

Discriminative Models for Multi-Class Object Layout

Chaitanya Desai · Deva Ramanan · Charless C. Fowlkes

Received: 2 March 2010 / Accepted: 18 March 2011
© Springer Science+Business Media, LLC 2011

Abstract Many state-of-the-art approaches for object recognition reduce the problem to a 0-1 classification task. This allows one to leverage sophisticated machine learning techniques for training classifiers from labeled examples. However, these models are typically trained independently for each class using positive and negative examples cropped from images. At test-time, various post-processing heuristics such as non-maxima suppression (NMS) are required to reconcile multiple detections within and between different classes for each image. Though crucial to good performance on benchmarks, this post-processing is usually defined heuristically.

We introduce a unified model for multi-class object recognition that casts the problem as a structured prediction task. Rather than predicting a binary label for each image window independently, our model simultaneously predicts a structured labeling of the entire image (Fig. 1). Our model learns statistics that capture the spatial arrangements of various object classes in real images, both in terms of which arrangements to suppress through NMS and which arrangements to favor through spatial co-occurrence statistics.

We formulate parameter estimation in our model as a max-margin learning problem. Given training images with ground-truth object locations, we show how to formulate learning as a convex optimization problem. We employ the cutting plane algorithm of Joachims et al. (Mach. Learn. 2009) to efficiently learn a model from thousands

of training images. We show state-of-the-art results on the PASCAL VOC benchmark that indicate the benefits of learning a global model encapsulating the spatial layout of multiple object classes (a preliminary version of this work appeared in ICCV 2009, Desai et al., IEEE international conference on computer vision, 2009).

Keywords Object recognition · Context · Structured prediction · Cutting plane

1 Introduction

A contemporary and successful approach to object recognition is to formulate it as a classification task, e.g. “Does an image window at location i contain a given object o ?”. The classification formulation allows immediate application of a variety of sophisticated machine learning techniques in order to learn optimal detectors from training data. Such methods have the potential to encapsulate those subtle statistical regularities of the visual world which separate object from background. As a result, learning approaches have often yielded detectors that are more robust and accurate than their hand built counterparts for a range of applications, from edge and face detection to general purpose object recognition (see e.g., Rowley et al. 1996; Viola and Jones 2004).

In contrast to the well founded techniques used for classification of individual image patches, the problem of correctly detecting and localizing multiple objects from multiple classes within an image of a scene has generally been approached in a far more ad-hoc manner. For example, *non-max suppression* (NMS) is required to remove some detections returned by a classifier based on overlap criteria or more complicated heuristics (e.g. the mode finding approach

C. Desai (✉) · D. Ramanan · C.C. Fowlkes
Department of Computer Science, UC Irvine, Irvine, CA, USA
e-mail: desaic@ics.uci.edu

D. Ramanan
e-mail: dramanan@ics.uci.edu

C.C. Fowlkes
e-mail: fowlkes@ics.uci.edu

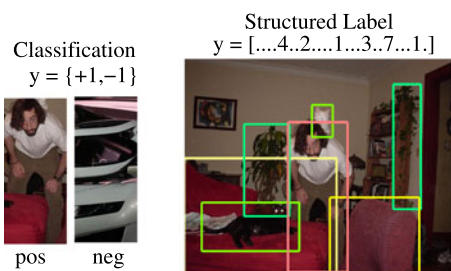


Fig. 1 Our framework. Classification-based approaches for recognition predict a binary label for a cropped window (*left*). We formulate the recognition problem as predicting a sparse, structured label vector specifying which windows, if any, contain particular objects in an entire input image. The latter allows our model to capture a wide range of contextual constraints among objects as described in Table 1 and Fig. 2

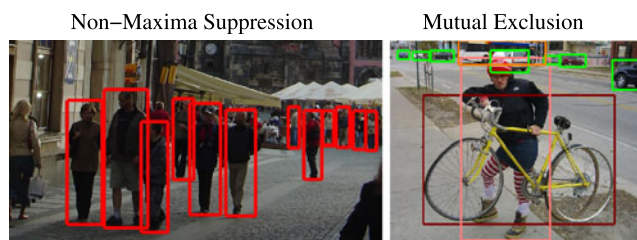


Fig. 2 Our novel contributions include the ability to learn inhibitory intra-class constraints (NMS) and inhibitory inter-class constraints (Mutual Exclusion) in a single unified model along with contextual cueing and spatial co-occurrence. Naïve methods for NMS or mutual exclusion may fail for objects that tend to overlap themselves (*left*) and other objects (*right*). In contrast, our framework *learns* how best to enforce such constraints from training data. We formulate the tasks of NMS and Mutual Exclusion using the language of structured prediction. This allows us to compute an optimal model by minimizing a convex objective function

of Dalal and Triggs 2005). Such tricks of the trade are essential to good performance on benchmarks designed to penalize multiple non-localized detections, however, they highlight a clear disconnect between training and testing phases. The objective optimized during learning only characterizes a sub-component of the final system used at runtime.

Furthermore, there is a wide range of possible interactions between object detections which is not fully captured by ad-hoc approaches. In street-level views, pedestrians are likely to occur standing next to each other, nearly overlapping, but unlikely to occur directly above or below each other (Fig. 2). In general, spatial object-object interactions may be arbitrarily complex and depend on latent information which is not readily available from single image. As an extreme example, studies of *proxemics* (Hall 1966), the body spacing and pose of people as they interact, shows that physical spacing between people depends in complicated ways on their “social distance”. While such complex interactions are difficult to encode, we argue there does exist useful information that is being ignored by current ad-hoc approaches to NMS.

Table 1 A taxonomy of interactions captured in our model. *Within* a single object class, our model can favor typical spatial layouts of objects (people often stand in crowds) while directly learning how to inhibit overlapping detections in such cases (NMS). Our model also captures long-range interactions between objects, such as the constraint that there exists at most one object instance (counting). Analogous interactions exist *between* object classes, including typical spatial relations between objects (bottles sit on tables), mutual exclusion (dog and cat detectors should not respond to the same image region), and co-occurrence (couches and cars do not commonly co-occur)

	Within-class	Between-class
Activation	Textures of objects	Spatial cueing
Inhibition	NMS	Mutual exclusion
Global	Expected counts	Co-occurrence

NMS is generally described in terms of intra-class inhibition, but can be generalized to suppression of overlapping detections between different classes. We refer to this more general constraint, that two objects cannot occupy the same 3D volume at the same time, as *mutual exclusion*. As seen in a 2D image projection, the exact nature of this constraint depends on the object classes. Figure 2 (right) shows an example of ground-truth labelings in the PASCAL VOC (Everingham et al. 2007) dataset in which strict mutual-exclusion would produce sub-optimal performance.

Object detections can also serve to enhance rather than inhibit other detections within a scene. This has been an area of active research in object recognition over the last few years (Torralba et al. 2004; Murphy et al. 2003; Galleguillos et al. 2008; He et al. 2004; Hoiem et al. 2008; Baur et al. 2008; Kumar and Hebert 2005). For example, different object classes may be likely to co-occur in a particular spatial layout. People ride on bikes, bottles rest on tables, and so on. In *contextual cueing*, a confident detection of one object (a bike) provides evidence that increases the likelihood of detecting another object (a person above the bike) (Baur et al. 2008; Galleguillos et al. 2008; Kumar and Hebert 2005). Contextual cueing can also occur within an object category, e.g., a crowd of pedestrians reinforcing each other’s detection responses. An extreme example of this phenomena is *near-regular texture* in which the spatial locations of nearly identical elements provides a strong prior on the expected locations of additional elements, lowering their detection threshold (Liu et al. 2004).

