



Parameter Estimation and Energy Minimization for Region-based Semantic Segmentation

M. Pawan Kumar, Haithem Turki, Dan Preston, Daphne Koller

► To cite this version:

M. Pawan Kumar, Haithem Turki, Dan Preston, Daphne Koller. Parameter Estimation and Energy Minimization for Region-based Semantic Segmentation. [Technical Report] 2013. hal-00857918

HAL Id: hal-00857918

<https://inria.hal.science/hal-00857918>

Submitted on 4 Sep 2013

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.

Parameter Estimation and Energy Minimization for Region-based Semantic Segmentation

M. Pawan Kumar, Haithem Turki, Dan Preston, and Daphne Koller

Abstract—We consider the problem of parameter estimation and energy minimization for a region-based semantic segmentation model. The model divides the pixels of an image into non-overlapping connected regions, each of which is assigned a label indicating its semantic class. In the context of energy minimization, the main problem we face is the large number of putative pixel-to-region assignments. We address this problem by designing an accurate linear programming based approach for selecting the best set of regions from a large dictionary. The dictionary is constructed by merging and intersecting segments obtained from multiple bottom-up over-segmentations. The linear program is solved efficiently using dual decomposition. In the context of parameter estimation, the main problem we face is the lack of fully supervised data. We address this issue by developing a principled framework for parameter estimation using diverse data. More precisely, we propose a latent structural support vector machine formulation, where the latent variables model any missing information in the human annotation. Of particular interest to us are three types of annotations: (i) images segmented using generic foreground or background classes; (ii) images with bounding boxes specified for objects; and (iii) images labeled to indicate the presence of a class. Using large, publicly available datasets we show that our methods are able to significantly improve the accuracy of the region-based model.



1 INTRODUCTION

SEMANtic segmentation offers a useful representation of the scene depicted in an image by assigning each pixel to a specific semantic class (for example, ‘person’, ‘building’ or ‘tree’). It is a long standing goal of computer vision, and is an essential building block of many of the most ambitious applications such as autonomous driving, robot navigation, content-based image retrieval and surveillance. Its importance can be assessed by its long history in the computer vision literature.

Encouraged by their success in low-level vision applications such as image denoising and stereo reconstruction [1], earlier efforts focused on pixel-level models, where each pixel was assigned a label using features extracted from a regularly shaped patch around it [2], [3], or at an offset from it [4]. However, the features extracted from such patches are not reliable in the presence of background clutter. For example, a patch around a boundary pixel of a tree may contain sky or building pixels. This may inhibit an algorithm from using the fact that trees are mostly green.

To avoid the problem of pixel-based methods, researchers have started to develop region-based models. Such models divide an image into regions, where

each region is a set of connected pixels. The reasoning is that if the regions are large enough to provide reliable discriminative features, but small enough so that all the pixels within a region belong to the same semantic class, then we can obtain an accurate segmentation of the image. The first region-based models [5], [6], [7] for high-level vision defined the regions of the image as the segments obtained using a standard bottom-up over-segmentation approach [8], [9]. However, since over-segmentation approaches are agnostic to the task at hand, these regions may not capture the boundaries between the scene entities accurately. To address this issue, some works have suggested heuristics for selecting a good over-segmentation [5], [6]. However, even the best over-segmentation approach is unlikely to provide regions of sufficient accuracy.

In order to obtain regions that capture the boundaries accurately, some efforts have been made to combine multiple over-segmentations. For example, Pontafaru *et al.* [10] suggested taking the intersection of multiple over-segmentation. However, such an approach results in very small regions. Other models suggest using overlapping regions [11], [12], [13], [14]. However, the pixels within the overlap can support two contradicting hypotheses (that is, they can belong to two different semantic classes) thereby overcounting the data.

In this work, we use the region-based model of Gould *et al.* [15], which consists of two layers. The variables of the first layer correspond to the image pixels. A labeling of the first layer (that is, an assignment of values to the variables) divides the

- M. Pawan Kumar is with the Center for Visual Computing, Ecole Centrale Paris, Châtenay-Malabry, 92295, France.
E-mail: pawan.kumar@ecp.fr
- Haithem Turki, Dan Preston and Daphne Koller are with the Computer Science Department, Stanford University, Stanford, CA, 94305.
E-mail: {hturki,dpreston,koller}@cs.stanford.edu

pixels into contiguous, non-overlapping regions. The variables of the second layer correspond to the regions defined by the first layer. A labeling of the second layer assigns each region to a unique semantic class, thereby providing the segmentation of the image. The real-valued energy function of the model provides the desirability of an overall labeling (that is, the combined labeling of the first and second layers), that is, the lower the energy the better the labeling. Given an image, its semantic segmentation is inferred by minimizing the energy over all possible labelings. The advantage of the above formulation is two-fold: (i) in addition to the segmentation, the regions themselves would be inferred (including the number of regions), which implies that unlike other models, the regions would be task dependent; and (ii) the regions would be connected and non-overlapping, thereby avoiding data overcounting.

While the region-based model of Gould *et al.* [15] possesses several desirable qualities, its practical deployment poses two formidable problems. First, minimizing the energy associated with this model is extremely challenging due to the large number of possible pixel-to-region assignments. Furthermore, the regions are required to be connected, a difficult constraint to impose [16], [17]. Second, the energy function consists of several thousand parameters that need to be estimated accurately in order to obtain good segmentations via energy minimization.

In order to address the difficulty posed by energy minimization, Gould *et al.* [15] proposed a method that constructs a large dictionary of putative regions using multiple over-segmentations obtained by changing the parameters of a bottom-up approach. Specifically, the putative regions are defined as sets of pixels obtained by merging and intersecting the segments with each other. While merging segments together provides large regions, their intersection with small segments ensures that they align well with the boundary. The set of non-overlapping regions of the image are selected from the dictionary by minimizing the energy function using a simple greedy algorithm. However, while we would expect the dictionary itself to contain accurate regions, the greedy algorithm is susceptible to getting stuck in a bad local minimum. In order to alleviate this issue, we formulate the problem of selecting the best set of regions from a large dictionary as an integer program and design an accurate linear programming (LP) relaxation for it. Furthermore, we show how the LP relaxation can be solved efficiently by suitably modifying the dual decomposition framework [18].

In order to estimate the parameters of the region-based model, Gould *et al.* [15] devised a piecewise learning approach, called closed loop learning, which relies on a fully supervised training dataset. However, we argue that their approach suffers from two significant drawbacks. First, piecewise learning can

result in a bad solution. Second, and more important, the collection of fully supervised data is an onerous task as reflected by the small size of the current datasets for semantic segmentation [15], [19]. In order to overcome these problems, we propose the use of diverse data, where the level of annotation varies among the training samples, from pixelwise segmentation to bounding boxes and image-level labels. We design a principled framework for learning with diverse data, with the aim of exploiting the varying degrees of information in the different datasets to the fullest. Specifically, we formulate the parameter learning problem using a latent structural support vector machine (LSVM) [20], [21], where the latent variables model any missing information in the annotation. In order to optimize the objective function of LSVM using the self-paced learning algorithm [22], we modify our LP relaxation based energy minimization algorithm such that it can complete the annotation of weakly supervised images.

We demonstrate the efficacy of our novel energy minimization and parameter estimation methods using large, publicly available datasets. Earlier versions of this article have appeared as [23], [24].

2 THE REGION-BASED MODEL

We begin by providing a formal description of the two-layer region-based model of Gould *et al.* [15]. The first layer of the model consists of random variables \mathcal{V}^P that correspond to the pixels of a given image \mathbf{X} . Each random variable can take one label from the set $\mathcal{L}^P = \{1, 2, \dots, R\}$ that represents the region to which it belongs. Here R is the total number of regions in a given image (which has to be inferred automatically). A labeling of the first layer is a vector \mathbf{Y}^P that divides the pixels into connected, non-overlapping regions. The second layer consists of random variables $\mathcal{V}^R(\mathbf{Y}^P)$ that correspond to the regions defined by \mathbf{Y}^P . Each random variable can take one label from the set $\mathcal{L}^R = \{1, 2, \dots, L\}$ where L is the total number of semantic classes we consider. The desired segmentation is provided by a labeling \mathbf{Y}^R of the second layer. A labeling of the entire model is denoted by $\mathbf{Y} = (\mathbf{Y}^P, \mathbf{Y}^R)$. The energy of the model consists of two types of potentials, unary and pairwise, which we defined in terms of the image features and the model parameters. The potentials and the energy function are described in detail in the following subsections.

2.1 Unary Potential

For each variable $v_r \in \mathcal{V}^R(\mathbf{Y}^P)$ (corresponding to a region r) we define a unary potential $\theta_r(\mathbf{Y}_r^R; \mathbf{X})$ for assigning it the label \mathbf{Y}_r^R . The value of the unary potential depends on the parameters as well as the features extracted from the image. In more detail, let $\mathbf{u}_r(\mathbf{X})$ denote the features extracted from the pixels belonging to the region r , which can capture shape, appearance and texture information (for example,

green regions are likely to be grass or tree, while blue regions are likely to be sky). We refer the interested reader to [15] for details regarding the features used in this paper. The unary potential of assigning the region r to a semantic class i is defined as

$$\theta_r(i; \mathbf{X}) = \mathbf{w}_i^\top \mathbf{u}_r(\mathbf{X}), \quad (1)$$

where \mathbf{w}_i is the unary parameter corresponding to the class i .

