



Bilinear Classifiers for Visual Recognition

Hamed Pirsiavash Deva Ramanan Charless Fowlkes

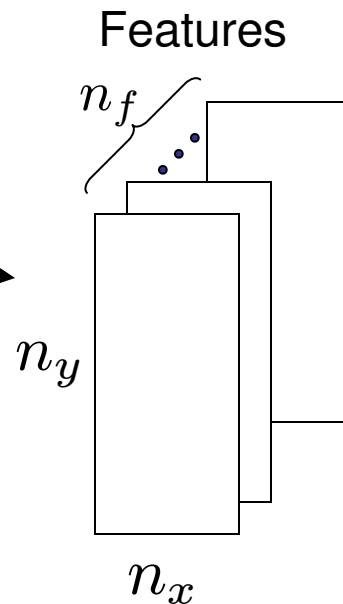
Computational Vision Lab.
University of California Irvine

To be presented in NIPS 2009

Introduction

■ Linear model for visual recognition

A window on
input image



Reshape
(:)

Feature vector

$$\mathbf{x} = \begin{bmatrix} \overleftrightarrow{1} \\ \vdots \end{bmatrix}_{n_y n_x n_f \times 1}$$

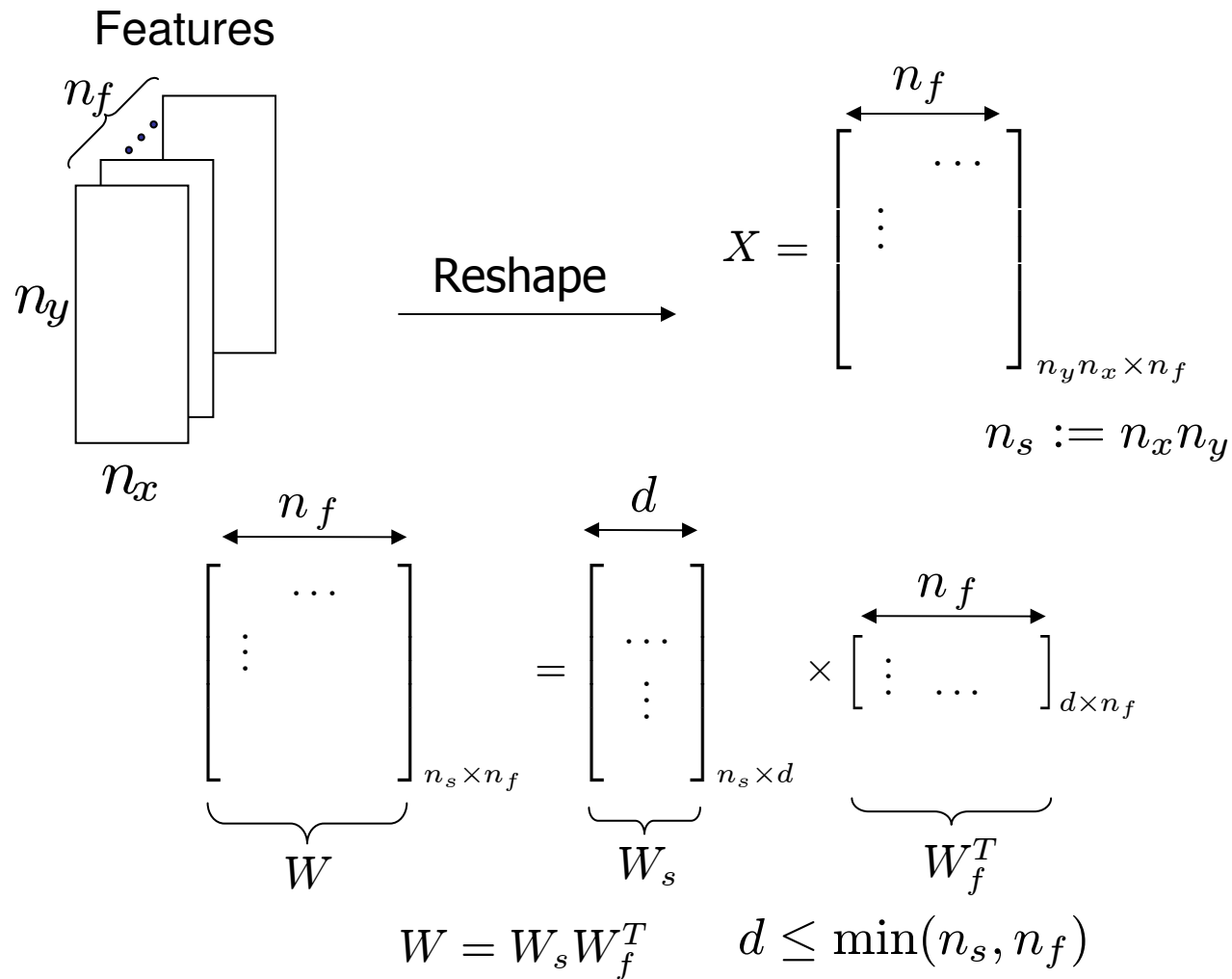


Introduction

- Linear classifier
 - Learn a template
 - Apply it to all possible windows

$$\underbrace{\begin{bmatrix} \vdots \end{bmatrix}^T}_{\text{Template (Model)}} \underbrace{\begin{bmatrix} \vdots \end{bmatrix}}_{\text{Feature vector}} > 0 \iff w^T x > 0$$

Introduction





Introduction

- Motivation for bilinear models
 - Reduced rank: less number of parameters
 - Better generalization: reduced over-fitting
 - Run-time efficiency
 - Transfer learning
 - Share a subset of parameters between different but related tasks



Outline

- Introduction
- Sliding window classifiers
- Bilinear model and its motivation
- Extension
- Related work
- Experiments
 - Pedestrian detection
 - Human action classification
- Conclusion

