

Big Data - Lecture 1

Optimization reminders

S. Gadat

Toulouse, Octobre 2014

Big Data - Lecture 1

Optimization reminders

S. Gadat

Toulouse, Octobre 2014

Schedule

- 1 Introduction
 - Major issues
 - Examples
 - Mathematics
- 2 Standard Convex optimisation procedures
 - Convex functions
 - Example of convex functions
 - Gradient descent method
- 3 Constrained Convex optimisation
 - Definition
 - Equality constraint
 - Inequality constraint
 - Lagrangian in general settings
 - KKT Conditions
- 4 Conclusion

I Introduction - Major issues in data science

- **Data science:** Extract from data some knowledge for industrial or academic exploitation.
- Involves:
 - ① **Signal Processing** (how to record the data and represent it?)
 - ② **Modelization** (What is the problem, what kind of mathematical model and answer?)
 - ③ **Statistics** (reliability of estimation procedures?)
 - ④ **Machine Learning** (what kind of efficient optimization algorithm?)
 - ⑤ **Implementation** (software needs)
 - ⑥ **Visualization** (how can I represent the resulting knowledge?)
- In its whole, this sequence of questions are at the core of Artificial Intelligence and may also be referred to as Computer Science problems.
- In this lecture, we will address some issues raised in **red** items. Each time, practical examples will be provided
- Most of our motivation comes from the *Big Data* world, encountered in image processing, finance, genetics and many other fields where knowledge extraction is needed, when facing to many observations described by many variables.
- n : number of observations - p : number of variables per observations

$$p \gg n \gg O(1).$$

I Introduction - Spam classification - Signal Processing datasets

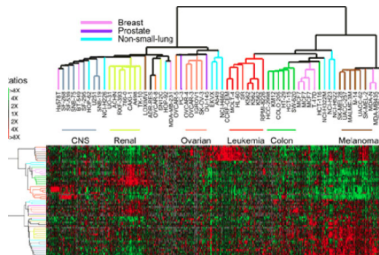
From a set of labelled messages (spam or not), build a classification for automatic spam rejection.

Variable	Mot ou Carac.	Modalités P/A	Variable	Mot ou Carac.	Modalités
make	make	make / Nmk	X650	650	650 / N65
address	address	addr / Nal	lab	lab	lab / Nlb
all	all	all / Nal	labs	labs	labs / Nls
X3d	3d	3d / N3d	telnet	telnet	teln / Ntl
our	our	our / Nou	X857	857	857 / N87
over	over	over / Nov	data	data	data / Nda
remove	remove	remo / Nrm	X415	415	415 / N41
internet	internet	inte / Nin	X85	85	85 / N85
order	order	orde / Nor	technology	technology	tech / Ntc
mail	mail	mail / Nma	X1999	1999	1999 / N19
receive	receive	rece / Nrc	parts	parts	part / Npr
will	will	will / Nwi	pm	pm	pm / Npm
people	people	peop / Npp	direct	direct	dire / Ndr
report	report	repo / Nrp	cs	cs	cs / Ncs
addresses	addresses	adds / Nas	meeting	meeting	meet / Nmt
free	free	free / Nfr	original	original	orig / or
business	business	busi / Nbs	project	project	proj / Npj
email	email	emai / Nem	re	re	re / Nre
you	you	you / Nyo	edu	edu	edu / Ned
credit	credit	cred / Ncr	table	table	tabl / Ntb
your	your	your / Nyr	conference	conferenc	e conf / Ncf
font	font	font / Nft	CsemiCol	;	Csci / NCs
X000	000	000 / N00	Cpar	(Cpar / NCp
money	money	money / Nmn	Ccroch	[Ccro / NCc
hp	hp	hp / Nhp	Cexclam	!	Cexc / NCe
hpl	hpl	hpl / Nhpl	Cdollar	\$	Cdol / NCd
george	george	geor / Nge	Cdie	#	Cdie / NCi

- Select among the words meaningful elements?
- Automatic classification?

I Introduction - Micro-array analysis - Biological datasets

One measures micro-array datasets built from a huge amount of profile genes expression. Number of genes p (of order thousands). Number of samples n (of order hundred).



Diagnostic help: healthy or ill?

- Select among the genes meaningful elements?
- Automatic classification?

I Introduction - Fraud detection - Industrial datasets

Set of individual electrical consumption for some people in Medellin (Colombia).

