

**PASCAL Workshop**  
**Thurnau March 16-18, 2005**

**MODERN CONVEX OPTIMIZATION**

**Arkadi Nemirovski**

**Technion - Israel Institute of Technology**

**nemirovs@ie.technion.ac.il**

- **Convex Programs**
- **Efficient solvability of generic convex programs**
- **Interior-point polynomial time algorithms for well-structured convex programs**
- **Beyond the scope of interior-point algorithms: simple methods for extremely large-scale convex programs**

♣ Geometrically, a convex program is

$$\min_x \{c^T x : x \in X\} \quad [X \subset \mathbf{R}^n : \text{convex set}]$$

♠ In applications, optimization programs usually arise in the form of

$$\min_{u \in \mathbf{R}^n} \{f_0(u) : f_i(u) \leq 0, i = 1, \dots, m\}. \quad (*)$$

When all functions  $f_0(u), \dots, f_m(u)$  are convex,  $(*)$  can be rewritten as

$$\min_{x=(t,u) \in X} t, \quad \underbrace{X = \{(t, u) : f_0(u) \leq t, f_1(u) \leq 0, \dots, f_m(u) \leq 0\}}_{\text{convex set!}}$$

and thus  $(*)$  is a convex program.

$$\min_{u \in \mathbf{R}^n} \{f_0(u) : f_i(u) \leq 0, i = 1, \dots, m\}. \quad (*)$$

♠ The simplest case of a convex program is a Linear Programming program where all  $f_i$  are linear. Such a problem can be posed in the canonical form

$$\min_x \{c^T x : Ax - b \in \mathbf{K} \equiv \mathbf{R}_+^n\}. \quad (\text{LP})$$

◇ Nonlinearity in (\*), if any, “sits” in the objective and in the constraints. It is easily seen that a *convex* program (\*) can be represented in the *conic* form similar to (LP):

$$\min_x \{c^T x : Ax - b \in \mathbf{K}\} \quad (\text{CP})$$

where  $\mathbf{K} \subset \mathbf{R}^N$  is a cone (closed, pointed, convex and with a nonempty interior), and  $x \mapsto Ax : \mathbf{R}^n \rightarrow \mathbf{R}^N$  is a linear embedding.

$$\min_x \{c^T x : Ax - b \in \mathbf{K}\} \quad (\text{CP})$$

♠ An extremely wide variety of convex programs “is captured” by just three generic cones  $\mathbf{K}$ :

♡ Nonnegative orthant  $\mathbf{R}_+^n = \{x \in \mathbf{R}^n : x \geq 0\} \Rightarrow \text{LP}$

♡ Direct products of Lorentz cones  $\mathbf{L}^n = \{(x, t) \in \mathbf{R}^n \times \mathbf{R} : \|x\|_2 \leq t\} \Rightarrow \text{conic quadratic programs}$

$$\min_x \{c^T x : \|A_i x - b_i\|_2 \leq c_i^T x - d_i\}$$

♡ Semidefinite cone  $\mathbf{S}_+^n = \{X \in \mathbf{S}^n : X \succeq 0\} \Rightarrow \text{semidefinite programs}$

$$\min_x \left\{ c^T x : \sum_j x_j A_j - B \succeq 0 \right\}.$$

♠ Specific nice structure of the three generic cones in question underlies modern powerful *primal-dual interior point methods* for Linear, Conic Quadratic and Semidefinite Programming.

♠ There exists a kind of calculus of “conic quadratic and semidefinite representable sets and functions” which offers a systematic way to recognize the possibility to reformulate a convex program as a conic quadratic or semidefinite program.

minimize  $c^T x$  subject to

$$\begin{aligned}
 (a) \quad & \begin{cases} Ax = b \\ x \geq 0 \end{cases} \\
 (b) \quad & \begin{cases} \left( \sum_{i=1}^8 |x_i|^3 \right)^{1/3} \leq x_2^{1/7} x_3^{2/7} x_4^{3/7} + 2x_1^{1/5} x_5^{2/5} x_6^{1/5} \\ 5x_2 \geq \frac{1}{x_1^{1/2} x_2} + \frac{2}{x_2^{1/3} x_3^{5/8} x_4} \\ \begin{pmatrix} x_2 & x_1 & & & \\ x_1 & x_4 & x_3 & & \\ & x_3 & x_6 & x_3 & \\ & & x_3 & x_8 & \end{pmatrix} \preceq 5I \end{cases} \\
 (c) \quad & \begin{cases} \begin{pmatrix} x_1 & x_2 - x_1 & x_3 - x_2 & x_4 - x_3 \\ x_2 - x_1 & x_2 & x_3 - x_2 & x_4 - x_3 \\ x_3 - x_2 & x_3 - x_2 & x_3 & x_4 - x_3 \\ x_4 - x_3 & x_4 - x_3 & x_4 - x_3 & x_4 \end{pmatrix} \succeq 0 \\ x_1 + x_2 \sin(\phi) + x_3 \sin(2\phi) + x_4 \sin(4\phi) \geq 0 \quad \forall \phi \in [0, \frac{\pi}{2}] \end{cases}
 \end{aligned}$$

♠ The problem can be converted, in a systematic way, into SDP

♠ Removing constraints (c), the resulting problem can be converted, in a systematic way, into CQP, and can be approximated, in a polynomial time fashion, by LP.

## Why convex programming?

♠ Convex Programming admits nice and powerful *Duality Theory* capable to certify optimality and to quantify reliably the quality of an approximate solution.

In contrast to this, in a typical nonconvex problems there are no easy ways to certify (global) optimality and to quantify the “level of non-optimality” of an approximate solution.

While typically Convex Duality does not allow to get a solution in a closed analytic form, it still allows for deep understanding of the problem of interest, for building its highly instructive and nontrivial reformulations, etc.

♠ Generic convex programs, under mild computability and boundedness assumptions, are *computationally tractable* - they admit theoretically (and to some extent, also practically) efficient solution methods. In contrast to this, nonconvex optimization programs for which a global solution can be efficiently found numerically are rare exceptions (mainly due to “hidden convexity” which we are lucky to recognize).