Sliding window classifiers

- Extract some visual features from a spatio-temporal window
 - e.g., histogram of gradients (HOG) in Dalal and Triggs' method
- Train a linear SVM using annotated positive and negative instances $w^T x > 0$

$$\min_w L(w) = \frac{1}{2} w^T w + C \sum_n \max(0, 1 - y_n w^T x_n)$$

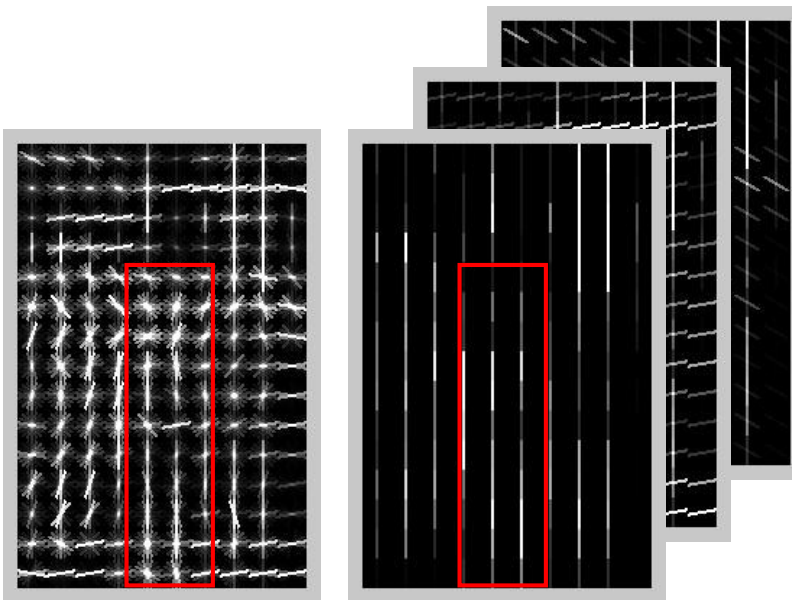
- Detection: evaluate the model on all possible windows in space-scale domain
 - Use convolution since the model is linear



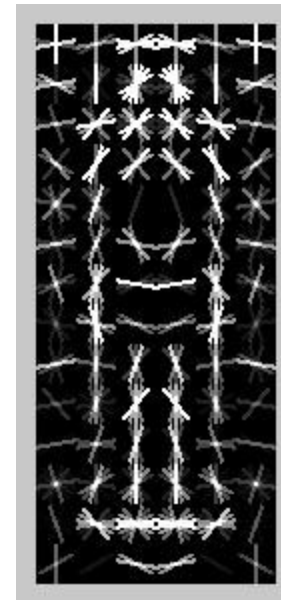
Sliding window classifiers



Sample image



Features



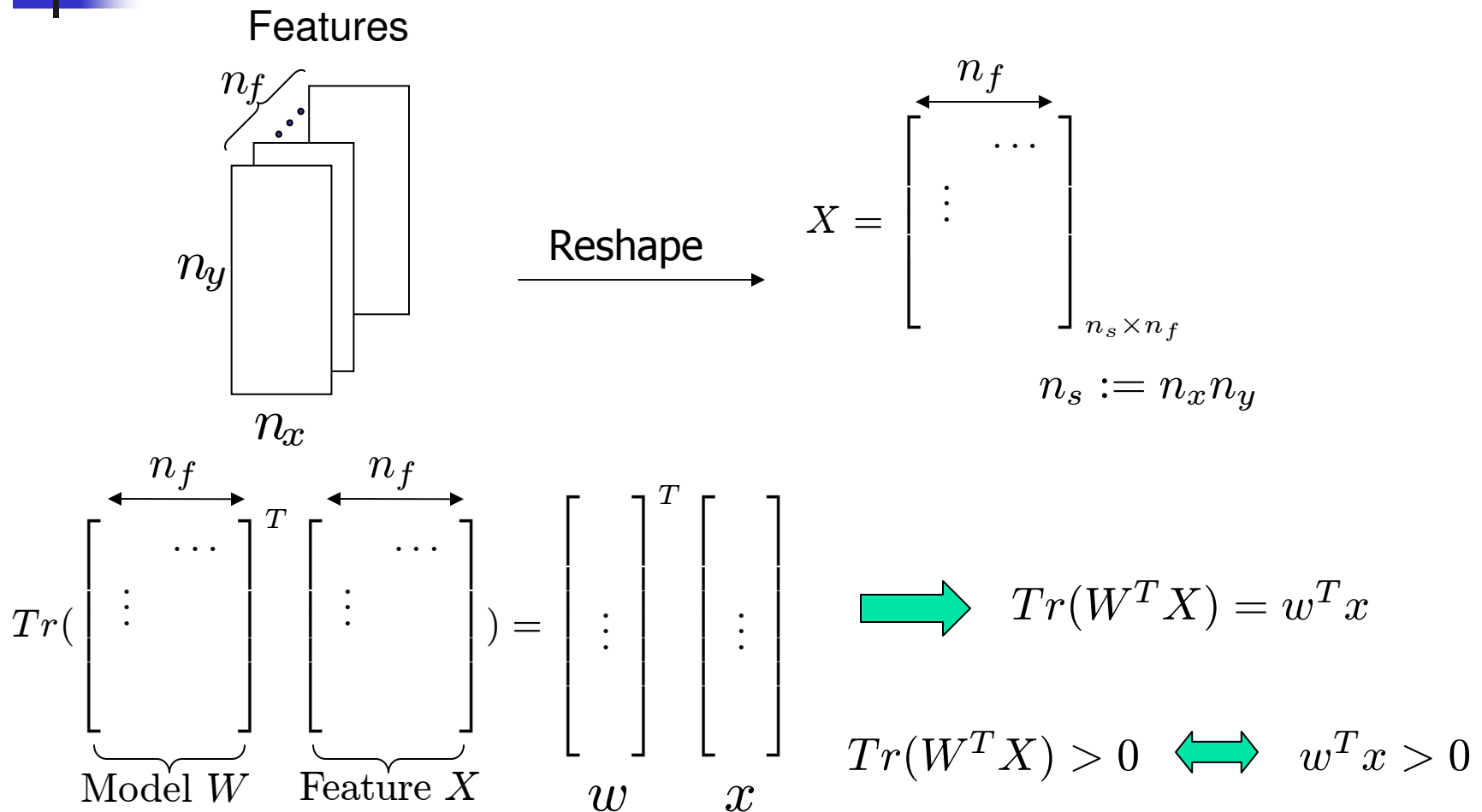
Sample
template W



Bilinear model (Definition)

- Visual data are better modeled as matrices/tensors rather than vectors
 - Why not use the matrix structure

Bilinear model (Definition)



Bilinear model (Definition)