- Each individual provides a monthly electrical consumption.
- The electrical firm measures the whole consumption for important hubs (they are formed by a set of clients).

Want to detect eventual fraud.

Problems:

- Missing data: completion of the table. How?
- Noise in the several measurements: how does it degrades the fraud detection?
- Can we exhibit several monthly behaviour of the clients?

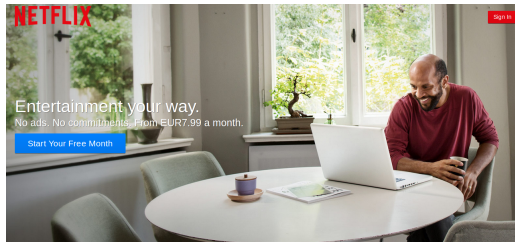
I Introduction - Data Completion & Recommendation systems - Advertisement and e-business datasets

Recommndation problems:



[Your Amazon.com](#) [Today's Deals](#) [Gift Cards](#) [Sell](#) [Help](#)

And more recently:



- What kind of database?
- Reliable recommendation for clients?
- Online strategy?

1 Introduction - Credit scoring - Actuaries datasets

Build an indicator (Q score) from a dataset for the probability of interest in a financial product (Visa premier credit card).

TABLE 1 - Liste des variables et de leur libellé

Identif.	Libellé
matric	Matricule (identifiant client)
depts	Département de résidence
pro	Point de vente
sexseq	Sexe (qualitatif)
ager	Âge en années
famiq	Situation familiale (Poux : marié, Pcel : célibataire, Filv : divorcé, Foll : sans lien, Fsep : séparé de corps, Poux : veuf)
riat	Ancienneté de relation en mois
pcsqp	Catégorie socio-professionnelle (code nani)
qual5	Code "qualité" client évalué par la banque
GxvGxvS	plusieurs variables caractérisant les intérêts bancaires
impeds	Nombre d'impôts en cours
rejets	Montant total des rejets en francs
opqab	Nombre d'opérations par gauchet dans le mois
meqr	Moyenne des mouvements nets créditeurs des 3 mois en Kf
uvep	Total des avoirs épargne monétaire en francs
endet	Taux d'endettement
gaqet	Total des engagements en francs
gaqec	Total des engagements court terme en francs
gaqem	Total des engagements moyen terme en francs
lvuah	Nombre de comptes à vue
qmeqr	Moyenne des soldes moyens sur 3 mois
qmeqr	Moyenne des mouvements créditeurs en Kf
dmqr	Âge du dernier mouvement (en jours)

TABLE 2 - Liste des variables et de leur libellé — suite

Identif.	Libellé
teppn	Nombre d'opérations à M-1
factn	Montant facturé dans l'année en francs
lgagt	Engagement long terme
vienb	Nombre de produits contrats vie
vienm	Montant des produits contrats vie en francs
ueenmb	Nombre de produits épargne monétaire
ueenmb	Montant des produits d'épargne monétaire en francs
slqph	Nombre de produits d'épargne logement
slqph	Montant des produits d'épargne logement en francs
ylvrb	Nombre de comptes sur livret
ylvrb	Montant des comptes sur livret en francs
rlvrb	Nombre de produits d'épargne long terme
rlvrb	Montant des produits d'épargne long terme en francs
rlvrb	Nombre de produits épargne à terme
rlvrb	Montant des produits épargne à terme
abbees	Nombre de produits bons et certificats
abbees	Montant des produits bons et certificats en francs
meqr	Nombre de paiements par carte bancaire à M-1
meqr	Nombre total de cartes
meqr	Nombre de cartes point argent
segr2s	Ségmentation version 2
itarc	Total des avoirs sur tous les comptes
haver	Total des avoirs épargne financière en francs
judb1s	Nombre de jours à débit à M
judb2s	Nombre de jours à débit à M-1
judb3s	Nombre de jours à débit à M-2
carvp	Possession de la carte VISA Premier

- 1 Define a model, a question?
- 2 Use a supervised classification algorithm to rank the best clients.
- 3 Use logistic regression to provide a score.

I Introduction - What about maths?

