

Variable selection and outlier detection as a Mixed Integer Program (MIP) and Image segmentation

Stéphane Canu

asi.insa-rouen.fr/enseignants/~scanu

S4D 2018 : Research Summer School on Statistics for Data Science

Caen, June 18, 2018

Road map

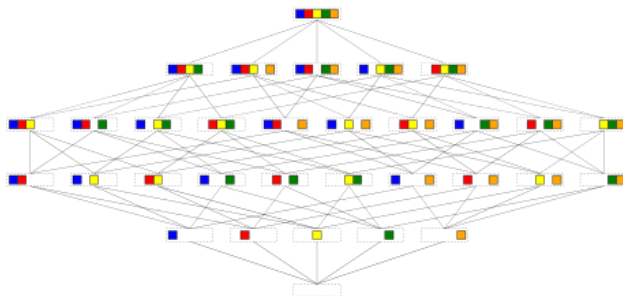
- 1 Examples of combinatorial problems in machine learning
- 2 Mixed integer (binary) programming (MIP)
- 3 L_0 proximal algorithm
- 4 Implementation
- 5 MIP for image processing

NP Hard =



Variable selection

$$f(x_1, \dots, x_j, \dots, x_p) = \sum_{j=1}^{p=5} x_j w_j$$

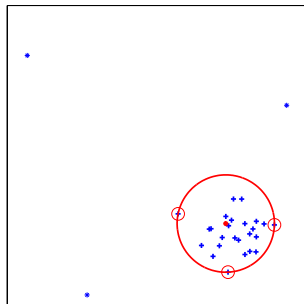
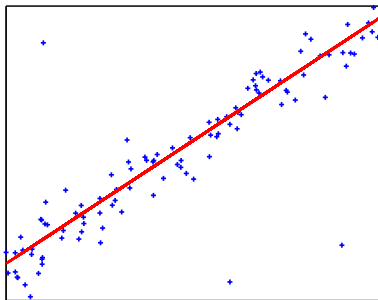


e.g. gene selection from microarray data, find w sparse such that $Xw = y$

Fit the data **and** remove useless variables

Enumerate of all possible **combinations** and choose

Outlier detection



Fit the data **and** remove useless observations (outliers)

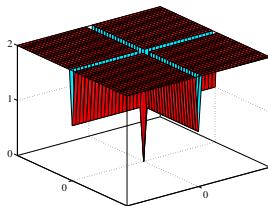
Enumerate of all possible **point configurations** and choose

What's common?

Example (The counting function)

$$\begin{aligned} c : \mathbb{R}^p &\longrightarrow \mathbb{R} \\ w &\longmapsto c(w) = \text{the number of nonzero components } w_i \text{ of } w \end{aligned}$$

It is often called the 0-norm denoted by $c(w) = \|w\|_0$.



Minimize a nonconvex nonsmooth target function or constraint

Nonconvex Nonsmooth problems in machine learning

Many **lattice based** problems

- variable selection, outlier detection, clustering,
- image processing, total variations,
- discrete artificial vision,
- sensor placement,
- distribution factorization,
- low rank factorization.

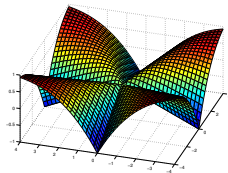
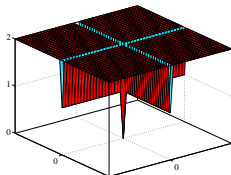


3 ways of solving combinatorial problems

- local optimization (in general)
 - ▶ Continuous relaxations (L_1 penalty, DC...)
 - ▶ Combinatorial algorithms (greedy search, spanning tree...)
- global optimization
 - ▶ **Mixed integer programming** (difficult to scale to large problems)

Road map

- 1 Examples of combinatorial problems in machine learning
- 2 Mixed integer (binary) programming (MIP)
- 3 L_0 proximal algorithm
- 4 Implementation
- 5 MIP for image processing



Variable selection with binary variables

Definition (the least square variable selection problem)

$$\begin{cases} \min_{w \in \mathbb{R}^p} & \|Xw - y\|^2 & \leftarrow \text{fit the data} \\ \text{s.t.} & \|w\|_0 \leq k & \leftarrow \text{with } k \text{ variables} \end{cases}$$

- introduce p new binary variable $z \in \{0, 1\}^p$
- for useless variables: $z_j = 0 \Leftrightarrow w_j = 0$ \rightarrow a coupling mechanism
- $\|w\|_0 = \sum_{j=1}^p z_j$

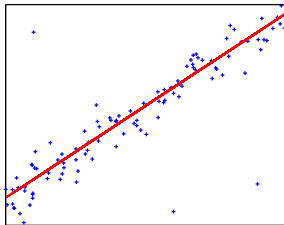
Definition (the LS variable selection problem with binary variables)

$$\begin{cases} \min_{w \in \mathbb{R}^p, z \in \{0, 1\}^p} & \|Xw - y\|^2 \\ \text{s.t.} & \|w\|_0 = \sum_{i=1}^p z_i \leq k \\ & z_j = 0 \Leftrightarrow w_j = 0, \quad j = 1, p \end{cases}$$

Outlier detection with binary variables

Introducing outliers variables $o \in \mathbb{R}^n$

$$y = Xw + \varepsilon + o, \quad o_i = \begin{cases} y_i - x_i^t w & \text{if } i \text{ outlier} \\ 0 & \text{else} \end{cases}$$



The least square (trimmed) regression problem with k outliers [?]