## ♣ Convex Duality [conic programming case]

$$\text{Opt} = \min_x \{c^T x : Ax - b \in \mathbf{K}\} \quad (\text{P})$$

♠ The origin of duality is in the desire to find a systematic way to bound from below the optimal value in (CP). In other words, we are interested in a mechanism for deriving from the constraints of (P) their consequences of the form

$$c^T x \geq a.$$

• The simplest way to build a linear inequality which is a consequence of the constraints of (P) is *linear aggregation*

$$Ax - b \in \mathbf{K} \Rightarrow \lambda^T (Ax - b) \geq 0 \quad (*)$$

• Question: What should be a “weight vector”  $\lambda$  in order to make the implication (\*) valid?

• Answer: The necessary and sufficient condition is that  $\lambda^T \xi \geq 0$  for every  $\xi \in \mathbf{K}$ :

$$\lambda \in \mathbf{K}_* \equiv \{\lambda : \lambda^T y \geq 0 \forall y \in \mathbf{K}\}.$$

• Fact: When  $\mathbf{K}$  is a closed, convex and pointed cone with a nonempty interior, so is the *dual cone*  $\mathbf{K}_*$ , and  $(\mathbf{K}_*)_* = \mathbf{K}$ .

$$\text{Opt} = \min_x \{c^T x : Ax - b \in \mathbf{K}\} \quad (\mathbf{P})$$

$$\mathbf{K}_* = \{\lambda : \lambda^T \xi \geq 0 \forall \xi \in \mathbf{K}\}$$

$$Ax - b \in \mathbf{K} \Rightarrow \lambda^T (Ax - b) \geq 0$$

• **Conclusion:** Let  $\lambda \in \mathbf{K}_*$  be such that  $A^T \lambda = c$ . Then  $\text{Opt} \geq b^T \lambda$ .

Equivalently: Let us associate with (P) the dual problem

$$\text{Opt}_* = \max_{\lambda} \{b^T \lambda : A^T \lambda = c, \lambda \in \mathbf{K}_*\} \quad (\mathbf{D})$$

Then  $\text{Opt}_*$  is a lower bound on  $\text{Opt}$ .

**Note:**  $x \mapsto Ax$  is an embedding, so that (P) can be rewritten equivalently as

$$(\mathbf{Pr}): \quad \text{Opt} = \min_{\xi} \{e^T \xi : \xi \in (\mathcal{L} - b) \cap \mathbf{K}\} \quad [e : A^T e = c; \mathcal{L} = \text{Im}A]$$

(D) is of the same geometric structure:

$$(\mathbf{D1}): \quad \text{Opt}_* = \max_{\lambda} \{b^T \lambda : \lambda \in (\mathcal{M} + e) \cap \mathbf{K}_*\} \quad [\mathcal{M} = \text{Ker } A^T = \mathcal{L}^\perp]$$

♡ Geometrically, conic problem is to minimize a linear functional over the intersection of affine plane and a cone. Such a problem is called strictly feasible, if the affine plane intersects the interior of the cone.



$$(P): \quad \text{Opt} = \min_x \{c^T x : Ax - b \in \mathbf{K}\}$$

$$\Updownarrow$$

$$(Pr): \quad \min_{\xi} \{e^T \xi : \xi \in (\mathcal{L} - b) \cap \mathbf{K}\} \quad [e : A^T e = c; \mathcal{L} = \text{Im}A]$$

$$(DI): \quad \max_{\lambda} \{b^T \lambda : \lambda \in (\mathcal{M} + e) \cap \mathbf{K}_*\} \quad [\mathcal{M} = \text{Ker } A^T = \mathcal{L}^\perp]$$

$$\Updownarrow$$

$$(D): \quad \text{Opt}_* = \max_{\lambda} \{b^T \lambda : A^T \lambda = c, \lambda \in \mathbf{K}_*\}$$

♣ **Conic Duality Theorem.** (i) *Conic duality is symmetric: (D) is a conic problem, and the problem dual to (D) is (equivalent to) (P).*

(ii) **[Weak Duality]:**  $\text{Opt}_* \leq \text{Opt}$ , so that whenever  $x$  is primal, and  $\lambda$  is dual feasible, the duality gap

$$c^T x - b^T \lambda = [Ax - b]^T \lambda = [c^T x - \text{Opt}] + \underbrace{[\text{Opt} - \text{Opt}_*]}_{\geq 0} + [\text{Opt}_* - b^T \lambda]$$

is nonnegative.

(iii) **[Strong Duality]** *Let one of the problems (P), (D) be strictly feasible and bounded. Then the other problem is solvable, and  $\text{Opt} = \text{Opt}_*$ .*

*In particular, if both (P), (D) are strictly feasible, both problems are solvable with equal optimal values.*

$$\begin{aligned}
 \text{(P):} \quad & \text{Opt} = \min_x \{c^T x : Ax - b \in \mathbf{K}\} \\
 \text{(D):} \quad & \text{Opt}_* = \max_{\lambda} \{b^T \lambda : A^T \lambda = c, \lambda \in \mathbf{K}_*\}
 \end{aligned}$$

♣ Optimality conditions: Let both (P), (D) be strictly feasible. Then a pair  $(x, \lambda)$  of primal-dual optimal solutions is comprised of optimal solutions to the respective problems

- if and only if

$$c^T x - b^T \lambda = 0 \quad \text{[zero duality gap]}$$

and

- if and only if

$$[Ax - b]^T \lambda = 0 \quad \text{[complementary slackness]}$$

♣ Note: Nonnegative orthant, Lorentz and Semidefinite cones are self-dual. As a result, program dual to an LP/CQP/SDP program is LP/CQP/SDP, respectively.

♣ A generic convex problem  $\mathcal{P}$  is a family of convex programs – instances

$$(p) : \quad \min_x \left\{ f_{(p)}(x) : x \in X_{(p)} \subset \mathbf{R}^{n(p)} \right\} \quad [f_{(p)} : \mathbf{R}^{n(p)} \rightarrow \mathbf{R}]$$

parameterized by finite-dimensional data vectors  $\text{Data}(p)$  and equipped with infeasibility measure  $I_{\mathcal{P}}(x, p)$ .

♡ Infeasibility measure quantifies the infeasibility of a candidate solution  $x \in \mathbf{R}^{n(p)}$  to an instance  $(p)$ . As a function of  $x \in \mathbf{R}^{n(p)}$ , the infeasibility measure  $I_{\mathcal{P}}(x, p)$  should be

- nonnegative everywhere and zero if and only if  $x$  is feasible for  $(p)$
- convex in  $x$ .

