

# Computation & optimization for Lasso - part 2

Luyang Han & Janosch Ott

ETH Zürich

22 October 2018

2018-10-14

Computation & optimization

Computation & optimization for Lasso - part 2

Luyang Han & Janosch Ott

ETH Zürich

22 October 2018

# Overview

1. Coordinate Descent
2. A Simulation Study
3. Least Angle Regression
4. Digression: Duality
5. ADMM
6. Screening Rules
7. Minor-Max Algorithms
8. Alternating Minimizations

2018-10-14

## Computation & optimization

### └ Overview

1. Coordinate Descent
2. A Simulation Study
3. Least Angle Regression
4. Digression: Duality
5. ADMM
6. Screening Rules
7. Minor-Max Algorithms
8. Alternating Minimizations

# Digression: Duality in optimization

2018-10-14

Computation & optimization

└ Digression: Duality

└ Digression: Duality in optimization

Digression: Duality in optimization

In various section, I came across terms like "dual" and "dual problem"

Primal	
Optimize	$\min f(x)$
Constraints	$g_i(x) \leq 0, h_j(x) = 0, x \in X$
Function	$L(x, \lambda, \mu) := f(x) + \sum_i \lambda_i g_i(x) + \sum_j \mu_j h_j(x)$
Dual	
Function	$q(\lambda, \mu) = \inf_{x \in X} L(x, \lambda, \mu)$
Constraints	$\lambda \geq 0$
Optimize	$\max q(\lambda, \mu)$

Why though? - **Dual problem is always convex!**

$x \in X$  for e.g. solutions in a cone or integer solutions

Terms: Primal problem, Lagrange function with dual variables/Lagrange-multipliers, dual function, dual problem

Dual problem is always convex! - I don't know much about optimization yet, but they really like convexity.

"(Convexity confers two advantages. The first is that, in a constrained problem, a convex feasible region makes it easier to ensure that you do not generate infeasible solutions while searching for an optimum.)

The second advantage is that all local optima are global optima. That allows local search algorithms to guarantee optimal solutions. And local search is often faster." [Rubin, 2016])

## Primal

Optimize  $\min f(x)$

Constraints  $g_i(x) \leq 0, h_j(x) = 0, x \in X$

Function  $L(x, \lambda, \mu) := f(x) + \sum_i \lambda_i g_i(x) + \sum_j \mu_j h_j(x)$

## Dual

Function  $q(\lambda, \mu) = \inf_{x \in X} L(x, \lambda, \mu)$

Constraints  $\lambda \geq 0$

Optimize  $\max q(\lambda, \mu)$

Why though? - **Dual problem is always convex!**

# Alternating Direction Method of Multipliers (ADMM)

## Problem

$$\underset{\beta \in \mathbb{R}^m, \theta \in \mathbb{R}^n}{\text{minimize}} \quad f(\beta) + g(\theta) \quad \text{subject to} \quad \mathbf{A}\beta + \mathbf{B}\theta - c = 0$$

## Lagrangian

$$f(\beta) + g(\theta) + \langle \mu, \mathbf{A}\beta + \mathbf{B}\theta - c \rangle$$

## Augmented Lagrangian

$$L_\rho(\beta, \theta, \mu) := f(\beta) + g(\theta) + \langle \mu, \mathbf{A}\beta + \mathbf{B}\theta - c \rangle + \frac{\rho}{2} \|\mathbf{A}\beta + \mathbf{B}\theta - c\|_2^2$$

2018-10-14

## Computation & optimization

└ ADMM

└ Alternating Direction Method of Multipliers  
(ADMM)Alternating Direction Method of  
Multipliers (ADMM)

Problem

$$\underset{\beta \in \mathbb{R}^m, \theta \in \mathbb{R}^n}{\text{minimize}} \quad f(\beta) + g(\theta) \quad \text{subject to} \quad \mathbf{A}\beta + \mathbf{B}\theta - c = 0$$

Lagrangian

$$f(\beta) + g(\theta) + \langle \mu, \mathbf{A}\beta + \mathbf{B}\theta - c \rangle$$