$$\begin{cases} \min_{w \in \mathbb{R}^p, o \in \mathbb{R}^n} & \frac{1}{2} \|Xw + o - y\|^2 \\ \text{s.t.} & \|o\|_0 \leq k \end{cases}$$

Introduce
binary
variables

$$i = 1, n \quad \begin{cases} t_i = 0 & (x_i, y_i) \text{ is an outlier} & o_i \neq 0 \\ t_i = 1 & (x_i, y_i) \text{ is NOT an outlier} & o_i = 0 \end{cases}$$

$$\|o\|_0 = \sum_{i=1}^n (1 - t_i)$$

Bi robust regression

Variable selection

AND

outlier detection

$$\left\{ \begin{array}{l} \min_{w \in \mathbb{R}^p} \quad \frac{1}{2} \|Xw - y\|^2 \\ \text{s.t.} \quad \|w\|_0 \leq k_v \end{array} \right. \quad \text{and} \quad \left\{ \begin{array}{l} \min_{w \in \mathbb{R}^p, o \in \mathbb{R}^n} \quad \frac{1}{2} \|Xw + o - y\|^2 \\ \text{s.t.} \quad \|o\|_0 \leq k_o \end{array} \right.$$

LS regression with variable selection AND outlier detection

Given k_v the number of variable required and k_o the number of outliers

$$\left\{ \begin{array}{l} \min_{w \in \mathbb{R}^p, o \in \mathbb{R}^n} \quad \|Xw - y - o\|^2 \\ \text{s.t.} \quad \|w\|_0 \leq k_v \\ \quad \quad \|o\|_0 \leq k_o. \end{array} \right. \quad (1)$$

Bi robust regression with binary variables

p binary variables $z_j \in \{0, 1\}$

- Variables
- Outliers

$$\|w\|_0 = \sum_{j=1}^p z_j \quad \text{and} \quad z_j = 0 \Leftrightarrow w_j = 0,$$

$$\left\{ \begin{array}{l} \min_{w, o} \|Xw - y - o\|^2 \\ \text{s.t.} \quad \|w\|_0 \leq k_v \\ \quad \|o\|_0 \leq k_o. \end{array} \right.$$

n binary variables $t_i \in \{0, 1\}$

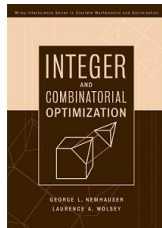
$$\|o\|_0 = \sum_{i=1}^n (1 - t_i) \quad \text{and} \quad 1 - t_i = 0 \Leftrightarrow o_i = 0,$$

Bi robust regression

$$\left\{ \begin{array}{l} \min_{w \in \mathbb{R}^p, o \in \mathbb{R}^n, z \in \{0,1\}^p, t \in \{0,1\}^n} \|Xw - y - o\|^2 \\ \text{s.t.} \quad \|w\|_0 = \sum z_j \leq k_v \\ \quad z_j = 0 \Leftrightarrow w_j = 0, \quad j = 1, p \\ \quad \|o\|_0 = \sum t_i \leq k_o \\ \quad 1 - t_i = 0 \Leftrightarrow o_i = 0, \quad i = 1, n \end{array} \right. \quad (2)$$

So far...

- combinatorial problems can be formulated using binary variables
- we have **mixed binary optimization** problem
- How to solve them?
 - ▶ reformulations
 - ▶ towards stronger relaxations
 - ▶ nice initialization



Mixed integer linear program (MILP)

- linear cost
- linear constraints
- **integer** and continuous variables

Definition (mixed integer linear program – MILP (canonical form))

$$\left\{ \begin{array}{ll} \min_{w \in \mathbb{R}^p, z \in \mathbb{N}^q} & f(w, z) = c^t w + d^t z \quad \leftarrow \text{linear} \\ \text{s.t.} & A w + B z \leq b \quad \leftarrow \text{linear} \\ & w \geq 0, \end{array} \right.$$

for some given $c \in \mathbb{R}^p$, $d \in \mathbb{R}^m$, $A \in \mathbb{R}^{m \times p}$, $B \in \mathbb{R}^{m \times q}$ and $b \in \mathbb{R}^m$.

- A **mixed binary linear program** is a MILP with $z \in \{0, 1\}^q$ binary.
- When its domain is not empty and bounded, a MILP admit a unique global minimum.

Mixed integer quadratic program (MIQP)

- quadratic cost
- linear constraints
- integer and continuous variables

Definition (mixed integer quadratic program – MIQP)

$$\left\{ \begin{array}{ll} \min_{x=(w,z) \in \mathbb{R}^p \times \mathbb{N}^q} & f(x) = \frac{1}{2}x^t Qx + c^t x \quad \leftarrow \text{quadratic} \\ \text{s.t.} & Ax \leq b \quad \leftarrow \text{linear} \\ & x \geq 0, \end{array} \right.$$

for some given symmetric matrix $Q \in \mathbb{R}^{(p+q) \times (p+q)}$

Mixed integer quadratically constrained quadratic program (MIQCP).

- quadratic cost, **quadratic constraints**, integer and continuous variables

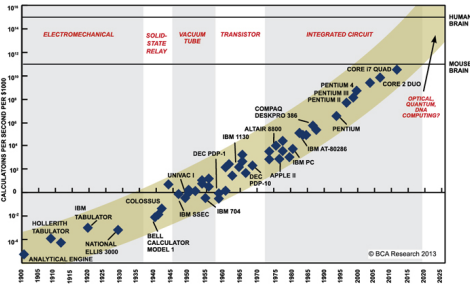
Problems hierarchy

MILP (= MIQP with $Q = 0$) \subset MIQP \subset MIQCP

Progresses in MILP

in 1989

MILP is a powerful modeling tool, “They are, however, theoretically complicated and computationally **cumbersome**” [?]



from 1996 to 2016 [? ?]

	improvement factor
machine solver formulation	$\times 2^{10} = 1000 - 1600$ $\times 1000 - 3600$???
global	$\times 1 - 5 \cdot 10^6$

a year to solve 10 – 20 years ago → now 30 seconds

“mixed integer linear techniques are nowadays **mature**, that is fast, robust, and are able to solve problems with up to millions of variables” [?]