■ Bilinear model

$$d \leq \min(n_s, n_f)$$

$$\underbrace{\begin{bmatrix} \cdots \\ \vdots \end{bmatrix}}_{W} \begin{matrix} \xleftarrow{n_f} \\ n_s \times n_f \end{matrix} = \underbrace{\begin{bmatrix} \cdots \\ \vdots \end{bmatrix}}_{W_s} \begin{matrix} \xleftarrow{d} \\ n_s \times d \end{matrix} \times \underbrace{\begin{bmatrix} \cdots \\ \vdots \end{bmatrix}}_{W_f^T} \begin{matrix} \xleftarrow{n_f} \\ d \times n_f \end{matrix}$$

$$W = W_s W_f^T$$

$$w^T x > 0 \quad \longleftrightarrow \quad f(x) > 0 \quad f(X) = \text{Tr}(W_f W_s^T X)$$

Bilinear model (Definition)

■ Bilinear model

$$d \leq \min(n_s, n_f)$$

$$\underbrace{\begin{bmatrix} \cdots \\ \vdots \end{bmatrix}}_{W} \begin{matrix} \xleftarrow{n_f} \\ n_s \times n_f \end{matrix} = \underbrace{\begin{bmatrix} \cdots \\ \vdots \end{bmatrix}}_{W_s} \begin{matrix} \xleftarrow{d} \\ n_s \times d \end{matrix} \times \underbrace{\begin{bmatrix} \cdots \\ \vdots \end{bmatrix}}_{W_f^T} \begin{matrix} \xleftarrow{n_f} \\ d \times n_f \end{matrix}$$

$$W = W_s W_f^T$$

$$w^T x > 0 \quad \longleftrightarrow \quad f(x) > 0$$

$$f(X) = \text{Tr}(W_f W_s^T X)$$

$$\longrightarrow f(X) = \text{Tr}(\underbrace{W_s^T X W_f}_{\text{Bilinear in } W_f \text{ and } W_s})$$

Bilinear in W_f and W_s 12

Bilinear model (Definition)

■ Bilinear model

$$d \leq \min(n_s, n_f)$$

$$\underbrace{\begin{bmatrix} \vdots & \dots \end{bmatrix}}_{W} \begin{matrix} \xleftarrow{n_f} \\ n_s \times n_f \end{matrix} = \underbrace{\begin{bmatrix} \vdots & \dots \end{bmatrix}}_{W_s} \begin{matrix} \xleftarrow{d} \\ n_s \times d \end{matrix} \times \underbrace{\begin{bmatrix} \vdots & \dots \end{bmatrix}}_{W_f^T} \begin{matrix} \xleftarrow{n_f} \\ d \times n_f \end{matrix}$$

$$W = W_s W_f^T$$

$$w^T x > 0 \quad \longleftrightarrow \quad f(x) > 0$$

$$f(X) = \text{Tr}(W_f W_s^T X)$$

$$\longrightarrow f(X) = \text{Tr}(\underbrace{W_s^T X W_f}_{\text{Bilinear in } W_f \text{ and } W_s})$$

Was Linear: $f(X) = \text{Tr}(W^T X)$

Bilinear in W_f and W_s 13



Bilinear model (Learning)

- Linear SVM for a given set of training pairs $\{X_n, y_n\}$

$$\min_W L(W) = \underbrace{\frac{1}{2} \text{Tr}(W^T W)}_{\text{Objective function}} + \underbrace{C}_{\substack{\downarrow \\ \text{Parameter}}} \underbrace{\sum_n \max(0, 1 - y_n \text{Tr}(W^T X_n))}_{\text{Constraints}}$$



Bilinear model (Learning)

- Linear SVM for a given set of training pairs $\{X_n, y_n\}$

$$\min_W L(W) = \underbrace{\frac{1}{2} \text{Tr}(W^T W)}_{\text{Regularizer}} + C \underbrace{\sum_n \max(0, 1 - y_n \text{Tr}(W^T X_n))}_{\text{Constraints}}$$

$W := W_s W_f^T$ ← Objective function

↓ Parameter



Bilinear model (Learning)

- For Bilinear formulation

$$\min L(W_s, W_f) = \frac{1}{2} \text{Tr}(W_f W_s^T W_s W_f^T) + C \sum_n \max(0, 1 - y_n \text{Tr}(W_f W_s^T X_n))$$

- Biconvex so solve by coordinate decent
 - By fixing one set of parameters, it's a typical SVM problem (with a change of basis)
 - Use off-the-shelf SVM solver in the loop



Motivation

- Regularization

- Similar to PCA, but not orthogonal and learned discriminatively and jointly with the template

$$W = W_s W_f^T$$

Reduced dimensional Template Subspace

- Run-time efficiency

- d convolutions instead of n_f

$$d \leq \min(n_s, n_f)$$



Motivation

- Transfer learning

$$W = W_s W_f^T$$

- Share the subspace W_f between different problems
 - e.g human detector and cat detector
- Optimize the summation of all objective functions
 - Learn a good subspace using all data



Extension

- Multi-linear

- High-order tensors

- $L(W_x, W_y, W_f)$ instead of just $L(W_s, W_f)$

- For 1D feature $L(W_x, W_y)$

- Separable filter for (Rank=1)

- Spatio-temporal templates

- $L(W_x, W_y, W_t, W_f)$



Related work (Rank restriction)

- Bilinear models
 - Often used in increasing the flexibility; however, we use them to reduce the parameters.
 - Mostly used in generative models like density estimation and we use in classification
- Soft Rank restriction
 - They used $Tr(W)$ rather than $Tr(W^T W)$ in SVM to regularize on rank
 - Convex, but not easy to solve
 - Decrease summation of eigen values instead of the number of non-zero eigen values (rank)
- Wolf et al (CVPR'07)
 - Used a formulation similar to ours with hard rank restriction
 - Showed results only for soft rank restriction
 - Used it only for one task (Didn't consider multi-task learning)



Related work (Transfer learning)

- Dates back to at least Caruana's work (1997)
 - We got inspired by their work on multi-task learning
 - Worked on: Back-propagation nets and k-nearest neighbor
- Ando and Zhang's work (2005)
 - Linear model
 - All models share a component in low-dimensional subspace (transfer)
 - Use the same number of parameters



