

A Study of Nesterov's Scheme for Lagrangian Decomposition and MAP Labeling

Anonymous CVPR submission

Paper ID 1549

Abstract

We study the MAP-labeling problem for graphical models by optimizing a dual problem obtained by Lagrangian decomposition. In this paper, we focus specifically on Nesterov's optimal first-order optimization scheme for non-smooth convex programs, that has been studied for a range of other problems in computer vision and machine learning in recent years. We show that in order to obtain an efficiently convergent iteration, this approach should be augmented with a dynamic estimation of a corresponding Lipschitz constant, leading to a runtime complexity of $O(\frac{1}{\epsilon})$ in terms of the desired precision ϵ . Additionally, we devise a stopping criterion based on a duality gap as a sound basis for competitive comparison and show how to compute it efficiently. We evaluate our results using the publicly available Middlebury database and a set of computer generated graphical models that highlight specific aspects, along with other state-of-the-art methods for MAP-inference.

1. Introduction

Problem We consider the problem of computing the most likely configuration x for a given graphical model, i.e. a distribution $p_{\mathcal{G}}(x; \theta) \propto \exp(-E_{\mathcal{G}}(\theta, x))$. We use the following standard notation [17]:

Let $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ be an undirected graph, where \mathcal{V} is a finite set of its nodes and $\mathcal{E} \subset \mathcal{V} \times \mathcal{V}$ is a set of edges. Let also \mathcal{X}_v , $v \in \mathcal{V}$ be a finite set of labels. The set $\mathcal{X} = \otimes_{v \in \mathcal{V}} \mathcal{X}_v$, where \otimes denotes the Cartesian product, will be called *labeling set* and its elements $x \in \mathcal{X}$ *labelings*. Thus each labeling is a collection $(x_v : v \in \mathcal{V})$ of labels. To shorten notation we will use x_{uv} for a pair of labels (x_u, x_v) and \mathcal{X}_{uv} for $\mathcal{X}_u \times \mathcal{X}_v$. Functions of the form $\theta_v : \mathcal{X}_v \rightarrow \mathbb{R}$, $v \in \mathcal{V}$, and $\theta_{uv} : \mathcal{X}_{uv} \rightarrow \mathbb{R}$, $uv \in \mathcal{E}$, are called *unary* and *pairwise potentials*, respectively. The collection of all potentials will be denoted by θ .

The problem to compute the most likely labeling x (*MAP labeling problem*) amounts to minimizing the energy function

$$\min_{x \in \mathcal{X}} E_{\mathcal{G}}(\theta, x) = \min_{x \in \mathcal{X}} \left\{ \sum_{v \in \mathcal{V}} \theta_v(x_v) + \sum_{uv \in \mathcal{E}} \theta_{uv}(x_{uv}) \right\}. \quad (1)$$

Background and Motivation Problem (1) is known to be NP-complete in general for graphs with cycles. We will concentrate mainly on the linear programming (LP) relaxation of the problem originally proposed by Schlesinger [14] – see [18] for a recent review.

Schlesinger [14] analysed also the dual LP as an upper bound of an integer solution and proposed two minimization algorithms: the DAG algorithm and a diffusion algorithm (cf. [18]). These algorithms decrease the value of the dual LP monotonically but do not attain its optima in general, since they can be interpreted as (block-)coordinate descent and thus can get stuck due to the non-smoothness of the dual objective.

Another algorithm, known as TRW-S, was proposed by Kolmogorov [5]. This algorithm computes the same fixed points as the diffusion algorithm and generalizes it by considering arbitrary sub-trees of the initial graph as elementary subproblems in contrast to separate nodes and neighboring edges in the diffusion. An alternative sub-gradient based scheme for dual function minimization was proposed by Komodakis [7]. Such sub-gradient iterations are guaranteed to compute the optimum of the dual function but have two drawbacks: i) no efficient convergence rate is backed by theory – to the best of our knowledge no improvement has been established with respect to the general convergence estimate $O(\frac{1}{\epsilon^2})$ [9] – and ii) absence of a stopping criterion that is sound from the optimization viewpoint.

Disadvantages of both approaches (TRW-S and sub-gradient iteration) are mainly caused by the non-smoothness of the dual objective. To overcome this problem, smoothing of the objective was proposed in a series of papers [2, 3, 12, 19]. However, questions concerning the worst-case complexity bound and theoretically sound stop-

ping conditions have remained open.

In a most recent work [4], smoothing of the dual objective was addressed with Nesterov’s optimal first-order optimization scheme. We will show in this paper, however, that without carefully modifying the generic scheme [9], the resulting complexity bound $O(\sqrt{|\mathcal{V}|}/\epsilon)$ is too loose for almost any real problem instance.

Contribution Our contribution is two-fold:

(i) We propose an algorithm for solving the dual LP problem with a guaranteed complexity of $O(\frac{1}{\epsilon})$ oracle calls (evaluations of the function or its gradient).

(ii) We formulate and analyse a general method for constructing an upper bound for algorithms maximizing the dual LP objective. The method is used in turn to devise a sound stopping criterion based on the duality gap.

Algorithm (i) is based on smoothing the dual objective and applying the optimal first-order optimization scheme by Nesterov [10]. Our approach is similar to the method described in [4] but differs from it in essential technical details:

(a) Instead of using a *fixed* Lipschitz constant for a gradient step of the algorithm, we *adaptively* estimate this constant during the iteration. This leads to a significantly smaller number of outer iterations of the algorithm necessary for convergence, at the cost of a few more oracle calls (no more than 4 on average) in the inner loop of the iteration. Overall, our algorithm is much faster.