Various mathematical fields we will talk about:

- Analysis: Convex optimization, Approximation theory
- Statistics: Penalized procedures and their reliability
- Probabilistic methods: concentration inequalities, stochastic processes, stochastic approximations

Famous keywords:

- Lasso
- Boosting
- Convex relaxation
- Supervised classification
- Support Vector Machine
- Aggregation rules
- Gradient descent
- Stochastic Gradient descent
- Sequential prediction
- Bandit games, minimax policies
- Matrix completion

Schedule

- 1 Introduction
 - Major issues
 - Examples
 - Mathematics
- 2 Standard Convex optimisation procedures
 - Convex functions
 - Example of convex functions
 - Gradient descent method
- 3 Constrained Convex optimisation
 - Definition
 - Equality constraint
 - Inequality constraint
 - Lagrangian in general settings
 - KKT Conditions
- 4 Conclusion

Convex functions

We recall some background material that is necessary for a clear understanding of how some machine learning procedures work. We will cover some basic relationships between convexity, positive semidefiniteness, local and global minimizers.

Definition (Convex sets, convex functions)

A set D is convex if for any $(x_1, x_2) \in D^2$ and all $\alpha \in [0, 1]$, $x = \alpha x_1 + (1 - \alpha)x_2 \in D$. A function f is convex if

- its domain D is convex
- $f(x) = f(\alpha x_1 + (1 - \alpha)x_2) \leq \alpha f(x_1) + (1 - \alpha)f(x_2)$.

Definition (Positive Semi Definite matrix (PSD))

A $p \times p$ matrix H is (PSD) if for all $p \times 1$ vectors z , we have $z^t H z \geq 0$.

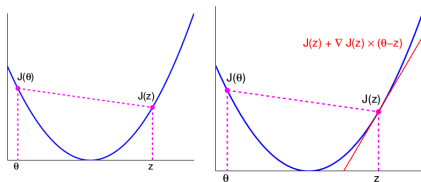
There exists a strong link between SDP matrix and convex functions, given by the following proposition.

Proposition

A smooth $C^2(D)$ function f is convex if and only if $D^2 f$ is SDP at any point of D .

The proof follows easily from a second order Taylor expansion.

Example of convex functions



- **Exponential function:** $\theta \in \mathbb{R} \mapsto \exp(a\theta)$ on \mathbb{R} whatever a is.
- **Affine function:** $\theta \in \mathbb{R}^d \mapsto a^t \theta + b$
- **Entropy function:** $\theta \in \mathbb{R}_+ \mapsto -\theta \log(\theta)$
- **p -norm:** $\theta \in \mathbb{R}^d \mapsto \|\theta\|_p := \sqrt[p]{\sum_{i=1}^d \|\theta_i\|^p}$.
- **Quadratic form:** $\theta \in \mathbb{R}^d \mapsto \theta^t P \theta + 2q^t \theta + r$ where P is symmetric and positive.

Why convex functions are useful?

From **external motivations**:

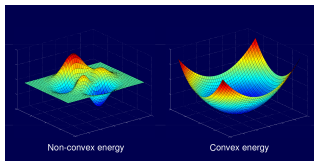
- Many problems in machine learning come from the minimization of a convex criterion and provide meaningful results for the statistical initial task.
- Many optimization problems admit a convex reformulation (SVM classification or regression, LASSO regression, ridge regression, permutation recovery, ...).

From a **numerical point of view**:

- Local minimizer = global minimizer. It is a powerful point since in general, descent methods involve $\nabla f(x)$ (or something related to), which is a local information on f .
- x is a local (global) minimizer of f if and only if $0 \in \partial f(x)$.
- Many fast algorithms for the optimization of convex function exist, and sometimes are independent on the dimension d of the original space.

Why convex functions are powerful?

Two kind of optimization problems:



- On the left: non convex optimization problem, use of Travelling Salesman type method. Greedy exploration step (simulated annealing, genetic algorithms).
- On the right: convex optimization problem, use local descent methods with gradients or subgradients.

Definition (Subgradient (nonsmooth functions?))

For any function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, and any x in \mathbb{R}^d , the subgradient $\partial f(x)$ is the set of all vectors g such that

$$f(x) - f(y) \leq \langle g, x - y \rangle.$$

This set of subgradients may be empty. Fortunately, it is not the case for convex functions.

Proposition

$f : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex if and only if $\partial f(x) \neq \emptyset$ for any x of \mathbb{R}^d .

