

Efficient SDP Inference for Fully-connected CRFs Based on Low-rank Decomposition

Peng Wang¹, Chunhua Shen^{1,2}, Anton van den Hengel^{1,2}

¹University of Adelaide, Australia

²Australian Centre for Robotic Vision

Abstract

Conditional Random Fields (CRFs) are one of the core technologies in computer vision, and have been applied to a wide variety of tasks. Conventional CRFs typically define edges between neighboring image pixels, resulting in a sparse graph over which inference can be performed efficiently. However, these CRFs fail to model more complex priors such as long-range contextual relationships. Fully-connected CRFs have thus been proposed. While there are efficient approximate inference methods for such CRFs, usually they are sensitive to initialization and make strong assumptions. In this work, we develop an efficient, yet general SDP algorithm for inference on fully-connected CRFs. The core of the proposed algorithm is a tailored quasi-Newton method, which solves a specialized SDP dual problem and takes advantage of the low-rank matrix approximation for fast computation. Experiments demonstrate that our method can be applied to fully-connected CRFs that could not previously be solved, such as those arising in pixel-level image co-segmentation.

1. Introduction

Semantic image segmentation, or pixel labelling, is a key problem in computer vision. Given an image, the task is to label every pixel against one or multiple pre-defined object categories. It is clear that to achieve satisfactory results, one must exploit contextual information. Scalability and speed of the algorithm are also of concerns, if we are to design an algorithm applicable to high-resolution images.

Conditional random fields (CRFs) have been one of the most successful approaches to semantic pixel labelling, which solves the problem as maximum a posteriori (MAP) estimation. Standard CRFs contain unary potentials that are typically defined on low-level local features of texture, color, and locations. Edge potentials, which are typically defined on 4- or 8-neighboring pixels, consist of smoothness terms that penalize label disagreement between similar pixels, and terms that model contextual relationships between different classes. Although these CRF models have

achieved encouraging results for segmentation, they fail to capture more complex priors such as long-range contextual information.

In the literature, fully-connected CRFs have been proposed for this purpose. The main challenge for inference on fully-connected CRFs stems from the computational cost. A fully-connected CRF over N image pixels has N^2 edges. Even for a small image with a few thousand pixels, the number of edges can be a few million. Although there have been a variety of methods for MAP estimation [1–10], they are usually computationally infeasible for such cases. The authors of [11, 12] have proffered an efficient mean field approximation method for MAP inference in multi-label CRF models with fully-connected pairwise terms. In their algorithms, the computational bottleneck can be expressed as the product of kernel matrices and column vectors. Given the assumption that the pairwise terms are in the form of a weighted mixture of Gaussian kernels, a filter-based method is used to accelerate the computation of the matrix-vector product in [11, 12]. Other work [13, 14] on the MAP inference in fully-connected CRFs also incorporate filter-based methods for fast computation, yet with different assumptions. The method proposed in [13] can be applied on generalized RBF kernels, instead of the original Gaussian kernels [11]. An efficient inference algorithm was developed in [14] for a special type of fully-connected CRF, in which the edge potentials are defined to capture spatial relationships among different objects, and only depend on their relative positions (that is they are spatially stationary). Note that the assumptions in [11–14] limit the practical value of these approaches. Note also that the mean field approximation adopted in [11–13] and the quadratic programming approach used in [14] may converge to local optima.

In general, semidefinite programming (SDP) relaxation provides accurate solutions for MAP estimation problems, but it is usually computationally inefficient (see [3] for a comparison of different relaxation methods). Standard interior-point methods require $\mathcal{O}(q^3 + qn^3 + q^2n^2)$ flops to solve a generic SDP problem in worst-case, where n and q are the semidefinite matrix dimension and the number of linear constraints respectively. Recently, several scal-

able SDP methods have been proposed for MAP estimation. Huang *et al.* [15] proposed an alternating direction methods of multipliers method (ADMM) to solve large-scale MAP estimation problems. Wang *et al.* [16] presented an efficient dual approach (refer to as SDCut), which can also be applied for MAP estimation. However, their methods still cannot be applied directly to large-scale fully-connected CRFs.

In this paper, a more efficient SDP algorithm is proposed for MAP estimation in fully-connected CRFs. There are two key contributions in this work: 1. An efficient low-rank SDP approach (based on SDCut) is proposed for MAP estimation in large-scale fully-connected CRFs. Several significant improvements over SDCut are presented, which makes SDCut much more scalable. The proposed SDP method solves a convex problem, and generally provides more stable and accurate solutions than mean field approximation. 2. Similar to the mean field approach [11], the most computationally expensive part of the proposed SDP method is calculating the product of kernel matrices and column vectors. Instead of the filter-based method used in [11], low-rank approximation methods for SPSD kernels (whose kernel matrix is symmetric positive semidefinite) are seamlessly integrated into our SDP method for fast computation. The use of low-rank approximation relaxes the limitation on the pairwise term from being (a mixture of) Gaussian kernels to all symmetric positive-semidefinite kernels.

As a result, our method is much more general and scalable, and so has a broader range of applications. The proposed SDP approach can handle fully-connected CRFs of $\#states \times \#variables$ up to 10^6 . In particular, we show that on an image co-segmentation application, the fast method of [11] is not applicable while our method achieves superior segmentation accuracy. To our knowledge, our method is the first pixel-level co-segmentation method. All previous co-segmentation methods have relied on super-pixel pre-processing in order to make the computation tractable. Wang *et al.* [17] and Frostig *et al.* [18] also proposed efficient approaches which find near-optimal solutions to SDP relaxation to MAP problems. The main difference is that their methods solve (generally *nonconvex*) quadratically constrained quadratic programs by projected gradient descent, while ours uses quasi-Newton methods to solve a *convex* semidefinite least-square problem. Notation is listed in Table 1. An extended version of this work is available at <http://arxiv.org/abs/1504.01492>.

2. Fully-connected Pairwise CRFs with SPSD Kernels

Consider a random field over N random variables $\mathbf{x} = [x_1, x_2, \dots, x_N]^T$ conditioned on the observation \mathbf{I} . Each variable can be assigned a label from the set $\mathcal{L} = \{1, \dots, L\}$. The energy function of a CRF (\mathbf{I}, \mathbf{x}) can be

\mathbf{X}	A matrix (bold upper-case letters).
\mathbf{x}	A column vector (bold lower-case letters).
S^n	The space of $n \times n$ symmetric matrices.
S_+^n	The cone of $n \times n$ symmetric positive semidefinite (SPSD) matrices.
\mathbb{R}^n	The space of real-valued $n \times 1$ vectors.
$\mathbb{R}_+^n, \mathbb{R}_-^n$	The non-negative and non-positive orthants of \mathbb{R}^n .
\mathbf{I}_n	The $n \times n$ identity matrix.
$\mathbf{0}$	An all-zero vector with proper dimension.
$\mathbf{1}$	An all-one vector with proper dimension.
\leq, \geq	Inequality between scalars or element-wise inequality between column vectors.
$\text{diag}(\mathbf{X})$	The vector of the diagonal elements of the input matrix \mathbf{X} .
$\text{Diag}(\mathbf{x})$	The $n \times n$ diagonal matrix whose main diagonal vector is the input vector \mathbf{x} .
$\text{trace}(\cdot)$	The trace of a matrix.
$\text{rank}(\cdot)$	The rank of a matrix.
$\delta(\text{cond})$	The indicator function which returns 1 if <i>cond</i> is true and 0 otherwise.
$\ \cdot\ _F$	Frobenius-norm of a matrix.
$\langle \cdot, \cdot \rangle$	Inner product of two matrices.
\circ	Hadamard product of two matrices.
\otimes	Kronecker product of two matrices.
$\nabla f(\cdot)$	The first-order derivative of function $f(\cdot)$.
$\nabla^2 f(\cdot)$	The second-order derivative of function $f(\cdot)$.
$n!$	The factorial of a non-negative integer n .