(b) Instead of *static* selection of a smoothing value, we select it *dynamically*, which usually gives a significant speed-up.

Method (ii) can be applied to any iterative scheme as soon as an approximate, but not necessarily feasible, primal solution can be computed. Hence, this contribution should be of wider interest. We use this method to define and evaluate a stopping condition for our algorithm. In contrast, no stopping criterion was specified in [4].

For the sake of clarity of our presentation, we consider here the special case of grid-graphs \mathcal{G} , mainly because the benchmark [15] conforms to this setting. Although none of our results is restricted to this special case, the quantitative evaluation is, of course. A generalization is mostly straightforward, and we add specific comments where issues might arise. All proofs of theoretical results are available as supplementary material, due to the space restriction.

2. Description of the Algorithm

Decomposition and Relaxation Our approach is based on the dual decomposition framework which was proposed for energy minimization by [16] and later on analysed by [7]. Let $\mathcal{G}^i = (\mathcal{V}^i, \mathcal{E}^i)$, $i = 1, 2$, be two *acyclic* subgraphs of the *master graph* \mathcal{G} . Let $\mathcal{V}^1 = \mathcal{V}^2 = \mathcal{V}$,

$\mathcal{E}^1 \cup \mathcal{E}^2 = \mathcal{E}$ and $\mathcal{E}^1 \cap \mathcal{E}^2 = \emptyset$ (e.g., \mathcal{E}^1 contains all horizontal edges of \mathcal{G} and \mathcal{E}^2 all vertical ones if \mathcal{G} is a grid graph). Then the overall energy becomes the sum of the energies corresponding to these sub-graphs,

$$\begin{aligned} E_{\mathcal{G}}(\theta, x) &= \sum_{i=1}^2 \sum_{v \in \mathcal{V}^i} \theta_v^i(x_v) + \sum_{uv \in \mathcal{E}^i} \theta_{uv}^i(x_{uv}) \\ &= E_{\mathcal{G}^1}(\theta^1, x) + E_{\mathcal{G}^2}(\theta^2, x), \end{aligned} \quad (2)$$

provided $\theta_{uv}^i = \theta_{uv}$, $uv \in \mathcal{E}^i$, $i = 1, 2$ and $\theta_v^1(x_v) + \theta_v^2(x_v) = \theta_v(x_v)$, $\forall v \in \mathcal{V}, x_v \in \mathcal{X}_v$. The latter condition can be represented in a parametric way as $\theta_v^1(x_v) = \frac{\theta_v(x_v)}{2} + \lambda_v(x_v)$ and $\theta_v^2(x_v) = \frac{\theta_v(x_v)}{2} - \lambda_v(x_v)$, $v \in \mathcal{V}$, $x_v \in \mathcal{X}_v$, where $\lambda_v(x_v) \in \mathbb{R}$. Thus we consider θ^i as a function of λ and obviously have

$$\min_{x \in \mathcal{X}} E_{\mathcal{G}}(\theta, x) \geq \max_{\lambda} \sum_{i=1}^2 \min_{x \in \mathcal{X}} E_{\mathcal{G}^i}(\theta^i(\lambda), x). \quad (3)$$

It is well-known [6] that all collections (of arbitrary cardinality) of acyclic sub-graphs covering the master graph are equivalent, in the sense that they lead to the same lower bound as the one presented on the right-hand side of equation (3). It is also well-known that this lower bound is equal to the solution of the following linear programming problem:

$$\begin{aligned} \min_{\mu} \sum_{v \in \mathcal{V}} \sum_{x_v \in \mathcal{X}_v} \theta_v(x_v) \mu_v(x_v) + \sum_{uv \in \mathcal{E}} \sum_{x_{uv} \in \mathcal{X}_{uv}} \theta_{uv}(x_{uv}) \mu_{uv}(x_{uv}) \\ \text{s.t.} \quad \sum_{x_v \in \mathcal{V}} \mu_v(x_v) = 1, v \in \mathcal{V} \\ \sum_{x_v \in \mathcal{V}} \mu_{uv}(x_{uv}) = \mu_u(x_u), x_u \in \mathcal{X}_u, uv \in \mathcal{E} \\ \sum_{x_u \in \mathcal{V}} \mu_{uv}(x_{uv}) = \mu_v(x_v), x_v \in \mathcal{X}_v, uv \in \mathcal{E} \\ \mu_{uv}(x_{uv}) \geq 0, x_{uv} \in \mathcal{X}_{uv}, uv \in \mathcal{E}. \end{aligned} \quad (4)$$

This formulation is based on the overcomplete representation commonly used for discrete graphical models [17], in terms of relaxed indicator vectors μ constrained to the *local polytope* $\mathcal{L}(\mathcal{G})$, that is defined by the constraints of (4). It is well-known that $\mathcal{L}(\mathcal{G})$ constitutes an outer bound (relaxation) of the convex hull of all indicator vectors of labelings (marginal polytope; cf. [17]). Consequently, (4) simply reads $\min_{\mu \in \mathcal{L}(\mathcal{G})} \langle \theta, \mu \rangle$.