Augmented Lagrangian

$$L_\rho(\beta, \theta, \mu) := f(\beta) + g(\theta) + \langle \mu, \mathbf{A}\beta + \mathbf{B}\theta - c \rangle + \frac{\rho}{2} \|\mathbf{A}\beta + \mathbf{B}\theta - c\|_2^2$$

# Alternating Direction Method of Multipliers (ADMM)

Problem - decomposable !

$$\underset{\beta \in \mathbb{R}^m, \theta \in \mathbb{R}^n}{\text{minimize}} \quad f(\beta) + g(\theta) \quad \text{subject to} \quad \mathbf{A}\beta + \mathbf{B}\theta - c = 0$$

Lagrangian

$$f(\beta) + g(\theta) + \langle \mu, \mathbf{A}\beta + \mathbf{B}\theta - c \rangle$$

Augmented Lagrangian

$$L_\rho(\beta, \theta, \mu) := f(\beta) + g(\theta) + \langle \mu, \mathbf{A}\beta + \mathbf{B}\theta - c \rangle + \frac{\rho}{2} \|\mathbf{A}\beta + \mathbf{B}\theta - c\|_2^2$$

2018-10-14

Computation &amp; optimization

└ ADMM

└ Alternating Direction Method of Multipliers  
(ADMM)

decomposable problem and constraints!

Alternating Direction Method of  
Multipliers (ADMM)

Problem - decomposable !

$$\underset{\beta \in \mathbb{R}^m, \theta \in \mathbb{R}^n}{\text{minimize}} \quad f(\beta) + g(\theta) \quad \text{subject to} \quad \mathbf{A}\beta + \mathbf{B}\theta - c = 0$$

Lagrangian

$$f(\beta) + g(\theta) + \langle \mu, \mathbf{A}\beta + \mathbf{B}\theta - c \rangle$$

Augmented Lagrangian

$$L_\rho(\beta, \theta, \mu) := f(\beta) + g(\theta) + \langle \mu, \mathbf{A}\beta + \mathbf{B}\theta - c \rangle + \frac{\rho}{2} \|\mathbf{A}\beta + \mathbf{B}\theta - c\|_2^2$$

# Alternating Direction Method of Multipliers (ADMM)

Problem - decomposable !

$$\underset{\beta \in \mathbb{R}^m, \theta \in \mathbb{R}^n}{\text{minimize}} f(\beta) + g(\theta) \quad \text{subject to } \mathbf{A}\beta + \mathbf{B}\theta - c = 0$$

Lagrangian - decomposable !

$$f(\beta) + g(\theta) + \langle \mu, \mathbf{A}\beta + \mathbf{B}\theta - c \rangle$$

Augmented Lagrangian

$$L_\rho(\beta, \theta, \mu) := f(\beta) + g(\theta) + \langle \mu, \mathbf{A}\beta + \mathbf{B}\theta - c \rangle + \frac{\rho}{2} \|\mathbf{A}\beta + \mathbf{B}\theta - c\|_2^2$$

2018-10-14

Computation &amp; optimization

└ ADMM

└ Alternating Direction Method of Multipliers (ADMM)

Alternating Direction Method of Multipliers (ADMM)

Problem - decomposable !

$$\underset{\beta \in \mathbb{R}^m, \theta \in \mathbb{R}^n}{\text{minimize}} f(\beta) + g(\theta) \quad \text{subject to } \mathbf{A}\beta + \mathbf{B}\theta - c = 0$$

Lagrangian - decomposable !

$$f(\beta) + g(\theta) + \langle \mu, \mathbf{A}\beta + \mathbf{B}\theta - c \rangle$$

Augmented Lagrangian

$$L_\rho(\beta, \theta, \mu) := f(\beta) + g(\theta) + \langle \mu, \mathbf{A}\beta + \mathbf{B}\theta - c \rangle + \frac{\rho}{2} \|\mathbf{A}\beta + \mathbf{B}\theta - c\|_2^2$$