In Table 1 we outline a simplified taxonomy of different types of object-object interactions, both positive and negative, within and between classes. The contribution of this paper is a single model that incorporates all interactions from Table 1 through the framework of structured prediction. Rather than returning a binary label for a each image window, our model simultaneously predicts a set of detections for multiple objects from multiple classes over the entire image. Given training images with ground-truth object

locations, we show how to formulate parameter estimation as a convex max-margin learning problem. We employ the cutting plane algorithm of Joachims et al. (2009) to efficiently learn globally optimal parameters from thousands of training images.

In the sections that follow we formulate the structured output model in detail, describe how to perform inference and learning, and detail the optimization procedures used to efficiently learn parameters. We show state-of-the-art results on the PASCAL 2007 VOC benchmark (Everingham et al. 2007), indicating the benefits of learning a global model that encapsulates the layout statistics of multiple object classes in real images. We conclude with a discussion of related work and future directions.

2 Model

We describe a model for capturing interactions across a family of object detectors. To do so, we will explicitly represent an image as a collection of overlapping windows at various scales. The location of the i th window is given by its center and scale, written as $l_i = (x, y, s)$. The collection of N windows are precisely the regions scored by a scanning-window detector. Write x_i for the features extracted from window i . For example, in our experiments x_i is a normalized histogram of gradient features (Dalal and Triggs 2005). The entire image can then be represented as the collection of feature vectors $X = \{x_i : i = 1, \dots, N\}$.

Assume we have K object models. We write $y_i \in \{0, \dots, K\}$ for the label of the i th window, where the 0 label designates the background. Let $Y = \{y_i : i = 1, \dots, N\}$ be the entire label vector for the set of all sub-windows in an image. We define the score of labeling image X with vector Y as:

$$S(X, Y) = \sum_{i,j} w_{y_i, y_j}^T d_{ij} + \sum_i w_{y_i}^T x_i \quad (1)$$

where w_{y_i, y_j} represent weights that encode valid geometric configurations of object classes y_i and y_j , and w_{y_i} represents a local template for object class i . d_{ij} is a spatial context feature that bins the relative location of window i and j into one of D canonical relations including above, below, overlapping, next-to, near, and far (Fig. 3). Hence d_{ij} is a sparse binary vector of length D with a 1 for the k th element when the k th relation is satisfied between the current pair of windows. w_{y_i, y_j} encodes the valid geometric arrangements of a single class. For example, if people occur beside one another but not above, the weight from w_{y_i, y_j} associated with next-to relations would then be large.

Local model. In our current implementation, rather than learning a local template, we simply use the output of the local detector as the single feature. To learn biases between

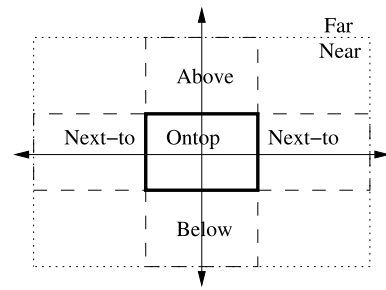


Fig. 3 A visualization of our spatial histogram feature d_{ij} . We consider the location of the center of window j with respect to a coordinate frame defined by window i , denoted by the thickly outlined box. The dashed and dotted rectangles represent regions over which the center of window j are binned. The relative location of j must either be far or near. For near windows, we consider above, ontop, below, and symmetric next-to bins as shown. To allow our model to reproduce the behavior of baseline modules that perform NMS with a criteria of 50% relative overlap, we also include a binary overlap feature. This makes d_{ij} a 7 dimensional sparse binary vector

different object classes, we append a constant 1 to make x_i two-dimensional.

Background class. Since we are concerned only with the relative difference in scores between labelings, we have an extra degree of freedom in defining the weights. We constrain local and pairwise background weights w_0 and w_{i0} and w_{0i} to be 0. Since the majority of windows in an image will be labelled as background, this significantly speeds up computations with the model.

3 Inference

Computing $\arg \max_Y S(X, Y)$ is NP hard unless the pairwise potentials happen to have some particular structure (e.g., super-modularity with $K = 1$). For more general cases, one must resort to search techniques such as branch-and-bound or A^* to find exact minima. In our experiments, we use a simple greedy forward search. We extensively evaluate the effectiveness of our greedy inference procedure in Sect. 7.

3.1 Greedy Forward Search

Our algorithm for optimizing (1) is analogous to greedy algorithms typically used for NMS (Leibe et al. 2004). (1) Initialize the label vector Y to the background class for each window. (2) Greedily select the single window that, when labelled as a non-background class, increases the score S by the largest amount. (3) Stop when instantiating any other detection decreases the total score. Naïvely re-computing the score at each step of the algorithm takes excessively long but we can track the potential gain of adding each detection incrementally.

We write I for a particular set of instanced window-class pairs $\{(i, c)\}$ and write $Y(I)$ for the associated label vector where $y_i = c$ for all pairs in I and $y_i = 0$ otherwise. We define the change in score obtained by adding window-class pair (i, c) to the set of instances I as

$$\Delta(i, c) = S(X, Y(I \cup \{(i, c)\})) - S(X, Y(I))$$

Initialize $I = \{\}$, $S = 0$ and $\Delta(i, c) = w_c^T x_i$ and repeat:

1. $(i^*, c^*) = \arg \max_{(i,c) \notin I} \Delta(i, c)$
2. $I = I \cup \{(i^*, c^*)\}$
3. $S = S + \Delta(i^*, c^*)$
4. $\Delta(i, c) = \Delta(i, c) + w_{c^*,c}^T d_{i^*,i} + w_{c,c^*}^T d_{i,i^*}$

until $\Delta(i^*, c^*) < 0$ or all windows are instanced. In Step 4, we update $\Delta(i, c)$ for un-instanced window-class pairs by adding the pairwise costs due to the newly instanced pair (i^*, c^*) . For additional speed ups, we ran the above algorithm on a set of windows that passed an initial minimal threshold and conservative NMS step. This substantially reduces the number of windows the algorithm must consider.

While this greedy selection procedure can produce sub-optimal results, we have found that in practice it yields quite good solutions. On our datasets, the greedy procedure produces globally optimal solutions on at least 97.6% of the test cases. In Sect. 7 we present an empirical comparison to other optimization techniques as well as discussing theoretical justifications for this good performance.

