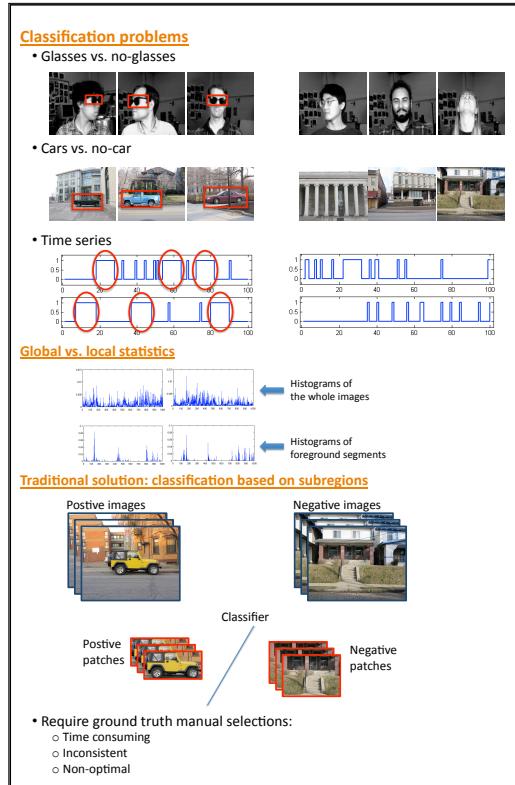


Weakly supervised discriminative localization and classification: a joint learning process

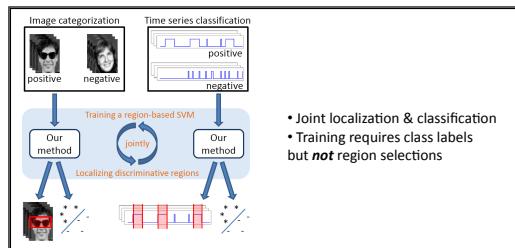
Minh Hoai Nguyen¹, Lorenzo Torresani², Fernando de la Torre¹, & Carsten Rother³

¹Carnegie Mellon University, ²Dartmouth College, ³Microsoft Research Cambridge

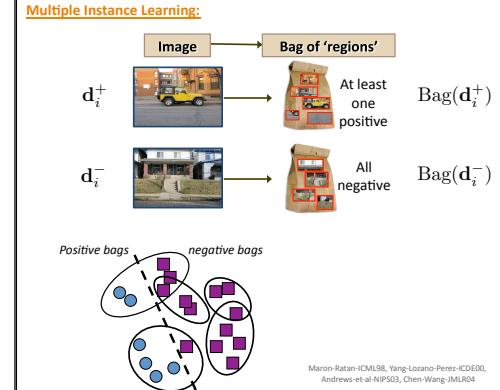
Motivation



Overview



Proposed method



Localization-classification SVM: the formulation

$$\begin{aligned} \text{minimize}_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t.} \quad & \max_{\mathbf{x} \in \text{Bag}(d_i^+)} \mathbf{w}^T \varphi(\mathbf{x}) + b \geq 1 \quad \forall i \\ & \max_{\mathbf{x} \in \text{Bag}(d_i^-)} \mathbf{w}^T \varphi(\mathbf{x}) + b \leq -1 \quad \forall i \end{aligned}$$

Optimization

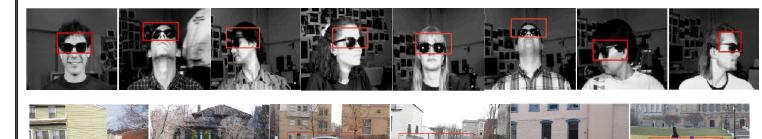
- Two iterative loops
 - Outer loop – coordinate descent: alternate between optimizing (\mathbf{w} , b) and the instances of positive bags that maximize the SVM scores.
 - Inner loop – constraint generation: add the most violated constraints into the constraint sets.
- Each iteration requires $\hat{\mathbf{x}} = \underset{\mathbf{x} \in \text{Bag}(d)}{\operatorname{argmax}} \mathbf{w}^T \varphi(\mathbf{x})$

Representation & Efficient localization

$$\begin{aligned} \varphi(\mathbf{x}) &= \text{Histogram of visual words} \\ \text{Bag}(\text{Image}) &= \text{all possible subwindows} \\ &\text{Efficient search using branch-and-bound} \\ &\text{Average 100ms/image (480*640 pixel)} \\ \text{Bag}(\text{Time series}) &= \left\{ \text{At most } k \text{ disjoint intervals} \right\} \\ &\text{Efficient search using dynamic programming} \\ &\text{Average 10ms/signal (15000 frames)} \end{aligned}$$

Experiments

Images – Localization results



Images – Quantitative results

Dataset	Measure	Bag of words	SVM with global statistics	SVM with human labels	Ours
Faces	Acc. (%)	80.11	82.97	86.79	90.0
	ROC Area	n/a	0.90	0.94	0.96
Cars	Acc. (%)	77.5	80.75	81.44	84.0
	ROC Area	n/a	0.86	0.88	0.90
Caltech-4 datasets	Acc. (%)	89.74	96.05	89.40	96.05
	ROC Area	n/a	0.99	0.95	0.99
	Acc. (%)	94.93	98.17	n/a	98.28
	ROC Area	n/a	1.00	n/a	1.00
	Acc. (%)	59.83	88.70	86.78	89.57
	ROC Area	n/a	0.95	0.91	0.95
Motorbikes	Acc. (%)	76.80	88.99	84.67	87.81
	ROC Area	n/a	0.95	0.92	0.94

Several Conclusions

- Human labels often are not optimal
- Tight bounding boxes often are not optimal; contextual information is important.
- Segmentation does not always help. Our method determines automatically the optimal support region for classification

Synthetic time series data

max # of disjoint intervals allowed		Using global statistics: ROC: 0.577, Acc: 66.5%						
k	1	2	3 to 7	8	12	16	20	
Acc. (%)	77.0	93.0	100	98.5	91.5	77.5	67.25	
ROC Area	.843	.980	1.00	.998	.933	.793	.613	

Several Conclusions

- Segmentation does help.
- Multiple disjoint intervals are necessary.
- Classification performance is not too sensitive to the number of maximum disjoint intervals allowed.

Mouse activity classification



Action	Dollár et al. [8]	1-NN	SVM	Ours
Drink	0.63	0.58	0.63	0.67
Eat	0.92	0.87	0.91	0.91
Explore	0.80	0.79	0.85	0.85
Groom	0.37	0.23	0.44	0.54
Sleep	0.88	0.95	0.99	0.99

$$F_1 = \frac{2 \times \text{Recall} \times \text{Precision}}{\text{Recall} + \text{Precision}}$$

F_1 score