Lagrangian problem can still be decomposed into  $\beta$  and  $\mu$  terms  
this has nice algorithm where we can execute some stuff in parallel, because  
we can decompose the Lagrangian

# Alternating Direction Method of Multipliers (ADMM)

Problem - decomposable !

$$\underset{\beta \in \mathbb{R}^m, \theta \in \mathbb{R}^n}{\text{minimize}} f(\beta) + g(\theta) \quad \text{subject to } \mathbf{A}\beta + \mathbf{B}\theta - c = 0$$

Lagrangian - decomposable !

$$f(\beta) + g(\theta) + \langle \mu, \mathbf{A}\beta + \mathbf{B}\theta - c \rangle$$

Augmented Lagrangian - NOT decomposable !

$$L_\rho(\beta, \theta, \mu) := f(\beta) + g(\theta) + \langle \mu, \mathbf{A}\beta + \mathbf{B}\theta - c \rangle + \frac{\rho}{2} \|\mathbf{A}\beta + \mathbf{B}\theta - c\|_2^2$$

2018-10-14

Computation & optimization

└ ADMM

└ Alternating Direction Method of Multipliers (ADMM)

Alternating Direction Method of Multipliers (ADMM)

Problem - decomposable !

$$\underset{\beta \in \mathbb{R}^m, \theta \in \mathbb{R}^n}{\text{minimize}} f(\beta) + g(\theta) \quad \text{subject to } \mathbf{A}\beta + \mathbf{B}\theta - c = 0$$

Lagrangian - decomposable !

$$f(\beta) + g(\theta) + \langle \mu, \mathbf{A}\beta + \mathbf{B}\theta - c \rangle$$

Augmented Lagrangian - NOT decomposable !

$$L_\rho(\beta, \theta, \mu) := f(\beta) + g(\theta) + \langle \mu, \mathbf{A}\beta + \mathbf{B}\theta - c \rangle + \frac{\rho}{2} \|\mathbf{A}\beta + \mathbf{B}\theta - c\|_2^2$$

Augmented: scalar product with  $\rho$  gets added,  
Method of Multipliers: is a way to make the algorithm more robust

advantage: better convergence

disadvantage: no longer parallel execution of subtasks due to l2-term, no longer decomposable in beta and theta terms, as l2 norm squares every entry of the vector

alternating direction: semi-decomposable, i.e. keeping one variable fixed while updating the other

$\rho$  is step length of iterative algorithm

All notes on this slide: see the slides by [Boyd]



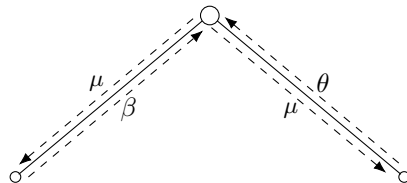
# Dual Variable Update

## Alternating Direction Method of Multipliers

$$\beta^{t+1} = \arg \min_{\beta \in \mathbb{R}^m} L_\rho(\beta, \theta^t, \mu^t)$$

$$\theta^{t+1} = \arg \min_{\theta \in \mathbb{R}^m} L_\rho(\beta^{t+1}, \theta, \mu^t)$$

$$\mu^{t+1} = \mu^t + \rho(\mathbf{A}\beta^{t+1} + \mathbf{B}\theta^{t+1} - c)$$



**Figure:** My own illustration of the dual ascent step in the ADMM algorithm utilising dual decomposition based on [Gordon and Tibshirani, 2012].

2018-10-14

## Computation & optimization

### ADMM

### Dual Variable Update

$$\begin{aligned}\beta^{t+1} &= \arg \min_{\beta \in \mathbb{R}^m} L_\rho(\beta, \theta^t, \mu^t) \\ \theta^{t+1} &= \arg \min_{\theta \in \mathbb{R}^m} L_\rho(\beta^{t+1}, \theta, \mu^t) \\ \mu^{t+1} &= \mu^t + \rho(\mathbf{A}\beta^{t+1} + \mathbf{B}\theta^{t+1} - c)\end{aligned}$$



**Figure:** My own illustration of the dual ascent step in the ADMM algorithm utilising dual decomposition based on [Gordon and Tibshirani, 2012].