3.2 Marginals

Many object recognition benchmarks such as PASCAL are scored by ranking detections with a precision-recall curve. This means we need to associate a score with each detected window. To obtain a score, we can appeal to a probabilistic version of our model, which would correspond to a conditional random field (CRF) written as $P(Y|X) = \frac{1}{Z(X)} e^{S(X,Y)}$. One natural score for an individual detection is to use the marginal posterior $P(y_i = c|X)$ however this requires marginalizing over an exponential number of configurations which is intractable in our model. Instead we develop an approximation based on the log-odds ratio¹

$$\begin{aligned} m(y_i = c) &= \log \frac{P(y_i = c|X)}{P(y_i \neq c|X)} \\ &= \log \frac{\sum_{y_r} P(y_i = c, y_r|X)}{\sum_{y_s, c' \neq c} P(y_i = c', y_s|X)} \end{aligned} \tag{2}$$

We write y_r and y_s for a $N - 1$ dimensional vector of labels for the remaining $N - 1$ windows other than i . Both sums

above still require marginalizing out an exponential number of labels, but let us assume the posterior mass inside each sum is dominated by the most probable label y_r^* and the second best label y_s^* with class c^* respectively.

$$\begin{aligned} y_r^* &= \arg \max_{y_r} S(X, y_i = c, y_r) \\ (y_s^*, c^*) &= \arg \max_{(y_s, c' \neq c)} S(X, y_i = c', y_s) \end{aligned} \tag{3}$$

Then the marginal log-odds ratio equation (2) can be approximated by

$$\begin{aligned} m(y_i = c) &\approx \log \frac{P(y_i = c, y_r^*|X)}{P(y_i = c^*, y_s^*|X)} \\ &= S(X, y_i = c, y_r^*) - S(X, y_i = c^*, y_s^*) \end{aligned}$$

It is straightforward to extend our greedy maximization procedure for optimizing (1) to solve (3). This is used for the per detection scoring presented in the result section. In practice, we approximate the marginal by $m(y_i = c) \approx \Delta(i, c)$ computed during the greedy optimization.

4 Learning

In order to describe the learning algorithm, we first re-write the score function from (1) in terms of a single linear parameter vector w . To do this, we encapsulate the effect of Y and X in a potential function, writing

$$S(X, Y) = \sum_{i,j} w_s^T \psi(y_i, y_j, d_{ij}) + \sum_i w_a^T \phi(x_i, y_i) \tag{4}$$

where w_s and $\psi()$ are vectors of length DK^2 , and w_a and $\phi()$ are vectors of length KF , where D is the number of spatial relations, K is the number of classes and F is the length of feature vector x_i . In general, each object class may use a feature vector of different length. The vector $\psi()$ will contain at most D nonzero entries and the vector $\phi()$ will contain only F nonzero entries. We can then write the score as $S(X, Y) = w^T \Psi(X, Y)$ where

$$w = \begin{bmatrix} w_s \\ w_a \end{bmatrix}, \quad \Psi(X, Y) = \begin{bmatrix} \sum_{ij} \psi(y_i, y_j, d_{ij}) \\ \sum_i \phi(x_i, y_i) \end{bmatrix} \tag{5}$$

where our greedy inference procedure solves

$$Y^* = \arg \max_Y w^T \Psi(X, Y) \tag{6}$$

4.1 Convex Training

Assume we are given a collection of training images X_i and labels Y_i . We want to find a model w that, given a new image X_i , tends to produce the true label vector $Y_i^* = Y_i$. We

¹The log-odds and marginals would give the same rank ordering of the detections if exact inference was feasible.

formulate this as a regularized learning problem:

$$\begin{aligned} & \arg \min_{w, \xi_i \geq 0} w^T w + C \sum_i \xi_i \\ & \text{s.t. } \forall_i, H_i \quad w^T \Delta \Psi(X_i, Y_i, H_i) \geq l(Y_i, H_i) - \xi_i \end{aligned} \tag{7}$$

where $\Delta \Psi(X_i, Y_i, H_i) = \Psi(X_i, Y_i) - \Psi(X_i, H_i)$. The constraint from (7) specifies the following: Consider the i th training image X_i and its true label Y_i . We want the true label to score higher than all other hypothesized labelings $\{H_i\}$. However not all incorrect labelings are equally bad. The loss function $l(Y_i, H_i)$ measures how incorrect H_i is and penalizes the slack variable ξ_i in proportion. This loss function formulation from (7) is often called margin-rescaling (Tsochantaridis et al. 2004).

We consider notions of loss that decompose across the N windows: $l(Y, H) = \sum_{i=1}^N l(y_i, h_i)$. One simple window-specific loss is 0-1:

$$l_{01}(y_i, h_i) = I(y_i \neq h_i)$$

Hence, the constraint from (7) requires that label Y scores much higher than those hypotheses H that differ from the ground-truth on many windows. However note that l_{01} incorrectly penalizes detections that overlap true positives as false positives. A more appropriate loss that handles overlap a bit better is:

$$l_{ov}(y_i, h_i) = \begin{cases} 1: & y_i \neq bg \wedge h_i \neq y_i \\ 1: & h_i \neq bg \wedge \neg \exists j \\ & \text{s.t. } [ov(i, j) > .5 \wedge y_j = h_i] \\ 0: & \text{otherwise} \end{cases} \tag{8}$$

The top condition corresponds to a missed detection, while the second corresponds to a false positive (where we check to make there does not exist an overlapping true detection). One may also define a soft loss that assigns a value between 0 and 1 for partially overlapping windows, as in Blaschko and Lampert (2008).

5 Cutting Plane Optimization

Consider the following unconstrained formulation that is equivalent to the constrained problem from (7):

$$\begin{aligned} w^* &= \arg \min_w L(w) \quad \text{where } L(w) = \frac{1}{2} \|w\|^2 + CR(w) \\ R(w) &= \sum_{i=1}^N \max_H (0, l(Y_i, H) - w^T \Delta \Psi(X_i, Y_i, H)) \end{aligned} \tag{9}$$

In the above formulation, $R(w)$ is a convex function since it is the maximum of a set of linear functions and N is the

total number of training examples. This proves that the overall objective function $L(w)$ is convex since it is the sum of two convex functions.

We follow the derivation from Teo et al. (2007) and call (7) the master problem. We define the following reduced problem

$$\begin{aligned} w_t &= \arg \min_w L_t(w) \\ \text{where } L_t(w) &= \frac{1}{2} \|w\|^2 + CR_t(w) \end{aligned} \tag{10}$$

where the convex hinge loss R is approximated by a piecewise linear function R_t . The approximation is constructed from a small set of lower-tangent planes called *cutting planes*. Each cutting plane will be a sub-gradient g of the function $R(w)$ computed at a particular point w_j . The sub-gradient is computed as:

$$\begin{aligned} g(w_j) &= - \sum_{i=1}^N \pi_i \Delta \Psi(X_i, Y_i, H_i^*) \\ \pi_i &= \begin{cases} 1 & \text{if } l(Y_i, H_i^*) - w_j^T \Delta \Psi(X_i, Y_i, H_i^*) \geq 0 \\ 0 & \text{otherwise} \end{cases} \\ H_i^* &= \arg \max_H l(Y_i, H) - w^T \Delta \Psi(X_i, Y_i, H) \end{aligned} \tag{11}$$

where H_i^* is the most violated constraint for image i under the current weight vector w . The subgradient provides a linear lower bound for $R(w)$.

$$R(w) \geq R(w_j) + g(w_j)^T (w - w_j) \quad \forall w \tag{12}$$

To obtain a tighter lower bound of $R(w)$, we will take the point-wise maximum of cutting planes computed at points w_1, \dots, w_{t-1} , adding the zero-plane to the set since the hinge loss R is nonnegative:

$$R_t(w) = \max \left(0, \max_{j=1, \dots, t-1} w^T g(w_j) + b_j \right) \quad \forall w \tag{13}$$

