CVPR 2013 TUTORIAL: A CRASH COURSE ON VISUAL SALIENCY MODELING: BEHAVIORAL FINDINGS AND COMPUTATIONAL MODELS

SALIENCY: INFORMATION THEORETIC, BAYESIAN, AND GRAPHICAL MODELS

Neil Bruce

Department of Computer University of Manitoba Winnipeg, MB, Canada **Contact Information:**

bruce@cs.umanitoba.ca www.cs.umanitoba.ca/~bruce



University of Manitoba



Overview

- Motivating ideas: A summary
- Information theoretic approaches
 - Some of my work (AIM)
- Coding based approaches
- Graph-based strategies
- Bayesian approaches
- Summary and Discussion

Some basic background

Entropy

$$H(X) = \mathbb{E}_X[I(x)] = -\sum_{x \in \mathbb{X}} p(x) \log p(x).$$

Joint entropy

$$H(X,Y) = \mathbb{E}_{X,Y}[-\log p(x,y)] = -\sum p(x,y)\log p(x,y)$$

Conditional entropy

$$H(X|Y) = \mathbb{E}_{Y}[H(X|y)] = -\sum_{y \in Y} p(y) \sum_{x \in X} p(x|y) \log p(x|y) = -\sum_{x,y} p(x,y) \log \frac{p(x,y)}{p(y)}$$

Mutual information

) p(y)

 $_{x,y}$

$$I(X;Y) = \mathbb{E}_{X,Y}[SI(x,y)] = \sum_{x,y} p(x,y) \log \frac{p(x,y)}{p(x)p(y)}$$

Can all be expressed in terms of KL-divergence

Some basic background

Other important ideas

Suspicious Coincidences (Barlow)

- One goal of the brain is detecting associations
- Find suspicious coincidences, and anticipate them
- Coding theory
 - Rate/Distortion
 - Data compression
 - Redundancy
- Bayes' Theorem

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}.$$

Computer Vision Classics



 KL-divergence across scale space representation
 Peaks provide a sense of scale of interest

Computer Vision Classics

- Kadir and Brady (IJCV 2001)
- Again appealing to peaks in scale space subject to local entropy
- Some extra steps (region clustering, etc...)
- Area has been further explored in detail by interest point/descriptor research community



Saliency

















Information Theoretic Saliency

A note on "categorization"

- Information theory, coding, Bayesian inference, Graphical models aren't easily separated
- Grouped thematically, but several may be present within any single model



Source: Unknown

AIM (Attention by Information Maximization)

- Appeals to role of coding, and information theory
- Key points:
 - Independent (sparse) coding
 - Want to quantify likelihood of observing local patch/ region of image
- \Box Likelihood related to self-information via $-\log(p(x))$

Computational Constraints

$$\square p(X = x_1, x_2, ..., x_n)$$

Problem of dimensionality: More general problem

X 1	x ₂	X 3	x ₄	x ₅
X 6	X ₇	X 8	X 9	X ₁₀
X ₁₁	X ₁₂	X ₁₃	x ₁₄	X ₁₅
X ₁₆	X ₁₇	X ₁₈	X ₁₉	X ₂₀
X ₂₁	X ₂₂	X ₂₃	x ₂₄	X ₂₅

- Seek a representation that makes the computation of p(X) tractable
- Content is not random but highly structured and exists in a lower dimensional manifold
- Solution: Take advantage of structure in the data -Sparse representation as one approach
- Information appeals to Shannon self-information

Computational Constraints

$$p(x_1, x_2, ..., x_n) = \prod_i p(x_i)?$$

□ Ideally $X_1, X_2, ..., X_n$ should be independent variables



The Model (AIM)



(Bruce and Tsotsos, NIPS 2005, JoV 2009)

Quantifying Performance

- There now exist a number of datasets, benchmarks, performance metrics, etc.
 - Benchmarking will be discussed later ③
- Different data sets, methodology and parameters
- Seems to hold up well in benchmarks, despite remaining largely untouched for 8 years...