Method of Multipliers: is a way to make the algorithm more robust, (if in second line  $\beta^t$  statt  $\beta^{t+1}$ )

alternating direction: semi-decomposable, i.e. keeping one variable fixed while updating the other

last step is called a dual variable update, this dual has nothing to do with two, but is connected to what is called a dual problem

dual variable update, we are working in the dual problem as "min L", thus convex problem, thus "dual decomposition" into subproblems which is possible by [Gordon and Tibshirani, 2012]

think of it as only the last line, sending  $\mu$  to the updaters for  $\beta$  and  $\theta$   
 $\rho$  in last line can be thought of as "step length"

All notes on this slide: see the slides by [Boyd]

## ADMM - Why?

- convex problems with nondifferentiable constraints
- blockwise computation
  - sample blocks
  - feature blocks

2018-10-14

Computation & optimization

└ ADMM

└ ADMM - Why?

Details for blockwise computation in Exercise 5.12.

# ADMM for the Lasso Problem

Problem in Lagrangian form

$$\underset{\beta \in \mathbb{R}^p, \theta \in \mathbb{R}^p}{\text{minimize}} \left\{ \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \|\theta\|_1 \right\} \quad \text{such that } \beta - \theta = 0$$

Augmented Lagrangian

$$L_\rho(\beta, \theta, \mu) := \left\{ \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \|\theta\|_1 \right\} + \langle \mu, \beta - \theta \rangle + \frac{\rho}{2} \|\beta - \theta\|_2^2$$

2018-10-14

Computation & optimization  
└ ADMM

└ ADMM for the Lasso

ADMM for the Lasso  
Problem

Problem in Lagrangian form

$$\underset{\beta \in \mathbb{R}^p, \theta \in \mathbb{R}^p}{\text{minimize}} \left\{ \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \|\theta\|_1 \right\} \quad \text{such that } \beta - \theta = 0$$

Augmented Lagrangian

$$L_\rho(\beta, \theta, \mu) := \left\{ \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \|\theta\|_1 \right\} + \langle \mu, \beta - \theta \rangle + \frac{\rho}{2} \|\beta - \theta\|_2^2$$

In the problem, I can decompose into beta and theta terms, i.e. show  $f(\beta)$  and  $g(\theta)$

the problem itself and the constraints,

A and B are unit matrices here

Computational cost: Initially  $\mathcal{O}(p^3)$ , which is a lot, for the SVD(singular value decomposition of  $\mathbf{X}$ ), after that comparable to coordinate descent or composite gradient from earlier

## ADMM for the Lasso

Update

Update

$$\beta^{t+1} = (\mathbf{X}^T \mathbf{X} + \rho \mathbf{I})^{-1} (\mathbf{X}^T \mathbf{y} + \rho \theta^t - \mu^t)$$

$$\theta^{t+1} = \mathcal{S}_{\lambda/\rho}(\beta^{t+1} + \mu^t/\rho)$$

$$\mu^{t+1} = \mu^t + \rho(\beta^{t+1} - \theta^{t+1})$$

where  $\mathcal{S}_{\lambda/\rho}(z) = \text{sign}(z)(|z| - \frac{\lambda}{\rho})_+$ .

2018-10-14

Computation &amp; optimization

└ ADMM

└ ADMM for the Lasso

ADMM for the Lasso  
Update

Update

$$\beta^{t+1} = (\mathbf{X}^T \mathbf{X} + \rho \mathbf{I})^{-1} (\mathbf{X}^T \mathbf{y} + \rho \theta^t - \mu^t)$$

$$\theta^{t+1} = \mathcal{S}_{\lambda/\rho}(\beta^{t+1} + \mu^t/\rho)$$

$$\mu^{t+1} = \mu^t + \rho(\beta^{t+1} - \theta^{t+1})$$

where  $\mathcal{S}_{\lambda/\rho}(z) = \text{sign}(z)(|z| - \frac{\lambda}{\rho})_+$ .

