

# Convex Optimization for Data Science

*Gasnikov Alexander*

[gasnikov.av@mipt.ru](mailto:gasnikov.av@mipt.ru)

## Lecture 1. Foundation of Convex analysis

October, 2016

## Structure of a course

- Lecture 1. Foundation of Convex analysis
- Lecture 2. Convex optimization and Big Data applications
  - Lectures 3. Complexity of optimization problems  
& Optimal methods for convex optimization problems
- Lecture 4. Stochastic optimization. Randomized methods
- Lectures 5. Primal-duality, regularization, restarts technique, mini-batch  
& Inexact oracle. Universal methods
  - Lecture 6. Gradient-free methods. Coordinate descent

### Projects/Examples:

<https://arxiv.org/find/all/1/all:+gasnikov/0/1/0/all/0/1>

## **Main books:**

*Polyak B.T.* Introduction to optimization. M. Nauka, 1983.

*Bertsekas D.P.* Nonlinear Programming. Belmont. MA: Athena Scientific, 1999.

*Boyd S., Vandenberghe L.* Convex optimization. – Cambridge University Press, 2004.

*Nesterov Yu.* Introduction Lectures on Convex Optimization. A Basic Course. Applied Optimization. – Springer, 2004.

*Nocedal J., Wright S.* Numerical optimization. – Springer, 2006.

*Nemirovski A.* Lectures on modern convex optimization analysis, algorithms, and engineering applications. – Philadelphia: SIAM, 2013.

*Bubeck S.* Convex optimization: algorithms and complexity // In Foundations and Trends in Machine Learning. – 2015. – V. 8. – no. 3-4. – P. 231–357.

*Evtushenko Yu.G.* <http://www.ccas.ru/personal/evtush/p/198.pdf>

<https://www.youtube.com/user/PreMoLab> (see courses of Yu.E. Nesterov and A.V. Gasnikov)

## Structure of Lecture 1

- Convex functions. Main properties. CVX
- Lagrange multipliers principle and vicinities
- Demyanov–Danskin’s formula, sensitivity in optimization
- Dual problem. Sion–Kakutani theorem ( $\min \max = \max \min$ )
- KKT-theorem. The role of convexity (sufficient condition)

### Examples:

- Conic duality – dual representation (Robust optimization)
- Constrained primal problem leads to unconstrained dual one
  - Implicit primal problem leads to explicit dual one
- Non separable large-dimensional primal problem leads to small dimensional dual problem on a bounded convex set (Slater’s arguments)

- **Convex set**

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

**Convex function**

$$\forall \alpha \in [0,1], x, y \rightarrow f(\alpha x + (1-\alpha)y) \leq \alpha f(x) + (1-\alpha)f(y)$$

**Examples of convex functions**

**i)**  $f(x) = \ln \left( \sum_{k=1}^n \exp(\langle a_k, x \rangle) \right)$  // by direct investigation of the Hessian

**ii)**  $f(x) = x_{[1]} + \dots + x_{[k]}$  (sum of largest  $k$  entries)

**iii)**  $f(X) = \lambda_{\max}(X)$ ,  $X = X^T$

**iv)**  $f(X) = \ln \det X^{-1}$ ,  $X \succ 0$

**v)**  $f(x, Y) = \langle x, Y^{-1}x \rangle$ ,  $Y \succ 0$

# Main Toolbox:

<http://cvxr.com/cvx/>

The screenshot shows a web browser window with the URL `cvxr.com/cvx/`. The page header indicates "Version 2.1, October 2016, Build 1112". There are two yellow callout boxes: the first mentions a video introduction by Professor Stephen Boyd, and the second announces "CVX 3.0 beta" with new features. The main text describes CVX as a Matlab-based modeling system for convex optimization, providing an example model:

$$\begin{aligned} & \text{minimize} && \|Ax - b\|_2 \\ & \text{subject to} && Cx = d \\ & && \|x\|_\infty \leq e \end{aligned}$$

Below the model, it states that the following code segment generates and solves a random instance of this model:

```
m = 20; n = 10; p = 4;
A = randn(m,n); b = randn(m,1);
C = randn(p,n); d = randn(p,1); e = rand;
cvx_begin
    variable x(n)
    minimize( norm( A * x - b, 2 ) )
    subject to
        C * x == d
        norm( x, Inf ) <= e
cvx_end
```

The browser's taskbar at the bottom shows several open files: `pdf_-_Markov_Roz...pdf`, `Convex optimizati...doc`, `OPT-Gasnikov (1).doc`, `Redko_10_22.pdf`, and `Redko_10_22.pdf`. The system tray shows the time as 20:32 on 22.10.2016.

- CVX can solve the following convex optimization problems

$$\langle c, x \rangle + \langle d, y \rangle \rightarrow \min_{(x,y): A \begin{bmatrix} x \\ y \end{bmatrix} = b, \begin{bmatrix} x \\ y \end{bmatrix} \in K},$$

where  $K$  – is a product of convex cones. Typically these cones:  $\mathbb{R}_+^n$ ,  $S_+^n$ ,  $L_2^n$  (positive cone, positive semidefinite cone, Lorentz cone).

