



**HAL**  
open science

# Cut Pursuit: fast algorithms to learn piecewise constant functions on general weighted graphs

Loic Landrieu, Guillaume Obozinski

► **To cite this version:**

Loic Landrieu, Guillaume Obozinski. Cut Pursuit: fast algorithms to learn piecewise constant functions on general weighted graphs. SIAM Journal on Imaging Sciences, 2017. hal-01306779v3

**HAL Id: hal-01306779**

**<https://hal.science/hal-01306779v3>**

Submitted on 31 Jul 2017 (v3), last revised 21 Aug 2017 (v4)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial 4.0 International License

# Cut pursuit: fast algorithms to learn piecewise constant functions on general weighted graphs\*

Loic Landrieu <sup>††</sup> and Guillaume Obozinski <sup>‡</sup>

**Abstract.** We propose working-set/greedy algorithms to efficiently solve problems penalized respectively by the total variation on a general weighted graph and its  $\ell_0$  counterpart the total level-set boundary size when the piecewise constant solutions have a small number of distinct level-sets; this is typically the case when the total level-set boundary size is small, which is encouraged by these two forms of penalization. Our algorithms exploit this structure by recursively splitting the level-sets of a piecewise-constant candidate solution using graph cuts. We obtain significant speed-ups over state-of-the-art algorithms for images that are well approximated with few level-sets.

**Key words.** working-set, total variation, sparsity, Mumford-Shah, greedy algorithm

**AMS subject classifications.** 90C25, 90C27, 65K10, 90C59, 68U10, 90C99, 62H11

**1. Introduction.** Estimation or approximation with piecewise constant functions has many applications in image and signal processing, machine learning and statistics. In particular, the assumption that natural images are well modeled by functions whose total variation is bounded motivates its use as a regularizer, which leads to piecewise constant images for discrete approximations. Moreover a number of models used in medical imaging [25] assume directly piecewise constant images. More generally, piecewise constant models can be used for compression, for their interpretability and finally because they are typically adaptive to the local regularity of the function approximated [69]. Piecewise constant functions display a form of structured sparsity since their gradient is sparse.

Both convex and non-convex formulations have been proposed to learn functions with sparse gradients. The most famous being the formulation of [62], hereafter referred to as ROF, which proposed to minimize the total variation subject to constraints of approximation of the noisy signal in the least squares sense, as well as the formulation of Mumford and Shah [46], which proposed to penalize the total perimeter of discontinuities of piecewise smooth functions. A fairly large literature is devoted to these formulations mainly in the fields of image processing and optimization. Although the connection between the total variation, the Mumford-Shah energy and graph cuts is today well-established, algorithms that leverage this connection are relatively recent. In particular for ROF, [13, 30] use the fact that the problem can be formulated as a parametric max-flow. [23] use graph cuts to solve the formulation of Mumford and Shah for the case of two constant components.

The literature on sparsity in computational statistics and machine learning has shown how the sparsity of the solution sought can be exploited to design algorithms which use parsimonious computations to solve the corresponding large-scale optimization problem with significant

---

\*Submitted to the editors 24/01/2017.

<sup>†</sup> Université Paris-Est, LASTIG MATIS, IGN, ENSG, F-94160 Saint-Mande, France ([loic.landrieu@ign.fr](mailto:loic.landrieu@ign.fr).)

<sup>‡</sup> Université Paris-Est, LIGM, Laboratoire d'Informatique Gaspard Monge (UMR 8049), CNRS, ENPC, ESIEE Paris, UPEM ([guillaume.obozinski@enpc.fr](mailto:guillaume.obozinski@enpc.fr)).

37 speed-ups [3]. Our work is motivated by the fact that this has to the best of our knowledge  
 38 not been fully leveraged to estimate and optimize with piecewise constant functions. In the  
 39 convex case, the algorithms proposed to exploit sparsity are working set<sup>1</sup> algorithms and the  
 40 related (fully corrective) Frank-Wolfe algorithm [31]. In the non-convex case, forward selection  
 41 algorithms such as OMP, FoBa and others have been proposed [45, 47, 70]<sup>2</sup>.

42 It is well understood that algorithms for the convex and non-convex cases are in fact fairly  
 43 related. In particular, for a given type of sparsity, the forward step of working set methods,  
 44 Frank-Wolfe and greedy algorithm is typically the same, and followed by the resolution of a  
 45 reduced problem.

46 Given their similarity, we explore in this paper both greedy and working set strategies. The  
 47 working set approach is used to solve optimization problems regularized by the total variation  
 48 while the greedy strategy solves problems penalized by the boundary size for piecewise constant  
 49 functions. In the convex case, our algorithms do not apply only to the cases in which the  
 50 data fitting term is the MSE or a separable smooth convex function, for which some efficient  
 51 algorithms implicitly exploiting sparsity exist [13, 2, 41], but also to a general smooth convex  
 52 term. Our algorithms are very competitive for deblurring and are applicable to the estimation  
 53 of piecewise constant functions on general weighted graphs.

54 **1.1. Notations.** Let  $G = (V, E, w)$  be an unoriented weighted graph whose edge set is of  
 55 cardinality  $m$  and  $V = [1, \dots, n]$ . For convenience of notations and proofs, we encode the  
 56 undirected graph  $G$ , as a directed graph with for each pair of connected nodes a directed edge  
 57 in each direction. Thus  $E$  denotes a collection of couples  $(i, j)$  of nodes, with  $(i, j) \in E$  if  
 58 and only if  $(j, i) \in E$ . We also have  $w \in \mathbb{R}^{2m}$  and  $w_{ij} = w_{ji}$ . For a set of nodes  $A \subset V$  we  
 59 denote  $\mathbf{1}_A$  the vector of  $\{0, 1\}^n$  such that  $[\mathbf{1}_A]_i = 1$  if and only if  $i \in A$ . For  $F \subset E$  a subset of  
 60 edges we denote  $w(F) = \sum_{(i,j) \in F} w_{ij}$ . By extension, for two subsets  $A$  and  $B$  of  $V$  we denote  
 61  $w(A, B) = w((A \times B) \cap E)$  the weight of the boundary between those two subsets. Finally we  
 62 denote  $\mathcal{C}$  the set of all partitions of  $V$  into connected components.

63 **1.2. General problem considered.**

64 **1.2.1. Problem formulation.** In this work we consider the problem of minimizing func-  
 65 tions  $Q$  of the form  $f(x) + \lambda\Phi(x)$  with  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  convex and differentiable, and  $\Phi : \mathbb{R}^n \rightarrow \mathbb{R}$   
 66 a penalty function that decomposes as  $\Phi(x) = \sum_{(i,j) \in E} w_{ij} \phi(x_i - x_j)$  with  $\phi : \mathbb{R} \rightarrow \mathbb{R}^+$  a  
 67 sparsity-inducing function such that  $\phi(0) = 0$ . The general problem writes  $\min_{x \in \mathbb{R}^n} Q(x)$   
 68 with

$$69 \quad (1) \quad Q(x) \doteq f(x) + \frac{\lambda}{2} \sum_{(i,j) \in E} w_{ij} \phi(x_i - x_j).$$

---

<sup>1</sup>We distinguish *working set* algorithms (aka column generation algorithm) that maintain an expansion of the solution which may have zero coefficients from *active set* algorithms that maintain an expansion using only non-zero coefficients and discard all other directions (or variables). This distinction can also be understood in the dual, where working set algorithms (which are dually cutting plane algorithms) maintain a superset of the active constraints, while active set algorithms maintain the exact set of active constraints.

<sup>2</sup>Proximal methods that perform soft-thresholding or the non-convex IHT methods maintain sparse solutions, but typically need to update a full dimensional vector at each iteration, which is why we do not cite them here. They blend however very well with active set algorithms.

70 Energies of this form were first introduced by [29] for image regularization, and are widely  
 71 used for their inducing spatial regularity as well as preserving discontinuities. In this paper,  
 72 we consider the case  $\phi$  equal to the absolute value, which corresponds to the total variation  
 73 (denoted TV), and the case  $\phi$  equal to one minus the Kronecker delta at 0, which leads to the  
 74 *total boundary size* penalty for piecewise constant functions. For these functions, the solution  
 75  $x^*$  of (1) has a sparse gradient  $\{x_i^* - x_j^* \mid (i, j) \in E\}$ . As a consequence, these solutions are  
 76 constant on the elements of a certain partition of  $V$  that is typically coarse, i.e. such that has  
 77 much fewer elements than  $|V|$ . We therefore reformulate the problem for candidate solutions  
 78 that have that property. We define the *support* of a vector  $x \in \mathbb{R}^n$  as the set  $S(x)$  of edges  
 79 supporting its gradients

$$80 \quad (2) \quad S(x) \doteq \{(i, j) \in E \mid x_i \neq x_j\},$$

81 and we will use  $S^c(x) \doteq E \setminus S(x)$  for the set on which the gradients are zero.

82 **1.2.2. Decomposition on a partition.** Any  $x \in \mathbb{R}^n$  can be written as  $x = \sum_{i=1}^k c_i \mathbf{1}_{A_i}$   
 83 with  $\Pi = \{A_1, \dots, A_k\} \in \mathcal{C}$  a partition of  $V$  into  $k$  connected components and  $c \in \mathbb{R}^k$ .  
 84 Conversely we say that  $x$  can be expressed by partition  $\Pi = (A_1, \dots, A_k)$  if it is in the set  
 85  $\text{span}(\Pi) = \text{span}(\mathbf{1}_{A_1}, \dots, \mathbf{1}_{A_k}) = \{\sum_{i=1}^k c_i \mathbf{1}_{A_i} \mid c \in \mathbb{R}^k\}$ . We denote

$$86 \quad (3) \quad x_\Pi \doteq \arg \min_{z \in \text{span}(\Pi)} Q(z),$$

87 the solution of (1) when  $x$  is constrained to be in  $\text{span}(\Pi)$ . Assuming that the regularization  
 88 strength is such that the solution  $x^*$  decomposes over a coarse partition, and that the con-  
 89 strained problem (3) is easy to solve for such a partition, problem (1) boils down to finding an  
 90 optimal partition  $\Pi^*$ :

$$91 \quad (4) \quad \Pi^* \doteq \arg \min_{\Pi \in \mathcal{C}} Q(x_\Pi).$$

92 An additional motivation to consider a sequence of partitions and solve sequentially problems  
 93 with  $x$  constrained to  $\text{span}(\Pi)$  is that the vectors of the form  $w(B, B^c)^{-1} \mathbf{1}_B$  are extreme points  
 94 of the set  $\{x \mid \text{TV}(x) \leq 1\}$ . In fact, the total variation is an *atomic gauge* in the sense of [17]  
 95 and the vectors of the form  $w(B, B^c)^{-1} \mathbf{1}_B$  are among the *atoms* of the gauge. We do not  
 96 develop this more abstract point of view in the paper, but provide a discussion in Appendix A.

97 Before presenting our approach we review some of the main relevant ideas in the related  
 98 literature.

99 **1.3. Related work.** [46] describe an image as *simple* if it can be expressed as a piecewise-  
 100 smooth function with few and small discontinuities, that is if the space can be partitioned in a  
 101 finite number of regions with short contours and such that the image varies smoothly in each of  
 102 these regions.

103 Given an observed noisy image modeled as a function  $J : \Omega \rightarrow \mathbb{R}$  whose domain  $\Omega$  is  
 104 an open, bounded and connected subset of  $\mathbb{R}^2$ , and assuming  $J \in L^\infty$ , Mumford and Shah  
 105 propose to obtain a denoised version  $I$  of the image via the minimization of an energy which  
 106 we can write as

$$107 \quad (\text{MS}) \quad \int_{\Omega} (I(x) - J(x))^2 dx + \mu \int_{\Omega \setminus \Gamma} \|\nabla I(x)\|^2 dx + \lambda \mathcal{H}_1(\Gamma),$$

108 where  $\mu$  and  $\lambda$  are two nonnegative regularization coefficients. It is composed of three terms: a  
 109 fidelity term quantifying the distortion between  $I$  and  $J$ , a term measuring the smoothness of  $I$   
 110 outside of a one-dimensional set of discontinuities  $\Gamma$ , and finally the one-dimensional Hausdorff  
 111 measure of this set  $\mathcal{H}_1(\Gamma)$ . David Mumford and Jayant Shah conjectured that this problem  
 112 admitted a solution  $(I^*, \Gamma^*)$  such that  $I^*$  was continuously differentiable on a finite number  $k$   
 113 of open sets  $R_i$  with  $\Gamma^* = \Omega \setminus \bigcup_i R_i$  a one dimensional set consisting of points connected by  
 114 rectifiable arcs.

115 In subsequent formalisations of the Mumford-Shah problem,  $I$  is constrained to the set  
 116  $\mathcal{C}^1(\Omega \setminus \Gamma)$  of continuously differentiable functions on  $\Omega \setminus \Gamma$ , where  $\Gamma$  is a closed set of Hausdorff  
 117 dimension 1. Ennio De Giorgi proposed a relaxed Mumford-Shah problem in which  $I$  is  
 118 constrained to the set  $\text{SBV}(\mathbb{R}^2)$  of special bounded total variation functions and  $\Gamma = \mathcal{S}I$  is the  
 119 *jump set* of  $I$  (for detailed presentations of the different formulations of the Mumford-Shah  
 120 problem and their connections, see [28, 5]). When  $\mu \rightarrow \infty$ , the smoothness term forces  $I$  to  
 121 be constant on the connected components of  $\Omega \setminus \Gamma$ .

122 If the number  $k$  of regions  $R_i$  (also called *phases*) on which  $I$  is constant is fixed to  $k$ , the  
 123 corresponding problem is referred to as the *piecewise constant Mumford-Shah problem* and can  
 124 be reformulated as:

$$125 \quad (\text{PC-MS}) \quad \min_{\Gamma, I} \sum_{i=1}^k \int_{R_i} (I_i - J(x))^2 dx + \lambda \mathcal{H}_1(\Gamma),$$

126 with  $I_i$  the constant value of  $I$  on  $R_i$  and  $\Omega = R_1 \cup \dots \cup R_k \cup \Gamma$ . Note that when  $k$  is fixed, the  
 127 sets  $R_i$  are not necessarily connected sets. Note that both (MS) and (PC-MS) extend naturally  
 128 to  $d$ -dimensional images by replacing  $\mathcal{H}_1$  by the  $d - 1$ -dimensional Hausdorff measure  $\mathcal{H}_{d-1}$ .

129 The setting in which  $k = 2$  is known as the Chan-Vese problem and was first approached  
 130 algorithmically using active contour methods [36, 1]. [16] propose a level-set based method for  
 131 the binary case, which has the advantage of foregoing edges and gradient completely, as they  
 132 are typically very sensitive to noise. This method has since been extended to the so called  
 133 *multiphase* setting where the number of *phases*, that is of level-sets of the function, is a power  
 134 of two [68]. The resolution of those problems is substantially sped up by the introduction of  
 135 graph-cut methods, for the binary phase [25] and in the multiphase setting [23].

136 Clearly, a counterpart of (PC-MS) in which the number of phases is not set a priori (and  
 137 can possibly be infinite) is also of interest. It has been introduced in the discrete setting by  
 138 [42] and has been studied in the continuous setting using the theory of *Caccioppoli partitions*  
 139 [66, 43].

140 Independently of the work of Mumford and Shah, [62] proposed the idea that the class of  
 141 functions with bounded variation is a good model for images, and relied on this idea to motivate  
 142 the minimization of the total variation under MSE approximation constraint as an approach  
 143 for image denoising. The introduction of the total variation had a lasting impact in imaging  
 144 sciences and was used for various tasks including denoising, deblurring and segmentation [12].  
 145 The variant<sup>3</sup> of the problem of Rudin, Osher and Fatemi (ROF) where the total variation is  
 146 used as a regularizer—corresponding to the proximal problem of the total variation—can be

---

<sup>3</sup>In [62] the TV is minimized under a constraint on the  $L_2$  distance between  $I$  and  $J$ .

147 written

148 (ROF) 
$$\min_{I \in \text{BV}} \int_{\Omega} (I(x) - J(x))^2 dx + \lambda \text{TV}(I),$$

149 where TV is the total variation and BV is the space of functions with bounded total variation.

150 In this paper we consider discretized versions of these formulations, in which the function  
 151 takes its value on the node set of a weighted graph  $G = (V, E, w)$ . Such discretizations are  
 152 for example naturally obtained if an a priori fine grained partition of the space in a collection  
 153 of elementary regions<sup>4</sup> is chosen and the image or function  $I$  is constrained to be constant  
 154 on each of these regions. The edge set  $E$  captures adjacencies between the elements, and the  
 155 weights  $w$  the size of the boundary between each pair of regions.

156 The ROF problem can be solved very efficiently for chain graphs using dynamic program-  
 157 ming [35] or exploiting the structure of the optimality conditions [19]. See [38] for a broader  
 158 discussion. In the general case, a first approach is to consider explicitly the set of edges pre-  
 159 senting discontinuities and iteratively update this set using calculus of variations based on the  
 160 Euler-Lagrange equations [1]. This class of methods is known as *active contours*. The level-  
 161 sets approach [54, 67] takes an opposite point of view and defines the discontinuity set as the  
 162 zero set of an auxiliary function. This allows for an indirect and continuous handling of the  
 163 evolution of the curve, thereby avoiding complications associated to making discrete changes  
 164 in the structure of the contours. In the recent literature, problems regularized with the total  
 165 variation are typically solved using proximal splitting algorithms [14, 57].

166 Some of the connections between graph-cuts and the total variation were already known  
 167 in [55] but some of these connections have been only fully exploited recently, when [13] and  
 168 [30] among others, exploited the fact that the ROF model can be reformulated as a para-  
 169 metric maximum flow problem, which, in these papers, is moreover shown to be solved by a  
 170 divide-and-conquer strategy: This algorithm entails to solve a sequence of max-flow problems  
 171 on the same graph, and the algorithm makes it possible to efficiently reuse partial computa-  
 172 tions performed in each max-flow problem with a push-relabel algorithm. These results on  
 173 the total variation are actually an instance of results that apply more generally to submodular  
 174 functions [2]. Indeed, the intimate relation existing between the total variation and graph-cuts  
 175 is due fundamentally to the fact that the former is the Lovász extension of the value of the  
 176 cut, which is a submodular function. Beyond the case of the total variation, [4] considers regu-  
 177 larizers that are obtained as Lovász extensions of symmetric submodular functions and recent  
 178 progress made on the efficient optimization of submodular functions produces simultaneously  
 179 new fast algorithms to compute proximal operators of the Lovász extension of submodular  
 180 function [41, 34].

181 In the discrete setting, problems regularized by the total variation or the total boundary  
 182 size are also related to the Potts model. Indeed, if the values of the level-set are quantized, the  
 183 corresponding energy to minimize is that of a discrete valued conditional random field (CRF),  
 184 with as many values as there are quantization levels [32, 67]. A number of optimization  
 185 techniques exist for CRFs [65]. One of the fastest is the  $\alpha$ -expansion algorithm of [10], which  
 186 relies on graph-cut algorithms [9].

---

<sup>4</sup>In the context of images these could be thought of as super-pixels, for example.

187 In the literature on sparsity, a number of algorithms have been proposed to take advantage  
 188 computationally of the sparsity of the solution. In the convex setting, these algorithms include  
 189 homotopy algorithms such as the LARS [21] or working set algorithms [52, 61, 26]. It should  
 190 be noted that the Frank-Wolfe algorithm [33], which has been revived and regained popular-  
 191 ity in recent years, is closely-related to working set methods and also provides a rationale to  
 192 algorithmically exploit the sparsity of solution of optimization problems. Although originally  
 193 designed to solve constrained optimization problems, [31] have shown how a variant can be nat-  
 194 urally constructed for the regularized setting, and can be applied to the case of total variation  
 195 regularization. The counterparts of these algorithms in the  $\ell_0$  setting are (a) greedy forward  
 196 selection approaches that compute a sequence of candidate solutions by iteratively decreasing  
 197 the sparsity of the candidate solutions, such as orthogonal matching pursuit [45], orthogonal  
 198 least squares [18] and related algorithms [47], (b) forward-backward selection approaches such  
 199 as the Single Best Replacement (SBR) algorithm [64], based on an  $\ell_0$  penalization or the FoBa  
 200 algorithm [70], which add backwards steps to remove previously introduced variables that are  
 201 no longer relevant. See [3] for a review. [2] proposes a number of algorithms to minimize sub-  
 202 modular functions, compute the associated proximal operators of the corresponding Lovász  
 203 extensions. In particular, generic primal and dual active set algorithms are proposed to solve  
 204 a linear regression problem regularized with the Lovász extension of a submodular function [2,  
 205 Chap. 7.12].

206 **2. A working set algorithm for total variation regularization.** In this section, we consider  
 207 the problem of solving the minimization of a convex, differentiable function  $f$  regularized  
 208 by a weighted total variation of the form  $\text{TV}(x) = \frac{1}{2} \sum_{(i,j) \in E} w_{ij} |x_i - x_j|$  with  $w_{ij}$  some  
 209 nonnegative weights<sup>5</sup>.

210 Based on the considerations of Section 1.2.2, we propose a working set algorithm which  
 211 alternates between solving a reduced problem of the form  $\min_{x \in \text{span}(\Pi)} Q(x)$  for  $Q(x) = f(x) +$   
 212  $\lambda \text{TV}(x)$ , and refining the partition  $\Pi$ . In Section 2.3, we will discuss how to solve the reduced  
 213 problem efficiently, but first we present a criterion for refining the partition  $\Pi$ .

