Biologically Plausible Detection of Amorphous Objects in the Wild

> Sunhyoung Han, Nuno Vasconcelos Statistical Visual Computing Lab. University of California, San Diego



PandaCam dataset



PandaCam dataset

- Video feed provided by San Diego Zoo showing real time movement of a Panda family in a natural habitat
- 24/7 surveillance type real application
 - Dramatic change of illumination
 - Low quality images
 - Video is a collection from multiple cameras and viewpoint
- Textureless amorphous object in highly structured background
 - Highly deformable shape of panda bear
 - Background have richer combination of structure, shape, texture



Outline

- Biologically plausible Discriminant saliency network

Connection with existing object recognition models

- HMAX, convolutional neural network
- Unique property of DiscSalNet
 - Detection of feature absence
- Experimental results & conclusion



Discriminant Saliency network

 Discriminant saliency implemented in standard neurophysiological model

- we are investigating a discriminant view of perception
 - brains are optimized for decision-making
 - this holds at all levels of neural circuitry from single cells to populations and high-level circuits
- Hierarchical network architecture



Neurophysiology of vision

- standard architecture of the primary visual cortex (V1)
 - a cascades of linear filter, divisive normalization, a quadratic nonlinearity, and spatial pooling [Carandini et al., 2005]



Natural stimuli statistics

- whenever we pass a natural image through a band-pass filter (Gabor, wavelet, DCT, PCA, etc.)
- response follows generalized Gaussian density (GGD)

$$p(x;lpha,eta)=rac{eta}{2lpha \Gamma(1/eta)} \exp\left(-(|x|/lpha)^eta
ight)$$

- parameters: α scale, β shape
- heavy tails for $\beta < 1$
- extremely consistent observation (any natural images, any filter)



Generalized Gaussian Density (GGD)

- GGDs have two interesting properties
 - negative log-likelihood is proportional to response

$$-\log p(x; lpha, eta) = \left(rac{|x|}{lpha}
ight)^{eta} + K$$

- parameter estimation:
 - using a conjugate (Gamma) prior for the inverse scale ($\theta = 1/\alpha^{\beta}$)

$$P_{\Theta}(\theta) = Gamma\left(\theta, 1 + \frac{\eta}{\beta}, \nu\right) = \frac{\nu^{1+\eta/\beta}}{\Gamma(1+\eta/\beta)} \theta^{\eta/\beta} e^{-\nu\theta}$$

• MAP estimate of α from sample of iid responses {x(1),...,x(n)}

$$\widehat{lpha} \propto \left(\sum_j |x(j)|^eta +
u
ight)^{1/eta}$$

is the norm of the vector of those responses

SVCL ₹UCSD

Simple cells

- even more interesting:
 - when we put the two together we obtain divisive normalization

 $-\log p(x;\{x(j)\}) \propto rac{|x|^eta}{\sum_j |x(j)|^eta+
u}$

- simple cell computations
 - probability of stimulus in cell's receptive field
 - under the GGD of the "normalizing responses" {x(j)}



simple cell

[Hubel & Wiesel 1962, Heeger 1992, Schwartz & Simoncelli 2001]

- from this
 - we can build all components of statistical inference
- classification/decision rules:
 - Log-likelihood ratio

$$g(x) \;\;=\;\; \log rac{P_{X|Y}(x|1)}{P_{X|Y}(x|-1)}$$

- inhibitory connections:
 learning of null
 hypothesis (Y=-1)
- excitatory connections: learning alternative hypothesis (Y=1)



• classification: • classification: -Log-likelihood ratio -posterior probability $\log \frac{1}{P_{X|Y}'}$ $P_{X|Y}(x|1)$ $P_{Y|X}(1|x) =$ g(x)= $P_{X|Y}(x|1)$ $\overline{P_{X|Y}(x|1) + P_{X|Y}(x|-1)}$ = $\frac{1}{1+\frac{P_{X|Y}(x|-1)}{P_{X|Y}(x|-1)}}$ $(1 + \exp\{-g(x)\})^{-1} = \sigma(x)$ _ $\sigma(x)=(1+e^{-x})$ $\sigma(.)$ Τŀ $g[x_j]$

SVCL ₹UCSD





Risks

- given a sample $\mathcal{D} = \{x_1, ..., x_n\}$
 - these can be used to compute various risks

$$\tilde{\psi}(x) = \begin{cases} \frac{1}{2} \log \frac{x}{1-x}, & x \ge .5\\ 0, & x < .5 \end{cases}$$