♡ The dimension of the data vector of an instance  $(p) \in \mathcal{P}$  is called the size  $\text{Size}(p)$  of the instance.

- **Linear Programming** is a generic problem with instances

$$\min_x \{c^T x : Ax - b \geq 0\},$$

the data vectors being  $(\dim c, \dim b, c, A, b)$ . As an infeasibility measure, one can take

$$I_{LP}(x, p) = \min \{t \geq 0 : Ax - b + t\mathbf{1} \geq 0\}, \quad \mathbf{1} = (1, \dots, 1)^T;$$

- **Conic Quadratic Programming** is a generic problem with instances

$$\min_x \{c^T x : \|A_i x - b_i\|_2 \leq c_i^T x - d_i, i = 1, \dots, m\},$$

the data vectors being  $(\dim x, m, \dim b_1, \dots, \dim b_m, c, \{A_i, b_i, c_i, d_i\}_{i=1}^m)$ . As an infeasibility measure, one can take

$$I_{CQP}(x, p) = \min \{t \geq 0 : \|A_i x - b_i\|_2 \leq c_i^T x - d_i + t, i = 1, \dots, m\};$$

- **Semidefinite Programming** is a generic problem with instances

$$\min_x \left\{ c^T x : \sum_j x_j A_j - B \succeq 0 \right\},$$

data vectors being  $(\dim x, \dim B, c, \{A_j\}_{j=1}^{\dim x}, B)$ . As an infeasibility measure, one can take

$$I_{SDP}(x, p) = \min \left\{ t \geq 0 : \sum_j x_j A_j - B + tI \succeq 0 \right\};$$

- Geometric Programming is a generic problem with instances

$$\min_x \{f_0(x) : f_i(x) \leq 0, i = 1, \dots, m\},$$
$$f_\ell(x) = \ln \left( \sum_{j=1}^k \exp\{a_{ij} + \alpha_{ij}^T x\} \right)$$

data vectors being  $(\dim x, m, k, \{a_{ij}, \alpha_{ij}\}_{i,j})$ . As an infeasibility measure, one can take

$$I_{\mathcal{GP}}(x, p) = \min \{t \geq 0 : f_i(x) \leq t, i = 1, \dots, m\}.$$

♣ A solution algorithm  $\mathcal{A}$  for a generic problem  $\mathcal{P}$  is a code for a Real Arithmetic Computer which, given on input

- the data vector  $\text{Data}(p)$  of an instance

$$(p) : \quad \min_x \{ f_{(p)}(x) : x \in X_{(p)} \subset \mathbf{R}^{n(p)} \}$$

of  $\mathcal{P}$ , and

- a required accuracy  $\epsilon > 0$ ,

returns on output

- either an  $\epsilon$ -solution to  $(p)$ , that is, a vector  $x_\epsilon$  which is  $\epsilon$ -optimal and  $\epsilon$ -feasible for  $(p)$ :

$$f_{(p)}(x_\epsilon) - \text{Opt}(p) \leq \epsilon \ \& \ I_{\mathcal{P}}(x_\epsilon, p) \leq \epsilon,$$

- or a correct conclusion “ $(p)$  is infeasible”,
- or a correct conclusion “ $\text{Opt}(p) = -\infty$ ”.

♣ The  $\epsilon$ -complexity of  $(p) \in \mathcal{P}$  w.r.t. a solution algorithm  $\mathcal{A}$  is the number  $\mathcal{T}_{\mathcal{A}}(\epsilon, p)$  of Real Arithmetic operations required to process the input  $(\text{Data}(p), \epsilon)$

♣ A solution algorithm  $\mathcal{A}$  for a generic problem  $\mathcal{P}$  is called *polynomial time* (“theoretically efficient”), if

$$\mathcal{T}_{\mathcal{A}}(\epsilon, p) \leq \text{Poly} \left( \text{Size}(p), \text{Digits}_{(p)}(\epsilon) \right),$$

where

$$\text{Digits}_{(p)}(\epsilon) = \ln \left( \frac{\text{Size}(p) + \|\text{Data}(p)\|_{\infty} + \epsilon^2}{\epsilon} \right)$$

is the *number of accuracy digits* in an  $\epsilon$ -solution to  $(p)$ .

Interpretation: When a solution algorithm is polynomial time, 10-fold increase in computer power allows

- to increase by an absolute constant factor the *sizes*  $\text{Size}(p)$  of instances which can be solved within a given accuracy in a given time, or
- to increase by an absolute constant factor the number of accuracy digits which can be obtained in a given time on instances of a given size.

♣ A generic problem  $\mathcal{P}$  is called *polynomially solvable* (“computationally tractable”), if it admits a polynomial time solution algorithm.

• **Claim:** *Under mild computability and boundedness assumptions, a generic convex problem is polynomially solvable.*

**In particular, adding to instances of LP/CQP/SDP/GP box constraints**

$$\|x\|_{\infty} \leq R$$

**and treating  $R$  as part of the data vectors, we make the corresponding generic problems polynomially solvable.**



♣ The basic fact underlying polynomial time solvability of generic convex problems is as follows:

Theorem. Consider an optimization program

$$\min_x \{f(x) : x \in B_n \equiv \{x \in \mathbf{R}^n : \|x\|_2 \leq 1\}\},$$

with convex and continuous on  $B_n$  objective  $f$ . There exists an explicit algorithm which, for every  $\epsilon > 0$ , is capable to find a feasible  $\epsilon$ -solution to the problem at the cost of computing the values  $f(x_i)$  and subgradients  $f'(x_i)$  of the objective at

$$N(\epsilon) = O(1)n^2 \ln \left( n + \frac{nV(f)}{\epsilon} \right),$$
$$V(f) = \max_{B_n} f - \min_{B_n} f$$

recursively generated points  $x_i \in \text{int } B_n$ , with  $O(1)n^2$  additional arithmetic operations per every one of these computations.

♠ “Universal” polynomial time algorithms like the one mentioned in the Theorem use *local* information (values and subgradients) on the objective and the constraints and therefore cannot utilize a priori knowledge of problem’s structure. For “well-structured” convex problems, like LP/CQP/SDP, there exist *Interior Point Polynomial Time* algorithms which do utilize problem’s structure and, as a result, yield better complexity bounds and exhibit much better practical performance.