- Many convex functions have cone representation (Nesterov–Nemirovski)

$$f(x) = \min_{y: A \begin{bmatrix} x \\ y \end{bmatrix} = b, \begin{bmatrix} x \\ y \end{bmatrix} \in K} \langle c, x \rangle + \langle d, y \rangle \quad ! \text{ (e.g. i) – v)}$$

### Other Toolboxes:

- <https://www-01.ibm.com/software/commerce/optimization/cplex-optimizer/> CPLEX (IBM product): Large-scale LP, QP
- <https://www.mosek.com/> Large-scale convex optimization problems

## How to obtain convex functions?

**Lemma 1.** Let  $f(x)$ ,  $g(x)$  convex,  $h(z)$  convex increasing. Then  $\alpha f(x) + \beta g(x)$  ( $\alpha, \beta \geq 0$ ),  $f(Ay + b)$ ,  $h(f(x))$  are also convex functions.

**Lemma 2.** Let  $G(x, y)$  – is convex function as a function of  $x$  for all  $y \in Y$ . Assume that problem  $\max_{y \in Y} G(x, y)$  is solvable for all  $x$ . Then  $f(x) = \max_{y \in Y} G(x, y)$  is convex. **Example.**  $f(x) = \|Ax - b\|_\infty$  is convex.

**Lemma 3.** Let  $G(x, y)$  – is convex function as a function of  $x$  and  $y$  on the convex set  $Q$ . Assume that problem  $\min_{y: (x, y) \in Q} G(x, y)$  is solvable for all  $x$ . Then  $f(x) = \min_{y: (x, y) \in Q} G(x, y)$  is convex. **Example.**  $f(x) = \inf_{y \in Q} \|x - y\|$  (where  $Q$  is a convex set) is convex. **How can be calculated  $\partial f(x)$ ?**



## Demyanov–Danskin’s formula

Let  $G(x, y)$  and

$$\boxed{f(x) = \max_y G(x, y)} \left( f(x) = \min_y G(x, y) \right)$$

are smooth enough functions. Assume that there exists  $y(x)$  such that

$$G(x, y(x)) = \max_y G(x, y) \left( G(x, y(x)) = \min_y G(x, y) \right).$$

Then  $\boxed{\nabla f(x) = \nabla_x G(x, y(x))} = \left\{ \frac{\partial G(x, y)}{\partial x_i} \right\}_i \Big|_{y=y(x)} \quad . // \quad f_x = G_x + \underbrace{G_y}_{0} y_x = G_x$

$$\boxed{\partial f(x) = \text{convex hull} \bigcup_{\tilde{x}: y(\tilde{x})=y(x)} \partial_x G(\tilde{x}, y(\tilde{x}))} \text{ (convex case).}$$

## Schur's complement

Using Lemma 3 and Demyanov–Danskin's formula

$$G(x, y) = \langle x, Cx \rangle + \langle y, Ay \rangle + 2\langle Bx, y \rangle,$$

one can show that if  $A \succ 0$  (strictly, i.e.  $\forall x \neq 0 \rightarrow \langle x, Ax \rangle > 0$ ) and

$$\begin{bmatrix} A & B \\ B^T & C \end{bmatrix} \succ 0. // \text{ this means that } G(x, y) \text{ is convex}$$

Then

$$C - B^T A^{-1} B \succ 0.$$

Indeed here we can find explicitly  $y(x) = \arg \min_y G(x, y) = -A^{-1} Bx$ . Hence

$$f(x) = \min_y G(x, y) = G(x, y(x)) = \langle x, (C - B^T A^{-1} B)x \rangle$$

is a convex function due to Lemma 3.

## Lagrange multipliers principle and Implicit function theorem

We have a sufficiently smooth optimization problem

$$f(x, y) \rightarrow \min_{g(x, y)=0},$$

where  $g : \mathbb{R}^{n+m} \rightarrow \mathbb{R}^m$ . Assume that implicit function theorem can be applied for  $g(x, y)=0$ . That is, there exists smooth  $y(x)$  such that  $g(x, y(x)) \equiv 0$ . If  $(x_*, y_*)$  is a solution of initial optimization problem then

$$\exists \lambda : \nabla_x L(x_*, y_*, \lambda) = 0, \nabla_y L(x_*, y_*, \lambda) = 0, \quad (*)$$

where

$$L(x, y, \lambda) = f(x, y) + \langle \lambda, g(x, y) \rangle,$$

$\lambda$  can be found from  $g(x_*(\lambda), y_*(\lambda)) = 0$  where  $x_*(\lambda), y_*(\lambda)$  satisfy (\*).

## Implicit function theorem

$$g(x, y(x)) \equiv 0 \Rightarrow g_x + g_y y_x \equiv 0 \Rightarrow y_x = -g_y^{-1} g_x$$

## Fermat principle

$$\frac{d}{dx} f(x, y(x)) = 0 \Rightarrow f_x + f_y y_x = 0 \Rightarrow f_x - f_y g_y^{-1} g_x = 0$$

## Lagrange multipliers principle

$$\left. \begin{array}{l} 0 = L_x = f_x + \lambda^T g_x \\ 0 = L_y = f_y + \lambda^T g_y \end{array} \right\} \times \left( -g_y^{-1} g_x \right) \Bigg\} + \Rightarrow$$

$$f_x - f_y g_y^{-1} g_x = 0. \quad (**)$$

Equality (\*\*) with  $g(x, y) = 0$  allows us to find extremum (optimal point).

