

Distribution-based Semi-Supervised Learning for Activity Recognition (AAAI'19)

Hangwei Qian, Sinno Jialin Pan, Chunyan Miao

Nanyang Technological University, Singapore

January 30, 2019



Outline

- 1 Problem Overview
- 2 Kernel Mean Embeddings for Feature Extraction
- 3 The Proposed DSSL for Semi-Supervised Learning
- 4 Experiments
- 5 Conclusion

Human Activity Recognition

Tremendous applications:

- elderly assistant
- healthcare
- fitness coaching
- smart building
- gaming



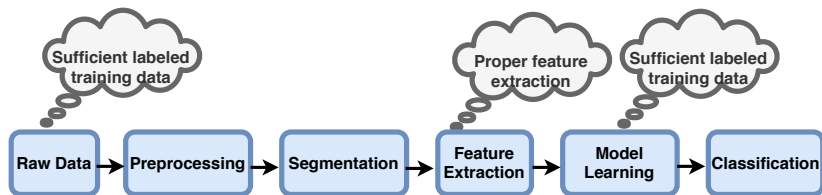
Human Activity Recognition

A multi-class classification problem

- Input: wearable onbody sensor data
- Output: activity labels



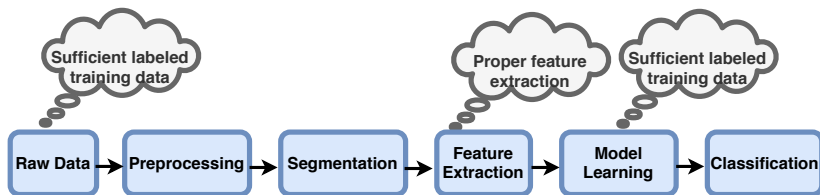
Problem Overview



Two key prerequisites:

- 1 expressive feature extraction → discriminate activities
- 2 sufficient labeled training data → build a precise model

Problem Overview



Two key prerequisites:

- 1 expressive feature extraction → discriminate activities → dependent on domain knowledge
- 2 sufficient labeled training data → build a precise model → require a huge amount of human annotation effort

Motivation

- 1 **Can we extract as many discriminative features as possible, in an automatic fashion?**
 - kernel mean embedding of distributions, with NO information loss
 - novel supervised methods **SMM**_{AR} and **R-SMM**_{AR} [7]¹
- 2 **Can we utilize labeled data as few as possible to alleviate human annotation effort?**
 - Distribution-based Semi-Supervised Learning (DSSL)

¹Hangwei Qian, Sinno Jialin Pan, and Chunyan Miao. **Sensor-based activity recognition via learning from distributions**. In AAAI'18 (oral).

Existing Feature Extraction Methods

Frame-level → vectorial-based

- Manual feature engineering, statistics of each frame

	time_1	time_2	time_3	time_4	time_5
feature_1	0.9134	0.2785	0.9649	0.9572	0.8147
feature_2	0.9058	0.6324	0.5469	0.1576	0.4854
feature_3	0.127	0.0975	0.9575	0.9706	0.8003

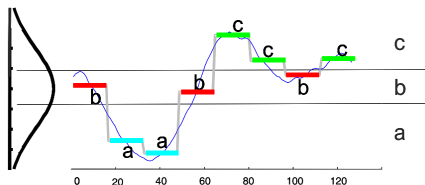
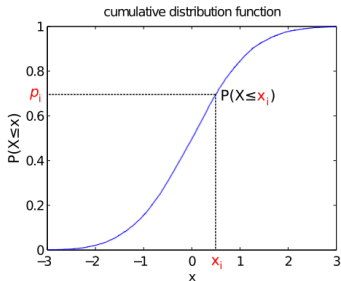
Existing Feature Extraction Methods

Frame-level \rightarrow vectorial-based

- Manual feature engineering, statistics of each frame

Segment-level \rightarrow matrix-based

- Statistical, i.e., moments of each segment
- Structural
 - The ECDF method [4]
 - The SAX method [3, 6]



