

Machine Learning approaches for Computer Vision applications

Alexandros Iosifidis

Aarhus University, Department of Engineering, ECE

alexandros.iosifidis@eng.au.dk

Short bio

Academic positions

- 2017 – Assistant Prof. of Machine Learning and Computer Vision, Aarhus University
- 2018 – Adjunct Prof. (Docent) of Machine Learning, Tampere University of Technology
- 2017 – Adjunct Prof. (Docent) of Data Management & Data Analytics, Lappeenranta University of Technology
- 2016-2017 Academy of Finland Postdoctoral Research Fellow
- 2015-2017 TTY Foundation Postdoctoral Research Fellow

Education

- 2014 PhD in Informatics (Computer Science), Aristotle University of Thessaloniki
- 2010 M.Eng. in Electrical & Computer Engineering, Democritus University of Thrace
- 2008 Diploma in Electrical & Computer Engineering, Democritus University of Thrace

Research interests

- Statistical Machine Learning and Artificial Neural Networks
- Image/Video/Signal analysis, Computer Vision, Computational Finance

Overview

Computer Vision applications

- › Image analysis/recognition/segmentation
- › Video analysis/segmentation
- › Face recognition/verification and affective computing
- › Human action recognition/localization/segmentation
- › Domain-specific image/video analysis

Recent contributions

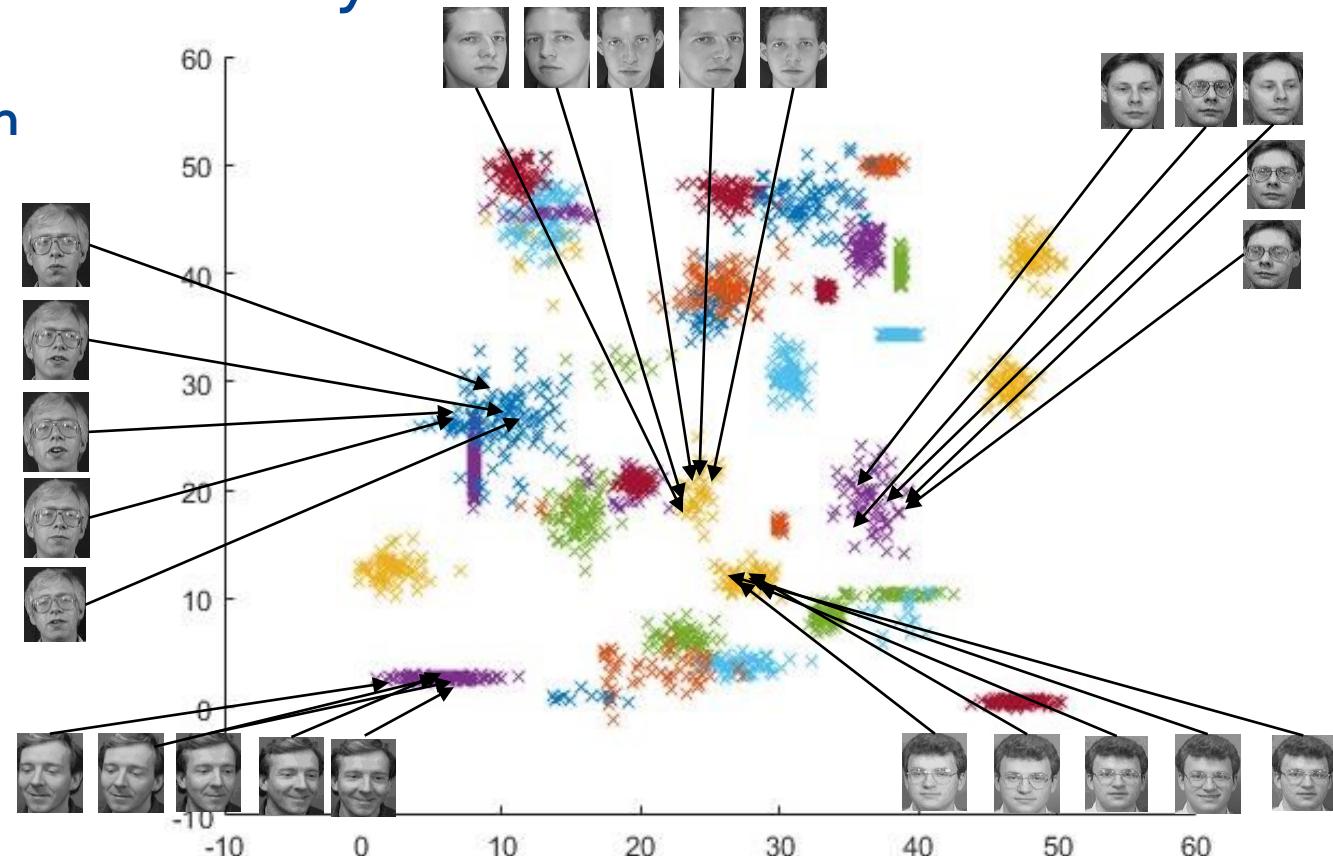
- › Graph-based analysis/recognition/clustering
- › Max-margin Classification
- › Discriminant Learning
- › Kernel-based learning
- › Multi-view/modal Data Analysis
- › Neural Network (Deep Learning) acceleration
- › Data-driven (Deep) Architecture Learning

Extensions to applications of other domains

Computer Vision applications

Human Face Analysis

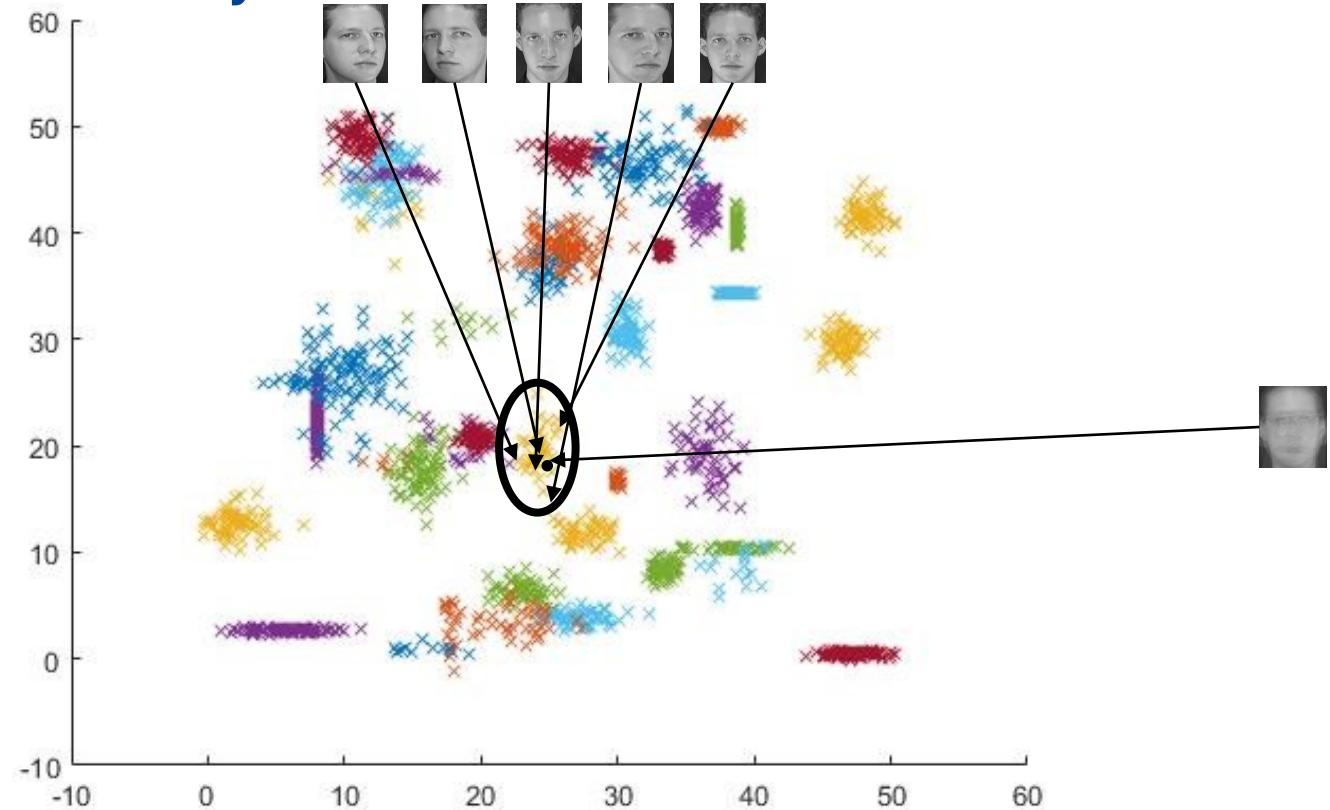
Face recognition



- A. Iosifidis, A. Tefas and I. Pitas, "Kernel Reference Discriminant Analysis", Pattern Recognition Letters, 2014
A. Iosifidis, A. Tefas and I. Pitas, "Class-specific Reference Discriminant Analysis with application in Human Behavior Analysis", IEEE Transactions on Human-Machine Systems, 2015

Human Face Analysis

Face verification
Face identification

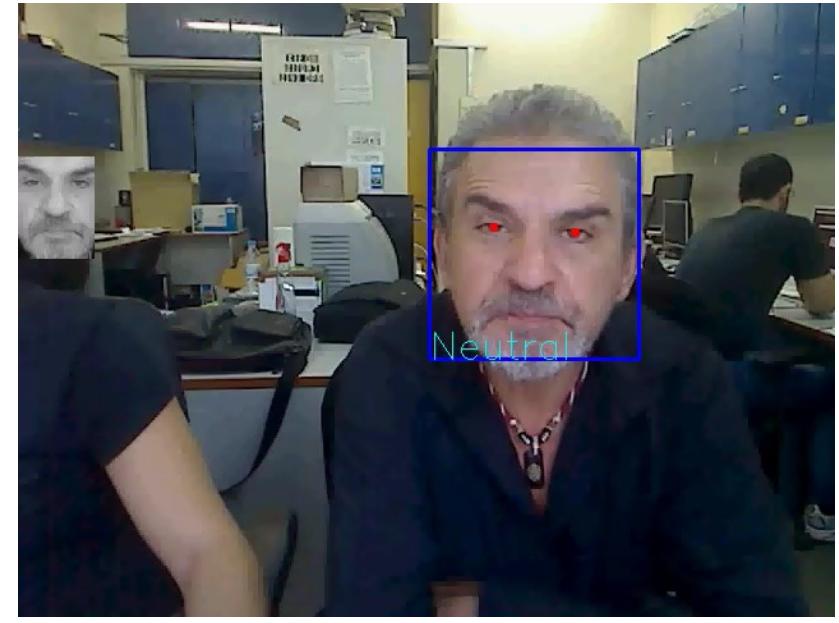
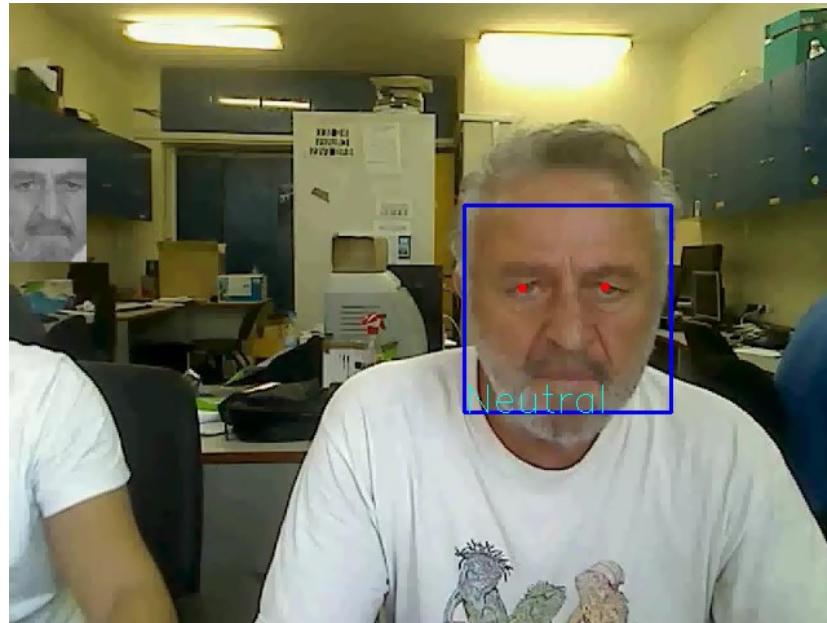