214 **2.1. Steepest binary cut.** Given a current partition  $\Pi$  and the solution of the associated  
 215 reduced problem  $x_\Pi = \arg \min_{x \in \text{span}(\Pi)} Q(x)$ , our goal is to compute a finer partition  $\Pi_{new}$   
 216 leading to the largest possible decrease of  $Q$ . To this end we consider updates of  $x$  of the  
 217 form  $x_\Pi + h u_B$  with  $u_B = \gamma_B \mathbf{1}_B - \gamma_{B^c} \mathbf{1}_{B^c}$  for some set  $B \subset V$  and some scalars  $h, \gamma_B$  and  
 218  $\gamma_{B^c}$  such that  $\|u_B\|_2 = 1$ . We postpone to Section 2.2 the precise discussion of how the choice  
 219 of  $B$  leads to a new partition and focus first on a rationale for choosing  $B$ , but essentially,  
 220 introducing  $u_B$  in the expansion of  $x$  will lead to a new partition in which the elements of  $\Pi$   
 221 are split along the boundary between  $B$  and  $B^c$ . A natural criterion is to choose the set  $B$   
 222 such that  $u_B$  is a descent direction which is as steep as possible, in the sense that  $Q$  decreases  
 223 the most, at first order. We denote  $Q'(x, v) = \lim_{h \rightarrow 0} h^{-1}(Q(x + hv) - Q(x))$  so that, when  
 224  $d \in \mathbb{R}^n$  is a unit vector,  $Q'(x, d)$  denotes the directional derivative of  $Q$  at  $x \in \mathbb{R}^n$  in the  
 225 direction  $d$ . Consequently, choosing  $B$  for which the direction  $u_B$  is steepest requires solving  
 226  $\min_{B \subset V} Q'(x_\Pi, u_B)$ .

---

<sup>5</sup>In particular, this is the form taken by the anisotropic total variation for images if the weights are deter-  
 mined by the Cauchy-Crofton formula (see e.g. [30]).

To further characterize  $Q'$  we decompose the objective function: Since the absolute value is differentiable on  $\mathbb{R}_*$ , setting  $S \doteq S(x_\Pi)$  allows us to split  $Q$  into two parts  $Q_S$  and  $\text{TV}|_{S^c}$  which are respectively differentiable and non-differentiable at  $x_\Pi$ :

$$\begin{cases} Q_S(x) & \doteq f(x) + \frac{\lambda}{2} \sum_{(i,j) \in S} w_{ij} |x_i - x_j|, \\ \text{TV}|_{S^c}(x) & \doteq \frac{\lambda}{2} \sum_{(i,j) \in S^c} w_{ij} |x_i - x_j|. \end{cases}$$

227  $\text{TV}|_{S^c}$  is a weighted total variation on the graph  $G$  but with weights  $w_{S^c}$  such that  $[w_{S^c}]_{i,j} \doteq w_{ij}$   
 228 for  $(i,j) \in S^c$  and 0 for  $(i,j) \in S$ . We extend the previous notations and define  $w_{S^c}(A, B) \doteq$   
 229  $w_{S^c}(A \times B) = w((A \times B) \cap S^c)$ .

**Proposition 1.** For  $x \in \mathbb{R}^n$ , if we set  $S = S(x)$  then the directional derivative in the direction of  $\mathbf{1}_B$  is

$$Q'(x, \mathbf{1}_B) = \langle \nabla Q_S(x), \mathbf{1}_B \rangle + \lambda w_{S^c}(B, B^c).$$

Moreover if  $\langle \nabla f(x), \mathbf{1}_V \rangle = 0$  then

$$Q'(x, u_B) = (\gamma_B + \gamma_{B^c}) Q'(x, \mathbf{1}_B).$$

230 *Proof.* See Appendix B. ■

Considering the case  $x = x_\Pi$ , then for  $S = S(x_\Pi)$ ,  $\nabla f(x_\Pi)$  is clearly orthogonal to  $\text{span}(\Pi)$  and thus to  $\mathbf{1}_V$ . Therefore, by the previous proposition, finding the steepest descent direction of the form  $u_B$  requires solving

$$\min_{B \subset V} (\gamma_B + \gamma_{B^c}) Q'(x_\Pi, \mathbf{1}_B)$$

231 To keep a formulation which remains amenable to efficient computations, we will ignore the  
 232 factor<sup>6</sup>  $\gamma_B + \gamma_{B^c}$ . This leads us to define a *steepest binary cut* as any cut  $(B_\Pi, B_\Pi^c)$  such that

$$233 \quad (5) \quad B_\Pi \in \arg \min_{B \subset V} \langle \nabla Q_S(x_\Pi), \mathbf{1}_B \rangle + \lambda w_{S^c}(B, B^c).$$

234 Note that since  $Q'(x, \mathbf{1}_\emptyset) = 0$ , we have  $\min_{B \subset V} Q'(x, \mathbf{1}_B) \leq 0$ . If  $\emptyset$  is a solution to (5), we  
 235 set  $B_\Pi = \emptyset$ . As formulated, it is well-known, at least since [55], that problem (5) can be  
 236 interpreted as a minimum cut problem in a suitably defined flow graph.

237 Indeed consider the graph  $G_{flow} = (V \cup \{s, t\}, E_{flow})$  illustrated in Figure 1, where  $s$  and  
 238  $t$  are respectively a source and sink nodes, and where the edge set  $E_{flow}$  and the associated  
 239 nonzero (undirected) capacities  $c \in \mathbb{R}^{|S^c|+n}$  are defined as follows

$$240 \quad (6) \quad E_{flow} = \begin{cases} (s, i), \forall i \in \nabla_+, & \text{with } c_{si} = \nabla_i Q_S(x), \\ (i, t), \forall i \in \nabla_-, & \text{with } c_{it} = -\nabla_i Q_S(x), \\ (i, j), \forall (i, j) \in S^c, & \text{with } c_{ij} = \lambda w_{ij}, \end{cases}$$

---

<sup>6</sup> $\gamma_B$  and  $\gamma_{B^c}$  could otherwise be determined by requiring that  $\langle \mathbf{1}_V, u_B \rangle = 0$ . More rigorously, descent directions considered could be required to be orthogonal to  $\text{span}(\Pi)$ , but this leads to even less tractable formulations, that we therefore do not consider here.

241 where  $\nabla_+ \doteq \{i \in V \mid \nabla_i Q_S(x) > 0\}$  and  $\nabla_- \doteq V \setminus \nabla_+$ . The vector  $\nabla Q_S(x)$  is directly computed  
 242 as  $\nabla Q_S(x) = \nabla f(x) + \frac{1}{2} \lambda D^\top y$ , with  $D \in \mathbb{R}^{2m \times n}$  the weighted edge incidence matrix whose  
 243 entries are equal to  $D_{(i,j),k} \doteq w_{ij}(1_{\{i=k\}} - 1_{\{j=k\}})$  and  $y \in \mathbb{R}^{2m}$  is the vector whose entries  
 244 are indexed by the elements of  $E$  and such that  $y_{(i,j)} \doteq \text{sign}(x_i - x_j)$  with the convention that  $\text{sign}(0) = 0$ . As stated in the next proposition, finding a minimal cut in this graph provides

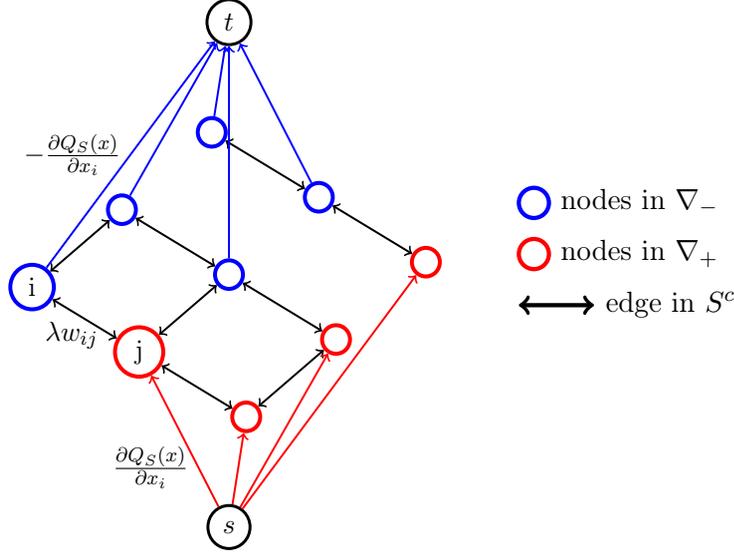


Figure 1: Directed graph for which finding a maximal flow is equivalent to solving (5). Neighboring nodes with different values of  $x$  in the original graph are linked by an undirected edge with capacity  $\lambda w_{ij}$ , nodes with non-negative gradient are linked to the source, and nodes with negative gradient to the sink with capacity  $|\nabla Q_S(x)|$ .

245 us with the desired steepest binary cut.  
 246

247 **Proposition 2.** *Let  $S = S(x)$  then  $(C, V_{flow} \setminus C)$  is a minimal cut in  $G_{flow}$  if and only if*  
 248  *$C \setminus \{s\}$ , and its complement in  $V$  are minimizers of  $B \mapsto Q'(x, \mathbf{1}_B)$ .*

249 This result is a well-know result which was first discussed in [55]. We refer the reader to [39]  
 250 for a proof.

Note that the min-cut/max-flow problem of Figure 1 decouples on each of the connected components of the graph  $G|_{S^c} \doteq (V, S^c)$  and that as a result solving (5) is equivalent to solving separately

$$\min_{C \subset A} \langle \nabla Q_S(x_\Pi), \mathbf{1}_C \rangle + \lambda w(C, A \setminus C)$$

251 for each set  $A$  that is a connected components of  $G|_{S^c}$ . The binary steepest cut thus actually  
 252 reduces to computing a steep cut in each connected component of the graph, and they can all  
 253 be computed in parallel. Let us insist that the connected components of  $G|_{S^c}$  are often but  
 254 not always the elements of  $\Pi$  since they can be unions of adjacent elements of  $\Pi$  when they  
 255 share the same value.

256 We can now characterize the optimality of  $x_\Pi$  or of the corresponding partition  $\Pi$ , based  
 257 on the value of the steepest binary partition:

258 **Proposition 3.** *We have  $x = \arg \min_{z \in \mathbb{R}^n} Q(z)$  if and only if  $\min_{B \subset V} Q'(x, \mathbf{1}_B) = 0$  and*  
 259  *$Q'(x, \mathbf{1}_V) = 0$ .*

260 *Proof.* See Appendix B ■

261 Note that the rationale we propose to choose the new direction  $\mathbf{1}_B$  is different than the one  
 262 typically used for working-set algorithms in the sparsity literature and variants of Frank-Wolfe.  
 263 When considering the minimization of an objective of the form  $f(x) + \lambda \Omega(x)$ , where  $f : \mathbb{R}^n \rightarrow \mathbb{R}$   
 264 is a differentiable function and  $\Omega$  is a norm, the optimality condition in terms of subgradient  
 265 is  $-\frac{1}{\lambda} \nabla f(x) \in \partial \Omega(x)$ , where  $\partial \Omega(x)$  is the subgradient of the norm  $\Omega$  at  $x$ . A classical result  
 266 from convex analysis is that  $\partial \Omega(x) = \{s \in \mathbb{R}^n \mid \langle s, x \rangle = \Omega(x) \text{ and } \Omega^\circ(s) \leq 1\}$  where  $\Omega^\circ$   
 267 denotes the dual norm [60, Thm. 23.5]. In particular, the subgradient condition is not satisfied  
 268 if  $\Omega^\circ(-\nabla f(x)) \geq \lambda$  and since  $\Omega^\circ(s) = \max_{\Omega(\xi) \leq 1} \langle s, \xi \rangle$  then  $\arg \max_{\Omega(\xi) \leq 1} \langle -\nabla f(x), \xi \rangle$  provides  
 269 a direction in which the inequality constraint is most violated. This direction is the same  
 270 as the Frank-Wolfe direction for the optimization problem  $\min_{x: \Omega(x) \leq \kappa} f(x)$ , also the same  
 271 as the direction proposed in a variant of the Frank-Wolfe algorithm proposed by [31] for the  
 272 regularized problem, and again the same as the direction that would be used in the primal  
 273 active set algorithm of [2, Chap. 7.12] for generic Lovász extensions of submodular function,  
 274 which is essentially a fully corrective and active-set version of the algorithm of [31]. This  
 275 rationale extends to the case where  $\Omega$  is more generally a gauge and is most relevant when it  
 276 is an atomic norm or gauge [17], which we discuss in Appendix A. For decomposable atomic  
 277 norms [48] that have atoms of equal Euclidean norm, one can check that the steepest descent  
 278 direction that we propose and the Frank-Wolfe direction are actually the same. However,  
 279 for the total variation the two differ. The Frank-Wolfe direction leads to the choice  $B^* =$   
 280  $\arg \max_{B \subset V} -w(B, B^c)^{-1} \langle \nabla f(x_\Pi), \mathbf{1}_B \rangle$ . We show in Section 2.7 and via results presented in  
 281 Figure 6 that using the steepest cut direction outperforms the Frank-Wolfe direction.

282 **2.2. Induced new partition in connected sets and new reduced problem.** For  $\Pi =$   
 283  $(A_1, \dots, A_k)$ ,  $B_\Pi$  is chosen so that the addition of a term of the form  $h u_B = h \gamma_B \mathbf{1}_B - h \gamma_{B^c} \mathbf{1}_{B^c}$   
 284 to  $x = \sum_{i=1}^k c_i \mathbf{1}_{A_i}$  decreases the objective function  $Q$  the most. At the next iteration, we could  
 285 thus consider solving a reduced problem that consists of minimizing  $Q$  under the constraint  
 286 that  $x \in \text{span}(\mathbf{1}_{A_1}, \dots, \mathbf{1}_{A_k}, \mathbf{1}_B)$  with  $B = B_\Pi$ . But there is in fact a simpler and more  
 287 relevant choice. Indeed, on the set  $\text{span}(\mathbf{1}_{A_1}, \dots, \mathbf{1}_{A_k}, \mathbf{1}_B)$ , the values  $x_{i_1}, x_{i_2}, x_{i_3}$  and  $x_{i_4}$  with  
 288  $i_1 \in A_j \cap B$ ,  $i_2 \in A_j \cap B^c$ ,  $i_3 \in A_{j'} \cap B$  and  $i_4 \in A_{j'} \cap B^c$  are a priori coupled; also, if  $A_j \cap B$  has  
 289 several connected components  $i \mapsto x_i$  must take the same value on these components. These  
 290 constraints seem unnecessarily restrictive.

291 Consider  $S_\Pi \doteq \bigcup_{(A, A') \in \Pi^2} \partial(A, A')$  with  $\partial(A, A') \doteq (A \times A') \cap E$ . With the notion of  
 292 support  $S(x)$  that we defined in (2) we actually have  $\text{span}(\Pi) = \{x \in \mathbb{R}^n \mid S(x) \subset S_\Pi\}$ .  
 293 Now, if  $x \in \text{span}(\mathbf{1}_{A_1}, \dots, \mathbf{1}_{A_k}, \mathbf{1}_B)$ , we have in general  $S(x) \subset S_{\text{new}} \doteq S_\Pi \cup \partial(B, B^c)$ , which  
 294 corresponds to allowing a larger support. But then it makes sense to allow  $x$  to remain in the  
 295 largest set with this maximal support  $S_{\text{new}}$ , that is equivalent to staying in the vector space  
 296  $\mathcal{X}_{S_{\text{new}}} \doteq \{x' \mid S(x') \subset S_{\text{new}}\}$ . But, if we now define  $\Pi_{\text{new}}$  as the partition of  $V$  defined as  
 297 the collection of all connected components in  $G$  of all sets  $A_j \cap B_\Pi$  and  $A_j \cap B_\Pi^c$  for  $A_j \in \Pi$ ,

298 then it is relatively immediate that  $\text{span}(\Pi_{\text{new}}) = \mathcal{X}_{S_{\text{new}}}$ . The construction of  $\Pi_{\text{new}}$  from  $\Pi$  is  
 299 illustrated in Figure 2.

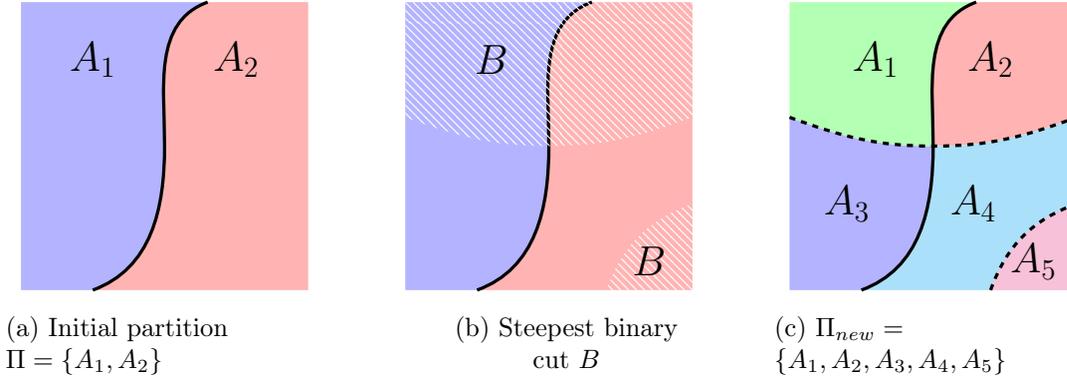


Figure 2: Illustration of the induced new partition. From an initial partition  $\Pi$ , the steepest binary cut  $B$  induced a new partition  $\Pi_{\text{new}}$ . The solid line  $\text{—}$  represent the initial contours  $S$ , and the dashed line  $\text{- - - -}$  the new contours  $S_{\text{new}} \setminus S$  introduced by  $B$ . Note that the binary partition induced by  $B$  can more than double the number of resulting components.

299

300 We therefore set  $\Pi_{\text{new}}$  to be the new partition and solve the reduced problem constrained  
 301 to  $\text{span}(\Pi_{\text{new}})$ . Note that in general we do not have  $S(x_{\Pi}) = S_{\Pi}$ , because the total variation  
 302 regularization can induce that the value of  $x_{\Pi}$  on several adjacent elements of  $\Pi$  is the same.  
 303 The following result shows that if a non-trivial cut  $(B_{\Pi}, B_{\Pi}^c)$  was obtained as a solution to  
 304 (5) then the new reduced problem has a solution  $x_{\Pi_{\text{new}}} = \arg \min_{x \in \text{span}(\Pi_{\text{new}})} Q(x)$  which is  
 305 strictly better than the previous one.

306 **Proposition 4.** *If  $B_{\Pi} \neq \emptyset$ ,  $Q(x_{\Pi_{\text{new}}}) < Q(x_{\Pi})$ .*

*Proof.* We clearly have

$$\text{span}(\Pi) \subset \text{span}(\mathbf{1}_{A_1}, \dots, \mathbf{1}_{A_k}, \mathbf{1}_{B_{\Pi}}) \subset \text{span}(\Pi_{\text{new}}),$$

so that

$$Q(x_{\Pi_{\text{new}}}) = \min_{x \in \text{span}(\Pi_{\text{new}})} Q(x) \leq \min_{x \in \text{span}(\Pi)} Q(x) = Q(x_{\Pi}).$$

307 Moreover, if  $B_{\Pi} \neq \emptyset$ , then  $Q'(x_{\Pi}, \mathbf{1}_B) < 0$ , which entails that there exists  $\varepsilon > 0$  such that  
 308  $Q(x_{\Pi_{\text{new}}}) \leq Q(x_{\Pi} + \varepsilon \mathbf{1}_B) < Q(x_{\Pi})$ . This completes the proof.  $\blacksquare$

309 We summarize the obtained working set scheme as Algorithm 1, and illustrate its two first  
 310 steps on a ROF problem in Figure 3. The following proposition provides a formal proof of  
 311 convergence.

312 **Proposition 5.** *The scheme presented in Algorithm 1 converges to the a global minimum  $x^*$   
 313 of  $Q$  in a finite a finite amount of steps bounded by  $n$ .*

**Algorithm 1** Cut pursuit

---

```

Initialize  $\Pi \leftarrow \{V\}$ ,  $x_\Pi \in \arg \min_{c \in \mathbb{R}} Q(c\mathbf{1}_V)$ ,  $S \leftarrow \emptyset$ 
while  $\min_{B \subset V} \langle \nabla Q_S(x_\Pi), \mathbf{1}_B \rangle + \lambda w_{S^c}(B, B^c) < 0$  do
    Pick  $B_\Pi \in \arg \min_{B \subset V} \langle \nabla Q_S(x_\Pi), \mathbf{1}_B \rangle + \lambda w_{S^c}(B, B^c)$ 
     $\Pi \leftarrow \{B_\Pi \cap A\}_{A \in \Pi} \cup \{B_\Pi^c \cap A\}_{A \in \Pi}$ 
     $\Pi \leftarrow$  connected components of elements of  $\Pi$ 
    Pick  $x_\Pi \in \arg \min_{z \in \text{span}(\Pi)} Q(z)$ 
     $S \leftarrow S(x_\Pi)$ 
end while
return  $(\Pi, x_\Pi)$ 

```

---

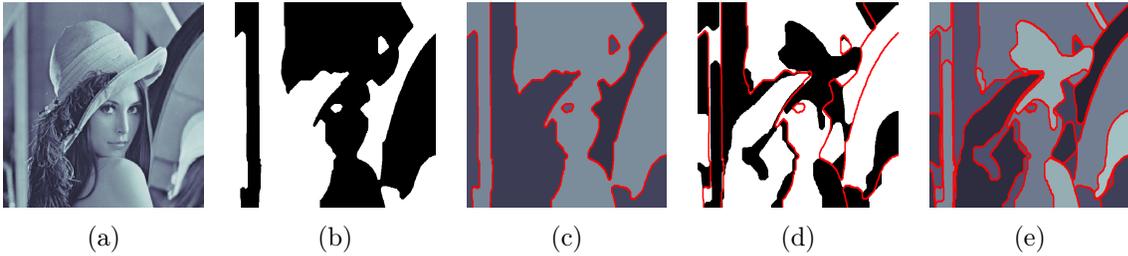


Figure 3: Two first iterations of cut pursuit for the ROF problem on the picture in (a). Images (b) and (d) represent the new cut at iterations 1 and 2 with  $B_\Pi$  and  $B_\Pi^c$  respectively in black and white, and (c) and (e) represent the partial solution in levels of gray, with the current set of contours  $S$  in red. The contours induced by the cut in (b) (resp. (d)) are superimposed on (c) (resp. (e)).