Experiments: Pedestrian detection

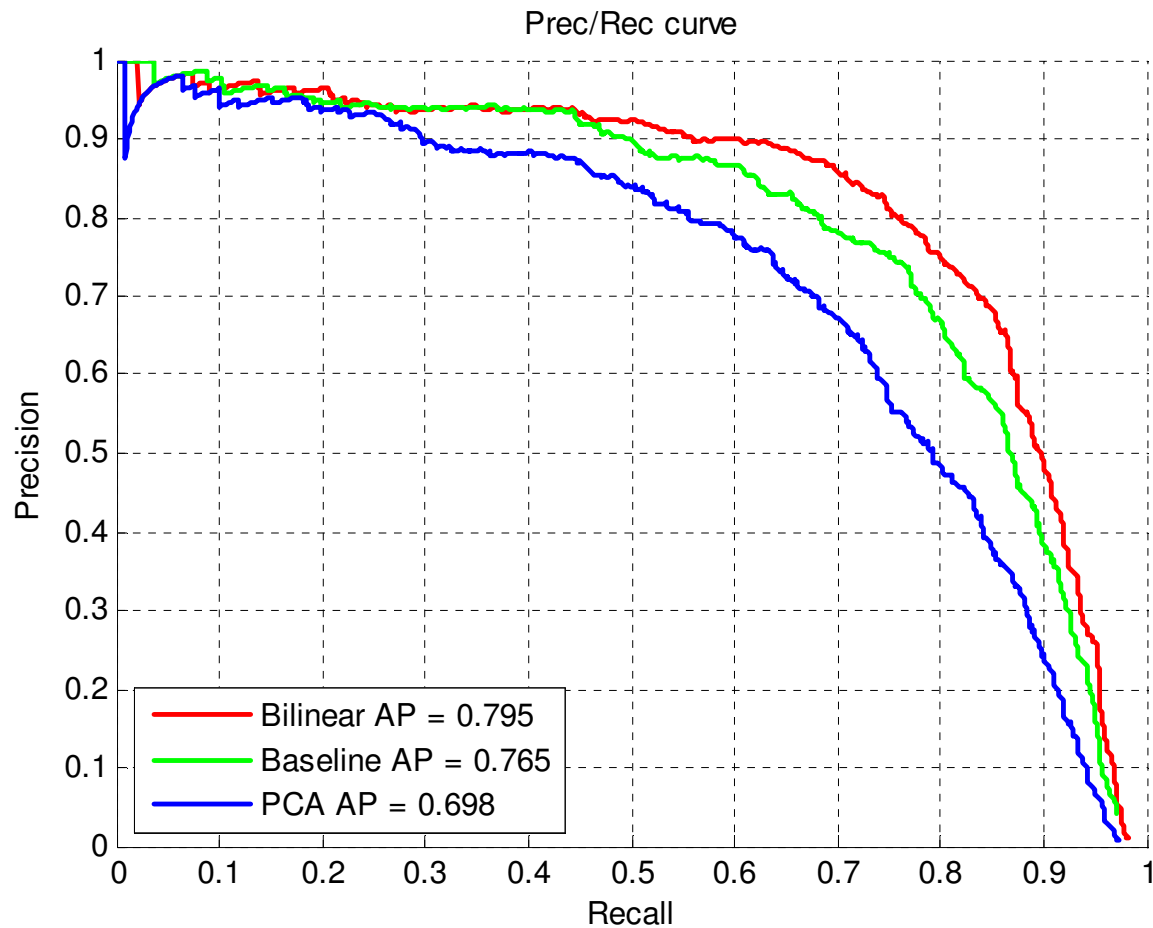
- Baseline: Dalal and Triggs' spatio-temporal classifier (ECCV'06)
 - Linear SVM on features: (84 for each 8×8 cell)
 - Histogram of gradient (HOG)
 - Histogram of optical flow
 - Made sure that the spatiotemporal is better than the static one by modifying the learning method



Experiments: Pedestrian detection

- Dataset: INRIA motion and INRIA static
 - 3400 video frame pairs
 - 3500 static images
- Typical values:
 - $n_s = 14 \times 6$, $n_f = 84$, $d = 5$ or 10
- Evaluation
 - Average precision
- Initialize with PCA in feature space
- Ours is 10 times faster