$\mathcal{S}$  is a soft-thresholding parameter

Computational cost: Initially  $\mathcal{O}(p^3)$ , which is a lot, for the SVD(singular value decomposition of  $\mathbf{X}$ ), after that comparable to coordinate descent or composite gradient from earlier

# Screening Rules

- very big data set, esp. huge number of predictors
- maybe too big to load into memory
- Screening rules eliminate predictors with minor calculation
- and very high / safe certainty (i.e. eliminated predictors would not show up in lasso model based on full data)

2018-10-14

Computation & optimization

└ Screening Rules

└ Screening Rules

Imagine a big data set, a very big data set, with such a huge design matrix, that you cannot load it into memory (RAM). Wh

- very big data set, esp. huge number of predictors
- maybe too big to load into memory
- Screening rules eliminate predictors with minor calculation
- and very high / safe certainty (i.e. eliminated predictors would not show up in lasso model based on full data)

# What is a good predictor?

includegraphics

correlation is covariance with some factors

covariance is an inner product on a vector space

high absolute correlation (=large absolute inner product)  $\Rightarrow$

high predictive power (see plots)  $\Rightarrow$   $x_j$  with largest inner  
product has predictive power, thus for that  $j$  we are most  
willing to accept some penalty from  $\lambda$

2018-10-14

Computation & optimization

└ Screening Rules

└ What is a good predictor?

What is a good predictor?

includegraphics  
correlation is covariance with some factors  
covariance is an inner product on a vector space  
high absolute correlation ( $\Rightarrow$  large absolute inner product)  $\Rightarrow$   
high predictive power (see plots)  $\Rightarrow$   $x_j$  with largest inner  
product has predictive power, thus for that  $j$  we are most  
willing to accept some penalty from  $\lambda$

# SAFE Rules

2018-10-14

Computation &amp; optimization

└ Screening Rules

└ SAFE Rules

# Dual Polytope Projection (DPP)

Suppose we want to calculate a lasso solution at  $\lambda < \lambda_{\max}$ .  
The DPP rule discards the  $j^{\text{th}}$  variable if

$$\left| \mathbf{x}_j^T \mathbf{y} \right| < \lambda_{\max} - \|\mathbf{x}_j\|_2 \|\mathbf{y}\|_2 \frac{\lambda_{\max} - \lambda}{\lambda}$$

## Sequential DPP rule

Suppose we have the lasso solution  $\hat{\beta}(\lambda')$  at  $\lambda'$  and want to  
screen variables for solutions at  $\lambda < \lambda'$ . We discard the  $j^{\text{th}}$   
variable if

$$\left| \mathbf{x}_j^T (\mathbf{y} - \mathbf{X} \hat{\beta}(\lambda')) \right| < \lambda' - \|\mathbf{x}_j\|_2 \|\mathbf{y}\|_2 \frac{\lambda_{\max} - \lambda}{\lambda}$$

$$\left| \mathbf{x}_j^T \mathbf{y} \right| < \lambda_{\max} - \|\mathbf{x}_j\|_2 \|\mathbf{y}\|_2 \frac{\lambda_{\max} - \lambda}{\lambda}$$

$$\left| \mathbf{x}_j^T (\mathbf{y} - \mathbf{X} \hat{\beta}(\lambda')) \right| < \lambda' - \|\mathbf{x}_j\|_2 \|\mathbf{y}\|_2 \frac{\lambda_{\max} - \lambda}{\lambda}$$



# Global Strong Rule

Suppose we want to calculate a lasso solution at  $\lambda < \lambda_{\max}$ .