## Sensitivity in optimization

We have a sufficiently smooth optimization problem

$$f(x) \rightarrow \min_{g(x)=b},$$

where  $g : \mathbb{R}^n \rightarrow \mathbb{R}^m$  ( $m < n$ ). Let

$$L(x, \lambda) = f(x) + \langle \lambda, b - g(x) \rangle.$$

For optimal solution  $x(b)$  there exists such  $\lambda(b)$  that

$$L_x(x(b), \lambda(b)) = 0. // g(x(b)) = b \quad (*)$$

$$\boxed{\nabla F(b) = \lambda(b)}, \quad \boxed{F(b) = \min_{g(x)=b} f(x)}.$$

Indeed,  $F(b) = L(x(b), \lambda(b), b)$ ;  $F_b = L_x x_b + L_\lambda \lambda_b + L_b = L_b = \lambda(b)$  due to (\*).

## Lagrange multipliers principle and Separation theorem (Convex case)

We have a convex optimization problem ( $g : \mathbb{R}^n \rightarrow \mathbb{R}^m$ )

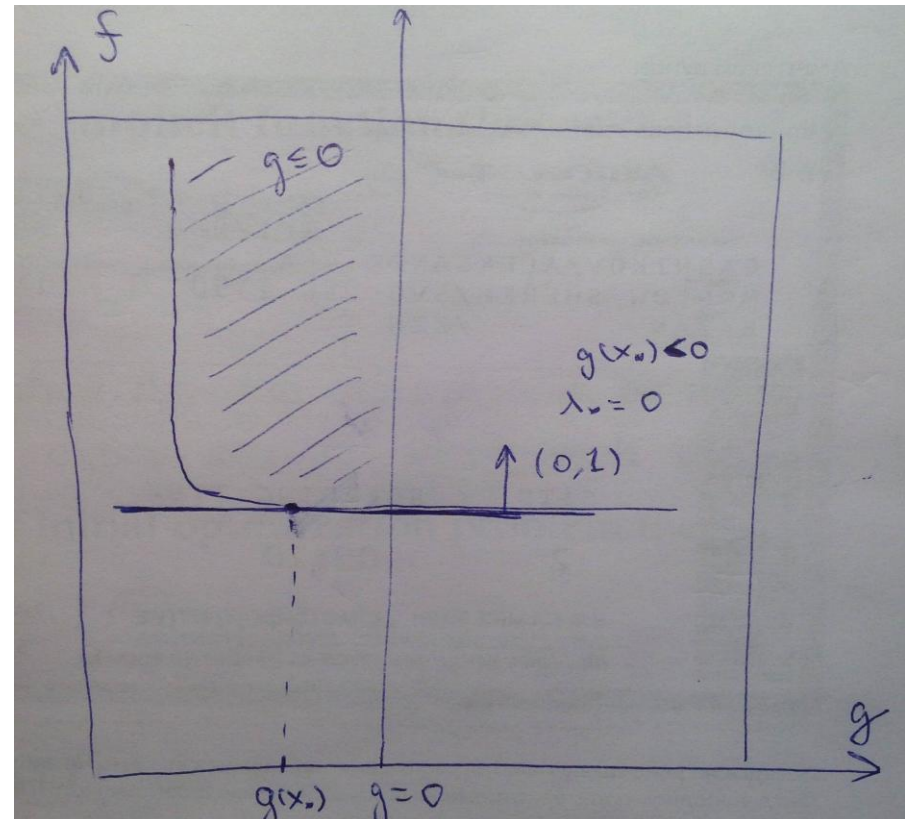
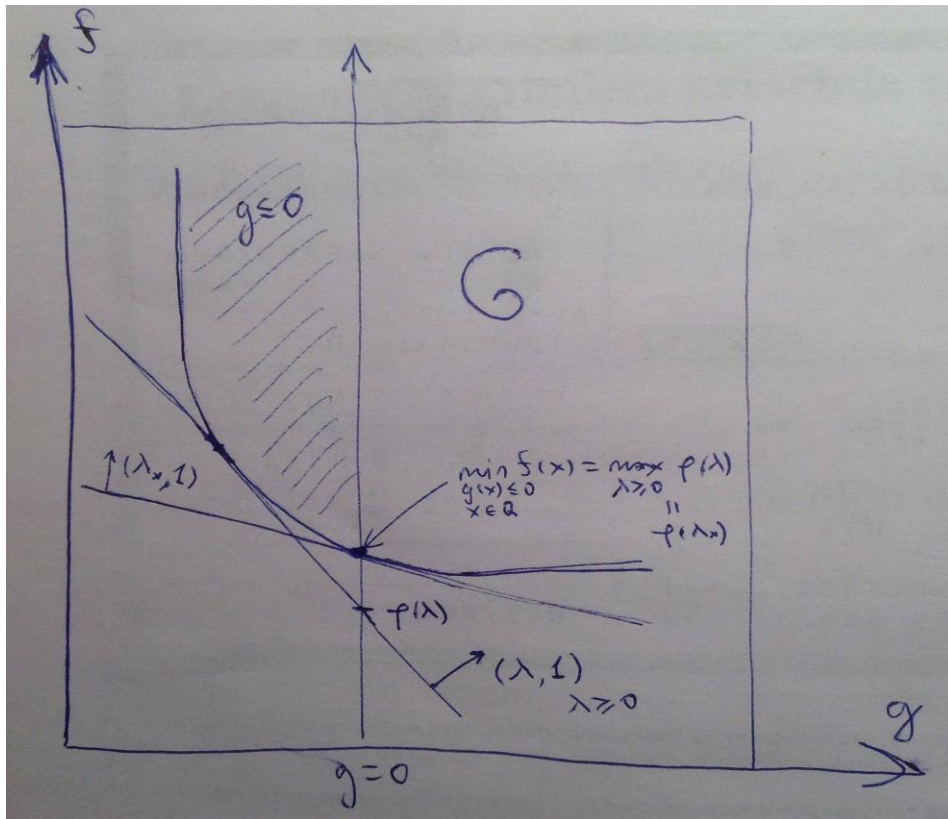
$$f(x) \rightarrow \min_{g(x) \leq 0, x \in Q}.$$