2.2 Pairwise Potential

For each pair of neighboring variables, whose corresponding regions share at least one boundary pixel, we define a pairwise potential $\theta_{rr'}(\mathbf{Y}_r^R, \mathbf{Y}_{r'}^R; \mathbf{X})$ for assigning labels \mathbf{Y}_r^R and $\mathbf{Y}_{r'}^R$ to v_r and $v_{r'}$ respectively. Similar to the unary potential, the value of the pairwise potential depends on the parameters and the image features. In more detail, let $\mathbf{p}_{rr'}(\mathbf{X})$ refer to the features extracted using the pixels belonging to regions r and r' , which can capture contrast and contextual information (for example, boats are likely to be above water, while cars are likely to be above road). We refer the interested reader to [15] for details regarding the features used in this paper. The pairwise potential for assigning the regions r and r' to semantic classes i and j respectively is defined as

$$\theta_{rr'}(i, j; \mathbf{X}) = \mathbf{w}_{ij}^\top \mathbf{p}_{rr'}(\mathbf{X}), \quad (2)$$

where \mathbf{w}_{ij} is the pairwise parameter corresponding to the classes i and j .

2.3 Energy Function

Given an image \mathbf{X} and a labeling \mathbf{Y} , we define the unary feature for a semantic class i as

$$\Psi_i(\mathbf{X}, \mathbf{Y}) = \sum_{r=1}^R \delta(\mathbf{Y}_r^R = i) \mathbf{u}_r(\mathbf{X}), \quad (3)$$

that is, it is the sum of the features over all the regions that are assigned the class i in the labeling \mathbf{Y} . Similarly, we define the pairwise feature for the semantic classes i and j as

$$\Psi_{ij}(\mathbf{X}, \mathbf{Y}) = \sum_{(r, r') \in \mathcal{E}^R(\mathbf{Y}^P)} \delta(\mathbf{Y}_r^R = i) \delta(\mathbf{Y}_{r'}^R = j) \mathbf{p}_{rr'}(\mathbf{X}), \quad (4)$$

that is, it is the sum of the features over all pairs of neighboring regions that are assigned the classes i and j respectively in the labeling \mathbf{Y} . Here, $\mathcal{E}^R(\mathbf{Y}^P)$ is the set of pairs of regions that share at least one boundary pixel, as defined by \mathbf{Y}^P . Using the unary and pairwise features, we define the joint feature vector of the input \mathbf{X} and output \mathbf{Y} of the model as

$$\Psi(\mathbf{X}, \mathbf{Y}) = [\Psi_i(\mathbf{X}, \mathbf{Y}), \forall i; \Psi_{ij}(\mathbf{X}, \mathbf{Y}), \forall i, j]. \quad (5)$$

In other words, the joint feature vector of the input and the output is the concatenation of the unary features for all semantic classes and the pairwise

features for all pairs of semantic classes. Similarly, we define the parameter of the model as

$$\mathbf{w} = [\mathbf{w}_i, \forall i; \mathbf{w}_{ij}, \forall i, j], \quad (6)$$

that is, the parameter is the concatenation of the unary parameters for all semantic classes and the pairwise parameters for all pairs of semantic classes.

The energy of the model can then be written concisely using the parameter and the joint feature vector as follows:

$$\begin{aligned} E(\mathbf{Y}; \mathbf{X}, \mathbf{w}) &= \mathbf{w}^\top \Psi(\mathbf{X}, \mathbf{Y}) \\ &= \sum_{v_r \in \mathcal{V}^R(\mathbf{Y}^P)} \theta_r(\mathbf{Y}_r^R) + \sum_{(v_r, v_{r'}) \in \mathcal{E}^R(\mathbf{Y}^P)} \theta_{rr'}(\mathbf{Y}_r^R, \mathbf{Y}_{r'}^R), \end{aligned} \quad (7)$$

Note that we have dropped the input \mathbf{X} from the notation of the individual potentials to avoid clutter.

3 ENERGY MINIMIZATION

For the model described above, we now consider the problem of obtaining the labeling \mathbf{Y}^* that minimizes the energy function. This will provide us with the best set of regions for the task (\mathbf{Y}^P) as well as their labels (\mathbf{Y}^R). In this section, we assume that the parameters \mathbf{w} have been provided. In the next section, we will describe a max-margin framework for estimating the parameters using a diversely annotated training dataset.

As noted earlier, the main difficulty in energy minimization arises due to the fact that there are many possible labelings \mathbf{Y}^P that group pixels into regions. Specifically, for a given $H \times W$ image there can be as many as HW regions (that is, each pixel is a region). Hence the total number of possible labelings \mathbf{Y}^P is $(HW)^{(HW)}$. Furthermore, we also need to ensure that the inferred labeling \mathbf{Y}^P provides connected regions, which is well-known to be a difficult constraint to impose [16], [17].

In order to overcome these problems, we make use of bottom-up over-segmentation approaches. Specifically, we minimize the energy using the following two steps: (i) construct a large dictionary of connected putative regions using multiple over-segmentations; and (ii) select the set of regions that minimize the energy (that is, infer \mathbf{Y}^P and \mathbf{Y}^R). We begin by describing our algorithm for selecting regions from a dictionary. We then provide details of the dictionary that we found to be effective in our experiments.

3.1 Region Selection as Optimization

Given a dictionary of regions, we wish to select a subset of regions such that: (i) the entire image is explained by the selected regions; (ii) no two selected regions overlap with each other; and (iii) the energy $E(\mathbf{Y}; \mathbf{X}, \mathbf{w})$ is minimized. Note that the dictionary itself may contain overlapping regions of any shape

or size. We do not place any restrictions on it other than the assumption that it contains at least one set of disjoint regions that explains the entire image. We formulate the above task as an integer program and provide an accurate linear programming (LP) relaxation for it.

3.1.1 Integer Programming Formulation

Before describing the integer program, we need to set up some notation. We denote the dictionary of regions by \mathcal{D} . The intersection of all the regions in \mathcal{D} defines a set of super-pixels \mathcal{S} . The set of all regions that contain a super-pixel $s \in \mathcal{S}$ is denoted by $\mathcal{C}(s) \subseteq \mathcal{D}$. Finally, the set of neighboring regions is denoted by \mathcal{E} , where two regions r and r' are considered neighbors of each other (that is, $(r, r') \in \mathcal{E}$) if they do not overlap and share at least one boundary pixel.

To formulate our problem as an integer program we define binary variables $y_r(i)$ for each region $r \in \mathcal{D}$ and $i \in \mathcal{L}^I = \mathcal{L}^R \cup \{0\}$. These variables indicate whether a particular region is selected and if so, which label it takes. Specifically, if $y_r(0) = 1$ then the region r is not selected, else if $y_r(i) = 1$ where $i \in \mathcal{L}^R$ then the region r is assigned the label i . Similarly, we define binary variables $y_{rr'}(i, j)$ for all neighboring regions $(r, r') \in \mathcal{E}$ such that $y_{rr'}(i, j) = y_r(i)y_{r'}(j)$. Furthermore, we define unary and pairwise potentials corresponding to the augmented label set \mathcal{L}^I as

$$\begin{aligned}\bar{\theta}_r(i) &= \begin{cases} \theta_r(i) & \text{if } i \in \mathcal{L}^R, \\ 0 & \text{otherwise,} \end{cases} \\ \bar{\theta}_{rr'}(i, j) &= \begin{cases} \theta_{rr'}(i, j) & \text{if } i, j \in \mathcal{L}^R, \\ 0 & \text{otherwise.} \end{cases} \end{aligned} \quad (8)$$

The problem of obtaining the desired subset of regions can then be formulated as the following integer program:

$$\begin{aligned} \min_{\mathbf{y} \in \text{SELECT}(\mathcal{D})} \quad & \sum_{r \in \mathcal{D}, i \in \mathcal{L}^I} \bar{\theta}_r(i) y_r(i) + \\ & \sum_{(r, r') \in \mathcal{E}, i, j \in \mathcal{L}^I} \bar{\theta}_{rr'}(i, j) y_{rr'}(i, j), \end{aligned} \quad (9)$$

where the feasible set $\text{SELECT}(\mathcal{D})$ is specified using the following constraints:

$$\begin{aligned} y_r(i), y_{rr'}(i, j) &\in \{0, 1\}, \forall r, r' \in \mathcal{D}, i, j \in \mathcal{L}^I, \\ \sum_{i \in \mathcal{L}^I} y_r(i) &= 1, \forall r \in \mathcal{D}, \\ \sum_{j \in \mathcal{L}^I} y_{rr'}(i, j) &= y_r(i), \forall (r, r') \in \mathcal{E}, i \in \mathcal{L}^I, \\ \sum_{i \in \mathcal{L}^I} y_{rr'}(i, j) &= y_{r'}(j), \forall (r, r') \in \mathcal{E}, j \in \mathcal{L}^I, \\ \sum_{r \in \mathcal{C}(s)} \sum_{i \in \mathcal{L}^R} y_r(i) &= 1, \forall s \in \mathcal{S}. \end{aligned} \quad (10)$$

The first set of constraints ensure that the variables \mathbf{y} are binary. The second constraint implies that each region r should be assigned one label from the set \mathcal{L}^I .

The third constraint enforces $y_{rr'}(i, j) = y_r(i)y_{r'}(j)$. The final constraint, which we call covering constraint, restricts each super-pixel to be covered by exactly one selected region.