314 *Proof.* At the beginning of each iteration, if  $\min_{B \subset V} Q'(x_\Pi, \mathbf{1}_B) < 0$  then the steepest  
315 binary partition is not trivial, that is  $B_\Pi \neq \emptyset$ . Consequently the new partition  $\Pi_{\text{new}}$  will have  
316 at least one more component than  $\Pi$ , and Proposition 4 states that the solution associated  
317 with  $\Pi_{\text{new}}$  will be strictly better than  $x_\Pi$ . This ensures that the objective function is strictly  
318 decreasing along iterations of the algorithm. If  $\min_{B \subset V} Q'(x_\Pi, \mathbf{1}_B) = 0$ , then Proposition 3  
319 ensures that optimality is reached, because for each value of  $\Pi$ , by construction  $x_\Pi$  is such that  
320  $Q'(x_\Pi, \mathbf{1}_V) = 0$ . Since the number of components of  $\Pi$  is strictly increasing and bounded by  
321  $n$ , the algorithm converges in at most  $n$  steps, in the worst case scenario. Provided that each  
322 constrained problem  $x_\Pi \in \arg \min_{z \in \text{span}(\Pi)} Q(z)$  is solved exactly in finite time, this proves  
323 that  $x_\Pi$  converges to the optimum  $x^*$ . In the next section we discuss how to exploit the  
324 sparse structure of  $x_\Pi$  to solve the reduced problem efficiently. ■

325 **Case of a non-convex function  $f$ .** We assumed in all this section that  $f$  is a differ-  
326 entiable convex function. However, from a theoretical point of view, a number of results still  
327 hold even if  $f$  is non-convex provided it is assumed *strictly differentiable* in the sense of [8,  
328 Chapter 6.2], or more simply if  $f$  is assumed to be continuously differentiable, since continuous  
329 differentiability implies *strict differentiability*. Indeed, it can be shown in that case that the

330 calculations on subgradient and directional derivative that prove our results are still valid for  
 331 such a function  $f$  for an appropriate generalization of the subgradient. As discussed in more  
 332 details in Appendix C, Propositions 1 and 2 then still hold. In the non-convex case, Algo-  
 333 rithm 1 has to be modified since it is no longer reasonable to assume that a global optimum  
 334 can be found when solving the reduced problem (3), and we could assume instead that the  
 335 solver called on the reduced problem finds a local optimum which strictly reduces the value of  
 336 the objective. In the previous sections, proofs of Proposition 3 and Proposition 5 essentially  
 337 showed that some first order subgradient conditions hold and relied on the fact that first order  
 338 subgradient conditions are sufficient to characterize minima of convex functions. For non-convex  
 339 functions, the same first order subgradient conditions still hold (although they are no longer  
 340 sufficient to characterize global minima) and these proposition can be extended, but new suf-  
 341 ficient conditions are needed to guarantee that the algorithm converge to a local minimum of  
 342 the objective (see the appendix for details).

343 From a practical point of view, we however do not recommend to use the algorithm for  
 344 non-convex functions, because the low dimensional constraints of the active set algorithm could  
 345 lead to find very suboptimal local minima of the function. Instead, we would recommend when  
 346 possible to use majorization-minimization (MM) algorithms, based on convex upper bounds  
 347 of  $f$ . For instance, it is of interest to be able to solve problem (1) for non-convex functions  
 348  $\phi$  and in particular so-called concave penalties such as MCP, SCAD and others; for these  
 349 formulations, MM schemes requiring to solve a sequence of TV are efficient ([53]) and can be  
 350 advantageously combined with cut pursuit, since the latter will leverage the partition of the  
 351 previous iterate as a warm-start for the next iteration. This is the scheme we use in Section 3.2.

352 **2.3. A reduced graph for the reduced problem.** Let  $\Pi$  be a coarse partition of  $V$  into  
 353 connected components. We argue that  $\min_{z \in \text{span}(\Pi)} Q(z)$  can be solved efficiently on a smaller  
 354 weighted graph whose nodes are associated with the elements of partition  $\Pi$ , and whose edges  
 355 correspond to pairs of adjacent elements in the original graph. Indeed, consider the graph  
 356  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  with  $\mathcal{V} = \Pi$  and  $\mathcal{E} = \{(A, B) \in \mathcal{V}^2 \mid \exists(i, j) \in (A \times B) \cap E\}$ . Figure 4 shows an  
 357 example of graph reduction on a small graph. For  $x \in \text{span}(\Pi)$  we can indeed express  $\text{TV}(x)$   
 358 simply:

**Proposition 6.** For  $x = \sum_{A \in \Pi} c_A \mathbf{1}_A$  we have  $\text{TV}(x) = \text{TV}_{\mathcal{G}}(c)$  with

$$\text{TV}_{\mathcal{G}}(c) \doteq \frac{1}{2} \sum_{(A, B) \in \mathcal{E}} w(A, B) |c_A - c_B|.$$

*Proof.*

$$\begin{aligned} 359 \quad 2\text{TV}(x) &= \sum_{(i, j) \in E} w_{ij} |x_i - x_j| = \sum_{(i, j) \in E} w_{ij} \sum_{(A, B) \in \Pi^2} \mathbf{1}_{\{i \in A, j \in B\}} |c_A - c_B| \\ 360 \quad &= \sum_{(A, B) \in \Pi^2} |c_A - c_B| \sum_{(i, j) \in E \cap (A \times B)} w_{ij}, \end{aligned}$$

361 hence the result using the definition of  $w(A, B)$ . ■

362 Note that if  $\text{TV}$  is the total variation associated with the weighted graph  $G$  with weights  
 363  $(w_{ij})_{(i, j) \in E}$  then  $\text{TV}_{\mathcal{G}}$  is the total variation associated with the weighted graph  $\mathcal{G}$  and the

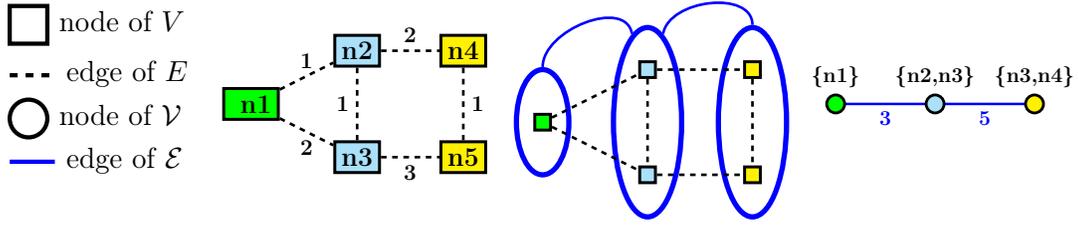


Figure 4: Example of reduced graph. Left: graph  $G$  with weights  $(w_{ij})_{(i,j) \in E}$  on the edges, middle: partition  $\Pi$  of  $G$  into connected components, right: reduced graph  $\mathcal{G}$  with weights  $(w_{AB})_{(A,B) \in \mathcal{E}}$  on the edges.

364 weights  $(w(A, B))_{(A,B) \in \mathcal{E}}$ . Denoting  $\tilde{f} : c \mapsto f(\sum_{A \in \Pi} c_A \mathbf{1}_A)$ , the reduced problem is equivalent  
 365 to solving  $\min_{c \in \mathbb{R}^k} \tilde{f}(c) + \lambda \text{TV}_{\mathcal{G}}(c)$  on  $\mathcal{G}$ . If  $\Pi$  is a coarse partition, we have  $|\mathcal{E}| \ll 2m$  and  
 366 computations involving  $\text{TV}_{\mathcal{G}}$  are much cheaper than those involving  $\text{TV}$ . As illustrated in  
 367 Section 2.4, the structure of  $\tilde{f}$  can often be exploited as well to reduce the computational cost  
 368 on the reduced problem. The construction of the reduced graph itself  $\mathcal{G}$  is cheap compared  
 369 to the speed-ups allowed, as it is obtained by computing the connected components of the  
 370 graph  $(V, E \setminus S(x))$ , which can be done in linear time by depth-first search. Note that once the  
 371 reduced problem is solved, if  $c_{\Pi} \in \arg \min_c \tilde{f}(c) + \lambda \text{TV}_{\mathcal{G}}(c)$ , then  $S(x_{\Pi})$  is directly computed  
 372 as  $S(x_{\Pi}) = \bigcup \{ \partial(A, A') \mid (A, A') \in \mathcal{E}, c_A \neq c_{A'} \}$ .

373 **2.4. Solving linear inverse problems with TV.** A number of classical problems in image  
 374 processing such as deblurring, blind deconvolution, and inpainting are formulated as ill-posed  
 375 linear inverse problems [15], where a low TV prior on the image provides appropriate regular-  
 376 ization. Typically if  $x_0 \in \mathbb{R}^n$  is the original signal,  $H$  a  $p \times n$  linear operator,  $\epsilon$  additive noise,  
 377 and  $y = Hx_0 + \epsilon \in \mathbb{R}^p$  the degraded observed signal, this leads to problems of the form

378 (7) 
$$x^* = \arg \min_{x \in \mathbb{R}^n} \frac{1}{2} \|Hx - y\|^2 + \lambda \text{TV}(x)$$

379 First order optimization algorithms, such as proximal methods, only require the computation  
 380 of the gradient  $H^T Hx - H^T y$  of  $f$  and can be used to solve (7) efficiently. However the reduced  
 381 problem can be computed orders of magnitude faster provided that the current partition is  
 382 coarse. Indeed, for a  $k$ -partition  $\Pi$  of  $V$ , we denote  $K \in \{0, 1\}^{n \times k}$  the matrix whose columns  
 383 are the vectors  $\mathbf{1}_A$  for  $A \in \Pi$ . Any  $x \in \text{span}(\Pi)$  can be rewritten as  $Kc$  with  $c \in \mathbb{R}^k$  and  
 384 the gradient of the discrepancy function with respect to  $c$  then writes:  $\nabla_c 1/2 \|HKc - y\|^2 =$   
 385  $K^T H^T HKc - K^T H^T y$ .

386 As a result, the reduced problem can be solved by a similar first-order scheme of much  
 387 smaller size, with parameters  $K^T H^T HK$  and  $K^T H^T y$ , which are of size  $k \times k$  and  $k$  respectively.  
 388 Given the sparsity of the matrix  $K$ ,  $HK$  is computed in time  $O(pn)$ ; consequently  $K^T H^T HK$   
 389 can be precomputed in  $O(k^2 p + pn)$  and  $K^T H^T y$  in  $O(pn)$ . Solving the reduced problem is  
 390 then very quick provided  $k$  is small compared to  $n$ .

391 In the case of a blur operator  $H$  with adequate symmetry, for which  $p = n$  is large, manipulating  
 392 the matrices  $H$  or  $H^T$  directly should be avoided. However  $x \mapsto Hx$  being a convolution, it

393 can be computed quickly using the fast Fourier transform and, in that case,  $K^\top H^\top H K$  and  
 394  $K^\top H y$  can be precomputed in  $\mathcal{O}(n \log n)$  time.

395 **2.5. Complexity analysis.** The computational bottlenecks of the algorithm could a priori  
 396 be (a) the computation of the steepest binary cut which requires to solve a min cut/max flow  
 397 problem, (b) the cost of solving the reduced problem, (c) the computation of the reduced graph  
 398 itself, (d) the number of global iterations needed.

399 (a) The steepest binary cut is obtained as the solution of a max-flow/min-cut optimization  
 400 problem. It is well-known that there is a large discrepancy between the theoretical  
 401 upper bound on the complexity of many graph-cut algorithms and the running times  
 402 observed empirically, the former being too pessimistic. In particular, the algorithm of  
 403 [10] has a theoretical exponential worst case complexity, but scales essentially linearly  
 404 with respect to the graph size in practice. In fact, it is known to scale better than  
 405 some algorithms with polynomial complexity, which is why we chose it.

406 (b) Solving the reduced problem can be done with efficient proximal splitting algorithms  
 407 such as [58], which is proved to reach a primal suboptimality gap of  $\varepsilon$  in  $\mathcal{O}(1/\varepsilon^2)$  it-  
 408 erations; in practice, the observed convergence rate is almost linear. Preconditioning  
 409 greatly speeds up convergence in practice. Moreover, the problems induced on the  
 410 reduced graph can typically be solved at a significantly reduced cost: in particular,  
 411 as discussed in section 2.4, for a quadratic data fitting term and  $H$  a blurr operator,  
 412 the gradient in the subgraph can be computed in  $\mathcal{O}(k^2)$  time, based on a single ef-  
 413 ficient FFT-based computation of the Hessian per global iteration which itself takes  
 414  $\mathcal{O}(k^2 n \log n)$  time. For problems with coarse solutions, this algorithm is only called for  
 415 small graphs so that this step only contributes to a small fraction of the the running  
 416 time.

417 (c) Computing the reduced graph, requires computing the connected components of the  
 418 graph obtained when removing the edges in  $S$ , and the weights  $w(A, B)$  between all  
 419 paris of components  $(A, B)$ . This can be efficiently performed in  $\mathcal{O}(m + n)$  through a  
 420 depth-first exploration of the nodes of the original graph.

421 (d) The main factor determining the computation time is the number of global iterations  
 422 needed. In the worst case scenario, this is  $\mathcal{O}(n)$ . In practice, the number of global  
 423 iterations seems to grow logarithmically with the number of constant regions at the  
 424 optimum. If for simple images or strongly regularized natural images 4 or 5 cuts seems  
 425 to suffice, a very complex image with very weak regularization might need many more.  
 426 In the end, our algorithm is only efficient on problems whose solutions do not have too  
 427 many components. E.g. in the deblurring task, it is competitive for solutions with up  
 428 to 10,000 components for a  $512 \times 512$  image.

429 **2.6. Regularization path of the total variation.** Since the regularization coefficient  $\lambda$  is  
 430 difficult to choose a priori, it is typically useful to compute an approximate regularization path,  
 431 that is the collection of solutions to (1) for a set of values  $\lambda_0 > \dots > \lambda_j > 0$ . For  $\ell_1$  sparsity,  
 432 [21] showed how a fraction of the exact regularization path can be computed in a time of the  
 433 same order of magnitude as the time need to compute of the last point. In general, when  
 434 the path is not piecewise linear, the exact path cannot be computed, but similar results have  
 435 been shown for group sparsity [61, 52]. The case of total variation has been studied as well for

436 1-dimensional signals in [7]. We propose a warm-start approach to compute an approximate<sup>7</sup>  
 437 solution path for the total variation.

438 The rationale behind our approach is that, if  $\lambda_i$  and  $\lambda_{i+1}$  are close, the associated solutions  
 439  $x_i^*$  and  $x_{i+1}^*$  should also be similar, as well as their associated optimal partition, which we will  
 440 refer to as  $\Pi_i^*$  and  $\Pi_{i+1}^*$ . Consequently, it is reasonable to use a warm-start technique which  
 441 consists of initializing Algorithm 1 with  $\Pi_i^*$  to solve the problem associated with  $\lambda_{i+1}$  and  
 442 to expect that it will converge in a small number of binary cuts. It is important to note  
 443 that while our algorithm lends itself naturally to warm-starts, to the best of our knowledge  
 444 similar warm-start techniques do not exist for proximal splitting approaches such as [57] or  
 445 [14]. Indeed solutions whose primal solutions are close can have vastly different auxiliary/dual  
 446 solutions, and in our experiments no initialization heuristics consistently outperformed a naive  
 447 initialization.

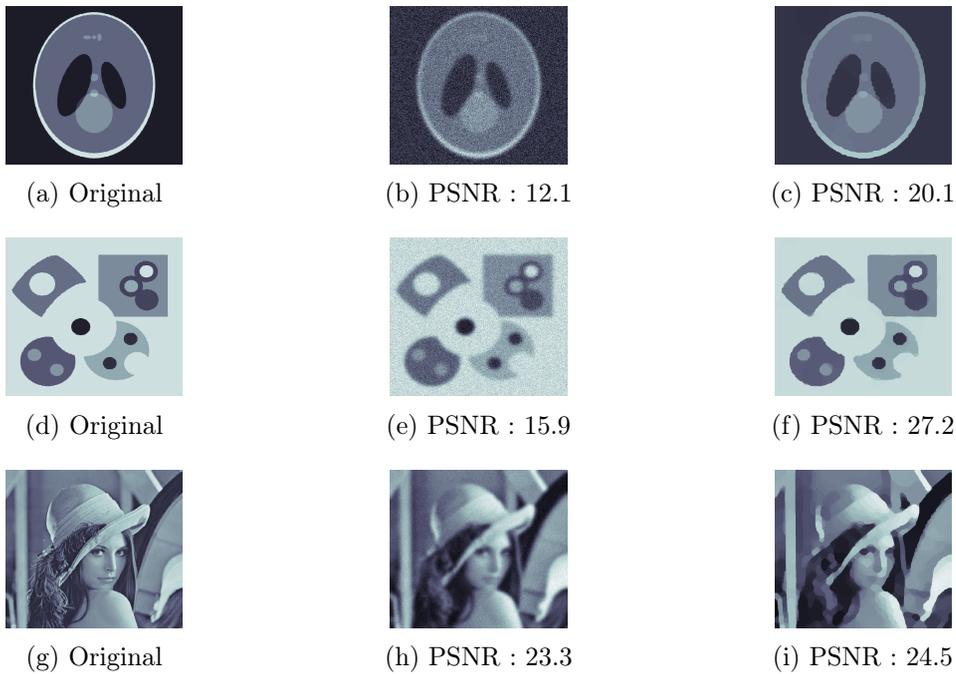


Figure 5: Benchmark on the deblurring task. Left column: original images, Middle column: blurred images, Right column: images retrieved by cut pursuit (CP)

448 **2.7. Numerical experiments: deblurring with TV.** To assess the performance in terms  
 449 of speed of our working set algorithm for the total variation regularization, we compare it with

---

<sup>7</sup>In fact for a quadratic data fitting term regularized by the total variation, the regularization path is piecewise linear and could thus in theory computed exactly, with a scheme similar to the LARS algorithm [21]. It should however be expected that this path has many point of discontinuity of the gradient, which entails that the cost of computation of the whole path is likely to be prohibitively high. We therefore do not consider further this possibility.

several state-of-the-art algorithms on a deblurring task of the form presented in section 2.4. Specifically, given an image  $x$ , we compute  $y = Hx + \epsilon$ , where  $H$  is a Gaussian blur matrix, and  $\epsilon$  is some Gaussian additive noise, and we solve (1) with a total variation regularization based on the 8-neighborhood graph built on image pixels. We use three  $512 \times 512$  images of increasing complexity to benchmark the algorithms: the Shepp-Logan phantom, a simulated example, and Lena, all displayed in Figure 5. For all images the standard deviation of the blur is set to 5 pixels.

A C++ implementation of the cut pursuit algorithm is available on the first author’s page<sup>8</sup>.

### 2.7.1. Competing methods.

**Preconditioned Generalized Forward Backward (PGFB).** As a general baseline, we consider a recent preconditioned generalized forward-backward splitting algorithm by [58] whose prior non-preconditioned version was shown to outperform state-of-the-art convex optimization on deblurring tasks in [57], including among others the algorithm of [14]. [58] demonstrate the advantages of the preconditioning strategy used over other adaptive metric approaches, such as the preconditioning proposed in [56] and the inertial acceleration developed in [44].

**Accelerated forward-backward with parametric max-flows (FB+).** Since efficient algorithms that solve the ROF problem have been the focus of recent work, and given that the ROF problem corresponds to the computation of the proximal operator of the total variation, we also compare with an implementation of the accelerated forward-backward algorithm of [49]. To compute the proximal operator, we use an efficient solver of the ROF problem based on a reformulation as a parametric max-flow proposed by [13]. The solver we use is the one made publicly available by the authors, which is based on a divide and conquer approach that works through the resolution of a parametric max-flow problem. This implies computing a sequence of max-flow problems, whose order make it possible to re-use the search trees in the [10] algorithm, thereby greatly speeding up computations.

**Cut pursuit with Frank-Wolfe descent direction (CPFWD).** We consider an alternative to the steepest binary partition to split the existing components of the partial solution: Inspired by the conditional gradient algorithm for regularized problems proposed by [31], consider a variant of cut pursuit in which we replace the steepest binary cut by the cut  $(B, B^c)$  such that  $\mathbf{1}_B$  is the Frank-Wolfe direction for the total variation, i.e. minimizing  $w(B, B^c)^{-1} \langle \nabla f(x), \mathbf{1}_B \rangle$  (see the discussion at the end of Section 2.1 and Appendix A). Note that the corresponding minimization of a ratio of combinatorial functions can in this setting be done efficiently using a slight modification of the algorithm of [20]. See Appendix D for more details. We chose not to make direct comparisons with the algorithms of [31] and of [2, Chap. 7.12], since it is clear that these algorithms will be outperformed by CPFWD. Indeed, these algorithms include a single term of the form  $\mathbf{1}_A$  in the expansion of  $x$  at each iteration, while CP and CPFWD grow much faster the subspace in which  $x$  is sought (its dimension typically more than doubles at each iteration). This entails that these algorithms must be slower than CPFWD, because for the former and for the latter, a single iteration requires to compute a Frank-Wolfe step, which requires solving several graph-cuts on the whole graph, and, as we discuss in Section 2.7.2 and

---

<sup>8</sup><https://github.com/loicland/cut-pursuit>

491 illustrate in Figure 7, the cost of graph-cuts already dominates the per iteration cost of CP  
 492 and CPFW.

493 **Cut pursuit.** To implement our algorithm (CP), we solve min-cut problems using the [37]  
 494 solver, which itself is based on [10] and [39]. The problems on the reduced graph are solved  
 495 using the PGFB algorithm. This last choice is motivated by the fact that the preconditioning  
 496 is quite useful as it compensates for the fact that the weights on the reduced graph can be  
 497 quite imbalanced.