Problem Smoothing Consider a single summand on the right-hand side of (3). It can be expressed as inner product of the local potential vector θ^i and a correspondingly chosen binary indicator vector $\phi(x)$:

$$U^i(\lambda) := \min_{x \in \mathcal{X}} E_{\mathcal{G}^i}(\theta^i(\lambda), x) = \min_{x \in \mathcal{X}} \langle \theta^i(\lambda), \phi(x) \rangle. \quad (5)$$

Since this is a non-smooth function, the objective – right-hand side of (3) – is also non-smooth. Applying the well-known approximation of the min (or – max) function

by the log-exponential function (cf. [10, 13]) leads to the smooth version

$$\hat{U}_\rho^i(\lambda) = -\rho \log \sum_{x \in \mathcal{X}} \exp \langle -\theta^i(\lambda)/\rho, \phi(x) \rangle \quad (6)$$

with smoothing parameter ρ , that uniformly approximates U^i , that is

$$\hat{U}_\rho^i(\lambda) \leq U^i(\lambda) \leq \hat{U}_\rho^i(\lambda) + \rho \log |\mathcal{X}|. \quad (7)$$

Thus, for $\hat{U}_\rho = \sum_{i=1}^2 \hat{U}_\rho^i$ and $U = \sum_{i=1}^2 U^i$,

$$\hat{U}_\rho(\lambda) \leq U \leq \hat{U}_\rho(\lambda) + 2\rho \log |\mathcal{X}|. \quad (8)$$

We will call a gradient ∇f of a differentiable function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ Lipschitz continuous with Lipschitz constant L if

$$\|\nabla f(z) - \nabla f(w)\| \leq L\|z - w\|, \quad \forall z, w \in \mathbb{R}^n, \quad (9)$$

where $\|\cdot\|$ is the ℓ_2 -norm in \mathbb{R}^n .

Defining vectors $D\hat{U}_\rho^i(\lambda) \in \mathbb{R}^{\sum_{v \in \mathcal{V}} |\mathcal{X}_v|}$ by

$$D\hat{U}_\rho^i(\lambda)_{v, x_v} := \frac{\sum_{x \in \mathcal{X}(v, x_v)} \exp \langle -\theta^i(\lambda)/\rho, \phi(x) \rangle}{\exp(-\hat{U}_\rho^i(\lambda)/\rho)}, \quad (10)$$

where $\mathcal{X}(v, x_v) = \{x' \in \mathcal{X}: x'_v = x_v\}$, we have:

Lemma 1 (follows from Theorem 1 in [10]) *The function $\hat{U}_\rho(\lambda) = \sum_{i=1}^2 \hat{U}_\rho^i(\lambda)$ is well-defined and continuously differentiable at any $\lambda \in \mathbb{R}^{\sum_{v \in \mathcal{V}} |\mathcal{X}_v|}$. Moreover, this function is concave, and its gradient*

$$\nabla \hat{U}_\rho(\lambda) = D\hat{U}_\rho^1(\lambda) - D\hat{U}_\rho^2(\lambda) \quad (11)$$

is Lipschitz-continuous with constant $L_\rho = 2\frac{|\mathcal{V}|}{\rho}$.

This lemma is analogous to the "Computing Lipschitz" lemma in [4] with the significant difference that we consider the ℓ_2 -norm instead of the ℓ_1 -norm. Jojic et al. [4] inconsistently apply an algorithm based on the ℓ_2 -norm, however. Therefore, the role of the ℓ_1 -Lipschitz estimate (which reads $L_\rho = \frac{2}{\rho}$, see [4]) for the algorithm design remains unclear in [4].

Optimal First-Order Iterative Optimization It is known [9] that concave continuously differentiable (with Lipschitz constant L) functions can be maximized by iterative first-order optimization methods in $O(\sqrt{\frac{L}{\epsilon}})$ iterations, where ϵ determines the absolute precision of achieved objective value. Thus, by virtue of Lemma 1, the number of iterations can grow as $\sqrt{|\mathcal{V}|}$ with the size of a model, in the worst case.

Next, we present such an algorithm, omitting technical details which can be found in Nesterov's book ([9] p. 76) and in the original paper [8].

Algorithm 1 (Variant of Algorithm 2.2.6 in [9]) *In addition to the Lipschitz-constant L_ρ , we introduce variables $\gamma^t, \alpha^t, \omega \in \mathbb{R}$ and vectors $\lambda^t, v^t, y^t \in \mathbb{R}^{\sum_{v \in \mathcal{V}} |\mathcal{X}_v|}$. Superscript t indexes the iteration.*

1. Choose $\lambda^0 = v^0 \in \mathbb{R}^{\sum_{v \in \mathcal{V}} |\mathcal{X}_v|}$ and set $\gamma^0 = L_\rho$.

2. t -th iteration ($t \geq 0$):

(a) Compute $\hat{U}_\rho(\lambda^t)$ and $\nabla \hat{U}_\rho(\lambda^t)$.

(b) Find $\omega^t \leq L_\rho$ as small as possible, such that

$$\hat{U}_\rho(y^t) \geq \hat{U}_\rho(\lambda^t) + \frac{1}{2\omega^t} \|\nabla \hat{U}_\rho(\lambda^t)\|^2, \quad (12)$$

where $y^t = \lambda^t + \frac{1}{\omega^t} \nabla \hat{U}_\rho(\lambda^t)$.

(c) Compute $\alpha^t \in (0, 1)$ from $\omega^t(\alpha^t)^2 = (1 - \alpha^t)\gamma^t$ and set $\gamma^{t+1} = (1 - \alpha^t)\gamma^t$.