3.1.2 Linear Programming Relaxation

Problem (9) is NP-hard since \mathbf{y} is constrained to be binary (which specifies a non-convex feasible region). However, we can obtain an approximate solution to the above problem by relaxing the constraint on \mathbf{y} such that $y_r(i)$ and $y_{rr'}(i, j)$ take (possibly fractional) values in the interval $[0, 1]$. The resulting LP relaxation is similar to the standard relaxation for energy minimization in pairwise random fields [25], [26], with the exception of the additional covering constraints. However, this relaxation is very weak when the pairwise potentials are not submodular (roughly speaking, when they encourage neighboring regions to take different labels) [27]. For example, consider the case where each region is either selected or not (that is, $|\mathcal{L}^I| = 2$). For two neighboring regions r and r' , the pairwise potential $\bar{\theta}_{rr'}(\cdot, \cdot)$ is 0 if one or both regions are not selected and $\bar{\theta}_{rr'}(1, 1)$ otherwise (as defined by equation (8)). If $\bar{\theta}_{rr'}(1, 1) > 0$ then the neighboring regions are encouraged to take different labels; that is, the pairwise potentials are non-submodular. This results in frustrated cycles for which the standard LP relaxation provides a weak approximation [28] (we tested this empirically for our problem, but do not include the results due to space limitations).

There are two common ways to handle non-submodular problems in the literature: (i) applying the roof duality relaxation [28], [29]; and (ii) using message passing algorithms [30], [31], [32] based on cycle inequalities [33]. Unfortunately, both these methods are not directly applicable in our case. Specifically, roof duality does not allow us to incorporate the covering constraints. Adding cycle inequalities still results in a weak approximation. For example, see Fig. 1 that shows a clique formed by three neighboring regions (r_1 , r_2 and r_3) along with all the regions that overlap with at least one of them. Consider the case where $|\mathcal{L}^I| = 2$ and $\bar{\theta}_{rr'}(1, 1) > 0$, thereby resulting in non-submodular pairwise potentials (shown in Fig. 1(b)). Since each super-pixel needs to be covered by exactly one selected region, it follows that the energy of the optimal assignment for this clique will be strictly greater than 0 (as at least two neighboring regions have to be selected). However, the LP relaxation, with all the cycle inequalities added in, still provides a fractional solution whose objective function value is 0 (shown in Fig. 1(c)).

Our example shows that cycle inequalities are not sufficient to make the relaxation tight. Instead, we require all the constraints that define the marginal polytope [26] (convex hull of the valid integral labelings) of the entire clique. We note here that the recent work of Sontag *et al.* [32] also advocates the use of

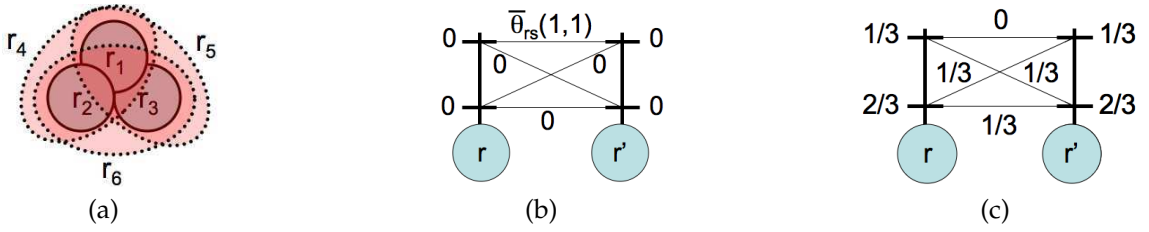


Fig. 1. (a) Neighboring regions r_1 , r_2 and r_3 , shown using solid lines, consist of a single super-pixel. The regions shown using dashed lines are formed by two super-pixels. Specifically, $r_4 = r_1 \cup r_2$, $r_5 = r_1 \cup r_3$ and $r_6 = r_2 \cup r_3$. (b) The potentials corresponding to the clique of size 6 formed by the regions. The branches (horizontal lines) along the trellis (the vertical line on top of a region) represent the different labels that each region may take. We consider a two label case here. The unary potential $\bar{\theta}_r(i)$ is shown next to the i^{th} branch of the trellis on top of region r . The pairwise potential $\bar{\theta}_{rr'}(i, j)$ is shown next to the connection between the i^{th} branch of r and the j^{th} branch of r' . The only non-zero potential $\bar{\theta}_{rr'}(1, 1) > 0$ corresponds to selecting both the regions r and r' . The optimal labeling of the clique must have an energy greater than 0 since at least two neighboring regions must be selected. (c) The optimal solution of the LP relaxation. The value of $y_r(i)$ is shown next to the i^{th} branch of r and the value of $y_{rr'}(i, j)$ is shown next to the connection between the i^{th} and j^{th} branches of r and r' respectively. Note that the solution satisfies all cycles inequalities, that is, $\sum_{(r,r') \in \mathcal{E}_C} y_{rr'}(0, 0) + y_{rr'}(1, 1) \geq 1$, where \mathcal{E}_C is a cycle. Hence the solution lies within the feasible region of the relaxation. However, it can be easily verified that its objective function value is 0, thereby proving that the relaxation is not tight.

such constraints. However, in their experiments they found cliques of size three (which are also a cycle of size three) to be sufficient. In contrast, we use cliques that are formed by three neighboring regions along with all the regions that overlap with at least one of the three regions. To the best of our knowledge, ours is one of the first examples in vision where constraints on cliques are essential for obtaining a good solution. Although a large number of constraints are required to specify the marginal polytope of a clique (their exact form is not important for this work), we show how the overall relaxation can be solved efficiently.

3.2 Solving the Relaxation

We use the dual decomposition framework that is well-known in the optimization community [18] and has recently been introduced in computer vision [34]. We begin by describing the general framework and then specify how it can be modified to efficiently solve our relaxation.

3.2.1 Dual Decomposition

Consider the following convex optimization problem: $\min_{\mathbf{z} \in \mathcal{F}} \sum_{k=1}^M g_k(\mathbf{z}_k)$, where \mathcal{F} represents the convex feasible region of the problem. The above problem is equivalent to the following: $\min_{\mathbf{z}_k \in \mathcal{F}, \mathbf{z}} \sum_k g_k(\mathbf{z}_k)$, s.t. $\mathbf{z}_k = \mathbf{z}$. Introducing the additional variables \mathbf{z}_k allows us to obtain the dual problem as

$$\max_{\boldsymbol{\lambda}_k} \min_{\mathbf{z}_k \in \mathcal{F}, \mathbf{z}} \sum_k g_k(\mathbf{z}_k) + \sum_k \boldsymbol{\lambda}_k(\mathbf{z}_k - \mathbf{z}), \quad (11)$$

where $\boldsymbol{\lambda}_k$ are the Lagrange multipliers. Differentiating the dual function with respect to \mathbf{z} we obtain the constraint that $\sum_k \boldsymbol{\lambda}_k = \mathbf{0}$, which implies that we can discard \mathbf{z} from the above problem. The simplified form of the dual suggests the following strategy for solving it. We start by initializing $\boldsymbol{\lambda}_k$ such that $\sum_k \boldsymbol{\lambda}_k = \mathbf{0}$. Keeping the values of $\boldsymbol{\lambda}_k$ fixed, we solve the following slave problems: $\min_{\mathbf{z}_k \in \mathcal{F}} (g_k(\mathbf{z}_k) + \boldsymbol{\lambda}_k \mathbf{z}_k)$. Upon obtaining the optimal solutions \mathbf{z}_k^* of the slave problems,

we update the values of $\boldsymbol{\lambda}_k$ by projected subgradient descent where the subgradient with respect to $\boldsymbol{\lambda}_k$ is \mathbf{z}_k^* . In other words, we update $\boldsymbol{\lambda}_k \leftarrow \boldsymbol{\lambda}_k + \eta_t \mathbf{z}_k^*$ where η_t is the learning rate at iteration t . In order to satisfy the constraint $\sum_k \boldsymbol{\lambda}_k = \mathbf{0}$ we project the value of $\boldsymbol{\lambda}_k$ to $\boldsymbol{\lambda}_k \leftarrow \boldsymbol{\lambda}_k - (\sum_k \boldsymbol{\lambda}_k) / M$. Under fairly general conditions, this iterative strategy known as dual decomposition converges to the globally optimal solution of the original problem. We refer the reader to [18] for details.

3.2.2 Dual Decomposition for Selecting Regions

When using dual decomposition, it is crucial to select slave problems whose optimal solutions can be computed quickly. With this in mind, we choose three types of slave problems. Below we describe each of these slave problems and justify their use by providing an efficient method to optimize them.

The first type of slave problems is similar to the one used for energy minimization in pairwise random fields. Specifically, each slave problem is defined using a subset of regions $\mathcal{D}_T \subseteq \mathcal{D}$ and edges $\mathcal{E}_T \subseteq \mathcal{E}$ that form a tree. For each such graph $(\mathcal{D}_T, \mathcal{E}_T)$ we define the following problem:

$$\begin{aligned} \min_{\mathbf{y} \geq 0} \quad & \sum_{r \in \mathcal{D}_T, i} \left(\frac{\bar{\theta}_r(i)}{n_r} + \lambda_r^1(i) \right) y_r(i) + \\ & \sum_{(r,r') \in \mathcal{E}_T, i, j} \left(\frac{\bar{\theta}_{rr'}(i, j)}{n_{rr'}} + \lambda_{rr'}^1(i, j) \right) y_{rr'}(i, j), \\ \text{s.t.} \quad & \mathbf{y} \geq 0, \sum_{i \in \mathcal{I}} y_r(i) = 1, \\ & \sum_{j \in \mathcal{I}} y_{rr'}(i, j) = y_r(i), \sum_{i \in \mathcal{I}} y_{rr'}(i, j) = y_{r'}(j), \end{aligned} \quad (12)$$

where n_r and $n_{rr'}$ are the total number of slave problems involving $r \in \mathcal{D}$ and $(r, r') \in \mathcal{E}$ respectively. As the LP relaxation is tight for trees (that is, it has integral solutions) [25], [26], the above problem can be solved efficiently using belief propagation [35].