The global strong rule discards the  $j^{\text{th}}$  variable if

$$\left| \mathbf{x}_j^T \mathbf{y} \right| < \lambda - (\lambda_{\max} - \lambda) = 2\lambda - \lambda_{\max}$$

## Sequential Strong Rule

Suppose we have the lasso solution  $\hat{\beta}(\lambda')$  at  $\lambda'$  and want to screen variables for solutions at  $\lambda < \lambda'$ . We discard the  $j^{\text{th}}$  variable if

$$\left| \mathbf{x}_j^T (\mathbf{y} - \mathbf{X} \hat{\beta}(\lambda')) \right| < 2\lambda - \lambda'$$

2018-10-14

## Computation & optimization

### └ Screening Rules

### └ Global Strong Rule

#### Global Strong Rule

Suppose we want to calculate a lasso solution at  $\lambda < \lambda_{\max}$ .  
The global strong rule discards the  $j^{\text{th}}$  variable if

$$\left| \mathbf{x}_j^T \mathbf{y} \right| < \lambda - (\lambda_{\max} - \lambda) = 2\lambda - \lambda_{\max}$$

#### Sequential Strong Rule

Suppose we have the lasso solution  $\hat{\beta}(\lambda')$  at  $\lambda'$  and want to screen variables for solutions at  $\lambda < \lambda'$ . We discard the  $j^{\text{th}}$  variable if

$$\left| \mathbf{x}_j^T (\mathbf{y} - \mathbf{X} \hat{\beta}(\lambda')) \right| < 2\lambda - \lambda'$$

# Minorization-Maximization Algorithms (MMA)

- Problem: minimize  $f(\beta)$  over  $\beta \in \mathbb{R}^p$   
for  $f$  possibly non-convex
- Introduce additional variable  $\theta$
- Use  $\theta$  to majorize (bound from above) the objective  
function to be minimized

Majorization-Minimization Algorithms work analogously.

2018-10-14

Computation & optimization

└ Minor-Max Algorithms

└ Minorization-Maximization Algorithms (MMA)

Minorization-Maximization  
Algorithms (MMA)

- Problem: minimize  $f(\beta)$  over  $\beta \in \mathbb{R}^p$   
for  $f$  possibly non-convex
- Introduce additional variable  $\theta$
- Use  $\theta$  to majorize (bound from above) the objective  
function to be minimized

Majorization-Minimization Algorithms work analogously.

## MMA visually

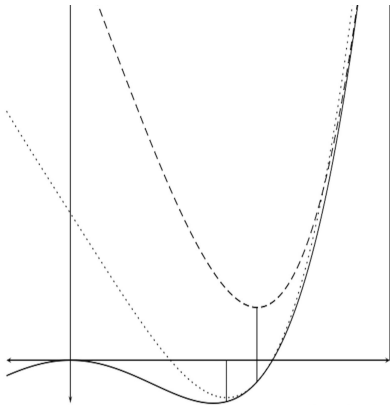


Figure: Figure from [de Leeuw, 2015]

2018-10-14

Computation & optimization

└ Minor-Max Algorithms

└ MMA visually

MMA visually

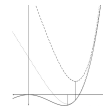


Figure: Figure from [de Leeuw, 2015]

## MMA analytically I

Def.  $\Psi : \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}$  **majorizes**  $f$  at  $\beta \in \mathbb{R}^p$  if

$$\forall \theta \in \mathbb{R}^p \quad \Psi(\beta, \theta) \geq f(\beta)$$

with equality for  $\theta = \beta$ .

Minor-Maxxalgorithm

- initialize  $\beta^0$
- update with  $\beta^{t+1} = \arg \min_{\beta \in \mathbb{R}^p} \Psi(\beta, \beta^t)$

2018-10-14

Computation &amp; optimization

└ Minor-Max Algorithms

└ MMA analytically I

MMA analytically I

Def.  $\Psi : \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}$  **majorizes**  $f$  at  $\beta \in \mathbb{R}^p$  if  
 $\forall \theta \in \mathbb{R}^p \quad \Psi(\beta, \theta) \geq f(\beta)$

with equality for  $\theta = \beta$ .