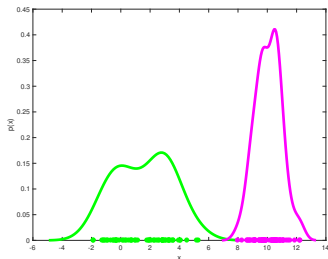
Existing Semi-Supervised Methods

- LapSVM [1]: manifold learning
- ∇ T SVM [2]: transductive
- SSKLR [5]: kernel logistic regression with Expectation-Maximization algorithm
- GLSVM [8]: multi-graph based

Outline

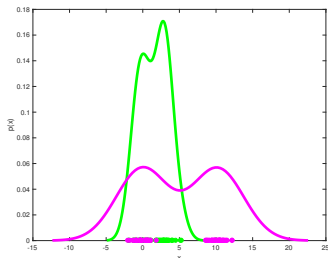
- 1 Problem Overview
- 2 Kernel Mean Embeddings for Feature Extraction**
- 3 The Proposed DSSL for Semi-Supervised Learning
- 4 Experiments
- 5 Conclusion

Intuition of Kernel Mean Embedding

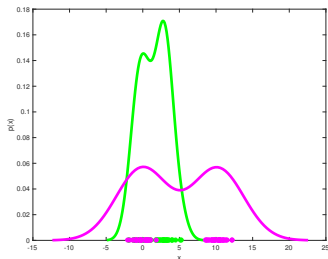


$(\mathbb{E}[x])$ as features

problem: many distributions have the same mean!



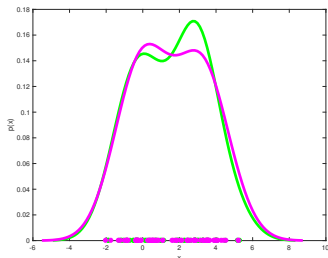
Intuition of Kernel Mean Embedding



$(\mathbb{E}[x])$ as features

problem: many distributions have the same mean!

$\left(\begin{array}{c} \mathbb{E}[x] \\ \mathbb{E}[x^2] \end{array} \right)$ as features



problem: many distributions have the same mean and variance!

Intuition of Kernel Mean Embedding

$(\mathbb{E}[x])$ as features

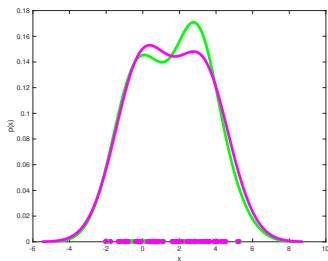
problem: many distributions have the same mean!

$\left(\begin{array}{c} \mathbb{E}[x] \\ \mathbb{E}[x^2] \end{array} \right)$ as features

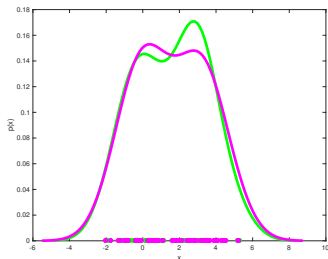
problem: many distributions have the same mean and variance!

$\left(\begin{array}{c} \mathbb{E}[x] \\ \mathbb{E}[x^2] \\ \mathbb{E}[x^3] \end{array} \right)$ as features

problem: many distributions still have the same first 3 moments!



Intuition of Kernel Mean Embedding



$$\mu[P_x] = \begin{pmatrix} \mathbb{E}[x] \\ \mathbb{E}[x^2] \\ \mathbb{E}[x^3] \\ \dots \\ \dots \end{pmatrix}$$

The **infinite dimensional features** should be able to discriminate different distributions!

Kernel Mean Embedding of Distributions

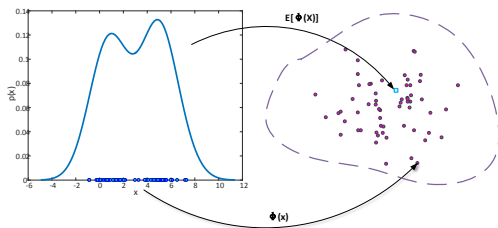


Figure 1:
Illustrations of kernel mean embeddings of a distribution and embeddings of empirical examples

$$\mu[P_X] = E_X[k(\cdot, x)] \in \mathcal{H} \quad (1)$$

$$\mu[X] = \frac{1}{m} \sum_{i=1}^m k(\cdot, x_i) \in \mathcal{H} \quad (2)$$

Here $X = \{x_1, \dots, x_m\} \stackrel{i.i.d.}{\sim} P_X$, \mathcal{H} is the RKHS associated with k .