Mixed integer software (available with matlab)

Software package	Matlab function
Open source	
GLPK	<u>glpk</u> for mixed integer linear programming
COIN OR	(not matlab yet)
Commercial	
Matlab	<u>intlinprog</u> for mixed integer linear programming
CPLEX	<u>cplexmilp</u> for mixed integer linear programming <u>cplexmiqp</u> for mixed integer quadratic programming <u>cplexmiqcp</u> for mixed integer quadratically constrained pg
GUROBI	<u>gurobi</u> for MILP, MIQP and MIQCQP
Mosek	<u>mosekopt</u> for MILP, MIQP and MIQCQP
SCIP	<u>opti_scip</u> for MILP, MIQP and MIQCQP

Mixed Integer Linear Programming Benchmark (MILP2010)

recommend CPLEX, GUROBI and Mosek (NOT intlinprog)

<http://plato.asu.edu/ftp/milpc.html>

So far

- The MIP hierarchy: MILP (linear) \subset MIQP (quadratic) \subset MIQCQP
- use standard MIP software
- the joy of having an exact solution
- formulate ML problems as a MILP or MIQP (if possible)



Bi robust regression

Variable selection

AND

outlier detection

$$\left\{ \begin{array}{l} \min_{w \in \mathbb{R}^p} \quad \frac{1}{2} \|Xw - y\|^2 \\ \text{s.t.} \quad \|w\|_0 \leq k_v \end{array} \right. \quad \text{and} \quad \left\{ \begin{array}{l} \min_{w \in \mathbb{R}^p, o \in \mathbb{R}^n} \quad \frac{1}{2} \|Xw + o - y\|^2 \\ \text{s.t.} \quad \|o\|_0 \leq k_o \end{array} \right.$$

LS regression with variable selection AND outlier detection

Given k_v the number of variable required and k_o the number of outliers

$$\left\{ \begin{array}{l} \min_{w \in \mathbb{R}^p, o \in \mathbb{R}^n} \quad \|Xw - y - o\|^2 \\ \text{s.t.} \quad \|w\|_0 \leq k_v \\ \quad \quad \|o\|_0 \leq k_o. \end{array} \right. \quad (3)$$

LS with fixed cardinality as a MIQP: the big M constraint

Assuming we know an **upper bound** M for w

$$\|w\|_0 \leq k \quad \Leftrightarrow \quad \begin{cases} z_j \in \{0, 1\}, & j = 1 : p \\ \sum_{i=1}^p z_j \leq k \\ |w_j| \leq z_j M \end{cases}$$

For useless variables:
 $z_j = 0 \Rightarrow w_j = 0$

LS with fixed cardinality as a MIQP [?]

$$\begin{cases} \min_{w \in \mathbb{R}^p, z \in \{0,1\}^p} & \frac{1}{2} \|Xw - y\|_2^2 \\ \text{s.t.} & \sum_{j=1}^p z_j \leq k \\ \text{and} & |w_j| \leq z_j M \quad j = 1, p \end{cases}$$

Variable selection AND outlier detection as a MILP

$$q \in \{1, 2\} \quad \left\{ \begin{array}{ll} \min_{w \in \mathbb{R}^p, o \in \mathbb{R}^n} & \|Xw - y - o\|_q^q \\ \text{s.t.} & \|w\|_0 \leq k_v \\ & \|o\|_0 \leq k_o. \end{array} \right.$$

$$q = 1$$

$$\left\{ \begin{array}{ll} \min_{w \in \mathbb{R}^p, o, \varepsilon^+, \varepsilon^- \in \mathbb{R}^n, z \in \{0,1\}^p, t \in \{0,1\}^n} & \sum_{i=1}^n \varepsilon_i^+ + \varepsilon_i^- \\ \text{s.t.} & \varepsilon_i^+ - \varepsilon_i^- = x_i^t w + o_i - y_i \quad i = 1, n \\ & \sum_{j=1}^p z_j \leq k_v \\ & |w_j| \leq z_j M_v \quad j = 1, p \\ & \sum_{i=1}^n (1 - t_i) \leq k_o \\ & |o_i| \leq t_i M_o \quad i = 1, n \\ & 0 \leq \varepsilon_i^+, 0 \leq \varepsilon_i^- \quad i = 1, n. \end{array} \right.$$

LSE with fixed cardinality as a MIQP with SOS constraints

Variable selection: $z_j = 0 \Rightarrow w_j = 0$ either $w_j = 0$ or $1 - z_j = 0$

Special ordered set (SOS) of type 1: at most one variable in the set can take a nonzero value,

$$w_j = 0 \text{ or } 1 - z_j = 0 \Leftrightarrow (w_j, 1 - z_j) : \text{SOS}$$

MIQP using special ordered set (SOS) of type 1

$$\left\{ \begin{array}{ll} \min_{w \in \mathbb{R}^p, e \in \mathbb{R}^n, z \in \{0,1\}^p} & \sum_{i=1}^n \frac{1}{2} (X_i^t w - y_i)^2 \quad \leftarrow \text{data loss} \\ \text{s.t.} & \sum_{j=1}^p z_j \leq k \quad \leftarrow \text{at most } k \text{ non 0 variables} \\ & (w_j, 1 - z_j) : \text{SOS} \quad j = 1, p \end{array} \right.$$