♣ Consider a *canonical cone*

$$\begin{aligned} \mathbf{K} &= \mathbf{L}^{n_1} \times \dots \times \mathbf{L}^{n_p} \times \mathbf{S}_+^{n_{p+1}} \times \dots \times \mathbf{S}_+^{n_{p+q}} \\ &\subset E = \mathbf{R}^{n_1+1} \times \dots \times \mathbf{R}^{n_p+1} \times \mathbf{S}^{n_{p+1}} \times \dots \times \mathbf{S}^{n_{p+q}} \end{aligned}$$

•  $E$  is equipped with the inner product

$$\langle \xi, \lambda \rangle = \xi_1^T \lambda_1 + \dots + \xi_p^T \lambda_p + \text{Tr}(\xi_{p+1} \lambda_{p+1}) + \dots + \text{Tr}(\xi_{p+q} \lambda_{p+q})$$

•  $\mathbf{K}$  is equipped with *canonical barrier*

$$K(\xi) = K_1(\xi_1) + \dots + K_{p+q}(\xi_{p+q}) : \text{int } \mathbf{K} \rightarrow \mathbf{R},$$

where

• for a Lorentz factor  $\mathbf{L}^{n_i} = \{(x, t) \in \mathbf{R}^{n_i} \times \mathbf{R} : \|x\|_2 \leq t\}$ , we set

$$K_i(x, t) = -\ln(t^2 - x^T x),$$

• for a semidefinite factor  $\mathbf{S}_+^{n_i} = \{x \in \mathbf{S}^{n_i} : x \succeq 0\}$  we set

$$K_i(x) = -\ln \text{Det}(x).$$

• We define the *parameter*  $\vartheta(K)$  of barrier  $K$  as twice the number of Lorentz faactors plus the total row size of the semidefinite factors in  $\mathbf{K}$ .

$$\mathbf{K} = \mathbf{L}^{n_1} \times \dots \times \mathbf{L}^{n_p} \times \mathbf{S}_+^{n_{p+1}} \times \dots \times \mathbf{S}_+^{n_{p+q}}$$

$$K(\xi) = K_1(\xi_1) + \dots + K_{p+q}(\xi_{p+q}) : \text{int } \mathbf{K} \rightarrow \mathbf{R}$$

Important facts:

- $\mathbf{K}$  is self-dual;
- The anti-gradient mapping  $\xi \rightarrow -\nabla K(\xi)$  is a one-to-one mapping of  $\text{int } \mathbf{K}$  onto itself which is *self-inverse*:

$$-\nabla K(-\nabla K(\xi)) = \xi \quad \forall \xi \in \text{int } \mathbf{K}.$$

♣ Given a primal-dual pair of strictly feasible conic problems

$$(\mathbf{Pr}): \min_{\xi} \{ \langle e, \xi \rangle : \xi \in [\mathcal{L} - b] \cap \mathbf{K} \} \quad (\mathbf{Dl}): \max_{\lambda} \{ \langle b, \lambda \rangle : \lambda \in [\mathcal{L}^\perp + e] \cap \mathbf{K} \}$$

we associate with the pair

- primal central path  $\xi_*(t) = \operatorname{argmin}_{\xi} \{ t \langle e, \xi \rangle + K(\xi) : \xi \in [\mathcal{L} - b] \cap \mathbf{K} \} \quad [t > 0]$
- dual central path  $\lambda_*(t) = \operatorname{argmin}_{\lambda} \{ -t \langle b, \lambda \rangle + K(\xi) : \lambda \in [\mathcal{L}^\perp + e] \cap \mathbf{K} \} \quad [t > 0]$

♠ Important facts:

- both paths are well-defined on  $(0, \infty)$  and belong to the sets of strictly feasible solutions to the respective problems
- both paths are images of each other under scaled anti-gradient mapping:

$$\lambda_*(t) = -t^{-1} \nabla K(\xi_*(t)); \quad \xi_*(t) = -t^{-1} \nabla K(\lambda_*(t))$$

- Along the primal-dual central path  $(\xi_*(t), \lambda_*(t))$ , the duality gap is  $\vartheta(K)/t$ :

$$\text{DualityGap}(\xi, \lambda) \equiv \langle \xi, \lambda \rangle \equiv [\langle e, \xi \rangle - \text{Opt}(\mathbf{Pr})] + [\text{Opt}(\mathbf{Dl}) - \langle b, \lambda \rangle] \Big|_{\substack{\xi = \xi_*(t) \\ \lambda = \lambda_*(t)}} = \frac{\vartheta(K)}{t}$$

♠ We have seen that as  $t \rightarrow \infty$ , the duality gap along the primal-dual central path goes to 0, so that the path approaches the set of primal-dual optimal solutions. *In primal-dual path-following methods, one traces the path as  $t \rightarrow \infty$ , staying in an appropriately defined neighbourhood of the path and thus approaching primal-dual optimality.*

♠ “Feasible start” generic primal-dual path-following algorithm:

- A pair  $(\xi = \xi_*(t), \lambda = \lambda_*(t))$  on the central path is fully characterized by the following conditions:

$$\begin{array}{l}
 \xi \in [\mathcal{L} - b] \cap \text{int } \mathbf{K} \quad [\text{strict primal feasibility}] \\
 \lambda \in [\mathcal{L}^\perp + e] \cap \text{int } \mathbf{K} \quad [\text{strict dual feasibility}] \\
 \underbrace{t\lambda + \nabla K(\xi)}_{B(\xi, \lambda, t)} = 0 \quad \left[ \begin{array}{l} \text{augmented complementary slackness} \\ \text{in LP: } \lambda_i \xi_i = t^{-1} \quad \forall i \end{array} \right] \quad (*)
 \end{array}$$

- In order to trace the path, one iterates the following scheme:

Given a current triple  $(\xi, \lambda, t)$  with strictly primal-dual feasible  $(\xi, \lambda)$ , we

- choose a target value  $t_+ \geq t$  of the penalty;
- find a *primal-dual search direction*  $(\Delta\xi, \Delta\lambda)$  by linearizing (\*):

$$\Delta\xi \in \mathcal{L}; \quad \Delta\lambda \in \mathcal{L}^\perp; \quad B(\xi, \lambda, t_+) + \frac{\partial B(\xi, \lambda, t_+)}{\partial \xi} \Delta\xi + \frac{\partial B(\xi, \lambda, t_+)}{\partial \lambda} \Delta\lambda = 0$$

- Update the triple

$$(\xi, \lambda, t) \leftarrow (\xi_+ = \xi + \Delta\xi, \lambda_+ = \lambda + \Delta\lambda, t_+)$$

Various primal-dual path-following IP's follow the outlined scheme, utilizing two additional options:

- incorporating line search at the updating step
- exploiting various equivalent forms of the augmented complementary slackness condition.