Table 1: Notation.

expressed by the following Gibbs distribution:

$$P(\mathbf{x}|\mathbf{I}) := \frac{1}{Z(\mathbf{I})} \exp(-E(\mathbf{x}|\mathbf{I})), \quad (1)$$

where $E(\mathbf{x}|\mathbf{I})$ denotes the Gibbs energy function w.r.t. a labelling $\mathbf{x} \in \mathcal{L}^N$, and $Z(\mathbf{I}) := \sum_{\mathbf{x} \in \mathcal{L}^N} \exp(-E(\mathbf{x}|\mathbf{I}))$ is the partition function. In the rest of the paper, the conditioning w.r.t. \mathbf{I} is dropped for simplicity of notation.

Assuming that $E(\mathbf{x})$ only contains unary and pairwise terms, the MAP inference problem for the CRF (\mathbf{I}, \mathbf{x}) is equivalent to the following energy minimization problem:

$$\min_{\mathbf{x} \in \mathcal{L}^N} E(\mathbf{x}) := \sum_{i \in \mathcal{N}} \psi_i(x_i) + \sum_{i,j \in \mathcal{N}, i < j} \psi_{i,j}(x_i, x_j), \quad (2)$$

where $\mathcal{N} := \{1, \dots, N\}$. $\psi_i : \mathcal{L} \rightarrow \mathbb{R}$ and $\psi_{i,j} : \mathcal{L}^2 \rightarrow \mathbb{R}$ correspond to the unary and pairwise potentials respectively. The pairwise potentials considered in this paper can be written as:

$$\psi_{i,j}(x_i, x_j) := \mu(x_i, x_j) \sum_{m=1}^M w^{(m)} k^{(m)}(\mathbf{f}_i, \mathbf{f}_j), \quad (3)$$

where $\mathbf{f}_i, \mathbf{f}_j \in \mathbb{R}^D$ indicate D -dimensional feature vectors corresponding to variables x_i and x_j respectively. $k^{(m)} : \mathbb{R}^D \times \mathbb{R}^D \rightarrow \mathbb{R}$ denotes the function of the m -th SPSD kernel and $w^{(m)} \in \mathbb{R}_+$ is the associated weight. Following the term in [11], $\mu : \mathcal{L}^2 \rightarrow [0,1]$ is used to represent a symmetric label compatibility function, which has the properties that $\mu(l, l') = \mu(l', l), \forall l, l' \in \mathcal{L}$ and $\mu(l, l) = 0, \forall l \in \mathcal{L}$. The label compatibility function penalizes similar pixels being assigned with different/incompatible labels. A simple label compatibility function would be given by Potts model, that is $\mu(l, l') = \delta(l \neq l')$. The form of pairwise potential in (3)

is very general, and can be used to represent many potentials of practical interest.

Mean field approximation is used in [11] for solving problem (2), which is considered to be state-of-the-art. A filter-based method [19] is adopted in [11] to accelerate the computation. In the following two sections, we will briefly revisit mean field approximation and the filter-based method, especially their respective limitations.

2.1. Mean Field Approximation

In mean field approximation [20], a distribution $Q(\mathbf{x})$ is introduced to approximate the Gibbs distribution $P(\mathbf{x}|\mathbf{I})$, in which the marginals are supposed to be independent to each other such that $Q(\mathbf{x}) = \prod_{i \in \mathcal{N}} Q_i(x_i)$. The KL-divergence $D(Q||P)$ is minimized by iteratively applying the following update equation:

$$Q_i(l) = \frac{1}{Z_i} \exp \left(-\psi_i(l) - \sum_{l' \in \mathcal{L}} \mu(l, l') \sum_{m=1}^M w^{(m)} \sum_{j \in \mathcal{N}, j \neq i} k^{(m)}(\mathbf{f}_i, \mathbf{f}_j) Q_i(l') \right), \quad (4)$$

where Z_i is the normalization factor.

The computational bottleneck in updating the above equation can be expressed as the matrix-vector product $\mathbf{K}^{(m)} \mathbf{q}_l, \forall l \in \mathcal{L}, m = 1, \dots, M$, where $\mathbf{K}^{(m)} \in \mathcal{S}_+^N$ denotes the kernel matrix corresponding to $k^{(m)}$, that is $K_{ij}^{(m)} = k^{(m)}(\mathbf{f}_i, \mathbf{f}_j)$, and $\mathbf{q}_l = [Q_1(l), \dots, Q_N(l)]^\top, \forall l \in \mathcal{L}$. The naive implementation of the matrix-vector product needs $\mathcal{O}(N^2)$ time. Krähenbühl and Koltun [11] proposed to use a filter-based approach to compute the matrix-vector product in $\mathcal{O}(N)$ time, which will be discussed in the next section.

One significant limitation of mean field approximation is that it may converge to one of potentially many local optima, because the variational problem to be optimized may be non-convex. A consequence of this non-convexity is that mean field is often sensitive to the initialization of Q .

2.2. Filter-based Matrix-vector Product

Filter-based methods [19] have been used in [11] to speed up the above matrix-vector product. The method in [11] is based on the assumption that pairwise potentials are Gaussian kernels:

$$k^{(m)}(\mathbf{f}_i, \mathbf{f}_j) = \exp \left(-\frac{1}{2} (\mathbf{f}_i - \mathbf{f}_j)^\top \boldsymbol{\Lambda}^{(m)} (\mathbf{f}_i - \mathbf{f}_j) \right), \quad (5)$$

where $\boldsymbol{\Lambda}^{(m)} \in \mathcal{S}_+^D, m = 1, 2, \dots, M$. The product of a Gaussian kernel matrix and an arbitrary column vector can be expressed as a Gaussian convolution w.r.t. $\boldsymbol{\Lambda}^{(m)}$ in feature space (see [11, 19] for more details). From the viewpoint of signal processing, the Gaussian convolution can be seen as a low-pass filter over the feature space. Then the convolution result can be recovered from a set of samples whose spacing is proportional to the standard deviation

of the filter. A number of filtering methods [19, 21] can be used to compute the convolution efficiently, in which the computational complexity and memory requirement are both linear in N .

Filter-based approaches have a number of limitations, however:

1. In general, the pairwise potentials are limited to Gaussian kernels over a Euclidean feature space.
2. The feature dimension cannot be very high. The bilateral filtering method in [21] has an exponential complexity w.r.t. the dimension D . The time complexity using a permutohedral lattice [19] is quadratic in D , which works well only when the input dimension is $5 \sim 20$. Because it does not create new lattice points during the blur step, an accuracy penalty is accumulated with the growth of feature dimension.

In the following sections, we propose alternative methods for the matrix-vector product and MAP inference to overcome the aforementioned limitations.

3. Matrix-vector Product Based on Low-rank Approximation

One key contribution of this paper is the use of a low-rank approximation to the positive semidefinite kernel matrix, based on which low-rank quasi-Newton methods are developed for large-scale SDP CRF inference. We propose to approximate an SPSD kernel matrix $\mathbf{K} \in \mathcal{S}_+^N$ by a low-rank representation: $\mathbf{K} \approx \boldsymbol{\Phi} \boldsymbol{\Phi}^\top$, where $\boldsymbol{\Phi} \in \mathbb{R}^{N \times R_K}$ and $R_K \ll N$, such that both of the computational complexity and memory requirement for computing the aforementioned matrix-vector product are linear in N . Compared to [19, 21], the pairwise potential function is generalized to any positive semidefinite kernel function and there is no restriction on the input feature dimension.

The optimal low-rank approximation can be obtained by eigen-decomposition, while it is computationally inefficient whose computational complexity is generally cubic in N . There are a number of low-rank approximation methods achieving linear complexity in N , including Nyström methods [22–24], incomplete Cholesky decomposition [25, 26], random Fourier features [27, 28], and homogeneous kernel maps [29]. For detailed discussion, please refer to the review papers [30–32]. We adopt Nyström methods [29] in this paper for the low-rank approximation of kernel matrices.