The (Pareto) set  $G = \{(g(x), f(x)), x \in Q\} \oplus \mathbb{R}_+^{m+1}$  is close convex set in the space of pairs  $(g, f)$  ([http://stanford.edu/~boyd/cvxbook/bv\\_cvxslides.pdf](http://stanford.edu/~boyd/cvxbook/bv_cvxslides.pdf)).

i)  $f + g\lambda = \varphi(\lambda)$ , where  $\lambda \geq 0$  – be a tangent hyperplane to  $G$  iff

$$\varphi(\lambda) = \min_{x \in Q} \{f(x) + \langle \lambda, g(x) \rangle\} \text{ (dual function).}$$

ii) Lagrange multipliers principle can be obtained from the Separation theorem for the set  $G$  and hyperplane from i).



iii) For arbitrary  $x \in Q$ ,  $g(x) \leq 0$ ,  $\lambda \geq 0$

$$f(x) \geq \varphi(\lambda)$$

$$\underbrace{\min_{g(x) \leq 0, x \in Q} f(x)}_{\text{primal problem}} \geq \underbrace{\max_{\lambda \geq 0} \varphi(\lambda)}_{\text{dual problem}}. \text{ (weak duality)}$$

Typically for convex problem we have

$$\boxed{\min_{g(x) \leq 0, x \in Q} f(x) = \max_{\lambda \geq 0} \varphi(\lambda)}. \text{ (strong duality)}$$

**Example (LP).**  $\langle c, x \rangle \rightarrow \min_{Ax=b, x \geq 0}$ . Dual problem:  $\langle b, \lambda \rangle \rightarrow \max_{c-A^T \lambda \geq 0}$ .

$$\text{Hint: } \min_{x \geq 0} \left\{ \langle c, x \rangle + \max_{\lambda} \langle b - Ax, \lambda \rangle \right\} = \max_{\lambda} \left\{ \langle b, \lambda \rangle + \min_{x \geq 0} \langle c - A^T \lambda, x \rangle \right\}.$$

If primal problem is compatible then we have strong duality (otherwise dual problem reach infinity).

**What are the sufficient conditions for there is no duality gap (strong duality)?**

**Strong duality (for convex problem):** there exists non-vertical supporting hyperplane at  $(g, f) = (0, f_*)$ .



## Slater's condition

We introduce

$$Q_{\bar{\lambda}} = \left\{ \lambda \in \mathbb{R}_+^m : \varphi(\lambda) \geq \varphi(\bar{\lambda}) \right\}.$$

Assume that Slater's condition is true (sufficient condition in KKT):

*there exists such  $\bar{x} \in Q$  that  $g(\bar{x}) < 0$  ( $\gamma = \min_{i=1, \dots, m} \{-g_i(\bar{x})\}$ ).*

Then

$$\|\lambda_*\|_1 \leq \max_{\lambda \in Q_{\bar{\lambda}}} \|\lambda\|_1 = \max_{\lambda \in Q_{\bar{\lambda}}} \sum_{i=1}^m |\lambda_i| \leq \frac{1}{\gamma} (f(\bar{x}) - \varphi(\bar{\lambda})).$$

**Hint:**  $\varphi(\bar{\lambda}) \leq \varphi(\lambda) = \min_{x \in Q} \left\{ f(x) + \sum_{i=1}^m \lambda_i g_i(x) \right\} \leq f(\bar{x}) + \sum_{i=1}^m \lambda_i g_i(\bar{x}).$

## Quadratic Programming

Let's consider quadratic convex ( $C \succ 0$ ) optimization problem ( $Ax \leq b$  is solvable)

$$\langle x, Cx \rangle + \langle d, x \rangle \rightarrow \min_{Ax \leq b}.$$

Dual problem (we have strong duality because of affine restrictions)

$$\langle \lambda, \tilde{C}\lambda \rangle + \langle \tilde{d}, \lambda \rangle + \tilde{c} \rightarrow \min_{\lambda \geq 0},$$

$$\tilde{C} = AC^{-1}A^T, \tilde{d} = AC^{-1}d + b, \tilde{c} = \frac{1}{2} \langle d, C^{-1}d \rangle.$$

The solutions of primal and dual problems ( $x_*$  and  $\lambda_*$ ) satisfy

$$Cx_* + A^T \lambda_* + d = 0.$$

Typically, one can explicitly build dual problem iff one can explicitly connect primal and dual variables.

## Dual relaxation (strong duality)

For the problem ( $A \neq 0, A^T = A$ )

$$f(x) = \langle x, Ax \rangle \rightarrow \min_{\|x\|_2^2 \leq 1}$$

we have the following dual problem ( $C \succ 0$  means  $\forall x \rightarrow \langle x, Cx \rangle \geq 0$ )

$$\varphi(\lambda) = \begin{cases} -\lambda, & A + \lambda I \succ 0 \\ -\infty, & \text{otherwise} \end{cases} \rightarrow \max_{\lambda}$$