Variable selection AND outlier detection as a MIQP

$$q \in \{1, 2\} \quad \left\{ \begin{array}{ll} \min_{w \in \mathbb{R}^p, o \in \mathbb{R}^n} & \|Xw - y - o\|_q^q \\ \text{s.t.} & \|w\|_0 \leq k_v \\ & \|o\|_0 \leq k_o. \end{array} \right.$$

$q = 2$

$$\left\{ \begin{array}{ll} \min_{w \in \mathbb{R}^p, o \in \mathbb{R}^n, z \in \{0,1\}^p, t \in \{0,1\}^n} & (y - Xw - o)^t (y - Xw - o) \\ \text{s.t.} & \sum_{j=1}^p z_j = k_v \\ & \sum_{i=1}^n t_i \leq k_o \\ & (w_j, 1 - z_j) : \text{SOS} & j = 1, p \\ & (o_i, 1 - t_i) : \text{SOS} & i = 1, n. \end{array} \right.$$

Balls and Triks: the convex hull of the feasible set

$$\text{Conv} \left(\left\{ w \mid |w_j| \leq z_j M \text{ and } \sum_{j=1}^p z_j \leq k \right\} \right) = \left\{ w \mid \|w\|_\infty \leq M \text{ and } \|w\|_1 \leq kM \right\}$$

MIQP: a more structured representation [?]

$$\left\{ \begin{array}{ll} \min_{w \in \mathbb{R}^p, e \in \mathbb{R}^n, z \in \{0,1\}^p} & \sum_{i=1}^n \frac{1}{2} (X_i^t w - y_i)^2 \quad \leftarrow \text{data loss} \\ \text{s.t.} & \sum_{j=1}^p z_j \leq k \quad \leftarrow \text{at most } k \text{ non 0 variables} \\ & (w_j, 1 - z_j) : \text{SOS} \quad \begin{array}{l} j = 1, p \\ j = 1, p \end{array} \\ & |w_j| \leq M_\infty \\ & \sum_{j=1}^p |w_j| \leq M_1 \end{array} \right.$$

[?] claim: Adding these bounds typically leads to improved performance of the MIO, especially in delivering lower bound certificates

Balls and Triks: the convex hull of the feasible set

$$\mathcal{S} = \left\{ w, o \mid \sum_{j=1}^p z_j \leq k_v, |w_j| \leq z_j M_v, \sum_{i=1}^n (1 - t_i) \leq k_o, |\tau_i| \leq t_i M_o \right\},$$

$$\text{Conv}(\mathcal{S}) = \left\{ w, o \mid \|w\|_\infty \leq M_v, \|o\|_\infty \leq M_o, \|w\|_1 \leq k_v M_v, \|o\|_1 \leq k_o M_o \right\}$$

$$\left\{ \begin{array}{ll} \min_{w, o, z \in \{0,1\}^p, t \in \{0,1\}^n} & \frac{1}{q} \|Xw + o - y\|_q^q \\ \text{s.t.} & \sum_{j=1}^p z_j \leq k_v, \quad (w_j, 1 - z_j) : \text{SOS} \quad j = 1, p \\ & \sum_{i=1}^n (1 - t_i) \leq k_o, \quad (\tau_i, 1 - t_i) : \text{SOS} \quad i = 1, n \\ & \|w\|_1 \leq k_v M_v, \quad \|w\|_\infty \leq M_v \\ & \|o\|_1 \leq k_o M_o, \quad \|o\|_\infty \leq M_o, \end{array} \right. \quad (4)$$

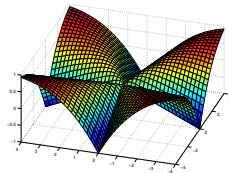
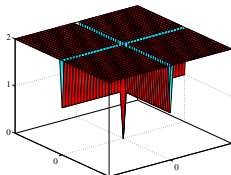
with problem-dependent constants M_v and M_o .

So far...

- birobust regression as a MIP
 - ▶ for **variable selection AND outlier detection** in regression
 - ▶ and in quantile regression, SVM, logistic regression
 - ▶ reformulation (practical matter)
- efficient software for **moderate size problem**
- for large size: use **first order algorithms**

Road map

- 1 Examples of combinatorial problems in machine learning
- 2 Mixed integer (binary) programming (MIP)
- 3 L_0 proximal algorithm**
- 4 Implementation
- 5 MIP for image processing



Variable selection: a specific case with a closed-form solution

Definition (the least square variable selection problem with $X = Id$)

given $k < p$

$$\begin{cases} \min_{\mathbf{u} \in \mathbf{R}^p} & \|\mathbf{u} - \mathbf{w}\|^2 & \longleftarrow \text{fit the data} \\ \text{s.t.} & \|\mathbf{u}\|_0 \leq k & \longleftarrow \text{with } k \text{ variables} \end{cases}$$

sort $|\mathbf{w}|$: $|w_{(1)}| \geq |w_{(2)}| \geq \dots |w_{(j)}| \geq \dots |w_{(p)}|$

Closed-form solution: the hard thresholding operator

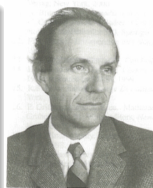
$$u_i^* = H_k(\mathbf{w}) = \begin{cases} w_j & \text{if } j \in \{(1), \dots, (k)\} \\ 0 & \text{else} \end{cases}$$

Proximity operator

Definition (Proximity operator [?])