Nyström methods can be used to approximate a positive semidefinite matrix $\mathbf{K} \in \mathcal{S}_+^N$, by sample $R_0 \ll N$ columns of \mathbf{K} (refer to as landmarks). Firstly \mathbf{K} is expressed as:

$$\mathbf{K} = \begin{bmatrix} \mathbf{W} & \mathbf{K}_{2,1}^\top \\ \mathbf{K}_{2,1} & \mathbf{K}_{2,2} \end{bmatrix}, \quad (6)$$

where $\mathbf{W} \in \mathcal{S}^{R_0}$ denotes the intersection of the sampled R_0 columns and rows. The matrix $\mathbf{K}_{2,2} \in \mathcal{S}^{N-R_0}$ can be

approximated as:

$$\mathbf{K}_{2,2} \approx \mathbf{K}_{2,1} \mathbf{\Gamma}_R \mathbf{\Sigma}_R^{-1} \mathbf{\Gamma}_R^\top \mathbf{K}_{2,1}^\top, \quad (7)$$

where $R \leq R_0$ and $\mathbf{\Sigma}_R = \text{Diag}([\lambda_1, \dots, \lambda_R]^\top)$. $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_R > 0$ are the R -largest eigenvalues of \mathbf{W} and $\mathbf{\Gamma}_R$ contains the corresponding (column) eigenvectors. Note that $\mathbf{\Gamma}_R \mathbf{\Sigma}_R \mathbf{\Gamma}_R^\top$ is the best rank- R approximation to \mathbf{W} . Then we have a rank- R approximation to \mathbf{K} :

$$\mathbf{K} \approx \left(\begin{bmatrix} \mathbf{W} \\ \mathbf{K}_{2,1} \end{bmatrix} \mathbf{\Gamma}_R \mathbf{\Sigma}_R^{-\frac{1}{2}} \right) \left(\begin{bmatrix} \mathbf{W} \\ \mathbf{K}_{2,1} \end{bmatrix} \mathbf{\Gamma}_R \mathbf{\Sigma}_R^{-\frac{1}{2}} \right)^\top, \quad (8)$$

which is proved to have a bounded error to the optimal rank- R approximation given by the eigen-decomposition [31].

There are several strategies to sample representative landmarks, *i.e.*, columns of \mathbf{K} , including the standard uniform sampling [23], non-uniform sampling [24] and k -means clustering [33]. In this paper, we adopt the k -means method in [33] to select landmarks. At each round of k -means, only R columns of \mathbf{K} , rather than the entire matrix \mathbf{K} , is required to be instantiated. Note that for Nyström methods, the positive semidefinite matrix \mathbf{K} to be approximated can be any kernel matrix $\mathbf{K}^{(m)}$ or the linear combination $\sum_{m=1}^M w^{(m)} \mathbf{K}^{(m)}$.

4. SDP Approach to MAP Estimation

SDP Relaxation In this section, we introduce an SDP relaxation to the problem (2). Throughout the rest of this paper, the label compatibility function is assumed to be given by Potts model, that is $\mu(l, l') = \delta(l \neq l')$ ¹.

By defining $\mathbf{X} \in \{0, 1\}^{N \times L}$, $\mathbf{H} \in \mathbb{R}^{N \times L}$ and $\mathbf{K} \in \mathcal{S}_+^N$ as $X_{i,l} = \delta(x_i = l)$, $H_{i,l} = \psi_i(l)$ and $K_{i,j} = \sum_{m=1}^M w^{(m)} k^{(m)}(\mathbf{f}_i, \mathbf{f}_j)$, the problem (2) can be expressed as the following binary quadratic problem (BQP)²:

$$\min_{\mathbf{X} \in \{0,1\}^{N \times L}} \tilde{\mathbf{E}}(\mathbf{X}) := \langle \mathbf{H}, \mathbf{X} \rangle - \frac{1}{2} \langle \mathbf{X} \mathbf{X}^\top, \mathbf{K} \rangle \quad (9a)$$

$$\text{s.t.} \quad \sum_{l=1}^L X_{i,l} = 1, \quad \forall i \in \mathcal{N}, \quad (9b)$$

Note that $\mathbf{E}(\mathbf{x}) = \tilde{\mathbf{E}}(\mathbf{X}) + \frac{1}{2} \mathbf{1}^\top \mathbf{K} \mathbf{1}$ for equivalent \mathbf{x} and \mathbf{X} .

By introducing $\mathbf{Y} := \begin{bmatrix} \frac{1}{x} \\ \frac{1}{x} \end{bmatrix} \begin{bmatrix} \frac{1}{x} \\ \frac{1}{x} \end{bmatrix}^\top$, the corresponding SDP relaxation to problem (9) can be expressed as:

$$\min_{\mathbf{Y} \in \mathcal{S}_+^{N+L}} \langle \mathbf{Y}, \frac{1}{2} \begin{bmatrix} \mathbf{0} & \mathbf{H}^\top \\ \mathbf{H} & -\mathbf{K} \end{bmatrix} \rangle, \quad (10a)$$

$$\text{s.t.} \quad Y_{l,l} = 1, \quad l \in \mathcal{L}, \quad (10b)$$

$$\frac{1}{2} (Y_{l,l'} + Y_{l',l}) = 0, \quad l \leq l', l, l' \in \mathcal{L}, \quad (10c)$$

$$\frac{1}{2} \sum_{l=1}^L (Y_{i+L,l} + Y_{l,i+L}) = 1, \quad i \in \mathcal{N}, \quad (10d)$$

$$Y_{i+L,i+L} = 1, \quad i \in \mathcal{N}. \quad (10e)$$

Clearly we have $\text{trace}(\mathbf{Y}) = N + L$ which is implicitly encoded by the linear constraints. The non-convex constraint

¹The SDP relaxation corresponding to an arbitrary label compatibility function is discussed in the extended version.

²The derivation from (2) to (9) can be found in the extended version.

$\text{rank}(\mathbf{Y}) = L$ is dropped by the SDP relaxation.

In the above formulation, all the constraints (10b), (10c), (10d), (10e) are linear w.r.t. \mathbf{Y} . Therefore they can be rewritten in terms of inner products, that is $\langle \mathbf{Y}, \mathbf{B}_i \rangle = b_i$, $i = 1, 2, \dots, q$, where $\mathbf{B}_i \in \mathcal{S}^{N+L}$, $\mathbf{b} \in \mathbb{R}^q$ and $q = 2N + L(L+1)/2$ is the total number of linear constraints in (10). The problem (10) can thus be expressed as the following general form:

$$\min_{\mathbf{Y} \in \mathcal{S}_+^{N+L}} p(\mathbf{Y}) := \langle \mathbf{Y}, \mathbf{A} \rangle, \quad (11a)$$

$$\text{s.t.} \quad \langle \mathbf{Y}, \mathbf{B}_i \rangle = b_i, \quad i = 1, 2, \dots, q, \quad (11b)$$

$$\text{where } \mathbf{A} = \frac{1}{2} \begin{bmatrix} \mathbf{0} & \mathbf{H}^\top \\ \mathbf{H} & -\mathbf{K} \end{bmatrix}.$$

4.1. Low-rank Quasi-Newton Methods

In this section, *several major improvements are proposed to make SDCut [16] scalable to the large-scale energy minimization problem (9), which is another key contribution of this work.*

4.1.1 SDCut Formulation

SDCut [16] solves the following approximation of (11) using quasi-Newton methods:

$$\min_{\mathbf{Y} \in \mathcal{S}_+^n} p_\gamma(\mathbf{Y}) := \langle \mathbf{Y}, \mathbf{A} \rangle + \frac{1}{2\gamma} \|\mathbf{Y}\|_F^2 - \frac{n^2}{2\gamma}, \quad (12a)$$

$$\text{s.t.} \quad \langle \mathbf{Y}, \mathbf{B}_i \rangle = b_i, \quad i = 1, 2, \dots, q, \quad (12b)$$

where $\gamma > 0$ is a parameter and $n = N + L$. Given a sufficiently large γ , the difference between the optimal solutions to (12) and (11) can be very small [16]. Note that $-\frac{n^2}{2\gamma}$ is a constant and it can be removed from the optimization problem.