There is no duality gap in this situation. That is we have strong duality

$$\min_{x \in Q} f(x) = \max_{\lambda} \varphi(\lambda)$$

**Note:**  $f_* = \lambda_{\min}$ ,  $x_*$  – eigen vector that corresponds  $\lambda_{\min}$ .

## Dual relaxation (weak duality)

Consider NP-hard two-way partitioning problem

$$f(x) = \langle x, Wx \rangle \rightarrow \min_{x_i^2=1, i=1, \dots, n} .$$

The dual problem is

$$\varphi(\lambda) = \begin{cases} -\sum_{i=1}^n \lambda_i, W + \text{Diag} \{ \lambda_i \}_{i=1}^n \succ 0 & \rightarrow \max_{\lambda} . \\ -\infty, \text{ otherwise} \end{cases}$$

Using weak duality one can show that

$$\min_{x_i^2=1, i=1, \dots, n} \langle x, Wx \rangle \geq n \lambda_{\min}(W).$$

## Necessary and sufficient conditions for convex programming

We have a convex optimization problem

$$f(x) \rightarrow \min_{g(x) \leq 0, Ax=b, x \in Q}.$$

Introduce:  $L(x, \lambda_0, \lambda, \mu) = \lambda_0 f(x) + \langle \lambda, g(x) \rangle + \langle \mu, Ax - b \rangle.$

Initial problem can be equivalently reformulated as

$$\sup_{\lambda \geq 0, \mu} L(x, 1, \lambda, \mu) \rightarrow \min_{x \in Q}.$$

Assume that:  $\min_{x \in Q} \sup_{\lambda \geq 0, \mu} L(x, 1, \lambda, \mu) = \sup_{\lambda \geq 0, \mu} \min_{x \in Q} L(x, 1, \lambda, \mu).$

Then  $\min_{x \in Q} L(x, 1, \lambda, \mu) \rightarrow \sup_{\lambda \geq 0, \mu} .$  (dual problem)

## Complementary slackness

$$L(x, \lambda_0, \lambda, \mu) = \lambda_0 f(x) + \langle \lambda, g(x) \rangle + \langle \mu, Ax - b \rangle,$$

$$\varphi(\lambda, \mu) = \min_{x \in Q} L(x, 1, \lambda, \mu) = L(x(\lambda, \mu), 1, \lambda, \mu) \rightarrow \sup_{\lambda \geq 0, \mu} .$$

From the optimality conditions:  $\varphi_\lambda = 0$  or  $\lambda_* = 0$ ,  $\varphi_\lambda \leq 0$ ;  $\varphi_\mu = 0$ . Since

$$\varphi_\lambda = L_x x_\lambda + L_\lambda = L_\lambda = g; \varphi_\mu = L_x x_\mu + L_\mu = L_\mu = Ax(\lambda, \mu) - b$$

one can obtain that

$$A \cdot x(\lambda, \mu) = b, \quad g(x(\lambda, \mu)) \leq 0.$$

Solution of the dual problem  $\lambda_* > 0$  only if  $g(x(\lambda_*, \mu_*)) = 0$ .

Let's remove  $Ax = b$ ,  $\lambda_0 = 0$  means that supporting hyperplane to  $G$  at  $(g, f) = (0, f_*)$  is vertical.

## Karush–Kuhn–Tucker theorem (KKT)

$$f(x) \rightarrow \min_{g(x) \leq 0, Ax=b, x \in Q} (g : \mathbb{R}^n \rightarrow \mathbb{R}^m)$$

**Necessary condition.** Let  $x_* \in Q$  – is a solution of the problem.

Then there exist  $\lambda_0 \geq 0, \lambda_i \geq 0, i = 1, \dots, m, \mu$  such that:

i)  $Ax_* = b, \lambda_i g_i(x_*) = 0, i = 1, \dots, m;$  (complementary slackness)

ii)  $\min_{x \in Q} L(x, \lambda_0, \lambda, \mu) = L(x_*, \lambda_0, \lambda, \mu).$

**Sufficient condition.** Let  $x_*$  satisfy  $g(x_*) \leq 0, Ax_* = b, x_* \in Q$  and i), ii) with  $\lambda_0 > 0$  (equivalent  $\lambda_0 = 1$ ) then  $x_*$  – is a solution of the problem.

## Saddle point (min max = max min)

We say that  $(x_*, \lambda_*)$  is a saddle point of  $L(x, \lambda)$  iff for all admissible  $(x, \lambda)$

$$L(x, \lambda_*) \geq L(x_*, \lambda_*) \geq L(x_*, \lambda).$$

$$\bar{L}(x) = \sup_{\lambda \in \Lambda} L(x, \lambda) \rightarrow \min_{x \in X} \quad \text{(i)} \quad \underline{L}(\lambda) = \inf_{x \in X} L(x, \lambda) \rightarrow \max_{\lambda \in \Lambda} \quad \text{(ii)}$$

- 1)  $\max_{\lambda \in \Lambda} \underline{L}(\lambda) \leq \min_{x \in X} \bar{L}(x)$  (if there exists saddle point then we have “=”);
- 2) Saddle point set of  $L(x, \lambda)$  coincide with pairs of solutions (i), (ii).