The second type of slave problems correspond to the covering constraints. Specifically, for each $s \in \mathcal{S}$ we define:

$$\min_{\mathbf{y} \geq 0} \sum_{r \in \mathcal{C}(s), i} \left(\frac{\bar{\theta}_r(i)}{n_r} + \lambda_r^2(i) \right) y_r(i) \quad (13)$$

$$\text{s.t.} \quad \sum_{r \in \mathcal{C}(s)} \sum_{i \in \mathcal{L}^R} y_r(i) = 1, \sum_{i \in \mathcal{L}^I} y_r(i) = 1. \quad (14)$$

It can be verified that the constraint matrix for the above problem is totally unimodular. In other words, the optimal solution is integral and in fact, can be found efficiently in $O(L|\mathcal{C}(s)|)$ time (where L is the number of region labels and $|\mathcal{C}(s)|$ is the number of regions in \mathcal{D} that cover s).

The third type of slave problems correspond to the clique constraints defined in §3.1. Specifically, for every clique defined by $\mathcal{D}_Q \subseteq \mathcal{D}$ and $\mathcal{E}_Q \subseteq \mathcal{E}$, we specify:

$$\min_{\mathbf{y} \in \mathcal{M}_Q} \sum_{r \in \mathcal{D}_Q, i} \left(\frac{\bar{\theta}_r(i)}{n_r} + \lambda_r^3(i) \right) y_r(i) + \sum_{(r, r') \in \mathcal{E}_Q, i, j} \left(\frac{\bar{\theta}_{rr'}(i, j)}{n_{rr'}} + \lambda_{rr'}^3(i, j) \right) y_{rr'}(i, j), \quad (15)$$

where \mathcal{M}_Q is the marginal polytope of the clique. Note that since \mathcal{M}_Q is the convex hull of all valid integral labelings and the objective function is linear, it follows that the above problem has an integral optimal solution. Furthermore, it corresponds to the minimization of a sparse higher order function [36], [37]. In other words, although there are $O((L+1)^{|\mathcal{D}_Q|})$ possible labelings for the clique, only a small fraction of them are valid. For example, consider the clique shown in Fig. 1. If region r_4 is selected, then regions r_1 and r_2 cannot be selected since they overlap with r_4 . Furthermore, we also know that exactly one of r_1 and r_4 must be selected. Using similar observations, the number of valid labelings of a general clique used in our relaxation (that is, a clique formed by three neighboring regions along with all the regions that overlap with at least one of them) can be shown to be $O((L|\mathcal{C}(s_Q)|)^3)$. Here, s_Q is the super-pixel s that is covered by at least one region in \mathcal{D}_Q and has the largest corresponding set $\mathcal{C}(s)$. Note that we can also use cliques defined by m neighboring regions that can be optimized in $O((L|\mathcal{C}(s_Q)|)^m)$ time (we omit details on the derivation of the time complexity due to lack of space). However, we found $m = 3$ to be sufficient to obtain a tight relaxation.

We iteratively solve the above slave problems and update λ . Upon convergence, we obtain the primal solution (that is, the subset of regions and their labels) in a similar manner to [30], [34], [36]. Briefly, this involves sequentially considering each super-pixel (in some arbitrary order) and picking the best region for it (according to the values of the dual variable λ). The primal-dual gap provides us with an estimate of how

good our solution is. Typically, we obtain a very small gap by using our approach.

3.3 Generating the Dictionaries

Ideally, we would like to form a dictionary that consists of all the regions obtained by merging and intersecting the segments provided by a bottom-up approach. However, this will clearly result in a very large dictionary that cannot be handled even using our efficient algorithm. Instead, we iteratively search over the regions using the following strategy. We initialize our dictionary with the regions obtained from one (very coarse) over-segmentation. In the subsequent iterations, we consider two different types of dictionaries (similar to [15]). The first dictionary consists of the current set of regions \mathcal{D}_{cur} (those that have provided the best explanation of the image until now, in terms of the energy) as well as all the regions obtained by merging two neighboring regions in \mathcal{D}_{cur} . The second type of dictionary uses multiple over-segmentations to define the regions. Specifically, in addition to \mathcal{D}_{cur} , it also consists of all the regions obtained by merging every segment from the over-segmentations with all its overlapping and neighboring regions in \mathcal{D}_{cur} . Similarly, it also contains all the regions obtained by intersecting every segment with all its overlapping regions in \mathcal{D}_{cur} . While using the first dictionary results in larger regions (by merging neighboring regions together), the second dictionary allows us to correct any mistakes (merging two incompatible regions) by considering intersections of segments and regions.

It is worth noting that, unlike most of the previous works [7], [12], [38], [39], [40], our dictionaries define regions that are not just segments obtained by a bottom-up approach. This additional degree of freedom is important for obtaining regions that provide discriminative features while respecting scene boundaries.

Although dual decomposition is in general an efficient strategy for solving convex problems, when the size of the dictionary is too large, it may not be practically feasible. In order to improve the efficiency of our overall energy minimization algorithm, we use two heuristics based on a shrinking strategy [41] and a move-making strategy [42] respectively. Since these heuristics are not central to the understanding of the remainder of the paper, we provide their description in the appendices.

4 PARAMETER ESTIMATION

Given a training dataset that consists of images with different types of ground-truth annotations, our goal is to learn accurate parameters for the region-based model. To this end, we design a principled framework for learning with diverse data, with the aim of exploiting the varying degrees of information in the different datasets to the fullest: from the cleanliness

of pixelwise segmentations, to the vast availability of bounding boxes and image-level labels. Specifically, we formulate the parameter learning problem using a latent structural support vector machine (LSVM) [20], [21], where the latent variables model any missing information in the annotation. For this work, we focus on three types of missing information: (i) the specific class of a pixel labeled using a generic foreground or background class (for example, images from the VOC segmentation dataset [19] or the Stanford background dataset [15]; (ii) the segmentation of an image annotated with bounding boxes of objects (for example, images from the VOC detection dataset [19]); and (iii) the segmentation of an image labeled to indicate the presence of a class (for example, images from the ImageNet dataset [43]). We show how the energy minimization algorithm described in the previous section can be suitably modified to obtain an accurate local minimum to the optimization problem corresponding to LSVM using the self-paced learning algorithm [22].

4.1 Learning with Generic Classes

To simplify the discussion, we first focus on the case where the ground-truth annotation of an image specifies a pixelwise segmentation that includes generic foreground or background labels. As will be seen in subsequent subsections, the other cases of interest, where the ground-truth only specifies bounding boxes for objects or image-level labels, will be handled by solving a series of LSVM problems that deal with generic class annotations.

4.1.1 Notation

As mentioned earlier in section 2, we denote an image by \mathbf{X} and a labeling (that is, a segmentation) by \mathbf{Y} . The joint feature vector is denoted by $\Psi(\mathbf{X}, \mathbf{Y})$, and the energy of a segmentation is equal to $\mathbf{w}^\top \Psi(\mathbf{X}, \mathbf{Y})$, where \mathbf{w} are the parameters of the model. The best segmentation of an image is obtained using energy minimization, that is, $\mathbf{Y}^* = \arg\min_{\mathbf{Y}} \mathbf{w}^\top \Psi(\mathbf{X}, \mathbf{Y})$.

We denote the training dataset as $\mathcal{T} = \{(\mathbf{X}_k, \mathbf{A}_k), k = 1, \dots, n\}$, where \mathbf{X}_k is an image and \mathbf{A}_k is the corresponding annotation. For each image \mathbf{X} with annotation \mathbf{A} , we specify a set of latent, or hidden, variables \mathbf{H} such that $\{\mathbf{A}, \mathbf{H}\}$ defines a labeling \mathbf{Y} of the model. In other words, for each pixel p labeled using the generic foreground (background) class, the latent variable \mathbf{H}_p models its specific foreground (background) class.

4.1.2 Learning as Risk Minimization

Given the dataset \mathcal{T} , we learn the parameters \mathbf{w} by training an LSVM. Briefly, an LSVM minimizes an upper bound on a user-specified risk, or loss, $\Delta(\mathbf{A}, \{\hat{\mathbf{A}}, \hat{\mathbf{H}}\})$. Here, \mathbf{A} is the ground-truth and $\{\hat{\mathbf{A}}, \hat{\mathbf{H}}\}$ is the predicted segmentation for a given set of parameters. In this work, we specify the loss using the overlap score, which is the measure of accuracy

for the VOC challenge [19]. For an image labeled using specific foreground classes and a generic background (the label set denoted by \mathcal{F}), we define the loss function as

$$\Delta(\mathbf{A}, \{\hat{\mathbf{A}}, \hat{\mathbf{H}}\}) = 1 - \frac{1}{|\mathcal{F}|} \sum_{i \in \mathcal{F}} \frac{|\mathcal{P}_i(\mathbf{A}) \cap \mathcal{P}_i(\{\hat{\mathbf{A}}, \hat{\mathbf{H}}\})|}{|\mathcal{P}_i(\mathbf{A}) \cup \mathcal{P}_i(\{\hat{\mathbf{A}}, \hat{\mathbf{H}}\})|}, \quad (16)$$

where the function $\mathcal{P}_i(\cdot)$ returns the set of all the pixels labeled using class i . Note that when i is the generic background, then $\mathcal{P}_i(\{\hat{\mathbf{A}}, \hat{\mathbf{H}}\})$ is the set of all pixels labeled using any specific background class. A similar loss function can be defined for images labeled using specific background classes and a generic foreground (the label set \mathcal{B}). Formally, the parameters are learned by solving the following non-convex optimization problem:

$$\begin{aligned} \min_{\mathbf{w}, \xi_k \geq 0} \quad & \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{n} \sum_{k=1}^n \xi_k, \\ \text{s.t.} \quad & \mathbf{w}^\top \Psi(\mathbf{X}_k, \{\bar{\mathbf{A}}_k, \bar{\mathbf{H}}_k\}) - \min_{\mathbf{H}} \Psi(\mathbf{X}_k, \{\mathbf{A}_k, \mathbf{H}\}) \\ & \geq \Delta(\mathbf{A}_k, \{\bar{\mathbf{A}}_k, \bar{\mathbf{H}}_k\}) - \xi_k, \forall \bar{\mathbf{A}}_k, \bar{\mathbf{H}}_k. \end{aligned} \quad (17)$$

Intuitively, the problem requires that for every image the energy of the ground-truth annotation, augmented with the best value of its latent variables, should be less than the energy of all other labelings. The desired margin between the two energy values is proportional to the loss.