Remark 1. *The Lagrangian dual problem of (12) can be simplified to*

$$\max_{\mathbf{u} \in \mathbb{R}^q} d_\gamma(\mathbf{u}) := -\frac{\gamma}{2} \|(\mathbf{C}(\mathbf{u}))_+\|_F^2 - \mathbf{u}^\top \mathbf{b} - \frac{n^2}{2\gamma}, \quad (13)$$

where $\mathbf{C}(\cdot) : \mathbb{R}^q \rightarrow \mathcal{S}^n$ is defined as $\mathbf{C}(\mathbf{u}) := -\mathbf{A} - \sum_{i=1}^q u_i \mathbf{B}_i$, and $(\cdot)_+ : \mathcal{S}^n \rightarrow \mathcal{S}_+^n$ is defined as $(\mathbf{Y})_+ = \mathbf{\Gamma} \text{Diag}(\max(\mathbf{0}, \boldsymbol{\lambda})) \mathbf{\Gamma}^\top$. $\boldsymbol{\lambda} := [\lambda_1, \dots, \lambda_n]^\top$ and $\mathbf{\Gamma}$ stand for the respective eigenvalues and eigenvectors of \mathbf{Y} , that is $\mathbf{Y} = \mathbf{\Gamma} \text{Diag}(\boldsymbol{\lambda}) \mathbf{\Gamma}^\top$. The relationship between the optimal solution to the primal (12), \mathbf{Y}^* , and the solution to the dual (13), \mathbf{u}^* , is $\mathbf{Y}^* = \gamma(\mathbf{C}(\mathbf{u}^*))_+$.

The proof of the above results can be found in [16]. Assuming that $\text{trace}(\mathbf{Y}) = n$, it is shown in [16] that the objective function value of the dual (13) for any $\mathbf{u} \in \mathbb{R}^q$ yields a lower-bound to the optimal objective value of the BQP (9).

The objective function of (13), d_γ , is differentiable but not necessarily twice differentiable, and its gradient is:

$$\nabla d_\gamma(\mathbf{u}) = -\gamma \left[\langle (\mathbf{C}(\mathbf{u}))_+, \mathbf{B}_1 \rangle, \dots, \langle (\mathbf{C}(\mathbf{u}))_+, \mathbf{B}_q \rangle \right]^\top - \mathbf{b}. \quad (14)$$

Algorithm 1 LR-SDCut algorithm for MAP estimation.

Input: \mathbf{A} , $\{\mathbf{B}_i\}_{i=1,2,\dots,q}$, \mathbf{b} , γ , K_{\max} , $\tau > 0$, $r \ll N$.

- 1 **Initialization:** $\mathbf{u}^{(0)} = \mathbf{0}$, $\tilde{\mathbf{E}}^* = +\inf$, $\mathbf{A} = \mathbf{A} - \nu \mathbf{I}_N$ where ν is the r -th smallest eigenvalue of \mathbf{A} .
 - 2 **for** $k = 0, 1, 2, \dots, K_{\max}$ **do**
 - 3 **Step1:** $\mathbf{u}^{(k+1)} = \mathbf{u}^{(k)} - \rho \mathbf{H} \nabla d_\gamma(\mathbf{u}^{(k)})$, where \mathbf{H} is updated to approximate $(\nabla^2 d_\gamma(\mathbf{u}^{(k)}))^{-1}$ and $0 < \rho \leq 1$ is the step size.
 Step2: $\mathbf{X}^{(k+1)} = \text{Round}(\gamma(\mathbf{C}(\mathbf{u}^{(k+1)}))_+)$.
 Step3: If $\tilde{\mathbf{E}}(\mathbf{X}^{(k+1)}) < \tilde{\mathbf{E}}^*$, $\mathbf{X}^* = \mathbf{X}^{(k+1)}$.
 Step4: Exit, if $\frac{(d_\gamma(\mathbf{u}^{(k+1)}) - d_\gamma(\mathbf{u}^{(k)}))}{\max\{|d_\gamma(\mathbf{u}^{(k+1)})|, |d_\gamma(\mathbf{u}^{(k)})|, 1\}} \leq \tau$.
 - 4 **end**
- Output:** \mathbf{X}^* , $\tilde{\mathbf{E}}^*$.
-

Such that Wang *et al.* [16] adopted quasi-Newton methods to solve the dual problem (13). At each iteration of quasi-Newton methods, only the objective function d_γ and its gradient (14) need to be computed, where the computational bottleneck is the calculation of $(\mathbf{C}(\mathbf{u}))_+$, which needs all the positive eigenvalues and the corresponding eigenvectors of $\mathbf{C}(\mathbf{u})$.

Although it is shown in [16] that SDCut already runs much faster than standard interior-point methods, there are still several issues to be addressed for the problem to be solved in this work:

1. It is shown in [16] that $\text{rank}((\mathbf{C}(\mathbf{u}))_+)$ drops significantly in the first several iterations, and Lanczos methods [34] can be used to efficiently compute a few leading eigenpairs. However, because $(\mathbf{C}(\mathbf{u}))_+$ is not necessarily low-rank in the initial several iterations, much of time may be spent on the first several eigen-decompositions. In the CRFs considered in this paper, there are up to 681,600 variables. Using the original SDCut method, the time spent on the first several iterations can be prohibitive.
2. In general, a BFGS-like method has a superlinear convergence speed under the condition that the objective function is twice continuously differentiable. However, the dual objective function (13) is not necessarily twice differentiable. *So the convergence speed of SDCut is unknown.* In practice, SDCut usually needs more than 100 iterations to converge.

In the next two sections, we introduce two improvements to the SDCut method, which address the above two problems and increase the scalability of SDCut significantly. The improved method is referred to as LR-SDCut and its procedure is summarized in Algorithm 1.

4.1.2 A Low-rank Initial Point

If the initialization of the dual variable $\mathbf{u}^{(0)}$ is $\mathbf{0}$, then we have $\mathbf{C}(\mathbf{u}^{(0)}) = -\mathbf{A}$. Without affecting the optimal solution to (11), \mathbf{A} can be perturbed so as to reduce $\text{rank}((\mathbf{C}(\mathbf{u}^{(0)}))_+)$ to a small integer, based on:

1. For $\mathbf{Y} \in \mathcal{S}_+^n \cap \{\text{trace}(\mathbf{Y}) = n\}$, $\langle \mathbf{Y}, \mathbf{A} + \nu \mathbf{I}_n \rangle = \langle \mathbf{Y}, \mathbf{A} \rangle + \nu n$. So the matrix \mathbf{A} in the problem (11) can be equivalently replaced by $\mathbf{A} + \nu \mathbf{I}_n$, $\forall \nu \neq 0$.

2. Suppose that $\lambda \neq 0$ and $\mathbf{x} \in \mathbb{R}^n$ is an eigenpair of $\mathbf{A} \in \mathcal{S}^n$, *i.e.*, $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$, then $\mathbf{A} + \nu \mathbf{I}_n$ has an eigenpair: $\lambda + \nu$ and \mathbf{x} , $\forall \nu \neq 0$.

To decrease the rank of $(\mathbf{C}(\mathbf{u}^{(0)}))_+$ to $r \ll n$, we can equivalently replace \mathbf{A} by $\mathbf{A} - \nu \mathbf{I}_n$, where ν is the r -th smallest eigenvalue of \mathbf{A} .

