Learning Transferable Distance Functions For Human Action Recognition and Detection

Weilong Yang Simon Fraser University



THINKING OF THE WORLD



Action Recognition and Detection











Applications

Action related video search

Sports and Dancing video search

Event Detection

 Automatic abnormality detection in surveillance videos





Motivation

- On KTH & Weizmann action datasets, almost 100% accurancy is achieved. [Jhuang et al. ICCV07, Fathi & Mori CVPR08]
- Most of methods rely on a large amout of training set.
 - Half-half split or Leave-one-out cross validation
- It is unrealistic to collect this many training samples for some action.





Related Works

- One shot learning of object categories [Fei-Fei et al.] ICCV031
- Visual Object Identification [Ferencz et al. IJCV07]

Transfer Learning:

Inter-class Transfe The ability of a system to recognize and apply knowledge and skills learned in previous tasks to novel tasks .[Pan & Yang, TKDE 2009]



Patch based Action comparison



Query

Template



M $D_{qt} = \sum_{i=1}^{M} \min_{j \in [1,N]} \left\{ \sum_{s=1}^{S} \min_{r \in R_s} d(q_{is}, t_{jr}) \right\}$ Frame-to-Frame

Distance



 Fb_y^+

 Fb_y^-

Motion Descriptor [Efros et al. ICCV 03]





Local Distance Function

$$D_{qt} = \sum_{s=1}^{\infty} d_{qt,s}$$
$$D_{qt} = \sum_{s=1}^{\infty} w_{q,s} d_{qt,s} = \langle \mathbf{w}_q, \mathbf{d}_{qt} \rangle$$

[Frome et al. NIPS06]

Local Distance Function

$$D_{qt} = \sum_{s=1} w_{q,s} d_{qt,s} = \langle \mathbf{w}_q \cdot \mathbf{d}_{qt} \rangle$$

Transferable Distance Function

Transferable Distance Function

Max-Margin Formulation

$$\min_{\mathbf{P},\xi} \quad \frac{1}{2} \|\mathbf{P}\|^2 + C \sum_{ijk} \xi_{ijk}$$

$$s.t.: \quad \forall i, j, k: \quad \langle \mathbf{P}^T \mathbf{F}_i \cdot \mathbf{x}_{ijk} \rangle \ge 1 - \xi_{ijk}$$

$$\forall m: \quad \langle \mathbf{P} \cdot \mathbf{f}_m \rangle \ge 0$$

$$\forall i, j, k: \quad \xi_{ijk} \ge 0$$

• Triplet $\mathbf{x}_{ijk} = \mathbf{d}_{ij} - \mathbf{d}_{ik}$

It is convex and similar to the primal problem of SVM

Hyper-Features

Codebook representation

- Descriptor for each patch $h = \{g, x, y\}$
 - HOG + Positions
- Obtaining codebook with the size of |V|
 - K-means clustering
- Hyper-feature for each patch
 - \circ A |V| dimensional vector ${f f}$

 $\mathbf{f}_m(v_i) = \frac{K_\sigma(D(v_i, h_m))}{\sum_{j=1}^{|V|} K_\sigma(D(v_j, h_m))}$

Summary of Features

Patch Matching

Motion Cue

Patch Weighting

Experiments on Action Recognition

 Train the transferable distance function on Weizmann, and test on KTH.

Visualization

 $w_i = \langle \mathbf{P} \cdot \mathbf{f}_i \rangle$

Learnt Weights on Testing Actions Codeword Ranking

Five Rounds of Experiments

For each round, we randomly select one actor, then choose one clip per action from this actor as the template.

Confusion Matrix of the Round 2

Efficiency

- With the learnt distance function, we can sort the patches on each frame by their saliency.
- Instead of using all patches, we can choose the top N patches with high weights for matching.

Human Action Detection

• • •

• • •

• • •

28

Cascade Structure

• • •

• • •

•••

• • •

Reject

Cascade Stage 1

Cascade Stage 2

Reject

• •

Cascade Stage N

Cascade Structure

Efficient Action Detection

Contributions

Transferable distance function Learning

- Hyper-features based on appearance and positions
- Max-margin Learning framework
- Action recognition from one clip
 - Template Matching based on motion
- Efficient action detection from one clip
 - Cascade structure

Thank You !

THINKING OF THE WORLD