Convexity and gradient: Constrained case

In either constrained or unconstrained problems, descent methods are powerful with convex functions. In particular, consider constrained problems in $\mathcal{X} \subset \mathbb{R}^d$. The most famous local descent method relies on

$$y_{t+1} = x_t - \eta g_t \quad \text{where} \quad g_t \in \partial f(x_t),$$

and

$$x_{t+1} = \Pi_{\mathcal{X}}(y_{t+1}),$$

where $\eta > 0$ is a fixed step-size parameters.

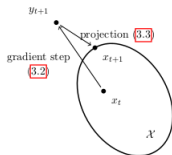


Fig. 3.2 Illustration of the Projected Subgradient Descent method.

Théorème (Convergence of the projected gradient descent method, fixed step-size)

If f is convex over \mathcal{X} with $\mathcal{X} \subset B(0, R)$ and $\|\partial f\|_{\infty} \leq L$, the choice $\eta = \frac{R}{L\sqrt{t}}$ leads to

$$f\left(\frac{1}{t} \sum_{s=1}^t x_s\right) - \min f \leq \frac{RL}{\sqrt{t}}$$

Convexity and gradient: Smooth unconstrained case

Results can be seriously improved with smooth functions with bounded second derivatives.

Definition

f is β smooth if ∇f is β Lipschitz:

$$\|\nabla f(x) - \nabla f(y)\| \leq \beta \|x - y\|.$$

Standard gradient descent over \mathbb{R}^d becomes

$$x_{t+1} = x_t - \eta \nabla f(x_t),$$

Théorème (Convergence of the gradient descent method, β smooth function)

If f is a convex and β -smooth function, then $\eta = \frac{1}{\beta}$ leads to

$$f\left(\frac{1}{t} \sum_{s=1}^t x_s\right) - \min f \leq \frac{2\beta \|x_1 - x_0\|^2}{t-1}$$

Remark

- Note that the two past results *do not depend on the dimension of the state space d* .
- The last result can be extended to the constrained situation.

Schedule

- 1 Introduction
 - Major issues
 - Examples
 - Mathematics
- 2 Standard Convex optimisation procedures
 - Convex functions
 - Example of convex functions
 - Gradient descent method
- 3 **Constrained Convex optimisation**
 - **Definition**
 - **Equality constraint**
 - **Inequality constraint**
 - **Lagrangian in general settings**
 - **KKT Conditions**
- 4 Conclusion

Constrained optimisation: Definition

Elements of the problem:

- θ unknown vector of \mathbb{R}^d to be recovered
- $J : \mathbb{R}^d \mapsto \mathbb{R}$ function to be minimized
- f_i and g_i differentiable functions defining a set of constraints.

Definition of the problem:

- $\min_{\theta \in \mathbb{R}^d} J(\theta)$ **such that:**
- $f_i(\theta) = 0, \forall i = 1, \dots, n$
- $g_i(\theta) \leq 0, \forall i = 1, \dots, m$

Set of admissible vectors:

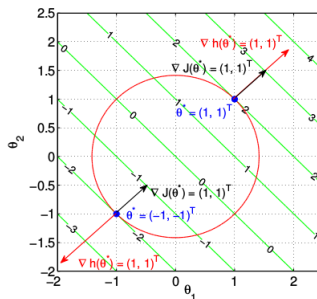
$$\Omega := \left\{ \theta \in \mathbb{R}^d \mid f_i(\theta) = 0, \forall i \text{ and } g_j(\theta) \leq 0, \forall j \right\}$$

Constrained optimisation: Example

Typical situation:

Example

$$\begin{aligned} \min_{\theta \in \mathbb{R}^2} \quad & \theta_1 + \theta_2 \\ \text{s.t.} \quad & \theta_1^2 + \theta_2^2 - 2 = 0 \end{aligned}$$



Ω : circle of radius $\sqrt{2}$

Optimal solution: $\theta^* = (-1, -1)^T$ and $J(\theta^*) = -2$.

Important restriction: we will restrict our study to **convex functions** J .