Empirical risks based on sample $\mathcal{R} = \{x_1, \dots, x_n\}$									
expected NLL	$E_X\left[-\log P_X(x)\right]$	$-rac{1}{n}\sum_{i=1}^n l_{lpha,eta}(x_i)$	entropy $H[X]$						
expected LLR	$E_X \left[\log \frac{P_{X Y}(x 1)}{P_{X Y}(x 0)} \right]$	$rac{1}{n}\sum_{i=1}^n g(x_i)$	$\begin{aligned} KL[P_X(x) P_{X Y}(x 0)] \\ -KL[P_X(x) P_{X Y}(x 1)] \end{aligned}$						
MI	$E_{X}\left[I(Y;X=x)\right]$	$\frac{1}{2}\sum_{i=1}^{n} \xi\{\sigma[q(z_i)]\}$	I(Y;X)						
expected confidence (LLR)	$E_X[LLR(x)]$	$\frac{1}{n}\sum_{i=1}^{n}\tilde{\psi}\{\sigma[g(x_i)]\}$	$KL[P_{X Y}(x 1) P_{X Y}(x 0)]$						
expected confidence (MI)	$E_X[IC(x)]$	$rac{1}{n}\sum_{i=1}^n ilde{\xi}\{\sigma[g(x_i)]\}$							

- non-linearity plus pooling
- computations of the complex cell





Discriminant saliency Layer

- SC network:
 - Gabor decomposition into orientation channels
 - simple units compute target probabilities P_{Y|X}(1|X)
 - complex units compute risk



SVCL ₹UCSD

HMAX units

- similar to DS units but:
 - simple units do trivial Gabor filtering
 - complex units more like feature detectors (max pooling) than risk assessors



SVCL ₹UCSD

Comparison with other algorithms

Discriminant Saliency layer vs. extended neural network



Discriminant Saliency network



- two layers of simple + complex units
 - layer 1: simple featurebased DS units (gabor, DCT basis)
 - layer 2:
 - templates sampled randomly from L1 training responses
 - template-based DS units
- Hierarchical architecture shared with HMAX, convNN, SPMK





SVCL ₹UCSD

Experiment: localization



Results: localization performance



Precision Recall curve



Results: Panda detection

Post processing

- box filter is applied to 7 scale saliency maps
- center location and scale of largest saliency is detected
- non-maximum suppression is applied



VJ: Viola & Jones CVPR 2001 ScSPM : Yang et al. CVPR 2009 SPMK: Lazebnik et al. CVPR 2006 partModel: Felzenszwalb et al CVPR 2008

Results: Panda detection



SVCL ₹UCSD

Results: comparison on Caltech101 and Natural scene 15 classification dataset

Method	# Layer 2 units		Recognition rate		# Layer 2 units		Recognition rate	
	N15	C101	N15	C101	N15	C101	N15	C101
V1 model[1]	-	4000	-	42 ± 0.5	-	86K	-	65
HMAX[2]	-	4075	-	56	-	-	-	-
convNN[3]	-	4096	-	65.5	-	-	-	-
BoF[4]	-	4200	-	64.6	8400	-	81.4 ± 0.5	-
ScSPM[5]	5120	5120	75.3 ± 0.5	64.8±0.7	21504	21504	80.28 ± 0.9	73.2 ± 0.5
NBNN[6]	-	-	-	-	-	10M	-	70.4
HGMM[7]	-	-	-	-	46080	310272	85.2	73.1
DiscSalNet	600	4040	82 ± 0.5	70 ± 0.5	22500	20200	85.4±0.3	73.1±0.6

[1] Pinto et al. PLoS 2008 [2] Mutch & Lowe CVPR 2006 [3] Jarrett et al. ICCV 2009
[4] Lazebnik et al. CVPR 2006 [5] Yang et al. CVPR 2009 [6] Boiman et al. CVPR 2008
[7] Zhou et al. ICCV 2009

- Performances are **highly dependent** to the # of L2 units
- With similar L2 unit number, Saliency Net achieves the best performance

Conclusion

- top-down discriminant saliency model for amorphous object detection
 - Cell level optimal tuning for given problem
 - Complex feature in higher level (increase selectivity)
 - Absence of feature can also be exploited
 - Achieves best performance among state-of-the-arts
 computer vision algorithms on standard test benchmarks
- Edge based (SIFT, HOG, etc) feature may not be enough
 - State-of-the-arts method could be tuned too much to the specific dataset

Thank you!