4.1.3 Learning an LSVM

Algorithm 1 describes the main steps of learning an LSVM using our recently proposed self-paced learning (SPL) algorithm [22]. Unlike the concave convex procedure (CCCP) [20], [21] used in previous works, which treats all images equally, SPL automatically chooses a set of easy images at each iteration (in step 4), and uses only those images to update the parameters.

Following our earlier work [22], we choose the initial threshold σ_0 such that half the images are considered easy in the first iteration, and set the annealing factor $\mu = 1.3$. These settings have been shown to work well in practice for a large number of applications [22]. In order to avoid learning from images whose latent variables are never imputed correctly, we measure the accuracy of the model after each iteration using a validation set (different from the test set). We report test results using the parameters that are the most accurate on the validation set.

Let us take a closer look at what is needed to learn the parameters of our model. The first step of SPL requires us to impute the latent variables of each training image \mathbf{X} , given the ground-truth annotation \mathbf{A} , by solving problem (18). In other words, this step completes the segmentation to provide us with a positive example. We call this annotation-consistent inference. The second step requires us to find a segmentation with low energy but high loss (a negative

Algorithm 1 *The self-paced learning algorithm for LSVM.*

input $\mathcal{T}, \mathbf{w}_0, \sigma_0, \mu, \epsilon$.

 1: $\mathbf{w} \leftarrow \mathbf{w}_0, \sigma \leftarrow \sigma_0$.

 2: **repeat**

3: Impute the latent variables as

$$\mathbf{H}_k = \underset{\mathbf{H}}{\operatorname{argmin}} \mathbf{w}^\top \Psi(\mathbf{X}_k, \{\mathbf{A}_k, \mathbf{H}\}). \quad (18)$$

 4: Compute the slack variables ξ_k as

$$\begin{aligned} \xi_k = \max_{\bar{\mathbf{A}}_k, \bar{\mathbf{H}}_k} & \Delta(\mathbf{A}_k, \{\bar{\mathbf{A}}_k, \bar{\mathbf{H}}_k\}) - \\ & \mathbf{w}^\top \Psi(\mathbf{X}_k, \{\bar{\mathbf{A}}_k, \bar{\mathbf{H}}_k\}) \\ & + \mathbf{w}^\top \Psi(\mathbf{X}_k, \{\mathbf{A}_k, \mathbf{H}_k\}). \end{aligned} \quad (19)$$

 Using ξ_k , define variables $v_k = \delta(\xi_k \leq \sigma)$, where $\delta(\cdot) = 1$ if its argument is true and 0 otherwise.

5: Update the parameters by solving the following problem that only considers easy images:

$$\begin{aligned} \min_{\mathbf{w}, \xi_k \geq 0} & \quad \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{n} \sum_{k=1}^n v_k \xi_k, \\ \text{s.t.} & \quad \mathbf{w}^\top (\Psi(\mathbf{X}_k, \{\bar{\mathbf{A}}_k, \bar{\mathbf{H}}_k\}) - \\ & \quad \Psi(\mathbf{X}_k, \{\mathbf{A}_k, \mathbf{H}_k\})) \\ & \quad \geq \Delta(\mathbf{A}_k, \{\bar{\mathbf{A}}_k, \bar{\mathbf{H}}_k\}) - \xi_k, \forall \bar{\mathbf{A}}_k, \bar{\mathbf{H}}_k. \end{aligned} \quad (20)$$

 6: Change the threshold $\sigma \leftarrow \sigma \mu$.

 7: **until** Decrease in objective (17) is below tolerance ϵ and all images have been labeled as easy.

example) by solving problem (19). We call this loss-augmented inference. The third step requires solving the convex optimization problem (20). In this work, we solve it using stochastic subgradient descent [44], where the subgradient for a given image \mathbf{X} is the joint feature vector $\Psi(\mathbf{X}, \{\mathbf{A}^*, \mathbf{H}^*\})$, where $\{\mathbf{A}^*, \mathbf{H}^*\}$ are obtained by solving problem (19). We refer the interested reader to [44] for details.

To summarize, in order to learn from generic class segmentations, we require two types of inference algorithms—annotation-consistent inference and loss-augmented inference. In the remainder of this subsection, we describe these inference algorithms for the region-based model. However, we note that both these types of inference algorithms can be designed for several existing models by suitably modifying their energy minimization algorithm, which would allow enable us to learn their parameters using datasets labeled with generic classes.

4.1.4 Annotation-Consistent Inference

The goal of annotation-consistent inference is to impute the latent variables that minimize the energy under the constraint that they do not contradict the ground-truth annotation (which specifies a pixelwise segmentation using generic classes). In other words, a pixel marked as a specific class must belong to a region labeled as that class. Furthermore, a pixel marked

as generic foreground (background) must be labeled using a specific foreground (background) class.

For a given dictionary of regions \mathcal{D} , annotation-consistent inference is equivalent to the following integer program (IP):

$$\min_{\bar{\mathbf{y}} \in \operatorname{SELECT}(\mathcal{D})} \boldsymbol{\theta}^\top \bar{\mathbf{y}} \quad \text{s.t.} \quad \Delta(\mathbf{A}, \bar{\mathbf{y}}) = 0. \quad (21)$$

The set $\operatorname{SELECT}(\mathcal{D})$ refers to the set of all valid assignments to $\bar{\mathbf{y}}$, as specified in equation (10). The constraint $\Delta(\mathbf{A}, \bar{\mathbf{y}}) = 0$ (where we have overloaded Δ for simplicity to consider $\bar{\mathbf{y}}$ as its argument) ensures that the imputed latent variables are consistent with the ground-truth. Similar to the energy minimization problem discussed in the previous section, the above IP is solved approximately by relaxing the elements of $\bar{\mathbf{y}}$ to take values between 0 and 1, resulting in a linear program (LP).

Fig. 2(a) shows examples of the segmentation obtained using the above annotation-consistent inference over different iterations of SPL. Note that the segmentations obtained are able to correctly identify the specific classes of pixels labeled using generic classes. The quality of the segmentation, together with the ability of SPL to select correct images to learn from, results in an accurate set of parameters.

4.1.5 Loss-Augmented Inference

The goal of loss-augmented inference is to find a labeling that minimizes the energy while maximizing the loss (as shown in problem (19)), which can be formulated as the following IP:

$$\min_{\bar{\mathbf{y}} \in \operatorname{SELECT}(\mathcal{D})} \boldsymbol{\theta}^\top \bar{\mathbf{y}} - \Delta(\mathbf{A}, \bar{\mathbf{y}}). \quad (22)$$

Unfortunately, relaxing $\bar{\mathbf{y}}$ to take fractional values in the interval $[0, 1]$ for the above problem does not result in an LP. This is due to the dependence of Δ on the labeling $\bar{\mathbf{y}}$ in both its numerator and denominator (see equation (16)). We address this issue by adopting a two stage strategy: (i) replace Δ by another loss function that results in an LP relaxation; and (ii) using the solution of the first stage as an accurate initialization, solve problem (22) via iterated conditional modes (ICM). For the first stage, we use the macro-average error over all classes as the loss function, that is

$$\Delta'(\mathbf{A}, \{\hat{\mathbf{A}}, \hat{\mathbf{H}}\}) = 1 - \frac{1}{|\mathcal{L}|} \sum_{i \in \mathcal{L}} \frac{|\mathcal{P}_i(\mathbf{A}) \cap \mathcal{P}_i(\{\hat{\mathbf{A}}, \hat{\mathbf{H}}\})|}{|\mathcal{P}_i(\mathbf{A})|}, \quad (23)$$

where \mathcal{L} is the appropriate label set (\mathcal{F} for images labeled using specific foreground and generic background, \mathcal{B} for images labeled using specific background and generic foreground). Note that the denominator of Δ' does not depend on the predicted labeling. Hence, it can be absorbed into the unary potentials, leading to a pairwise energy minimization problem, which can be solved using our LP relaxation.