Minor-Maxalgorithm

- initialize  $\beta^0$
- update with  $\beta^{t+1} = \arg \min_{\beta \in \mathbb{R}^p} \Psi(\beta, \beta^t)$

# MMA analytically II

This scheme generates a sequence of  $\beta$ 's for which the cost  $f(\beta^t)$  is nonincreasing, because

$$f(\beta^t) \stackrel{(i)}{=} \Psi(\beta^t, \beta^t) \stackrel{(ii)}{\geq} \Psi(\beta^{t+1}, \beta^t) \stackrel{(iii)}{\geq} f(\beta^{t+1})$$

where

(i) & (iii) Definiton of majorize

(ii)  $\beta^{t+1}$  is a minimizer of  $\beta \mapsto \Psi(\beta, \beta^t)$

2018-10-14

Computation & optimization

└ Minor-Max Algorithms

└ MMA analytically II

MMA analytically II

This scheme generates a sequence of  $\beta$ 's for which the cost  $f(\beta^t)$  is nonincreasing, because

$$f(\beta^t) \stackrel{(i)}{=} \Psi(\beta^t, \beta^t) \stackrel{(ii)}{\geq} \Psi(\beta^{t+1}, \beta^t) \stackrel{(iii)}{\geq} f(\beta^{t+1})$$

where

(i) & (iii) Definiton of majorize

(ii)  $\beta^{t+1}$  is a minimizer of  $\beta \mapsto \Psi(\beta, \beta^t)$

for inequalities: show previous slide

# Biconvexity

Let's consider an example . . .

$$f(\alpha, \beta) = (1 - \alpha\beta)^2$$

2018-10-14

Computation & optimization  
└ Alternating Minimizations  
└ Biconvexity

Biconvexity

Let's consider an example . . .

$$f(\alpha, \beta) = (1 - \alpha\beta)^2$$

Mathematica: 3D plot  $(1-xy)^2$ ,  $x$  in  $[-2,2]$ ,  $y$  in  $[-2,2]$

The formula is a link.

# Biconvexity

Let's consider an example . . .

$$f(\alpha, \beta) = (1 - \alpha\beta)^2$$

Def. A function  $f(\alpha, \beta) : \mathbb{R}^m \times \mathbb{R}^n \rightarrow \mathbb{R}$  is **biconvex**, if for each  $\alpha \in \mathbb{R}^m$  the function  $\alpha \mapsto f(\alpha, \beta)$  is convex and for each  $\beta \in \mathbb{R}^n$  the function  $\beta \mapsto f(\alpha, \beta)$  is convex. Analogously, a set  $\mathcal{C} \subseteq \mathcal{A} \times \mathcal{B}$ , for  $\mathcal{A}, \mathcal{B}$  convex sets, is called biconvex, if it is convex

2018-10-14

## Computation & optimization

- Alternating Minimizations
  - Biconvexity

Mathematica: 3D plot  $(1-xy)^2$ ,  $x$  in  $[-2,2]$ ,  $y$  in  $[-2,2]$   
The formula is a link.

Let's consider an example . . .

$$f(\alpha, \beta) = (1 - \alpha\beta)^2$$

Def. A function  $f(\alpha, \beta) : \mathbb{R}^m \times \mathbb{R}^n \rightarrow \mathbb{R}$  is biconvex, if for each  $\alpha \in \mathbb{R}^m$  the function  $\alpha \mapsto f(\alpha, \beta)$  is convex and for each  $\beta \in \mathbb{R}^n$  the function  $\beta \mapsto f(\alpha, \beta)$  is convex. Analogously, a set  $\mathcal{C} \subseteq \mathcal{A} \times \mathcal{B}$ , for  $\mathcal{A}, \mathcal{B}$  convex sets, is called biconvex, if it is convex

# Alternate Convex Search

Block coordinate descent applied to  $\alpha$  and  $\beta$  blocks

1. Initialize  $(\alpha^0, \beta^0)$  at some point in the biconvex set to minimize over
2. For  $t = 0, 1, 2, \dots$ 
  - (i) Fix  $\beta = \beta^t$  and update  $\alpha^{t+1} \in \arg \min_{\alpha \in \mathcal{C}_{\beta^t}} f(\alpha, \beta^t)$
  - (ii) Fix  $\alpha = \alpha^{t+1}$  and update  $\beta^{t+1} \in \arg \min_{\beta \in \mathcal{C}_{\alpha^{t+1}}} f(\alpha^{t+1}, \beta)$

For a function bounded from below, the algorithm converges to a partial optimum (i.e. as biconvexity, only optimal in one coordinate if the other coordinate is fixed).