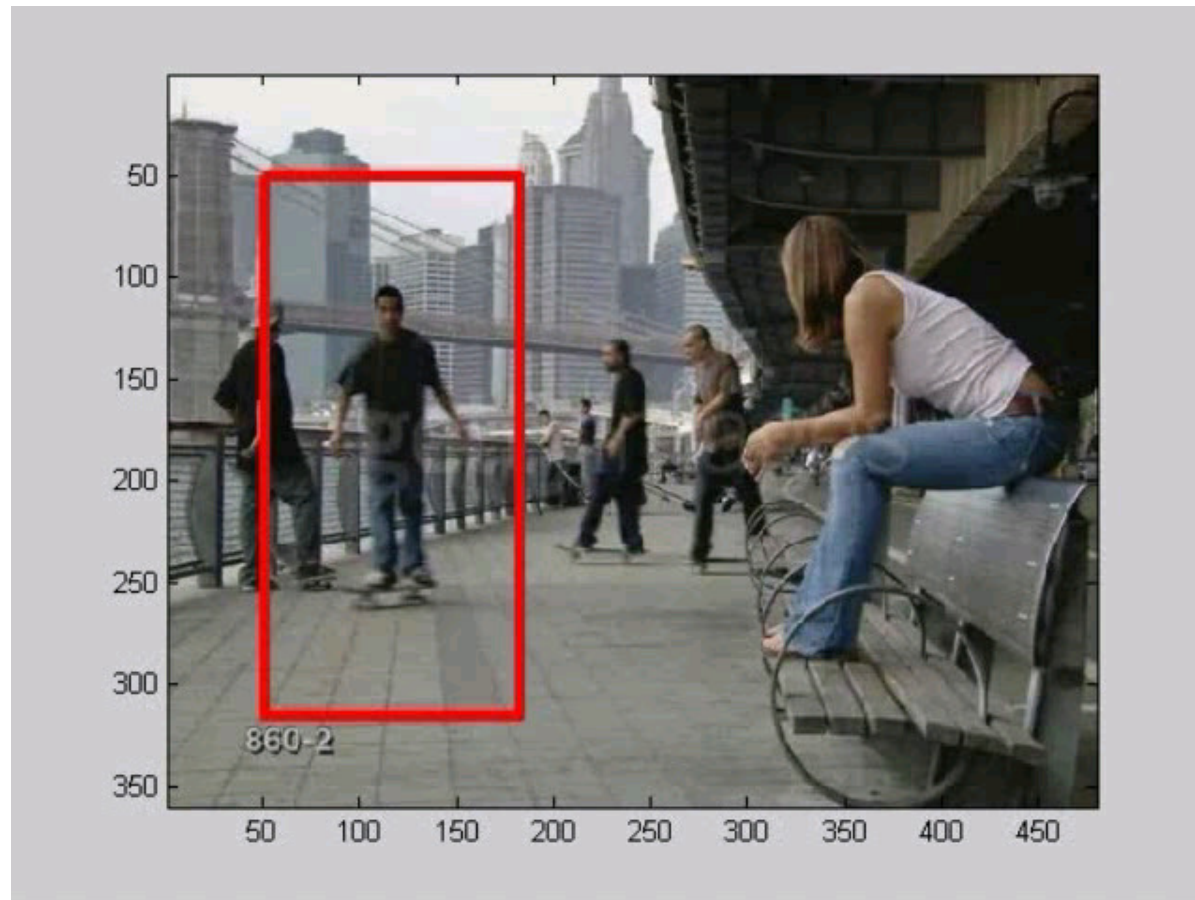
Experiments: Pedestrian detection



Experiments: Pedestrian detection



Experiments: Pedestrian detection



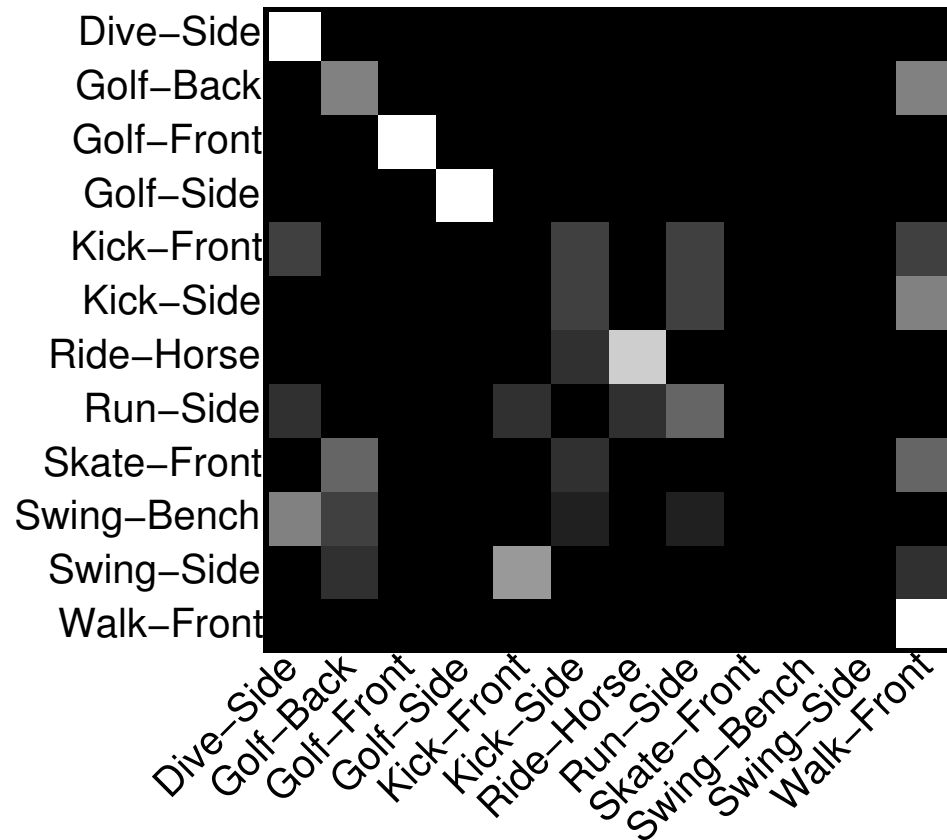
[Link](#)

Experiments:

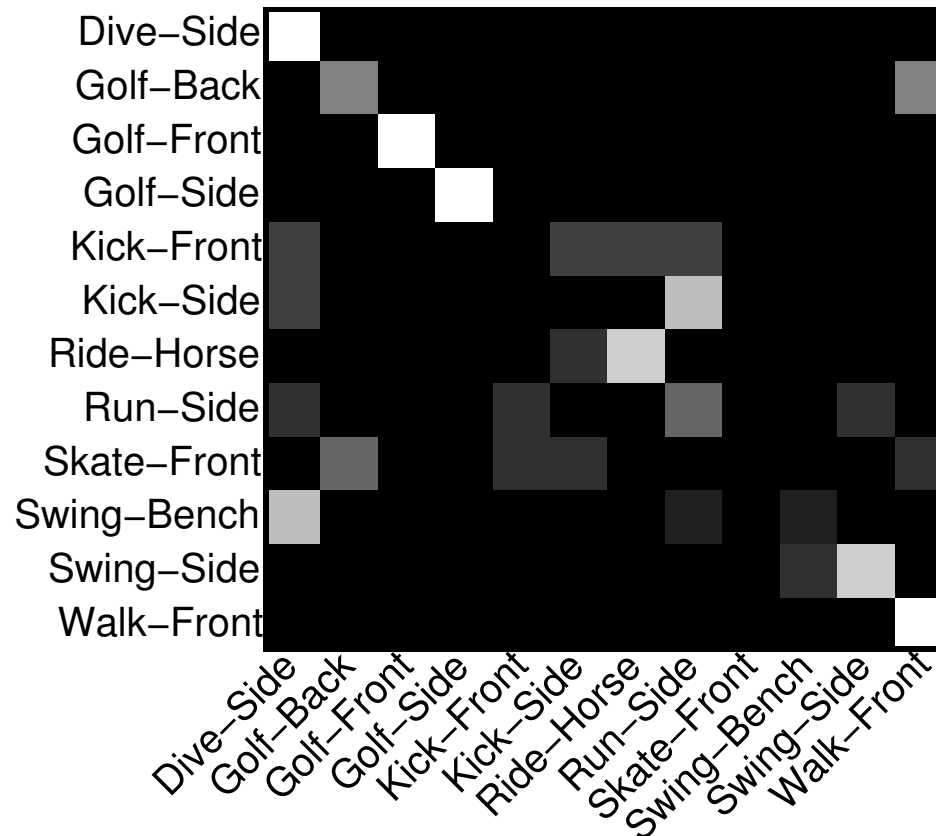
Human action classification 1

- 1 vs all action templates
- Voting:
 - A second SVM on confidence values
- Dataset:
 - UCF Sports Action (CVPR 2008)
 - They obtained 69.2%
 - We got 64.8% **but**
 - More classes: 12 classes rather than 9
 - Smaller dataset: 150 videos rather than 200
 - Harder evaluation protocol: 2-fold vs. LOOCV
 - 87 training examples rather than 199 in their case

