xrightarrow{K \rightarrow \infty} \infty$  so  $h(x) \rightarrow 0$

## METHOD OF MULTIPLIERS

### (AUGMENTED LAGRANGIAN METHOD)

$$\mathcal{L}_A(x; \lambda, \mu) := f(x) - \lambda^T h(x) + \frac{\mu}{2} \|h(x)\|^2$$

↳ which is composed by  $\begin{cases} \mathcal{L}(x, \lambda) = f(x) - \lambda^T h(x) \\ Q(x, \mu) = f(x) + \frac{\mu}{2} \|h(x)\|^2 \end{cases}$

↳ the basic IDEA is:

- fix  $\mu_k > 0$ ,  $\lambda_k \in \mathbb{R}^n$  and approximately solve  $\min \mathcal{L}_A(x_k, \lambda_k; \mu_k) \rightarrow x_k$
- update  $\mu_k, \lambda_k$  and repeat until convergence

↳ how do you update  $x_k$ ?

↳ the optimality condition for  $\min_x \mathcal{L}_A(x, \lambda_k, \mu_k)$  (P<sub>2</sub>)

$$\text{is } \nabla_x \mathcal{L}_A = 0$$

$$\bullet \nabla_x \mathcal{L}_A = \nabla f(x_k) - \sum_{i=1}^m ((\lambda_k)_i - \mu_k h_i(x_k)) \nabla h_i(x_k) \approx 0$$

↳ while the optimality conditions for  $[\min f(x) \text{ st } h(x) = 0]$  (P<sub>1</sub>)

$$\mathcal{L}(x, \lambda) = f(x) - \lambda^T h(x)$$

$$\text{KKT: } \begin{cases} \nabla_x \mathcal{L} = 0 \\ \nabla_\lambda \mathcal{L} = 0 \end{cases} \rightarrow \begin{cases} \nabla f(x) - \sum_{i=1}^m \lambda_i \nabla h_i(x) \\ -h(x) = 0 \end{cases}$$

↳ let  $x^* = \arg \min f(x)$  and  $\lambda^*$  sat. KKT conditions

$$\text{then } \bullet \nabla f(x^*) - \sum \lambda_i^* \nabla h_i(x) = 0$$

↳ let's look closely at both optimality conditions

$$\bullet \nabla f(x^*) - \sum \lambda_i^* \nabla h_i(x) = 0$$
$$\bullet \nabla f(x_k) - \sum_{i=1}^m ((\lambda_k)_i - \mu_k h_i(x_k)) \nabla h_i(x_k) \approx 0$$

↳ we want that  $\lambda_i^* \approx (\lambda_k)_i - \mu_k h_i(x_k)$   
at least asymptotically

$$h_i(x_k) \approx \frac{(\lambda_k)_i - \lambda_i^*}{\mu_k} \rightarrow \begin{array}{l} \text{the more } (\lambda_k)_i \rightarrow \lambda_i^* \\ \text{the more } h_i(x_k) = 0 \end{array}$$

↳ this suggests the following rule:

$$\lambda_{k+1} = \lambda_k - \mu_k h(x_k)$$

### METHOD OF MULTIPLIERS ALGORITHM

- Given  $\mu_0 > 0$ ,  $x_0^{(0)} \in \mathbb{R}^n$ ,  $\lambda_0 \in \mathbb{R}^m$

- For  $k \geq 0$

- Build  $\mathcal{L}_A(x, \lambda_k; \mu_k)$  and find

$$x_k \approx \operatorname{argmin} \mathcal{L}_A(x; \lambda_k, \mu_k)$$

- Update  $\lambda_{k+1} = \lambda_k - \mu_k h(x)$

$$\mu_{k+1} = \rho \mu_k$$

choose  $x_{k+1}^{(0)}$

## ALTERNATE DIRECTION METHOD OF MULTIPLIERS (ADMM)

↳ augmented Lagrangian for separable problems.