A. Iosifidis and M. Gabbouj, "Class-Specific Kernel Discriminant Analysis revisited: further analysis and extensions", IEEE Transactions on Cybernetics, 2017

A. Iosifidis and M. Gabbouj, "Scaling up Class-Specific Kernel Discriminant Analysis for large-scale Face Verification", IEEE Transactions on Information Forensics and Security, 2016

Human Face Analysis

Facial expression recognition (Affective Computing)



Human Face Analysis

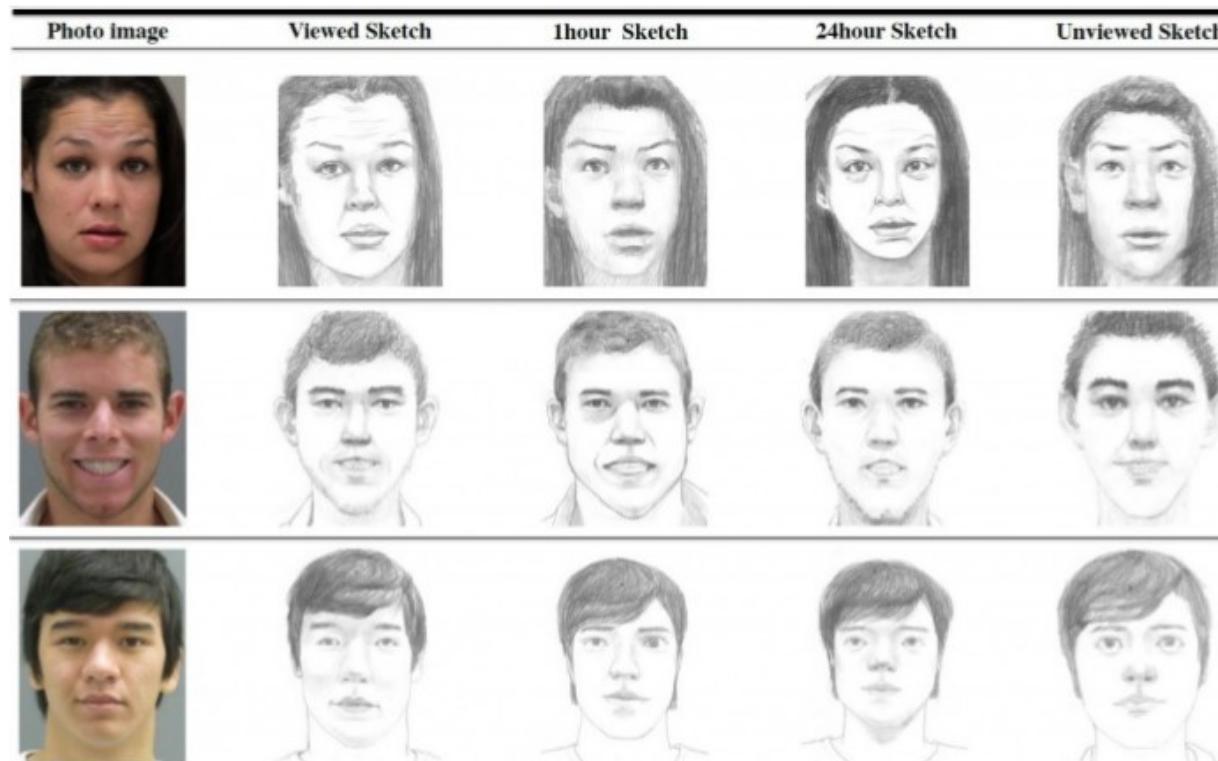
Visual Voice Activity Detection → Assign the correct face to the observed voice



3DTV FP7-ICT

Human Face Analysis

Face - Sketch recognition



G. Cao, A. Iosifidis and M. Gabbouj, "Multi-modal Subspace Learning with Dropout regularization for Cross-modal Recognition and Retrieval ", IPTA 2016 (*Best Student Paper Award*)

Image Analysis

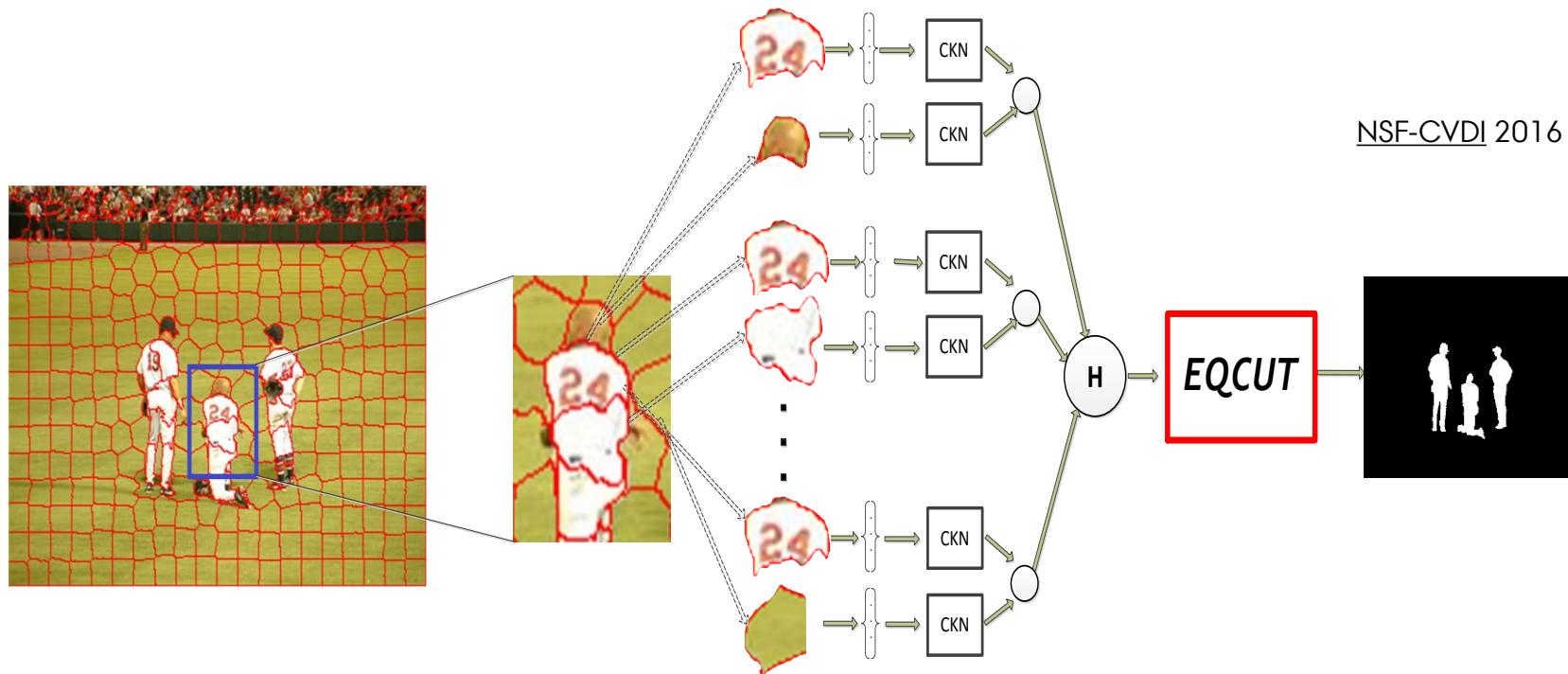
Salient object segmentation → Unsupervised (generic) case

NSF-CVBI 2016



Image Analysis

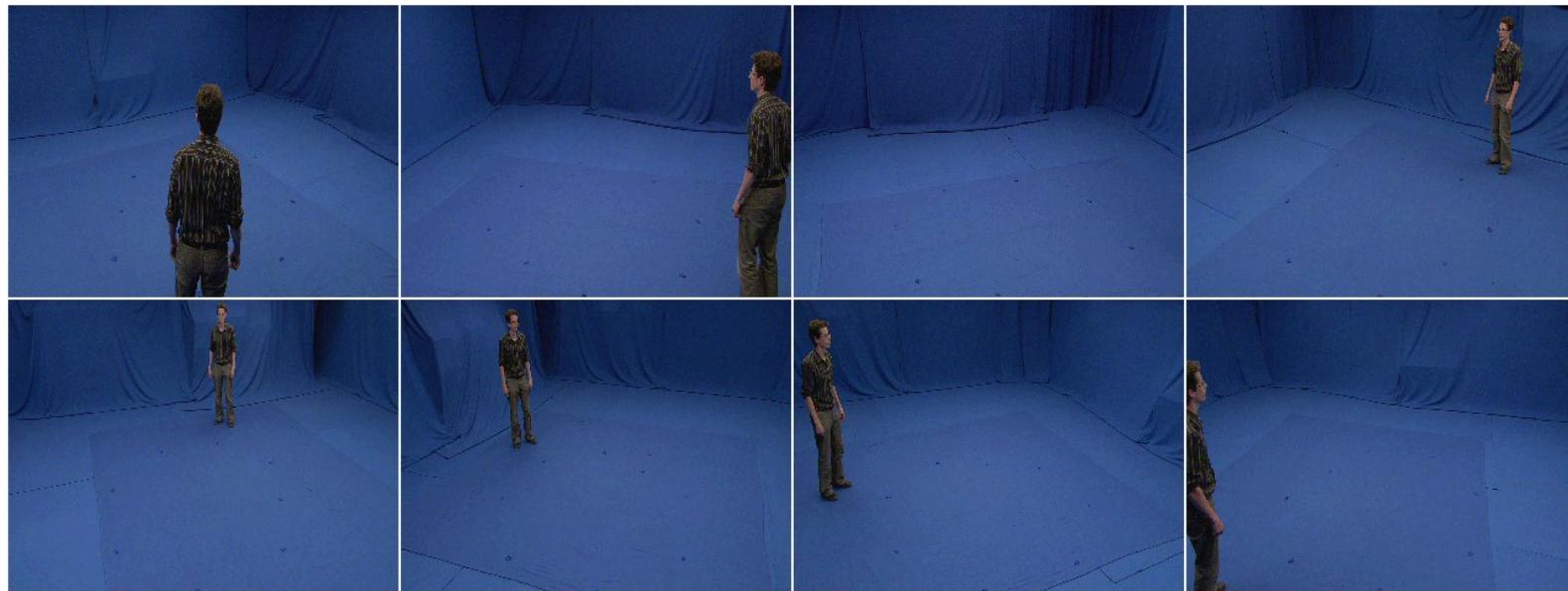
Salient object segmentation → Supervised (User-directed case)



Human Action Recognition

How to combine action observations from various views/cameras
Restricted application scenario → **Movie production**

i3DPost FP6



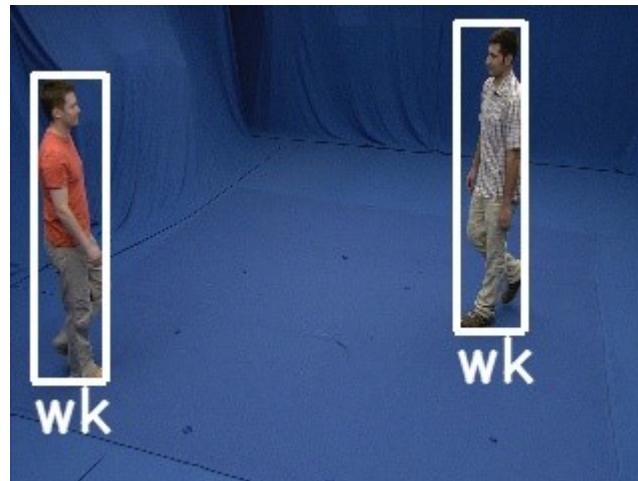
- A. Iosifidis, A. Tefas and I. Pitas, "View-invariant action recognition based on Artificial Neural Networks", IEEE Transactions on Neural Networks and Learning Systems, 2012
- A. Iosifidis, A. Tefas, N. Nikolaidis and I. Pitas, "Multi-view Human Movement Recognition based on Fuzzy Distances and Linear Discriminant Analysis", Computer Vision and Image Understanding, 2012.