## Sion–Kakutani minimax theorem

*Let continuous function  $L(x, \lambda)$  is convex in  $x$  and concave in  $\lambda$ . The sets  $X$  and  $\Lambda$  are convex and  $X$  is compact. Then  $\max_{\lambda \in \Lambda} \underline{L}(\lambda) = \min_{x \in X} \bar{L}(x)$ .*



## Neyman–Pearson’s Lemma

Assume we have random samples  $x = (x_1, \dots, x_n)$  and two hypotheses about the probability nature of these samples ( $L(\cdot)$  is likelihood function)

$$H_0 : L(x|H_0); H_1 : L(x|H_1).$$

Let’s introduce decision rule:  $0 \leq \varphi(x) \leq 1$  – probability to decide in favor of hypothesis  $H_1$  if one observe vector of samples  $x$ . We’d like to find the best  $\varphi(x)$  in terms of the following infinite dimensional LP-problem:

$$\beta = P(H_0|H_1) = 1 - \int \varphi(x)L(x|H_1)dx \rightarrow \min_{P(H_1|H_0) \leq \alpha};$$

$$P(H_1|H_0) = \int \varphi(x)L(x|H_0)dx \quad |\lambda > 0, 0 \leq \varphi(x) \leq 1.$$

We can solve this problem by means of Lagrange multipliers principle

$$L(\varphi(\cdot), \lambda) = 1 - \int \varphi(x) L(x|H_1) dx + \lambda \cdot \left( \int \varphi(x) L(x|H_0) dx - \alpha \right) \rightarrow \min_{0 \leq \varphi(x) \leq 1},$$

$$\int \varphi(x) \left( \lambda L(x|H_0) - L(x|H_1) \right) dx \rightarrow \min_{0 \leq \varphi(x) \leq 1},$$

$$\varphi(x) = \begin{cases} 1, & \Lambda(x) > \lambda \\ p(x), & \Lambda(x) = \lambda, \Lambda(x) = \frac{L(x|H_1)}{L(x|H_0)}, \\ 0, & \Lambda(x) < \lambda \end{cases}$$

where

$$\int_{\Lambda(x) > \lambda} L(x|H_0) dx + \int_{\Lambda(x) = \lambda} p(x) L(x|H_0) dx = \alpha.$$

Multiplier  $\lambda$  is determined from here in a unique manner, and  $\beta$  doesn't depend on the choice of  $p(x)$ . For more general facts about hypothesis testing via optimization see [http://www2.isye.gatech.edu/~nemirovs/StatOpt\\_LN.pdf](http://www2.isye.gatech.edu/~nemirovs/StatOpt_LN.pdf)

## Robust Optimization (A. Nemirovski)

$$\langle c, x \rangle \rightarrow \min_{\langle \alpha, x \rangle \leq \beta, \alpha \in A} .$$

The set of vectors  $A$  can be infinitely big (e.g. in robust optimization  $A$  is a box). So in general we have a LP problem with infinitely many constraints. Can we solve this problem efficiently? The answer is YES if the set  $A$  has a Fourier–Motzkin representation. This mean the following

$$A = \{ \alpha : \exists u : A\alpha + Bu + b = 0; C\alpha + Du + e \in K \},$$

where  $K$  is a convex cone (with simple enough dual cone  $K^*$ ).

Note that all reasonable convex set have such representation, e.g. boxes,

$$X = \{ x \in \mathbb{R}_+^3 : x_1 x_2 x_3 x_4 \geq 1 \}, X = \{ X \in \mathbb{R}^{m \times n} : \|X\|_{\text{nuclear}} \leq 1 \}.$$

## Farkas' Lemma (& Conic duality)

Assume  $\exists (\bar{\alpha}, \bar{u}) : A\bar{\alpha} + B\bar{u} + b = 0; C\bar{\alpha} + D\bar{u} + e \in \text{icr } K$  or  $K = \mathbb{R}_+^m$ . Then

$\langle \alpha, x \rangle \leq \beta$  for all  $\alpha \in A$  *iff*

$$A^T \mu + C^T \lambda + x = 0, B^T \mu + D^T \lambda = 0, \langle b, \mu \rangle + \langle e, \lambda \rangle \leq \beta$$

is compatible in  $(\mu, \lambda)$ , where  $\lambda \in K^* = \{ \lambda : \langle \lambda, y \rangle \geq 0 \text{ for all } y \in K \}$ .

**Show the transition “ $\Leftarrow$ ”.** For all  $\alpha \in A, \lambda \in K^*$  we have

$$0 = \langle A\alpha + Bu + b, \mu \rangle, 0 \leq \langle C\alpha + Du + e, \lambda \rangle.$$

Hence

$$\underbrace{\langle -A^T \mu - C^T \lambda, \alpha \rangle}_x + \underbrace{\langle -B^T \mu - D^T \lambda, u \rangle}_0 \leq \underbrace{\langle b, \mu \rangle + \langle e, \lambda \rangle}_{\leq \beta}.$$

Using Farkas' lemma we can reformulate initial LP optimization problem:

$$\langle c, x \rangle \rightarrow \min_{(x, \mu, \lambda) \in X},$$

where the convex set  $X$  has the following Fourier–Motzkin representation (reduced to cone representation, since that the problem can be efficiently solved for example by S. Boyd's CVX)

$$X = \left\{ (x, \mu, \lambda) : \begin{aligned} &A^T \mu + C^T \lambda + x = 0, \\ &B^T \mu + D^T \lambda = 0, \langle b, \mu \rangle + \langle e, \lambda \rangle \leq \beta, \lambda \in K^* \end{aligned} \right\}.$$

Details can be found in the book:

*Ben-Tal A., Ghaoui L.El., Nemirovski A. Robust optimization. – Princeton University Press, 2009.*

## Arbitrage free theorem

**Auxiliary fact:** One and only one is true

$$\exists x: Ax \geq 0, Ax \neq 0 \text{ (i)} \quad \text{vs} \quad \exists y > 0: y^T A = 0. \text{ (ii)}$$

**Hint:** (i)  $\Rightarrow \forall y > 0 \rightarrow 0 < \langle y, Ax \rangle = \langle A^T y, x \rangle \neq 0 \Rightarrow$  (ii) is false; (ii)  $\Rightarrow$

Convex hull of rows of  $A$  is hyperplane or whole space  $\Rightarrow$  (i) is false.

$$\neg \exists x: Ax \geq 0 (Ax \neq 0), \quad A = \begin{pmatrix} S \cdot u & C_u \\ S \cdot d & C_d \\ S & C \\ -S & -C \end{pmatrix} \Leftrightarrow C = \frac{1-d}{u-d} C_u + \frac{u-1}{u-d} C_d,$$

where we know  $d < 1 < u$ ,  $C_u, C_d$ .