$$\begin{aligned} \min \quad & f_1(x) + f_2(z) \\ \text{st} \quad & Ax + Bz = C \end{aligned}$$

↳ can rewrite constraints as  $D = [A \mid B]$ ,  $w = \begin{bmatrix} x \\ z \end{bmatrix}$   
thus  $Dw = C$

↳ exploit problem structure to simplify iterations

ADMM uses the Augmented Lagrangian Method in 3 steps:

- ①  $\min \mathcal{L}_A(x, z; \lambda, \mu)$  wrt  $x$
- ②  $\min \mathcal{L}_A(x, z; \lambda, \mu)$  wrt  $z$
- ③ update  $\lambda$

$$\mathcal{L}_A(x, z; \lambda, \mu) = f_1(x) + f_2(z) - \lambda^T (Ax + Bz - C) + \frac{\mu}{2} \|Ax + Bz - C\|^2$$

## ADMM ALGORITHM

- Given  $x_0, z_0, \lambda_0, \mu_0 \geq 0$

- For  $k \geq 0$

- Find  $x_{k+1} = \arg \min \mathcal{L}_A(x_k, z_k; \lambda_k, \mu_k)$   
 $= \arg \min f_1(x) + \cancel{f_2(z)} - \lambda_k^T (Ax + Bz - C) + \frac{\mu}{2} \|Ax + Bz - C\|^2$

- Find  $z_{k+1} = \arg \min \mathcal{L}_A(x_k, z_k; \lambda_k, \mu_k)$   
 $= \arg \min \cancel{f_1(x)} + f_2(z) - \lambda_k^T (Ax + Bz - C) + \frac{\mu}{2} \|Ax + Bz - C\|^2$

- Update  $\lambda_{k+1} = \lambda_k - \mu_k (Ax_{k+1} + Bz_{k+1} - C)$

$$\mu_{k+1} = \rho \mu_k$$

$$\text{Choose } x_{k+1}, z_{k+1}$$

## ADMM CONVERGENCE PROPERTIES

↳ under mild assumptions, for  $\kappa \rightarrow \infty$

- $Ax + Bz - c \rightarrow 0$
- $f_1(x_{k+1}) + f_2(z_{k+1}) \rightarrow f_1(x^*) + f_2(z^*)$
- $\lambda_{k+1} \rightarrow \lambda^*$

↳ stronger assumptions are needed for  $x_{k+1} \rightarrow x^*$  and  $z_{k+1} \rightarrow z^*$

↳ rate of convergence: ADMM behaves like a 1<sup>st</sup> order method

## ADMM IN SCALED FORM

↳ using  $\rho = \frac{\lambda}{\mu}$  you can replace:

$$x_{k+1} = \operatorname{argmin}_x f_1(x) - \lambda^T (Ax + Bz - c) + \frac{\mu}{2} \|Ax + Bz - c\|^2$$

with

$$\begin{aligned} x_{k+1} &= \operatorname{argmin}_x f_1(x) + \frac{\mu_k}{2} \|Ax + Bz - c - p_k\|^2 + \frac{\rho}{2\mu_k} \|\lambda_k\|^2 \\ &= \operatorname{argmin}_x f_1(x) + \frac{\mu_k}{2} \|Ax + Bz - c - p_k\|^2 \end{aligned}$$

not def on  $x$

(and similarly for  $z_{k+1}$ )

↳ then only update  $p_{k+1} = p_k - (\lambda x_{k+1} + Bz_{k+1} - c)$

## ADMM APPLICABLE TO LASSO

(LEAST ABSOLUTE SHRINKING AND SELECTION OPERATOR)

$$\min \frac{1}{2} \|Ex - bu\|_2^2 + \gamma \|x\|_1$$

↓

$$\min \frac{1}{2} \|Ex - bu\|_2^2 + \gamma \|z\|_1$$

$$\text{s.t. } x - z = 0$$