Human Motion Categorization & Detection

Juan Carlos Niebles^{1,2} Advisor : Prof. Fei-Fei Li¹





Juan Carlos Niebles^{1,2}





BARRANQUILLA - COLOMBIA



capture

game interfaces

In this talk



- Where are the humans?
 - Which pixels?
- What motions are they doing?
 - Learn models and classify novel sequences

In this talk





- Where are the humans?
 - Which pixels?
- What motions are they doing?
 - Learn models and classify novel sequences

Identification of human actions in video



Challenges:

- Camera Motion
- Complex Background
- Viewpoint Change

Localizing human actions in long sequences



Multiple actions by different people in one single sequence

Camera Motion



Dynamic Background

Previous Work

- 1. Recognition by template correlation/matching
 - Efros et al, ICCV '03
 - Shechtman et al, CVPR '05
 - Ramanan et al, NIPS '03

2. Action analysis using graphical models

- Song et al, PAMI '03
- Fanti et al, ICCV '05
- Boiman et al, ICCV '05
- Bergler, CVPR '97

3. Spatial-temporal Interesting Points

- Laptev et al, ICCV '03; Schuldt et al, ICPR '04
- Dollar et al, ICCV VS-PETS '05
- Ke et al, ICCV '05

Motivation



Sparse and *Local* representation

Spatio-temporal information

Johansson, 1973

Spatio-Temporal Interest Points

Detection

$$\mathbf{R} = (\mathbf{I} * g_x * g_y * h_{ev})^2 + (\mathbf{I} * g_x * g_y * h_{od})^2$$

Local optimum of R define the position of features











[Dollar et al '05]

Spatial-Temporal Interest Points

- walking
- running
- jogging
- boxing
- handclapping
- handwaving



Spatial-Temporal Interest Points

Detection

$$\mathbf{R} = (\mathbf{I} * g_x * g_y * h_{ev})^2 + (\mathbf{I} * g_x * g_y * h_{od})^2$$

Local optimum of R define the position of features



Description

Spatial-temporal cube



Experiments showed that optical flow descriptor is equally effective

[Dollar et al '05]





Codebook and Representation



Codebook and Representation



Experiment I: Codebook





Unsupervised learning using pLSA



pLSA Model



$$p(w_i \mid d_j) = \sum_{k=1}^{K} p(w_i \mid z_k) p(z_k \mid d_j)$$

action category vectors

Word distribution per action category

action category weights

Action category distribution per video

pLSA Model



$$p(w_i | d_j) = \sum_{k=1}^{n} p(w_i | z_k) p(z_k | d_j)$$

Unsupervised Learning

$$L = \prod_{i=1}^{M} \prod_{j=1}^{N} p(w_i \mid d_j)^{n(w_i, d_j)}$$



Recognition



Experiment I:

KTH dataset [Schuldt et al., 2004]:



25 persons, indoors and outdoors, 4 long sequences per person

Experiment I: Performance

- Leave-one person out cross validation
- Average performance: 81.50%
- Unsupervised training
- Handle multiple motions





Experiment I: Caltech dataset



Only words from the corresponding action are shown

Experiment I: A longer sequence





Trained with the KTH data

Tested with our own data

Experiment I: Multiple motions







handwaving

Trained with the KTH data

Tested with our own data

Experiment II:

Figure Skating data set: [Y.Wang, G.Mori et al, CVPR 2006]



7 persons, 3 action classes: camel spin, stand spin, sit spin

Experiment II: Examples

Figure skating actions







Camel spin Sit spin Stand spin

Experiment II: Long Sequences









references

• Juan Carlos Niebles, Hongcheng Wang and Li Fei-Fei. Unsupervised Learning of Human Action Categories Using Spatial-Temporal Words. **IJCV**. September. 2008.



- with
 - Hongcheng Wang

but...