Prediction of fixation patterns



Behavioral phenomena





В	こういいしゃくいいいいいいいいい		いろう ひょう ふらん		今日後ので ちのて ごりくの かの からで ちゅう たいてい かんてい かんてい			シバーション マンシー・ション シャンション ションション
	4	9	ц	6	9	٢	ę	4
	5	q	1	"	`	4	2	4
	6	۶	٩	5	>	١	7	2
	9	8	7	3	5	2	2	9



(Bruce and Tsotsos, 2009, Bruce and Tsotsos 2011)

Spatiotemporal Cells



Examples





Many facets of AIM

- Role of feature representation (performance and neuroscience)
- JoV 2009, CRV 2011, Front. Comp. Neurosci. 2011
- Role of context in defining algorithm output
- ICIP 2009, ICPR 2004
- Recurrence/hierarchical processing (CRV 2012, ICIAR 2012, CVPR Workshop Biol.
- Consistent Vision, 2011)

□ Hou and Zhang (NIPS 2008)

- Measure entropy gain of each feature
- Maximize entropy across sample features
- Select features with large coding length increment



Activity ratio p_i for ith feature:
$$p_i = \frac{\sum_k |\mathbf{w}_i \mathbf{x}^k|}{\sum_i \sum_k |\mathbf{w}_i \mathbf{x}^k|}$$

with
$$\mathbf{X} = [\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^k, \dots]$$

 Most efficient strategy is to make equal use of all features – i.e. maximize H(p)

New excitation of feature i:
$$\hat{p}_j = \begin{cases}
\frac{p_j + \varepsilon}{1 + \varepsilon}, & j = i \\
\frac{p_j}{1 + \varepsilon}, & j \neq i
\end{cases}$$

Changes entropy of feature activities:

 $\frac{\partial H(\mathbf{p})}{\partial p_i} = -\frac{\partial p_i \log p_i}{\partial p_i} - \frac{\partial \sum_{j \neq i} p_j \log p_j}{\partial p_i}$

Which boils down to:

$$ICL(p_i) = \frac{\partial H(\mathbf{p})}{\partial p_i} = -H(\mathbf{p}) - p_i - \log p_i - p_i \log p_i$$

Salient feature set:

 $S = \{i \mid ICL(p_i) > 0\}$ i.e. Do subsequent activations of feature i increase entropy of the system?

Allows redistribution of energy among salient features:

$$d_i = \frac{\text{ICL}(p_i)}{\sum_{j \in \mathcal{S}} \text{ICL}(p_j)}$$



Finally, salience may be computed, with $M = [m_1, m_2, ..., m_n]$

$$\mathbf{m}_k = \sum_{i \in \mathcal{S}} d_i \mathbf{w}_i \mathbf{x}^k$$





A: input sample







Spatiotemporal data

Dynamic Visual Attention

- 1. At time t, calculate feature ICL based on p^t
- 2. Given current eye fixation, generate a saliency map with foveal bias.
- 3. By a saccade, move eye to the global maximum of the saliency map.
- 4. Sample top N "informative" (largest ICL) features in fixation neighborhood. (In our experiment, N = 10)
- 5. Calculate $\hat{\mathbf{p}}^t$, update \mathbf{p}^{t+1} , and go to Step. 1.



Conditional Entropy

- □ Li, Zhou, Yan and Yang (ACCV 2009)
 - Saliency based on conditional entropy
 - Minimum uncertainty of local region given surround
 - Conditional entropy given by coding length (assuming lossy distortion) modeled as multivariate Gaussian data
 - Segmentation to detect proto-objects
 - Extended by Yan et al. to multi-resolution

Conditional Entropy



Probability/Clutter

- Rosenholtz (Vis. Res. 1999, JoV 2007, ACM TAP 2011)
 - Distribution of features determined (e.g. in color-space)
 - Mean and covariance of distractors computed
 - Target saliency given by Mahalanobis distance given target, and mean/covariance of distractor distribution
 - Later versions also account for role of "clutter"

Rarity Based Saliency

Mancas (2007)

Considers rarity of features (both local and global, including subject to self-information)

Multi-scale approach reminiscent of Itti et al.