The Proximity operator of a function h is:

$$\begin{aligned} \mathbf{prox}_h : \mathbb{R}^p &\longrightarrow \mathbb{R} \\ w &\longmapsto \mathbf{prox}_h(w) = \arg \min_{u \in \mathbb{R}^p} h(u) + \frac{1}{2} \|u - w\|^2 \end{aligned}$$



Example

$h(w) = 0$	$\mathbf{prox}_h(w) = w$	
$h(w) = \rho \mathbf{pen}_\lambda(w)$	$\mathbf{prox}_h(w) = \mathbf{shr}_{\rho\lambda}(w)$	shrinkage
$h(w) = \mathbb{I}_C(w)$	$\mathbf{prox}_h(w) = \arg \min_{u \in C} \frac{1}{2} \ u - w\ ^2$	projection

The proximity operator as a projection

$$\mathbf{prox}_{\mathbb{I}_{\{\|w\|_0 \leq k\}}}(w) = \arg \min_{\|u\|_0 \leq k} \frac{1}{2} \|u - w\|^2 = H_k(w) = \begin{cases} w_i & \text{if } i \in \{(1), \dots, (k)\} \\ 0 & \text{else} \end{cases}$$

The projected gradient (L_0 projection or proximal)

$$\text{for solving } \begin{cases} \min_{w \in \mathbb{R}^p} & \frac{1}{2} \|Xw - y\|^2 \\ \text{s.t.} & \|w\|_0 \leq k_v \end{cases}$$

Algorithm 1 L_0 gradient projection algorithm [?]

Data: X, y, w initialization

Result: w

while *not converged* **do**

$$g \leftarrow \nabla g(w) = X^\top (Xw - y),$$

the gradient

$$\rho \leftarrow \text{choose a stepsize}$$

$$d \leftarrow w - \rho g,$$

forward (explicit)

$$w \leftarrow H_k(d),$$

the projection-proximal step, backward (implicit)

end

if $\varepsilon \leq \rho \leq \frac{1}{\|X^\top X\|}$, it converges towards a local minimum [?] since its objective function satisfies the Kurdyka-Lojasiewicz inequality.

Proximal alternating linearized minimization (PALM)

$$\left\{ \begin{array}{ll} \min_{w \in \mathbb{R}^p, o \in \mathbb{R}^n} & \frac{1}{2} \|Xw + o - y\|^2 \\ \text{s.t.} & \|w\|_0 \leq k_v \\ & \|o\|_0 \leq k_o \end{array} \right.$$

given o

$$\left\{ \begin{array}{ll} \min_{w \in \mathbb{R}^p} & \frac{1}{2} \|Xw + o - y\|^2 \\ \text{s.t.} & \|w\|_0 \leq k_v \end{array} \right.$$

given w

$$\left\{ \begin{array}{ll} \min_{o \in \mathbb{R}^n} & \frac{1}{2} \|o - (y - Xw)\|^2 \\ \text{s.t.} & \|o\|_0 \leq k_o \end{array} \right.$$

Algorithm 2 Proximal alternating linearized minimization (PALM) [?]

Data: X, y initialization $w, o = 0$

Result: w, o

while *not converged* **do**

$$d \leftarrow w - \rho_v X^T (Xw + o - y),$$

variable selection

$$w \leftarrow H_{k_v}(d),$$

$$\delta \leftarrow o - \rho_o (Xw + o - y),$$

eliminating outliers

$$o \leftarrow H_{k_o}(\delta),$$

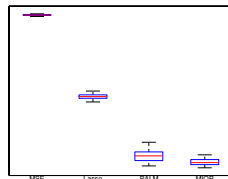
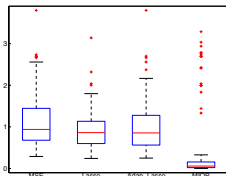
end

Prox summary

- PALM is fast and scalable
- convergence proofs towards a local minimum
- improvement:
 - ▶ accelerations: FISTA and others
 - ▶ Newton proximal
 - ▶ more improvement with randomization

Road map

- 1 Examples of combinatorial problems in machine learning
- 2 Mixed integer (binary) programming (MIP)
- 3 L_0 proximal algorithm
- 4 Implementation**
- 5 MIP for image processing



Combine the best of the two worlds

Combine the best of the two worlds [?]

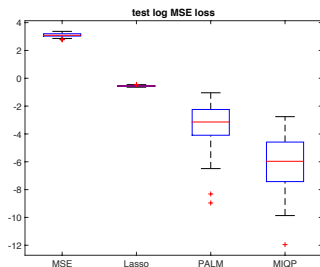
1. $(w, o) \leftarrow$ PALM alternating proximal gradient method
2. use w and o as a warm start for MIP (with Cplex)
 - ▶ $(w, o) \leftarrow$ Polish coefficients on the active set
 - ▶ initialize the constants M_w, M_o

Experimental setup

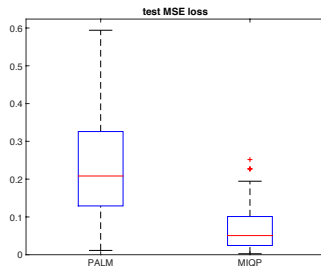
- My mac
- Matlab
- Cplex 12.6.1 (cplexmiqp)
- time out = 5 min