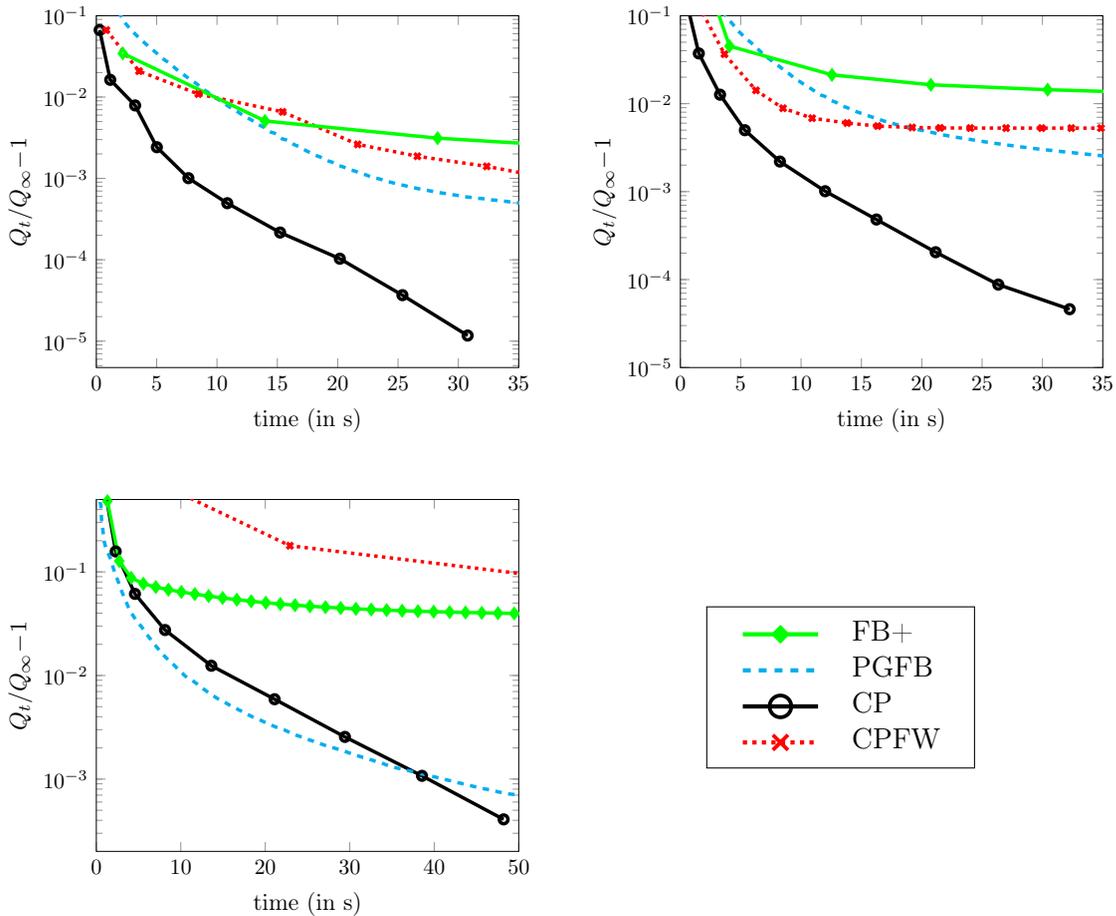


Figure 6: Relative primal suboptimality gap  $Q_t/Q_\infty - 1$  at time  $t$  (in seconds) for different algorithms on the deblurring task: accelerated forward backward (FB+), Preconditioned Generalized Forward Backward (PGFB), Cut pursuit (CP) and a variant using Frank-Wolfe directions (CPFW), and for different  $512 \times 512$  images and different regularization values: Shepp-Logan phantom (left), our simulated example (middle) and Lena (right). The marks in (FB+), (CP) and (CPFW) corresponds to one iteration.

498 **2.7.2. Results.** Figure 6 presents the convergence speed of the different approaches on the  
 499 three test images on a quad-core CPU at 2.4 Ghz. Precisely, we represent the relative primal  
 500 suboptimality gap  $(Q_t - Q_\infty)/Q_\infty$  where  $Q_\infty$  is the lowest value obtained by CP in 100 seconds.  
 501 We can see that our algorithm significantly speeds up the direct optimization approach PGFB  
 502 when the solution is sparse, and that it remains competitive in the case of a natural image  
 503 with strong regularization. Indeed since the reduced problems are of a much smaller size than  
 504 the original, our algorithm can perform many more forward-backward iterations in the same  
 505 allotted time.

506 The variant of cut pursuit using Frank-Wolfe directions (CPFW) is as efficient over the  
 507 first few iterations but then stagnates. The issue is that the computation of a new Frank-Wolfe  
 508 direction does not take into account the current support  $S(x)$  which provides a set of edges  
 509 that are “free”; this means that the algorithm overestimates the cost of adding new boundaries,  
 510 resulting in overly-conservative updates.

511 Accelerated forward-backward with parametric max-flow (FB+) is also slower than the cut  
 512 pursuit approach in this setting. This can be explained by the fact that the calls to max-flow  
 513 algorithms, represented by a mark on the curve, are better exploited in the cut pursuit setting.  
 514 Indeed in the forward-backward algorithm, the solutions of parametric max-flow problems  
 515 are exploited by performing one (accelerated) proximal gradient step. By contrast, in the  
 516 cut pursuit setting, the solution of each max-flow problem is used to optimize the reduced  
 517 problem. Since the reduced graph is typically much smaller than the original, a precise solution  
 518 can generally be obtained very quickly, yet providing a significant decrease in the objective  
 519 function. Furthermore, as the graph is split into smaller and smaller independent connected  
 520 components by cut pursuit, the call to the max-flow solver of [10] are increasingly efficient  
 521 because the augmenting paths search trees are prevented from growing too wide, which is the  
 522 main source of computational effort.

523 Figure 7 presents the breakdown of computation time for each algorithm over 60 seconds of  
 524 computation. In PGFB, the forward-backward updates naturally dominate the computation  
 525 time, as well as the fast Fourier transform needed to compute the gradient at each iteration.  
 526 In FB+, the computation of the proximal operator of the partial solution through parametric  
 527 maximum flows is by far the costliest. Our approach and CPFW share a similar breakdown  
 528 of computation time as their structures are similar. The maximum flow represents the highest  
 529 cost, with the fast Fourier transform needed to compute  $K^\top H^\top H K$  a close second. Finally  
 530 diverse operations such as computing the reduced graph takes a small fraction of the time.  
 531 More interestingly, solving the reduced problem (with the PGFB subroutine of CP) takes  
 532 comparatively very little time (roughly 3%) when this is the only step that actually decreases  
 533 the objective function. This is expected as, even at the last iteration, the reduced graph had  
 534 only 300 components so that the associated problem is solved very rapidly.

535 **2.8. Numerical experiments: approximate TV regularization path.** We now present the  
 536 computation of an approximate regularization path for the ROF minimization, using warm-  
 537 starts as described in Section 2.6. We consider the task of ROF-denoising on three natural  
 538 images presented in Figure 9. For each image we pick 20 values of  $\lambda$  evenly distributed  
 539 logarithmically in the range of parameters inducing from coarse to perfect reconstructions.

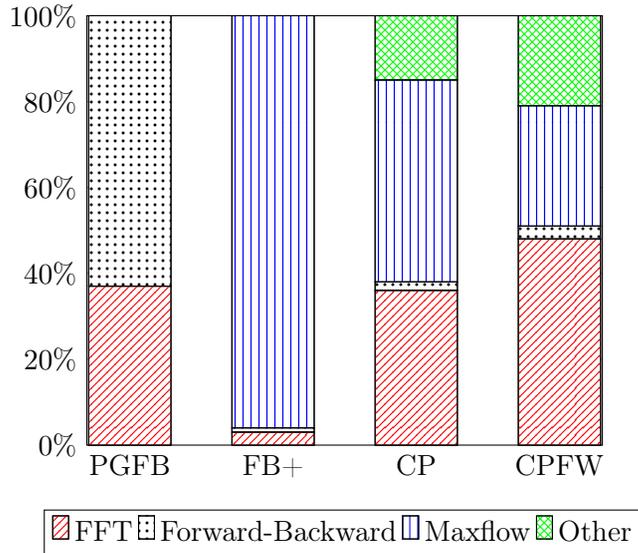


Figure 7: Time breakdown for the different algorithms over 60 seconds of optimization.

540 **2.8.1. Competing methods. Parametric max-flows (PMF).** We use the parametric  
 541 max-flow based ROF solver of [13] to compute each value. In our numerical experiments, it  
 542 was the fastest of all available solvers, and moreover returns an exact solution.

543 **Cut pursuit (CP).** We use the algorithm presented in this paper to separately compute the  
 544 solutions for each parameter value. The algorithm stops when it reaches a relative primal  
 545 suboptimality gap  $Q_t/Q_\infty - 1$  of  $10^{-5}$ , with  $Q_\infty$  the exact solution given by PMF.

546 **Cut pursuit path (CPP).** We use the warm start approach proposed in Section 2.6, with  
 547 the same stopping criterion.

548 **2.8.2. Results.** We report in Figure 9 the time in seconds necessary to reach a primal  
 549 suboptimality gap of  $10^{-5}$  for the different approaches. We observe that, in general, cut  
 550 pursuit (CP) is slightly faster than the parametric max-flow. It should be noted, however,  
 551 that the latter finds an exact solution and remains from that point of view superior. Warm-  
 552 starts allow for a significant acceleration, needing at most two calls to the max-flow code to  
 553 reach the desired gap. Unlike the deblurring task, for high noise levels, cut pursuit remains  
 554 here very competitive for natural images which are not sparse, as illustrated in Table 10 and  
 555 Figure. 8.

556 As the regularization strength decreases, the coarseness of the solution decreases, and as  
 557 a consequence the cut pursuit approaches CP and CPP become less and less efficient. This is  
 558 because as the number of components increases, so does the time needed to solve the reduced  
 559 problem. We note however that for the values provided with the peak PSNR, the warm-start  
 560 approach is faster than PMF.

561 PMF and CP perform significantly worse on sparse images and for high values of  $\lambda$ . This  
 562 can be explained by the inner workings of the max-flow algorithm of [10]. Indeed for high

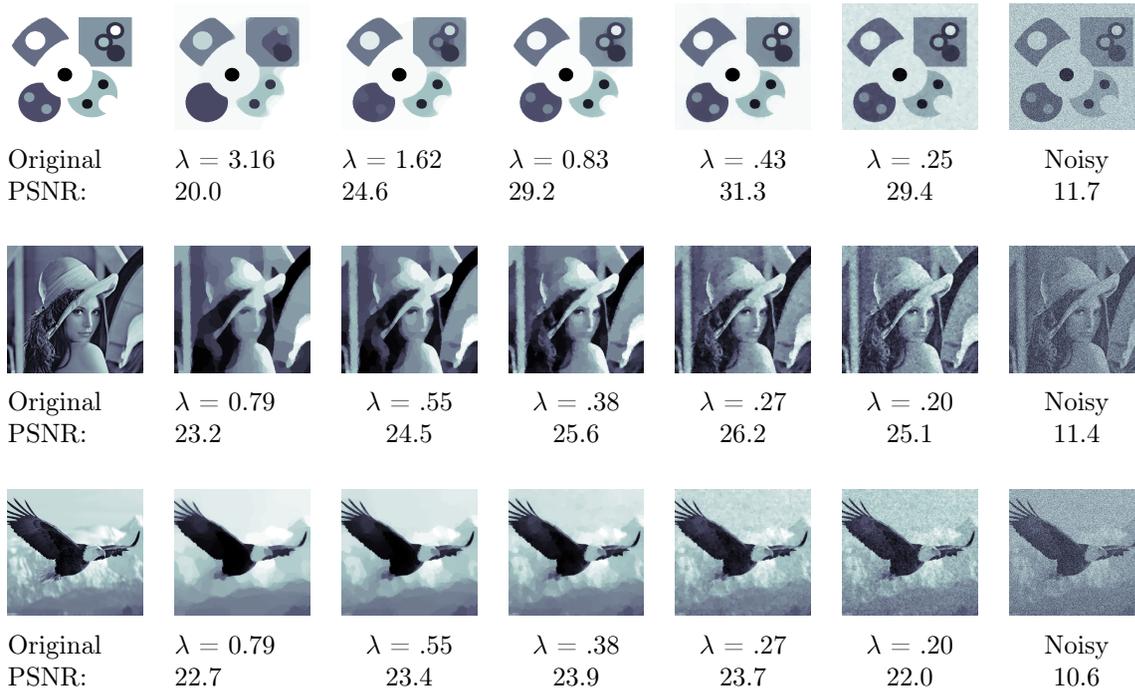


Figure 8: Illustration of the regularization path for the three images in the data set for 5 of the 20 values in the regularization parameters in the path. The peak PSNR is reached for  $\lambda = 0.53$ ,  $0.28$  and  $0.34$  respectively.

563 values of  $\lambda$  or sparse images, the pairwise term of the corresponding Potts model will dominate,  
 564 which forces the algorithm to build deep search trees to find augmenting paths. Indeed as the  
 565 size of the regions formed by the cut increase, the combinatorial exploration of all possible  
 566 augmenting paths drastically increases as well. The warm-started path approach does not  
 567 suffer from this problem because the graph is already split in smaller components at the  
 568 warm-start initialization, which prevents the search trees from growing too large.

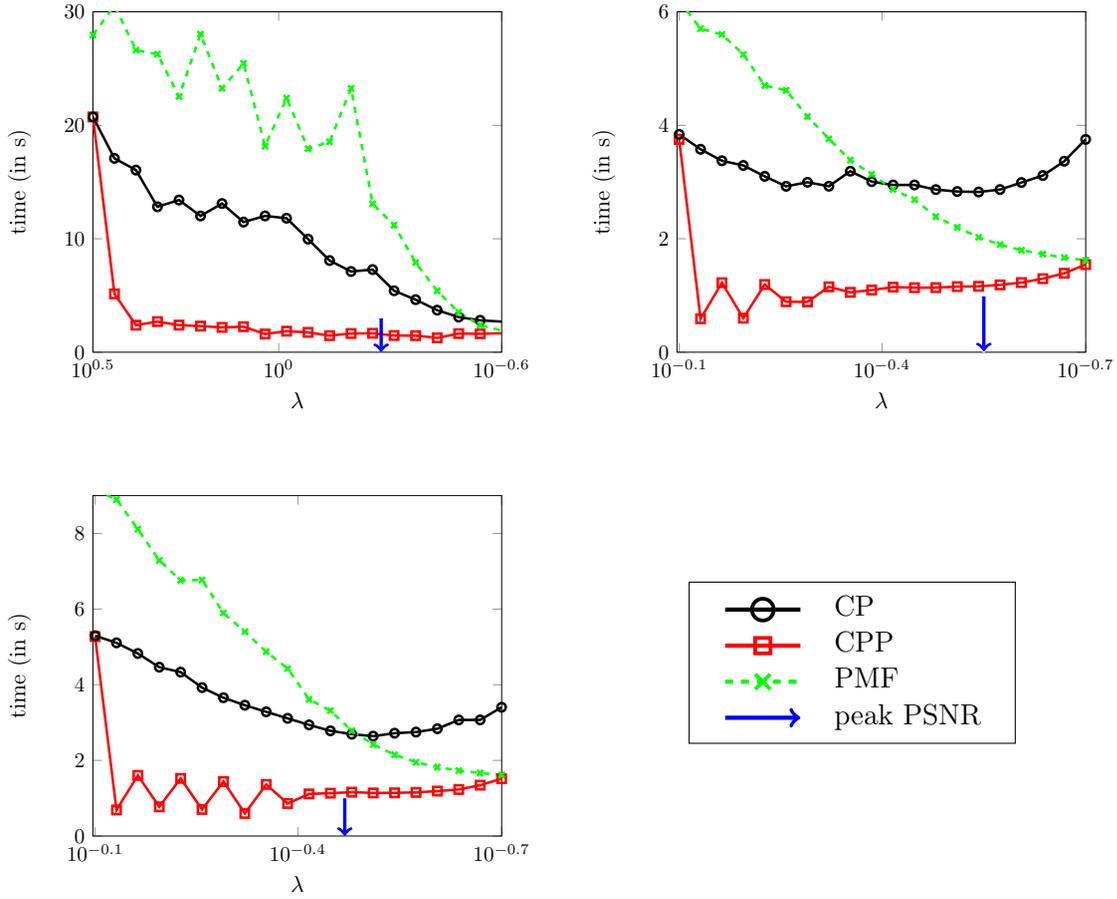


Figure 9: Time in seconds necessary to solve the problem regularized with a given  $\lambda$  (from the warm-start initialization when applicable) with a relative primal suboptimality gap of  $10^{-5}$ , for regularly sampled values of  $\lambda$  along the regularization path. The competing methods are cut pursuit (CP), cut pursuit with warm-start (CPP) and the parametric max-flow solver (PMF) for different  $512 \times 512$  noisy images: simulated example (left), Lena (middle) and eagle (right). The computation times are averaged over 10 random degradations of the images by uniform noise. The blue arrow indicates the best PSNR value.

569 **3. Generalized minimal partition.** We consider now a generalization of the minimal parti-  
 570 tion problem  $\min_{x \in \mathbb{R}^n} Q(x)$  with  $Q(x) = f(x) + \lambda \Gamma(x)$  where  $\Gamma(x) \doteq \frac{1}{2} \sum_{(i,j) \in S(x)} w_{ij}$  the total  
 571 boundary size penalty for piecewise constant functions. This non-convex non-differentiable  
 572 problem being significantly harder than the previous one, we restrict the functions  $f$  we con-  
 573 sider to be separable functions of the form  $f(x) = \sum_{i \in V} f_i(x_i)$  with  $f_i : \mathbb{R} \mapsto \mathbb{R}$  continuous<sup>9</sup>.

<sup>9</sup>The algorithmic scheme we propose in this section does not require the functions  $f_i$  to be convex, but convexity will make subproblems easier to solve, and, as discussed later, can be helpful to establish sufficient conditions for convergence (see Section 3.1.1 and Appendix E.1.2)

Method	Simulated	Lena	Eagle
CPP	59	25	27
CP	194	62	70
PMF	356	67	91

Figure 10: Time in seconds necessary to compute the entire approximate regularization path at a relative primal suboptimality gap of  $10^{-5}$  for the different algorithms, averaged over 10 samplings of the noise.

574 Our formulation, like [42], but unlike most instances of the minimal partition problem in the  
 575 literature, does not constrain the number of components in advance. We call the corresponding  
 576 problem *generalized minimal partition problem*.

577 Inspired by greedy feature selection algorithms in the sparsity literature and by the working  
 578 set algorithm we presented for TV regularization, we propose to exploit the assumption that  
 579 the optimal partition  $\Pi^*$  is not too large to construct an algorithm that greedily optimizes the  
 580 objective by adding and removing cuts in the graph.

581 Indeed, the problem that we consider has a fixed regularization coefficient  $\lambda$ , and so its  
 582 natural counterpart for classical sparsity is the problem of minimizing an objective of the form  
 583  $f(x) + \lambda\|x\|_0$  which subsumes AIC, BIC and other information criteria. The algorithmic ap-  
 584 proach we consider is thus the counterpart of a very natural greedy algorithm to minimize the  
 585 former objective, which surprisingly is almost absent from the literature, perhaps for the fol-  
 586 lowing reasons: On the one hand, work on *stagewise regression* and forward-backward greedy  
 587 algorithms, which both add and remove variables, goes back to the 60ies [22], but the algo-  
 588 rithms then considered were based on sequences of tests as opposed to a greedy minimization  
 589 of a penalized criterion.

590 On the other hand, the literature on greedy algorithms for sparse models has almost ex-  
 591 clusively focused on solving the constrained problem  $\min_x f(x)$  s.t.  $\|x\|_0 \leq k$ , with algorithms  
 592 such as OMP, Orthogonal least squares (OLS), FoBa, and CoSamp, which can alternatively  
 593 be viewed as algorithms that are greedily approximating the corresponding Pareto frontier. A  
 594 notable exception is IHT.

595 A very natural variant of OLS solving  $\min_x f(x) + \lambda\|x\|_0$  can however be obtained by  
 596 adding the  $\ell_0$  penalty to the objective. This algorithm was formally considered in [64] under the  
 597 name Single Best Replacement (SBR), in reference to the similar Single Maximum Likelihood  
 598 Replacement (SMLR) of [40]. At each iteration, the algorithm considers adding (forward step)  
 599 or removing (backward step) a single variable, whichever reduces the value of the objective  
 600 most. It should be noted that while the similar OLS and OMP are forward algorithms, SBR is  
 601 a forward-backward algorithm, which can remove a variable provided doing so only increases  
 602  $f$  by less than  $\lambda$ .

603 We argue in the following section that a similar algorithm can be designed for the general-  
 604 ized minimal partition problem, using a general scheme which is similar to that of cut pursuit.  
 605 We thus call this algorithm  $\ell_0$ -cut pursuit. In particular, it follows a similar structure, in which  
 606 a partition is successively split into its constant connected components. The main differences is

607 an adapted rationale to split elements of the partition, and the addition of a explicit backward  
608 step.

609 **3.1. A greedy algorithm for generalized minimal partition.** As in cut pursuit, we propose  
610 an algorithm which greedily splits the elements of the current partition  $\Pi = (A_1, \dots, A_k)$  in  
611 forward steps, reoptimizes the value taken by  $x$  on each of the  $A_j$ , then, in backward steps,  
612 possibly merges some regions (or moves some of the boundaries between regions), and iterates.

**3.1.1. Forward step.** Assume that we split the set of existing regions  $(A_j)_{1 \leq j \leq k}$  by introducing a global cut  $(B, B^c)$  for some set  $B \subset V$ , so as to minimize the global objective, i.e.

$$\min_{B \subset V} \min_{(h_j, h'_j)_{1 \leq j \leq k}} \sum_{j=1}^k \left[ \sum_{i \in A_j \cap B} f_i(h_j) + \sum_{i \in A_j \cap B^c} f_i(h'_j) \right] + \lambda \sum_{j=1}^k w(A_j \cap B, A_j \cap B^c)$$

613 This cut induces a cut on each element  $A_j$  of the form  $(A_j \cap B, A_j \cap B^c)$ . Two simple  
614 properties should be noted: (a) the additional boundary perimeter incurred with the cut is  
615 simply the sum of the perimeters of the cuts induced within each element  $A_j$  and is precisely  
616 of the form  $\sum_{j=1}^k w(A_j \cap B, A_j \cap B^c)$  — the boundary between pre-existing components is  
617 “free” (cf Figure 2), (b) if the value of  $x$  is re-optimized under the constraint that it should be  
618 constant on each of the elements  $A_j \cap B$  and  $A_j \cap B^c$  of the new partition, then the separability  
619 of  $f$  and the fact that  $\Gamma(x)$  stays constant when the value of each of the regions is modified  
620 together entail that the optimization can be done separately on each set  $A_j$ . So the choice of  
621 an optimal cut reduces to independent choices of optimal cut on each set  $A_j$  as defined by the  
622 objective

$$623 \quad (8) \quad \min_{B_j \subset A_j} \min_{(h, h')} \sum_{i \in B_j} f_i(h) + \sum_{i \in A_j \setminus B_j} f_i(h') + \lambda w(B_j, A_j \setminus B_j).$$

624 This optimization problem is difficult to solve globally, because even if the functions  $f_i$   
625 were assumed convex, it would not be a convex optimization problem. However, for  $B_j$  fixed,  
626 the partial minimization with respect to  $h$  and  $h'$  is an optimization problem in  $\mathbb{R}^2$ , and, for  
627  $(h, h')$  fixed, the optimisation with respect to  $B_j$  is solved as a min-cut/max-flow problem  
628 very similar to the one for the steepest binary cut of Section 2.1. We therefore propose  
629 the alternating minimization algorithm presented in pseudo-code as Subroutine 2. Under  
630 appropriate hypotheses on  $f$  detailed in Appendix E.1.2, this algorithm finds a local minimum  
631 of the objective. In particular, these hypotheses hold if each  $f_i$  is strictly convex and in general  
632 position so to as to avoid ties in assignments of  $i$  to  $B$  or  $B^c$ , for example if  $f_i(\cdot) = (\cdot - x_i)^2$   
633 with  $x_i$  drawn i.i.d. from a continuous distribution, which corresponds to our case of interest.