(d) Set $v^{t+1} = \frac{(1 - \alpha^t)\gamma^t v^t + \alpha^t \nabla \hat{U}_\rho(\lambda^t)}{\gamma^{t+1}}$.

(e) Choose $\lambda^{t+1} = \frac{\alpha^t \gamma^t v^{t+1} + \gamma^{t+1} y^t}{\gamma^t}$.

Lemma 2 [9] *Condition (12) is fulfilled for any $\omega^t \geq L_\rho$.*

Theorem 1 (modified Theorem 2.2.2 in [9]) *Algorithm 1 has the following bound on worst-case sub-optimality:*

$$\hat{U}_\rho^* - \hat{U}_\rho(\lambda^t) \leq \frac{L}{\left(1 + \frac{t}{2} \sqrt{\frac{L}{\omega^{*t}}}\right)^2} \|\lambda^0 - \lambda^*\|^2, \quad (13)$$

where \hat{U}_ρ^* and λ^* are optimal function and variable values and $\omega^{*t} = \max_{k \leq t} \omega^k$.

The estimate (13) shows that ω^t should be as small as possible. One possible way to achieve this is to perform exact linear search in the direction of $\nabla \hat{U}_\rho(\lambda^t)$ at each iteration of Algorithm 1, which is not particularly efficient however. A simple alternative is to set $\omega^t = L_\rho$ in view of Lemma 2, as done in [4]. Our experiments however show that the worst-case estimate of L_ρ according to Lemma 1 is quite loose and leads to a poor convergence rate.

Instead, we applied backtracking linear search ([1] p. 464, [11]), which consistently leads to speed-up factors up to 100 for our datasets. The analysis of backtracking linear search (see [11], eqn. (4.12), for a detailed proof) shows that for t iterations one needs no more than

$$N_k \leq 2 \left[1 + \frac{\ln d}{\ln u} \right] (t + 1) + \frac{1}{\ln u} \ln \frac{2uL}{dL_0} \quad (14)$$

oracle calls, where $d, u, L_0 \in \mathbb{R}$ are parameters of the search procedure. Using $d = u = L_0 = 2$, for example, each iteration of Algorithm 1 requires about 4 oracle calls on average (empirically, for our datasets, 3 or 4 oracle calls per iteration in most cases).

Selecting the Smoothing Parameter Inequality (8) and Lemma 1 show that selection of the smoothing value ρ is a trade-off between accuracy of the approximation and speed of the algorithm. The following lemma describes how to optimally select ρ for any algorithm \mathcal{A} that satisfies some conditions.

Lemma 3 Let \mathcal{A} be any algorithm depending on a smoothing parameter $\rho > 0$ with convergence rate $\hat{U}_\rho^* - \hat{U}_\rho(\lambda^t) \leq \frac{1}{\rho\tau(t)}$, where $\tau(t)$ is a monotonously non-decreasing function of the number of iterations t . Suppose that $U(\lambda) - \hat{U}_\rho(\lambda) \leq \rho\Delta$, for some $\Delta > 0$ and $\forall \rho, \lambda$, where \hat{U}_ρ is the smoothed objective function U . Let ϵ be the prescribed precision. Then, selection of the smoothing parameter ρ as

$$\rho = \frac{\epsilon}{2\Delta} \quad (15)$$

minimizes the worst-case bound on the number of iterations to achieve precision ϵ .

Our empirical results show $\omega^{*t} \propto \frac{1}{\rho}$ for ω^{*t} defined in Theorem 1. Thus, according to (13) and (8), this lemma can be directly applied to Algorithm 1, as done in [10] and [4]. In these papers, an upper bound $\Delta = 2 \log |\mathcal{X}|$ was used, leading to

$$\rho = \frac{\epsilon}{4 \log |\mathcal{X}|}. \quad (16)$$

This bound, however, can be rather loose in practice, that slows down convergence.

In contrast to this worst-case approach, we adapt ρ so as to allow for stronger smoothing in the initial and intermediate phase of the iteration, while still achieving the precision ϵ at convergence. We select ρ such that $\Delta \approx U(\lambda^0) - \hat{U}_\rho(\lambda^0) \lesssim \epsilon/2$, increasing ρ by the factor 2 if necessary. Usually, 3 to 6 computations of $\hat{U}_\rho(\lambda^0)$ suffice until $U(\lambda^0) - \hat{U}_\rho(\lambda^0) > \epsilon/2$. Such adaptive estimation of ρ leads to a speed up of the overall algorithm of order $2^2 \dots 2^5$. We check the inequality $U(\lambda^t) - \hat{U}_\rho(\lambda^t) < \epsilon/2$ during the iteration. If it does not hold (e.g. when convergence slows down close to the optima of \hat{U}_ρ), we decrease ρ by the factor 2.

3. Stopping Criterion

The stopping criterion we propose is based on a duality gap between the value of the primal LP, given by (4), and its dual $U(\lambda)$, given by right-hand side of (3). Since we optimize the dual problem and thus know its value, we focus in this section on estimating the value of the primal function, whose objective we will denote by P . We further denote by $\mathbb{R}_+(\mathcal{G}) = \mathbb{R}_+^{|\otimes_{v \in \mathcal{V}} \mathcal{X}_v| + |\otimes_{uv \in \mathcal{E}} \mathcal{X}_{uv}|}$ a nonnegative linear half-space containing the local polytope $\mathcal{L}(\mathcal{G})$. Finally, we denote the optimal primal value over the local polytope by $P^* = \min_{\mu \in \mathcal{L}(\mathcal{G})} P(\mu) = \min_{\mu \in \mathcal{L}(\mathcal{G})} \langle \theta, \mu \rangle$.