4.1.3 Rounding Schemes and Early Stop

Traditionally, a feasible binary solution \mathbf{X} to the BQP problem (9) is obtained by rounding the optimal solution \mathbf{Y}^* to the corresponding SDP formulation (12). The rounding procedure will be carried out until the quasi-Newton algorithm converges. In contrast, we perform the rounding procedure on the non-optimal solution $\mathbf{Y}^{(k)} := \gamma(\mathbf{C}(\mathbf{u}^{(k)}))_+$ at each iteration k of the quasi-Newton algorithm (Step2 in Algorithm 1). In practice, we find that the dual objective value of (13), *i.e.* the lower-bound to the optimal value of $\tilde{\mathbf{E}}(\mathbf{X})$, increases dramatically in the first several iterations. Simultaneously, the value of $\tilde{\mathbf{E}}(\mathbf{X}^{(k)})$ also drops significantly for the first several k s. This observation inspires us to stop the quasi-Newton algorithm long before convergence, with little decline in the quality of the final binary solution to (9).

In this work, we adopt the random rounding scheme proposed in [35] to derive \mathbf{X} from $\mathbf{Y}^{(k)} := \gamma(\mathbf{C}(\mathbf{u}^{(k)}))_+$. Note that because $\mathbf{Y}^{(k)}$ is positive semidefinite, it can be decomposed to $\mathbf{Y}^{(k)} = \Psi\Psi^\top$, where $\Psi \in \mathbb{R}^{N \times R_Y}$ and $R_Y = \text{rank}(\mathbf{Y}^{(k)})$. The rounding scheme can be expressed in the following two steps:

1. *Random Projection:* $\hat{\mathbf{X}} = \Psi\mathbf{P}$, where $\mathbf{P} \in \mathbb{R}^{R_Y \times L}$ and each entry $P_{i,j}$ is independently sampled from the standard Gaussian distribution with mean 0 and variance 1.
2. *Discretization:* Obtain $\mathbf{X} \in \{0,1\}^{N \times L}$ by discretizing the above $\hat{\mathbf{X}}$, that is, $X_{i,l} = \delta(\hat{X}_{i,l} > \hat{X}_{i,l'}, \forall l' \in \mathcal{L}, l' \neq l)$.

4.2. Computational Cost and Memory Requirement

The computational bottleneck of LR-SDCut is the eigen-decomposition of $\mathbf{C}(\mathbf{u})$ at each iteration, which is performed by Lanczos methods [34] in this paper. Lanczos methods only require users to implement the matrix-vector product $\mathbf{C}(\mathbf{u})\mathbf{d} = -\mathbf{A}\mathbf{d} - (\sum_{i=1}^q u_i \mathbf{B}_i)\mathbf{d}$, where $\mathbf{d} \in \mathbb{R}^n$ denotes a so-called ‘‘Lanczos vector’’ produced by Lanczos algorithms iteratively. In this section, we will show how to accelerate the computation of this matrix-vector product by exploiting the specific structures of \mathbf{A} and $\{\mathbf{B}_i\}_{i=1,\dots,q}$, and then give the computational cost and memory requirement of LR-SDCut.

For the problem (10), $\mathbf{A} = \frac{1}{2} \begin{bmatrix} \mathbf{0} & \mathbf{H}^\top \\ \mathbf{H} & -\mathbf{K} \end{bmatrix}$ and $\sum_{i=1}^q u_i \mathbf{B}_i = \begin{bmatrix} \text{Diag}(\mathbf{u}_1) + \frac{1}{2} \text{LTri}(\mathbf{u}_2) & \frac{1}{2} \mathbf{u}_3^\top \otimes \mathbf{1} \\ \frac{1}{2} \mathbf{u}_3 \otimes \mathbf{1}^\top & \text{Diag}(\mathbf{u}_4) \end{bmatrix}$, where $\mathbf{u}_1 \in \mathbb{R}^L$, $\mathbf{u}_2 \in \mathbb{R}^{L(L-1)/2}$, $\mathbf{u}_3, \mathbf{u}_4 \in \mathbb{R}^N$ denote the respective dual variables w.r.t. constraints (10b), (10c), (10d), (10e) and such that $\mathbf{u} = [\mathbf{u}_1^\top, \mathbf{u}_2^\top, \mathbf{u}_3^\top, \mathbf{u}_4^\top]^\top$. $\text{LTri}(\mathbf{u}) :$

$\mathbb{R}^{L(L-1)/2} \rightarrow \mathcal{S}^L$ produces an $L \times L$ symmetric matrix whose lower triangular part is made up of the elements of the input vector $\mathbf{u} \in \mathbb{R}^{L(L-1)/2}$, that is $\text{LTri}(\mathbf{u}) = \begin{cases} 0 & \text{if } i = j \\ u_{(L-1)!/j!+i-j} & \text{if } i > j \\ u_{(L-1)!/i!+j-i} & \text{if } i < j \end{cases}$. Then $\mathbf{C}(\mathbf{u})\mathbf{d}$ is expressed as:

$$\mathbf{C}(\mathbf{u})\mathbf{d} = -\frac{1}{2} \left[\begin{array}{c} \mathbf{H}^\top \mathbf{d}_2 \\ \mathbf{H} \mathbf{d}_1 - \mathbf{K} \mathbf{d}_2 \end{array} \right] - \underbrace{\left[\begin{array}{c} \mathbf{u}_1 \circ \mathbf{d}_1 + \frac{1}{2} \text{LTri}(\mathbf{u}_2) \mathbf{d}_1 + \frac{1}{2} (\mathbf{u}_3^\top \mathbf{d}_2) \mathbf{1} \\ \frac{1}{2} (\mathbf{1}^\top \mathbf{d}_1) \mathbf{u}_3 + \mathbf{u}_4 \circ \mathbf{d}_2 \end{array} \right]}_{(\sum_{i=1}^q u_i \mathbf{B}_i) \mathbf{d}: \mathcal{O}(L^2+N)}, \quad (15)$$

Ad: $\mathcal{O}(NL+NR_K)$ $(\sum_{i=1}^q u_i \mathbf{B}_i) \mathbf{d}: \mathcal{O}(L^2+N)$

where $\mathbf{d}_1 \in \mathbb{R}^L$, $\mathbf{d}_2 \in \mathbb{R}^N$ and such that $\mathbf{d} = [\mathbf{d}_1^\top, \mathbf{d}_2^\top]^\top$. Note that the calculation of $\mathbf{K} \mathbf{d}_2$ can be accelerated by Nyström methods. Accordingly, the computational cost of solving (10) by LR-SDCut is:

$$\underbrace{\left(\mathcal{O} \left((N+L)R_Y^2 + \underbrace{(NR_K + NL + L^2)}_{\text{matrix-vector product (15)}} R_Y \right) \right)}_{\text{Lanczos factorization}} \times \# \text{Lanczos-Iters} \times \# \text{Descent-Iters}, \quad (16)$$

and the memory requirement is $\mathcal{O}(N(L + R_Y + R_K) + LR_Y)$, where R_K and R_Y denotes the rank of \mathbf{K} and $(\mathbf{C}(\mathbf{u}))_+$ respectively. As mean field approximation, the computational complexity is also linear in N .

5. Applications

To show the superiority of the proposed method, we evaluate it and other methods on two applications in this section: image segmentation and image co-segmentation. In the following our experiments, the maximum number of iterations K_{\max} for LR-SDCut is set to 10; the initial rank r is set to 20; and the penalty parameter γ is set to 1000.

5.1. Application 1: Image Segmentation

Following the work in [11], pairwise potentials for image segmentation are expressed in the following form:

$$K_{i,j}^{(1)} = \exp \left(-\frac{|\mathbf{p}_i - \mathbf{p}_j|^2}{2\theta_\alpha^2} - \frac{|\mathbf{c}_i - \mathbf{c}_j|^2}{2\theta_\beta^2} \right), \quad (17)$$

where \mathbf{p}_i and \mathbf{c}_i are the position and color value of pixel i respectively, and similarly for \mathbf{p}_j and \mathbf{c}_j . The matrix defined in (17) corresponds to the appearance kernel which penalizes the case that two adjacent pixels with similar color and different labels. The label compatibility function is given by the Potts model $\mu(l, l') = \delta(l \neq l')$.