634 In this algorithm, since the minimization with respect to  $B_j$  can lead to more than two  
635 connected components, we use the same idea as presented in Section 2.2 and illustrated on  
636 Figure 2, which is to treat each connected component as a new element of the partition.

637 Further details on Subroutine 2 and initialization strategies are discussed in Appendix E.1.

638 **3.1.2. Saturated sets.** A particular situation occurs when the optimal solution  $B_j$  of  
639 problem (8) is equal to  $\emptyset$  or  $A_j$ : in that case, any split of  $A_j$  would increase the objective.

640 We then say that the component  $A_j$  is *saturated*. The overall algorithm maintains a set  $\Sigma$  of  
 641 *saturated* components which do not need to be processed anymore in the splitting steps.

642 For cut pursuit (i.e. in the TV case), it was essentially sufficient to design the splitting step  
 643 to specify the algorithm: indeed, after splitting with a steepest binary cut, the problem solved  
 644 on the reduced graph involved in that case a total variation term penalizing the difference of  
 645 values between adjacent regions (cf Proposition 6), and this TV term could thus induce the  
 646 merge of two adjacent regions. By contrast, for  $\ell_0$  cut pursuit, given that the optimization of  
 647 the values on each region is independent and without any incidence on the definition of their  
 648 contours, merge steps and other steps to modify the shape of the regions should be added  
 649 explicitly. We discuss them in the next two sections.

650 **3.1.3. Backward steps.** In greedy algorithms for plain sparsity, backward steps remove  
 651 variables to reduce the support of the solution. In our case, the appropriate notion of support  
 652 is  $S(x)$  (cf Equation 2), which is formed as the union of the boundaries between pairs of  
 653 components. A backward step is a step that reduces the total boundary perimeter. The most  
 654 natural way to obtain this is by merging two adjacent components.

**Simple merge step:** For a region  $C$ , let  $f_C^* := \min_h \sum_{i \in C} f_i(h)$ . If a pair of adjacent  
 components  $(A, B)$  is merged into a single constant component, and the value of  $A \cup B$  is,  
 reoptimized, the objective  $Q$  increases by

$$\delta_-(A, B) := f_A^* + f_B^* - f_{A \cup B}^* + \lambda w(A, B).$$

655 A merge effectively decreases the value of the objective and is thus worth it if  $\delta_-(A, B) > 0$   
 656 i.e. if  $f_{A \cup B}^* - (f_A^* + f_B^*) < \lambda w(A, B)$ .

657 It should be noted that the merge step considered does not, in general, correspond to  
 658 canceling exactly a previous cut, but can merge adjacent subregions that have each been  
 659 obtained by splitting different regions. The merge step is described in Subroutine 4.

660 A shortcoming of the simple merge step is that while the removal of boundaries between  
 661 components is considered, a simple change of the shape of the created boundaries that could  
 662 reduce total boundary length is not possible. However, since the optimal binary computation  
 663 only considers binary partitions, the shape of the components might be suboptimal. We  
 664 therefore propose another kind of step.

**Merge-resplit:** This step is a combination of a merge step immediately followed by a  
 new split step on the merged components. It is a “backward-then-forward” step, which can  
 be worth it even if the corresponding backward step taken individually is not decreasing the  
 objective. Given  $h_A := \arg \min \sum_{i \in A} f_i(A)$  and  $h_B := \arg \min \sum_{i \in B} f_i(B)$ , the merge resplit  
 step amounts to solve the corresponding

$$\min_{A', B'} \sum_{i \in A'} f_i(h_A) + \sum_{i \in B'} f_i(h_B) + \lambda w(A', B') \quad \text{s.t.} \quad B' = (A \cup B) \setminus A'.$$

665 But this problem can again be solved as a min-cut/max-flow problem on the region  $A \cup B$ .

666 Note that this merge-resplit step is very similar to what [10] call an  $\alpha$ - $\beta$  *swap* in the context  
 667 of energy minimization in Markov random fields: nodes assigned to other components<sup>10</sup> than

---

<sup>10</sup>In the context of MRFs the components correspond to a number of different classes fixed in advance and are in general not connected.

668  $A$  or  $B$  keep their current assignments to components, but the nodes of  $A \cup B$  are reassigned  
 669 to  $A$  or  $B$  so that the boundary between  $A$  and  $B$  minimizes the above energy.

670 The merge-resplit step includes the possibility of a simple merge step (without resplitting),  
 671 since all elements can be “swapped” in the same set by the  $\alpha$ - $\beta$  swap, so that the new boundary  
 672 is effectively empty. Finally, note that during the merge-resplit step the values of  $x_A$  and  $x_B$   
 673 are held constant and only updated upon completion of the step. In fact, in a number of cases,  
 674 it might be possible to iterate such steps for a given pair  $(A, B)$ . We do not consider this  
 675 computationally heavier possibility.

676 **3.1.4. The  $\ell_0$  cut pursuit algorithm.** Given definitions of forward and backward steps,  
 677 different algorithms can be obtained by iterating and alternating these steps differently. We  
 678 propose to alternate between splitting all components at once (possibly in parallel) and then  
 679 iterating backward steps over all adjacent pairs of components. This allows for the splitting  
 680 to be done in parallel directly on the original flow graph, thus avoiding the memory overheads  
 681 associated with constructing a new flow graph for each new component. This leads to two  
 682 variants for the main algorithm which are presented as Algorithms 5 and 6, depending on  
 683 whether only simple merge or merge-resplit steps are used. Implementation details of the  
 684 algorithms and other possible variants are discussed in Appendix E.2.

685 Under mild assumptions, Algorithm 5 converges in a finite number of iterations and yields  
 686 a partition  $\Pi = (A_1, \dots, A_n)$  such that  $x_\Pi \doteq \arg \min_{z \in \text{span}(\Pi)} Q(z)$  is a local minimum of  $Q$ .  
 687 See in Appendix E.3 for a precise statement and a proof.

---

**Subroutine 2**  $[\Pi, \mathcal{E}, \Sigma] \leftarrow \text{split}(\Pi, \mathcal{E}, \Sigma, A)$

---

*[Splits the component  $A$  with a binary cut: updates the current partition  $\Pi$ , the component adjacency structure  $\mathcal{E}$  and the set of saturated components  $\Sigma$ ]*

**for**  $A \in \Pi$  **do**

$\Pi \leftarrow \Pi \setminus \{A\}$

$B \leftarrow \arg \min_{B \subset A, h, h'} \sum_{i \in B} f_i(h) + \sum_{i \in B^c} f_i(h')$

**while not\_converged do**

$x \leftarrow \arg \min_h \sum_{i \in B} f_i(h)$

$x' \leftarrow \arg \min_h \sum_{i \in A \setminus B} f_i(h)$

$B \leftarrow \arg \min_{B \subset A} \sum_{i \in B} f_i(x) + \sum_{i \in B^c} f_i(x') + \lambda w(B, B^c)$

**end while**

**if**  $B \in \{\emptyset, A\}$  **then**

$\Sigma \leftarrow \Sigma \cup \{A\}$

**end if**

$[B_1, \dots, B_k] \leftarrow$  connected components of  $B$  and  $A \setminus B$

$\Pi \leftarrow \Pi \cup \{B_1, \dots, B_k\}$

$\mathcal{E} \leftarrow$  updated adjacency structure

**end for**

---

---

**Subroutine 3**  $[\Pi, \mathcal{E}, \Sigma] \leftarrow \text{simple\_merge}(\Pi, \mathcal{E}, \Sigma, A, B)$ 


---

*[Merges components A and B]*  
 $\Pi \leftarrow \Pi \setminus \{A, B\} \cup \{A \cup B\}$   
 $\mathcal{E} \leftarrow \mathcal{E} \setminus \{\{A, B\}\}$   
 $\Sigma \leftarrow \Sigma \setminus \{A, B\}$   
**for**  $C$  neighbors of  $A$  or  $B$  **do**  
     $\mathcal{E} \leftarrow \mathcal{E} \cup \{\{A \cup B, C\}\}$   
**end for**

---



---

**Subroutine 4**  $[\Pi, \mathcal{E}, \Sigma] \leftarrow \text{resplit}(\Pi, \mathcal{E}, \Sigma, A, B)$ 


---

*[Performs a merge-resplit step on components A and B.]*  
 $x_A \leftarrow \arg \min_h \sum_{i \in A} f_i(h)$   
 $x_B \leftarrow \arg \min_h \sum_{i \in B} f_i(h)$   
 $C \leftarrow \arg \min_{C \subset A \cup B} \sum_{i \in C} f_i(x_A) + \sum_{i \in A \cup B \setminus C} f_i(x_B) + \lambda w(C, A \cup B \setminus C)$   
**if**  $C \notin \{A, B\}$  **then**  
     $\Sigma \leftarrow \Sigma \setminus \{A, B\}$   
**else**  
     $[C_1, \dots, C_k] \leftarrow$  connected components of  $C$  and  $A \cup B \setminus C$   
     $\Pi \leftarrow \Pi \setminus \{A, B\} \cup \{C_1, \dots, C_k\}$   
     $\mathcal{E} \leftarrow$  updated adjacency structure  
**end if**

---



---

**Algorithm 5** Simple merge variant  
 $(\ell_0\text{-CPm})$ 


---

**Initialization:**  $\Pi_0 = \{V\}, \mathcal{E} = \Sigma = \emptyset$   
**while**  $\Pi \neq \Sigma$  **do**  
    **for**  $A \in \Pi \setminus \Sigma$  in parallel **do**  
         $[\Pi, \mathcal{E}, \Sigma] \leftarrow \text{split}(\Pi, \mathcal{E}, A, \Sigma)$   
    **end for**  
    Compute  $\delta_-(A, B)$  for all  $(A, B) \in \mathcal{E}$   
    **while**  $\max_{\{A, B\} \in \mathcal{E}} \delta_-(A, B) > 0$  **do**  
         $\{A, B\} = \arg \max_{\{A', B'\} \in \mathcal{E}} \delta_-(A', B')$   
         $[\Pi, \mathcal{E}', \Sigma] \leftarrow \text{merge}(\Pi, \mathcal{E}, \Sigma, A, B)$   
        Update  $\delta_-(A, B)$  for  $\{A, B\} \in \mathcal{E}' \setminus \mathcal{E}$   
         $\mathcal{E} \leftarrow \mathcal{E}'$   
    **end while**  
**end while**

---



---

**Algorithm 6** Merge-resplit variant  
 $(\ell_0\text{-CPs})$ 


---

**Initialization:**  $\Pi_0 = \{V\}, \mathcal{E} = \Sigma = \emptyset$   
**while**  $\Pi \neq \Sigma$  **do**  
    **for**  $A \in \Pi \setminus \Sigma$  in parallel **do**  
         $[\Pi, \mathcal{E}, \Sigma] \leftarrow \text{split}(\Pi, \mathcal{E}, \Sigma, A)$   
    **end for**  
     $\mathcal{E}' \leftarrow \mathcal{E}$   
    **for**  $\{A, B\} \in \mathcal{E}'$  **do**  
        **if**  $\{A, B\} \in \mathcal{E}$  **then**  
             $[\Pi, \mathcal{E}, \Sigma] \leftarrow \text{resplit}(\Pi, \mathcal{E}, \Sigma, A, B)$   
        **end if**  
    **end for**  
**end while**

---

688

689 **3.2. Numerical experiments: denoising with  $\ell_0$  cut pursuit.** We now present experiments  
690 empirically demonstrating the superior performance of the  $\ell_0$ -cut pursuit algorithm presented  
691 in section 3. We assess its performance against two state-of-the art algorithms to minimize the  
692 total boundary size of two noisy  $512 \times 512$  images: the Shepp-Logan phantom [63] and another

693 simulated example. In order to illustrate the advantage of our algorithm over alternatives which  
 694 discretize the value range, we add a small random shift of grey values to both images. We also  
 695 test the algorithms on a spatial statistic aggregation problem using open-source data<sup>11</sup> which  
 696 consists of computing the statistically most faithful simplified map of the population density  
 697 in the Paris area over a regular grid represented in Figure 12. The raster is triangulated to  
 698 obtain a graph with 252,183 nodes and 378,258 edges. We use the squared loss weighted by  
 699 the surface of each triangle as a fidelity term.

700 A C++ implementation of the  $\ell_0$ -cut pursuit algorithm is available<sup>12</sup>.

### 701 3.2.1. Competing methods.

702  **$\alpha$ -expansions on quantized models (CRF $i$ ).** If the range of values of  $x_i$  is quantized,  
 703 the MPP and TV problems reduce to a Potts model, in which each class  $c$  is associated with  
 704 a (non necessarily connected) level-set [32]. In the MPP case, the pairwise terms are of the  
 705 form  $1_{\{c_i \neq c_j\}} w_{ij}$ . We use  $\alpha$ -expansions [10] to approximately minimize the corresponding  
 706 energy. More precisely, we use the  $\alpha$ -expansions implementation of [27], which uses the same  
 707 max-flow code [9] as our algorithm. We denote the resulting algorithm CRF $i$  where  $i$  is the  
 708 number of levels of quantization of the observed image value range. While this algorithm is  
 709 not theoretically guaranteed to converge, it does in practice and the local minima are shown  
 710 by [10] to be within a multiplicative constant of the global optimum.

711 **Non-convex relaxation (TV $_{0.5}$ ).** We considered a non-convex counterpart of the total  
 712 variation, similar to the formulations considered in [51] or [71], but with  $t \mapsto (\epsilon + t)^{\frac{1}{2}}$  in lieu  
 713 of  $t \mapsto |t|$ . The resulting functional can be minimized locally using a reweighted TV scheme  
 714 described in [53]. We use our cut pursuit algorithm to solve each reweighted TV problem as  
 715 it is the fastest implementation.

716  **$\ell_0$ -cut pursuit** We implemented three versions of  $\ell_0$  cut pursuit with different backward steps.  
 717 In the simplest instantiation,  $\ell_0$ -CPf, no backward step is used and the reduced graph can only  
 718 increase in size. In  $\ell_0$ -CPm, described in Algorithm 5, the simple merge step is performed after  
 719 each round of cuts. Finally in  $\ell_0$ -CPs, described in Algorithm 6, merge steps are replaced by  
 720 merge-resplit steps but without priority queue.

721 After a few preliminary experiments, we chose not to include either level-set methods [16] or  
 722 active contour methods based on solving Euler-Lagrange equations [36] as their performances  
 723 were much lower than the algorithms we consider.

724 Comparing speed results of code is always delicate as the degree of code optimization varies  
 725 from one implementation to another. The  $\alpha$ -expansion code uses the implementation of [27]  
 726 which is a highly optimized code,  $\ell_0$ -CPf and  $\ell_0$ -CPm are implemented in C++, while  $\ell_0$ -CPs  
 727 and TV $_{0.5}$  are implemented in Matlab with a heavy use of mex-files. Even if minor improve-  
 728 ments could be obtained on the latter, we believe that it would not change the performances  
 729 significantly. In particular, a justification for direct time comparisons here is that computation  
 730 time for each of the algorithms is mostly spent computing min cuts which is done in all codes  
 731 using the same implementation of [9] and which accounts for most of the computation time.

---

<sup>11</sup><https://www.data.gouv.fr/fr/datasets/donnees-carroyees-a-200m-sur-la-population>

<sup>12</sup><https://github.com/loicland/cut-pursuit>

732 **3.2.2. Results.** Given that the MPP is hard, and that all the algorithms we consider only  
 733 find local minima, we compare the different algorithms both in terms of running time and in  
 734 terms of the objective value of the local minima found. The marks on the curves correspond  
 735 to one iteration of each of the considered algorithms: For  $TV_{0.5}$  there is a mark for each  
 736 reweighted TV problem to solve, for  $CRF_k$ , a mark corresponds to one  $\alpha$ -expansion step,  
 737 i.e. solving  $k$  max-flow problems. For  $\ell_0$ -CP this corresponds to one forward (split) and one  
 738 backward step. For clarity, the large number of marks were omitted in the third experiment,  
 739 as well as for  $\ell_0$ -CPs in the first experiment.

740 In Figure 11, we report the energy obtained by the different algorithms normalized by  
 741 the energy of the best constant approximation. We can see that our algorithms find local  
 742 optima that are essentially as good or better than  $\alpha$ -expansions for the discretized problems  
 743 in less time, as long as the solutions are sufficiently sparse. For the population density data,  
 744 the implementation  $\ell_0$ -CPm with simple merge is faster and finds a better local minimum  
 745 than CRF40, but is outperformed by CRF60. The implementation with swaps merge-resplit  
 746 ( $\ell_0$ -CPs) is on par with CRF60 when it comes to speed, and finds a slightly better minimum.

747 The simple merge step provides a better solution than the purely forward approach at the  
 748 cost of a slight increase in computational time. The merge-resplit backward step improves the  
 749 quality of the solution further, but comes with a significant increase in computation.

750 We report in Table 13 performance in PSNR that shows that  $\ell_0$ -CP outperforms the CRF  
 751 formulations for quantization levels that lead to comparable running time.

752 The comparison with CRF formulations is investigated in more details in Appendix F,  
 753 where we report the performance of approximations with CRFs solved with iterative  $\alpha$ -expansions  
 754 for different numbers of quantization levels, as compared with the performance of  $\ell_0$ -CPm. The  
 755 results show that the running time for the CRF formulations grows linearly with the num-  
 756 ber of classes, although the performance in PSNR does not increase monotonically, and has  
 757 oscillations which lead to results that are worse than  $\ell_0$ -CPm for some number of classes.

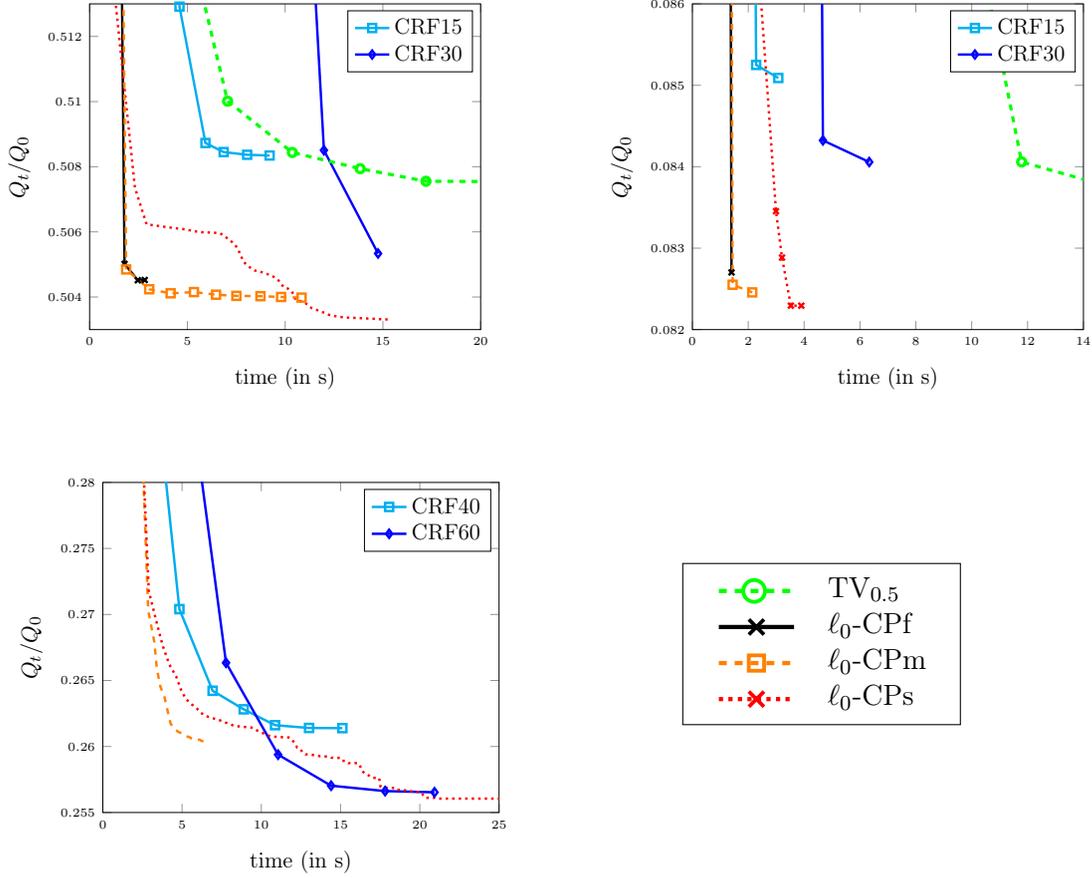


Figure 11: *Generalized minimal partition energy* at time  $t$  (in seconds) divided by the same energy for the best constant approximation obtained by different algorithms: Non-convex relaxation ( $TV_{0.5}$ ),  $\ell_0$ -CPf with no backward step,  $\ell_0$ -CPm with simple merge step,  $\ell_0$ -CPs with merge-resplit steps, and finally,  $\alpha$ -expansions with different number of levels of quantization (see image legends), for different images: the Shepp-Logan phantom (left), our simulated example (middle) and the map simplification problem (right). Markers correspond respectively to one reweighting, one  $\alpha$ -expansion cycle and one cut for ( $TV_{0.5}$ ), (CRF) and ( $\ell_0$ -CP).