A typical issue for many algorithms which optimize a dual problem (3) is that, during the iteration, one can only get *infeasible* primal points $\tilde{\mu}$, that is $\tilde{\mu}$ does not satisfy the constraints of (4). In this connection, we propose to construct a mapping $\chi: \mathbb{R}_+(\mathcal{G}) \rightarrow \mathcal{L}(\mathcal{G})$ yielding primal feasible points, which enjoys the following properties:

Lemma 4 Let $\tilde{\mu}^t \in \mathbb{R}_+(\mathcal{G})$ be any sequence such that $P(\tilde{\mu}^t) \rightarrow P^*$. Let also $\min_{\mu \in \mathcal{L}(\mathcal{G})} \|\tilde{\mu}^t - \mu\| \rightarrow 0$. Then $P(\chi(\tilde{\mu}^t)) \rightarrow P^*$.

We define the shorthand $\mu' := \chi(\tilde{\mu})$ and the set $\mathcal{L}(\mathcal{G}, \mu'(\mathcal{V})) = \{\mu \in \mathcal{L}(\mathcal{G}) : \mu_v = \mu'_v, v \in \mathcal{V}\}$. A mapping χ as characterized by Lemma 4 can be constructed in the following two-steps way:

$$\mu''_v = \frac{\tilde{\mu}_v}{\sum_{x_v \in \mathcal{X}_v} \tilde{\mu}_v(x_v)}, v \in \mathcal{V}, \quad (17)$$

$$\mu' = \arg \min_{\mu \in \mathcal{L}(\mathcal{G}, \mu''(\mathcal{V}))} \langle \theta, \mu \rangle. \quad (18)$$

It is easy to see that problem (18) decomposes into $|\mathcal{E}|$ independent optimization problems (for each $uv \in \mathcal{E}$) of the form

$$\begin{aligned} \mu'_{uv} &= \arg \min_{\mu_{uv}} \sum_{x_{uv} \in \mathcal{X}_{uv}} \theta_{uv}(x_{uv}) \mu_{uv}(x_{uv}), \\ \text{s.t. } \sum_{x_v \in \mathcal{V}} \mu_{uv}(x_{uv}) &= \mu''_u(x_u), x_u \in \mathcal{X}_u \\ \sum_{x_u \in \mathcal{V}} \mu_{uv}(x_{uv}) &= \mu''_v(x_v), x_v \in \mathcal{X}_v \\ \mu_{uv}(x_{uv}) &\geq 0, x_{uv} \in \mathcal{X}_{uv}. \end{aligned} \quad (19)$$

Such linear programs are well-studied and known as *transportation problems*. Since the size of each individual problem is small, they can be easily solved by any appropriate method of linear programming.

We point out that the existence of a sequence $\tilde{\mu}^t$ satisfying the conditions of Lemma 4 is important for the theoretical properties of $\chi(\tilde{\mu}^t)$ to hold. But to compute $\chi(\tilde{\mu}^t)$, one only needs a subset of coordinates of the sequence, namely $\tilde{\mu}_v^t, v \in \mathcal{V}$. We show existence of a sequence $\tilde{\mu}_v^t$ by construction.

Theorem 2 When $\rho \rightarrow 0, t \rightarrow \infty$, for the sequence

$$\mu_v^{\rho,t} = \frac{D\hat{U}_\rho^1(\lambda^t)_v + D\hat{U}_\rho^2(\lambda^t)_v}{2}, v \in \mathcal{V} \quad (20)$$

a sequence $\tilde{\mu}^{\rho,t} \in \mathbb{R}_+(\mathcal{G})$ exists such that $\tilde{\mu}_v^{\rho,t} = \mu_v^{\rho,t}, v \in \mathcal{V}$, and $\tilde{\mu}^{\rho,t}$ satisfies the conditions of Lemma 4, namely $\forall \delta > 0 \exists \rho > 0: \exists t^* : \forall t > t^* \|P(\tilde{\mu}^{\rho,t}) - P^*\| < \delta$ and $\min_{\mu \in \mathcal{L}(\mathcal{G})} \|\tilde{\mu}^{\rho,t} - \mu\| < \delta$. Here λ^t is computed by Algorithm 1 for a given ρ , and $D\hat{U}_\rho^i(\lambda)_v, i = 1, 2$, are vectors with coordinates $D\hat{U}_\rho^i(\lambda)_{v,x_v}$ given by (10).

This theorem basically says that for ρ small enough, values $\mu_v^{\rho,t}$ plugged into formula (17) in place of $\tilde{\mu}_v$ would yield primal objective values which will converge with $t \rightarrow \infty$ to a value close enough to P^* .

4. Experiments

In our experiments we study different grid structured models with potentials of first and second order. Exemplarily we will discuss two of them. The first one is a synthetic model with 20×20 nodes, five labels and potential functions sampled uniformly from the interval $[0; 0.5]$ (corresponding plot in Figure 4 and top plots in Figures 1-3), the second is the Tsukuba stereo problem from the Middlebury MRF-Benchmark [15] (bottom plots in Figures 1-3).

We compare different variants of the Nesterov’s method (NEST) among each other and with standard methods, namely TRW-S [5], Norm-Product Belief-Propagation (NPBP) [2] and sub-gradient methods [7]. Thanks to the authors we can use their original code for TRW-S and NPBP. Since we compare different implementations of these methods, on the time axis we plot the number of oracle calls (function or gradient evaluations) instead of direct time measurements.