Outline

- 1 Problem Overview
- 2 Kernel Mean Embeddings for Feature Extraction
- 3 The Proposed DSSL for Semi-Supervised Learning**
- 4 Experiments
- 5 Conclusion

Contribution

DSSL: Distribution-based Semi-Supervised Learning

- 1 All orders of statistical moments features are extracted implicitly and automatically
- 2 DSSL relaxes SMM_{AR} 's full supervision assumption, and exploit unlabeled instances to learn an underlying data structure
- 3 DSSL is the first attempt on semi-supervised learning with distributions, with rigorous theoretical proofs provided.
- 4 Extensive experiments to show the efficacy of DSSL.

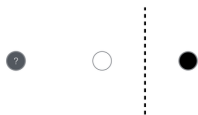
Intuition of DSSL

- Label annotation is time-consuming
- Unlabeled data is abundant and informative

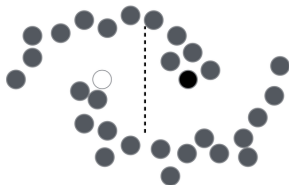


Intuition of DSSL

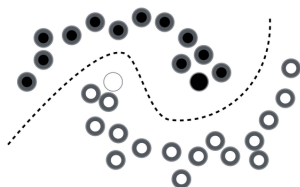
- Label annotation is time-consuming
- Unlabeled data is abundant and informative



what if unlabeled data is available?



Intuition of DSSL



Intuition: unlabeled data sheds light on the underlying manifolds of **data space**

Difficulty:

- Classical setting: $x \in \mathbb{R}^n$
- Our setting: $\mu[X] \in \mathcal{H}$

Distribution-based SSL: Main idea

- 1 map the activity segments into a RKHS \rightarrow sufficient features
- 2 wrap the RKHS space to reflect the manifold of the data \rightarrow modify the similarity measure $\langle f, g \rangle_{\tilde{\mathcal{H}}} \triangleq \langle f, g \rangle_{\tilde{\mathcal{H}}} + F(f, g)$
 - data within a manifold (**instead of closer Euclidean distance**) \rightarrow more similar
 - data with different labels \rightarrow less similar

Challenges

$$\langle f, g \rangle_{\check{\mathcal{H}}} \triangleq \langle f, g \rangle_{\tilde{\mathcal{H}}} + F(f, g) \quad (3)$$

$$f^* = \arg \min_{f \in \check{\mathcal{H}}} \frac{1}{I} \sum_{i=1}^I \ell([\mu_{\mathbb{P}_i}]_{\tilde{\mathcal{H}}}, y_i, [f]_{\check{\mathcal{H}}}) + \|f\|_{\check{\mathcal{H}}}^2, \quad (4)$$

- 1 How to construct the data-dependent kernel by incorporating unlabeled training data?
- 2 Is the new space valid? Since a RKHS is defined by inner product.
- 3 How to calculate the loss function given two items are not in the same space?

Challenge 1/3 Construction of kernel

$$\langle \mathbf{f}, \mathbf{g} \rangle_{\tilde{\mathcal{H}}} \stackrel{\Delta}{=} \langle \mathbf{f}, \mathbf{g} \rangle_{\tilde{\mathcal{H}}} + \langle S\mathbf{f}, S\mathbf{g} \rangle_{\mathcal{V}}, \quad (5)$$

where S is a bounded linear operator.

Denote $\mathbf{f}(\boldsymbol{\mu}) = (f(\boldsymbol{\mu}_{\mathbb{P}_1}), \dots, f(\boldsymbol{\mu}_{\mathbb{P}_n}))$,

$$\langle S\mathbf{f}, S\mathbf{f} \rangle_{\mathcal{V}} = \mathbf{f}(\boldsymbol{\mu}) M \mathbf{f}(\boldsymbol{\mu})^{\top} \quad (6)$$

In our case, $M = rL^2$, where L is the Laplacian matrix

Challenge 2/3 Validity of the new space

Theorem 1

$\tilde{\mathcal{H}}$ is a valid RKHS.