Human Action Recognition

How to combine action observations from various views/cameras

Restricted application scenario → **Movie production**

Localization of people and classification of their actions



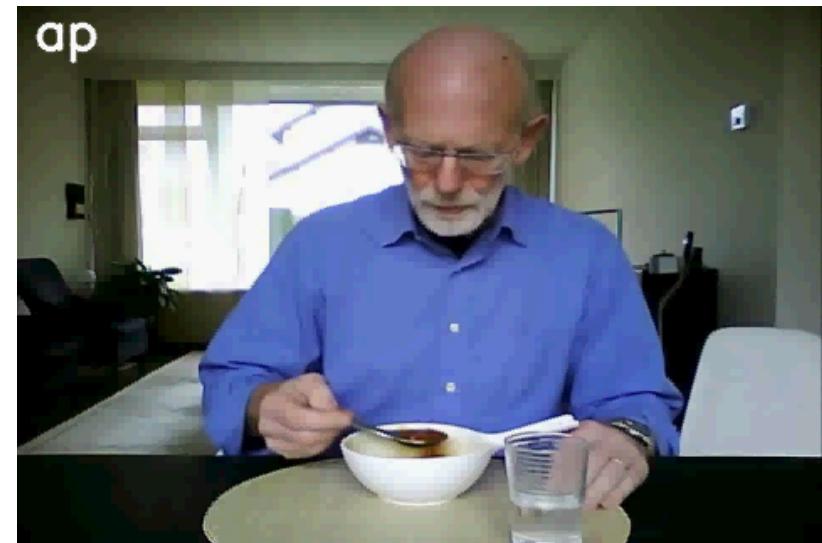
i3DPost FP6

- A. Iosifidis, A. Tefas and I. Pitas, "View-invariant action recognition based on Artificial Neural Networks", IEEE Transactions on Neural Networks and Learning Systems, 2012
- A. Iosifidis, A. Tefas, N. Nikolaidis and I. Pitas, "Multi-view Human Movement Recognition based on Fuzzy Distances and Linear Discriminant Analysis", Computer Vision and Image Understanding, 2012.

Human Action Recognition

Human action recognition for assisted living of the elderly

MOBISERV FP7



A. Iosifidis, E. Marami, A. Tefas, I. Pitas and K. Lyroudia, “The MOBISERV-AIIA Eating and Drinking multi-view database for vision-based assisted living”, J-IHMSP, 2015.

Human Action Recognition

What can we do for complex actions in complex/cluttered scenes?



Human Action Recognition

What can we do for complex actions in complex/cluttered scenes?

We focus on Space-Time Interest Points (STIPs) and follow the similar approaches



A. Iosifidis, A. Tefas and I. Pitas, "Discriminant Bag of Words based Representation for Human Action Recognition", Pattern Recognition Letters, 2014

A. Iosifidis, A. Tefas and I. Pitas, "Distance-based Human Action Recognition using optimized class representations", Neurocomputing, 2015

Human Action Recognition

When enriched visual information is available → Stereo Cameras

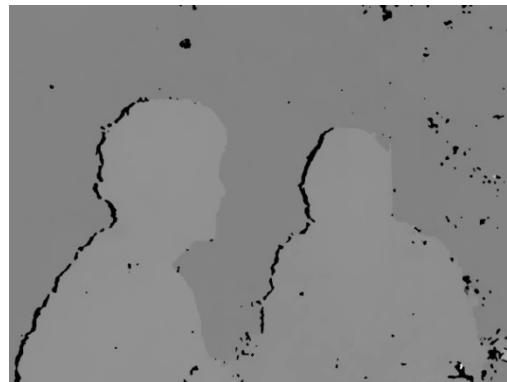
Left

channel



Disparity

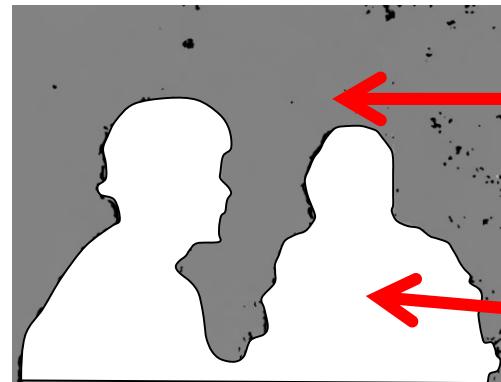
map



3DTV-S FP7-ICT

Right

channel



Disparity zone 2

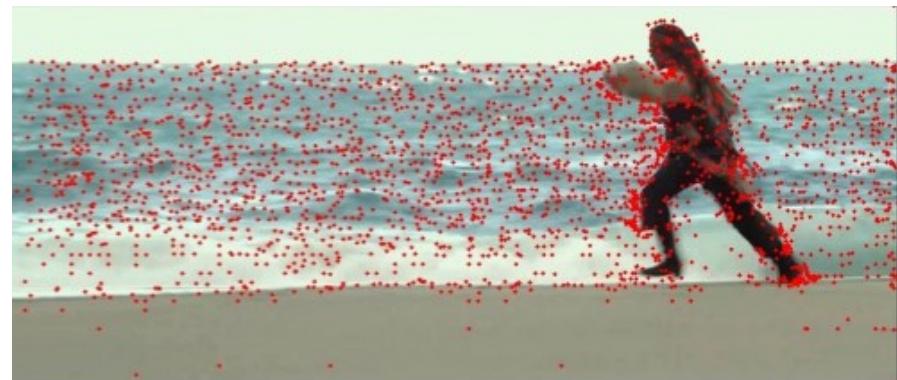
Disparity zone 1

Human Action Recognition

When enriched visual information is available → Stereo Cameras

3DTV-S FP7-ICT

Original DT-based
description

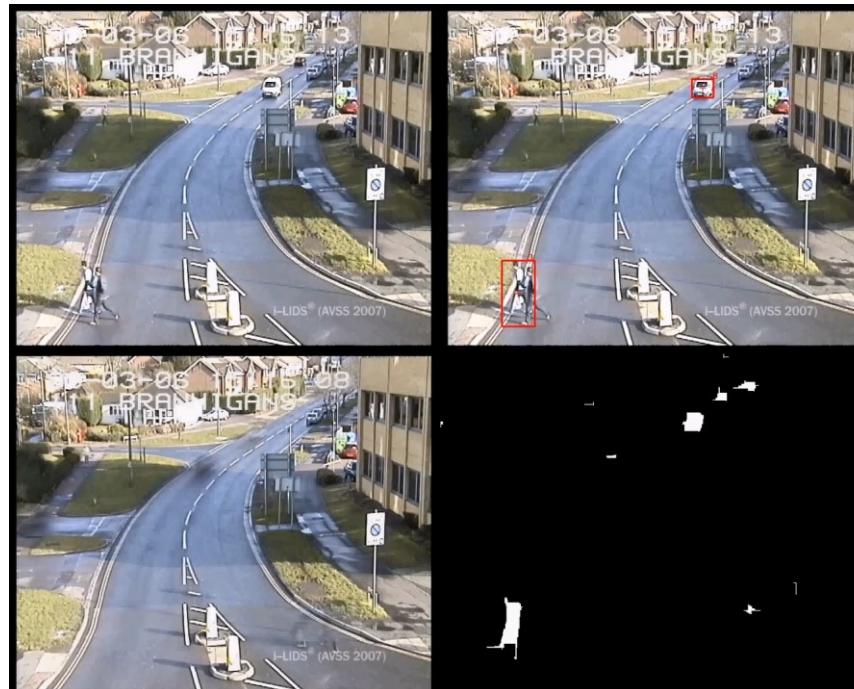


Disparity-enhanced DT-
based description



Video analysis

Object tracking

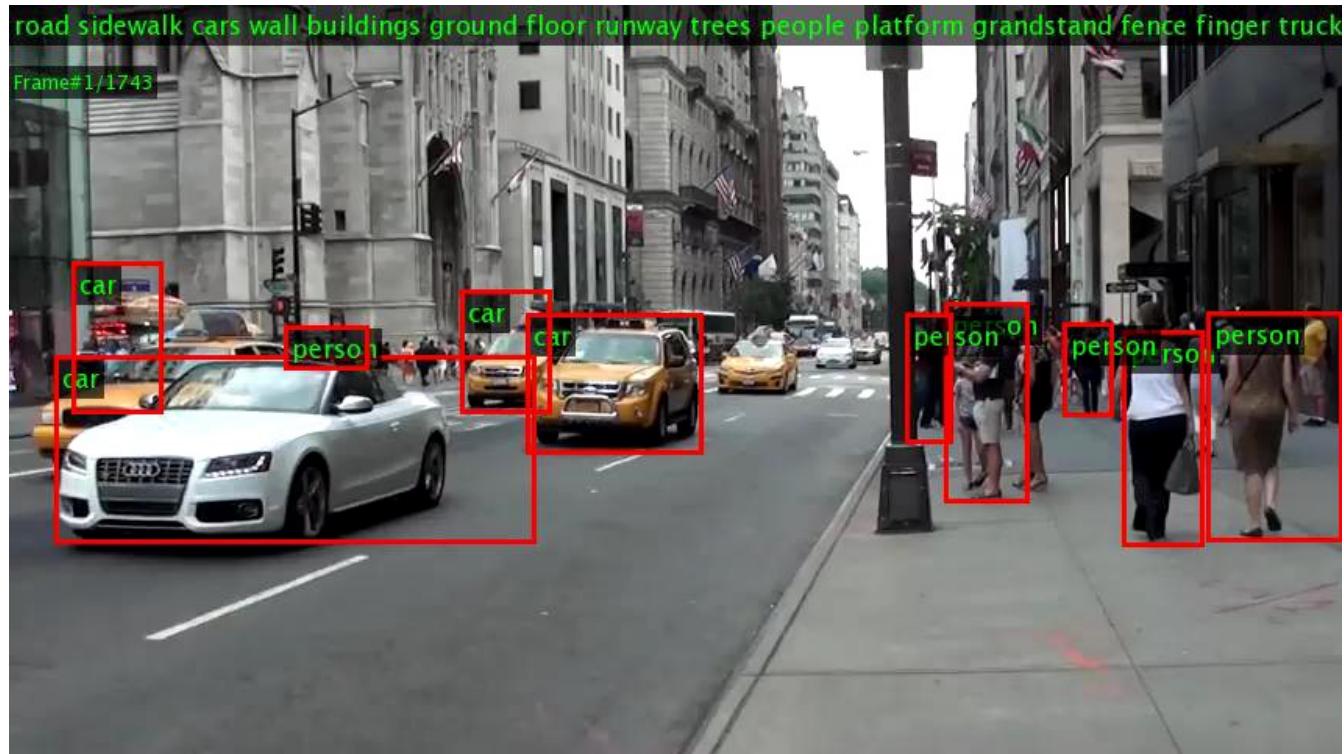


Video analysis

Object detection/recognition

Scene analysis

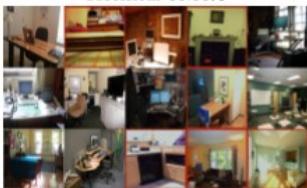
NSF-CVDI 2015



Multi-view/modal Data analysis

Image and Text (I2T and T2I) Retrieval

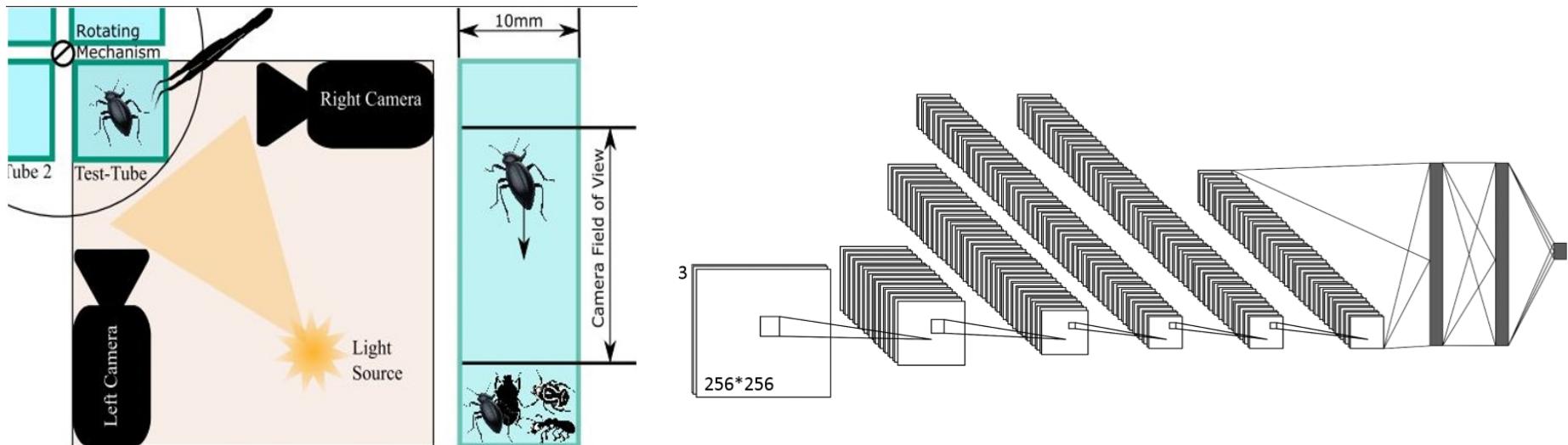
NSF-CVDI 2016

Image Query	Text Query
	1. A very big building with many windows and a clock on it. 2. A very old tall building with a large clock tower sticking out of it. 3. The clock tower stands high above the city. 4. A clock that is on the side of a large building. 5. The bridge is in front of a huge building with a clock tower in the middle of it.
Precision: 53.33%	Precision: 86.67%
	