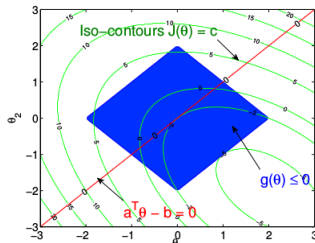
Constrained convex optimisation

A constrained optimization problem is "convex" if:

- J is a convex function
- f_i are linear or affine functions
- g_i are convex functions

Exemple

$$\begin{aligned} \min_{\theta \in \mathbb{R}^2} \quad & \theta_1^2 + \theta_2^2 + \theta_1\theta_2 \\ & -2\theta_1 + 2\theta_2 - 2 \\ \text{s.c.} \quad & \theta_1 - \theta_2 = 0 \\ & \|\theta\|_1 - 2 \leq 0 \end{aligned}$$



Case of a unique equality constraint

$$\min_{\theta} J(\theta) \quad \text{such that} \quad a^t \theta - b = 0$$

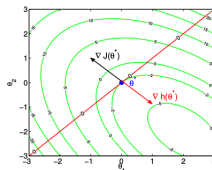
- Descent direction h : $\nabla J(\theta)^t h < 0$.
- Admissible direction h : $a^t(\theta + h) - b = 0 \iff a^t h = 0$.

Optimality θ^* is optimal if there is no admissible descent direction starting from θ^* . The only possible case is when $\nabla J(\theta^*)$ and a are linearly dependent:

$$\exists \lambda \in \mathbb{R} \quad \nabla J(\theta^*) + \lambda a = 0.$$

Example

$$\begin{aligned} \min_{\theta \in \mathbb{R}^2} \quad & \theta_1^2 + \theta_2^2 + \theta_1 \theta_2 \\ & -2\theta_1 + 2\theta_2 - 2 \\ \text{s.t.} \quad & \theta_1 - \theta_2 = 0 \end{aligned}$$



In this situation:

$$\nabla J(\theta) = \begin{pmatrix} 2\theta_1 + \theta_2 - 2 \\ \theta_1 + 2\theta_2 + 2 \end{pmatrix} \quad \text{and} \quad a = \begin{pmatrix} 1 \\ -1 \end{pmatrix}$$

Hence, we are looking for θ such that $\nabla J(\theta) \propto a$. Computations lead to $\theta_1 = -\theta_2$. Optimal value reached for $\theta_1 = 1/2$ (and $J(\theta^*) = -15/4$).

Lagrangian function

$$\min_{\theta} J(\theta) \quad \text{such that} \quad f(\theta) := a^t \theta - b = 0$$

We have seen the important role of the scalar value λ above.

Definition (Lagrangian function)

$$L(\lambda, \theta) = J(\theta) + \lambda f(\theta)$$

λ is the Lagrange multiplier. The optimal choice of (θ^*, λ^*) corresponds to

$$\nabla_{\theta} L(\lambda^*, \theta^*) = 0 \quad \text{and} \quad \nabla_{\lambda} L(\lambda^*, \theta^*) = 0.$$

Argument: θ^* is optimal if there is no admissible descent directions h . Hence, ∇J and ∇f are linearly dependent. As a consequence, there exists λ such that

$$\nabla_{\theta} L(\lambda^*, \theta^*) = \nabla J(\theta) + \lambda \nabla f(\theta) = 0 \quad (\text{Dual equation})$$

Since θ must be admissible, we have

$$\nabla_{\lambda} L(\lambda^*, \theta^*) = f(\theta^*) = 0 \quad (\text{Primal equation})$$

Case of a unique inequality constraint

$$\min_{\theta} J(\theta) \quad \text{such that} \quad g(\theta) \leq 0$$

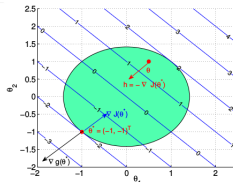
- Descent direction h : $\nabla J(\theta)^t h < 0$.
- Admissible direction h : $\nabla g(\theta)^t h \leq 0$ guarantees that $g(\theta + \alpha h)$ is decreasing with α .

Optimality θ^* is optimal if there is no admissible descent direction starting from θ^* . The only possible case is when $\nabla J(\theta^*)$ and $\nabla g(\theta^*)$ are linearly dependent and opposite:

$$\exists \lambda \in \mathbb{R} \quad \nabla J(\theta^*) = -\mu \nabla g(\theta^*) \quad \text{with} \quad \mu \geq 0.$$

Exemple

$$\begin{aligned} \min_{\theta \in \mathbb{R}^2} \quad & \theta_1 + \theta_2 \\ \text{s.t.} \quad & g(\theta) = \theta_1^2 + \theta_2^2 - 2 \leq 0 \end{aligned}$$



We can check that $\theta^* = (-1, -1)$.