The kernel matrix $\mathbf{K}^{(1)}$ can be decomposed to the Hadamard product of two independent kernel matrices: $\mathbf{K}^{(1)} = \mathbf{K}_p^{(1)} \circ \mathbf{K}_c^{(1)}$, where $k_p^{(1)}(\mathbf{f}_i, \mathbf{f}_j) = \exp \left(\frac{-|\mathbf{p}_i - \mathbf{p}_j|^2}{2\theta_\alpha^2} \right)$ and $k_c^{(1)}(\mathbf{f}_i, \mathbf{f}_j) = \exp \left(\frac{-|\mathbf{c}_i - \mathbf{c}_j|^2}{2\theta_\beta^2} \right)$.

Nyström methods are performed on $\mathbf{K}_p^{(1)}$ and $\mathbf{K}_c^{(1)}$ individually: $\mathbf{K}_p^{(1)} \approx \Phi_p \Phi_p^\top$ and $\mathbf{K}_c^{(1)} \approx \Phi_c \Phi_c^\top$, where

	Unary	MF+filter	MF+Nys.	LR-SDCut
Time(s)	NA	0.29	6.6	74
Accu.	0.79	0.83	0.83	0.83
Energy	$1.29 \cdot 10^5$	$9.79 \cdot 10^4$	$1.15 \cdot 10^5$	$9.02 \cdot 10^4$

Table 2: Quantitative results of image segmentation. Our method runs slower than mean field methods but gives significantly lower energy. Unfortunately, the lower energy does not lead to better segmentation accuracy.

$\Phi_p \in \mathbb{R}^{N \times R_p}$ and $\Phi_c \in \mathbb{R}^{N \times R_c}$. Then we have:

$$\mathbf{K}^{(1)} \mathbf{d} = (\mathbf{K}_p^{(1)} \circ \mathbf{K}_c^{(1)}) \mathbf{d} \quad (18a)$$

$$= \text{diag} \left(\Phi_p \Phi_p^\top \text{Diag}(\mathbf{d}) \Phi_c \Phi_c^\top \right) \quad (18b)$$

$$= \left(\left(\Phi_p \Phi_p^\top (\text{Diag}(\mathbf{d}) \Phi_c) \right) \circ \Phi_c \right) \mathbf{1}. \quad (18c)$$

This computation requires $\mathcal{O}(NR_c R_p)$ operations (R_c and R_p are set to 20 and 10 respectively). Performing Nyström methods on $\mathbf{K}_p^{(1)}$ and $\mathbf{K}_c^{(1)}$ separately instead of on $\mathbf{K}^{(1)}$ directly brings two benefits: 1) The memory requirement may be reduced from $R_c R_p$ to $R_c + R_p$; 2) For multiple images with the same size, we only need to perform Nyström on $\mathbf{K}_p^{(1)}$ once, as the input features (positions $\mathbf{p}_i, i = 1, \dots, N$) are the same for these images.

The improved Nyström method [33] is adopted to obtain the low rank approximation of $\mathbf{K}_c^{(1)}$ and $\mathbf{K}_p^{(1)}$. k -means clustering is used in [33] to select representative landmarks.

Experiments The proposed algorithm is compared with mean field on MSRC 21-class database. The test data are 93 representative images with accurate ground truth provided by [11]. The unary potentials are also obtained from [11]. The parameters θ_α , θ_β and $w^{(1)}$ are set to 60, 20 and 10 respectively. The iteration number limit for mean field inference is set to 20. All experiments are conducted using a single CPU with 10GB memory. As for the matrix-vector product in the mean field method, both the filter-based and Nyström-based approaches are evaluated (refer to as MF+filter and MF+Nys. respectively). The evaluated images have around 60,000 pixels and so the number of MRF variables is also around 60,000 for each image.

Our method achieves similar segmentation results to the mean field approach. In Table 2, quantitative results are demonstrated. Although the computational complexity of mean field and our method are both linear in N , mean field is still faster than ours in this experiment. This is partially because the code of mean field is highly optimized using C++, while ours is unoptimized. A speed up is expected if our code is further optimized and parallelized. Note that the filter-based method [19] can be also incorporated into our algorithm to compute matrix-vector products, which is likely to be faster than Nyström methods but limited to Gaussian kernels in general.

Despite the slower speed, *our method achieves significantly lower energy value than mean field*, which shows that our method is better from the viewpoint of energy minimization. Unfortunately, the lower energy of our solution

Data	#pics	N	LR-SDCut	MF+Nys.	N	SDLR	SDCut
Cow	10	681600	1415s	1965s	6713	9530s	307s
Sheep	8	545280	1066s	2045s	5375	6932s	583s
Tree	9	613440	1137s	1490s	6026	1090s	1316s

Table 3: Running times for image co-segmentation. Our method is slightly faster than mean field. The number of MRF variables N for two groups of evaluated methods are shown in the third and sixth columns. The problems solved by our approach are much larger than those of SDLR and SDCut.

	LR-SDCut	MF+Nys.	SDLR	SDCut
Cow	0.73 ($-1.59 \cdot 10^5$)	0.67($-1.58 \cdot 10^5$)	0.66	0.69
Sheep	0.74 ($-8.07 \cdot 10^4$)	0.49($-6.87 \cdot 10^4$)	0.57	0.58
Tree	0.83 ($-2.23 \cdot 10^5$)	0.65($-2.03 \cdot 10^5$)	0.66	0.68

Table 4: Segmentation accuracy (energy) of image co-segmentation. Our method and mean field work on original pixels, while SDLR and SDCut work on superpixels. For all the three evaluated datasets, our method achieves the lowest energies and highest segmentation scores.

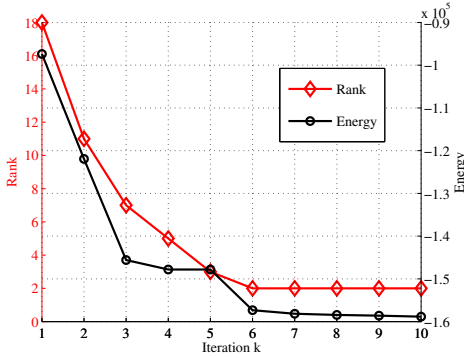


Figure 2: Rank and energy at each iteration for co-segmentation on the “cow” data set. Both of the rank of $(\mathbf{C}(\mathbf{u}^{(k)}))_+$ and the energy of binary solution $\mathbf{y}^{(k)}$ decrease significantly in the first several iterations.

does not lead to better segmentation performance. Actually, all of the evaluated methods have similar segmentation accuracy.

5.2. Application 2: Image Co-segmentation

The image co-segmentation problem requires that the same object be segmented from multiple images. There are two optimization criteria: the color and spatial consistency within one image and the separability between foreground and background over all images. There is no unary potentials for image co-segmentation and the pairwise potentials are shown in the following:

$$K_{i,j}^{(1)} = \varphi_{ij} \exp \left(-\frac{|\mathbf{p}_i - \mathbf{p}_j|^2}{2\theta_\alpha^2} - \frac{|\mathbf{c}_i - \mathbf{c}_j|^2}{2\theta_\beta^2} \right), \quad (19a)$$

$$\mathbf{K}^{(2)} = \Omega_N (\kappa N \mathbf{I}_N + \tilde{\mathbf{K}}^{(2)})^{-1} \Omega_N, \quad (19b)$$

where $\varphi_{ij} = 1$ if pixels i and j locate in the same image; $\varphi_{ij} = 0$, otherwise. $\kappa > 0$ is a regularization parameter. $\mathbf{K}^{(1)}$ is a block-diagonal matrix, and the matrix-vector product for $\mathbf{K}^{(1)}$ can be computed using the method described in Section 5.1. $\mathbf{K}^{(2)}$ is the inter-image discriminative clustering cost matrix (see [36] for details). $\Omega_N = \mathbf{I}_N - \frac{1}{N} \mathbf{1}\mathbf{1}^\top$ is the centering projection matrix, and $\tilde{\mathbf{K}}^{(2)}$ is the χ^2 kernel matrix of sift features. $\tilde{\mathbf{K}}^{(2)}$ can be approxi-