Variable selection AND outlier detection on a toy dataset

- $y = Xw + o + \varepsilon$
- $n = 300$ observations with $p = 25$ variables
- linear model with ε a centered gaussian noise with $\text{SNR} \approx 1$
- $k_o = 50$ outliers and $k_v = 5$ non zeros variables
- 100 repetitions



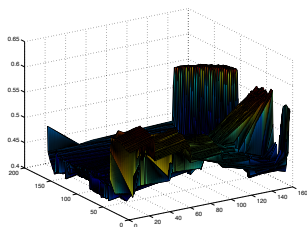
(log) performances



zoom

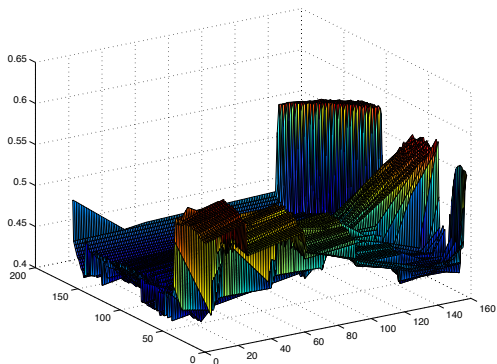
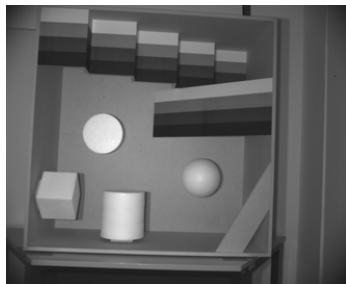
Road map

- 1 Examples of combinatorial problems in machine learning
- 2 Mixed integer (binary) programming (MIP)
- 3 L_0 proximal algorithm
- 4 Implementation
- 5 MIP for image processing



A Second Order Total Variation MIP Model for Image Segmentation Using
the L_0 Norm
(*on going work*)

The problem: depth estimation



Fundamental hypothesis

piecewise linear model

The second order TV Potts model in the 1d case

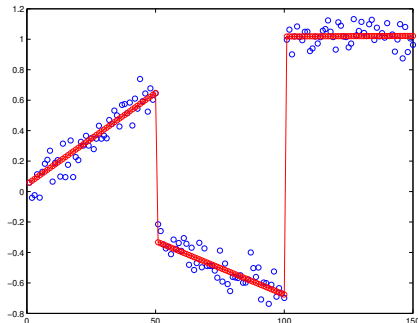
The objective: given n observations (z_1, \dots, z_n)

Retrieve the best linear fit with k discontinuities

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^n} \quad & \|\mathbf{w} - \mathbf{z}\|_1 \\ \text{s.t.} \quad & \|\nabla_x^2 \mathbf{w}\|_0 \leq k \end{aligned}$$

piece wise linear model

$$\begin{aligned} w_i &= a_\ell i + b_\ell \\ \ell &= 1, \dots, k^* + 1; \end{aligned}$$



$\|\nabla_x^2 \mathbf{w}\|_0 \leq k$ with k small

- $\|\nabla_x^2 \mathbf{w}\|_0 = 0$ impose linear model
- $\|\nabla_x^2 \mathbf{w}\|_0 \neq 0$ allow discontinuity

The second order TV Potts model as a MIP

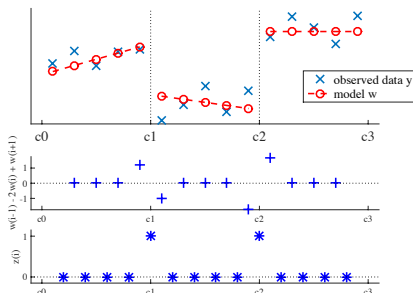
$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^n} \quad & \|\mathbf{w} - \mathbf{z}\|_1 \\ \text{s.t.} \quad & \|\nabla_x^2 \mathbf{w}\|_0 \leq k \end{aligned}$$

$$\nabla_x^2 \mathbf{w} \propto w_{i-1} - 2w_i + w_{i+1}$$

$$\begin{cases} w_{i-1} - 2w_i + w_{i+1} = 0 & \Rightarrow x_i = 1, \quad i = 2, \dots, n-1, \\ w_{i-1} - 2w_i + w_{i+1} \neq 0 & \Rightarrow x_i = 0, \quad i = 2, \dots, n-1 \end{cases}$$

For a given M large enough (the big M trick)

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^n, \mathbf{x} \in \{0,1\}^{n-1}} \quad & \sum_{i=1}^n |w_i - z_i| \\ \text{s.t.} \quad & |w_{i-1} - 2w_i + w_{i+1}| \leq M(x_{i-1} + x_i), \quad i = 2, \dots, n-1, \\ & \sum_{i=1}^{n-1} x_i \leq k. \end{aligned}$$



Eliminating the absolute values: positive and negative parts

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^n, \mathbf{x} \in \{0,1\}^{n-1}} \quad & \sum_{i=1}^n |w_i - z_i| \\ \text{s.t.} \quad & |w_{i-1} - 2w_i + w_{i+1}| \leq M(x_{i-1} + x_i), \quad i = 2, \dots, n-1, \\ & \sum_{i=1}^{n-1} x_i \leq k. \end{aligned}$$

Positive and negative parts: $z_+, z_- \geq 0$

$$z = z_+ - z_- \quad \text{and} \quad |z| = z_+ + z_-$$

$$\begin{aligned} \minimize_{\mathbf{w}, \varepsilon^+, \varepsilon^- \in \mathbb{R}^n, \mathbf{x} \in \{0,1\}^{n-1}} \quad & \sum_{i=1}^n (\varepsilon_i^+ + \varepsilon_i^-) \\ \text{s.t.} \quad & \mathbf{w} - \mathbf{z} = \varepsilon^+ - \varepsilon^- \\ & w_{i-1} - 2w_i + w_{i+1} \leq M(x_{i-1} + x_i), \quad i = 2, \dots, n-1 \\ & -w_{i-1} + 2w_i - w_{i+1} \leq M(x_{i-1} + x_i), \quad i = 2, \dots, n-1 \\ & \sum_{i=1}^{n-1} x_i \leq k \\ & 0 \leq \varepsilon^+, 0 \leq \varepsilon^- \end{aligned}$$

How to accelerate the solver?

