Pictorial Structures for Object Recognition

Pedro F. Felzenszwalb Presented by Hanlin Tang COS/PSY 598b

Feature detection



╢



Felzenswab Talk, 2007

Pictorial Recognition





Felzenswab Talk, 2007

1

The Task











The Task

 $\Theta = \{u, c\}$ I $L = \{I_1, I_2, I_3...\}$

 $L^* = \operatorname*{argmax}_L p(L|I, \theta)$

Bayes Rule to the Rescue...

$p(L|I,\theta) \propto p(I|L,\theta)p(L|\theta)$

Assuming part independence,

$$P(L|I,\theta) \propto \left(\prod_{i=1}^{n} p(I|l_i, u_i) \prod_{(v_i, v_j) \in E} p(l_i, l_j|c_{ij})\right)$$

A more intuitive formulation:

$$L^* = \operatorname*{argmax}_{L} p(L|I, \theta)$$
$$= \operatorname*{argmin}_{L} - \log[p(L|I, \theta)]$$

Fischler and Eschlager (1972)

The Grand Assumption

Want to define $d_{ij}(l_i, l_j)$

Let:
$$x = l_1 - l_2$$

In 1dimension:

$$\left(\frac{x-\mu}{\sigma}\right)^2$$

In N-dimensions: $(x - \mu)^T \Sigma^{-1} (x - \mu)$



Iconic Models - Faces



 $\Theta = \{u, c\}$ $L = \{I_1, I_2, I_3...\}$

 $L^* = \operatorname*{argmax}_L p(L|I, \theta)$

What an eye should look like

• The naïve way:



$$u = \{\mu, \sigma\}$$

$$\mu = \begin{pmatrix} 0.2 & 0.4 & 0.4 & 0.4 & 0.2 \\ 0.1 & 0.5 & 0.6 & 0.5 & 0.1 \\ 0.1 & 0.1 & 0.9 & 0.1 & 0.1 \\ 0.1 & 0.3 & 0.8 & 0.4 & 0.1 \\ 0.2 & 0.3 & 0.3 & 0.3 & 0.2 \end{pmatrix}$$

But can reduce dimensions!

• The intuition:



More intuition-building:



COS424 notes

Applied to natural data





 $\alpha(l_i)$

 $u_i = (\mu_i, \Sigma_i)$

$p(I | l_i, u_i) \propto N(\alpha(l_i), u_i, \Sigma_i)$

Characterizing springs



 $p(l_i, l_i | c_{ii}) = N(l_i - l_i, s_{ii}, \Sigma_{ii})$

$$L^* = \arg\min_L \left(\sum_{i=1}^n m_i(l_i) + \sum_{(v_i,v_j) \in E} d_{ij}(l_i,l_j) \right)$$

Seek models to define:

- Location
- Appearance
- Connections
- $p(I|l_i, u_i)$
- $p(l_i, l_j | c_{ij})$

$$l_{i} = (x, y)$$

$$u_{i} = (\mu_{i}, \Sigma_{i})$$

$$c_{ij} = (S_{ij}, \Sigma_{ij})$$

$$p(I | l_{i}, u_{i}) \propto N(\alpha(l_{i}), u_{i}, \Sigma_{i})$$

$$p(l_{i}, l_{i} | c_{ij}) = N(l_{i} - l_{j}, s_{ij}, \Sigma_{ij})$$

Some Math: $p(l_i, l_i | c_{ii}) = N(l_i - l_i, s_{ii}, \Sigma_{ii})$

$p(l_i, l_j | c_{ij}) = \mathcal{N}(T_{ij}(l_i) - T_{ji}(l_j), 0, D_{ij})$

















Articulated Models - Humans



$$L^* = \arg\min_L \left(\sum_{i=1}^n m_i(l_i) + \sum_{(v_i,v_j) \in E} d_{ij}(l_i,l_j) \right)$$

Seek models to define:

- Location
- Appearance
- Connections
- $p(I|l_i, u_i)$
- $p(l_i, l_j | c_{ij})$

 $l_{i} = (x, y, s, \theta)$ $u_{i} = (q_{1}, q_{2})$ $c_{ij} = (x_{ij}, y_{ij}, x_{ji}, y_{ji}, \sigma_{x}, \sigma_{y}, \sigma_{s}, \theta_{ij}, k)$ $p(I | l_{i}, u_{i}) = q_{1}^{n_{1}} (1 - q_{1})^{n_{1}} q_{2}^{n_{2}} (1 - q_{2})^{n_{2}} (0.5)^{T - A_{1} - A_{2}}$ $p(l_{i}, l_{i} | c_{ij}) \propto N(T_{ji}(l_{i}) - T_{ij}(l_{j}), 0, D_{ij})$





Figure 13: Matching results (sampling 200 times).