mated by a low-rank decomposition: $\tilde{\mathbf{K}}^{(2)} \approx \tilde{\Phi} \tilde{\Phi}^\top$, where $\tilde{\Phi} \in \mathbb{R}^{N \times R_{K_2}}$. Based on the matrix inversion lemma, we have:

$$\begin{aligned} \mathbf{K}^{(2)} &= \frac{1}{\kappa N} \Omega_N \left(\mathbf{I}_N - \underbrace{\tilde{\Phi} (\kappa N \mathbf{I}_D + \tilde{\Phi}^\top \tilde{\Phi})^{-1} \tilde{\Phi}^\top}_{\text{decompose to } \tilde{\Phi} \tilde{\Phi}^\top \text{ and } \tilde{\Phi} \tilde{\Phi}^\top \mathbf{1}=0} \right) \Omega_N \\ &= \frac{1}{\kappa N} (\Omega_N - \tilde{\Phi} \tilde{\Phi}^\top). \end{aligned} \quad (20)$$

Through the above equation, the matrix-vector product for $\mathbf{K}^{(2)}$ can be computed efficiently in the complexity of $O(NR_{K_2})$ (R_{K_2} is set to 640 in the experiments). The pairwise potentials are not necessarily submodular, because entries of $\mathbf{K}^{(2)}$ may be negative. Note that *the matrix-vector product for $\mathbf{K}^{(2)}$ cannot be performed by the filter-based method of [11], because $\mathbf{K}^{(2)}$ may not be a Gaussian kernel.*

Experiments Three groups of images are selected from the MSRC dataset for image co-segmentation. Besides our approach and mean field, the SDP-based algorithms in [36] (denoted as SDLR) and [16] (denoted as SDCut) are also evaluated. Our method and mean field are evaluated at the original pixel level, while SDLR and SDCut are evaluated only on superpixels.

The code for SDLR and SDCut is provided by authors of the original papers. The default settings are used. The iteration limit for mean field is set to 100. To prevent mean field from converging to undesirable local optima, we randomly run the method 5 times. All experiments are conducted on a single CPU with 20GB memory. The *intersection-over-union* accuracy is used to measure the segmentation performance.

From the results illustrated in Fig. 1, we see that our approach achieves much more accurate co-segmentation results than both SDLR and SDCut. The performance of mean field is also worse than ours.

Table 3 demonstrates the number of variables and computational time for each method. The variable numbers of the problems solved by our method and mean field are around 100 times larger than those for SDLR and SDCut. Our approach is slightly faster than mean field, and significantly more scalable than SDLR and SDCut.

The quantitative performance is shown in Table 4. Our approach achieves significantly better co-segmentation accuracy than all the other methods. As for energy minimization, our approach also produces lower energies than mean field. Empirically, we found mean field is sensitive to initialization. Take “tree” as example, the difference is $5.3 \cdot 10^4$ between the best and worst energy in the 5 repeats of mean field with random initializations. If we repeat mean field 100 times, the best energy improves from $-2.03 \cdot 10^5$ to $-2.08 \cdot 10^5$, but still worse than ours ($-2.23 \cdot 10^5$).

Figure 2 shows the change of $\text{rank}((\mathbf{C}(\mathbf{u}^{(k)}))_+)$ and $\tilde{E}(\mathbf{y}^{(k)})$ w.r.t. iteration k . Both of the rank and energy drops

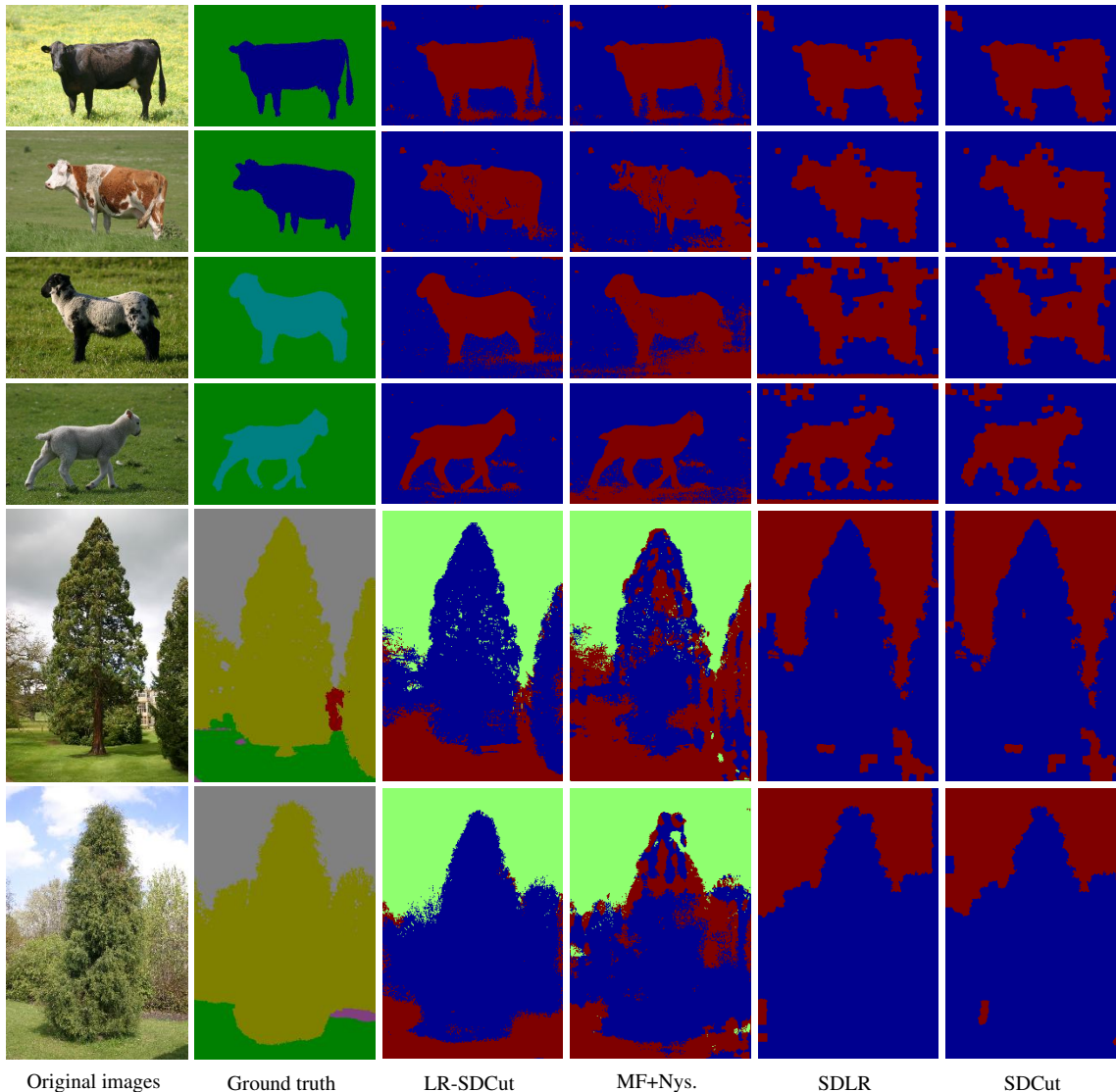


Figure 1: Qualitative results for image co-segmentation. Three classes of objects from MSRC datasets are used for the evaluation. Our approach and Mean Field (MF+Nys.) are performed on the original pixel-level images. Because SDLR [36] and SDCut [36] cannot scale up to pixel-level images, they are evaluated on superpixels. Our method performs best visually. We randomly repeat mean field approximation 5 times for each dataset and select the best result. Mean field is not stable at this task and sometimes converges to an undesirable local optimal point (see “tree” for example). SDLR and SDCut achieve worse results than our’s, since some image details are lost due to the use of superpixels.

quickly in the first several iterations. Simultaneously, the lower-bound of the optimal energy $\bar{E}(\mathbf{y})$ (*i.e.* the dual objective value) increases from $-8.09 \cdot 10^7$ to $-4.36 \cdot 10^5$.