*Ross S.* An elementary introduction to mathematical finance. Cambridge University Press, 2011.

## Computation of Wasserstein barycenter (M. Cuturi et al.)

$$\begin{aligned}
 H_W(L) &= \min_{\substack{\sum_{j=1}^n x_{ij}=L_i, \sum_{i=1}^n x_{ij}=W_j \\ x_{ij} \geq 0, i, j=1, \dots, n}} \left\{ \gamma \sum_{i,j=1}^n x_{ij} \ln x_{ij} + \sum_{i,j=1}^n c_{ij} x_{ij} \right\} = \\
 &= \max_{\lambda, \mu} \left\{ \langle \lambda, L \rangle + \langle \mu, W \rangle - \gamma \sum_{i,j=1}^n \exp\left(\frac{-c_{ij} + \lambda_i + \mu_j}{\gamma} - 1\right) \right\} = \\
 &= \max_{\lambda} \left\{ \underbrace{\langle \lambda, L \rangle - \gamma \sum_{j=1}^n W_j \ln \left( \frac{1}{W_j} \sum_{i=1}^n \exp\left(\frac{-c_{ij} + \lambda_i}{\gamma}\right) \right)}_{H_W^*(\lambda)} \right\}, \quad (*) \\
 L \in S_n(1) &= \left\{ L \geq 0 : \sum_{k=1}^n L_k = 1 \right\} \quad (W \in S_n(1), \gamma > 0).
 \end{aligned}$$

Due to Demyanov–Danskin theorem for  $L \in S_n(1)$  function  $H_W(L)$  is smooth with

$$\nabla H_W(L) = \lambda^*,$$

where  $\lambda^*$  is unique solution of (\*), that satisfy  $\langle \lambda^*, e \rangle = 0$ . Moreover

$$H_W^*(\lambda) = \max_{L \in S_n(1)} \{ \langle \lambda, L \rangle - H_W(L) \} = \gamma \sum_{j=1}^n W_j \ln \left( \frac{1}{W_j} \sum_{i=1}^n \exp \left( \frac{-c_{ij} + \lambda_i}{\gamma} \right) \right).$$

Wasserstein barycenter calculation problem has the following form:

$$\sum_{k=1}^m H_{W_k}(L) \rightarrow \min_{L \in S_n(1)}. \quad (**)$$

Unfortunately,  $H_{W_k}(L)$  and their gradients can't be calculate explicitly.



Let's reformulate (\*\*) in a dual (explicit) manner

$$-\sum_{k=1}^m H_{W_k}(L_k) \rightarrow \max_{\substack{L_1=L_m | \lambda^1 \\ \dots\dots\dots \\ L_{m-1}=L_m | \lambda^{m-1} \\ L_1, \dots, L_m \in S_n(1)}} ,$$

$$\sum_{k=1}^{m-1} \max_{L_k \in S_n(1)} \left\{ \langle \lambda^k, L_k \rangle - H_{W_k}(L_k) \right\} + \max_{L_m \in S_n(1)} \left\{ \left\langle -\sum_{k=1}^{m-1} \lambda^k, L_m \right\rangle - H_{W_m}(L_m) \right\} \rightarrow \min_{\lambda^1, \dots, \lambda^{m-1} \in \mathbb{R}^n} ,$$

$$\sum_{k=1}^{m-1} H_{W_k}^*(\lambda^k) + H_{W_m}^*\left(-\sum_{k=1}^{m-1} \lambda^k\right) \rightarrow \min_{\lambda^1, \dots, \lambda^{m-1} \in \mathbb{R}^n} , \quad (***)$$

$$L_* = \nabla H_{W_k}^*(\lambda_*^k) \text{ for all } k = 1, \dots, m-1,$$

where  $L_*$  and  $\{\lambda_*^k\}_{k=1}^{m-1}$  – unique solutions of problems (\*\*), (\*\*\*)

## When it is worth to solve dual problem instead of primal one?

Assume that  $Ax = b$  is compatible. Let (conditional optimization problem)

$$\frac{1}{2} \|x\|_2^2 \rightarrow \min_{Ax=b}.$$

We can build dual problem (strong duality)

$$\begin{aligned} \min_{Ax=b} \frac{1}{2} \|x\|_2^2 &= \min_x \max_{\lambda} \left\{ \frac{1}{2} \|x\|_2^2 + \langle b - Ax, \lambda \rangle \right\} = \\ &= \max_{\lambda} \min_x \left\{ \frac{1}{2} \|x\|_2^2 + \langle b - Ax, \lambda \rangle \right\} = \max_{\lambda} \left\{ \langle b, \lambda \rangle - \frac{1}{2} \|A^T \lambda\|_2^2 \right\}. \end{aligned}$$

Since  $Ax = b$  is compatible, then due to Fredholm's theorem there is no such  $\lambda$ , that  $A^T \lambda = 0$  and  $\langle b, \lambda \rangle > 0$ , hence dual problem has finite solution.