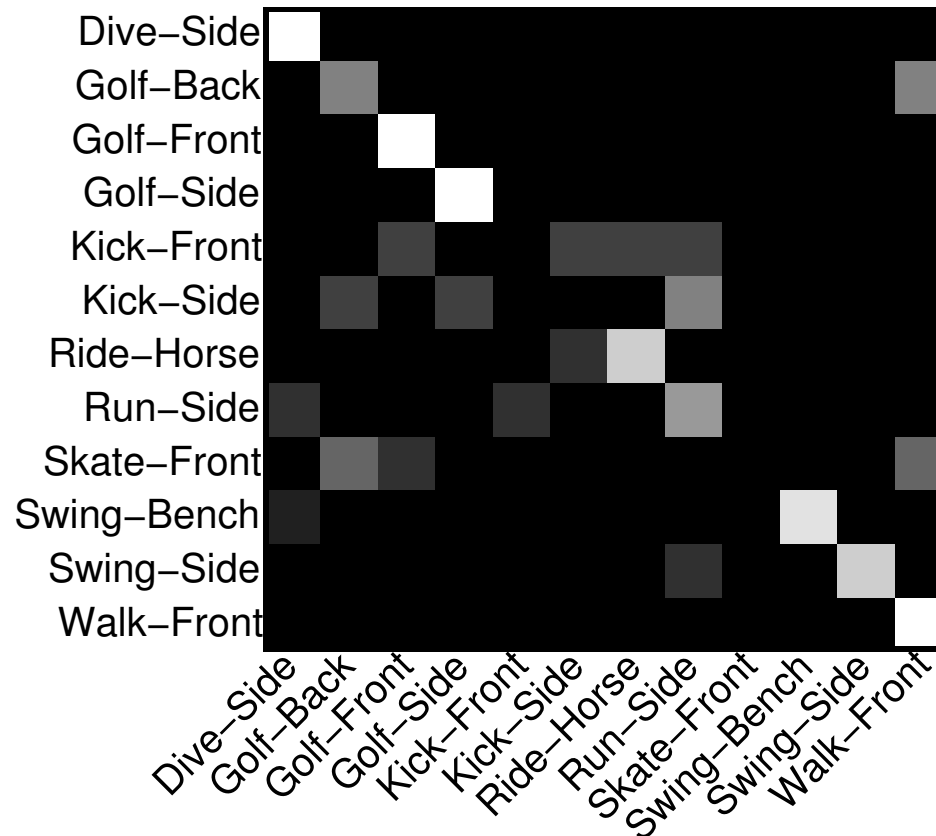
UCF action Results: PCA (0.444)



UCF action Results: Linear (0.518)

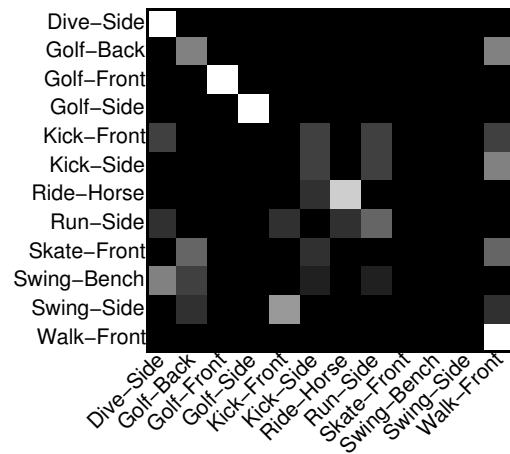


UCF action Results: Bilinear (0.648)

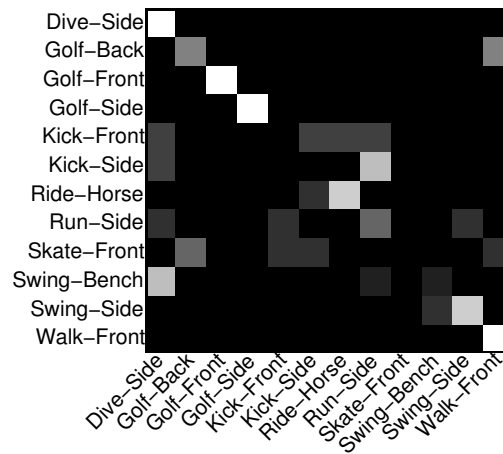


UCF action Results

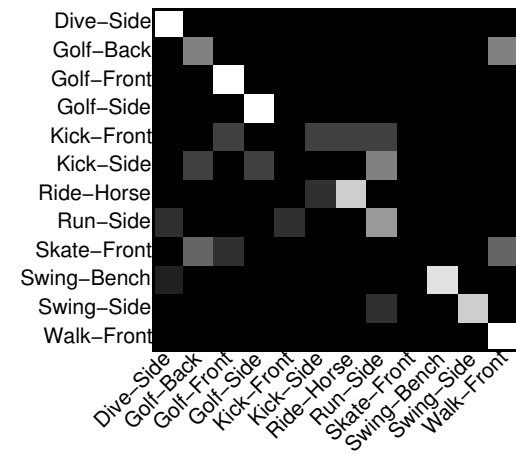
PCA on features (0.444)



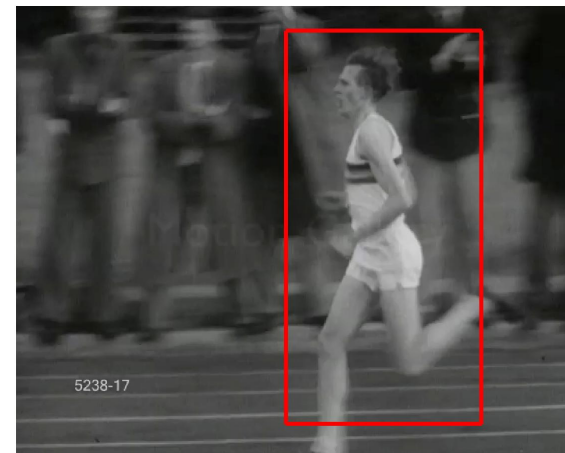
Linear (0.518)
(Not always feasible)