• We need models that exploit geometrical arrangements of features

in the object recognition world



- all have same P under bag-of-words model.
- part-based models that capture geometrical information
 - Constellation model
 - Pictorial structures
 - etc

Constellation Model



- Small number of features
- Strong shape representation

bags of features



Large number of features

No geometrical or shape information

Constellation of bags of features



© Large number of features © Strong shape representation
- P_p : Parts
- w : Observed patches (shape, motion and position)
- P's are *similar* to the parts in Constellation model
- Each Part is the parent of a "bag-of-features"

- Large number of features
- Strong shape representation











learned models











- Static features
 Edge map + Shape Context
 [Belongie '03]
- Motion features
 Interest Points + ST-gradients
 [Dollar '05]







p

- Y set of possible locations for the P parts (a discretized grid over the image)
- h index variable to select a specific position of the
- P parts among all the possibilities (in the grid)
- m chooses an specific assignment of features w to parts P



$$p(\mathbf{w}, \mathbf{Y}|\boldsymbol{\theta}) \approx \sum_{\omega=1}^{\Omega} \left[\pi_{\omega} \sum_{\mathbf{h} \in H} p(\mathbf{h}|\boldsymbol{\theta}_{\omega}) p(\mathbf{Y}|\mathbf{h}, \boldsymbol{\theta}_{\omega}) p(\mathbf{w}|\mathbf{Y}, \mathbf{m}^{*}, \mathbf{h}, \boldsymbol{\theta}_{\omega}) \right]_{Part \ layer} \sum_{Local \ feature \ layer} \left[\frac{1}{2} \sum_{\boldsymbol{h} \in H} p(\mathbf{h}|\boldsymbol{\theta}_{\omega}) p(\mathbf{Y}|\mathbf{h}, \boldsymbol{\theta}_{\omega}) p(\mathbf{w}|\mathbf{Y}, \mathbf{m}^{*}, \mathbf{h}, \boldsymbol{\theta}_{\omega}) \right]_{Local \ feature \ layer} \left[\frac{1}{2} \sum_{\boldsymbol{h} \in H} p(\mathbf{h}|\boldsymbol{\theta}_{\omega}) p(\mathbf{y}|\mathbf{h}, \boldsymbol{\theta}_{\omega}) p(\mathbf{w}|\mathbf{Y}, \mathbf{w}^{*}, \mathbf{h}, \boldsymbol{\theta}_{\omega}) \right]_{Local \ feature \ layer} \left[\frac{1}{2} \sum_{\boldsymbol{h} \in H} p(\mathbf{h}|\boldsymbol{\theta}_{\omega}) p(\mathbf{y}|\mathbf{h}, \boldsymbol{\theta}_{\omega}) p(\mathbf{w}|\mathbf{y}, \mathbf{w}^{*}, \mathbf{h}, \boldsymbol{\theta}_{\omega}) p(\mathbf{w}|\mathbf{y}, \mathbf{w}^{*}, \mathbf{h}, \boldsymbol{\theta}_{\omega}) \right]_{Local \ feature \ layer} \left[\frac{1}{2} \sum_{\boldsymbol{h} \in H} p(\mathbf{h}|\boldsymbol{\theta}_{\omega}) p(\mathbf{w}|\mathbf{y}, \mathbf{w}^{*}, \mathbf{h}, \boldsymbol{\theta}_{\omega}) p(\mathbf{w}|\mathbf{w}|\mathbf{y}, \mathbf{w}^{*}, \mathbf{h}, \boldsymbol{\theta}_{\omega}) p(\mathbf{w}|\mathbf{w}|\mathbf{w}^{*}, \mathbf{w}^{*}, \mathbf{w}, \mathbf{w}^{*}, \mathbf{w}, \mathbf{w}^{*}) p(\mathbf{w}|\mathbf{w}|\mathbf{w}^{*}, \mathbf{w}^{*}, \mathbf{w}, \mathbf{w}^{*}) p(\mathbf{w}|\mathbf{w}|\mathbf{w}^{*}, \mathbf{w}^{*}) p(\mathbf{w}|\mathbf{w}|\mathbf{w}^{*}, \mathbf{w}^{*}) p(\mathbf{w}|\mathbf{w}|\mathbf{w}^{*}, \mathbf{w}^{*}) p(\mathbf{w}|\mathbf{w}^{*}) p(\mathbf{w}|\mathbf$$