5.1 Dual QP for Cutting Planes

Consider the reduced problem from (10) $L_t(w) = \frac{1}{2} \|w\|^2 + CR_t(w)$, where $R_t(w)$ is as defined in (13). The primal QP can be written as:

$$\begin{aligned} & \arg \min_{w, \xi > 0} \frac{1}{2} \|w\|^2 + C\xi \\ & \text{s.t. } w \cdot g(w_i) + b_i \leq \xi, \quad \forall i = 1, \dots, t \end{aligned}$$

The full Lagrangian and the associated KKT conditions are:

$$\begin{aligned} L(w, \xi, \alpha, \mu) &= \frac{1}{2} \|w\|^2 + C\xi + \sum_{i=1}^t \alpha_i (w \cdot g(w_i) + b_i - \xi) - \mu \xi \end{aligned}$$

Taking the required partial derivatives for the KKT conditions gives:

$$\frac{\partial L}{\partial w} = 0 \implies w = \sum_{i=1}^t \alpha_i g(w_i)$$

$$\frac{\partial L}{\partial \xi} = 0 \implies C \geq \sum_{i=1}^t \alpha_i$$

Plugging in the KKT conditions into the Lagrangian yields the dual QP:

$$\begin{aligned} \arg \max_{\alpha > 0} & -\frac{1}{2} \sum_{i=1}^t \sum_{j=1}^t \alpha_i g(w_i)^T g(w_j) \alpha_j \\ \text{s.t.} & \sum_{i=1}^t \alpha_i \leq C \end{aligned}$$

The solution vector α to the QP is used to recover w using the 1st KKT condition. Note that solving the dual QP of the reduced problem is a function of t variables and is independent of the dimensionality of the feature vector $\Psi(\cdot)$. This makes the cutting plane approach easily scalable to learning from high dimensional features.

5.2 Standard Cutting Plane Algorithm

Initialize $t = 0$ and the set of cutting planes to be empty. Iterate:

1. Compute $w_t = \arg \min_w L_t(w)$ where $L_t(w) = \frac{1}{2} \|w\|^2 + CR_t(w)$. This can be solved with a dual QP with t variables. Since t is typically small (10-100), this can be solved with off-the-shelf solvers. We use the publicly available simplex solver from Franc (2006). Compute $L_t(w_t)$.
2. Compute the subgradient $g(w_t)$ and add the new cutting plane $w^T g(w_t) + b_t$ to the set. Compute $L(w_t)$.

As in Teo et al. (2007), we iterate until the stopping condition $L(w_t) - L_t(w_t) < \epsilon$. Define the optimal solution as $L^* = \min_w L(w)$. It is relatively straightforward to show that $\forall t$, we have the lower and upper bounds $L_t(w_t) \leq L^* \leq L(w_t)$. The iteration must terminate because the lower bound is non-decreasing $L_t(w_t) \geq L_{t-1}(w_{t-1})$.

We give the following intuition behind why the bounds hold: Since w^* is the globally optimal solution for problem (9), $L^* = L(w^*)$. By definition, $L(w^*) = \frac{1}{2} \|w^*\|^2 + CR(w^*)$ and $L_t(w^*) = \frac{1}{2} \|w^*\|^2 + CR_t(w^*)$. Since $R_t(w^*)$ is a point-wise max taken over a bunch of lower tangent planes to $R(w^*)$, we know that $R_t(w^*) \leq R(w^*)$. Therefore $L_t(w^*) \leq L(w^*)$. For any arbitrary w_t , such that $w_t \neq w^*$, the envelope of lower tangent planes will not be as tight as

that constructed using w^* . Mathematically this translates to $L_t(w_t) \leq L_t(w^*)$ whenever $w_t \neq w^*$. Thus

$$L_t(w_t) \leq L_t(w^*) \leq L(w^*)$$

Likewise, since $L()$ is the original cost function that we wish to minimize, when $w_t \neq w^*$, $L(w^*) \leq L(w_t)$ Therefore,

$$L_t(w_t) \leq L_t(w^*) \leq L(w^*) \leq L(w_t)$$

5.3 Online Cutting Plane Algorithm

Note that computing $g(w_t)$ in step 2 of Sect. 5.2 requires knowledge of H_i^* for all the N images in the training set. For large N , this is inefficient both, in terms of the number of computations required, as well as the memory needed to store all $\{H_i^*\}$ before the next true subgradient can be computed. This is in contrast to many online optimization techniques like perceptrons and stochastic gradient descent that are able to learn a reasonably good model without having to make a complete pass through the entire dataset. Motivated by such techniques, we observe that one can construct a cutting plane with a partial subgradient computed from a small number of examples $n \ll N$. We define a partial gradient $g(w_j)$, bias $b(w_j)$, and loss $L(w)$ computed on n examples as follows:

$$g^{(n)}(w_j) = - \sum_{i=1}^n \pi_i(w_j) \Delta \Psi(X_i, Y_i, H_i^*)$$

$$b^{(n)}(w_j) = \sum_{i=1}^n \pi_i(w_j)$$

$$L^{(n)}(w) = \frac{1}{2} \|w\|^2 + CR^{(n)}(w)$$

$$R^{(n)}(w) = \sum_{i=1}^n \max(0, l(Y_i, H_i^*) - w^T \Delta \Psi(X_i, Y_i, H_i^*))$$

We modify the standard cutting plane approach as follows: Initialize $t = 0$ and the set of cutting planes to be empty. Iterate:

1. Identical to step 1 in Sect. 5.2.
2. Iterate through the data in any order until one collects n examples for which $L^{(n)}(w_t) - L_t(w_t) > \epsilon$. Add in the partial cutting plane $w^T g^{(n)}(w_t) + b^{(n)}(w_t)$, and goto Step 1. If the condition is never met, stop.

Because $L^{(n)}(w_t) \leq L(w_t)$, we have that $L^{(n)}(w_t) - L_t(w_t) > \epsilon$ implies $L(w_t) - L_t(w_t) > \epsilon$. This means we only need n examples to discover that w_t is not ϵ -optimal. Once we discover this, we construct $w^T g^{(n)}(w_t) + b^{(n)}(w_t)$.

This cutting plane is a lower-bound to $R(w)$. However, because we have not used all N examples, it is no longer tight. Hence $L_t(w_t) \leq L^*$. If we cannot find N examples that violate ϵ -optimality, then $L^{(N)}(w_t) = L(w_t)$ and w_t is ϵ optimal.

During initial iterations of the algorithm, updates occur very frequently. This is because a single example often suffices to discover that w_t is not ϵ -optimal. Towards later iterations, when w_t is more accurate, updates are less common because n will need to be large to trigger a tolerance violation.