Bilinear (0.648)



UCF action Results



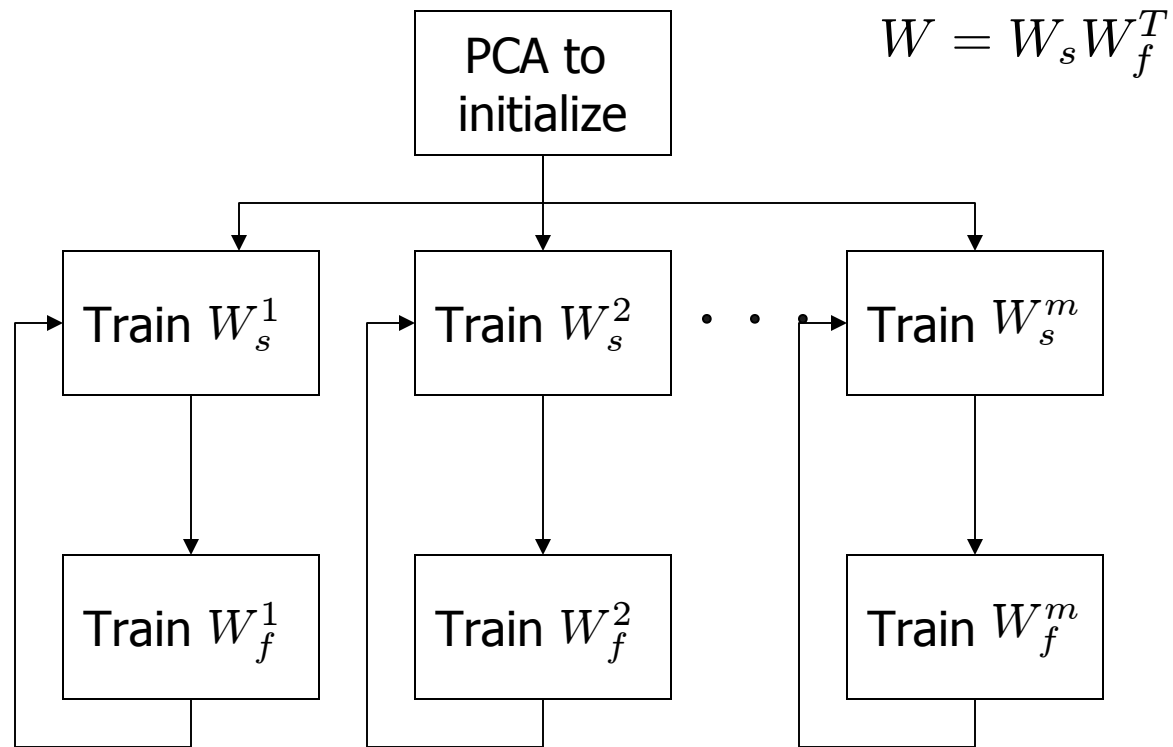
Experiments:

Human action classification 2

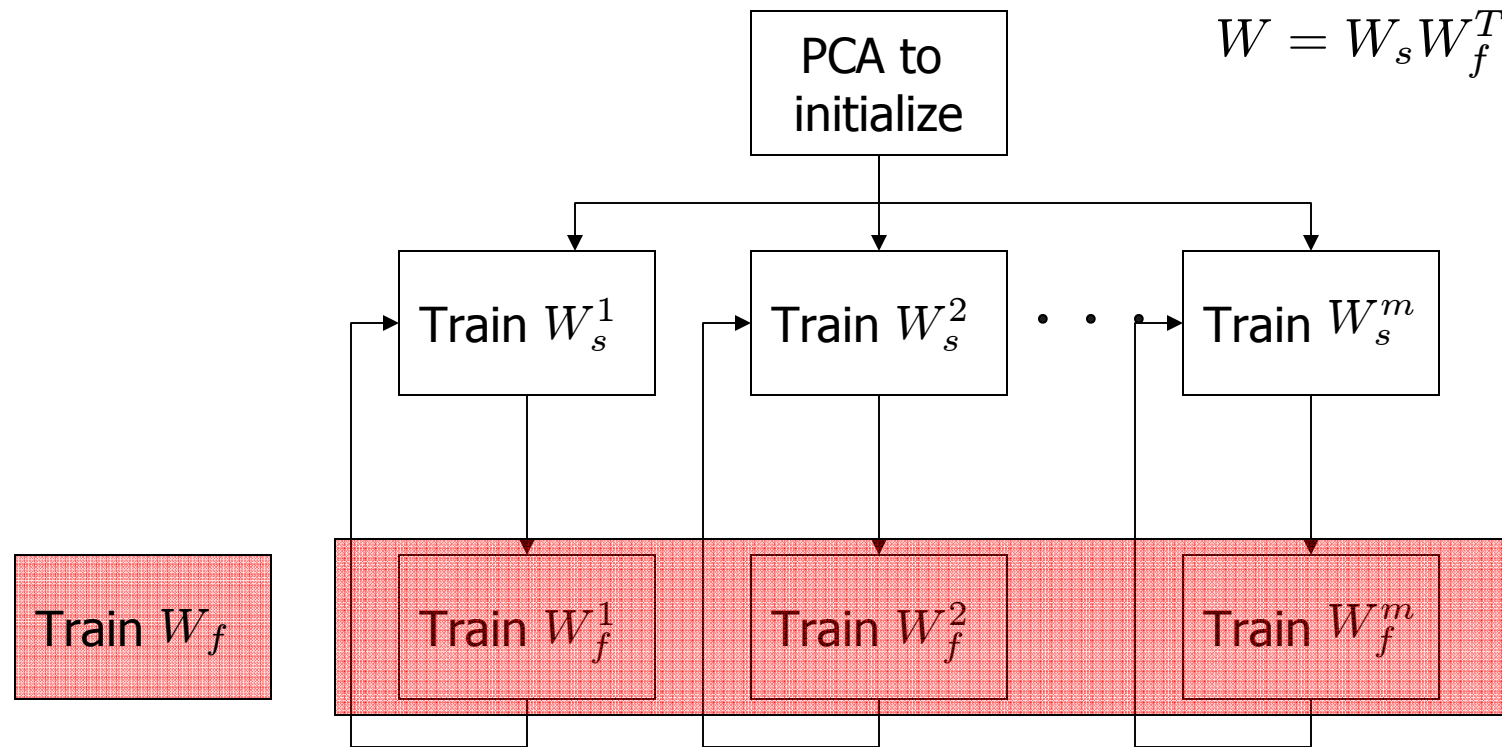
- Transfer
 - We used only two examples for each of 12 action classes
 - Once trained independently
 - Then trained jointly
 - Shared the subspace
 - Adjusted the C parameter for best result