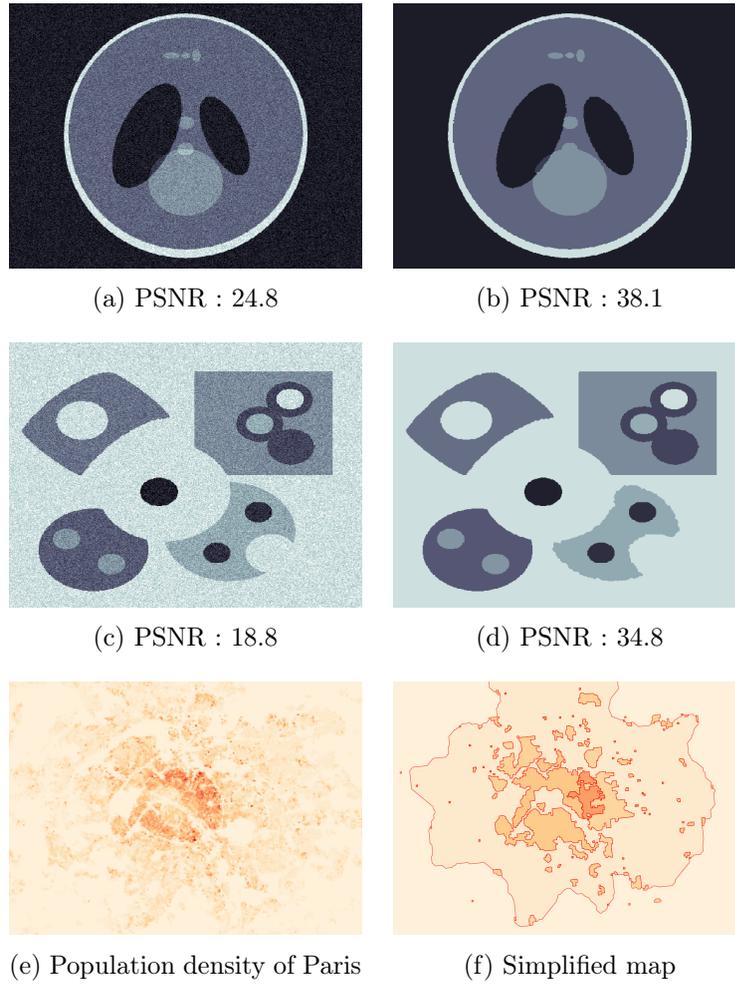


Figure 12: Benchmark on the denoising task. First two lines: (left) noisy images, (right) images retrieved by  $\ell_0$ -cut pursuit with simple merge steps ( $\ell_0$ -CPm). Last line: (left) rasterized population density of Paris area, (right) simplified map obtained by  $\ell_0$ -CPm: 69% of variance explained with 1.2% of contours perimeter.

Experiment	Phantom		Simulated	
	PSNR	time	PSNR	time
Noisy image	16.8	-	16.8	-
$\ell_0$ -CPm	<b>33.5</b>	<b>4.3</b>	<b>37.0</b>	4.6
CRF15	32.6	8.6	34.2	<b>4.0</b>
CRF30	33.3	25.3	34.8	11.4
$TV_{0.5}$	32.2	16.4	33.6	18.0

Figure 13: PSNR at convergence and time to converge in seconds for the four algorithms as well as the noisy image for the first two denoising experiments.

758 **4. Conclusion.** We proposed two algorithms to minimize functions penalized respectively  
759 by the total variation and by the total boundary size. They computationally exploit the fact  
760 that for sufficiently large regularization coefficients, the solution is typically piecewise constant  
761 with a small number of pieces, corresponding to a coarse partition. This is a consequence of the  
762 fact that, in the discrete setting, both the total variation and total boundary size penalize the  
763 size of the support of the gradient: indeed, functions with sparse gradients tend to have a small  
764 number of distinct level sets, which are moreover connected. The sparsity that is optimized is  
765 thus not exactly the same as the sparsity which is exploited computationally, although both  
766 are related.

767 By constructing a sequence of approximate solutions that are themselves piecewise constant  
768 with a small number of pieces, the proposed algorithms operate on reduced problems that can  
769 be solved efficiently, and perform only graph cuts on the original graph, which are thus the  
770 remaining bottleneck for further speed-ups. Like all working-set algorithms, the cut pursuit  
771 variants are not competitive if the solution has too many connected level-sets.

772 In the convex case, cut pursuit outperforms all proximal methods for deblurring images  
773 with simple solutions. For denoising with a ROF energy, it outperforms the parametric maxflow  
774 approach when computing sequences of solutions for different regularization strengths. In  
775 the  $\ell_0$  case, our algorithm can find a better solution in a shorter time than the non-convex  
776 continuous relaxation approach as well as the approach based on  $\alpha$ -expansions. Furthermore,  
777 while the performance of the latter hinges critically on setting an appropriate number of level-  
778 sets in advance, cut pursuit needs no such parametrization.

779 Future developments will consider the case of Lovász extensions of other symmetric sub-  
780 modular functions [4] and to the multivariate case. It would also be interesting to determine  
781 the conditions under which the alternating scheme presented in E.1 provides a globally optimal  
782 solution of (13), as it would be a necessary step in order to prove approximation guarantees  
783 to the solution of  $\ell_0$ -cut pursuit itself.

784

## REFERENCES

- 785 [1] G. AUBERT, M. BARLAUD, O. FAUGERAS, AND S. JEHAN-BESSON, *Image segmentation using active*  
 786 *contours: calculus of variations or shape gradients?*, SIAM Journal on Applied Mathematics, 63  
 787 (2003), pp. 2128–2154.
- 788 [2] F. BACH, *Learning with submodular functions: a convex optimization perspective*, Foundations and Trends  
 789 in Machine Learning, 6 (2013), pp. 145–373.
- 790 [3] F. BACH, R. JENATTON, J. MAIRAL, AND G. OBOZINSKI, *Optimization with sparsity-inducing penalties*,  
 791 Foundations and Trends in Machine Learning, 4 (2012), pp. 1–106.
- 792 [4] F. R. BACH, *Shaping level sets with submodular functions*, in Advances in Neural Information Processing  
 793 Systems, 2011, pp. 10–18.
- 794 [5] L. BAR, T. F. CHAN, G. CHUNG, M. JUNG, N. KIRYATI, R. MOHIEDDINE, N. SOCHEN, AND L. A.  
 795 VESE, *Mumford and Shah model and its applications to image segmentation and image restoration*,  
 796 in Handbook of Mathematical Methods in Imaging, Springer, 2011, pp. 1095–1157.
- 797 [6] R. BELLMAN, *A note on cluster analysis and dynamic programming*, Mathematical Biosciences, 18 (1973),  
 798 pp. 311–312.
- 799 [7] K. BLEAKLEY AND J.-P. VERT, *The group fused Lasso for multiple change-point detection*, arXiv preprint  
 800 arXiv:1106.4199, (2011).
- 801 [8] J. BORWEIN AND A. S. LEWIS, *Convex analysis and nonlinear optimization: theory and examples*,  
 802 Springer Science & Business Media, 2010.
- 803 [9] Y. BOYKOV AND V. KOLMOGOROV, *An experimental comparison of min-cut/max-flow algorithms for*  
 804 *energy minimization in vision.*, IEEE Transactions on Pattern Analysis and Machine Intelligence, 26  
 805 (2004), pp. 1124–1137.
- 806 [10] Y. BOYKOV, O. VEKSLER, AND R. ZABIH, *Fast approximate energy minimization via graph cuts*, IEEE  
 807 Transactions on Pattern Analysis and Machine Intelligence, 23 (2001), pp. 1222–1239.
- 808 [11] X. BRESSON, S. ESEDOĞLU, P. VANDERGHEYNST, J.-P. THIRAN, AND S. OSHER, *Fast global mini-*  
 809 *mization of the active contour/snake model*, Journal of Mathematical Imaging and Vision, 28 (2007),  
 810 pp. 151–167.
- 811 [12] A. CHAMBOLLE, V. CASELLES, D. CREMERS, M. NOVAGA, AND T. POCK, *An introduction to total*  
 812 *variation for image analysis*, in Theoretical foundations and numerical methods for sparse recovery,  
 813 De Gruyter, 2010, pp. 263–340.
- 814 [13] A. CHAMBOLLE AND J. DARBON, *On total variation minimization and surface evolution using parametric*  
 815 *maximum flows*, International Journal of Computer Vision, 84 (2009), pp. 288–307.
- 816 [14] A. CHAMBOLLE AND T. POCK, *A first-order primal-dual algorithm for convex problems with applications*  
 817 *to imaging*, Journal of Mathematical Imaging and Vision, 40 (2011), pp. 120–145.
- 818 [15] T. CHAN, S. ESEDOĞLU, F. PARK, AND A. YIP, *Recent developments in total variation image restoration*,  
 819 in Mathematical Models of Computer Vision, Springer Verlag, 2005, pp. 17–31.
- 820 [16] T. F. CHAN AND L. A. VESE, *Active contours without edges*, IEEE Transactions on Image Processing,  
 821 10 (2001), pp. 266–277.
- 822 [17] V. CHANDRASEKARAN, B. RECHT, P. A. PARRILO, AND A. S. WILLSKY, *The convex geometry of linear*  
 823 *inverse problems*, Foundations of Computational mathematics, 12 (2012), pp. 805–849.
- 824 [18] S. CHEN, C. F. COWAN, AND P. M. GRANT, *Orthogonal least squares learning algorithm for radial basis*  
 825 *function networks*, IEEE Transactions on Neural Networks, 2 (1991), pp. 302–309.
- 826 [19] L. CONDAT, *A direct algorithm for 1D total variation denoising*, IEEE Signal Processing Letters, 20  
 827 (2013), pp. 1054–1057.
- 828 [20] W. DINKELBACH, *On nonlinear fractional programming*, Management Science, 13 (1967), pp. 492–498.
- 829 [21] B. EFRON, T. HASTIE, I. JOHNSTONE, R. TIBSHIRANI, ET AL., *Least angle regression*, The Annals of  
 830 statistics, 32 (2004), pp. 407–499.
- 831 [22] M. EFROYMSON, *Multiple regression analysis*, Mathematical methods for digital computers, 1 (1960),  
 832 pp. 191–203.
- 833 [23] N. EL-ZEHIRY AND L. GRADY, *Discrete optimization of the multiphase piecewise constant Mumford-Shah*  
 834 *functional*, in Energy Minimization Methods in Computer Vision and Pattern Recognition, Springer,  
 835 2011, pp. 233–246.
- 836 [24] N. EL-ZEHIRY, P. SAHOO, AND A. ELMAGHRABY, *Combinatorial optimization of the piecewise constant*

- 837 Mumford-Shah functional with application to scalar/vector valued and volumetric image segmentation,  
 838 Image and Vision Computing, 29 (2011), pp. 365–381.
- 839 [25] N. Y. EL-ZEHIRY AND A. ELMAGHRABY, *Brain MRI tissue classification using graph cut optimization of*  
 840 *the Mumford–Shah functional*, in Proceedings of the International Vision Conference of New Zealand,  
 841 2007, pp. 321–326.
- 842 [26] J. FRIEDMAN, T. HASTIE, AND R. TIBSHIRANI, *Regularization paths for generalized linear models via*  
 843 *coordinate descent*, Journal of Statistical Software, 33 (2010), pp. 1–22.
- 844 [27] B. FULKERSON, A. VEDALDI, AND S. SOATTO, *Class segmentation and object localization with superpixel*  
 845 *neighborhoods*, in Proceedings of the International Conference on Computer Vision, IEEE, October  
 846 2009, pp. 670–677.
- 847 [28] N. FUSCO, *An overview of the Mumford-Shah problem*, Milan Journal of Mathematics, 71 (2003), pp. 95–  
 848 119.
- 849 [29] D. GEMAN AND G. REYNOLDS, *Constrained restoration and the recovery of discontinuities*, IEEE Trans-  
 850 actions on Pattern Analysis & Machine Intelligence, 14 (1992), pp. 367–383.
- 851 [30] D. GOLDFARB AND W. YIN, *Parametric maximum flow algorithms for fast total variation minimization*,  
 852 SIAM Journal on Scientific Computing, 31 (2009), pp. 3712–3743.
- 853 [31] Z. HARCHAOU, A. JUDITSKY, AND A. NEMIROVSKI, *Conditional gradient algorithms for norm-*  
 854 *regularized smooth convex optimization*, Mathematical Programming, 152 (2015), pp. 75–112.
- 855 [32] H. ISHIKAWA, *Exact optimization for Markov random fields with convex priors*, IEEE Transactions on  
 856 Pattern Analysis and Machine Intelligence, 25 (2003), pp. 1333–1336.
- 857 [33] M. JAGGI, *Revisiting Frank-Wolfe: projection-free sparse convex optimization*, in Proceedings of the 30th  
 858 International Conference on Machine Learning, 2013, pp. 427–435.
- 859 [34] S. JEGELKA, F. BACH, AND S. SRA, *Reflection methods for user-friendly submodular optimization*, in  
 860 Advances in Neural Information Processing Systems, 2013, pp. 1313–1321.
- 861 [35] N. A. JOHNSON, *A dynamic programming algorithm for the fused lasso and  $\ell_0$ -segmentation*, Journal of  
 862 Computational and Graphical Statistics, 22 (2013), pp. 246–260.
- 863 [36] M. KASS, A. WITKIN, AND D. TERZOPOULOS, *Snakes: Active contour models*, International Journal of  
 864 Computer Vision, 1 (1988), pp. 321–331.
- 865 [37] P. KOHLI AND P. H. TORR, *Efficiently solving dynamic Markov random fields using graph cuts*, in  
 866 International Conference on Computer Vision (ICCV), vol. 2, IEEE, 2005, pp. 922–929.
- 867 [38] V. KOLMOGOROV, T. POCK, AND M. ROLINEK, *Total variation on a tree*, SIAM Journal on Imaging  
 868 Sciences, 9 (2016), pp. 605–636.
- 869 [39] V. KOLMOGOROV AND R. ZABIH, *What energy functions can be minimized via graph cuts?*, IEEE Trans-  
 870 actions on Pattern Analysis and Machine Intelligence, 26 (2004), pp. 147–159.
- 871 [40] J. J. KORMYLO AND J. M. MENDEL, *Maximum likelihood detection and estimation of Bernoulli-Gaussian*  
 872 *processes*, IEEE Transactions on Information Theory, 28 (1982), pp. 482–488.
- 873 [41] K. KUMAR AND F. BACH, *Active-set methods for submodular optimization*, arXiv preprint  
 874 arXiv:1506.02852, (2015).
- 875 [42] Y. G. LECLERC, *Constructing simple stable descriptions for image partitioning*, International journal of  
 876 computer vision, 3 (1989), pp. 73–102.
- 877 [43] G. P. LEONARDI AND I. TAMANINI, *On minimizing partitions with infinitely many components*, Annali  
 878 dell’Università di Ferrara, 44 (1998), pp. 41–57.
- 879 [44] D. A. LORENZ AND T. POCK, *An inertial forward-backward algorithm for monotone inclusions*, Journal  
 880 of Mathematical Imaging and Vision, 51 (2014), pp. 311–325.
- 881 [45] S. MALLAT AND Z. ZHANG, *Adaptive time-frequency decomposition with matching pursuits*, in Time-  
 882 Frequency and Time-Scale Analysis, Proceedings of the IEEE-SP International Symposium, IEEE,  
 883 1992, pp. 7–10.
- 884 [46] D. MUMFORD AND J. SHAH, *Optimal approximations by piecewise smooth functions and associated vari-*  
 885 *ational problems*, Communications on pure and applied mathematics, 42 (1989), pp. 577–685.
- 886 [47] D. NEEDELL AND J. A. TROPP, *CoSaMP: Iterative signal recovery from incomplete and inaccurate*  
 887 *samples*, Applied and Computational Harmonic Analysis, 26 (2009), pp. 301–321.
- 888 [48] S. NEGAHBAN, B. YU, M. J. WAINWRIGHT, AND P. K. RAVIKUMAR, *A unified framework for high-*  
 889 *dimensional analysis of  $m$ -estimators with decomposable regularizers*, in Advances in Neural Informa-  
 890 tion Processing Systems, 2009, pp. 1348–1356.

- 891 [49] Y. NESTEROV, *Gradient methods for minimizing composite objective function*, tech. report, Université  
892 catholique de Louvain, Center for Operations Research and Econometrics (CORE), 2007.
- 893 [50] F. NIELSEN AND R. NOCK, *Optimal interval clustering: Application to Bregman clustering and statistical*  
894 *mixture learning*, Signal Processing Letters, 21 (2014), pp. 1289–1292.
- 895 [51] M. NIKOLOVA, M. K. NG, AND C.-P. TAM, *Fast nonconvex nonsmooth minimization methods for image*  
896 *restoration and reconstruction*, IEEE Transactions on Image Processing, 19 (2010), pp. 3073–3088.
- 897 [52] G. OBOZINSKI, B. TASKAR, AND M. JORDAN, *Multi-task feature selection*, Statistics Department, UC  
898 Berkeley, Technical report 743, (2006).
- 899 [53] P. OCHS, A. DOSOVITSKIY, T. BROX, AND T. POCK, *On iteratively reweighted algorithms for nonsmooth*  
900 *nonconvex optimization in computer vision*, SIAM Journal on Imaging Sciences, 8 (2015), pp. 331–372.
- 901 [54] S. OSHER AND J. A. SETHIAN, *Fronts propagating with curvature-dependent speed: algorithms based on*  
902 *Hamilton-Jacobi formulations*, Journal of Computational Physics, 79 (1988), pp. 12–49.
- 903 [55] J.-C. PICARD AND H. D. RATLIFF, *Minimum cuts and related problems*, Networks, 5 (1975), pp. 357–370.
- 904 [56] T. POCK AND A. CHAMBOLLE, *Diagonal preconditioning for first order primal-dual algorithms in convex*  
905 *optimization*, in Proceeding of the International Conference on Computer Vision (ICCV), IEEE, 2011,  
906 pp. 1762–1769.
- 907 [57] H. RAGUET, J. FADILI, AND G. PEYRÉ, *A generalized forward-backward splitting*, SIAM Journal on  
908 Imaging Sciences, 6 (2013), pp. 1199–1226.
- 909 [58] H. RAGUET AND L. LANDRIEU, *Preconditioning of a generalized forward-backward splitting and applica-*  
910 *tion to optimization on graphs*, SIAM Journal on Imaging Sciences, 8 (2015), pp. 2706–2739.
- 911 [59] N. RAO, P. SHAH, AND S. WRIGHT, *Forward-backward greedy algorithms for atomic norm regularization*,  
912 IEEE Transactions on Signal Processing, 63 (2015), pp. 5798–5811.
- 913 [60] R. T. ROCKAFELLAR, *Convex analysis*, Princeton University Press, 1970.
- 914 [61] V. ROTH AND B. FISCHER, *The group-lasso for generalized linear models: uniqueness of solutions and*  
915 *efficient algorithms*, in Proceedings of the 25th international conference on Machine learning, ACM,  
916 2008, pp. 848–855.
- 917 [62] L. I. RUDIN, S. OSHER, AND E. FATEMI, *Nonlinear total variation based noise removal algorithms*,  
918 Physica D: Nonlinear Phenomena, 60 (1992), pp. 259 – 268, [https://doi.org/http://dx.doi.org/10.](https://doi.org/http://dx.doi.org/10.1016/0167-2789(92)90242-F)  
919 [1016/0167-2789\(92\)90242-F](https://doi.org/http://dx.doi.org/10.1016/0167-2789(92)90242-F), <http://www.sciencedirect.com/science/article/pii/016727899290242F>.
- 920 [63] L. A. SHEPP AND B. F. LOGAN, *The Fourier reconstruction of a head section*, IEEE Transactions on  
921 Nuclear Science, 21 (1974), pp. 21–43.
- 922 [64] C. SOUSSEN, J. IDIER, D. BRIE, AND J. DUAN, *From Bernoulli–Gaussian deconvolution to sparse signal*  
923 *restoration*, IEEE Transactions on Signal Processing, 59 (2011), pp. 4572–4584.
- 924 [65] R. SZELISKI, R. ZABIH, D. SCHARSTEIN, O. VEKSLER, V. KOLMOGOROV, A. AGARWALA, M. TAPPEN,  
925 AND C. ROTHER, *A comparative study of energy minimization methods for Markov random fields*, in  
926 Proceeding of the European Conference in Computer Vision (ECCV), Springer, 2006, pp. 16–29.
- 927 [66] I. TAMANINI AND G. CONGEDO, *Optimal segmentation of unbounded functions*, Rendiconti del Seminario  
928 Matematico della Università di Padova, 95 (1996), pp. 153–174.
- 929 [67] Y.-H. R. TSAI AND S. OSHER, *Total variation and level set methods in image science*, Acta Numerica,  
930 14 (2005), pp. 509–573.
- 931 [68] L. A. VESE AND T. F. CHAN, *A multiphase level set framework for image segmentation using the*  
932 *Mumford and Shah model*, International Journal of Computer Vision, 50 (2002), pp. 271–293.
- 933 [69] Y.-X. WANG, J. SHARPBACK, A. SMOLA, AND R. J. TIBSHIRANI, *Trend filtering on graphs*, Journal of  
934 Machine Learning Research, 17 (2016), pp. 1–41.
- 935 [70] T. ZHANG, *Adaptive forward-backward greedy algorithm for sparse learning with linear models*, in Advances  
936 in Neural Information Processing Systems, 2009, pp. 1921–1928.
- 937 [71] H. ZOU, *The adaptive lasso and its oracle properties*, Journal of the American Statistical Association,  
938 101 (2006), pp. 1418–1429.