W	shape, motion and location features
Y	set of possible locations for the P parts (a discretized grid over the image)
h P	index variable to select a specific position of the parts among all the possibilities (in the grid)
m	chooses an specific assignment of features w to parts P



Part layer term: (constellation)

$$p(\mathbf{Y} | \mathbf{h}, \boldsymbol{\theta}_{\omega}) = N(\mathbf{Y}_{\mathbf{T}}(\mathbf{h}) | \boldsymbol{\mu}_{L, \omega}, \boldsymbol{\Sigma}_{L, \omega})$$

W	shape, motion and location features
Y	set of possible locations for the P parts (a discretized grid over the image)
h P	index variable to select a specific position of the parts among all the possibilities (in the grid)
m	chooses an specific assignment of features w to parts P



Local feature layer term:

$$p(\mathbf{w}|\mathbf{Y},\mathbf{m}^{*},\mathbf{h},\boldsymbol{\theta}_{\omega}) = \prod_{\mathbf{w}_{j}\in B_{g}} \underbrace{p(x_{j}^{r}|\boldsymbol{\theta}_{0}^{X})}_{B_{g} Shape} \underbrace{p(a_{j}|\boldsymbol{\theta}_{0}^{A})}_{B_{g} Appearance} \prod_{p=1}^{P} \prod_{\mathbf{w}_{i}\in P_{p}} \underbrace{p(x_{i}^{r}|\mathbf{Y},h_{p},\boldsymbol{\theta}_{p}^{X})}_{Part Shape} \underbrace{p(a_{i}|\boldsymbol{\theta}_{p}^{A})}_{Part Appearance}$$

learning



Parameter learning with $\mathbf{E}_{\omega}^{\mathbf{M}} = \left\{ \boldsymbol{\mu}_{L,\omega}, \boldsymbol{\Sigma}_{L,\omega}, \boldsymbol{\Sigma}_{p,\omega}^{X}, \boldsymbol{\theta}_{p,\omega}^{A}, \boldsymbol{\theta}_{0}^{X}, \boldsymbol{\theta}_{0}^{A} \right\} \quad p = 1 \dots P$ $\boldsymbol{\omega} = 1 \dots \Omega$

Dataset



"Actions as Space-Time Shapes" M. Blank, L. Gorelick, E. Shechtman, M. Irani, and R. Basri. IEEE International Conference on Computer Vision (ICCV), Beijing, October 2005.

experimental results

- 9 action classes, performed by 9 subjects [Blank et al 2005]
- Leave one out cross-validation
- Video Classification performance: 72.8%

bend	1.0	.00	.00	.00	.00	.00	.00	.00	.00
pjump	.00	1.0	.00	.00	.00	.00	.00	.00	.00
jack	.00	.00	1.0	.00	.00	.00	.00	.00	.00
wave1	.22	.11	.11	.44	.11	.00	.00	.00	.00
wave2	.00	.00	.11	.22	.67	.00	.00	.00	.00
jump	.00	.00	.00	.00	.00	.78	.00	.11	.11
run	.00	.00	.11	.00	.00	.11	.56	.11	.11
side	.00	.00	.00	.00	.00	.33	.11	.56	.00
walk	.00	.00	.00	.00	.00	.11	.00	.33	.56
	ben		jack	Wal	Wal ez	jum es	run	Side	Walk
					•				

Results



Results



references

 Juan Carlos Niebles and Li Fei-Fei. A Hierarchical Model of Shape and Appearance for Human Action Classification. CVPR, 2007



In this talk





- Where are the humans?
 - Which pixels?
- What motions are they doing?
 - Learn models and classify novel sequences





- detect moving people in YouTube videos
- extract the spatio-temporal volume that contains each person