6. Conclusions

In this paper, we have proposed an efficient, general method for solving fully-connected CRFs. The proposed SDP approach is more stable and accurate than mean field approximation, which is also more scalable than previous SDP methods. The use of low-rank approximation of the kernel matrix to perform matrix-vector products makes our approach even more efficient and applicable for any sym-

metric positive semidefinite kernel. In contrast, previous filter-based methods assume pairwise potentials to be based on a Gaussian or generalized RBF kernel. The computational complexity of our approach is linear in the number of CRF variables. The experiments on image co-segmentation validate that our approach can be applied to more general problems than previous methods.

As for future works, the proposed method can be parallelized to achieve even faster speed. The core of our method is quasi-Newton (or gradient descent) and eigen-decomposition, both of which can be parallelized on GPUs. Matrix-vector products, the main computational cost, can be implemented using CUDA function “cublasSgemm”.

Acknowledgements

This work was in part funded by the Data to Decisions Cooperative Research Centre, Australia. Correspondence should be addressed to C. Shen (e-mail: chhshen@gmail.com).

References

- [1] J. H. Kappes, B. Andres, F. A. Hamprecht, C. Schnorr, S. Nowozin, D. Batra, S. Kim, B. X. Kausler, J. Lellmann, N. Komodakis, et al., “A comparative study of modern inference techniques for discrete energy minimization problems,” in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2013.
- [2] 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.*, vol. 30, no. 6, pp. 1068–1080, 2008.
- [3] M. P. Kumar, V. Kolmogorov, and P. H. Torr, “An analysis of convex relaxations for MAP estimation of discrete MRFs,” *J. Mach. Learn. Res.*, vol. 10, pp. 71–106, 2009.
- [4] V. Kolmogorov, “Convergent tree-reweighted message passing for energy minimization,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 28, no. 10, pp. 1568–1583, 2006.
- [5] J. H. Kappes, B. Savchynskyy, and C. Schnorr, “A bundle approach to efficient MAP-inference by Lagrangian relaxation,” in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2012.
- [6] C. Rother, V. Kolmogorov, V. Lempitsky, and M. Szummer, “Optimizing binary MRFs via extended roof duality,” in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2007.
- [7] Y. Boykov, O. Veksler, and R. Zabih, “Fast approximate energy minimization via graph cuts,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 23, no. 11, pp. 1222–1239, 2001.
- [8] V. Kolmogorov and R. Zabih, “What energy functions can be minimized via graph cuts?,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 26, no. 2, pp. 147–159, 2004.
- [9] P. F. Felzenszwalb and D. P. Huttenlocher, “Efficient belief propagation for early vision,” *Int. J. Comp. Vis.*, vol. 70, no. 1, pp. 41–54, 2006.
- [10] M. J. Wainwright, T. S. Jaakkola, and A. S. Willsky, “MAP estimation via agreement on trees: message-passing and linear programming,” *IEEE T. Information Theory*, vol. 51, no. 11, pp. 3697–3717, 2005.
- [11] P. Krähenbühl and V. Koltun, “Efficient inference in fully connected CRFs with Gaussian edge potentials,” in *Proc. Adv. Neural Inf. Process. Syst.*, 2011.
- [12] P. Krähenbühl and V. Koltun, “Parameter learning and convergent inference for dense random fields,” in *Proc. Int. Conf. Mach. Learn.*, 2013.
- [13] N. D. Campbell, K. Subr, and J. Kautz, “Fully-connected CRFs with non-parametric pairwise potentials,” in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2013.
- [14] Y. Zhang and T. Chen, “Efficient inference for fully-connected CRFs with stationarity,” in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2012.
- [15] Q. Huang, Y. Chen, and L. Guibas, “Scalable semidefinite relaxation for maximum a posterior estimation,” in *Proc. Int. Conf. Mach. Learn.*, 2014.
- [16] P. Wang, C. Shen, and A. Hengel, “A fast semidefinite approach to solving binary quadratic problems,” in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2013.
- [17] S. I. Wang, R. Frostig, P. Liang, and C. D. Manning, “Relaxations for inference in restricted Boltzmann machines,” in *International Conference on Learning Representations*, 2014.
- [18] R. Frostig, S. I. Wang, P. S. Liang, and C. D. Manning, “Simple MAP inference via low-rank relaxations,” in *Proc. Adv. Neural Inf. Process. Syst.*, 2014.
- [19] A. Adams, J. Baek, and M. A. Davis, “Fast high-dimensional filtering using the permutohedral lattice,” in *Eurographics*, 2010.
- [20] D. Koller and N. Friedman, *Probabilistic graphical models: principles and techniques*, MIT press, 2009.
- [21] S. Paris and F. Durand, “A fast approximation of the bilateral filter using a signal processing approach,” in *Proc. Eur. Conf. Comp. Vis.*, 2006.
- [22] C. Williams and M. Seeger, “The effect of the input density distribution on kernel-based classifiers,” in *Proc. Int. Conf. Mach. Learn.*, 2000.
- [23] C. Williams and M. Seeger, “Using the Nyström method to speed up kernel machines,” in *Proc. Adv. Neural Inf. Process. Syst.*, 2001.
- [24] P. Drineas and M. W. Mahoney, “On the Nyström method for approximating a gram matrix for improved kernel-based learning,” *J. Mach. Learn. Res.*, vol. 6, pp. 2153–2175, 2005.
- [25] S. Fine and K. Scheinberg, “Efficient SVM training using low-rank kernel representations,” *J. Mach. Learn. Res.*, vol. 2, pp. 243–264, 2002.
- [26] F. R. Bach and M. I. Jordan, “Kernel independent component analysis,” *J. Mach. Learn. Res.*, vol. 3, pp. 1–48, 2003.
- [27] A. Rahimi and B. Recht, “Random features for large-scale kernel machines,” in *Proc. Adv. Neural Inf. Process. Syst.*, 2007.
- [28] A. Rahimi and B. Recht, “Weighted sums of random kitchen sinks: Replacing minimization with randomization in learning,” in *Proc. Adv. Neural Inf. Process. Syst.*, 2008.

- [29] A. Vedaldi and A. Zisserman, “Efficient additive kernels via explicit feature maps,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 34, no. 3, pp. 480–492, 2012.
- [30] F. Bach, “Sharp analysis of low-rank kernel matrix approximations,” *J. Mach. Learn. Res.*, vol. 30, pp. 185–209, 2013.
- [31] A. Gittens and M. W. Mahoney, “Revisiting the Nyström method for improved large-scale machine learning,” in *Proc. Int. Conf. Mach. Learn.*, 2013.
- [32] T. Yang, Y.-F. Li, M. Mahdavi, R. Jin, and Z.-H. Zhou, “Nyström method vs random Fourier features: A theoretical and empirical comparison,” in *Proc. Adv. Neural Inf. Process. Syst.*, 2012.
- [33] K. Zhang, I. W. Tsang, and J. T. Kwok, “Improved Nyström low-rank approximation and error analysis,” in *Proc. Int. Conf. Mach. Learn.*, 2008.
- [34] D. C. Sorensen, *Implicitly restarted Arnoldi/Lanczos methods for large scale eigenvalue calculations*, Springer, 1997.
- [35] J. Briët, F. M. de Oliveira Filho, and F. Vallentin, “The positive semidefinite Grothendieck problem with rank constraint,” in *Automata, Languages and Programming*, pp. 31–42. Springer, 2010.
- [36] A. Joulin, F. Bach, and J. Ponce, “Discriminative clustering for image co-segmentation,” in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2010.