Also consider many applications

Mancas (2012) (RARE)

Normalization/Whitening across color inputs and across scale, weighted combination/fusion

Attention
$$(I_j) = -\log\left(\frac{1}{S * |I_j|} \sum_{i=1}^{S} n_i\right)$$

RARE 2012



Self-resemblance

Seo and Milanfar (Journal of Vision, 2009)

- Local structure represented by matrix of local descriptors (steering kernels robust to noise/image distortions)
- Matrix cosine similarity forms a metric for resemblance at pixel to surround
- Amounts to an estimate of likelihood of local feature matrix given feature matrix of pixels in surround

$$\begin{split} K(\mathbf{x}_l - \mathbf{x}_i) = & \frac{\sqrt{\det(\mathbf{C}_l)}}{h^2} \exp\left\{\frac{(\mathbf{x}_l - \mathbf{x}_i)^T \mathbf{C}_l(\mathbf{x}_l - \mathbf{x}_i)}{-2h^2}\right\}, \quad \mathbf{C}_l \in \mathbb{R}^{2 \times 2} \\ S_i = & \frac{1}{\sum_{j=1}^N \exp\left(\frac{-1+\rho(\mathbf{F}_i, \mathbf{F}_j)}{\sigma^2}\right)} \end{split}$$
Self-resemblance



Site-Entropy Rate

Wang et al. CVPR 2011 (following Wang et al. CVPR 2010)



Site-Entropy Rate

- Wang et al. CVPR 2011 (following Wang et al. CVPR 2010)
- Average total information transmitted from location
 I to other nearby locations

$$S_i = \sum_k SER_{ki} = -\sum_k \left(\pi_{ki} \sum_j P_{kij} \log P_{kij}\right)$$

 π_i - Stationary distribution term (frequency with which random walker visits node i / frequency with which node i communicates with other nodes)

Site-Entropy Rate



Random walks: See also Achanta et al. 2009

Information Gain

Najemnik and Geisler



□ The "ideal observer"

- Subject to simulated constraints/uncertainty on perception
- Wish to maximize the information gain, or minimize uncertainty with respect to defined target location in making a saccade

Information Gain

Najemnik and Geisler (Nature 2005)





 $p_{i} = \frac{prior(i)\exp(d_{i}^{\prime 2}W_{i})}{\sum_{j=1}^{n} prior(j)\exp(d_{j}^{\prime 2}W_{j})}$

Information Gain

Butko and Movellan, ICDL 2008, IEEE TAMD 2010







Discriminant / Decision Theoretic Saliency

Spatial definition for "c"

$$S(l) = I_l(\mathbf{X}; Y)$$

= $\sum_c \int p_{\mathbf{X}(l), Y(l)}(\mathbf{x}, c) \log \frac{p_{\mathbf{X}(l), Y(l)}(\mathbf{x}, c)}{p_{\mathbf{X}(l)}(\mathbf{x}) p_{Y(l)}(c)} d\mathbf{x}.$



Decision Theoretic Saliency

Diagrams images

$$P_X(x;\alpha,\beta) = \frac{\beta}{2\alpha\Gamma(1/\beta)} \exp\left\{-\left(\frac{|x|}{\alpha}\right)^{\beta}\right\}$$

$$I(X;Y) = \sum_{c} P_Y(c) KL[P_{X|Y}(x|c) || P_X(x)],$$

$$\operatorname{KL}[P_X(x;\alpha_1,\beta_1) || P_X(x;\alpha_2,\beta_2)] = \log\left(\frac{\beta_1 \alpha_2 \Gamma(1/\beta_2)}{\beta_2 \alpha_1 \Gamma(1/\beta_1)}\right) + \left(\frac{\alpha_1}{\alpha_2}\right)^{\beta_2} \frac{\Gamma((\beta_2+1)/\beta_1)}{\Gamma(1/\beta_1)} - \frac{1}{\beta_1}$$



Discriminant / Decision Theoretic Saliency

- Derived explicitly from a minimum Bayes error definition
- "c" applicable to centre/surround, but also other classes (e.g. face vs. null hypothesis)
- Specific mathematical relationship can be shown to:
 - Suspicious coincidences, decision theory, neural computation/complex cells/circuitry, tracking
- See: Han and Vasconcelos Vis. Res. 2010, Mahadevan and Vasconcelos, TPAMI 2010, Gao et al. IEEE TPAMI 2009, Gao and Vasconcelos Neur. Comp 2009, Gao, Mahadevan and Vasconcelos, 2007, Gao and Vasconcelos ICCV 2007, Gao, Mahadevan and Vasconcelos NIPS 2007

Suspicious coincidences

 See also: Choe and Sarma AAAI 2006 (On relation between orientation filter responses and natural image statistics)