Indeed, if there exists such  $x$  that  $Ax = b$  then for all  $\lambda$ :  $\langle Ax, \lambda \rangle = \langle b, \lambda \rangle$ . Hence,  $\langle x, A^T \lambda \rangle = \langle b, \lambda \rangle$ . Assume that there exists such a  $\lambda$ , that  $A^T \lambda = 0$  and  $\langle b, \lambda \rangle > 0$ . If it is so we observe a contradiction:

$$0 = \langle x, A^T \lambda \rangle = \langle b, \lambda \rangle > 0.$$

So instead of conditional primal optimization problem one can solve

$$\langle b, \lambda \rangle - \frac{1}{2} \|A^T \lambda\|_2^2 \rightarrow \max_{\lambda} \text{ (dual problem – unconditional!)}$$

and reestablish solution of the primal problem from  $x(\lambda) = A^T \lambda$ .

If matrix  $A$  is not a full rank matrix then dual problem have many solutions (affine manifold). But all of them lead to the same (unique) primal solution  $x(\lambda) \equiv x_*$ .

**When it is worth to solve dual problem instead of primal one?**

$$\tilde{F}(x) = \langle c, x \rangle + \|x\|_a^2 + \gamma \sum_{k=1}^n x_k \ln x_k \rightarrow \min_{x \in S_n(1)}, // \text{ see Lecture 5}$$

or equivalently

$$\langle c, x \rangle + t + \gamma \sum_{k=1}^n x_k \ln x_k \rightarrow \min_{\substack{x \in S_n(1), \|x\|_a^a \leq t^{a/2}, \\ 0 \leq t \leq n^{2/a}, 0 \leq x_k \leq 1, k=1, \dots, n}} .$$

Using Sion–Kakutani theorem we can build dual problem (strong duality)

$$\tilde{G}(\lambda) = \min_{\substack{0 \leq t \leq n^{2/a}, \\ 0 \leq x_k \leq 1, k=1, \dots, n}} \left\{ \sum_{k=1}^n c_k x_k + t + \lambda_1 \cdot \left( \sum_{k=1}^n x_k - 1 \right) + \lambda_2 \cdot \left( \sum_{k=1}^n x_k^a - t^{a/2} \right) + \gamma \sum_{k=1}^n x_k \ln x_k \right\} \rightarrow \max_{\lambda_1 \in \mathbb{R}, \lambda_2 \geq 0} .$$

We can connect primal variables  $(t, x)$  with dual  $\lambda = (\lambda_1, \lambda_2)$  ( $1 < a < 2$ ):

$$t(\lambda) = \min \left\{ \left( \frac{\lambda_2 a}{2} \right)^{\frac{2}{2-a}}, n^{\frac{2}{a}} \right\},$$

and to find  $x(\lambda)$  we have to solve  $n$  one-dimensional convex optimization sub-problems on finite line segments. So we can obtain  $x(\lambda)$  with precision  $\varepsilon$  after  $O(n \ln(n/\varepsilon))$  arithmetic operation using half partition line segment method. Using Demyanov–Danskin’s formula one can obtain

$$\frac{\partial \tilde{G}}{\partial \lambda_1} = \sum_{k=1}^n x_k(\lambda) - 1, \quad \frac{\partial \tilde{G}}{\partial \lambda_2} = \sum_{k=1}^n x_k(\lambda)^a - t(\lambda)^{a/2}.$$

We’d like to solve dual problem with ellipsoid method. Since that we have to bound somehow the dual variables. We use Slater’s approach

$$\begin{aligned} -\|c\|_\infty - \gamma \ln n \leq \tilde{F}_* = \tilde{G}_* \leq \sum_{k=1}^n c_k \bar{x}_k + \bar{t} + \lambda_1 \cdot \left( \sum_{k=1}^n \bar{x}_k - 1 \right) + \\ + \lambda_2 \cdot \left( \sum_{k=1}^n \bar{x}_k^a - \bar{t}^{a/2} \right) + \gamma \sum_{k=1}^n \bar{x}_k \ln \bar{x}_k. \end{aligned}$$

If  $\lambda_1 \geq 0$ , then put  $\bar{t} = 1$ ,  $\bar{x}_k = 1/(2n)$ ,  $k = 1, \dots, n$ . Hence

$$\frac{1}{2}\lambda_1 + \frac{1}{2}\lambda_2 \leq 2\|c\|_\infty + 2\gamma \ln(n) + 1,$$

If  $\lambda_1 < 0$ , then put  $\bar{t} = 8$ ,  $\bar{x}_k = 2/n$ ,  $k = 1, \dots, n$ . Hence

$$|\lambda_1| + \frac{1}{2}\lambda_2 \leq 3\|c\|_\infty + 2\gamma \ln(2n) + 8.$$

Anyway, we have

$$\|\lambda_*\|_1 \leq 6\|c\|_\infty + 4\gamma \ln(2n) + 16 \stackrel{def}{=} C.$$

So we can restrict ourselves by solving bounded dual problem

$$-\tilde{G}(\lambda) \rightarrow \min_{\lambda_1 \in \mathbb{R}, \lambda_2 \geq 0, \|\lambda\|_1 \leq C}.$$

Using ellipsoid method (Lecture 3) we can find  $\varepsilon$ -solution of this problem after  $O(r^2 \ln(C/\varepsilon))$  iterations, the cost of one iteration  $O(r^2 + n \ln(n/\varepsilon))$ ,  $r = 2$ .

To be continued...