Block coordinate descent applied to  $\alpha$  and  $\beta$  blocks

1. Initialize  $(\alpha^0, \beta^0)$  at some point in the biconvex set to minimize over
2. For  $t = 0, 1, 2, \dots$ 
  - (i) Fix  $\beta = \beta^t$  and update  $\alpha^{t+1} \in \arg \min_{\alpha \in \mathcal{C}_{\beta^t}} f(\alpha, \beta^t)$
  - (ii) Fix  $\alpha = \alpha^{t+1}$  and update  $\beta^{t+1} \in \arg \min_{\beta \in \mathcal{C}_{\alpha^{t+1}}} f(\alpha^{t+1}, \beta)$

For a function bounded from below, the algorithm converges to a partial optimum (i.e. as biconvexity, only optimal in one coordinate if the other coordinate is fixed).



# References I



Trevor Hastie, Robert Tibshirani, and Martin Wainwright (2015)

Statistical learning with sparsity: the Lasso and generalizations

*CRC Press*; Boca Raton, FL



Jan De Leeuw (2015)

Block Relaxation Methods in Statistics

[doi.org/10.13140/RG.2.1.3101.9607](https://doi.org/10.13140/RG.2.1.3101.9607) (last accessed: 02.10.18)



S. Boyd

Alternating Direction Method of Multipliers

[https://web.stanford.edu/~boyd/papers/pdf/admm\\_slides.pdf](https://web.stanford.edu/~boyd/papers/pdf/admm_slides.pdf)  
(last accessed: 14.10.18)



Geoff Gordon and Ryan Tibshirani (2012)

Uses of Duality

<https://www.cs.cmu.edu/~ggordon/10725-F12/slides/18-dual-uses.pdf> (last accessed: 14.10.18)

2018-10-14

Computation & optimization  
└ Alternating Minimizations

└ References

References I

-  Trevor Hastie, Robert Tibshirani, and Martin Wainwright (2015)  
Statistical learning with sparsity: the Lasso and generalizations  
*CRC Press*; Boca Raton, FL
-  Jan De Leeuw (2015)  
Block Relaxation Methods in Statistics  
[doi.org/10.13140/RG.2.1.3101.9607](https://doi.org/10.13140/RG.2.1.3101.9607) (last accessed: 02.10.18)
-  S. Boyd  
Alternating Direction Method of Multipliers  
[https://web.stanford.edu/~boyd/papers/pdf/admm\\_slides.pdf](https://web.stanford.edu/~boyd/papers/pdf/admm_slides.pdf)  
(last accessed: 14.10.18)
-  Geoff Gordon and Ryan Tibshirani (2012)  
Uses of Duality  
<https://www.cs.cmu.edu/~ggordon/10725-F12/slides/18-dual-uses.pdf> (last accessed: 14.10.18)

## References II



Paul Rubin (2016)

What are the advantages of convex optimization compared to more general optimization problems?

[https://www.quora.com/](https://www.quora.com/What-are-the-advantages-of-convex-optimization-compared-to-more-general-optimization-problems?m=1)

[What-are-the-advantages-of-convex-optimization-compared-to-more-general-optimization-problems?m=1](https://www.quora.com/What-are-the-advantages-of-convex-optimization-compared-to-more-general-optimization-problems?m=1)  
(last accessed: 14.10.18)

2018-10-14

Computation & optimization  
└ Alternating Minimizations  
└ References

Comments . . .  
Questions . . .  
Suggestions . . .

2018-10-14

Computation & optimization  
└ Alternating Minimizations

Comments . . .  
Questions . . .  
Suggestions . . .

That's it.  
Thanks for listening.

Fill out your feedback sheets!

2018-10-14

Computation & optimization  
└ Alternating Minimizations

That's it.  
Thanks for listening.

Fill out your feedback sheets!