Figure 14: In this case, the binary image doesn't provide enough information to estimate the position of one arm.



Figure 15: This example illustrates how our method works well with noisy images.









- Restriction on d_{ij} allows linear running time for Finding L*
- Efficient ways to sample from the Posterior $p(L|I, \theta)$

Learning







 $I = \{I^1, I^2, ...\}$

 $L = \{L^1, L^2, ...\}$ $\Theta = \{u, c\}$

Learning the Model



 $\theta^* = \arg \max_{\theta} \prod_{k=1}^m p(I^k | L^k, \theta) \prod_{k=1}^m p(L^k | \theta).$

Learning Appearance and Connections



$$c_{ij}^{*} = \arg \max_{c_{ij}} \prod_{k=1}^{m} p(l_{i}^{k}, l_{j}^{k} | c_{ij}).$$

$$q(v_i, v_j) = \prod_{k=1}^{m} p(l_i^k, l_j^k | c_{ij}^*).$$

$$E^* = \arg\max_E \prod_{(v_i, v_j) \in E} q(v_i, v_j) = \arg\min_E \sum_{(v_i, v_j) \in E} -\log q(v_i, v_j)$$

General Framework



Felzenswab Talk, 2007

Key Points

- Pictorial models bring context to recognition
- Robust to noise, scale, lighting effects
- Generalized structure
- Heavily dependent on prior training of model
- Require robust definitions of (u,v)

1000	12345-709012345-799012345-739012345-7990
1.	
2	
3	
4	医囊骨骨骨骨骨骨骨骨骨骨骨骨骨骨骨骨骨骨骨骨骨骨骨骨骨骨骨骨骨骨骨骨骨骨骨骨
5	***************************************
	343449999844088666888888433384446646K86
7	弟后后于老老老老老老爸的这些的说明我们的那些的问题。 第二十十十十十十十十十十十十十十十十十十十十十十十十十十十十十十十十十十十十
	基本专业的专业主要的现在分词有限的保持的资源的保持的公司并不可能行用。
9	EV=+************************************
10	####################################
11	EEEF18248888888888888888888888888897684423
12	◆● ++ 予 + + は自身自己自然的的自身を必ず上的自己自然的子人" -> A + 5 -
13	医单子子氏单位的复数的复数形式 医丁丁乙基苯基甲基苯基乙基 化化化化合金
14	\$\$\$fff(B686668884)+===+)278888994348429
1.5	E也有此他的问题中国的方面。 # +1 A目目目目的方式中的分析中
16	##3#?#VV#3VF#X1+ -+1A###9AAA4999

18	<pre>#####EG####HXZ1+ -1AB#XXXAAMBBBA</pre>
19	EEBERMY####################################
27	4.6944.4.9.74.82.8.018.5.2.3.64.8.64.6.3.4.5.1.3.3.7.7.4.5.4.7
. 21	BAN ATTENT BULLING AN ALBERT TE SEE SEE TE ATAUNT
22	
	BEETARRAWSIIIIIII = + + + + + + + + + + + + + + +
22	######################################
6 V 24	00000000000000000000000000000000000000
	NWAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA
28	************
29	
	XXX777771X87X66474/174889X1111117777
31	X X X 7 7 7 7 7 7 MAGA 3 7 X 7 1 1 7 ABGAAGAAX 1 7 7 7 7 7 7 7 7 7 X X X
32	***************************************
-33-	XXX7230E-FEGEGXXXAGEVTA076EGENVHENAX
34	XXXX本集集合 X目前目前最终的目前在XX本型1集集中建立中国中国目
35	XXABBER JEMANMAXXXXXABJABREBERXEBHMM

HAIR WAS LOCATED AT (11, 21) L/EDGE WAS LOCATED AT (25, 11) R/EDGE WAS LOCATED AT (25, 24) L/EYE WAS LOCATED AT (21, 15) R/EYE WAS LOCATED AT (21, 21) NOSE WAS LOCATED AT (26, 18) MOUTH WAS LOCATED AT (29, 17)

Fischler and Eschlager (1972)

123456789012345678901234567890123456789

Original picture.