For the lower and upper bounds shown in our plots, we used values of the non-smooth dual objective U (see its definition after eq. (7)) and primal objective P , evaluated by means of (19), respectively.

Lipschitz Constant Estimation First we compare the performance of Nesterov’s method for different estimates of the Lipschitz constant. Adaptive selection of the Lipschitz constant leads to a significantly faster convergence than the fixed one. We also applied the calculation of the Lipschitz constant L_ρ as suggested in [4]. The top plot in Figure 1 shows that for the synthetic model the algorithm does not converge to the optimum, as their estimation of the Lipschitz constant does not yield valid bounds, as empirically observed by checking criterion (12). For the Tsukuba model, this effect is not so pronounced, which explains good applied results reported in [4]. As can be seen in Figure 1 (bottom), the remaining gap is not significantly larger in this case, however we observed violations of (12) for Jojic’s method here as well. On the Tsukuba model example one can also see, that a gradient step size, inferred from a fixed L_ρ given by Lemma 1, is so small, that there is almost no improvement of the objective function during iteration. Due to the smoothing, a gap between upper and lower bounds remains for any $\rho > 0$ and decreases with the smoothing (see Theorem 2).

Smoothing selection Next we compare Nesterov’s method with fixed smoothing to a method with adaptive smoothing for the same precision. In the first case, precision is selected according to (16). Both methods use adaptive estimation of the Lipschitz constant. Results are shown in Figure 2. Adaptive smoothing often works faster, as can be seen in the top plot of Figure 2, since it

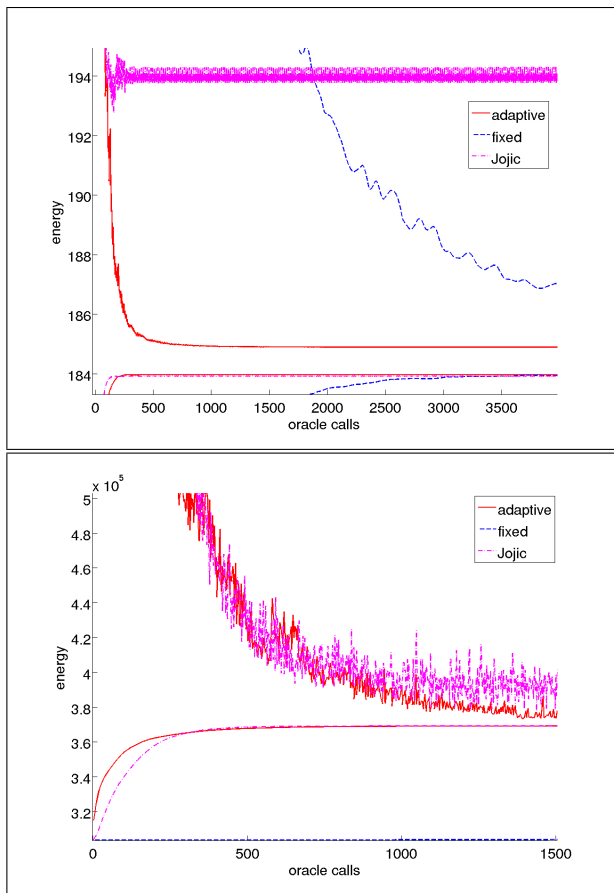


Figure 1. Nesterov’s method for synthetic (top) and Tsukuba (bottom) models with 3 different ways of Lipschitz constant L_ρ selection: (a) fixed, (b) adaptive, (c) L_ρ selected according to [4]. Smoothing value ρ is fixed. While adaptive estimation outperforms the fixed setting, the method suggested in [4] produces invalid values of the Lipschitz constant and does not converge to the optimum.

leads to a smoother function and thus to smaller values of the Lipschitz constant. However, since the adaptive smoothing depends on an actual gap between smoothed and non-smoothed functions, in cases where this gap is close to the upper bound given by (8), adaptive and fixed smoothing lead to similar results, as shown in the bottom plot of Figure 2.

Comparison TRW-S and Sub-Gradient Compared to TRW-S and sub-gradient methods, the proposed method gives better lower bound than TRW-S and converges significantly faster than the sub-gradient ascent. As an update rule for the sub-gradient ascent we use $\lambda^{t+1} = \lambda^t + \frac{\partial U(\lambda^t)}{2\sqrt{t+1}}$, where $\partial U(\lambda^t)$ denotes a sub-gradient of the dual function U . TRW-S is enormously fast, but can get stuck in local fixed points, as shown in the top plot of Figure 3. Unlike