5.4 Finding Most-Violated Constraint

In step (2) of Sect. 5.3, we need to compute the partial subgradient of $R(w)$ at the current w_t . To do so, we need to compute the most violated constraint H_i^* for an image i in (11). Dropping the i subscript notation, we can rewrite (11) as

$$\begin{aligned} H^* &= \arg \max_H l(Y, H) + w^T \Psi(X, H) \\ &= \arg \max_H \sum_{i,j} w_{h_i, h_j}^T d_{ij} + \sum_i (w_{h_i}^T x_i + l(h_i, y_i)) \end{aligned}$$

Since the loss function decomposes into a sum over windows, solving for H^* is very similar to the original maximization (1) except that the local match costs have been augmented by the loss function. Using the loss function in (8), the local scores for invalid object labels for a given window are incremented by one. This makes these labels more attractive in the maximization, and so they are more likely to be included in the most-violated constraint H^* . We can compute an approximation to H^* with a greedy forward search as in Sect. 3.1.

Our algorithm is an under-generating approximation (Finley and Joachims 2008), so there are not formal guarantees optimality. However, as stated in Sect. 3.1, greedy forward search tends to produce scores similar to the brute-force solution, and so we suspect our solutions are close to optimal. A detailed empirical evaluation of our greedy approach is presented in Sect. 7.

6 Results

We have focused our experimental results for multiclass object recognition on the PASCAL Visual Object Challenge. It is widely regarded as the most difficult available benchmark for recognition. We use the 2007 data which is the latest for which test annotations are available. The data consists of 10000 images spanning 20 object classes with a 50% test-train split. The images are quite varied, making this an especially difficult testbed for high-level contextual reasoning.

Baseline: State-of-the-art approaches tend to be scanning window detectors (Everingham et al. 2007). We use the publicly available code (Felzenszwalb 2008) as a baseline. It implements an intra-class NMS post-processing step. The code is an improved version of Felzenszwalb et al. (2008) that out-scores many of the previous best performers from the 2007 competition, suggesting it is a strong baseline for comparison.

Per-class scores: We follow the VOC protocol for reporting results (Everingham et al. 2007). A putative detection is considered correct if the intersection of its bounding box with the ground-truth bounding box is greater than 50% of their union. Multiple detections for the same ground-truth are considered false positives. We compute Precision-Recall (PR) curves and score the average precision (AP) across classes in Table 2. For twelve of the twenty classes, we achieve the best score when compared to the 2007 competition and the baseline model. We also compare to a version of Felzenszwalb et al. (2008) in which detections from multiple classes are pooled before apply-

Table 2 Per-class AP scores on PASCAL 2007 (Everingham et al. 2007). We show the winning score from the 2007 challenge in the *first* data column. This column is composed of various state-of-the-art recognition algorithms. The *second* column is our baseline obtained by running the code from (Felzenszwalb 2008). It outperforms many of the 2007 entries, suggesting it is a strong baseline for comparison. The *third* column pools detections across multiple classes before applying NMS procedure from (Felzenszwalb 2008) (MC-NMS). The *third* column is our approach, which provides a stark improvement over MC-NMS and generally improves performance over classification-trained approaches

Class		Baseline	MC-NMS	Our model
Plane	.262	0.278	0.270	0.288
Bike	.409	0.559	0.444	0.562
Bird	.098	0.014	0.015	0.032
Boat	.094	0.146	0.125	0.142
Bottle	.214	0.257	0.185	0.294
Bus	.393	0.381	0.299	0.387
Car	.432	0.470	0.466	0.487
Cat	.240	0.151	0.133	0.124
Chair	.128	0.163	0.145	0.160
Cow	.140	0.167	0.109	0.177
Table	.098	0.228	0.191	0.240
Dog	.162	0.111	0.091	0.117
Horse	.335	0.438	0.371	0.450
Motbike	.375	0.373	0.325	0.394
Person	.221	0.352	0.342	0.355
Plant	.120	0.140	0.091	0.152
Sheep	.175	0.169	0.091	0.161
Sofa	.147	0.193	0.188	0.201
Train	.334	0.319	0.318	0.342
TV	.289	0.373	0.359	0.354

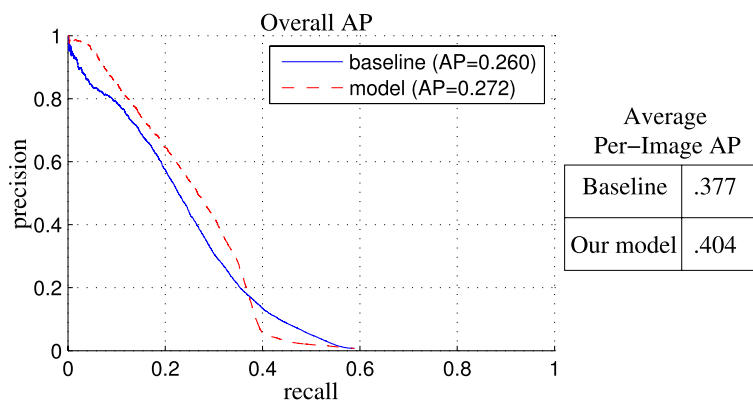
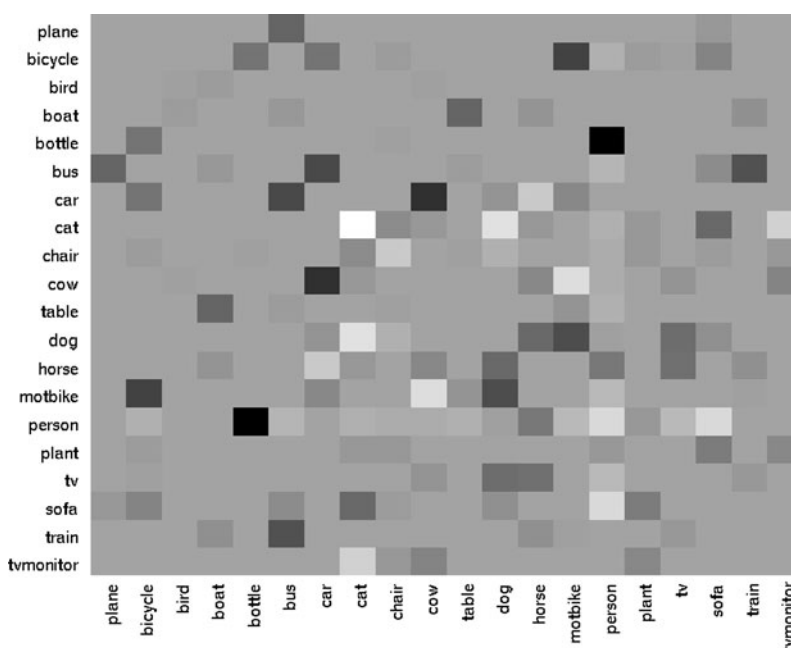


Fig. 4 Multi-class AP scores on PASCAL 2007. On the *left*, we score overall AP. We construct the baseline curve by pooling detections across classes and images when computing PR curves. Our global model clearly provides a noticeable boost in performance in the low-

recall high-precision regime. On the *right*, we pool detections on a per-image base, compute the per-image AP, and average the result over images. We see a noticeable improvement of 10% over our baseline (Felzenszwalb 2008)

Fig. 5 We visualize the weights for our overlap threshold across all our models. Light areas correspond to an increase in score. The structure in these weights indicate the subtlety required for applying mutual exclusion across classes. For example, because people and bottles have similar shapes, the local detectors we use (Felzenszwalb et al. 2008) can confuse them. Our global model learns to strongly compete such overlapping detections using a negative weight. However, people and sofas tend to overlap because people partially occlude sofas when sitting down. In this case, we learn a positive weight that reinforces both detections



ing NMS (MC-NMS). This tends to hurt performance, indicating the need for proper training of multiclass inhibition. The improvement over MC-NMS is generally large. In most cases, the improvement over the baseline is small, but for indoor classes such as tables and bottles and outdoor classes such as motorbikes and trains, the improvement is close to 10%.