Fast and robust algorithms for face detection

complexity of video sequences real-world sequences Compression artifacts & low quality •Videos contain multiple shots Unknown number of humans •Arbitrary human motion and Balan et al '07 poses •Unknown camera parameters and motion •Background clutter, motion and occlusions





extraction detail

System Overview



pedestrian detection

- Upright human detection is somewhat successful
 - Dalal & Triggs 05, Sabzmeydani & Mori 07, Laptev 06, and more



clustering ped detections



clustering ped detections



- Must-link: detections with similar appearance and consistent spatial locations
- Cannot-link: detections that occur within the same frame

Ø Clustering with constraints, Klein et al '02

First clustering step



- Clustering with constraints [Klein et al '02]
- Use simple & cheap descriptors: global color histogram
- Gives 'over-segmented' clusters Conservative clustering threshold

Second clustering step



- Clustering with constraints [Klein et al '02]
- More expensive descriptor: head & torso estimates
- Reject tracks
 Poor head torso estimates

Detection & clustering result



Output = One cluster per person in the sequence

System Overview



extracting the human region in a single frame [Ramanan 06]







Original image

Top-down person model

Estimated person region

- Pictorial structure representation
 - each part represents local visual properties
 - 'springs' capture spatial relationships
 - stretch & fit to find right configuration

images in this slide from: Deva Ramanan







Body part templates

- N-part model
- A configuration is given by $L = \{l_i\}; i = 1 \dots N$, with $l_i = \{x_i, y_i, \theta_i\}$
- We are interested in

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

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

Image evidence shape prior (measurement)

- efficient inference
- measurement is the bottle-neck expensive convolutions with nonseparable filters

reducing measurement computation

- Benefits
 - can be limited to a smaller space, reducing computation time
 - avoid distracting background observations
 - final estimation accuracy can be improved
- Use temporal consistency
 - Compute full model at first frame
 - Propagate part estimations to next frame







input image

original search space for torso (full frame)

reduced search space for torso

reducing measurement computation

	torso	upper-left arm			
	4	•			
input image	marginal Diffuse approxi	d GMM mation	marginal I a	Diffused GM	M on

- GMM approximation is done via Kernel Density Approximation (KDA) [Han et al. 'PAMI08]
- Densities are propagated in a Bayesian filtering framework

$$p(\mathbf{X}_{t} | \mathbf{Z}_{1:t}) \propto p(\mathbf{Z}_{t} | \mathbf{X}_{t}) p(\mathbf{X}_{t} | \mathbf{Z}_{1:t-1})$$
$$= \left(\sum_{i=1}^{N_{1}} \mathcal{N}(\kappa_{i}, \mathbf{X}_{i}, \mathbf{P}_{i})\right) \left(\sum_{j=1}^{N_{2}} \mathcal{N}(\tau_{j}, \mathbf{y}_{j}, \mathbf{Q}_{j})\right)$$

extracting the human region

Projecting the part posteriors into the image gives a rough segmentation of the body region


results



results



results



experiments

Precision Recall Comparison

	Detection only			Detection & Clustering			Full model		
	Prec	Rec	F	Prec	Rec	F	Prec	Rec	F
	0.89	0.31	0.46	0.89	0.30	0.45	0.83	0.73	0.78
	0.90	0.25	0.39	0.91	0.24	0.38	0.87	0.62	0.72
Rate	0.92	0.19	0.32	0.92	0.19	0.32	0.86	0.51	0.64
	0.93	0.16	0.27	0.94	0.15	0.27	0.92	0.43	0.58
	0.94	0.13	0.24	0.94	0.13	0.23	0.88	0.32	0.46

Computation Time

 \sim 1 order of magnitude speed up

more results









references

 Juan Carlos Niebles, Bohyung Han, Andras Ferencz and Li Fei-Fei. Extracting Moving People from Internet Videos. ECCV 2008



• With:



Bohyung Han



Andras Ferencz

Conclusions

- Spatio-temporal & spatial features + statistical classifier
 - Bag of words (unsupervised)
 - Constellation of bags of features
- Human Motion extraction
 - Real world sequences from YouTube

Interesting Research Issues

- Motion taxonomy/vocabulary
- higher level activities
- ...

In this talk



- Where are the humans?
 - Which pixels?
- What motions are they doing?
 - Learn models and classify novel sequences

Thank you