Lagrangian in general settings

We consider the minimization problem:

- $\min_{\theta} J(\theta)$ such that
- $g_j(\theta) \leq 0, \forall j = 1, \dots, m$
- $f_i(\theta) = 0, \forall i = 1, \dots, n$

Definition (Lagrangian function)

We associate to this problem the Lagrange multipliers $(\lambda, \mu) = (\lambda_1, \dots, \lambda_n, \mu_1, \dots, \mu_m)$.

$$L(\theta, \lambda, \mu) = J(\theta) + \sum_{i=1}^n \lambda_i f_i(\theta) + \sum_{j=1}^m \mu_j g_j(\theta)$$

- θ primal variables
- (λ, μ) dual variables

KKT Conditions

Definition (KKT Conditions)

If J and f, g are smooth, we define the Karush-Kuhn-Tucker (KKT) conditions as

- Stationarity: $\nabla_{\theta} L(\lambda, \mu, \theta) = 0$.
- Primal Admissibility: $f(\theta) = 0$ and $g(\theta) \leq 0$.
- Dual admissibility $\mu_j \geq 0, \forall j = 1, \dots, m$.

Theorem

A convex minimization problem of J under convex constraints f and g has a solution θ^ if and only if there exists λ^* and μ^* such that KKT conditions hold.*

Example:

$$J(\theta) = \frac{1}{2} \|\theta\|_2^2 \quad \text{s.t.} \quad \theta_1 - 2\theta_2 + 2 \leq 0$$

We get $L(\theta, \mu) = \frac{\|\theta\|_2^2}{2} + \mu(\theta_1 + 2\theta_2 + 2)$ with $\mu \geq 0$. Stationarity: $(\theta_1 + \mu, \theta_2 - 2\mu) = 0$.

$$\theta_2 = -2\theta_1 \quad \text{with} \quad \theta_2 \leq 0.$$

We deduce that $\theta^* = (-2/5, 4/5)$.

Dual problems (1)

We introduce the *dual* function:

$$\mathcal{L}(\lambda, \mu) = \min_{\theta} L(\theta, \lambda, \mu).$$

We have the following important result

Theorem

Denote the optimal value of the constrained problem $p^* = \min \{J(\theta) | f(\theta) = 0, g(\theta) \leq 0\}$, then

$$\mathcal{L}(\lambda, \mu) \leq p^*.$$

Remark:

- The dual function \mathcal{L} is lower than p^* , for any $(\lambda, \mu) \in \mathbb{R}^n \times \mathbb{R}_+^m$
- We aim to make this lower bound as close as possible to p^* : idea to maximize w.r.t. λ, μ the function \mathcal{L} .

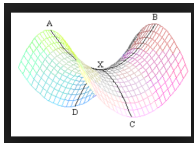
Definition (Dual problem)

$$\max_{\lambda \in \mathbb{R}^n, \mu \in \mathbb{R}_+^m} \mathcal{L}(\lambda, \mu).$$

$L(\theta, \lambda, \mu)$ affine function on λ, μ and thus **convex**. Hence, \mathcal{L} is **convex** and almost unconstrained.

Dual problems (2)

- Dual problems are easier than primal ones (because of almost constraints omissions).
- Dual problems are equivalent to primal ones: maximization of the dual \Leftrightarrow minimization of the primal (not shown in this lecture).
- Dual solutions permit to recover primal ones with KKT conditions (Lagrange multipliers).



Example:

- Lagrangian: $L(\theta, \mu) = \frac{\theta_1^2 + \theta_2^2}{2} + \mu(\theta_1 - 2\theta_2 + 2)$.
- Dual function $\mathcal{L}(\mu) = \min_{\theta} L(\theta, \mu) = -\frac{5}{2}\mu^2 + 2\mu$.
- Dual solution: $\max \mathcal{L}(\mu)$ such that $\mu \geq 0$: $\mu = 2/5$.
- Primal solution: $\text{KKT} \Rightarrow \theta = (-\mu, 2\mu) = (-2/5, 4/5)$.

Take home message

- Big Data problems arise in a large variety of fields. They are complicated for a computational reason (and also for a statistical one, see later).
- Many Big Data problems will be traduced in an optimization of a **convex problem**.
- Efficient algorithms are available to optimize them:
independently on the dimension of the underlying space.
- Primal - Dual formulations are important to overcome some constraints on the optimization.
- Numerical convex solvers are widely and freely distributed.