(a) Query by original image feature	(b) Query by projected image feature
	Precision: 100%
(c) Query by text	
Image Query	Text Query
	1. An open laptop sits on a desk in front of a window. 2. An Apple laptop sitting on a wooden desk. 3. An Apple laptop sitting on a wooden desk in an office. 4. An Apple laptop on a desk in an office. 5. A desk with a laptop sitting on top of it.
Precision: 60.00%	Precision: 86.67%
	
(a) Query by original image feature	(b) Query by projected image feature
	Precision: 66.67%
(c) Query by text	

Domain-specific Media Data Analysis

Classification of Aquatic Macroinvertebrates (bugs in lakes)

AoF DETECT



- J. Arje, S. Karkkainen, K. Meissner, A. Iosifidis, T. Ince, M. Gabbouj and S. Kiranyaz, “The effect of automated taxa identification errors on biological indices”, Expert Systems with Applications, 2017
- J. Raitoharju, E. Riabchenko, K. Meissner, I. Ahmad, A. Iosifidis, M. Gabbouj and S. Kiranyaz, “Data Enrichment in Fine-Grained Classification of Aquatic Macroinvertebrates”, ICPR, 2016

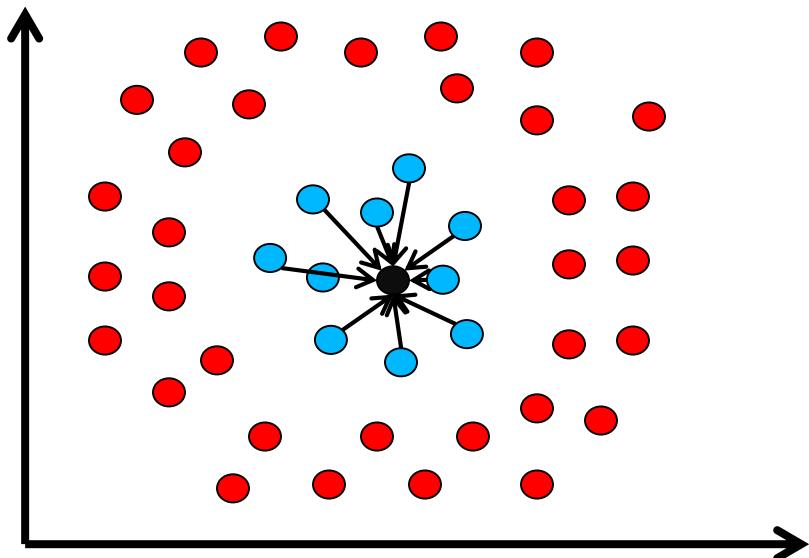
Recent Contributions

Discriminant Learning

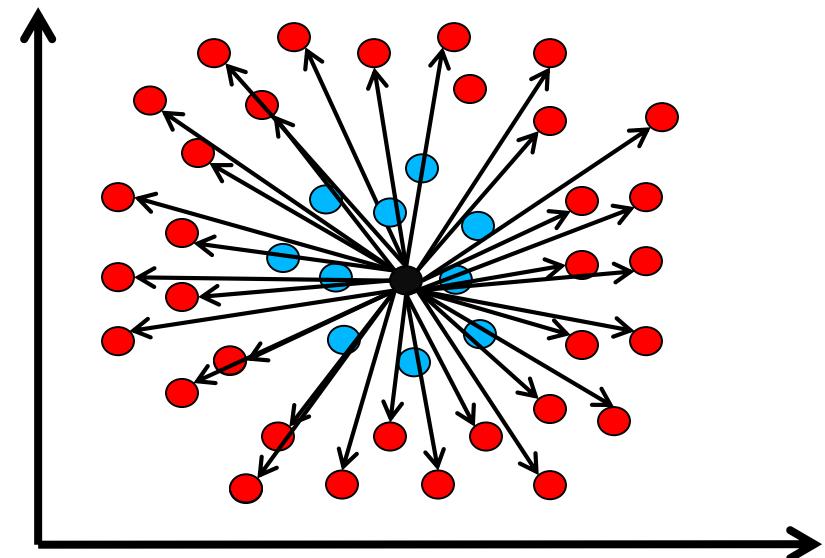
Class-specific kernel Discriminant Analysis

- › Data represented as vectors
- › Vector to vector transformation/mapping
- › Find a data mapping that maximizes discrimination of the class of interest from the rest of the world

Intra-class scatter S_i



Out-of-class scatter S_p



Discriminant Learning

➤ Traditional CSKDA:

1. Data mapping to the feature space:

$$\mathbf{x}_i \in \mathbb{R}^D \rightarrow \phi(\mathbf{x}_i) \in \mathcal{F}$$

2. Application of the linear projection in

$$\mathbf{S}_i = \Phi \mathbf{L}_i \Phi^T$$

$$\mathbf{S}_p = \Phi \mathbf{L}_p \Phi^T$$

$$\mathbf{W}^T \mathbf{S}_i \mathbf{W} = \mathbf{A}^T \Phi^T \Phi \mathbf{L}_i \Phi^T \Phi \mathbf{A} = \mathbf{A}^T \mathbf{K} \mathbf{L}_i \mathbf{K} \mathbf{A}$$

$$\mathbf{W}^T \mathbf{S}_p \mathbf{W} = \mathbf{A}^T \Phi^T \Phi \mathbf{L}_p \Phi^T \Phi \mathbf{A} = \mathbf{A}^T \mathbf{K} \mathbf{L}_p \mathbf{K} \mathbf{A}$$

- \mathbf{A} is calculated by applying eigen-analysis to the matrix $(\mathbf{K} \mathbf{L}_p \mathbf{K})^{-1} (\mathbf{K} \mathbf{L}_i \mathbf{K})$ ← scaling to big data issues!
stability issues!

Discriminant Learning

We showed that the standard CSKDA solution is equivalent to:

$$\hat{\mathcal{J}} = \|\mathbf{W}^T \Phi - \mathbf{T}\|_F^2, \quad s.t. : \text{rank}(W) \leq d$$

$$\hat{\mathcal{J}} = \|\mathbf{B}^T (\mathbf{Q}^T \Phi) - \mathbf{T}\|_F^2 = \|\mathbf{B}^T (\mathbf{A}^T \mathbf{K}) - \mathbf{T}\|_F^2$$

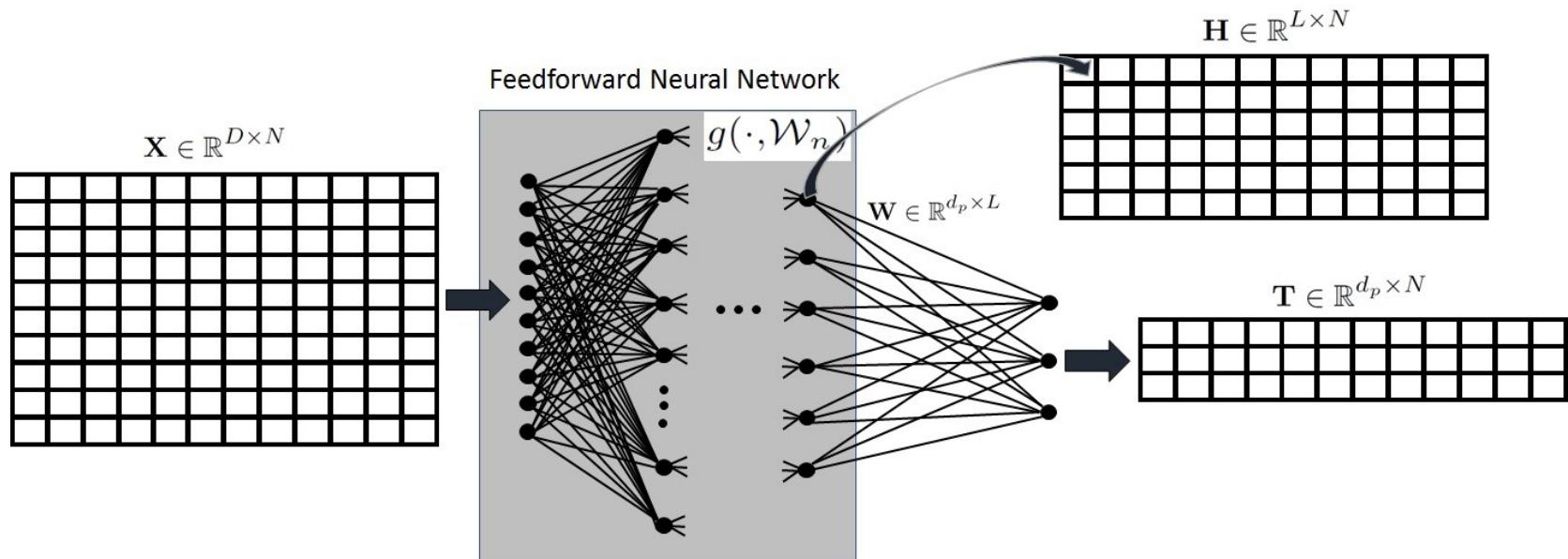
using data-independent targets \mathbf{T} !

Benefits

- › Much stable and fast solutions can be calculated
- › Approximation schemes are readily available and extremely efficient/effective
- › Incremental/Decremental solutions are available
- › Hierarchical (deep) models for class-specific learning are now possible

Discriminant Learning

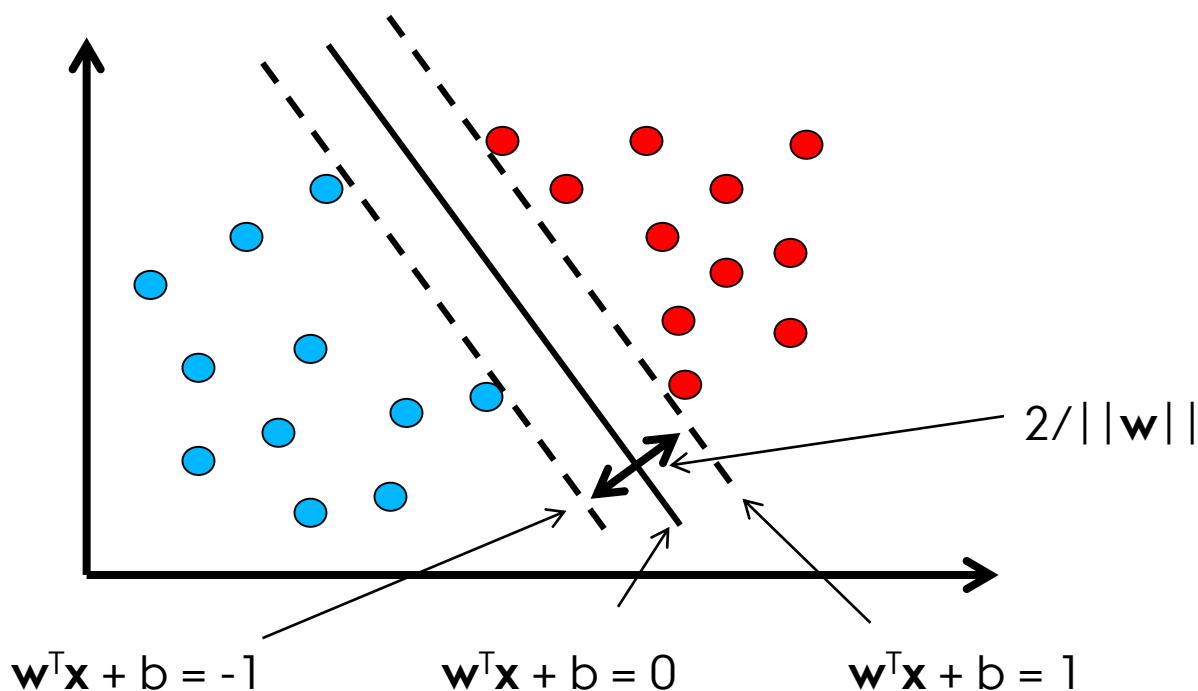
Deep Class-specific Discriminant Analysis model:



Classification

Max-margin based classification

- › Find the decision hyperplane discriminating the classes with maximum margin
- › Theoretical guarantees for generalization error
- › One (global) solution



Classification

- Binary classification:

$$\min_{\mathbf{w}, b} \frac{1}{2} \mathbf{w}^T \mathbf{w} + c \sum_{i=1}^N \xi_i,$$

s.t.:

$$y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i, \quad \xi_i \geq 0, \quad i = 1, \dots, N$$

$y_i \in \{-1, 1\}$ are the binary labels.