□ Torralba, Oliva, Castelhano and Henderson, Psych. Rev. 2006



$$p(O=1, X|L, G)$$

$$=\frac{1}{p(L|G)}p(L|O=1, X, G)p(X|O=1, G)p(O=1|G)$$

- This builds on several prior efforts, followed by some additional targeted efforts:
 - Context/Contextual priors:
 - Hidalgo-Sotelo, Oliva and Torralba, CVPR 2005
 - Torralba, NIPS 2001
 - and others...
 - Top-down control:
 - Oliva, Torralba, Castelhano and Henderson, ICIP 2003
 - Ehinger, Hidalgo-Sotelo, Torralba, Oliva, 2009
 - Oliva and Torralba, TICS 2007

□ Zhang et al., J. of Vision, 2008

• SUN
$$s_z = p(C = 1 | F = f_z, L = l_z)$$

= $\frac{p(F = f_z, L = l_z | C = 1)p(C = 1)}{p(F = f_z, L = l_z)}$

$$\log s_z = \underbrace{-\log p(F = f_z)}_{\text{Self-information:}} + \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log likelihood:}} + \underbrace{\log p(C = 1 \mid L = l_z)}_{\text{Log ation prior:}} \\ \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log likelihood:}} + \underbrace{\log p(C = 1 \mid L = l_z)}_{\text{Location prior:}} \\ \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log likelihood:}} + \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log ation prior:}} \\ \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log likelihood:}} + \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log ation prior:}} \\ \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log likelihood:}} + \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log ation prior:}} \\ \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log likelihood:}} + \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log ation prior:}} \\ \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log likelihood:}} + \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log ation prior:}} \\ \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log ation prior:}} + \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log ation prior:}} \\ \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log ation prior:}} + \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log ation prior:}} + \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log ation prior:}} \\ \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log ation prior:}} + \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log ation prior:}} + \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log ation prior:}} + \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log ation prior:}} + \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log ation prior:}} + \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log ation prior:}} + \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log ation prior:}} + \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log ation prior:}} + \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log ation prior:}} + \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log ation prior:}} + \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log ation prior:}} + \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log ation prior:}} + \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log ation prior:}} + \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log ation prior:}} + \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log ation prior:}} + \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log ation prior:}} + \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log ation prior:}} + \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log ation prior:}} + \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log ation prior:}} + \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log ation prior:}} + \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log ation prior:}} + \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log ation prior:}} + \underbrace{\log p(F = f_z \mid C = 1)}_{\text{Log ati$$

- □ Zhang et al. Proc. Cog. Sci. Soc. 2009,
 - SUNDAy, Dynamic analysis of scenes
- Kanan et al. 2009, Visual Cognition
 - Top down saliency
- □ Barrington et al. J. of Vision, 2008
 - NIMBLE: Saccade based visual memory
- Static model of natural image statistics, modeled as GGD lends itself to a very fast computational framework

Itti and Baldi, NIPS 2006, Vis. Res. 2009, Neural Netw, 2010





□ Family *M* of observer-dependent models or hypotheses about the world
 □ Observer beliefs: {*P*(*M*)}_{*M*∈*M*}

Bayesian foundation of probability: data is what changes a prior into a posterior:



$$P(M \mid D) = \frac{P(D \mid M)}{P(D)} P(M)$$





Surprise =
$$d[P(M | D), P(M)]_{M \in \mathcal{M}}$$



After a moment...

- Beliefs stabilize, prior and posterior become identical,
- and additional snow frames carry no surprise.

MTV CNN FOX BBC

 $P(\mathcal{M}),$

 $P(\mathcal{M}|D)$



Unit of surprise

□ Shannon: $I(D) = -\log P(D)$

1 bit observed when outcome has probability 0.5
E.g., toss a fair coin and observe "heads"

□ Surprise: $S(M) = \log \frac{P(M \mid D)}{P(M)}$

I wow experienced when belief in M changes by factor 2
E.g., from 20% belief that a coin is fair to 40%

Computing surprise

Discrete Case (Binomial/Dirichlet)

$$S(D, \mathcal{M}) \approx N \times KL\left[p, \frac{a_1}{a_1 + b_1}\right]$$

p =observed frequency over N binary data samples
a1 /(a1+b1) = expected frequency (Dirichlet)

$$S(D,\mathcal{M}) = \frac{N}{2} \left[\frac{1}{\kappa_1} + \frac{\overline{\sigma}^2}{s_1^2} + \log \frac{\upsilon_1 s_1^2}{2\overline{\sigma}^2} - \Psi\left(\frac{\upsilon_1}{2}\right) + \frac{(\overline{m} - \mu_1)^2}{s_1^2} \right]$$