Multi-class scores: Per-class APs do not score the consistency of detections across classes on an image, which is one of our goals for multi-class recognition. We consider two approaches for multiclass scores in Fig. 4. First we pool detections across classes and images (running the default NMS procedure in Felzenszwalb (2008) before pooling), and generate a single PR curve. Our model provides a noticeable

improvement, particularly in the high precision—low recall regime. We also pool detections on a per image bases, generating a per-image multi-class AP. We average this AP across all images. Our model again provided a strong improvement of 10% over the baseline. This is because the baseline does not correctly reconcile detections from various classes due to the fact that the detectors were trained independently.

Models: We visualize the pairwise weights learned in our models in both Figs. 5 and 6. These are trained discriminatively, taking into account the behavior of the local detector. For example, our model learns to aggressively compete bottle and person detections because local detectors confuse the two. This is contrast to simple co-occurrence weights that are trained by frequency counting as in Galleguillos et al.

Fig. 6 We visualize the pairwise spatial weights for each pair of classes as a 5×5 image (analogous to Fig. 3). Light areas indicate a favorable arrangement. We show a closeup for particular relations from classes where the global model helps performance. On the *top*, we see that bottles tend to sit above tables. In the *middle*, cars lie both near and far from trains, but rarely above or directly next to them. On the *bottom*, we see that motorbikes tend to occur next to one another in images

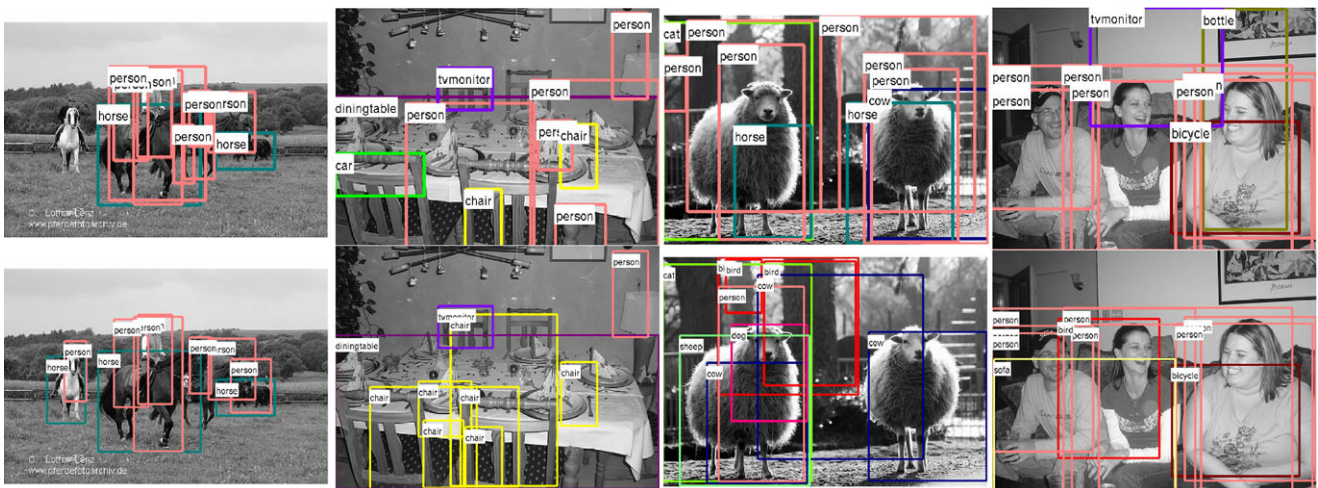
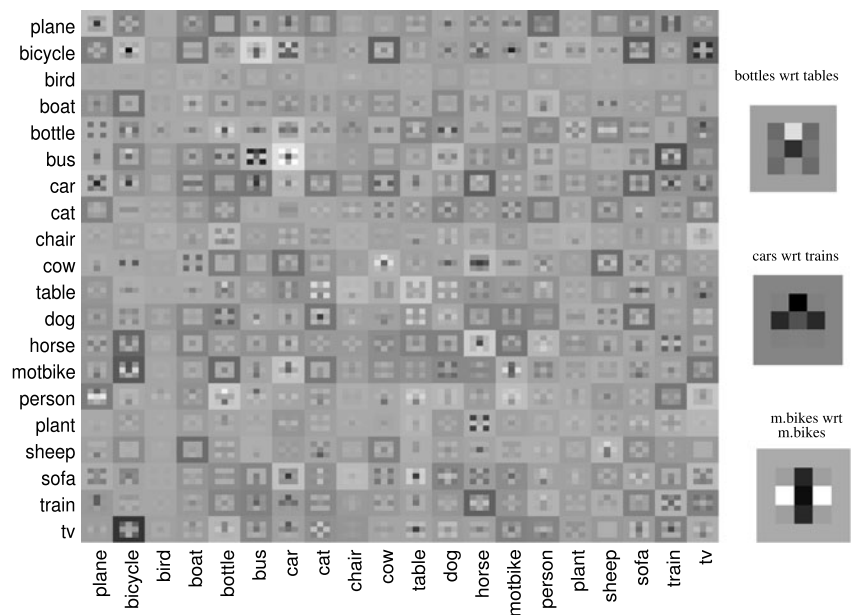


Fig. 7 Example test images. On the *top row*, we show the top 10 detections from the baseline model after standard NMS. On the *bottom row*, we show the top 10 marginal detections from our global model. On the *left*, we see that horse and person detections are better localized by the globally tuned NMS model. In the *left center*, our model seems to favor

patterns of chair detections that overlap, as maybe common in scenes of tables. In the *right center*, our model exploits co-occurrence cues favoring groups of animals. Finally, on the *right*, our model appears to be exploiting relational cues about sofas and people while enforcing mutual exclusion between the bottle and people detections

(2008), Baur et al. (2008). We also learn meaningful multiclass spatial layouts—e.g., bottles tend to occur above tables. We refer the reader to the captions for additional analysis. Figure 7 shows example multi-class detections from our model as compared to the baseline. Our model appears to produce better detections by understanding interactions between objects that spatially overlap, such as people when riding horses. It also learns how to correctly enforce mutual exclusion between classes, allowing people and sofas to overlap but not people and bottles.

Does context help? Our results suggest that the benefit from using context to improve per-class AP is only marginal

on PASCAL. We provide a couple of hypotheses as to why this is so:

1. *Contextual layout models are better suited for images with multiple objects.* The PASCAL dataset is somewhat impoverished as far as presence of sufficient inter-class and intra-class context is concerned. The PASCAL dataset contains 20 object classes. However, more than half the images contain only a single object class with two instances of that object class typically present in the image. We agree with the sentiment from Choi et al. (2010) that “contextual information is most useful

when many different object categories are present simultaneously in an image, with some object instances that are easy to detect (i.e. large objects) and some instances that are hard to detect (i.e. small objects).” Along similar lines, Park et al. (2010) suggest that context provides a stronger improvement for detecting small objects rather than large objects. We hypothesize that our models may similarly exhibit a stronger improvement on datasets containing such variety.