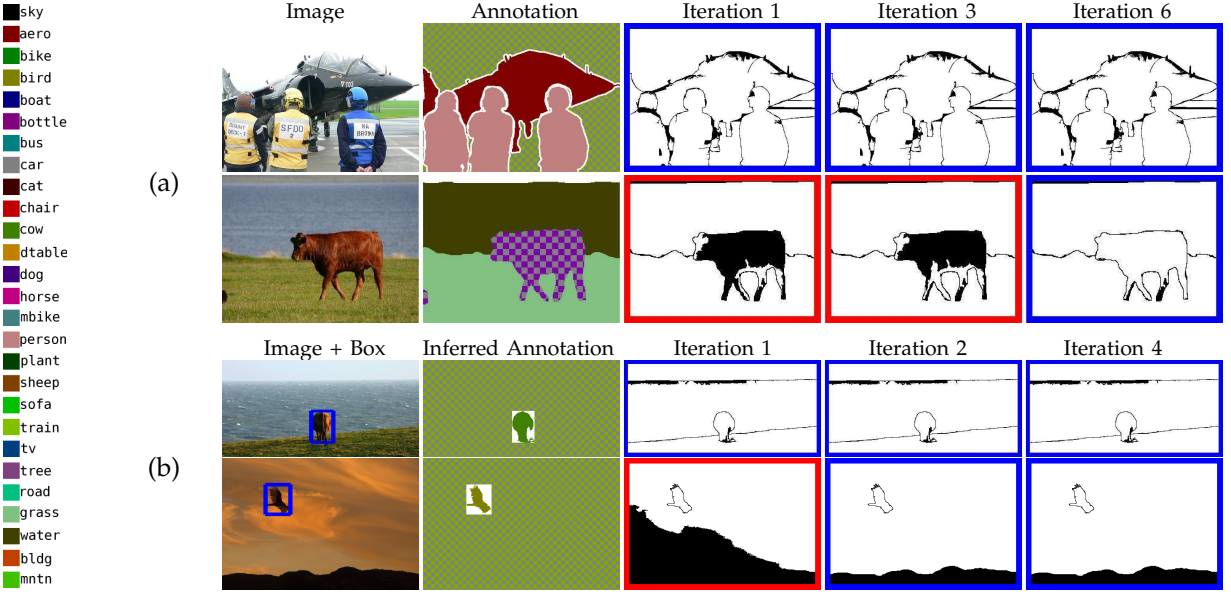


Fig. 2. Labelings obtained using annotation-consistent inference during different iterations of SPL. (a) Images annotated with generic classes. Column 2 shows the annotation (where the checkered patterns indicate generic classes). In columns 3 – 5, pixels labeled using the correct specific-class by annotation-consistent inference are shown in white, while pixel labeled using the wrong specific-class are shown in black (we labeled these images with specific-class annotations only for the purpose of illustration; these annotations were not used during training). A blue surrounding box on the labeling implies that the example was selected as easy by SPL, while a red surrounding box indicates that it wasn't selected during the specified iteration. Note that SPL discards the image where the cow (row 2) is incorrectly labeled. (b) Images annotated using bounding boxes. Column 2 shows the annotation obtained using bounding box inference. Note that the objects have been accurately segmented. Furthermore, SPL discards the image where the sky (row 2) is incorrectly labeled.

In our experiments, ICM converged within very few iterations (typically less than 5) when initialized in this manner. As will be seen in section 5, the approximate subgradients provided by the above loss-augmented inference were sufficient to obtain an accurate set of parameters.

4.2 Learning with Bounding Boxes

We now focus on learning specific-class segmentation from training images with user-specified bounding boxes for instances of some classes. To simplify our description, we make the following assumptions: (i) the image contains only one bounding box, which provides us with the location information for an instance of a specific foreground class; and (ii) all the pixels that lie outside the bounding box belong to the generic background class. We note that our approach can be trivially extended to handle cases where the above assumptions do not hold true (for example, bounding boxes for background or multiple boxes per image).

The major obstacle in using bounding box annotations is the lack of a readily available loss function that compares bounding boxes \mathbf{B} to pixelwise segmentations $(\mathbf{A}, \hat{\mathbf{H}})$. Note that it would be unwise to use a loss function that compares two bounding boxes (the ground-truth and the predicted one that can be derived from the segmentation), as this function would not be compatible with the overlap score loss used in the previous section. In other words, minimizing such a loss function would not necessarily improve the segmentation accuracy. We address this issue by adopting a simple, yet effective, strategy that solves

a series of LSVM problems for generic classes. Our approach consists of three steps:

- Given an image \mathbf{X} and its bounding box annotation \mathbf{B} , we infer its segmentation \mathbf{Y}^B using the current set of parameters (say, the parameters learned using generic class segmentation data). The segmentation is obtained by minimizing an objective function that augments the energy of the model with terms that encourage the segmentation to agree with the bounding box (see details below).
- Using the above segmentation, we define a generic class annotation \mathbf{A} of the image (see details below).
- Annotations \mathbf{A} are used to learn the parameters.

The new parameters then update the segmentation, and the entire process is repeated until convergence (that is, until the segmentations stop changing). Note that the third step simply involves learning an LSVM as described in the previous section. We describe the first two steps in more detail.

4.2.1 Using the Bounding Box for Segmentation.

We assume that the specified bounding box is tight (a safe assumption for most datasets) and penalize any row and column of the bounding box that is not covered by the corresponding class. Here, a row or column is said to be covered if it contains a sufficient number of pixels s that have been assigned to the corresponding class of the bounding box. Formally, given a bounding box \mathbf{B} of class c , we define an annotation \mathbf{A}' such that all the pixels p inside the bounding box

have no label specified in \mathbf{A}' (denoted by $\mathbf{A}'_p = 0$) and all the pixels outside the bounding box have a generic background label (consistent with our assumption). Furthermore, we define latent variables \mathbf{H} that model the specific semantic classes for each pixel. Using \mathbf{A}' and \mathbf{B} , we estimate \mathbf{H} as follows:

$$\mathbf{H} = \underset{\bar{\mathbf{H}}}{\operatorname{argmin}} \mathbf{w}^\top \Psi(\mathbf{X}, \{\mathbf{A}', \bar{\mathbf{H}}\}) + \sum_q \kappa_q I_q(\bar{\mathbf{H}}, c). \quad (24)$$

Here q indexes the rows and columns of \mathbf{B} , and I_q is an indicator function for whether q is covered by the latent variables $\bar{\mathbf{H}}$. The penalty κ_q has a high value κ_{max} for the center row and center column, and decreases at a linear rate to κ_{min} at the boundary. In our experiments, we set $\kappa_{max} = 10\kappa_{min}$ and cross-validated the value of κ_{min} using a small set of images. We found that our method produced very similar segmentations for a large range of κ_{min} .

We refer to problem (24) as bounding-box inference. Note that I_q adds a higher-order potential to the energy of the model since its value depends on the labels of all the pixels in a particular row or column. However, the potential is sparse (that is, it only takes a non-zero value for a small number of labelings). Hence, the above problem can be optimized efficiently for any semantic segmentation model [36]. In the following, we describe bounding-box inference for our region-based model.

4.2.2 Bounding-Box Inference

Given a dictionary of regions \mathcal{D} and a bounding box \mathbf{B} of class c , we obtain the segmentation by solving the LP relaxation of the following IP:

$$\begin{aligned} \min_{\bar{\mathbf{y}} \in \operatorname{SELECT}(\mathcal{D}), \bar{y}_q \in \{0,1\}} \quad & \theta^\top \bar{\mathbf{y}} + \sum_q \kappa_q (1 - \bar{y}_q) \\ \text{s.t.} \quad & \Delta(\mathbf{A}', \bar{\mathbf{y}}) = 0, \bar{y}_q \leq \sum_{r \in \mathcal{C}(q)} \bar{y}_c^r. \end{aligned} \quad (25)$$

Here \bar{y}_q is a boolean variable whose value is the complement of the indicator function I_q in problem (24). Note that, for our model, a row or a column is considered covered if at least one region overlapping with it is assigned the class c . The loss function Δ is measured over all pixels that lie outside the bounding box, which are assumed to belong to the generic background class. Fig. 2(b) shows some example annotations obtained from bounding-box inference, together with the results of annotation-consistent inference during different iterations of SPL. The quality of the annotations and the ability of SPL to select good images ensures that our model is trained without noise.

Lempitsky *et al.* [45] have suggested a method to obtain a binary segmentation of an image with a user-specified bounding box. However, our setting differs from theirs in that, unlike the low-level vision model used in [45] (likelihoods from RGB values, contrast dependent penalties), we use a more sophisticated

high-level model which encodes information about specific classes and their pairwise relationship using a region-based representation. Hence, we can resort to a much simpler optimization strategy and still obtain accurate segmentations.

4.2.3 From Segmentation to Annotation.

Let $\mathbf{Y}^B = \{\mathbf{A}', \mathbf{H}\}$ denote the labeling obtained from bounding-box inference. Using \mathbf{Y}^B we define an annotation \mathbf{A} as follows. For each pixel p inside the bounding box that was labeled as class c , that is, $\mathbf{Y}_p^B = c$, we define $\mathbf{A}_p = c$. For pixels p inside the bounding box such that $\mathbf{Y}_p^B \neq c$, we define $\mathbf{A}_p = 0$. In other words, during annotation-consistent inference these pixels can belong to any class, foreground or background. The reason for specifying \mathbf{A}_p in this manner is that, while we are fairly certain that the pixels labeled $\mathbf{Y}_p^B = c$ do belong to the class c , due to the lack of information in the annotation we are not sure of which class the other pixels belong to. Not labeling such pixels prevents using the mistakes made in bounding-box inference to learn the parameters. Finally, for all the pixels outside the bounding box we set $\mathbf{A}_p = \mathbf{A}'_p$, that is, they are labeled as generic background.