- Multi-class classification:

$$\min_{\mathbf{w}_k, b_k} \sum_{k=1}^K \frac{1}{2} \mathbf{w}_k^T \mathbf{w}_k + c \sum_{i=1}^N \sum_{k \neq l_i} \xi_i^k$$

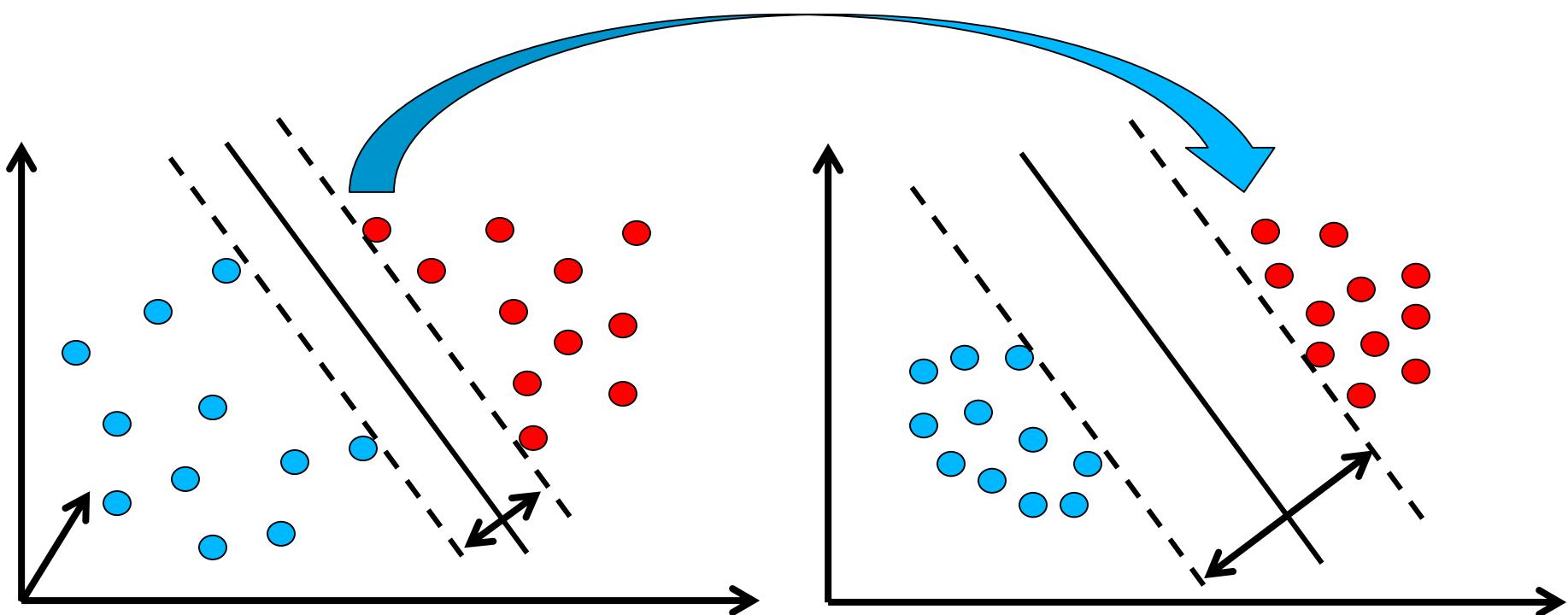
s.t.:

$$\mathbf{w}_{l_i}^T \mathbf{x}_i + b_{l_i} \geq \mathbf{w}_k^T \mathbf{x}_i + b_k + 2 - \xi_i^k, \quad \xi_i^k \geq 0, \quad i = 1, \dots, N, \quad k \neq l_i.$$

Classification

In order to increase class discrimination

- › Discriminant data mapping
- › Max-margin classification



Classification

Discriminant max-margin based classification

- › We proved that these two processing steps can be applied at once!
- › Max-margin to a discriminant (kernel) space

➤ Multi-class classification:

$$\text{s.t.: } \min_{\mathbf{w}_k, b_k} \sum_{k=1}^K \frac{1}{2} \mathbf{w}_k^T \mathbf{w}_k + c \sum_{i=1}^N \sum_{k \neq l_i} \xi_i^k + \sum_{k=1}^K \frac{\lambda}{2} \mathbf{w}_k^T \mathbf{S} \mathbf{w}_k$$

Discriminant Term

$$\mathbf{w}_{l_i}^T \mathbf{x}_i + b_{l_i} \geq \mathbf{w}_k^T \mathbf{x}_i + b_k + 2 - \xi_i^k, \quad \xi_i^k \geq 0, \quad i = 1, \dots, N, \quad k \neq l_i$$

➤ Joint optimization of K decision functions $\{\mathbf{w}_k, b_k\}$,
 $k=1, \dots, K$.

Probability-based Visual Saliency

Saliency model:

- › The probability of an image region (pixel/super-pixel) to be salient is obtained by:

$$\underset{\mathbf{P}(\mathbf{x})}{\operatorname{argmin}} \left(\sum_i (\mathbf{P}(\mathbf{x} = \mathbf{x}_i))^2 v_i + \frac{1}{2} \sum_{i,j} (\mathbf{P}(\mathbf{x} = \mathbf{x}_i) - \mathbf{P}(\mathbf{x} = \mathbf{x}_j))^2 w_{i,j} \right)$$

s.t. $\sum_i \mathbf{P}(\mathbf{x} = \mathbf{x}_i) = 1.$



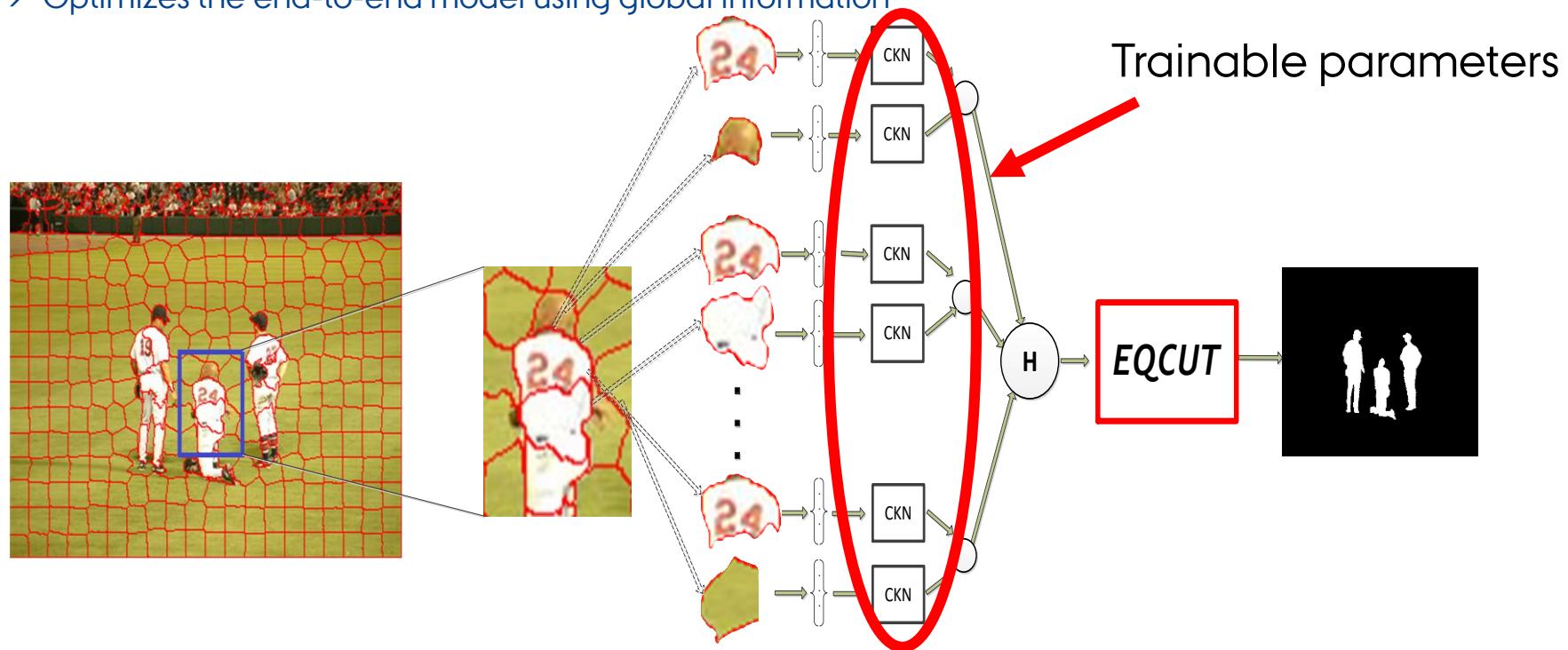
Properties

- › Global optimum solution
- › Generic framework for visual saliency:
 - › Diffusion methods are special cases of PSE
 - › PSE optimally refines the solution of QCut (SoA unsupervised saliency estimation method)
- › Now Saliency Estimation can also be modelled as an One-Class Classification Problem

Supervised Visual Saliency

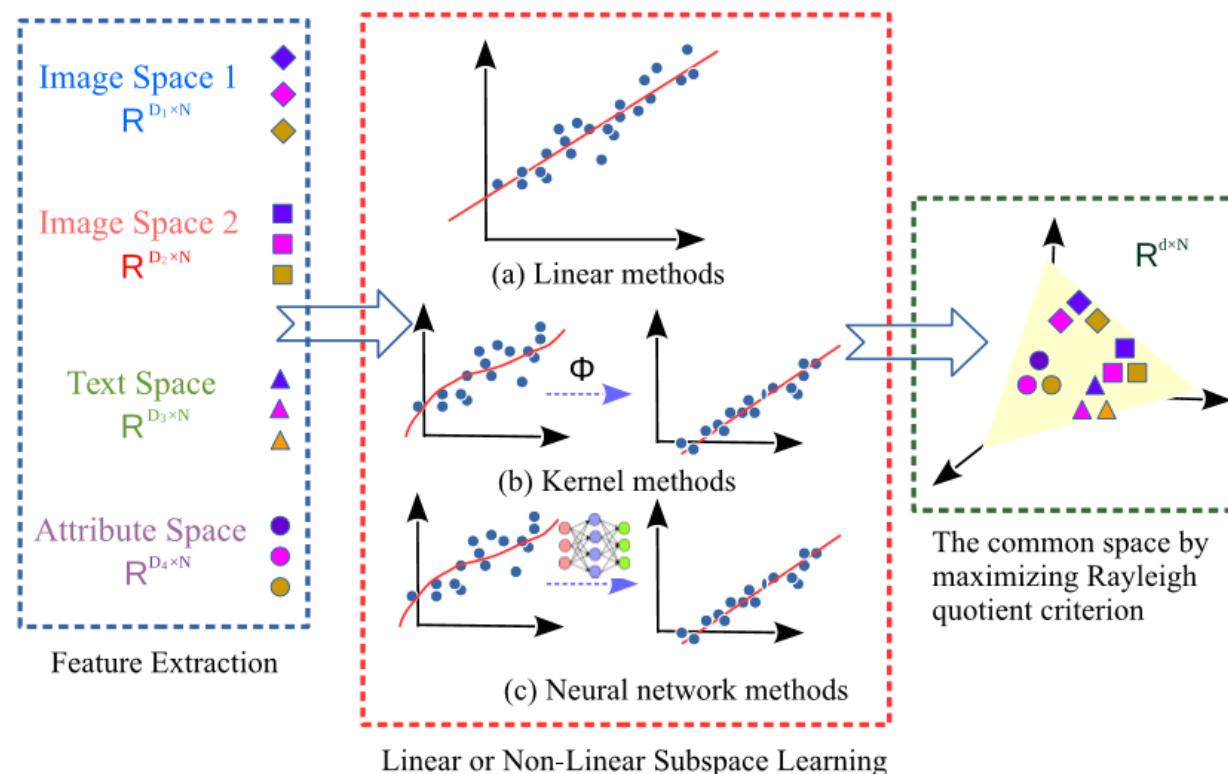
End-to-end learning for visual saliency that:

- › Exploits successive learnable feature transformations
- › Optimizes the end-to-end model using global information



Multi-view/modal Data Analysis

Generalized Multi-view Embedding



G. Cao, A. Iosifidis, K. Chen and M. Gabbouj, "Generalized Multi-view Embedding", IEEE Transactions on Cybernetics, 2018

Multi-view/modal Data Analysis

Generalized Multi-view Embedding

- › We showed that most multi-view embedding methods can be modeled using the Rayleigh Quotient

$$\mathcal{J} = \arg \max_{\mathbf{W}} \frac{\text{Tr}(\mathbf{W}^\top \mathbf{P} \mathbf{W})}{\text{Tr}(\mathbf{W}^\top \mathbf{Q} \mathbf{W})}$$

	P	Q
CCA	$\begin{bmatrix} 0 & \Sigma_{12} & \cdots & \Sigma_{1V} \\ \Sigma_{21} & 0 & \cdots & \Sigma_{2V} \\ \vdots & \vdots & \ddots & \vdots \\ \Sigma_{V1} & \Sigma_{V2} & \cdots & 0 \end{bmatrix}$	$\begin{bmatrix} \Sigma_{11} & 0 & \cdots & 0 \\ 0 & \Sigma_{22} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \Sigma_{VV} \end{bmatrix}$
PLS	$\begin{bmatrix} 0 & \Sigma_{12} & \cdots & \Sigma_{1V} \\ \Sigma_{21} & 0 & \cdots & \Sigma_{2V} \\ \vdots & \vdots & \ddots & \vdots \\ \Sigma_{V1} & \Sigma_{V2} & \cdots & 0 \end{bmatrix}$	$\begin{bmatrix} I & 0 & \cdots & 0 \\ 0 & I & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & I \end{bmatrix}$
LDA	$\begin{bmatrix} P_{11} & P_{12} & \cdots & P_{1V} \\ P_{21} & P_{22} & \cdots & P_{21} \\ \vdots & \vdots & \ddots & \vdots \\ P_{V1} & P_{V2} & \cdots & P_{VV} \end{bmatrix}$	$\begin{bmatrix} Q_{11} & 0 & \cdots & 0 \\ 0 & Q_{22} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & Q_{VV} \end{bmatrix}$

Multi-view/modal Data Analysis

Generalized Multi-view Embedding

- › We showed that most multi-view embedding methods can be modeled using the Rayleigh Quotient
- › Based on this observation, we proposed a new MvLDA criterion incorporating inter-view and intra-view variance criteria

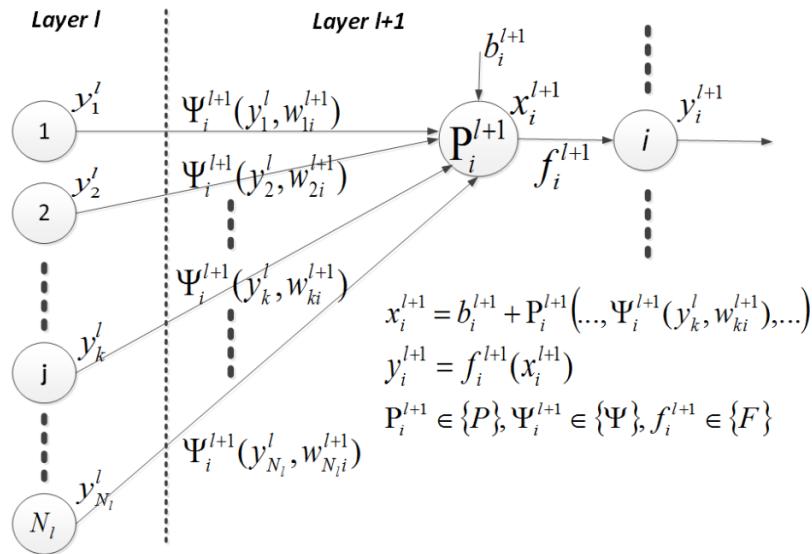
$$\begin{aligned}\mathbf{S}_W &= \sum_{i=1}^V \sum_{c=1}^C \mathbf{w}_i^\top \mathbf{x}_i \left(\mathbf{I} - \sum_{c=1}^C \frac{1}{N_c} \mathbf{e}_c \mathbf{e}_c^\top \right) \mathbf{x}_i^\top \mathbf{w}_i \\ &= \sum_{i=1}^V \sum_{j=1}^V \sum_{c=1}^C \mathbf{w}_i^\top \mathbf{Q}_{ii} \mathbf{w}_i\end{aligned}$$

$$\begin{aligned}\mathbf{S}_B &= \sum_{i=1}^V \sum_{j=1}^V \sum_{p=1}^C \sum_{q=1}^C \underset{p \neq q}{(\mathbf{m}_p^i - \mathbf{m}_q^i)(\mathbf{m}_p^i - \mathbf{m}_q^i)^\top} \\ &= \sum_{i=1}^V \sum_{j=1}^V \sum_{p=1}^C \sum_{q=1}^C \mathbf{w}_i^\top \mathbf{x}_i \mathbf{L}_B \mathbf{x}_j^\top \mathbf{w}_j \\ \mathbf{S}'_B &= \sum_{i=1}^V \sum_{j=1}^V \sum_{p=1}^C \sum_{q=1}^C \underset{p \neq q}{(\mathbf{m}_p^i - \mathbf{m}_q^i)(\mathbf{m}_p^j - \mathbf{m}_q^j)^\top} \\ &= \sum_{i=1}^V \sum_{j=1}^V \sum_{p=1}^C \sum_{q=1}^C \mathbf{w}_i^\top \mathbf{x}_i \mathbf{L}'_B \mathbf{x}_j^\top \mathbf{w}_j\end{aligned}$$

Neural Networks (Deep Learning)

Basic idea of neural network-based learning

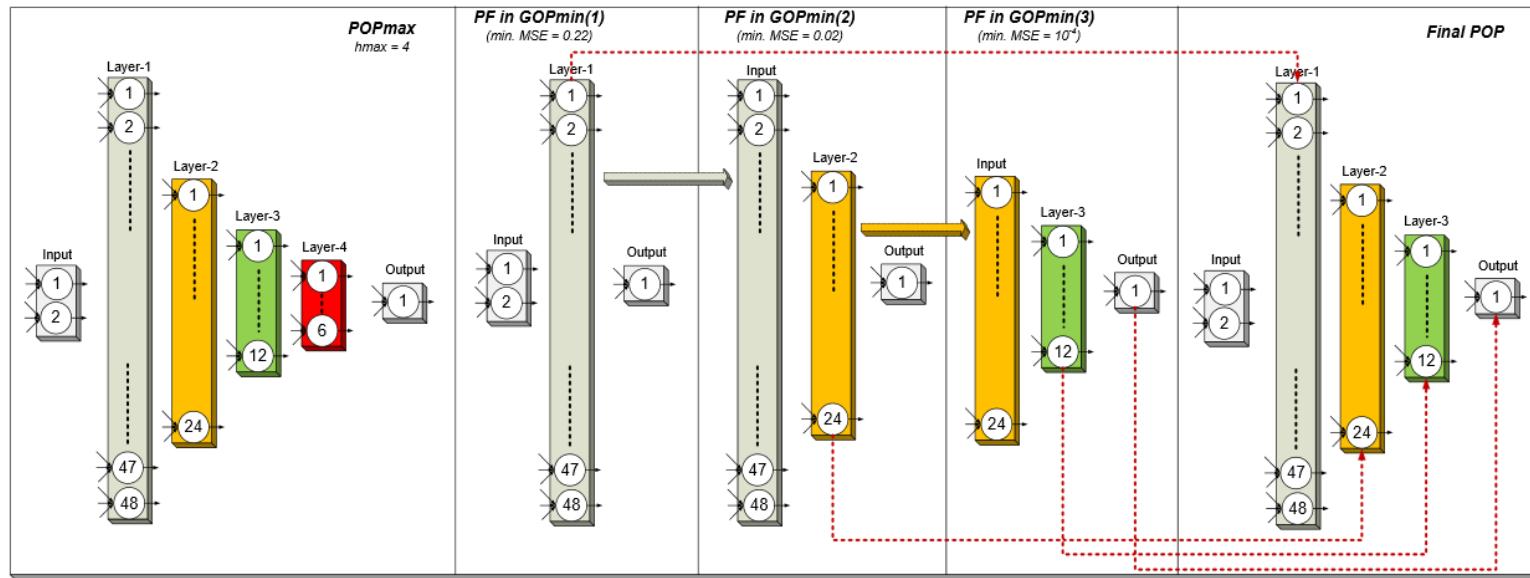
- › Data-driven fine-tuning the parameters of a network to regress inputs to targets using:
 - › User-defined number of layers
 - › User-defined activation functions (always the same for the entire network, except the last one)
 - › User-defined neuron pooling operator (always the same for all neurons)
 - › User-defined neuron nodal operator (always the same for all neurons)



Neural Networks (Deep Learning)

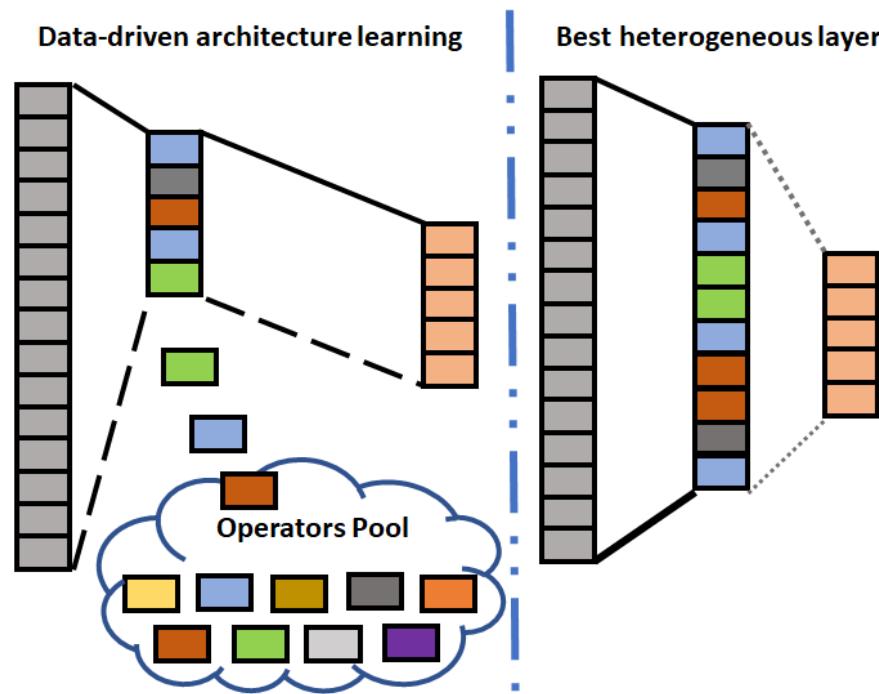
We proposed a Progressive Operational Feedforward Neural network learning approach

- › Data-driven network's architecture
- › Data-driven network's parameters tuning



Neural Networks (Deep Learning)