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^n, \mathbf{x} \in \{0,1\}^{n-1}} \quad & \sum_{i=1}^n |w_i - z_i| \\ \text{s.t.} \quad & |w_{i-1} - 2w_i + w_{i+1}| \leq M(x_{i-1} + x_i), \quad i = 2, \dots, n-1, \\ & \sum_{i=1}^{n-1} x_i \leq k. \end{aligned}$$

Initialization: use **first order algorithm** (ADMM, proximal, greedy...)

- initialize \mathbf{w} and \mathbf{x}
- initialize parameter M

Stronger constraints

- based on the convex hull of the actual constraints
- typically: $\|D_n \mathbf{w}\|_1 \leq k$
- local implied bound cuts to treat the big-M [?]

[?] claim: *Adding these bounds typically leads to improved performance of the MIO, especially in delivering lower bound certificates*

First order algorithm :toward a local minimum

$$D_n = \begin{pmatrix} 1 & -2 & 1 & 0 & \dots & 0 \\ 0 & 1 & -2 & 1 & \dots & 0 \\ & & & \ddots & & \\ 0 & \dots & 0 & 1 & -2 & 1 \end{pmatrix} \in \mathbb{R}^{n-2 \times n}$$

$$\begin{cases} \min_{\mathbf{w} \in \mathbb{R}^n} & \|\mathbf{w} - \mathbf{z}\|_1 \\ \text{s.t.} & \|D_n \mathbf{w}\|_0 \leq k \end{cases} \Rightarrow \begin{cases} \min_{\mathbf{w} \in \mathbb{R}^n, \mathbf{r} \in \mathbb{R}^{n-2}} & \|\mathbf{w} - \mathbf{z}\|_1 \\ \text{s.t.} & \|\mathbf{r}\|_0 \leq k \quad \text{and} \quad \mathbf{r} = D_n \mathbf{w} \end{cases}$$

The augmented Lagrangian with $\lambda > 0$

$$L(\mathbf{w}, \mathbf{r}, \Lambda) = \|\mathbf{w} - \mathbf{z}\|_1 + \Lambda^t (\mathbf{r} - D_n \mathbf{w}) + \frac{\lambda}{2} \|\mathbf{r} - D_n \mathbf{w}\|^2 + I_{\{\|\mathbf{r}\|_0 \leq k\}}$$

The ADMM algorithm:

$$\begin{aligned} \mathbf{w}^{k+1} &= \arg \min_{\mathbf{w}} L(\mathbf{w}^k, \mathbf{r}^k, \Lambda^k) \\ \mathbf{r}^{k+1} &= \arg \min_{\mathbf{r}} L(\mathbf{w}^{k+1}, \mathbf{r}^k, \Lambda^k) \\ \Lambda^{k+1} &= \Lambda^k + \rho (D_n \mathbf{w}^{k+1} - \mathbf{r}^{k+1}) \end{aligned}$$

Combine the best of the two worlds

The problem:

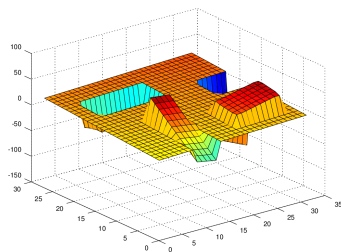
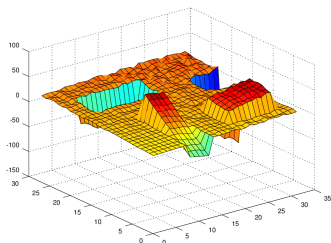
$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^n, \mathbf{x} \in \{0,1\}^{n-1}} \quad & \sum_{i=1}^n |w_i - y_i| \\ \text{s.t.} \quad & |w_{i-1} - 2w_i + w_{i+1}| \leq M(x_{i-1} + x_i), \quad i = 2, \dots, n-1, \\ & \sum_{i=1}^{n-1} x_i \leq k. \end{aligned}$$

The proposed solution [?]

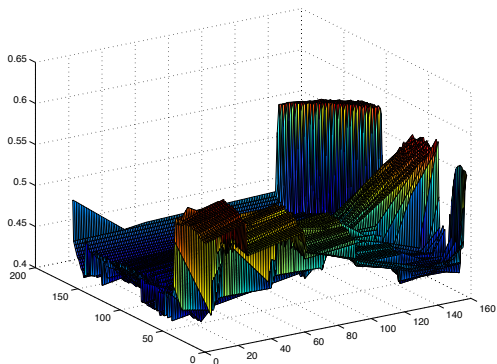
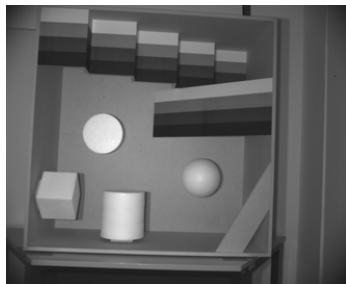
1. $(\mathbf{w}, \mathbf{x}) \leftarrow$ ADMM alternating proximal gradient method
2. use \mathbf{w} and \mathbf{x} as a warm start for MILP (with Cplex)
 - ▶ $(\mathbf{w}, \mathbf{x}) \leftarrow$ Polish coefficients on the active set
 - ▶ initialize the constant M