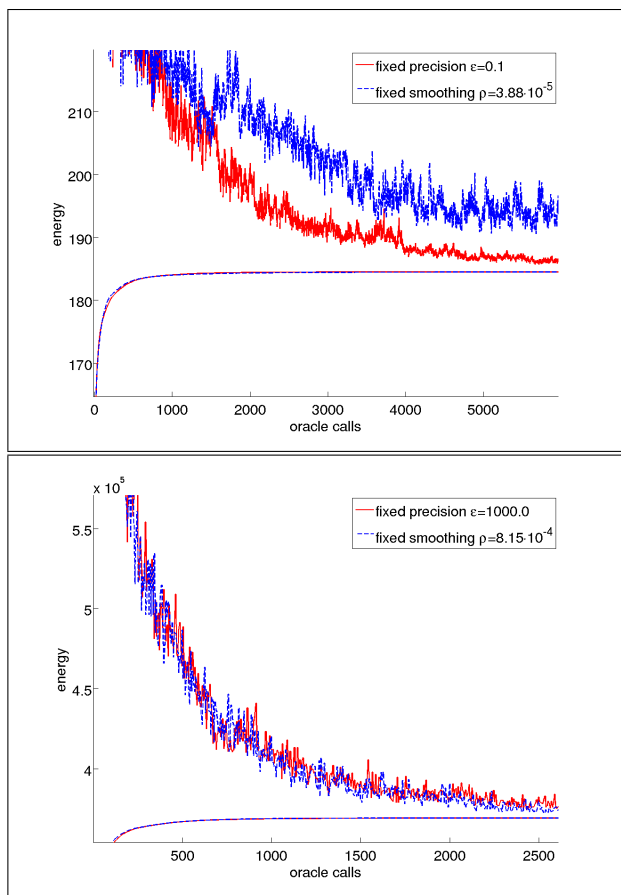


Figure 2. Nesterov’s method for synthetic (top) and Tsukuba (bottom) models with (a) fixed smoothing ρ and (b) adaptive smoothing calculated from a fixed precision ϵ . Parameters ρ and ϵ are connected by (16). Adaptive smoothing usually works faster, since it leads to a smoother function and thus to smaller values of the Lipschitz constant. However, since the adaptive smoothing depends on an actual gap between smoothed and non-smoothed functions, in cases where this gap is close to the upper bound given by (8), adaptive and fixed smoothing lead to similar results, as observed in the bottom plot.

TRW-S, the sub-gradient method is guaranteed to converge to the optimum, but its convergence is extremely slow.

Comparison to Smoothed NPBP Finally, we compare our method of solving a smoothed objective to NPBP, for which we use the entropy approximation as suggested in [2] and set $c_{ab} = 1$, $c_a = 0$ and $c_{ab,a} = 0$. We have selected different values of smoothing parameters ρ for these methods to guarantee, that upper bounds to a difference between smoothed and non-smoothed objectives coincide. For NPBP we apply additionally our method to construct a primal bound. This ends up in a mathematically sound stopping criterion for NPBP, which is lacking in [2]. However, since we optimize different smoothed functions, their opti-

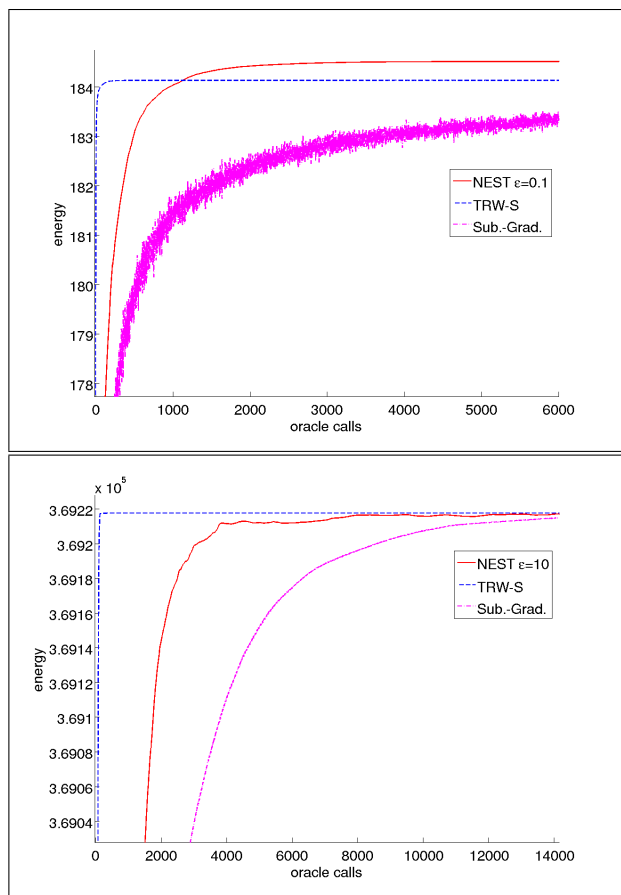


Figure 3. Comparison of (a) Nesterov’s, (b) TRW-S and (c) sub-gradient methods for a synthetic (top) and Tsukuba model (bottom). The plot shows LP lower bounds. TRW-S is the fastest one, but it gets stuck in a fixed point in the top plot, whereas Nesterov’s method calculates a tighter lower bound on the objective. The sub-gradient method is the slowest one.

mal values differ and a fair comparison is not obvious. With less smoothing we obtain tighter bounds for both methods as shown in Figure 4, while the speed of convergence decreases when the smoothing decreases.

5. Conclusion

We presented an in-depth study of Nesterov’s optimal first-order optimization scheme applied to the MAP labeling problem based on Lagrangian decomposition. Our study shows that a direct application of the scheme leads to poor convergence rates based on parameter settings governed by the worst-case optimality bounds. As a remedy, we proposed to modify the approach by i) adaptively estimating and selecting *both* the Lipschitz constant and the smoothing parameter, respectively, and ii) a sound termination condition based on the primal-dual gap.