Augmented Complementary Slackness  $t\lambda + \nabla K(\xi) = 0$  can be rewritten in many equivalent forms (e.g.,  $t\xi + \nabla K(\lambda) = 0$ ). *Outside of the central path, linearizations of various forms of A.C.S. are not equivalent to each other, so that different forms of A.C.S. lead to different path-following methods.*



♠ Example: “Primal” path-following method. With A.C.S. written as  $t\lambda + \nabla K(\xi) = 0$  and the primal-dual pair coming from the program

$$\min_x \{c^T x : Ax - b \in \mathbf{K}\},$$

the outlined scheme results in the essentially purely primal recurrence

$$\begin{aligned} t &\mapsto t_+ > t, \quad x \mapsto x_+ = x - [A^T \nabla^2 K(\underbrace{Ax - b}_{\xi}) A]^{-1} [t_+ c + A^T \nabla K(\xi)] \\ \xi &\mapsto \xi_+ = Ax_+ - b, \quad \lambda \mapsto \lambda_+ = -t_+^{-1} [\nabla K(\xi) + \nabla^2 K(\xi)[\xi_+ - \xi]] \end{aligned} \quad (*)$$

♡ Theorem [short-step primal path-following method] Let (\*) be initialized at a starting point  $x = x_0, t = t_0 > 0$  “close to the path”:

$$[tc + A^T \nabla K(\xi)]^T [A^T \nabla^2 K(\xi) A]^{-1} [tc + A^T \nabla K(\xi)] \leq 0.1^2 \quad [\xi = Ax - b]$$

Then, with the penalty updating policy  $t_+ = t[1 + 0.1\vartheta^{-1/2}(K)]$ , the algorithm is well-defined, keeps all iterates close to the path and, for every  $\epsilon > 0$ , results in primal-dual strictly feasible solution with duality gap  $\leq \epsilon$  in no more than

$$N(\epsilon) = O(1)\sqrt{\vartheta(K)} \ln \left( 2 + \frac{\vartheta(K)}{t_0 \epsilon} \right)$$

steps.

♡ Theorem assumes that the algorithm is started “close to the path”. This can be ensured, *keeping the complexity bound essentially intact*, by a suitable initialization scheme (based on the same path-following techniques as applied to an appropriate auxiliary problem).

♡ The theoretical complexity bounds stated by Theorem are the best known for LP/CQP/SDP. At the same time, the associated worst-case-oriented “short step” penalty updating policy  $t \mapsto t(1 + O(1)\vartheta^{-1/2}(K))$  is too slow from the practical viewpoint.

“Practical” IPm’s do not require “close to the path” (or even feasible) starting point, use more “aggressive” stepsize policies, operate in much wider neighbourhood of the central path and are “symmetric” with respect to primal and dual variables. As a result, practical IPM’s exhibit much better behaviour in actual computations than the short-step primal path-following method (although at the price of “spoiling” the theoretical complexity bound to  $O(1)\vartheta(K) \ln\left(\frac{\vartheta(K)}{\epsilon}\right)$ ).

♣ After two decades of Interior Point Revolution, the entire Convex Programming is “within the reach” of Interior Point Polynomial Time methods

⇒ theoretical (and to some extent – practical) possibility to solve convex programs to *high accuracy with low iteration count*

**However:** When solving programs with  $n$  variables, an IPM iteration requires assembling and solving  $n \times n$  Newton system of linear equations. With standard Linear Algebra, this costs  $O(n^3)$  operations, unless the matrix of the system is highly sparse with favourable sparsity pattern. As a matter of fact,

- typical LPs of real-world origin do lead to sparse Newton systems  
⇒ IPMs are capable to solve LPs with  $10^4 - 10^5$  and even  $10^6$  variables

- typical nonlinear convex programs (especially SDPs) lead to dense Newton systems

⇒ in reality, IPMs can fail to solve a nonlinear convex program with “just”  $10^4$  variables – the very first iteration will last forever...

(!) With design dimension  $n \sim 10^4$ , iteration cost  $O(n^3)$  becomes prohibitively large, and with  $n \sim 10^5$  and more, a linear in  $n$  iteration cost becomes a must...

♣ In the “extremely large-scale case” ( $n$  of order of tens and hundreds of thousands), (!) rules out all advanced convex optimization techniques, including all known polynomial time algorithms. As far as nonsmooth and/or constrained convex problems are concerned, at the present level of our knowledge (!) leaves us with the only option: *first order gradient-type methods*.

♣ A convenient framework for presenting gradient methods is *Convex Programming in saddle point form*, where the problem of interest is

$$\min_{x \in X} \left\{ f(x) = \max_{y \in Y} \phi(x, y) \right\} \quad (S)$$

- $X$  and  $Y$  – compact convex sets in Euclidean spaces
- $\phi(x, y) : X \times Y \rightarrow \mathbf{R}$  – Lipschitz continuous function convex in  $x \in X$  and concave in  $y \in Y$ .