2. *Context is more useful for higher-level semantic tasks.* Our baseline local detectors are state-of-the-art models that have consistently produced competitive results on PASCAL in terms of per class AP. We believe that contextual reasoning may only provide limited improvement over highly-tuned local detectors when scored for tasks such as object detection and localization. This view is corroborated by other empirical evaluations of context using tuned local detectors (Divvala et al. 2009) and (Galleguillos et al. 2008). However, we argue that context is helpful for higher level semantic inferences such as scene or action understanding. In the extreme case, given a perfect person and bottle detector, context cannot improve detection performance of either class. But even given such perfect detectors, one still requires a contextual layout model to recognize “drinking” actions because people and bottles must be simultaneously found in particular spatial relationships.

Our per-image AP and overall AP scores partially validate the second hypothesis. Per-image AP scores can be interpreted as a loose proxy for a holistic understanding of a “scene” since one must reconcile detections across multiple classes simultaneously. Under this criteria, our model does improve AP from 37% to 40%, which is noticeably stronger than the 1% improvement in overall AP. In subsequent work (Desai et al. 2010), we further investigate the effect of contextual layout models on the high-level semantic task of action recognition. We demonstrate that the contextual models developed here can be used to increase the accuracy of a static-image action classifier by 12%. Notably, this increase is obtained over a baseline using the exact same state-of-the-art local detectors used here. Hence we believe that our contextual layout model is more rewarding when used for higher level semantic tasks.

7 Analysis of Greedy Inference

We compare our greedy inference algorithm to two other approximate inference approaches: Loopy Belief Propagation (LBP) and Tree Re-Weighted Belief Propagation (TRW) (Wainwright et al. 2002; Kolmogorov 2006). Although LBP has been widely used for approximate inference in graphical models with cycles, LBP is not guaranteed to converge and

Table 3 Table comparing the average energy, precision and recall across different approximation techniques

	Av. score	Av. prec	Av. recall
Greedy	1.174	0.7939	0.4673
TRW-S	1.185	0.771	0.4707
LBP	1.185	0.771	0.4707

is susceptible to getting trapped at non-optimal fix points. In the TRW approach, the original MAP problem is initially formulated as an integer program, whose binary constraints are “relaxed” to give a linear program (LP). The TRW algorithm is a variant of Belief Propagation that solves the resulting LP and has been shown to be significantly faster for this problem structure than off-the-shelf LP solvers (Yanover and Meltzer 2006). The solution given by TRW provides an upper bound on the solution of our score maximization problem. Notably, if the solution to the LP relaxation is integral, then the bound is tight and the solution is guaranteed to be a global optimum.

We took the model learned using the approach discussed in Sect. 5 and ran the 3 approximation techniques: Greedy, LBP and TRW on the test set comprising 4952 images from PASCAL VOC 2007 dataset. We used the publicly available software from Meltzer (2006) for LBP and the software from MSR (2006) for TRW-S.² Table 3 compares the approximation techniques in terms of how well they maximize the score function, and their accuracy on PASCAL 2007 test set. For 4942 out of 4952 images (99.8%), Greedy and LBP yield identical results. More importantly, for 4832 of the 4952 images (97.6%), all the three schemes produce identical labels. For these cases, we verified that TRW-S produces integer solutions. This means that greedy produces the *provably globally optimal* solution in almost all images, while being two orders of magnitude faster than either approach.

One theoretical explanation for the near-optimal performance of the greedy search procedure comes from the study of maximizing *sub-modular set functions*. While such problems are NP hard in general, simple greedy heuristics can be shown to have strong approximation guarantees (Nemhauser et al. 1978). If the pairwise weights are all ≤ 0 , then one can show that $S(X, Y)$ from (1) is a submodular set function because it satisfies the *diminishing returns* property: consider two sets of instanced windows I_1 and I_2 , where $I_1 \subseteq I_2$, and a particular un-instanced window i . The increase in $S(X, Y)$ due to instancing i must be smaller for I_2 because all pairwise interactions are negative. This means that greedy inference algorithms enjoy strong theoretical guarantees for contextual layout models with solely negative interactions. In practice, we observe that 90% of all the pairwise weights

²We also tested QPBO (Rother et al. 2007) which gave similar results.

associated with the model trained on PASCAL 2007 images are ≤ 0 , which makes $S(X, Y)$ close to sub-modular and our greedy maximization algorithm theoretically well-motivated.

8 Discussion and Related Work

There has been a wide variety of work in the last few years on contextual modeling in image parsing (Torralba et al. 2004; Sudderth et al. 2005; Hoiem et al. 2008; Galleguillos et al. 2008; Shotton et al. 2006; He et al. 2004; Anguelov et al. 2005). These approaches have typically treated the problem as that of finding a joint labeling for a set of pixels, super-pixels, or image segments and are usually formulated as a CRF. Such CRFs for pixel/segment labeling use singleton potential features that capture local distributions of color, textons, or visual words. Pairwise potentials incorporate the labelings of neighboring pixels but in contrast to older work on MRFs these pairwise potentials may span a very large set of neighboring sites (e.g. Torralba et al. 2004; Tu 2008). Learning such complicated potentials is a difficult problem and authors have relied primarily on boosting (Shotton et al. 2006; Torralba et al. 2004; Tu 2008) to do feature selection in a large space of possible potential functions.

These approaches are appealing in that they can simultaneously produce a segmentation and detection of the objects in a scene. Thus they automatically enforce NMS and hard mutual exclusion (although as our examples show, this may not be entirely desirable). However, the discriminative power of these methods for detection seems limited. While local image features work for some object classes (grass, sky etc.), a clear difficulty with the pixel/segment labeling approach is that it is hard to build features for objects defined primarily by shape. It still remains to be shown whether such approaches are competitive with scanning window templates on object detection benchmarks.

In principle, one could define unary potentials for CRFs using, say, HOG templates centered on individual pixels. However, the templates must score well when centered on every pixel within a particular segment. Thus templates will tend to be overly-smoothed. Our method is fundamentally different in that the output is *sparse*. A complete object detection is represented by the activation of a single pixel and so the unary potential can be quite strong. Furthermore, a detection in our model represses detections corresponding to small translations while, in the pixel labeling model, exactly the opposite has to happen. We thus make a tradeoff, moving to more powerful discriminative unary features but sacrificing tractable pairwise potentials.

Alternatively, (Galleguillos et al. 2008; Kumar and Hebert 2005) group pixels into object-sized segments and

then define a CRF over the labels of the segments. This approach has the advantage that unary potentials can now be defined with object templates, say, centered on the segment. However, the initial segmentation must be fairly accurate and enforces NMS and mutual exclusion without object-level layout models.