Road map

- 1 Examples of combinatorial problems in machine learning
- 2 Mixed integer (binary) programming (MIP)
- 3 L_0 proximal algorithm
- 4 Implementation
- 5 MIP for image processing



The problem: depth estimation based on image Z

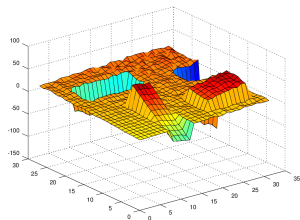


Fundamental hypothesis

piecewise linear model

ℓ_0 constraints on row and columns

$$\begin{aligned} \min_{W \in \mathbb{R}^{m \times n}} \quad & \|W - Z\|_1 \\ \text{s.t.} \quad & \|\nabla_x^2 W\|_0 \leq k_{r_i}, \quad \forall i = 1, \dots, n, \\ & \|\nabla_y^2 W\|_0 \leq k_{c_j}, \quad \forall j = 1, \dots, m. \end{aligned}$$



Fundamental hypothesis

piecewise 2d linear model

the problem as a MILP

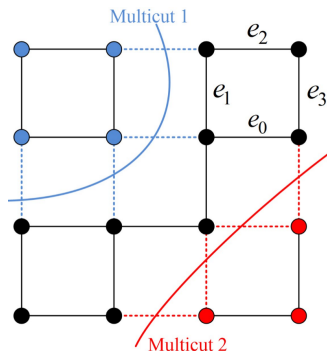
$$\begin{aligned} \min_{W \in \mathbb{R}^{m \times n}} \quad & \|W - Z\|_1 \\ \text{s.t.} \quad & \|\nabla_x^2 W\|_0 \leq k_{r_i}, \quad \forall i = 1, \dots, n, \\ & \|\nabla_y^2 W\|_0 \leq k_{c_j}, \quad \forall j = 1, \dots, m. \end{aligned}$$

$$\begin{aligned} \text{minimize} \quad & \sum_{i=1}^n \sum_{j=1}^m |W_{ij} - Z_{ij}| \\ W \in \mathbb{R}^{n \times m} & \\ x \in \{0,1\}^{n \times (m-1)} & \text{row} \\ y \in \{0,1\}^{(n-1) \times m} & \text{column} \\ \text{s.t.} \quad & |W_{i,j+1} + W_{i,j-1} - 2W_{ij}| \leq M_{r_i}(x_{i,j-1} + x_{ij}), \quad \begin{array}{l} i = 1, n \\ j = 2, m-1 \end{array} \\ & |W_{i+1,j} + W_{i-1,j} - 2W_{ij}| \leq M_{c_j}(y_{i-1,j} + y_{ij}), \quad \begin{array}{l} j = 1, m \\ i = 2, n-1 \end{array} \\ & \sum_{j=2}^{m-1} x_{ij} \leq k_{r_i}, \quad i = 1, \dots, n \\ & \sum_{i=2}^{n-1} y_{ij} \leq k_{c_j}, \quad j = 1, \dots, m \end{aligned}$$

A strongest formulation with multi-cut constraints

row cut: $x \in \{0, 1\}^{n \times (m-1)}$

column cut: $y \in \{0, 1\}^{(n-1) \times m}$



$4 \times (n - 2)^2$ additional constraints

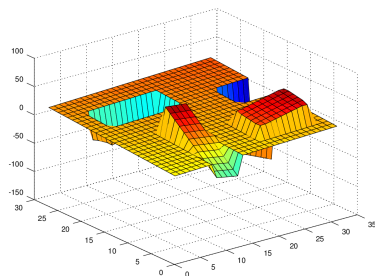
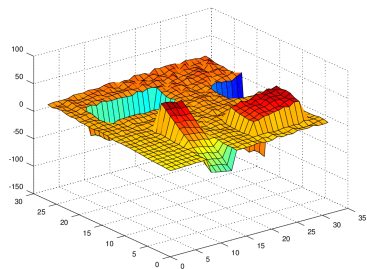
$$y_1 + x_2 + y_3 \geq x_0,$$

$$x_0 + x_2 + y_3 \geq y_1,$$

$$x_0 + y_1 + y_3 \geq x_2,$$

$$x_0 + y_1 + x_2 \geq y_3.$$

The matlab simulation



```
function [ W,X,Y,fval ] = Snd_Order_TV_MILP( Z,kr,kc,Tmax,disp )
```

```
Z          28x28  = 784
```

Reduced MIP has 5283 rows, 4541 columns, and 22807 **nonzeros**.

Reduced MIP has 1104 binaries, 0 generals, 0 SOSs, and 0 indicators.

10 sec to find the minimum
more than 10 minutes to prove it
can be reduced [?]

Road map (done)

- 1 Examples of combinatorial problems in machine learning
- 2 Mixed integer (binary) programming (MIP)
- 3 L_0 proximal algorithm
- 4 Implementation
- 5 MIP for image processing

Conclusion

- Machine learning with MIP [?]

pros	cons
it works global optimum flexible that is what we want to do	it does not scale only linear or quadratic show some instability it's not what we want to do

- Future work
 - ▶ efficient generic solver
 - ▶ efficient implementation: parallelization, randomisation, GPU
 - ▶ efficient hyper parameter calibration