Note: A convex program  $\min_{x \in X} f(x)$  can be in many ways reduced to (S):

- ◇ there always is a trivial possibility  $\phi(x, y) \equiv f(x)$
- ◇ when  $f$  is “well-structured”, other options are possible, e.g.:

- $[x \in \mathbf{R}^n]$   $f(x) \equiv \max_{1 \leq j \leq m} [a_j^T x + b_j] = \max_y \left\{ \sum_i y_i [a_i^T x + b_i] : y \succeq 0, \sum_i y_i = 1 \right\}$

- $[x \in \mathbf{S}^m]$   $f(x) \equiv \lambda_{\max}(x) = \max_y \{ \text{Tr}(yx) : y \succeq 0, \text{Tr}(y) = 1 \}$

$[\lambda_{\max}(x) \equiv \lambda_1(x) \geq \lambda_2(x) \geq \dots \geq \lambda_m(x)]$  are eigenvalues of  $x \in \mathbf{S}^m$

- $[x \in \mathbf{S}^m]$   $f(x) \equiv \lambda_1(x) + \dots + \lambda_k(x) = \max_y \{ \text{Tr}(xy) : 0 \preceq y \preceq I, \text{Tr}(y) = k \}$

- $[x \in \mathbf{S}^m]$   $f(x) \equiv \|x\| = \max_{u, v} \{ \text{Tr}(x[u - v]) : u, v \succeq 0, \text{Tr}(u) = \text{Tr}(v) = 1 \}$

♠ Here,  $f$  is highly nonsmooth and nonlinear, while  $\phi$  is just bilinear!

$$\min_{x \in X} \left\{ f(x) = \max_{y \in Y} \phi(x, y) \right\} \quad (S)$$

♣ In contrast to IPMs, all known first order methods for solving (S) are unable to utilize detailed a priori knowledge of problem's structure and use *black box representation* of (S). In this representation,

- the sets  $X, Y$  are known in advance
- $\phi$  is known to belong to a given family  $\mathcal{F}$  of convex-concave Lipschitz continuous functions on  $Z = X \times Y$
- *quantitative* information is obtained during the solution process via subsequent calls to the *First Order Oracle* which, given on input  $(x, y) \in X \times Y$ , returns  $\phi(x, y)$  along with a subgradient  $\phi'_x(x, y)$  of  $\phi$  in  $x$  and a supergradient  $\phi'_y(x, y)$  of  $\phi$  in  $y$ .

$$f_\phi^* = \min_{x \in X} \left\{ f_\phi(x) = \max_{y \in Y} \phi(x, y) \right\}, \quad \phi \in \mathcal{F} \quad (S)$$

♣ Limits of performance of *black-box-oriented* methods for (S) are given by **Information-Based Complexity Theory**.

♠ In IBC, one considers a **class**  $(X, Y, \mathcal{F})$  of **problems** of the form of (S) associated with  $X, Y$  and a given family  $\mathcal{F}$  of convex-concave functions  $\phi$ . A **solution method**  $\mathcal{B}$  is a collection of rules for generating

- **search points**  $z_t = (x_t, y_t) \in X = X \times Y$  where  $\phi, \phi'$  are computed;
- **approximate solutions**  $x^t \in X$ .

The only restriction on rules is **casuality**:  $z_t$  and  $x^t$  should depend solely on the “past information”  $\{\phi(z_\tau), \phi'(z_\tau)\}_{\tau < t}$ .

♠ The  **$\epsilon$ -complexity** of  $(X, Y, \mathcal{F})$  w.r.t.  $\mathcal{B}$  is the minimal number of steps  $N$  in which  $\mathcal{B}$  solves *all* problems from the class within accuracy  $\epsilon$ :

$$\text{Compl}^{\mathcal{B}}(\epsilon) = \min \{N : f_\phi(x^t) - f_\phi^* \leq \epsilon \forall (\phi \in \mathcal{F}, t \geq N)\}.$$

♠ The  **$\epsilon$ -complexity** of  $(X, Y, \mathcal{F})$  is the best of complexities w.r.t.  $\mathcal{B}$ 's:

$$\text{Compl}(\epsilon) = \min_{\mathcal{B}} \text{Compl}^{\mathcal{B}}(\epsilon).$$

$$f_\phi^* = \min_{x \in X} \{f_\phi(x) = \max_{y \in Y} \phi(x, y)\}; \quad \phi \mapsto \Phi(x, y) = \begin{bmatrix} \phi'_x(x, y) \\ -\phi'_y(x, y) \end{bmatrix} \quad (S)$$

♣ **Information-Based Complexity**  $\text{Compl}(\cdot)$  yields “ultimate” limits on worst-case performance of black-box-oriented methods as applied to a given class of problems  $(S)$ . These limits are known for all “non-parametric” classes of convex problems, including large-scale ones. Here are “large-scale” results related to the case of “Euclidean geometry”:

Let  $X, Y$  be subsets of the unit Euclidean ball  $B_n$  in  $\mathbb{R}^n$ . Then

**A. Nonsmooth case:** If  $\mathcal{F}$  is the set of all convex-concave functions  $\phi$  such that  $\|\Phi(z') - \Phi(z'')\|_2 \leq L$  for all  $z', z'' \in X \times Y$ , then

$$10 (L/\epsilon)^2 \geq \text{Compl}(\epsilon) \geq \underbrace{\frac{1}{10} (L/\epsilon)^2}_{\substack{\text{provided } X = Y = B_n \\ \text{and } n \geq (L/\epsilon)^2}}$$

The lower bound remains valid even in the *minimization* case, i.e., when  $\phi$ 's are further restricted to be independent of  $y$  and to be as simple as  $\phi(x, y) = \max_{1 \leq i \leq n} [\epsilon_i x_i + a_i]$ ,  $\epsilon_i = \pm 1$ .



**B. Smooth saddle point case:** If  $\mathcal{F}$  is the set of all convex-concave functions  $\phi$  such that  $\|\Phi(z') - \Phi(z'')\|_2 \leq L\|z' - z''\|_2$  for all  $z', z'' \in X \times Y$ , then

$$10(L/\epsilon) \geq \text{Compl}(\epsilon) \geq \underbrace{\frac{1}{10}(L/\epsilon)}_{\substack{\text{provided } X = Y = B_n \\ \text{and } n \geq L/\epsilon}}$$

The lower bound remains valid when  $\phi$ 's are further restricted to be bilinear:  $\phi(x, y) = a^T x + x^T A y + b^T y$ .

**C. Smooth minimization case:** If, in addition to smoothness from B,  $\phi$ 's are restricted to be independent of  $y$ , then

$$10(L/\epsilon)^{1/2} \geq \text{Compl}(\epsilon) \geq \underbrace{\frac{1}{10}(L/\epsilon)^{1/2}}_{\substack{\text{provided } X = B_n \text{ and} \\ n \geq (L/\epsilon)^{1/2}}}$$

The lower bound remains valid when  $\phi$ 's are further restricted to be convex quadratic forms of  $x$ :  $\phi(x, y) = x^T A x - 2b^T x$ ,  $A \succeq 0$ .