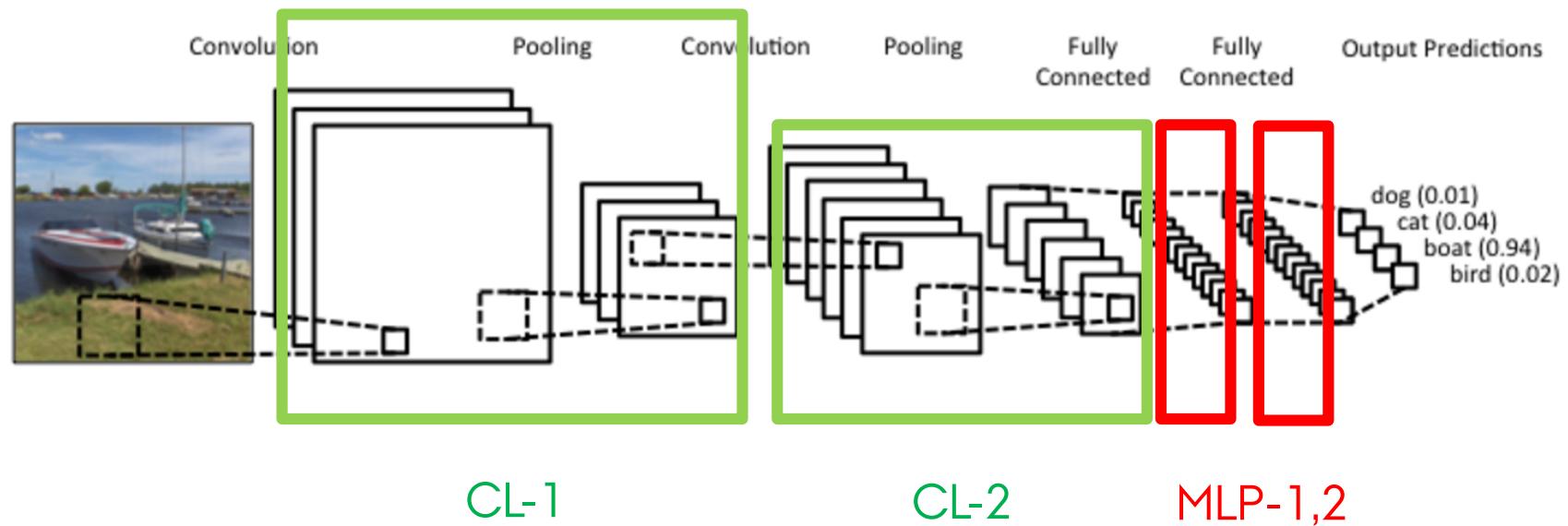
We proposed a data-driven neural network architecture learning scheme
› Fully heterogeneous layers



Convolutional Neural Networks (Deep Learning)

CNN architecture:

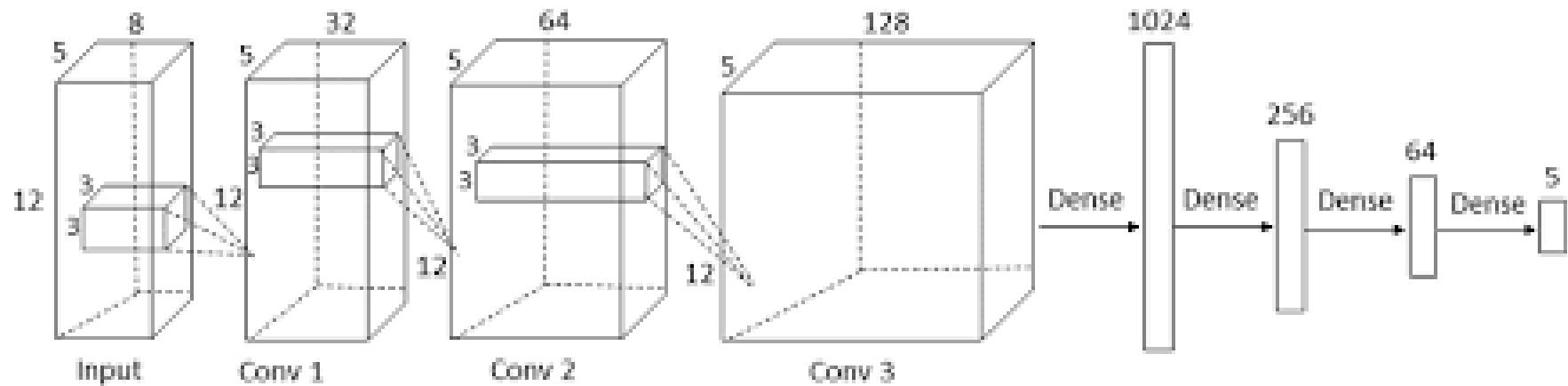
- › Convolutional layers
- › Multilayer Perceptron (vector) layers



Convolutional Neural Networks (Deep Learning)

CNN architecture:

- › Convolutional layers
- › Multilayer Perceptron (vector) layers

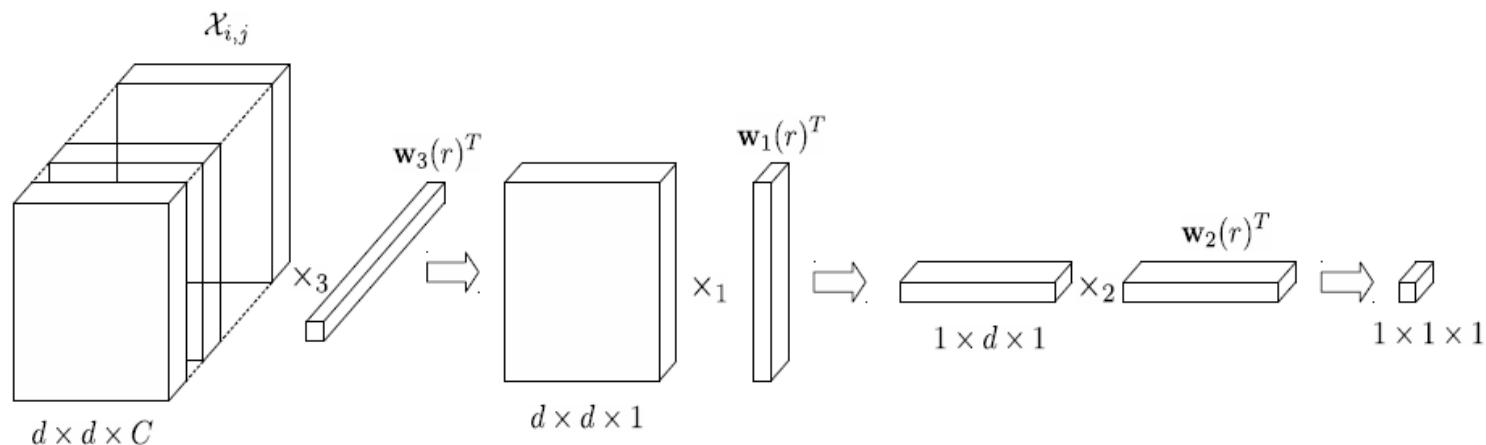


Real CNN architecture: CLs are tensors!

Convolutional Neural Networks (Deep Learning)

We proposed a tensor-based CNN filters modeling that can lead to:

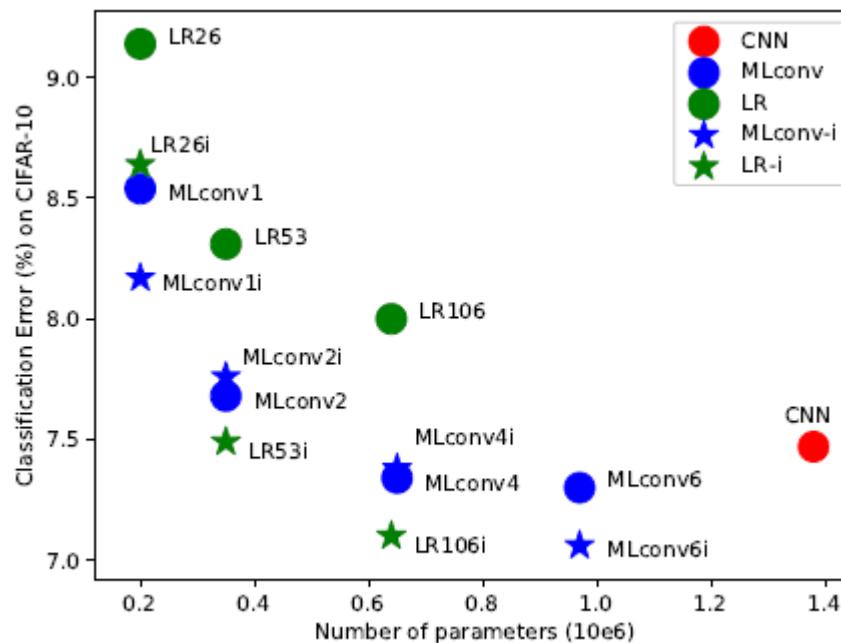
- › Lower number of network parameters (reduced memory footprint)
- › Faster classification (reduced computational cost)



Convolutional Neural Networks (Deep Learning)

We proposed a tensor-based CNN filters modeling that can lead to:

- › Lower number of network parameters (reduced memory footprint)
- › Faster classification (reduced computational cost)

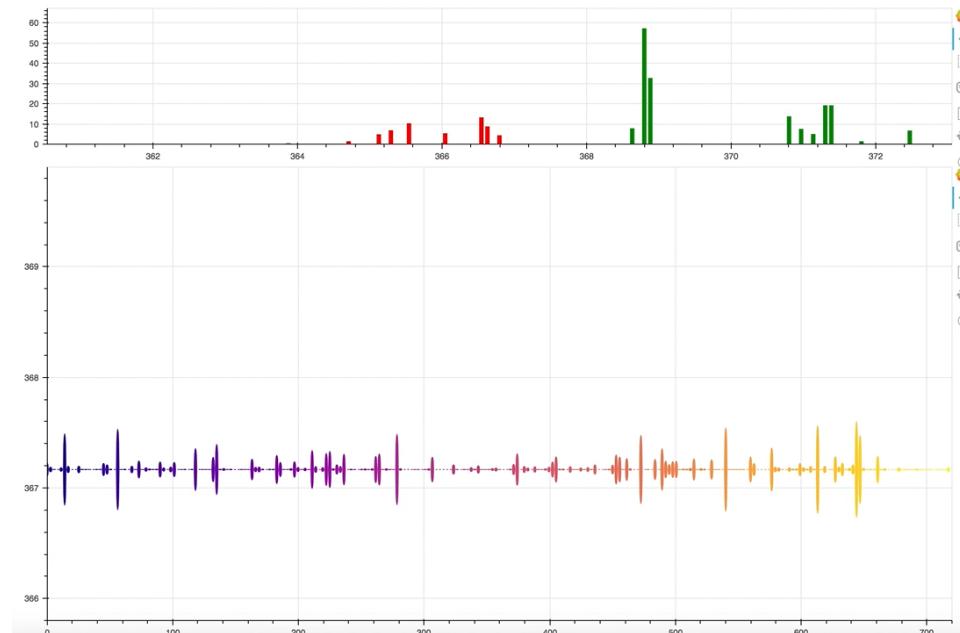


D.T. Thanh, A. Iosifidis and M. Gabbouj, "Improving Efficiency in Convolutional Neural Network with Multilinear Filters", under minor revisions in Neural Networks (arXiv 2017)

Extension to applications of other domains

Other applications

Stock prediction in financial markets



- A. Tsantekidis, N. Passalis, A. Tefas, J. Kanniainen, M. Gabbouj and A. Iosifidis, “Using Deep Learning to Detect Price Change Indications in Financial Markets”, EUSIPCO, 2017
- N. Passalis, A. Tsantekidis, A. Tefas, J. Kanniainen, M. Gabbouj and A. Iosifidis, “Time-series Classification using Neural Bag-of-Features”, EUSIPCO, 2017