4.3 Learning with Image-Level Labels

We use a similar three step process to the one described above for bounding boxes in order to take advantage of the numerous images with image-level labels, which indicate the presence of a class. In more detail, given an image containing class c , we define an annotation \mathbf{A}' that does not specify a label for any pixel of the image (that is, $\mathbf{A}'_p = 0$ for all p). We estimate the value of the latent variables \mathbf{H} that model the specific semantic class of each pixel by solving the following image-label inference problem:

$$\mathbf{H} = \underset{\bar{\mathbf{H}}}{\operatorname{argmin}} \mathbf{w}^\top \Psi(\mathbf{X}, \{\mathbf{A}', \bar{\mathbf{H}}\}) + \kappa_{max} I(\bar{\mathbf{H}}, c). \quad (26)$$

Here I is an indicator function for whether the image is covered by the latent variables $\bar{\mathbf{H}}$. Similar to a row or a column of a bounding box, an image is considered covered if a sufficient number of pixels s are assigned the label c . Once again, the function I introduces a sparse higher-order potential that can be handled efficiently [36]. This would allow us to learn a general semantic segmentation model using the method described above. For the region-based model, we perform image-label inference as follows.

Given a dictionary \mathcal{D} and image containing class c , we obtain a segmentation by solving the LP relaxation of the following IP:

$$\min_{\bar{\mathbf{y}} \in \operatorname{SELECT}(\mathcal{D}), \bar{y} \in \{0,1\}} \quad \theta^\top \bar{\mathbf{y}} - \kappa_{max} \bar{y}, \text{ s.t. } \bar{y} \leq \sum_{r \in \mathcal{D}} \bar{y}_c^r. \quad (27)$$

The value of y is the complement of the indicator function I in problem (26). Once again, SPL reduces

the noise during training by selecting images with correct annotations and latent variables.

To obtain an annotation \mathbf{A} from the above segmentation, we label all pixels p belonging to class c as $\mathbf{A}_p = c$. For the rest of the pixels, we define $\mathbf{A}_p = 0$. The annotations obtained in this manner are used to refine the parameters and the entire process is repeated until convergence.

5 EXPERIMENTS

We now demonstrate the efficacy of our energy minimization and parameter estimation approaches using large, publicly available datasets.

5.1 Energy Minimization

We compare our LP relaxation based approach for energy minimization with three baselines: (i) regions obtained by the intersection of multiple over-segmentations of an image; (ii) regions defined by the segments of the best single over-segmentation (in terms of the energy of the model); and (iii) regions selected from a dictionary similar to ours by the method of [15] (using the code provided by the authors). For each method, we use the same set of parameters, which are learned by the closed loop learning (CLL) technique of [15].

5.1.1 Dataset

We use the publicly available Stanford background dataset [15]. It consists of 715 images (collected from standard datasets such as PASCAL, MSRC and geometric context) whose pixels have been labeled as belonging to one of seven background classes or a generic foreground class. For each image we use three over-segmentations obtained by employing different kernels for the standard mean-shift algorithm [9]. Similar to [15], we split the dataset into 572 images for training and 143 images for testing. We report results on four different splits.

5.1.2 Results

We evaluate the accuracy of the labeling obtained by measuring the percentage of pixels whose labels matched the ground-truth. Here, the label of a pixel is the label assigned to the region to which it belongs. Table 1 shows the average energy and accuracy for the different approaches. Note that all region-based methods outperform the pixel-based approach described in [15] (that provides comparable results to [4]). However, the choice of the regions greatly affects the value of the energy and hence the accuracy of the segmentation. By using large dictionaries and an accurate LP relaxation, our approach provides a statistically significant improvement (using paired t-test with $p = 0.05$) over other methods, both in terms of energy and accuracy. Our algorithm takes less than 10 minutes per image on average on a 2.4 GHz processor.

	Energy	Accuracy
Pixel	-	76.65 ± 1.20
Intersection	16150 ± 2005	76.84 ± 1.34
Segmentation	6796 ± 833	77.85 ± 1.50
[15]	4815 ± 592	78.52 ± 1.40
Our Method	1630 ± 306	79.42 ± 1.41

TABLE 1

The mean and standard deviation of the energy and the pixel-wise accuracy over four folds is shown. The first row corresponds to a pixel-based model [15]. The second row uses regions obtained by intersecting the three over-segmentations. The third row shows the results obtained by using the best single over-segmentation (in terms of the energy function). The fourth row corresponds to the method described in [15]. The fifth row shows our method's results. Using multiple over-segmentations to define an accurate dictionary and selecting the regions by our LP relaxation based method results in a lower energy labeling that provides better accuracy.

Fig. 3 shows some example segmentations obtained by the various approaches. Note that the regions corresponding to the intersection of over-segmentations respect the boundaries of the scene entities. However, they are too small to provide reliable features. Even using the best single over-segmentation results in regions that are not large enough. The method of [15] overcomes this problem to some extent by using an accurate dictionary of regions. However, as it is prone to getting stuck in a bad local minima, it sometimes selects regions that result in a high energy labeling. Our approach addresses this deficiency by allowing us to make large moves at each iteration.

5.2 Parameter Estimation

Despite obtaining a substantial improvement in terms of the energy, our LP relaxation approach only results in a small improvement in terms of the segmentation accuracy. This is due to the fact that the parameters learned using CLL are not suited to our inference approach. To alleviate this, we use our LSVM formulation to estimate the parameters using large, publicly available datasets that specify varying levels of annotation for the training images.

5.2.1 Generic Class Annotations

Comparison. We show the advantage of the LSVM formulation over CLL, which was specially designed for the region-based model, for the problem of learning a specific-class segmentation model using generic class annotations.

Datasets. We use two datasets: (i) the VOC2009 segmentation dataset, which provides us with annotations consisting of 20 specific foreground classes and a generic background; and (ii) SBD, which provides us with annotations consisting of 7 specific background classes and a generic foreground. Thus, we consider 27 specific classes, which results in a harder learning problem compared to methods that use only the

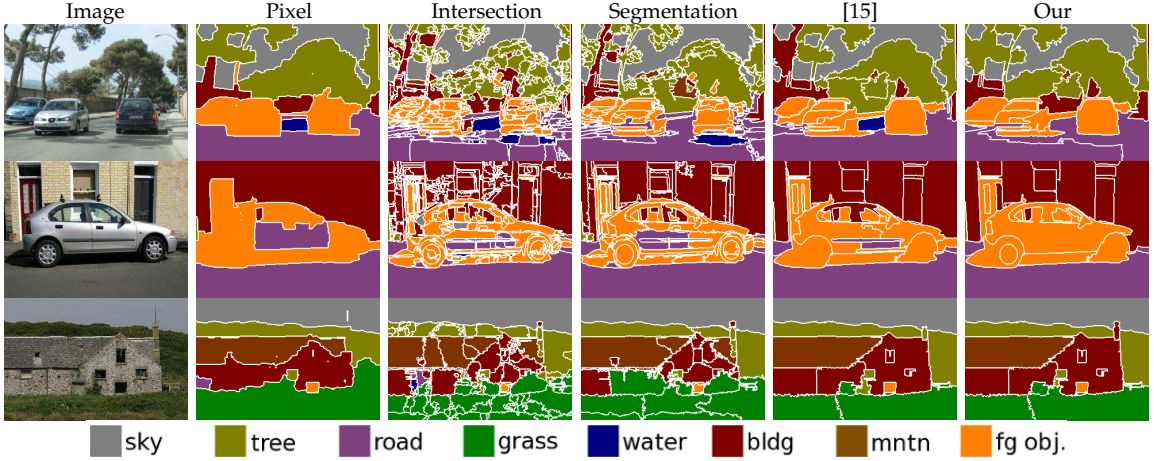


Fig. 3. Examples of semantic segmentation obtained using different types of regions. Our method provides large regions that align with the scene boundaries.

VOC2009 segmentation dataset or SBD. The total size of the dataset is 1846 training (1274 from VOC2009 and 572 from SBD), 278 validation (225 from VOC2009 and 53 from SBD) and 840 test (750 from VOC2009 and 90 from SBD) images. For CLL, the validation set is used to learn the pairwise potentials and several hyper-parameters while for LSVM it is used for early stopping (see section 4.1).

Results. Tables 2 and 3 (rows 1 and 2) show the accuracies obtained for SBD and VOC2009 test images respectively. The accuracies are measured using the overlap score, that is, $1 - \Delta(\mathbf{A}, \hat{\mathbf{A}}, \hat{\mathbf{H}})$, where \mathbf{A} is the ground-truth and $(\hat{\mathbf{A}}, \hat{\mathbf{H}})$ is the predicted segmentation. While both CLL and LSVM produce specific-class segmentations of all the test images, we use generic classes while measuring the performance due to the lack of specific-class ground-truth annotations. Note that LSVM provides better accuracies for nearly all the object classes in VOC2009 (17 of 21 classes). For SBD, LSVM provides a significant boost in performance for ‘sky’, ‘road’, ‘grass’ and ‘foreground’. With the exception of ‘building’, the accuracies for other classes is comparable. The reason for poor performance in the ‘mountain’ class is that several ‘mountain’ pixels are labeled as ‘tree’ in SBD (which confuses both the learning algorithms). Our results convincingly demonstrate the advantage of using LSVM.

5.2.2 Bounding Box Annotations

Comparison. We now compare the model learned using only generic class annotations with the model that is learned by also considering bounding box annotations. In keeping with the spirit of SPL, we use the previous LSVM model (learned using *easier* examples) as initialization for learning with additional bounding boxes.

Datasets. In addition to VOC2009 and SBD, we use some of the bounding box annotations that were introduced in the VOC2010 detection dataset. Our criteria for choosing the images is that (i) they were

not present in the VOC2009 detection dataset (which were used to obtain detection-based image features for the model); and (ii) none of their bounding boxes overlapped with each other. This provides us with an additional 1564 training images that have previously not been used to learn a segmentation model.