**Appendix A. The total variation as an atomic gauge.** It is well known that the total variation is the Lovász extension of the submodular function  $F : B \mapsto w(B, B^c)$  [2, chap. 6.2]. The *base polytope* associated with  $F$  is the set  $\mathcal{B}_F \doteq \{s \in \mathbb{R}^n \mid s(B) \leq F(B), B \subset V, s(V) = F(V)\}$ , where  $s(B) \doteq \sum_{i \in B} s_i$ . For any submodular function  $F$  such that  $F(\emptyset) = F(V) = 0$ , which is true in particular for all symmetric submodular functions, the Lovász extension  $\gamma_F$  is a *gauge function* which is the *support function*<sup>13</sup> of  $\mathcal{B}_F$ :  $\gamma_F(x) = \max_{s \in \mathcal{B}_F} \langle s, x \rangle$  and its *polar gauge* is the gauge of  $\mathcal{B}_F$  [4]. The total variation is thus a gauge function and its polar gauge is  $\text{TV}^\circ$  with

$$\text{TV}^\circ(s) = \begin{cases} \max_{\emptyset \subsetneq B \subsetneq V} \frac{s(B)}{w(B, B^c)} & \text{if } s(V) = 0 \\ +\infty & \text{else.} \end{cases}$$

Chandrasekaran et al. [17] have recently introduced the concept of *atomic gauge*. Given a closed set  $\mathcal{A} \subset \mathbb{R}^n$  whose elements are called *atoms*, the associated atomic gauge is the gauge  $\gamma_{\mathcal{A}}$  of the convex hull  $C_{\mathcal{A}}$  of  $\mathcal{A} \cup \{0\}$ , i.e.  $\gamma_{\mathcal{A}}(x) \doteq \inf\{t \mid x \in tC_{\mathcal{A}}\}$ . The polar gauge is the support function of  $\mathcal{A} \cup \{0\}$ , that is  $\gamma_{\mathcal{A}}^\circ(s) = \sup_{a \in \mathcal{A} \cup \{0\}} \langle a, s \rangle$ . Given that  $\mathcal{A} \subset \mathbb{R}^n$ , using Caratheodory's theorem, we have that

$$\gamma_{\mathcal{A}}(x) = \inf \left\{ \sum_{a \in \mathcal{A}} c_a \mid \forall a \in \mathcal{A}, c_a \geq 0, \sum_{a \in \mathcal{A}} c_a a = x \right\}.$$

939 Regularizing with an atomic gauge thus favors solutions that are sparse combinations of  
 940 atoms, which motivated the use of algorithms that exploit the sparsity of the solution com-  
 941 putationally [33, 59]. It is clear from previous definitions that Lovász extensions are atomic  
 942 gauges. In particular the total variation is the atomic gauge associated with the set of atoms  
 943  $\mathcal{A} = \{w(B, B^c)^{-1} \mathbf{1}_B + \mu \mathbf{1}_V\}_{B \notin \{\emptyset, V\}, \mu \in \mathbb{R}}$  or equivalently the set  $\mathcal{A}' = \{\frac{1}{2}w(B, B^c)^{-1}(\mathbf{1}_B -$   
 944  $\mathbf{1}_{B^c}) + \mu' \mathbf{1}_V\}_{B \notin \{\emptyset, V\}, \mu' \in \mathbb{R}}$ . Expressing solutions to problem regularized with the total varia-  
 945 tion as combinations of set indicators or cuts as we propose to do in this paper is thus very  
 946 natural from this perspective.

947 For the total variation, the Frank-Wolfe direction associated to  $s = -\nabla f(x)$  such that  
 948  $\langle s, \mathbf{1}_V \rangle = 0$  is

$$949 \quad (9) \quad \arg \max_{\xi: \text{TV}(\xi) \leq 1} \langle s, \xi \rangle = \arg \max_{\mathbf{1}_B: B \notin \{\emptyset, V\}} \frac{1}{w(B, B^c)} \langle s, \mathbf{1}_B \rangle,$$

950 since the maximizer is necessarily an extreme point of the set  $\{\xi \mid \text{TV}(\xi) \leq 1\}$  and therefore  
 951 among the atoms.

## 952 **Appendix B. Proof of Propositions 1 and 3.**

**Proposition 1.** For  $x \in \mathbb{R}^n$ , if we set  $S = S(x)$  then

$$Q'(x, \mathbf{1}_B) = \langle \nabla Q_S(x), \mathbf{1}_B \rangle + \lambda w_{S^c}(B, B^c).$$

Moreover if  $\langle \nabla f(x), \mathbf{1}_V \rangle = 0$  then

$$Q'(x, u_B) = (\gamma_B + \gamma_{B^c}) Q'(x, \mathbf{1}_B).$$

---

<sup>13</sup>See [60] for definitions of *gauge*, *polar gauge* and *support function* of a set.

*Proof.* For  $B \subset V$  we have that  $Q'(x, \mathbf{1}_B) = \langle \nabla Q_S(x), \mathbf{1}_B \rangle + \sup_{\epsilon \in \partial \text{TV}|_{S^c}(x)} \langle \epsilon, \mathbf{1}_B \rangle$ . This can be shown using the chain rule for subgradients that we have:

$$\partial \text{TV}|_{S^c}(x) = \left\{ \frac{1}{2} D^\top \delta \mid \delta_S = 0, \|\delta_{S^c}\|_\infty \leq 1, \forall (i, j) \in E, \delta_{ij} = -\delta_{ji} \right\},$$

953 with  $D \in \mathbb{R}^{2m \times n}$  the matrix whose only non-zero entries are  $D_{(i,j),i} = w_{ij}$  and  $D_{(i,j),j} = -w_{ij}$   
 954 for all  $(i, j) \in E$ , and with the notations  $\delta_S \in \mathbb{R}^{2m}$  and  $\delta_{S^c} \in \mathbb{R}^{2m}$  for the vectors whose entries  
 955 are equal to those of  $\delta$  respectively on  $S$  and  $S^c$  and equal to zero otherwise.

Therefore if  $\epsilon = \frac{1}{2} D^\top \delta_{S^c}$  then

$$\langle \epsilon, \mathbf{1}_B \rangle = \left\langle \frac{1}{2} \delta_{S^c}, D \mathbf{1}_B \right\rangle = \frac{1}{2} \sum_{(i,j) \in S^c} \delta_{ij} w_{ij} ([\mathbf{1}_B]_i - [\mathbf{1}_B]_j).$$

The supremum is reached for  $\delta_{ij} = \text{sign}([\mathbf{1}_B]_i - [\mathbf{1}_B]_j)$  for  $(i, j) \in S^c$ , so that  $\sup_{\epsilon \in \partial \text{TV}|_{S^c}(x)} \langle \epsilon, \mathbf{1}_B \rangle = w_{S^c}(B, B^c)$ .

For the second statement, we have that

$$Q'(x, u_B) = \langle \nabla Q_S(x), u_B \rangle + \sup_{\epsilon \in \partial \text{TV}|_{S^c}(x)} \langle \epsilon, u_B \rangle.$$

Letting  $g = \nabla Q_S(x)$ , and since  $\langle \nabla f, \mathbf{1} \rangle = 0$ , we have  $\langle g, \mathbf{1} \rangle = 0$ . Consequently  $\langle g, \mathbf{1}_{B^c} \rangle = \langle g, \mathbf{1} - \mathbf{1}_B \rangle = -\langle g, \mathbf{1}_B \rangle$ , and we have:

$$\langle g, u_B \rangle = \gamma_B \langle g, \mathbf{1}_B \rangle - \gamma_{B^c} \langle g, \mathbf{1}_{B^c} \rangle = (\gamma_B + \gamma_{B^c}) \langle g, \mathbf{1}_B \rangle.$$

956 Similarly,  $\langle \epsilon, u_B \rangle = \left\langle \frac{1}{2} \delta_{S^c}, D u_B \right\rangle = \frac{1}{2} \gamma_B \langle \delta_{S^c}, D \mathbf{1}_B \rangle - \frac{1}{2} \gamma_{B^c} \langle \delta_{S^c}, D \mathbf{1}_{B^c} \rangle = \frac{1}{2} (\gamma_B + \gamma_{B^c}) \langle \delta_{S^c}, D \mathbf{1}_B \rangle$   
 957 because  $D \mathbf{1}_B = -D \mathbf{1}_{B^c}$ . Taking the supremum over  $\epsilon$  then proves the result.  $\blacksquare$

958 **Proposition 3.** *We have  $x = \arg \min_{z \in \mathbb{R}^n} Q(z)$  if and only if  $\min_{B \subset V} Q'(x, \mathbf{1}_B) = 0$  and*  
 959  *$Q'(x, \mathbf{1}_V) = 0$ .*

960 *Proof.* ( $\Rightarrow$ ) If  $x$  is the solution of problem (1), the directional derivative of  $Q$  along  
 961 any direction must be nonnegative, which implies that  $Q'(x, \mathbf{1}_B) \geq 0$  for all  $B$ . But since  
 962  $\min_{B \subset V} Q'(x, \mathbf{1}_B) \leq Q'(x, \mathbf{1}_\emptyset) = 0$ , this proves the first part. Then since  $w(V, \emptyset) = 0$  we  
 963 have  $Q'(x, \mathbf{1}_V) = \langle \nabla Q_S(x), \mathbf{1}_V \rangle$ , and, in fact, since all elements of the subgradient of  $\text{TV}|_{S^c}$   
 964 are orthogonal to  $\mathbf{1}_V$  we also have  $Q'(x, -\mathbf{1}_V) = -\langle \nabla Q_S(x), \mathbf{1}_V \rangle$ . So  $0 \leq Q'(x, -\mathbf{1}_V) =$   
 965  $-Q'(x, \mathbf{1}_V) \leq 0$ .

966 ( $\Leftarrow$ ) Conversely we assume that  $\min_{B \subset V} Q'(x, \mathbf{1}_B) = 0$  and  $Q'(x, \mathbf{1}_V) = 0$ .  
 967 Since  $Q'(x, \mathbf{1}_V) = 0$  and since  $w_{S^c}(V, \emptyset) = 0$  we have  $\langle \nabla Q_S(x), \mathbf{1}_V \rangle = 0$ . Now, for any set  
 968  $A$  which is a maximal connected component of  $G|_{S^c} \doteq (V, S^c)$ , we also have  $w_{S^c}(A, A^c) =$   
 969  $0$  so that  $0 \leq Q'(x, \mathbf{1}_A) = \langle \nabla Q_S(x), \mathbf{1}_A \rangle$  but the same holds for the complement  $A^c$  and  
 970  $\langle \nabla Q_S(x), \mathbf{1}_A \rangle + \langle \nabla Q_S(x), \mathbf{1}_{A^c} \rangle = \langle \nabla Q_S(x), \mathbf{1}_V \rangle = 0$  so that  $\langle \nabla Q_S(x), \mathbf{1}_A \rangle = 0$ .

971 As a consequence the capacities of the graph  $G_{flow}$  defined in (6) of the article are such  
 972 that, for any set  $A$  which is a maximal connected component of  $G|_{S^c}$ , we have

$$974 \quad (10) \quad \sum_{i \in \nabla_+ \cap A} c_{si} = \sum_{i \in \nabla_- \cap A} c_{it}.$$

975 Then since  $Q'(x, \mathbf{1}_\emptyset) = 0$  and since  $\min_{B \subset V} Q'(x, \mathbf{1}_B) = 0$  it is a minimizing argument.  
 976 The characterization of the steepest partition as a minimal cut then guarantees that there  
 977 exists a minimal cut in  $G_{flow}$  which does not cut any edge in  $S^c$  and isolates the source or the  
 978 sink from the rest of the graph. Given equality (10), the set of minimal cuts are the cuts that  
 979 remove indifferently for each maximal connected component  $A$  either all edges  $\{(s, i)\}_{i \in A}$  or  
 980 the edges  $\{(i, t)\}_{i \in A}$ .

981 A consequence of the max-flow/min-cut duality is that to this cut corresponds a maximal  
 982 flow  $e \in \mathbb{R}^{2m}$  in  $G_{flow}$ . This flow is such that it is saturated at the minimal cut, and we thus  
 983 have  $e_{si} = c_{si}$  for all  $i \in \nabla_+$  and  $e_{it} = c_{it}$  for all  $i \in \nabla_-$ , again because of equation (10).

984 Writing flow conservation yields

$$985 \quad (11) \quad \begin{cases} e_{si} + \sum_{j \in N_i} (e_{ji} - e_{ij}) = 0 & \forall i \in \nabla_+ \\ -e_{it} + \sum_{j \in N_i} (e_{ji} - e_{ij}) = 0 & \forall i \in \nabla_-, \end{cases}$$

986 with  $N_i = \{j \mid (i, j) \in S^c\}$ .

987 By replacing  $e_{si}$  and  $e_{it}$  by their value, the flow conservation (11) at node  $i$  rewrites

$$988 \quad \nabla_i Q_S(x) + \sum_{j \in N_i} \lambda w_{ij} \delta_{ij} = 0$$

$$989 \quad (12) \quad \nabla_i Q_S(x) + \frac{1}{2} \sum_{j \in N_i} \lambda w_{ij} (\delta_{ij} - \delta_{ji}) = 0,$$

990 with  $\delta_{ij} = \frac{e_{ji} - e_{ij}}{\lambda w_{ij}}$  for  $(i, j) \in S^c(x)$  and  $\delta_{ij} = \delta_{ji} = 0$  for all edges  $(i, j) \in S(x)$ . The  
 991 flow  $e$  must respect the capacity at all edges and hence  $0 \leq e_{ij} \leq c_{ij} = \lambda w_{ij}$  for all edges  
 992 in  $S^c(x)$ . Since the flow is maximal, only one of  $e_{ij}$  or  $e_{ji}$  is non zero. Hence  $\delta$  we naturally  
 993 have  $\delta_{ij} = -\delta_{ji}$ , and  $|\delta_{ij}| \leq 1$ . But we can rewrite (12) as  $\nabla Q_S(x) = \frac{1}{2} \lambda D^T \delta$  with  $\delta_S = 0$   
 994 and  $\|\delta_{S^c}\| \leq 1$  with  $D$  as in the characterization of the subgradient of  $\text{TV}|_{S^c}$  which shows that  
 995  $-\frac{1}{\lambda} \nabla Q_S(x) \in \partial \text{TV}|_{S^c}(x)$  thus that  $0 \in \partial Q(x)$ , and finally that  $x$  minimizes  $Q$ . ■

996 **Remark:** We proved Proposition 3 using directly the flow formulation and the simplest  
 997 possible arguments. It is also possible to prove the result more directly using more abstract  
 998 results. We actually used the fact that  $x$  is a minimum of  $Q$  if and only if, for  $S = S(x)$ ,  
 999  $-\frac{1}{\lambda} \nabla Q_S(x) \in \partial \text{TV}|_{S^c}(x)$ . But it is possible to give another representation of  $\partial \text{TV}|_{S^c}(x)$   
 1000 using that the subgradient of a gauge  $\gamma$  at  $x$  is  $\partial \gamma(x) = \{s \mid \langle x, s \rangle = \gamma(x), \gamma^\circ(s) \leq 1\}$ .  
 1001 Indeed, for  $\gamma = \text{TV}$ , the set  $\{\gamma^\circ(s) \leq 1\}$  is simply the submodular polytope  $\mathcal{P}_F$  of  $F : B \mapsto$   
 1002  $w(B, B^c)$ . As a result  $\partial \text{TV}|_{S^c}(x) = \{s \in \mathbb{R}^n \mid \langle s, x \rangle = 1, \forall B, s(B) \leq w_{S^c}(B, B)\}$ . But  
 1003 having that  $\min_{B \subset V} \langle \nabla Q_S(x), \mathbf{1}_B \rangle + \lambda w_{S^c}(B, B^c) = 0$  is equivalent to having  $-\frac{1}{\lambda} \nabla Q_S(x) \in$   
 1004  $\{s \in \mathbb{R}^n \mid \forall B, s(B) \leq w_{S^c}(B, B)\}$ . There thus just remains to show that  $\langle \nabla Q_S(x), x \rangle =$   
 1005  $\text{TV}(x)$ . Let  $\Pi_S$  denote the set of maximal connected components of  $G|_{S^c} = (V, S^c)$ , so that  
 1006 we have  $x = \sum_{A \in \Pi_S} c_A \mathbf{1}_A$ . Since  $w_{S^c}(V, \emptyset) = 0$ , we have  $0 = Q'(x, \mathbf{1}_V) = \langle \nabla Q_S(x), \mathbf{1}_V \rangle$ .  
 1007 Similarly for  $A \in \Pi_S$ , we have  $w_{S^c}(A, A^c) = 0$ , which entails that  $\langle \nabla Q_S(x), \mathbf{1}_A \rangle \geq 0$ . But  
 1008 then  $-\langle \nabla Q_S(x), \mathbf{1}_A \rangle = \langle \nabla Q_S(x), \mathbf{1}_{A^c} \rangle \geq 0$  also, which proves  $\langle \nabla Q_S(x), \mathbf{1}_A \rangle = 0$ . Finally by  
 1009 linearity  $\langle \nabla Q_S(x), x \rangle = \sum_{A \in \Pi_S} c_A \langle \nabla Q_S(x), \mathbf{1}_A \rangle = 0 = \text{TV}|_{S^c}(x)$  which proves the result.

### Appendix C. Theoretical results for cut pursuit with a non-convex function $f$ .

This appendix discusses how the propositions of Section 2 can be extended to the case of non-convex functions  $f$ .

It relies on the fact that notions of directional derivative and subgradient can be extended to non-convex functions. This presents some difficulties in general and different definitions of directional derivatives and subgradient have been introduced by Dini, by Clarke, and by Michel and Penot [8, Chap. 6.1]. These extended subgradients do not behave like usual subgradients in general and some of the rules of the calculus of subgradient are no longer valid. Fortunately, for so-called *regular* functions, that is functions for which the Dini, Clarke and Michel-Penot subgradient all coincide, the usual subgradient calculus applies [8, Chap. 6.2]. In particular, a function  $Q = f + g$  with  $f$  *strictly differentiable*<sup>14</sup> and  $g$  convex is *regular* at any point  $x$  of the interior of its domain and  $\partial Q(x) = \nabla f(x) + \partial g(x)$ , where  $\partial \cdot$  denotes here the generalized subgradient for regular function (that coincides with the usual subgradient if the function is convex). This is in particular true for  $g = \text{TV}$ . As a consequence, the proof of Propositions 1 and 2 only require  $f$  to be *strictly differentiable*. Similarly, Proposition 3 no longer holds as stated because the first order subgradient condition is not sufficient for optimality, but we still have

**Proposition 7.** *For  $Q = f + \text{TV}$  with  $f$  strictly differentiable,  $0 \in \partial Q(x)$  if and only if  $\min_{B \subset V} Q'(x, \mathbf{1}_B) = 0$  and  $Q'(x, \mathbf{1}_V) = 0$ .*

*Proof.* Since  $f$  is strictly differentiable,  $Q$  is regular so that the usual subgradient calculus applies and the proof is the same as that of Proposition 3. ■

If  $f$  is non-convex, solving the subproblem on the reduced graph is more difficult, even if only a local minimum is sought. To extend Algorithm 1 to the non-convex setting, it seems appropriate to assume that reoptimizing on the reduced graph (at the end of the main loop) yields a vector  $x_{\Pi^t}$  which is a local minimum of the reduced objective and such that  $Q(x_{\Pi^t}) < Q(x_{\Pi^{t-1}})$ .

With that modification Proposition 4 remains true, and instead of Proposition 5, we have that the algorithm converges in a finite number of iterations to a point  $x^*$ , which is a local minimum of  $Q$  in the subspace  $\text{span}(\Pi)$  and satisfies  $0 \in \partial Q(x^*)$ . This is not sufficient in general for  $x^*$  to be a minimum of  $Q$ . However, if  $T(x^*)$  denotes the tangent cone of  $Q$  at  $x^*$ , that is  $T(x^*) := \{h \in \mathbb{R}^n \mid Q'(x^*, h) = 0\}$  (since there are no directions such that  $Q'(x^*, h) < 0$ ), and if  $\nabla^2 f(x^*)$  denotes the Hessian of  $f$  at  $x^*$  then, by standard arguments, the condition  $\forall h \in T(x^*), \langle h, \nabla^2 f(x^*)h \rangle > 0$  is sufficient to guarantee that  $x^*$  is a local minimum of  $Q$ .

**Appendix D. Computation of the Frank-Wolfe direction.** The computation of the Frank-Wolfe direction defined in (9) requires to optimize a ratio of combinatorial functions. More precisely, it requires to solve

$$\max_{B \notin \{\emptyset, V\}} \frac{N(B)}{D(B)} \quad \text{with} \quad N(B) \doteq -\langle \nabla f(x), \mathbf{1}_B \rangle, \quad \text{and} \quad D(B) \doteq w(B, B^c).$$

---

<sup>14</sup> $f$  is strictly differentiable at  $x$  if there exists  $\varphi \in \mathbb{R}^n$  such that  $\forall h \in \mathbb{R}^n, \lim_{y \rightarrow x, t \downarrow 0} \frac{f(y+th) - f(y)}{t} = \langle \varphi, h \rangle$ .

1044 Given that  $B \mapsto \frac{N(B)}{D(B)}$  it is the ratio of a supermodular function (in fact a modular function)  
 1045 and a nonnegative submodular function, it can be maximized efficiently by Algorithm 7 as  
 1046 proved in Proposition 8.