Other applications

Stock prediction in financial markets

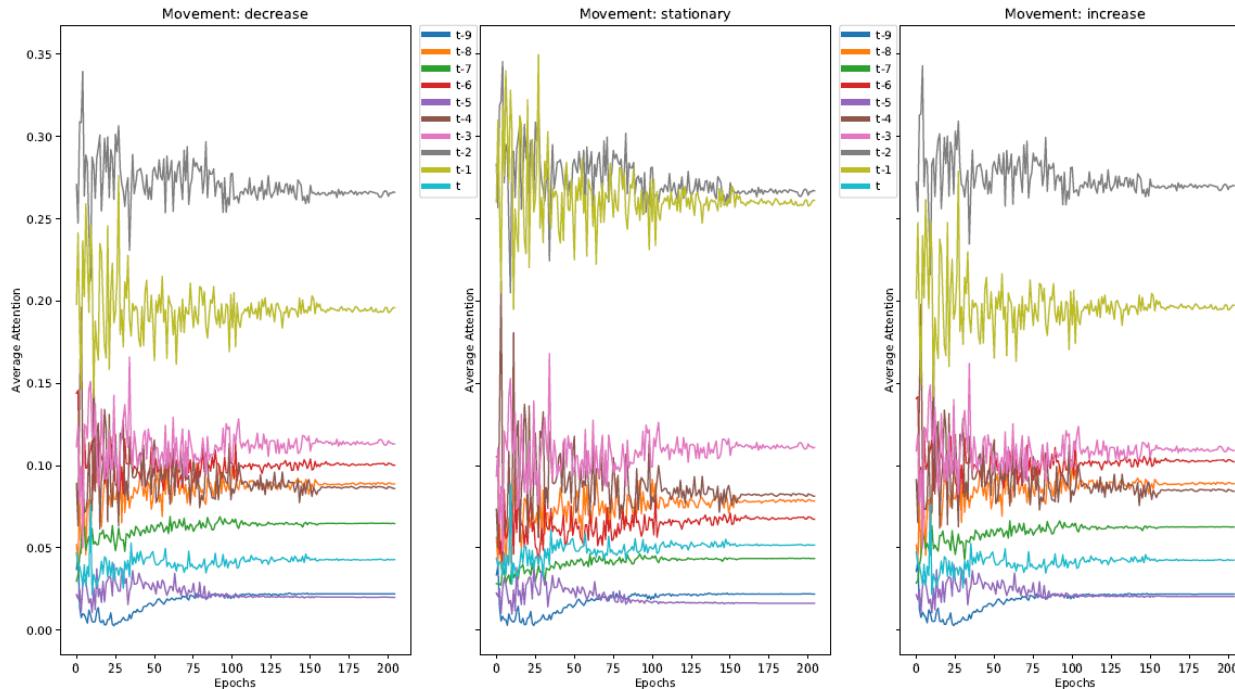
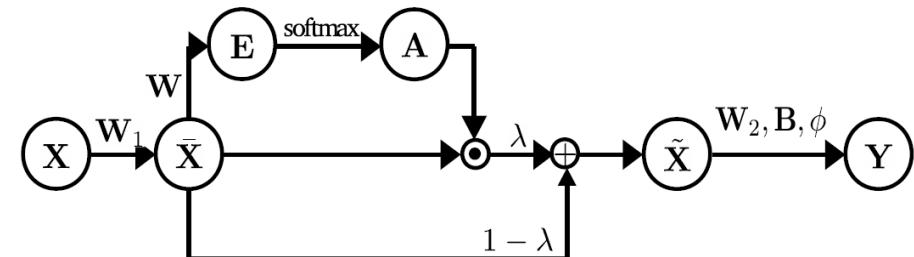


Fig. 3. Average attention of 10 temporal instances during training in 3 types of movement: decrease, stationary, increase. Values taken from configuration A(TABL) in Setup2, horizon $H = 10$

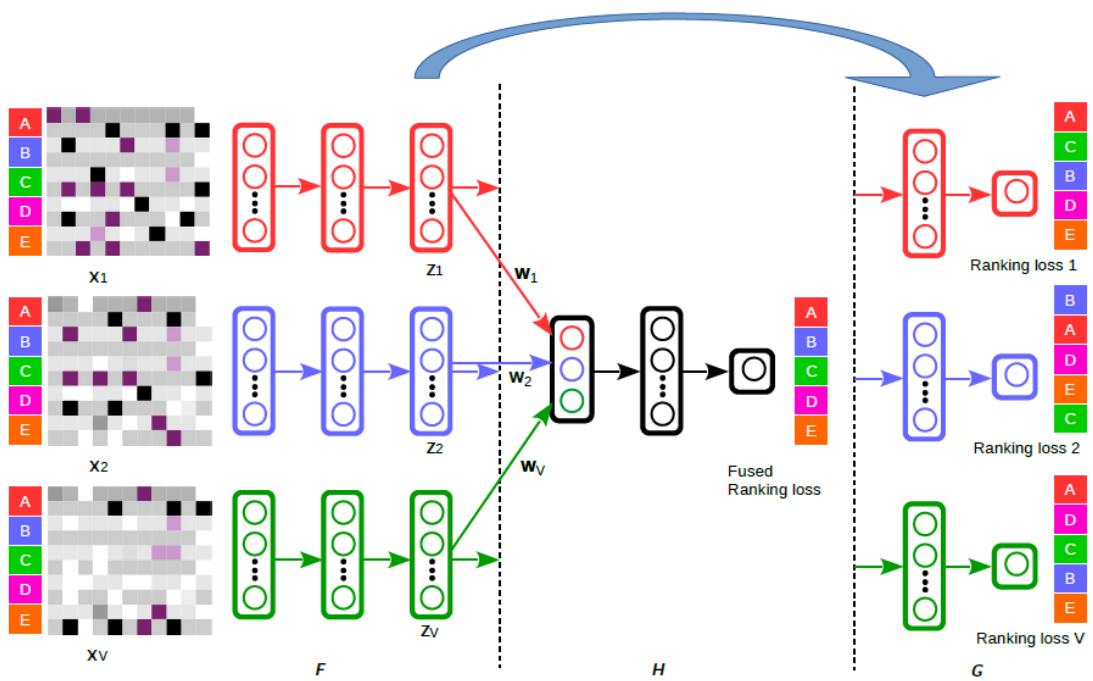


Models	Accuracy %	Precision %	Recall %	F1 %
Prediction Horizon $H = 10$				
SVM[48]	-	39.62	44.92	35.88
MLP[48]	-	47.81	60.78	48.27
CNN[47]	-	50.98	65.54	55.21
LSTM[48]	-	60.77	75.92	66.33
A(BL)	29.21	44.08	48.14	29.47
A(TABL)	70.13	56.28	58.26	56.03
B(BL)	78.37	67.73	68.89	67.71
B(TABL)	78.91	68.04	71.21	69.20
C(BL)	82.52	73.89	76.22	75.01
C(TABL)	84.70	76.95	78.44	77.63
Prediction Horizon $H = 20$				
SVM[48]	-	45.08	47.77	43.20
MLP[48]	-	51.33	65.20	51.12
CNN[47]	-	54.79	67.38	59.17
LSTM[48]	-	59.60	70.52	62.37
A(BL)	42.01	47.71	45.38	38.61
A(TABL)	62.54	52.36	50.96	50.69
B(BL)	70.33	62.97	60.64	61.02
B(TABL)	70.80	63.14	62.25	62.22
C(BL)	72.05	65.04	65.23	64.89
C(TABL)	73.74	67.18	66.94	66.93
Prediction Horizon $H = 50$				
SVM[48]	-	46.05	60.30	49.42
MLP[48]	-	55.21	67.14	55.95
CNN[47]	-	55.58	67.12	59.44
LSTM[48]	-	60.03	68.58	61.43
A(BL)	51.92	51.59	50.35	49.58
A(TABL)	60.15	59.05	55.71	55.87
B(BL)	72.16	71.28	68.69	69.40
B(TABL)	75.58	74.58	73.09	73.64
C(BL)	78.96	77.85	77.04	77.40
C(TABL)	79.87	79.05	77.04	78.44

D.T. Tran, A. Iosifidis, J. Kanniainen and M. Gabbouj, "Temporal Attention augmented Bilinear Network for Financial Time-Series Data Analysis", under minor revisions in IEEE Transactions on Neural Networks and Learning Systems (arXiv 2017)

Other applications

Multi-faceted and multi-objective data ranking



Safety/hazard Assessment

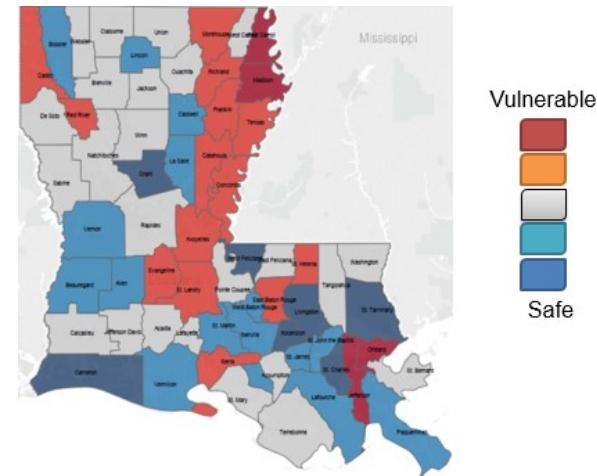
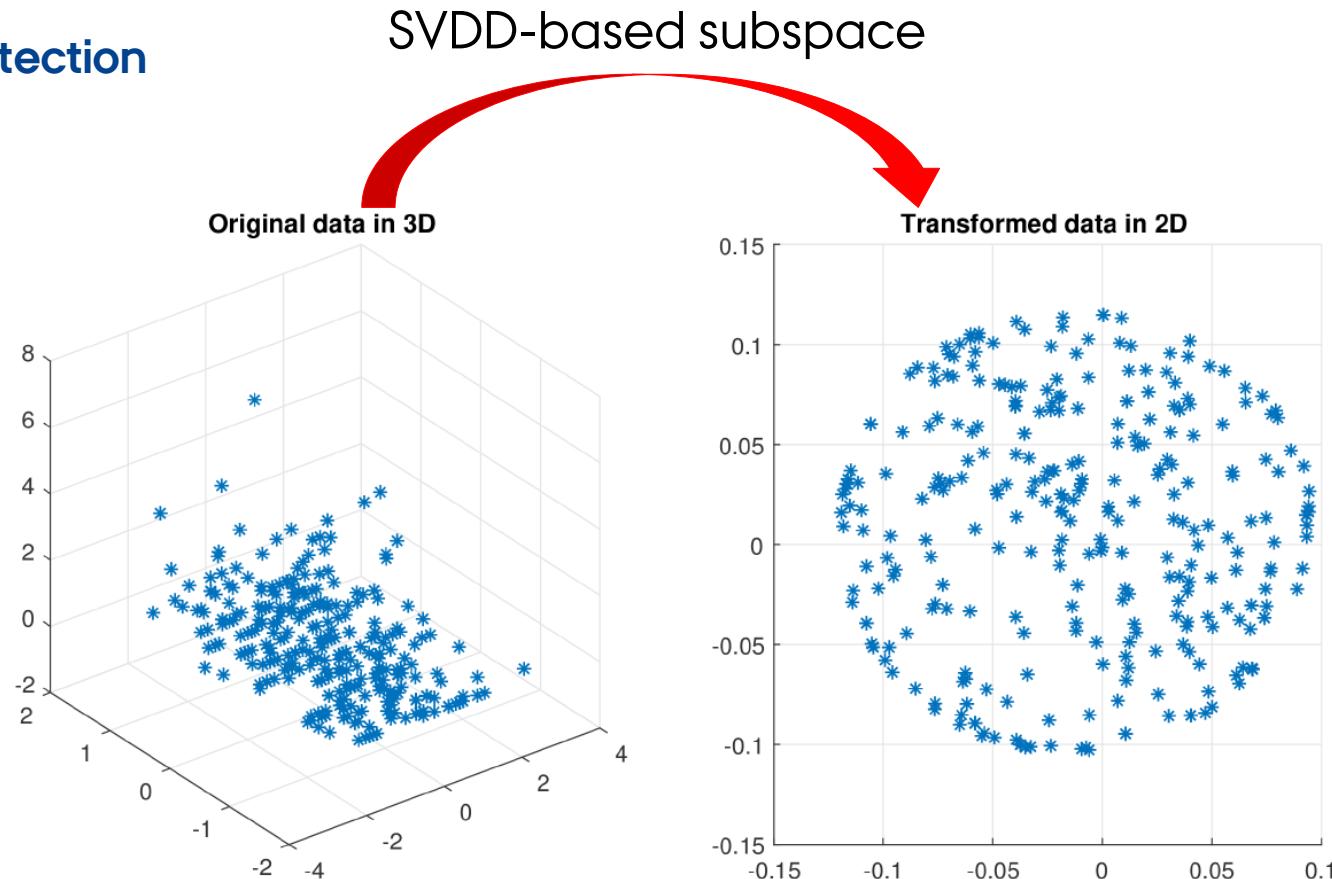


TABLE 1: Average Prediction Results (%) on 3 University Ranking Datasets in 2015.

Methods	Kendal's tau	Accuracy
Best Single View	65.38	-
Feature Concat	35.10	-
LMvCCA [5]	86.04	94.49
LMvMDA [5]	87.00	94.97
MvDA [6]	85.81	94.34
SmVR [8]	80.75	-
DMvCCA [5]	70.07	93.20
DMvMDA [5]	70.81	94.75
MvCCAE (<i>ours</i>)	75.94	94.01
MvMdae (<i>ours</i>)	81.04	94.85
DMvDR (<i>ours</i>)	89.28	95.30

Other applications

Novelty detection





Thank you for your attention!

Q/A