**Note:** The outlined complexity bounds are dimension-independent!

♣ The outlined IBCT results and their extensions/modifications provide us with

**Bad news:** In large-scale case, gradient type algorithms, same as all black-box-oriented methods, possess slow – sublinear – rate of convergence.

⇒ With gradient type methods, one cannot hope to get high accuracy solutions, provided that  $n$  is large.

However,

*In numerous applications, all we need are medium-accuracy solutions*

*⇒ What is of primary importance in large-scale optimization, is whether the rate of convergence does or does not deteriorate as the dimension grows, and not whether this rate is high or slow.*

**Good news:** In the case of problems with “favourable geometry”, the complexity is dimension-independent (or nearly so)

**More good news:** The (nearly) dimension-independent complexity is yielded by computationally cheap gradient-type methods.

⇒ When solving large-scale problems, gradient methods are natural candidates...

## ♣ Some History:

♠ The very first gradient-type method for solving  $\min_X f(x)$  with non-smooth convex  $f$  – **Subgradient Descent**

$$x_{t+1} = \operatorname{argmin}_{x \in X} \|[x_t - \gamma_t f'(x_t)] - x\|_2^2$$

originates from N. Shor ('63) and B. Polyak ('65) and over years was intensively modified, primarily via utilizing past information (*bundle methods*).

♠ SD is intrinsically adjusted to problems with Euclidean geometry. A substantial generalization of SD, the **Mirror Descent** method [Nem.&Yudin'77, Ben-Tal,Margalit,Nem.'01, Teboulle&Beck'03] allows to adjust the scheme to nonsmooth minimization and saddle point problems on

- simplexes

$$\Delta_n = \{x \in \mathbf{R}^n : x \geq 0, \sum_i x_i \leq 1\}$$

- spectahedrons

$$\Sigma_n = \{x \in \mathbf{S}^n : x \succeq 0, \operatorname{Tr}(x) \leq 1\},$$

etc.

♠ Until recently, common wisdom said that

- The only way to exploit problem's structure in nonlinear convex minimization is offered by IPMs (and thus is too computationally demanding in the extremely large-scale case);

- Computationally cheap gradient type methods are black-box oriented and thus must obey the IBCT limits of performance. Consequently, when solving a large-scale convex program

$$\min_{x \in X} f(x), \quad (*)$$

accuracy after  $t$  steps cannot be better than

- ◇  $O(t^{-1/2})$  for a nonsmooth Lipschitz continuous  $f$  (really slow...)

- ◇  $O(t^{-2})$  for  $C^{1,1}$ -smooth  $f$  (alas, smoothness is a rare commodity...)

♠ A breakthrough in understanding the situation is due to Yu. Nesterov ('03) who observed that if  $f$  is of nice analytic structure, it usually can be represented as  $f(x) = \max_{y \in Y} \phi(x, y)$  with  $C^{1,1}$ -smooth  $\phi$ , and therefore the resulting saddle point problem can be solved by simple methods at the rate  $O(t^{-1})$ .

**Note:** Accelerating convergence from the “common wisdom”  $O(t^{-1/2})$  to  $O(t^{-1})$  indeed makes a difference!

$$\min_{x \in X} \left\{ f(x) = \max_{y \in Y} \phi(x, y) \right\} \quad (S)$$

♣ Nesterov’s  $O(t^{-1})$ -method combines a large-scale optimal algorithm for minimizing *smooth* convex functions with smooth approximation of  $f$ , based on saddle point representation and adjusted from step to step.

We are about to present an alternative  $O(t^{-1})$ -method – **Mirror Prox** – which works directly with (S).

♣ **Setup for MP** is given by  $Z = X \times Y$ , a norm  $\|\cdot\|$  on the space  $E$  where  $Z$  lives and a  $C^1$  strongly convex *distance-generating function*  $\omega(\cdot) : Z \rightarrow \mathbf{R}$ . We associate with the setup data the parameters

$$\Theta = \max_{u, v \in Z} \underbrace{[\omega(u) - \omega(v) - \langle \omega'(v), u - v \rangle]}_{\text{local distance } \omega_v(u)}$$

$$\alpha = \max \left\{ c : \omega_v(u) \geq \frac{c}{2} \|u - v\|^2 \quad \forall u, v \in Z \right\}$$

and the prox mappings

$$P_v(\xi) = \operatorname{argmin}_{z \in Z} \{ \langle \xi, z \rangle + \omega_v(z) \}.$$

♠ “**Implementability Assumption**”:  $Z$  and  $\omega(\cdot)$  are simple and fit each other, so that prox mappings are easy to compute.

$$\left. \begin{array}{c} \min_{x \in X} \left\{ f(x) = \max_{y \in Y} \phi(x, y) \right\} \\ \Updownarrow \\ \max_{y \in Y} \left\{ g(y) = \min_{x \in X} \phi(x, y) \right\} \end{array} \right\} \mapsto \Phi(x, y) = \begin{bmatrix} \phi'_x(x, y) \\ -\phi'_y(x, y) \end{bmatrix} \quad (S)$$

♣ The basic Mirror Prox algorithm is

$$z_{t-1} \mapsto w_t := P_{z_{t-1}}(\gamma_t \Phi(z_{t-1})) \mapsto z_t := P_{z_{t-1}}(\gamma_t \Phi(w_t))$$

$$s^t \equiv (x^t, y^t) = \left( \sum_{\tau=\lfloor t/2 \rfloor}^t \gamma_\tau \right) \sum_{\tau=\lfloor t/2 \rfloor}^t \gamma_\tau w_\tau$$

where  $\gamma_t > 0$  are stepsizes.