Modification i) still enables to theoretically infer favor-

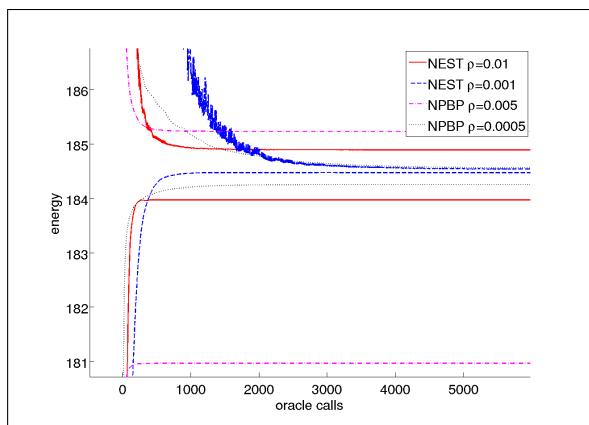


Figure 4. Comparison of (a) Nesterov's method and (b) NPBP for the synthetic model for two different smoothing values. Corresponding values of ρ for Nesterov's method and NPBP differ in two times due to different entropy approximations used in these methods. For smaller ρ both methods produce tighter bounds, but show slower convergence.

able complexity bounds and runtime guarantees. Contribution ii) removes ad-hoc thresholds for stopping the iteration and thus ensures comparability and reproducibility of results. It entails a method for constructing a primal feasible solution that should also be applicable to alternative approaches focusing on dual objective optimization. In our experiments we applied it to generate a primal solution for a Norm-Product Belief Propagation [2]. Our experiments also show that our method i) converges significantly faster than the sub-gradient ascent and ii) has a comparable convergence to the state-of-the-art smoothed Norm-Product Belief Propagation.

Our further work will focus on graphical models that are more general than the grid graphs considered in this paper. While such grid graphs naturally appear in standard low-level vision problems as current benchmarks show, less structured graphs are also of vital interest for various applications. Early experiments indicate that the relative performance of our method increase considerably in these cases, and that our contribution provides a solid basis for tackling such problems.

References

- [1] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, New York, NY, USA, 2004. 3
- [2] T. Hazan and A. Shashua. Norm-product belief propagation: Primal-dual message-passing for approximate inference. *IEEE Trans. on Information Theory*, Accepted for publication June 2010. 1, 5, 6, 7
- [3] J. K. Johnson, D. Malioutov, and A. S. Willsky. Lagrangian relaxation for MAP estimation in graphical models. In *45th*

- Annual Allerton Conference on Communication, Control and Computing*, 2007. 1
- [4] V. Jojic, S. Gould, and D. Koller. Accelerated dual decomposition for MAP inference. In *ICML*, pages 503–510, 2010. 2, 3, 4, 5
- [5] V. Kolmogorov. Convergent tree-reweighted message passing for energy minimization. *IEEE Trans. on Pattern Analysis and Machine Intelligence (PAMI)*, 28(10):1568–1583, 2006. 1, 5
- [6] N. Komodakis, N. Paragios, and G. Tziritas. MRF energy minimization and beyond via dual decomposition. *IEEE Trans. PAMI (in press)*. 2
- [7] N. Komodakis, N. Paragios, and G. Tziritas. MRF optimization via dual decomposition: Message-passing revisited. In *ICCV*, 2007. 1, 2, 5
- [8] Y. Nesterov. A method for solving a convex programming problem with convergence rate $1/k^2$. *Soviet Math. Dokl.*, 27(2):372–376, 1983. 3
- [9] Y. Nesterov. *Introductory Lectures on Convex Optimization: A Basic Course*. Kluwer Academic Publishers, Boston/Dordrecht/London, 2004. 1, 2, 3
- [10] Y. Nesterov. Smooth minimization of non-smooth functions. *Math. Program., Ser. A*(103):127–152, 2004. 2, 3, 4
- [11] Y. Nesterov. Gradient methods for minimizing composite objective function. *CORE Discussion Paper 2007/76*, page 30, 2007. 3
- [12] P. Ravikumar, A. Agarwal, and M. J. Wainwright. Message-passing for graph-structured linear programs: Proximal methods and rounding schemes. *J. Mach. Learn. Res.*, 11:1043–1080, 2010. 1
- [13] R. Rockafellar and R. J.-B. Wets. *Variational Analysis*. Springer, 2nd edition, 2004. 3
- [14] M. Schlesinger. Syntactic analysis of two-dimensional visual signals in the presence of noise. *Kibernetika*, (4):113–130, July-August 1976. 1
- [15] R. Szeliski, R. Zabih, D. Scharstein, O. Veksler, V. Kolmogorov, A. Agarwala, M. Tappen, and C. Rother. A comparative study of energy minimization methods for Markov random fields with smoothness-based priors. *IEEE Trans. Pattern Anal. Mach. Intell.*, 30:1068–1080, June 2008. 2, 5
- [16] M. Wainwright, T. Jaakkola, and A. Willsky. MAP estimation via agreement on (hyper)trees: message passing and linear programming approaches. In *Allerton Conf. on Communication, Control and Computing*, 2002. 2
- [17] M. J. Wainwright and M. I. Jordan. Graphical models, exponential families, and variational inference. *Found. Trends Mach. Learn.*, 1(1-2):1–305, 2008. 1, 2
- [18] T. Werner. A linear programming approach to max-sum problem: A review. *IEEE Trans. on Pattern Recognition and Machine Intelligence (PAMI)*, 29(7), July 2007. 1
- [19] T. Werner. Revisiting the decomposition approach to inference in exponential families and graphical models. Technical report, Center for Machine Perception, Czech Technical University, 2009. 1