Results. Tables 2 and 3 (row 3) show the accuracies obtained by training on the above dataset for VOC2009 and SBD respectively. Once again, we observe an improvement in the accuracies for nearly all the VOC2009 classes (18 of 21 classes) compared to the LSVM trained using only generic class annotations. For SBD, we obtain a significant boost for ‘tree’, ‘water’ and ‘foreground’, while the accuracies of ‘road’, ‘grass’ and ‘mountain’ remain (almost) unchanged.

		a	s	t	r	g	w	b	m	
		v	k	e	a	r	a	t	i	n
		g	y	e	d	s	r	g	n	g
CLL		53.1	77.7	48.4	70.1	73.5	55.6	<u>62.5</u>	00.0	36.0
LSVM		54.3	<u>79.1</u>	48.2	<u>75.5</u>	76.0	55.1	61.4	00.0	39.1
BOX		54.8	78.3	48.6	75.4	76.0	59.9	60.8	00.0	39.6
LABELS		55.3	78.1	<u>49.5</u>	<u>75.5</u>	<u>76.1</u>	<u>60.1</u>	62.0	00.0	<u>41.3</u>
CCCP		53.8	75.4	48.7	70.0	74.0	59.9	62.5	00.0	39.9

Table 3. Accuracies for the SBD test set. See caption of Fig. 2 for an explanation of the various methods.

5.2.3 Image-Level Annotations

Comparison. We compare the model learned using generic class and bounding box annotations with the model that is learned by also considering image-level labels. Once again, following the idea of SPL closely, we use the model learned in the previous subsection as an initialization for learning with the additional image-level labels.

Datasets. In addition to SBD, VOC2009 segmentation dataset and VOC2010 detection dataset, we use a sub-

Our LP relaxation approach can also be used to perform energy minimization in other region-based models, such as those developed for geometric reconstruction [39], object discovery [46] and object detection and segmentation [38], [40]. Furthermore, the accuracy of the method may also help in the development of a successful unified scene understanding model.

While our parameter estimation approach focused on only three types of annotations, it can be extended to handle other types of data. For example, instead of just a two-level hierarchy (the generic classes and the specific classes), we can consider a general hierarchy of labels, such as the one defined in ImageNet [43] (‘Ferrari’ and ‘Honda’ sub-classes for ‘car’). Such a hierarchy can be viewed as a tree over labels. Given a pixel p annotated with a non-leaf label l , we can specify a latent variable \mathbf{H}_p that models its leaf-level label (where the leaves are restricted to lie in the subtree rooted at l). The loss function and the inference algorithms for generic classes can be trivially modified to deal with this more general case.

Our ongoing work is in two directions: (i) dealing with noisy labels, such as the labels obtained from Google Images or Flickr; and (ii) improving the efficiency of SPL by exploiting the fact that very easy images can be discarded during later iterations. Both these directions are aimed towards learning a segmentation model from the millions of freely available images on the Internet.

REFERENCES

- [1] 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,” *PAMI*, 2008.
- [2] X. He, R. Zemel, and M. Carriera-Perpinan, “Multiscale conditional random fields for image labeling,” in *CVPR*, 2004.
- [3] S. Konishi and A. Yuille, “Statistical cues for domain specific image segmentation with performance analysis,” in *CVPR*, 2000.
- [4] J. Shotton, J. Winn, C. Rother, and A. Criminisi, “TextonBoost: Joint appearance, shape and context modeling for multi-class object recognition and segmentation,” in *ECCV*, 2006.
- [5] D. Larlus and F. Jurie, “Combining appearance models and Markov random fields for category level object segmentations,” in *CVPR*, 2008.
- [6] A. Rabinovich, A. Vedaldi, C. Galleguillos, E. Wiewiora, and S. Belongie, “Objects in context,” in *ICCV*, 2007.
- [7] A. Saxena, M. Sun, and A. Ng, “Make3D: Learning 3D scene structure from a single still image,” *PAMI*, 2008.
- [8] P. Arbelaez, M. Maire, C. Fowlkes, and J. Malik, “From contours to regions: An empirical evaluation,” in *CVPR*, 2009.
- [9] D. Comaniciu and P. Meer, “Mean shift analysis and applications,” in *ICCV*, 1997.
- [10] C. Pontafaru, C. Schmid, and M. Hebert, “Object recognition by integrating multiple segmentations,” in *ECCV*, 2008.
- [11] J. Gonfaus, X. Boix, J. Van de Weijer, A. Bagdanov, J. Serrat, and J. Gonzalez, “Harmony potentials for joint classification and segmentation,” in *CVPR*, 2010.
- [12] P. Kohli, L. Ladicky, and P. Torr, “Robust higher order potentials for enforcing label consistency,” in *CVPR*, 2008.
- [13] L. Ladicky, C. Russell, P. Kohli, and P. Torr, “Associative hierarchical CRFs for object class image segmentation,” in *ICCV*, 2009.
- [14] F. Li, J. Carreira, and C. Sminchisescu, “Object recognition as ranking holistic figure-ground hypotheses,” in *CVPR*, 2010.
- [15] S. Gould, R. Fulton, and D. Koller, “Decomposing a scene into geometric and semantically consistent regions,” in *ICCV*, 2009.
- [16] S. Nowozin and C. Lampert, “Global connectivity potentials for random field models,” in *CVPR*, 2009.
- [17] S. Vicente, V. Kolmogorov, and C. Rother, “Graph cut based image segmentation with connectivity priors,” in *CVPR*, 2008.
- [18] D. Bertsekas, *Nonlinear Programming*. Athena Scientific, 1999.
- [19] M. Everingham, L. Van Gool, C. Williams, J. Winn, and A. Zisserman, “The PASCAL visual object classes (VOC) challenge,” *IJCV*, 2010.
- [20] P. Felzenszwalb, D. McAllester, and D. Ramanan, “A discriminatively trained, multiscale, deformable part model,” in *CVPR*, 2008.
- [21] C.-N. Yu and T. Joachims, “Learning structural SVMs with latent variables,” in *ICML*, 2009.
- [22] M. P. Kumar, B. Packer, and D. Koller, “Self-paced learning for latent variable models,” in *NIPS*, 2010.
- [23] M. P. Kumar and D. Koller, “Efficiently selecting regions for scene understanding,” in *CVPR*, 2010.
- [24] M. P. Kumar, H. Turki, D. Preston, and D. Koller, “Learning specific-class segmentation from diverse data,” in *ICCV*, 2011.
- [25] C. Chekuri, S. Khanna, J. Naor, and L. Zosin, “A linear programming formulation and approximation algorithms for the metric labelling problem,” *Disc. Math.*, 2005.
- [26] M. Wainwright, T. Jaakkola, and A. Willsky, “MAP estimation via agreement on trees: Message passing and linear programming,” *IEEE Info. Theory*, 2005.
- [27] P. Hammer, “Some network flow problems solved with pseudo-Boolean programming,” *Operations Research*, 1965.
- [28] C. Rother, V. Kolmogorov, V. Lempitsky, and M. Szummer, “Optimizing binary MRFs via extended roof duality,” in *CVPR*, 2007.
- [29] E. Boros and P. Hammer, “Pseudo-Boolean optimization,” *Discrete Applied Mathematics*, 2002.
- [30] N. Komodakis and N. Paragios, “Beyond loose LP-relaxations: Optimizing MRFs by repairing cycles,” in *ECCV*, 2008.
- [31] M. P. Kumar and P. Torr, “Efficiently solving convex relaxations for MAP estimation,” in *ICML*, 2008.
- [32] D. Sontag, T. Meltzer, A. Globerson, T. Jaakkola, and Y. Weiss, “Tightening LP relaxations for MAP using message passing,” in *UAI*, 2008.
- [33] F. Barahona and A. Mahjoub, “On the cut polytope,” *Mathematical Programming*, 1986.
- [34] N. Komodakis, N. Paragios, and G. Tziritas, “MRF optimization via dual decomposition: Message-passing revisited,” in *ICCV*, 2007.
- [35] J. Pearl, *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. Morgan Kaufman, 1998.
- [36] N. Komodakis and N. Paragios, “Beyond pairwise energies: Efficient optimization for higher-order MRFs,” in *CVPR*, 2009.
- [37] C. Rother, P. Kohli, W. Feng, and J. Jia, “Minimizing sparse higher order functions of discrete variables,” in *CVPR*, 2009.
- [38] C. Gu, J. Lim, P. Arbelaez, and J. Malik, “Recognition using regions,” in *CVPR*, 2009.
- [39] D. Hoiem, A. Efros, and M. Herbert, “Geometric context from a single image,” in *ICCV*, 2005.
- [40] L.-J. Li, R. Socher, and L. Fei-Fei, “Towards total scene understanding: Classification, annotation and segmentation in an automatic framework,” in *CVPR*, 2009.
- [41] T. Joachims, “Making large-scale SVM learning practical,” in *Advances in Kernel Methods*. MIT Press, 1999.
- [42] Y. Boykov, O. Veksler, and R. Zabih, “Fast approximate energy minimization via graph cuts,” *PAMI*, 2001.
- [43] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “ImageNet: A large-scale hierarchical image database,” in *CVPR*, 2009.
- [44] S. Shalev-Shwartz, Y. Singer, and N. Srebro, “Pegasos: Primal estimated sub-gradient solver for SVM,” in *ICML*, 2009.
- [45] V. Lempitsky, P. Kohli, C. Rother, and T. Sharp, “Image segmentation with a bounding box prior,” in *ICCV*, 2009.
- [46] B. Russell, A. Efros, J. Sivic, W. Freeman, and A. Zisserman, “Using multiple segmentations to discover objects and their extent in image collections,” in *CVPR*, 2006.