Continuous Case (Gaussian, unknown mean and variance)











Itti and Baldi, NIPS 2006, Vis. Res. 2009, Neural Netw, 2010





Harel 2006

- Scale-space pyramid from intensity, color, orientation
- Fully connected graph over all grid locations
- Graph weights proportional to similarity of feature values, and spatial distance

$$d((i,j)||(p,q)) \triangleq \left|\log \frac{M(i,j)}{M(p,q)}\right|$$

$$w_1((i,j),(p,q)) \triangleq d((i,j)||(p,q)) \cdot F(i-p,j-q), \text{ where}$$

$$F(a,b) \triangleq \exp\left(-\frac{a^2+b^2}{2\sigma^2}\right).$$

□ Harel, NIPS 2006

- Treated as Markov chain that reflects expected time spent by a random walker (walking forever)
- Weights of outbound edges normalized to 1 with equivalence relation defined between nodes/states and edges/transition probabilities
- Saliency corresponds to equilibrium distribution

(a) Sample Picture With Fixation



Pang ICME 2008

- Stochastic model based on signal detection theory
- Dynamic Bayes net with 4 layers
 - Layer 1: Itti-like saliency determination
 - Layer 2: Gaussian state-space model (stochastic saliency map)
 - Layer 3: Overt shifts determined by HMM
 - Layer 4: Density map predicts positions
Graph Based techniques

Avraham and Lindenbaum (PAMI 2010)

The Esaliency Algorithm

- 1) Select candidates using some segmentation process.
- 2) Use the preference for a small number of expected targets (and possibly other preferences) to set the initial (prior) probability for each candidate to be a target.
- 3) Measure visual similarity between every two candidates and infer the correlations between the corresponding labels.
- 4) Represent the label dependencies using a Bayesian network.
- 5) Find the N most likely joint assignments.
- 6) Deduce the saliency of each candidate by marginalization.

$$I(l_i, l_j) = \sum_{l_i=0,1} \sum_{l_j=0,1} p(l_i, l_j) \log \frac{p(l_i, l_j)}{p(l_i)p(l_j)}$$

Labels are binary random variables:

$$p(l_i = 1, l_p = 1) = \gamma(d_{ij})\sqrt{\mu_i(1 - \mu_i)\mu_j(1 - \mu_j) + \mu_i\mu_j}$$

$$p(l_i = 1, l_j = 0) = \mu_i - p(l_i = 1, l_j = 1)$$

$$p(l_i = 0, l_j = 1) = \mu_j - p(l_i = 1, l_j = 1)$$

$$p(l_i = 0, l_j = 0) = 1 - \mu_i - \mu_j + p(l_i = 1, l_j = 1).$$

E-Saliency

Dependency on parent nodes for label

$$p(\bar{l}) = p(l_r) \prod_{i=1,\dots,n; i \neq r} p(l_i | l_{\text{par}(i)})$$

Marginalization considering most likely assignments:

$$p'(\bar{l}) = \frac{p(\bar{l})}{\sum_{j=1}^{N} p(\bar{l}^j)}$$

The saliencies are then:

$$p_T(c_i) = \sum_{j=1}^N p'(\bar{l^j}) \cdot l_i^j$$

E-saliency





Probabilistic and Bayesian models

Rao, NeuroReport 2005



- Bayesian, Integrate and Fire model
- Heavily inspired by biology, brain imaging
- See also Rao and Ballard, Nat. Neurosci. 1999

Probabilistic and Bayesian models

□ Chikkerur et al., Vis. Res. 2010, MIT Ph.D. Thesis



Graph Based techniques

Strongly inspired by biology



Learning/Object detection Methods

See also:

- What is an object? (Alexe et al. 2010)
- Deselaers et al. (ECCV 2010)
- Carreira and Sminchisescu (CVPR 2010)
- Gu et al. (CVPR 2009)
- van de Sande et al. (ICCV 2011)
- and many more...

Some take home points

Take home points...

- Much overlap in fundamental ideas that inspire techniques in this domain
- This isn't surprising (these are all fundamental principles in many efforts not just saliency)
- Reveals that the details are important (think of the tree)
- There are several benchmarks (which are important), but consider also application
- Saliency is very useful but won't solve everything