---

**Algorithm 7**      Computation      of

---

$\max_A N(A)/D(A)$

---

**Initialization:**  $\lambda_0 = 1, \lambda_{-1} = 0, t = 0$   
**while**  $\lambda_t \neq \lambda_{t-1}$  **do**  
 1047      $\mathcal{S}_t \leftarrow \text{Arg max}_{A \subset V} N(A) - \lambda_t D(A)$   
        $A_t \leftarrow \text{arg min}_{A \subset \mathcal{S}_t} D(A)$   
        $\lambda_{t+1} \leftarrow \frac{N(A_t)}{D(A_t)}$   
        $t \leftarrow t + 1$   
**end while**  
**return**  $A_t$

---

1048 **Proposition 8.** *The sequence  $(\lambda_t)_t$  generated by Algorithm 7 is monotonically increasing and*  
 1049 *converges in a finite number of iterations to  $\max_{\emptyset \subsetneq A \subset V} \frac{N(A)}{D(A)}$ .*

1050 *Proof.* As the maximum of a finite number non-increasing linear functions of a scalar  
 1051 argument, the function  $\varphi : \lambda \mapsto \max_{A \subset V} [N(A) - \lambda D(A)]$  is a non-increasing, continuous,  
 1052 piecewise linear convex function. It is also non negative because  $N(\emptyset) - \lambda D(\emptyset) = 0$ . It is  
 1053 immediate to check that  $\lambda^* := \min\{\lambda \mid \varphi(\lambda) = 0\} = \max_{\emptyset \subsetneq A \subset V} \frac{N(A)}{D(A)}$ . At each iteration, if  
 1054  $\varphi(\lambda_t) \neq 0$ , we must have  $\lambda_{t+1} > \lambda_t$ , because the function  $\lambda \mapsto N(A_t) - \lambda D(A_t)$  is strictly  
 1055 positive for  $\lambda = \lambda_t$  and equal to 0 for  $\lambda = \lambda_{t+1}$ . Moreover by construction, the sets  $A_t$  are all  
 1056 distinct, as long as  $\varphi(\lambda_t) \neq 0$ . As a consequence we must reach  $\varphi(\lambda_T) = 0$  after a finite number  
 1057 of iterations  $T$ . At the end of the algorithm,  $\varphi(\lambda_T) = 0$  entails that  $\forall A \subset V, N(A) \leq \lambda_T D(A)$ ,  
 1058 which entails that for all  $A \neq \emptyset, D(A)^{-1} N(A) \leq \lambda_T = D(A_{T-1})^{-1} N(A_{T-1})$ . This shows that  
 1059  $\lambda_T = \max_{\emptyset \subsetneq A \subset V} \frac{N(A)}{D(A)}$ . This concludes the proof. The choice of taking the maximizer with  
 1060 smallest value of  $D(A)$  on line 4 of the algorithm is not key to convergence of the algorithm,  
 1061 but aims at computing the right-derivative which maximizes the step size in  $\lambda$ . ■

1062 Note that this algorithm is closely related to the algorithm of [20] to maximize a ratio of  
 1063 functions, and in fact applies to any functions  $N$  and  $D$ ; but the minimization of the function  
 1064  $(A \mapsto \lambda D(A) - N(A))$  can be done in polynomial here because, since  $D$  and  $N$  are respectively  
 1065 sub- and super-modular, their difference is submodular. Moreover, when  $D$  is submodular and  
 1066  $N$  is modular, the number of iterations may be bounded by  $d$ , because the algorithm may be  
 1067 reinterpreted as the divide-and-conquer algorithm to maximise submodular functions over the  
 1068 submodular polytope [2, p.160] (for the general case, it may only be bounded in general by  
 1069  $2^d$ ).

1070 **Appendix E. Details of the derivation, technical elements and proofs for  $\ell_0$  cut pursuit.**  
 1071

1072 **E.1. Splitting step.** Since in Section 3.1.1 the problem of finding an optimal binary cut  
 1073 of the component  $A_j$  is decoupled from the same problem on other components and leads to  
 1074 formulation (8), we discuss the splitting step for the case of the optimal binary cut of the

1075 initial component  $V$ .

1076 In the same way that we defined the steepest binary cut in cut pursuit for the convex  
1077 formulation, we define the optimal binary partition  $(B, B^c)$  of  $V$  such that  $Q$  optimized over  
1078  $\text{span}(\mathbf{1}_B, \mathbf{1}_{B^c})$  is as small as possible. Ideally, we should impose that  $B$  and  $B^c$  have a single  
1079 connected component each, because as argued in section 2.3, it does not make sense to impose  
1080 that  $x_i$  should have the same values in different connected components. However, since this  
1081 constraint is too difficult to enforce, we first ignore it and address it later with post-processing  
1082 described in Section E.1.3. Note however that the penalization of the perimeter of the boundary  
1083 between  $B$  and  $B^c$  should strongly discourage the choice of sets  $B$  with many connected  
1084 components.

1085 **E.1.1. Optimal binary cut with alternating minimization.** Since  $\Gamma(h\mathbf{1}_B + h'\mathbf{1}_{B^c}) =$   
1086  $\Gamma(\mathbf{1}_B) = w(B, B^c)$ , and ignoring the connectedness constraint, the corresponding optimization  
1087 problem is of the form

$$1088 \quad (13) \quad \min_{B \subset V} \min_{h, h' \in \mathbb{R}} \sum_{i \in B} f_i(h) + \sum_{i \in B^c} f_i(h') + \lambda w(B, B^c).$$

1089 This problem is a priori hard to solve in general, because  $B \mapsto \min_{h, h' \in \mathbb{R}} f(h\mathbf{1}_B + h'\mathbf{1}_{B^c})$  is  
1090 not a submodular function. However, when  $h, h'$  are fixed, the assumption that  $f$  is separable  
1091 entails that  $B \mapsto f(h\mathbf{1}_B + h'\mathbf{1}_{B^c})$  is a modular function, so that the objective can be optimized  
1092 with respect to  $B$  by solving a max-flow problem. Similarly as for the flow problem (6) we  
1093 define the flow graph  $G_{flow} = (V \cup \{s, t\}, E_{flow})$  whose edge set and capacities are defined by:

$$1094 \quad (14) \quad E_{flow} = \begin{cases} (s, i), \forall i \in \nabla_+, & \text{with } c_{si} = f_i(h) - f_i(h'), \\ (i, t), \forall i \in \nabla_-, & \text{with } c_{it} = f_i(h') - f_i(h), \\ (i, j), \forall (i, j) \in E, & \text{with } c_{ij} = \lambda w_{ij}, \end{cases}$$

1095 where  $\nabla_+ \doteq \{i \in V \mid f_i(h) > f_i(h')\}$  and  $\nabla_- \doteq V \setminus \nabla_+$ .

1096 The regularity and convexity of  $f$  with respect to  $h$  and  $h'$  guarantee that the objective  
1097 can be minimized efficiently with respect to these variables. As suggested by [11] or [24],  
1098  $\psi(B, h, h') = \sum_{i \in B} f_i(h) + \sum_{i \in B^c} f_i(h') + \lambda w(B, B^c)$  can be efficiently minimized by alterna-  
1099 tively minimizing with respect to  $B$  and  $(h, h')$  separately.

1100 **E.1.2. Proof of convergence of the alternating minimization scheme.** The alternating  
1101 scheme used to compute the optimal binary cut provide a local minimum of  $\psi(B, h, h') =$   
1102  $\sum_{i \in B} f_i(h) + \sum_{i \in B^c} f_i(h') + \lambda w(B, B^c)$  with the following assumptions:

- 1103 • (A0): the functions  $f_i$  are continuous,
- 1104 • (A1): the solution of  $\min_{(h, h')} \psi(h, h', B)$  exists and is unique for all sets  $B$
- 1105 • (A2): the minimizer with respect to  $B$  of  $\psi(h_A, h'_A, B)$  is unique for all  $A$ .

1106 Note that (A1) holds if for example all functions  $f_i$  are strictly convex. (A2) can be shown to  
1107 hold with probability one if  $f_i$  is appropriately random, for example if  $f_i(\cdot) = (\cdot - x_i)^2$  with  
1108  $x_i$  drawn i.i.d. from a continuous distribution, which corresponds to our case of interest.

1109 **Proposition 9.** *Assuming that the assumptions (A0), (A1) and (A2) hold, the alternate min-*  
1110 *imization scheme converges in a finite number of iterations to a local minimum of  $\psi(h, h', B)$  in*

1111 *the sense that there exists a neighborhood  $\mathcal{N}_B$  of  $(h_B, h'_B)$  such that for all  $(h, h', A) \in \mathcal{N}_B \times 2^V$ ,*  
 1112 *we have  $\psi(h, h', A) \geq \psi(h_B, h'_B, B)$ .*

1113 *Proof.* Let  $\psi(B) = \min_{h, h'} \psi(h, h', B)$ . By construction and with assumption (A1), the  
 1114 sequence  $(\psi(B^t))_t$  is strictly decreasing until minimization with respect to either  $(h, h')$  or  $B$   
 1115 yields no progress, i.e. until a partial minimum with respect to both blocks is attained. Since  
 1116 the set  $2^V$  is finite, the algorithm must converge in a finite number of iterations.

1117 The point  $B$  attained must be a local minimum in the sense above: indeed for any set  $A$   
 1118 different than  $B$ , we must have  $\phi(h_B, h'_B, B) < \phi(h_B, h'_B, A)$  because the algorithm stopped  
 1119 (which excludes  $\phi(h_B, h'_B, B) > \phi(h_B, h'_B, A)$ ) and because an equality is excluded by (A2).  
 1120 But then by assumption (A0),  $\phi$  is continuous with respect to  $(h, h')$  so that in a neighborhood  
 1121  $\mathcal{N}_B$  of  $(h_B, h'_B)$  we must have  $\phi(h, h', A)$  sufficiently close to  $\phi(h_B, h'_B, A)$  for the inequality  
 1122 characterizing a local minimum to hold. ■

1123 **E.1.3. From binary cut to partition in connected components.** Like the working set  
 1124 algorithm proposed for the total variation,  $\ell_0$ -cut pursuit recursively splits the components of  
 1125 the current partition  $\Pi$ . The sets  $B$  and  $B^c$  obtained as a solution of (13) are not necessarily  
 1126 connected sets, but splitting  $B$  and  $B^c$  into their connected components and assigning each  
 1127 connected component its own value obviously does not change the contour perimeter  $\Gamma$  and  
 1128 can only decrease  $f$ . Given the collection of connected components  $A_1, \dots, A_k$  of  $B$  and  $B^c$   
 1129 we therefore set  $x = h_1 \mathbf{1}_{A_1} + \dots + h_k \mathbf{1}_{A_k}$  with  $h_j$  the minimizer of  $h \mapsto \sum_{i \in A_j} f_i(h)$ . Note  
 1130 that each  $h_i$  could possibly be computed in parallel given the separability of  $f$ .

1131 **E.2. Implementation.** As in the convex case,  $\ell_0$ -cut pursuit maintains a current partition  
 1132  $\Pi$  that is recursively split and computes optimal values for each of its components. It is  
 1133 comprised of three main steps: the splitting of the current partition, the computation of the  
 1134 connected components and their values, and a potential merging step, when necessary.

1135 **E.2.1. Splitting.** For each component an optimal binary partition  $(B, B^c)$  is obtained by  
 1136 solving (13) as described in section E.1.1: we alternatively minimize the objective with respect  
 1137 to  $B$  and with respect to  $(h, h')$  until either  $B$  does not change or a maximum number of  
 1138 iterations is reached. In practice, the algorithm converges in 3 steps most of the time. The  
 1139 choice of an appropriate initialization for  $B$  is non-trivial. Since the problem in which  $\lambda = 0$   
 1140 is often simpler, and can in a number of cases be solved analytically, we chose to use that  
 1141 solution to initialize our alternating minimization scheme. Indeed, for  $\lambda = 0$ , and when  $f$  is a  
 1142 squared Euclidean distance  $f : x \mapsto \|x - x_0\|_2^2$  the objective of (13) is the same as the objective  
 1143 of one-dimensional  $k$ -means with  $k = 2$ ; in this particular setting, the problem reduces to a  
 1144 change-point analysis problem, and an exact solution can be computed efficiently by dynamic  
 1145 programming [6]. This can be generalized to the case of Bregman divergences and beyond [50].

1146 As described in section E.1.3, the partition  $\Pi$  is updated by computing its connected com-  
 1147 ponents after it is split by  $(B, B^c)$ . Subroutine 2 gives the procedure algorithmically.  
 1148 It is important to note that this is the only operation that involves the original graph  $G$ , and  
 1149 hence will be the computational bottleneck of the algorithm. Fortunately since  $f$  is separable,  
 1150 this procedure can be performed on each component in parallel.

1151

**E.2.2. Simple merge..** This backward step consists of checking for each neighboring components  $A$  and  $B$  in  $\Pi$  whether merging them into a single component decreases the energy. If we denote  $\Pi_-(A, B)$  the partition obtained by merging  $A$  and  $B$ , the corresponding decrease in energy  $\delta_-(A, B)$  is

$$\delta_-(A, B) = f(x_\Pi) - f(x_{\Pi_-(A, B)}) + \lambda w(A, B),$$

1152 with  $\Pi_-(A, B) \doteq \Pi \setminus \{A, B\} \cup \{A \cup B\}$ .

1153 The exact implementation of the while loop described in Algorithm 5 is in fact based on a  
1154 priority-queue. The value  $\delta_-(A, B)$  is computed for each neighboring components, and stored  
1155 in the priority queue. Each pair that provides a nonnegative decrease is merged, and  $\delta_-$  is  
1156 updated for the neighbors of  $A$  and  $B$  to reflect the change in value and graph topology. This  
1157 operation scales with the size of the reduced graph only, and therefore can be performed effi-  
1158 ciently for problems in which the partition  $\Pi$  does not get too large.

1159

1160 **E.2.3. Merge-resplit..** This more complex backward step, already described in 3.1.3 is  
1161 significantly computationally more intensive as it is performed on the edges of the full graph,  
1162 by contrast with the simple merge which only considers the edges of the reduced graph. As a  
1163 consequence, while all potential simple merge steps can be precomputed and performed based  
1164 on a priority queue by merging first the pair of components yielding the largest decrease in  
1165 objective value, it would be too computationally heavy in the merge-resplit case and we thus  
1166 perform boundary changes only once for each pair of neighbors in the graph  $\mathcal{E}$ . The pseudocode  
1167 of the procedure is detailed in subroutine 4

1168 **E.2.4. Other algorithmic variants.** We discuss here the relevance of constructing more  
1169 greedy algorithms and of variants to tackle the problem in which the total boundary size is  
1170 constrained instead of penalized.

1171 It would have been theoretically possible to implement a more greedy version of  $\ell_0$  cut  
1172 pursuit in which one performs a single forward step (corresponding to splitting a single region)  
1173 at a time or a single backward step at a time by maintaining a global priority queue and one  
1174 greedily chooses the most beneficial, but the overhead costs would have been prohibitive.

1175 The  $\ell_0$  cut pursuit algorithms constructed in Section 3 is a greedy algorithm to solve  
1176 a formulation in which the total boundary size is penalized and not constrained. It is worth  
1177 pointing out that trying to solve directly the constrained case seems difficult: indeed, designing  
1178 algorithms that are only based on forward steps (e.g., in the style of OMP, OLS, etc) might  
1179 not succeed, because of the dependence between the cuts that need to be introduced to form  
1180 the final solution. Based on similar ideas as the ones used in  $\ell_0$ -cut pursuit, we designed and  
1181 tested an algorithm generalizing the FoBa algorithm [70]. The obtained algorithm tended to  
1182 remain trapped in bad local minima and yielded solutions that were much worse than the ones  
1183 based on the penalized formulation.

1184 **E.3. Convergence to a local minimum of the generalized MP problem.** We now prove  
1185 the local optimality of the solution provided by Algorithm 5.

1186 **Proposition 10.** *If assumption (A0) holds, then the  $\ell_0$  cut pursuit algorithm provides in a*

1187 finite number of iterations a partition  $\Pi = (A_1, \dots, A_n)$  such that  $x_\Pi \doteq \arg \min_{z \in \text{span}(\Pi)} Q(z)$   
 1188 is a local minimum of  $Q$ .

1189 *Proof.* The fact that  $f$  is separable ensures that  $x_\Pi$  can be minimized separately over each  
 1190 connected component. We denote  $x_{A_i} \in \arg \min_z \sum_{i \in A_i} f_i(z)$ .

1191 We denote  $\Pi^t$  the partition at iteration  $t$ , and  $x_\Pi^t$  the associated solution. We first prove  
 1192 that the sequence  $Q(x_\Pi^t)$  is strictly decreasing. Indeed if the stopping criterion for the algorithm  
 1193 is not met, then there exists at least one component  $A_j$  which is not saturated, i.e. such  
 1194 that there exists a binary partitions  $B \subsetneq A_j$  such that  $\min_{h, h'} \sum_{i \in B} f_i(h) + \sum_{i \in B^c} f_i(h') +$   
 1195  $\lambda w(B, B^c) < \sum_{i \in A_j} f_i(x_{A_j})$ . Consequently this component will be split in the next partition  
 1196 to yield a strict decrease of the objective function  $Q$ , at least equal to the one provided by the  
 1197 minimizing arguments  $(h, h')$ . Since the set of all partition is a finite set, the algorithm stops  
 1198 in a finite number of steps. We now prove that the partition  $\Pi$  attained when the algorithm  
 1199 stops is such that the corresponding variable  $x_\Pi$  is a local minimum of  $Q$ . Let  $\mathcal{E}$  be the set of  
 1200 pairs of adjacent components of  $\Pi$ . We can assume that  $x_A \neq x_B$  for any  $(A, B) \in \mathcal{E}$ . If it  
 1201 is not the case we replace  $\Pi$  by the partition in which such components are merged, without  
 1202 changing  $x_\Pi$ . Consequently there exists  $\delta_1 > 0$  such that  $|x_A - x_B| > \delta_1$  for any  $(A, B) \in \mathcal{E}$ .

1203 As all edge weights are assumed strictly non negative we have that  $w_{\min} = \min_{(i,j) \in E} w_{i,j} >$   
 1204  $0$ . Since (A0) states that  $f$  is continuous, there exists an Euclidean ball centered at  $x_\Pi$  and of  
 1205 radius  $\delta_2$  in which all elements are strictly greater than  $f(x_\Pi) - \min_{(i,j) \in E} w_{i,j}$ .

1206 We now prove by contradiction that  $x_\Pi$  is a local minimum of  $Q$ . Let  $x'$  be an element  
 1207 of the euclidian ball  $\mathcal{B}$  centered at  $x_\Pi$  and of radius  $\min(\frac{1}{3}\delta_1, \delta_2)$  such that  $Q(x') < Q(x_\Pi)$ .  
 1208 We can first recognize that since the values of  $x_\Pi$  associated to each connected component  
 1209 differs by at least  $\delta$ ,  $x'$  cannot have two connected components of  $\Pi$  sharing a common value.  
 1210 Consequently the boundary perimeter can only increase  $\Gamma(x') \geq \Gamma(x_\Pi)$ .

1211 If we first assume that  $\Gamma(x') = \Gamma(x_\Pi)$ , then  $x'$  must be piecewise constant with respect  
 1212 to  $\Pi$ , and be such that  $f(x') < f(x_\Pi)$ , which is a contradiction with the definition of  $x_\Pi$ .  
 1213 We must then assume that  $\Gamma(x') > \Gamma(x)$ . Since the smallest increment in  $\Gamma$  is  $w_{\min}$ , we have  
 1214  $\Gamma(x') \geq \Gamma(x) + w_{\min}$ . Since the radius of  $\mathcal{B}$  is smaller than  $\delta_2$ , we have that  $f(x) \geq f(x_\Pi) - w_{\min}$ ,  
 1215 and consequently  $Q(x_\Pi) \geq Q(x)$ , which is a contradiction.

1216 **Appendix F. CRF formulation and number of quantization levels.** In this appendix, we  
 1217 report the performance of  $\ell_0$ -cut pursuit and  $\alpha$ -expansions for different numbers of quantization  
 1218 levels for denoising an image. The regularization strength is chosen by cross-validation to  
 1219 maximize the PSNR.

1220 The fact that  $\ell_0$ -CPm does not rely on an a priori quantized level leads to overall good  
 1221 performance, with significantly faster computation times. By contrast, the running time for  
 1222 the  $\alpha$ -expansions based algorithms has a complexity which empirically grows linearly with the  
 1223 number of classes, and the performance whether measured in terms of the original objective  
 1224 or in PSNR does not increase monotonically as a function of the number of classes.

1225 Plotting the corresponding PSNRs shows that the smaller local minima of the objective  
 1226 found correlate well with gains in PSNR, and that the corresponding gains can be quite  
 1227 substantial as illustrated as well in Table 13. The fact that the level of performance for CRF  
 1228 is highly sensitive to the exact number of classes is a shortcoming of the method, especially  
 1229 given its computational cost.

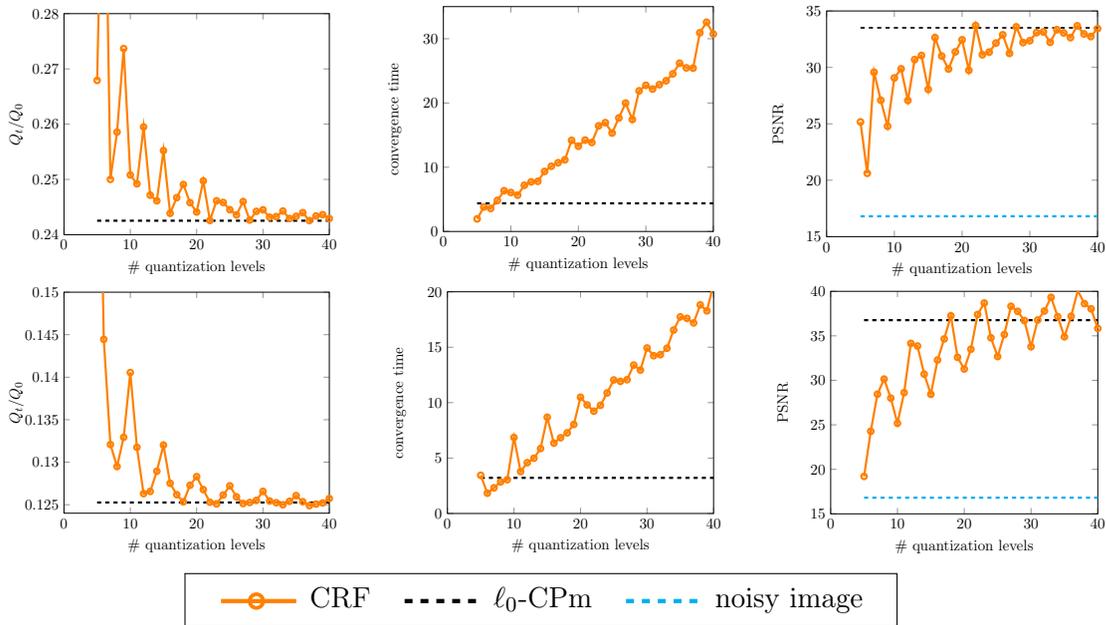


Figure 14: **Behavior of the  $\alpha$ -expansion based algorithm on CRF formulations for different number of quantization levels** for the phantom (top) and the simulated data (bottom) averaged on 10 denoising experiments: (left) ratio between the energy  $Q$  at convergence and the energy at time 0, (middle) running time, (right) corresponding PSNRs. The two algorithms represented are  $\alpha$ -expansions (CRF) for a varying number of quantization levels and  $\ell_0$ -CPm.