

1 In Brief	
Goal	
 Recover 3D human body pose - 3D pose = joint angles - use either individual images of 	e from monocular image silhouettes or video sequences
Applications	
 Human computer interaction Markerless motion capture 	 Gesture recognition Visual surveillance

Contributions

- "Model-free" learning based approach no explicit 3D model
- Mixture of kernel regressors trained using human motion capture data
- Multimodal probabilistic solutions in 3D pose space
- Temporal fusion using a particle filter style tracker

Silhouette Descriptors

Why Silhouettes

- Relatively simple and low-level
 - Capture most of the available pose information
 - Insensitive to surface attributes (clothing colour, texture..)
- Distortions caused by background subtraction, shadows • Ambiguity: hides internal details and depth ordering

Robust encoding of local shape — Shape Context Histograms



(a) extract silhouette



edge points







(c) find local (d) distribution (e) vector quantize shape contexts in SC space to get histogram

Training and Test Data 3

- Capture typical human movements, not just kinematically possible ones, using real human motion capture data
- Use both real silhouettes from motion capture and synthetic silhouettes from several human body models (POSER from Curious Labs)



Monocular Human Motion Capture with a Mixture of Regressors

Ankur Agarwal and Bill Triggs http://lear.inrialpes.fr

GRAVIR-INRIA-CNRS, Grenoble, France

Multimodal Pose Estimation

- The silhouette (z) to pose (x) problem is inherently multi-valued.
- Treating it as a function can lead to averaging or zig-zagging between different solutions.
- Introduce a discrete latent variable $k \in \{1, 2 \dots K\}$ to encode the information missing in the silhouette.
- Assume a mixture of experts model based on K underlying functional regression rules $\mathbf{x} \sim \mathbf{r}_k(\mathbf{z})$:

$$p(\mathbf{x} \mid \mathbf{z}) = \sum_{k=1}^{K} p(\mathbf{x} \mid \mathbf{z}, k) \, p(k \mid \mathbf{z}) , \qquad p(\mathbf{x} \mid \mathbf{z}, k) = \mathcal{N}(\mathbf{r}_k(\mathbf{z}), \mathbf{\Lambda}_k)$$

Mixture of Regressors by E-M 5

- Reduce dimensionality of silhouette data using kernel PCA: $z \rightarrow \phi(z)$
- Initialize clustering with local connected components analysis
- Fit a mixture of regressive Gaussians to the joint density ($\phi(\mathbf{z}), \mathbf{x}$):

$$\begin{pmatrix} \boldsymbol{\phi}(\mathbf{z}) \\ \mathbf{x} \end{pmatrix} \simeq \sum_{k=1}^{K} \pi_k \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Gamma}_k)$$

• Linear regressor within each component k

$$p(\mathbf{x} \mid \mathbf{z}) = \mathcal{N}(\mathbf{r}_k(\mathbf{z}), \mathbf{\Lambda}_k) \qquad \mathbf{r}_k(\mathbf{z}) \equiv \mathbf{A}_k \, \boldsymbol{\phi}(\mathbf{z}) + \mathbf{b}_k$$

• Special covariance structure enforces "regressive" noise model

$$\boldsymbol{\mu}_{k} = \begin{pmatrix} \boldsymbol{\phi}(\bar{\mathbf{z}}_{k}) \\ \mathbf{r}_{k}(\bar{\mathbf{z}}_{k}) \end{pmatrix}, \boldsymbol{\Gamma}_{k} = \begin{pmatrix} \boldsymbol{\Sigma}_{k} & \boldsymbol{\Sigma}_{k} \mathbf{A}_{k}^{\top} \\ \mathbf{A}_{k} \boldsymbol{\Sigma}_{k} & \mathbf{A}_{k} \boldsymbol{\Sigma}_{k} \mathbf{A}_{k}^{\top} + \boldsymbol{\Lambda}_{k} \end{pmatrix}$$



• M-step: Estimate A_k, b_k by weighted least squares regression, Λ_k from residual errors. Compute μ_k, Σ_k, π_k for each class.

• E-step: Reestimate class membership weights for each point.











Numbers of solutions and RMS joint angle reconstruction errors for 3 test sequences.



- Particle filter tracker, samples from dynamics $p(\mathbf{x}_t | \mathbf{x}_{t-1})$ as usual
- Uses regressive mixture $p(\mathbf{x}_t | \mathbf{z}_t)$ to assign posterior particle weights
- (Re)initializes by sampling from full mixture



RMS tracking error of individual joint angles on a 500 frame sequence (-1 when no person is detected).

{Ankur.Agarwal,Bill.Triggs}@inrialpes.fr

Pose from Static Images

 Provides multiple solutions for pose, with corresponding probabilities Most cases of ambiguity are identified

	% of frames with m solutions			Error in the	Error in best of
	m = 1	m = 2	$m \ge 3$	top solution	top 4 solutions
est person	62	28	10	6.14°	4.84 °
est motion	65	28	6	7.40 °	5.37°
ain subset	72	23	5	6.14°	4.55°

Self-Initializing 3D Tracking

$$p(\mathbf{x}_0 | \mathbf{z}_0) = \sum_{k=1}^{K} p(k | \mathbf{z}_0). \ \mathcal{N}(\mathbf{r}_k(\mathbf{z}_0), \mathbf{\Lambda}_k)$$

Potentially real time owing to closed form solution for posterior.

Automatic (re)initization

• Errors stabilize rapidly on (re)initialization





Detects the presence of a person and decides whether to wait, initialize or track using observed silhouette shape.

Upper Body Gesture Recognition 7.2



Test sequence labelling





Conclusion

- silhouettes

Work supported by an MENRT Doctoral Fellowship (French Education Ministry) and the European project LAVA







 Associate (by hand) different mixture components with gestures • Use posterior class probabilities to identify action

Training gestures (Basketball signals)

• "Model free" methods for recovering 3D human pose from monocular

 Multiple hypothesis pose estimates with associated probabilities • Stable pose recovery from static images and image sequences Action recognition using mixture components