**Theorem:** Let  $\Phi(\cdot)$  be Hölder continuous with exponent  $\sigma \in [0, 1]$  and constant  $L$ :

$$\|\Phi(z) - \Phi(z')\|_* \leq L \|z - z'\|^\sigma \quad \forall z, z' \in Z.$$

Then the stepsize policy  $\gamma_t = 0.7L^{-1} \left(\frac{\Theta}{t}\right)^{\frac{1-\sigma}{2}} \alpha^{\frac{1+\sigma}{2}}$  ensures that

$$\epsilon(s^t) \equiv \left[ f(x^t) - \min_X f \right] + \left[ \max_Y g - g(y^t) \right] \leq O(1)L \left( \frac{\Theta}{\alpha t} \right)^{\frac{1+\sigma}{2}}.$$

- **Setup:**  $(\omega(\cdot), Z = X \times Y, \|\cdot\|) \Rightarrow (\Theta, \alpha)$
- **Problem:**  $\min_{x \in X} \max_{y \in Y} \phi(x, y)$
- **Assumption:**  $\|\Phi(z) - \Phi(z')\|_* \leq L\|z - z'\|^\sigma, \Phi = \begin{bmatrix} \phi'_x \\ -\phi'_y \end{bmatrix}$



$$\epsilon(s^t) \leq O(1)L \left( \frac{\Theta}{\alpha t} \right)^{\frac{1+\sigma}{2}}$$

♠ **Nonsmooth case  $\sigma = 0$ :**  $\phi$  is Lipschitz continuous with constant  $L$  w.r.t.  $\|\cdot\| \Rightarrow O(t^{-1/2})$ -rate of convergence

♠ **Smooth case  $\sigma = 1$ :** gradient of  $\phi$  is Lipschitz continuous with constant  $L$  w.r.t.  $\|\cdot\|: \Rightarrow O(t^{-1})$ -rate of convergence

♣ **Adjusting  $(\omega(\cdot), \|\cdot\|)$ ,** one can optimize the efficiency estimate, thus adjusting the method to the geometry of problem in question.

♣ **Good setups for MP** are known when  $X$  and  $Y$  are direct products of (simple subsets) of

- **Euclidean balls**  $\{x \in \mathbf{R}^n : \|x\|_2 \leq R\}$ ,
- **boxes**  $\{x \in \mathbf{R}^n : \|x\|_\infty \leq R\}$ ,
- **“matrix boxes”**  $\{x \in \mathbf{S}^n : -RI \preceq x \preceq RI\}$ ,
- **simplexes**  $\{x \in \mathbf{R}^n : 0 \leq x, \sum_i x_i \leq R\}$ ,
- **spectahedrons**  $\{x \in \mathbf{S}^n : 0 \preceq x, \text{Tr}(x) \leq R\}$ .

♠ **The resulting efficiency estimate for MP** is nearly or fully dimension-independent, provided there are no box-type factors in  $X, Y$



♣ Example: SVM

♠ Problem: Given  $n$  points  $z^i \in \mathbf{R}^d$  partitioned into two sets  $B$  (“blue”) and  $R$  (“red”), find affine form

$$\alpha^T z + \beta = 0$$

which separates best of all the blue and the red points.

♠ Model:

$$\min_{\alpha, \beta} \left\{ \sum_{i=1}^n [1 - y_i(\alpha^T z^i + \beta)]_+ : \|\alpha\|_1 \leq \rho \right\} \quad (*)$$

• We have

$$\begin{aligned} f(\alpha, \beta) &\equiv \sum_i [1 - y_i(\alpha^T z^i + \beta)]_+ = \sum_i \max_{0 \leq \lambda_i \leq 1} \lambda_i [1 - y_i(\alpha^T z^i + \beta)] \\ &= \max_{0 \leq \lambda_i \leq 1} \{ \lambda^T [\mathbf{1} + A\alpha] - \beta \lambda^T y \} \quad [A_{ij} = -y_i z_j^i] \end{aligned}$$

⇓

$$F(\alpha) \equiv \min_{\beta} f(\alpha, \beta) = \max_{\lambda \in \Lambda} \lambda^T [\mathbf{1} + A\alpha], \quad \Lambda = \{ \lambda : 0 \leq \lambda_i \leq 1, \lambda^T y = 0 \}$$

• Thus, (\*) reduces to the saddle point problem

$$\min_{\|\alpha\|_1 \leq \rho} F(\alpha), \quad F(\alpha) = \max_{\lambda \in \Lambda} \lambda^T [\mathbf{1} + A\alpha].$$

♠ With MP, an  $\epsilon$ -solution “costs”  $O(1) \frac{\rho \sqrt{n \ln d} \max_j \|A_j\|_2}{\epsilon}$  multiplications of given vectors by  $A, A^T$ .

$$\min_{\alpha, \beta} \left\{ \sum_{i=1}^n [1 - y_i(\alpha^T z^i + \beta)]_+ : \|\alpha\|_1 \leq \rho \right\} \quad (*)$$

$$\begin{array}{c} \Updownarrow \\ \min_{\|\alpha\|_1 \leq \rho} F(\alpha), \quad F(\alpha) = \max_{\lambda \in \Lambda} \lambda^T [\mathbf{1} + A\alpha] \end{array} \quad (**)$$

$$A = [-y_1 z^1, \dots, -y_n z^n]^T : n \times d, \quad \Lambda = \{\lambda \in \mathbf{R}^n : 0 \leq \lambda \leq \mathbf{1}, y^T \lambda = 0\}$$

### ♣ Numerical illustration:

♠ Training set:  $n = 15,000$  randomly generated sparse vectors (7,506 blue and 7,494 red) in  $d = 5,000$ -dimensional space

$\Rightarrow$  3,142,764 nonzeros in  $15000 \times 5000$  matrix  $A$  (density 0.04)

♠ Generation ensures that the blue and the red vectors admit sparse “near-separator”  $\alpha_*^T z + \beta_*$  which misclassifies  $\approx 6\%$  of vectors.

♠ (\*) can be posed as an LP program with 25,001 variables and 25,001 inequality constraints and can be solved by IP methods.

(\*\*) can be solved by Mirror Prox.

♠ With  $\rho = 3.0$ , the results are as follows:

CPU time	# of matrix-vector multiplications	relative error in objective	classification error	
			training set	testing set

**Mirror Prox solutions:**

62'	114	3.4e-3	9.5%	8.9%
125'	230	1.3e-3	9.5%	10.1%
188'	346	5.0e-4	9.5%	9.7%
251'	462	2.7e-4	9.5%	9.7%
314'	578	1.8e-4	9.5%	9.4%
506'	916	1.8e-4	9.5%	10.0%
568'	1028	5.4e-5	9.5%	9.8%

**Interior Point solution:**

1782'	—	2.e-12	9.5%	10.0%
-------	---	--------	------	-------

Note: “True” separator has 71 nonzeros, MP and IP separators have 60 nonzeros.