



Hans Lobel

Machine Learning Group Pontificia Universidad Católica de Chile

Our Contribution

- Common recognition methods use a 2 level hierarchy: • A visual dictionary learned in an unsupervised fashion. • A top-level classifier, generally using a one-vs-all scheme.
- We present a method to learn this hierarchy, where both levels are jointly optimized to find a task-specific representation.
- Experimental evaluation shows an performance increase, leading to state-of-the-art performance using smaller dictionaries.



Hierarchical Joint Max-Margin Learning of Mid- and Top-Level Representations for Visual Recognition

René Vidal

Center for Imaging Science The Johns Hopkins University

Image Encoding

- We model each image with: \circ A set of L rectangular regions defined over an image I. • A set of local features, $\{v_j^I\}$, extracted from an image I.
- We encode a rectangular region *l* defined over *I* as follows:

$$v_{l,\Theta,z}(I) = [\langle \theta_1, v_{z_{(l,1)}}^I \rangle, \langle \theta_1, v_{z_{(l,1)}}^I \rangle]$$

where variables $z_{(l,k)}^{I} = \operatorname{argmax} \langle \theta_k, v_j^{I} \rangle$ induce max-pooling.

Learning Problem

a regularized max-margin learning problem:

$$\begin{split} \min_{W,\Theta,\{\xi_i\}} \quad &\frac{1}{2} \|W\|_F^2 + \frac{C_1}{2K} \|\Theta\|_F^2 + \frac{C_2}{N} \sum_{i=1}^N \xi_i \\ \text{s.t.} \quad &\sum_l^L \sum_k^K w_{y^i,l,k} \langle \theta_k, v_{z_{(l,k)}^{I^i}}^{I^i} \rangle - \sum_l^L \sum_k^K w_{y,l,k} \langle \theta_k, v_{z_{(l,k)}^{I_i}}^{I_i} \rangle \geq \Delta(y_i, y) - \xi_i \\ &\forall i \in \{1, \dots, N\} \land \forall y \in \{1, \dots, M\} \land \forall z. \\ &w_{y,l,k} \geq 0, \forall y \in \{1, \dots, M\} \land \forall l \in \{1, \dots, L\} \land \forall k \in \{1, \dots, K\}. \end{split}$$

• We use an alternating minimization approach: $_{\odot}$ When Θ is fixed, we solve a standard multi-class SVM for W. $_{\circ}$ When W is fixed, we solve a Latent Structural SVM for Θ .

important Initial Dictionary

• A visual dictionary $\Theta = [\theta_1 \ \theta_2 \ \theta_3 \ \dots \ \theta_K]$, formed by linear classifiers.

 $\theta_2, v_{z_{(l,2)}}^I \rangle, \ldots, \langle \theta_K, v_{z_{(l,K)}}^I \rangle]$

• Given a set of training examples $\{I^i, y^i\}_{i=1}^N$, we find object classifier parameters W and dictionary parameters Θ by solving

Machine Learning Group Pontificia Universidad Católica de Chile





Álvaro Soto

Results For Scene/Object Categorization

• Datasets: Caltech 101, 15 Scene Categories, MIT67 Indoor. • Images downsampled to no more than 300 pixels.

- HOG+LBP descriptors extracted densely.
- L rectangular regions selected randomly.



Performance increases with the number of regions up to a certain amount. Random regions give better performance than Spatial Pyramid at similar number of regions.



Maximum performance is achieved with a small dictionary size.

Comparison to Other Methods

		Dataset		
od	# Words	Caltech101	15 Scenes	MIT67
ine	400	63.9 ± 0.6	78.1 ± 0.3	33.2
/	400	64.6 ± 0.8	81.4 ± 0.5	-
	2048	73.4	80.5 ± 0.6	-
R	1024	73.2 ± 0.8	82.7 ± 0.5	-
Μ	1024	73.2 ± 0.5	80.3	-
argin	5250	_	82.7 ± 0.5	-
Bank	200 per class	_	80.9	37.6
le Models	200	_	78.6 ± 0.7	37.9
ve Patches	210 per class	_	-	38.1
sed	200 or 300	72.9 ± 0.6	84.6 ± 0.4	39.5

Compared to other related methods, we obtain state-of-the-art performance on two out of three datasets, using a much smaller dictionary size.