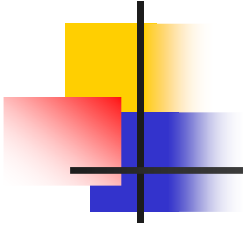
Transfer



Transfer



$$\min_{W_f} \sum_{i=1}^m L(W_f, W_s^i)$$

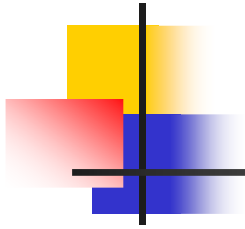


Results: Transfer

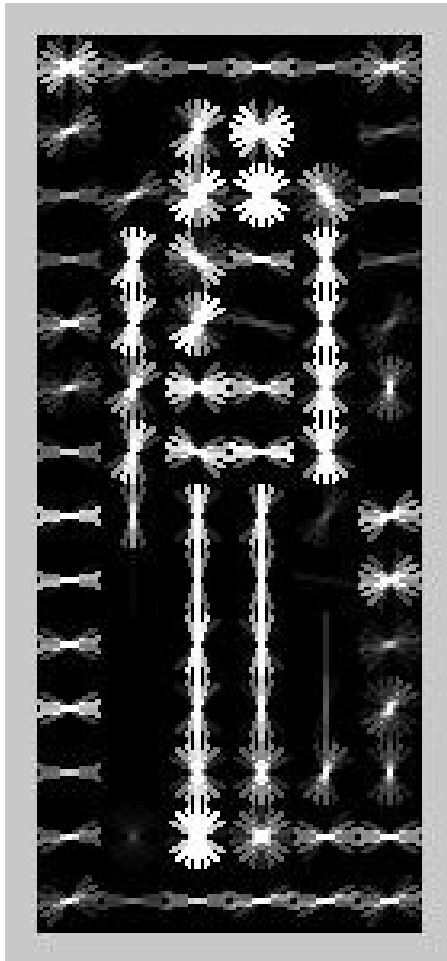
Average classification rate

	Coordinate decent iteration 1	Coordinate decent iteration 2
Independent bilinear (C=.01)	0.222	0.289
Joint bilinear (C=.1)	0.269	0.356

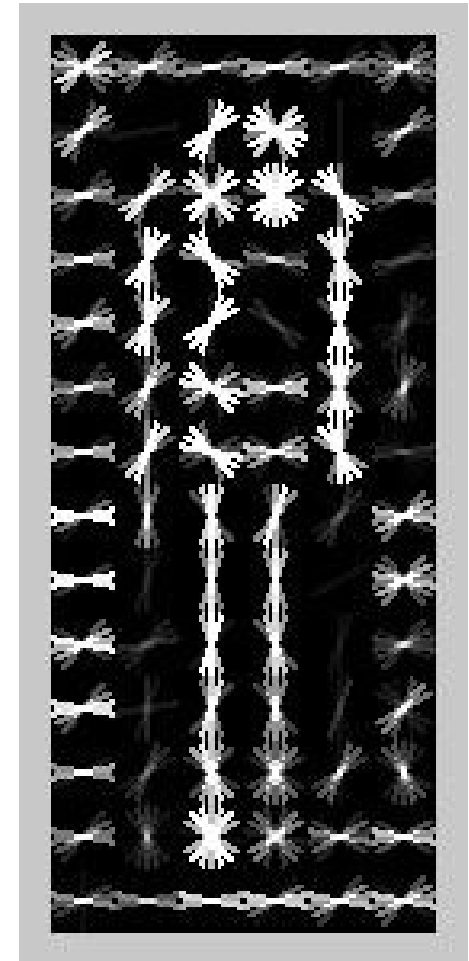
Results: Transfer (for “walking”)



Iteration 1



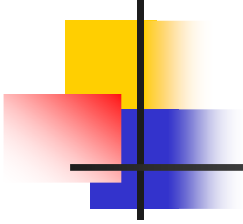
Refined at
Iteration 2





Conclusion

- Introduced multi-linear classifiers
 - Exploit natural matrix/tensor representation of spatio-temporal data
- Trained with existing efficient linear solvers
- Shared subspace for different problems
 - A novel form of transfer learning
- Got better performance and about 10X speed up in run-time compared to the linear classifier.
- Easy to apply to most high dimensional features (instead of dimensionality reduction methods like PCA)
- Simple: \sim 20 lines of Matlab code



Thanks!



Bilinear model (Learning details)

- Linear SVM for a given set of training pairs $\{x_n, y_n\}$

$$\min_W L(W) = \frac{1}{2} \text{Tr}(W^T W) + C \sum_n \max(0, 1 - y_n \text{Tr}(W^T X_n))$$

- For Bilinear formulation

$$\min L(W_f, W_s) = \frac{1}{2} \text{Tr}(W_f W_s^T W_s W_f^T) + C \sum_n \max(0, 1 - y_n \text{Tr}(W_f W_s^T X_n))$$

- It is **biconvex** so solve by coordinate decent



Bilinear model (Learning details)

- Each coordinate descent iteration:
freeze W_s

$$\min_{\tilde{W}_f} L(\tilde{W}_f, W_s) = \frac{1}{2}(\tilde{W}_f^T \tilde{W}_f) + C \sum_n \max(0, 1 - y_n \text{Tr}(\tilde{W}_f^T \tilde{X}_n))$$

$$\text{where } \tilde{W}_f = A_s^{\frac{1}{2}} W_f^T, \tilde{X}_n = A_s^{-\frac{1}{2}} W_s^T X_n, A_s = W_s^T W_s$$

then freeze W_f

$$\min_{\tilde{W}_s} L(W_f, \tilde{W}_s) = \frac{1}{2}(\tilde{W}_s^T \tilde{W}_s) + C \sum_n \max(0, 1 - y_n \text{Tr}(\tilde{W}_s^T \tilde{X}_n))$$

$$\text{where } \tilde{W}_s = W_s A_f^{\frac{1}{2}}, \tilde{X}_n = X_n W_f A_f^{-\frac{1}{2}}, A_f = W_f^T W_f$$