To our knowledge, the problem of end-to-end learning of multi-object detection (i.e. learning NMS) has not been explored. The closest work we know of is that of Blaschko and Lampert (2008) who use structured regression to predict the bounding box of a single detection within an image. Both models are trained using images rather than cropped windows. Both are optimized using the structural SVM formalism of Tsochantaridis et al. (2004). However, the underlying assumptions and resulting models are quite different. In the regression formalism of Blaschko and Lampert (2008), one assumes that each training image contains a single object instance, and so one cannot leverage information about the layout of multiple object instances, be it from the same class or not. The models may not perform well on images without the object because such images are never encountered during training. In our model, we can use all bounding-box labels from all training images, including those that do not contain any object, to train a model that will predict those very labels.

9 Conclusion

We have presented a system for multi-class object detection with spatial interactions that can be efficiently trained in a discriminative, end-to-end manner. This approach is able to fuse the outputs of state of the art template based object detectors with information about contextual relations between objects. Rather than resorting to post-processing to clean up detections, our model learns optimal non-max suppression parameters and detection thresholds for each class. The resulting system outperforms published results on the PASCAL VOC 2007 object detection dataset.

Acknowledgements This work was supported by NSF grant 0812428, a gift from Microsoft Research and the UC Lab Fees Research Program.

References

- Anguelov, D., Taskar, B., Chatalbashev, V., Koller, D., Gupta, D., Heitz, G., & Ng, A. (2005). Discriminative learning of Markov random fields for segmentation of 3d scan data. In *CVPR, II* (pp. 169–176).
- Baur, R., Efros, A. A., & Hebert, M. (2008). *Statistics of 3d object locations in images* (Tech. Rep. CMU-RI-TR-08-43). Robotics Institute, Pittsburgh, PA.

- Blaschko, M. B., & Lampert, C. H. (2008). Learning to localize objects with structured output regression. In *ECCV* (pp. 2–15). Berlin: Springer.
- Choi, M., Lim, J., Torralba, A., & Willsky, A. (2010). Exploiting hierarchical context on a large database of object categories. In *IEEE conference on computer vision and pattern recognition, CVPR*
- Dalal, N., & Triggs, B. (2005). Histograms of oriented gradients for human detection. In *CVPR I* (pp. 886–893).
- Desai, C., Ramanan, D., & Fowlkes, C. (2009). Discriminative models for multi-class object layout. In *IEEE international conference on computer vision*.
- Desai, C., Ramanan, D., & Fowlkes, C. (2010). Discriminative models for static human-object interactions. In *Workshop on structured prediction in computer vision, CVPR*.
- Divvala, S., Hoiem, D., Hays, J., & Efros, A. (2009). An empirical study of context in object detection. In *CVPR*.
- Everingham, M., Van Gool, L., Williams, C. K. I., Winn, J., & Zisserman, A. (2007). The PASCAL visual object classes challenge 2007 (VOC2007) results. <http://www.pascal-network.org/challenges/VOC/voc2007/workshop>.
- Felzenszwalb, P. (2008). <http://people.cs.uchicago.edu/pff/latent>.
- Felzenszwalb, P., McAllester, D., & Ramanan, D. (2008). A discriminatively trained, multiscale, deformable part model. In *CVPR*.
- Finley, T., & Joachims, T. (2008). Training structural svms when exact inference is intractable. In *Proceedings of the 25th international conference on machine learning* (pp. 304–311). New York: ACM.
- Franc, V. (2006). <http://cmp.felk.cvut.cz/xfrancv/libqp/html>.
- Galleguillos, C., Rabinovich, A., & Belongie, S. (2008). Object categorization using co-occurrence, location and appearance. In *CVPR*, Anchorage, AK.
- Hall, E. (1966). *The hidden dimension*. New York: Anchor Books.
- He, X., Zemel, R., & Carreira-Perpinan, M. (2004). Multiscale conditional random fields for image labeling. In *CVPR* (Vol. 2). Los Alamitos: IEEE Comput. Soc.
- Hoiem, D., Efros, A., & Hebert, M. (2008). Putting objects in perspective. *IJCV*, 80(1), 3–15.
- Joachims, T., Finley, T., & Yu, C. (2009). Cutting plane training of structural SVMs. *Machine Learning*, 77(1), 27–59.
- Kolmogorov, V. (2006). Convergent tree-reweighted message passing for energy minimization. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 28, 1568–1583. <http://doi.ieeecomputersociety.org/10.1109/TPAMI.2006.200>.
- Kumar, S., & Hebert, M. (2005). A hierarchical field framework for unified context-based classification. In *Tenth IEEE international conference on computer vision, ICCV, 2005* (Vol. 2).
- Leibe, B., Leonardis, A., & Schiele, B. (2004). Combined object categorization and segmentation with an implicit shape model. In *Workshop on statistical learning in computer vision, ECCV* (pp. 17–32).
- Liu, Y., Lin, W., & Hays, J. (2004). Near-regular texture analysis and manipulation. *ACM Transactions on Graphics*, 23(3), 368–376.
- Meltzer, T. (2006). <http://www.cs.huji.ac.il/talyam/inference.html>.
- MSR (2006). <http://research.microsoft.com/en-us/downloads/dad6c31e-2c04-471f-b724-ded18bf70fe3/>.
- Murphy, K., Torralba, A., & Freeman, W. (2003). Using the forest to see the trees: a graphical model relating features, objects and scenes. *NIPS* 16.
- Nemhauser, G., Wolsey, L., & Fisher, M. (1978). An analysis of approximations for maximizing submodular set functions. *Mathematical Programming*, 14(1), 265–294.
- Park, D., Ramanan, D., & Fowlkes, C. (2010). Multiresolution models for object detection. In *ECCV*.
- Rother, C., Kolmogorov, V., Lempitsky, V., & Szummer, M. (2007). Optimizing binary mrfs via extended roof duality. In *CVPR*.
- Rowley, H. A., Baluja, S., & Kanade, T. (1996). Neural network-based face detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20, 23–38.
- Shotton, J., Winn, J., Rother, C., & Criminisi, A. (2006). Textonboost: joint appearance, shape and context modeling for multi-class object recognition and segmentation. *Lecture Notes in Computer Science*, 3951, 1.
- Sudderth, E., Torralba, A., Freeman, W., & Willsky, A. (2005). Learning hierarchical models of scenes, objects, and parts. In *ICCV, II* (pp. 1331–1338).
- Teo, C., Smola, A., Vishwanathan, S., & Le, Q. (2007). A scalable modular convex solver for regularized risk minimization. In *SIGKDD*. New York: ACM.
- Torralba, A., Murphy, K., & Freeman, W. (2004). Contextual models for object detection using boosted random fields. *NIPS*.
- Tsochantaridis, I., Hofmann, T., Joachims, T., & Altun, Y. (2004). Support vector machine learning for interdependent and structured output spaces. In *ICML*. New York: ACM.
- Tu, Z. (2008). Auto-context and its application to high-level vision tasks. In *CVPR*.
- Viola, P. A., & Jones, M. J. (2004). Robust real-time face detection. *IJCV*, 57(2), 137–154.
- Wainwright, M., Jaakkola, T., & Willsky, A. (2002). Map estimation via agreement on (hyper)trees: message-passing and linear programming approaches. *IEEE Transactions on Information Theory*, 51, 3697–3717.
- Yanover, C., & Meltzer, T. Y. W. (2006). Linear programming relaxations and belief propagation—an empirical study. In *JMLR* (pp. 1887–1907).