Challenge 3/3 Loss function calculation

$$f^* = \arg \min_{f \in \check{\mathcal{H}}} \frac{1}{I} \sum_{i=1}^I \ell([\mu_{\mathbb{P}_i}]_{\check{\mathcal{H}}}, y_i, [f]_{\check{\mathcal{H}}}) + \|f\|_{\check{\mathcal{H}}}^2, \quad (7)$$

Proposition 1

$$\check{\mathcal{H}} = \tilde{\mathcal{H}}.$$

Proposition 2

$$\check{K} = (I + \check{K}M)^{-1} \check{K},$$

where \check{K} with $\check{K}_{ij} = \check{k}(\mu_{\mathbb{P}_i}, \mu_{\mathbb{P}_j})$ is the kernel matrix for $\tilde{\mathcal{H}}$ on $\mu_{\mathbb{P}_i}$'s, and \check{K} is the kernel matrix in the altered space $\check{\mathcal{H}}$.

Outline

- 1 Problem Overview
- 2 Kernel Mean Embeddings for Feature Extraction
- 3 The Proposed DSSL for Semi-Supervised Learning
- 4 Experiments**
- 5 Conclusion

Experimental Setup

- labeled training set, unlabeled training set and test set:
0.02:0.1:0.88
- evaluation: micro-F1 (miF), macro-F1 (maF)

Table 1: Statistics of datasets used in experiments.

Datasets	# Sample	# Instances per sample	# Feature	# Class
Skoda	1,447	68	60	10
HCI	264	602	48	5
WISDM	389	705	6	6

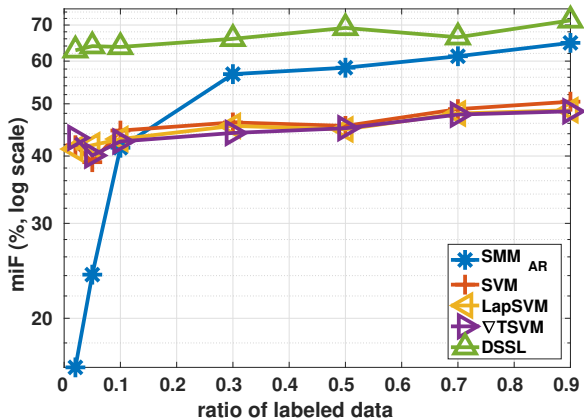
Experimental Results

Table 2: Experimental results on 3 activity datasets (unit: %).

Methods		Skoda		HCI		WISDM	
		miF	maF	miF	maF	miF	maF
Vectorial-based supervised	SVMs	85.7±1.8	42.5±0.9	69.7±9.6	69.6±9.4	41.5±5.2	39.6±6.8
	SAX_3	39.6±6.3	18.7±2.9	36.0±3.0	34.7±2.5	34.6±1.4	30.6±1.2
	SAX_6	37.2±6.1	18.6±2.8	39.7±7.3	38.4±7.9	34.9±3.0	30.5±5.0
	SAX_9	40.3±6.5	19.9±3.2	39.8±8.7	37.0±9.2	33.6±2.9	28.8±5.8
	ECDF_5	84.2±2.1	41.6±1.0	67.7±10.1	67.6±9.1	42.1±6.3	40.5±7.7
	ECDF_15	79.8±1.5	39.2±0.7	68.4±10.4	68.5±9.6	39.4±3.3	36.2±5.7
	ECDF_30	72.6±1.2	35.4±0.3	68.6±11.1	68.7±10.5	37.7±2.5	32.6±4.9
	ECDF_45	65.7±2.5	31.5±1.3	68.6±11.4	68.6±10.8	36.4±1.4	31.3±3.6
Vectorial-based semi-supervised	LapSVM	89.7±2.1	44.6±1.2	76.1±4.8	76.3±4.7	40.1±3.8	34.5±3.5
	▽TSVM	85.9±2.7	84.8±2.8	75.4±11.5	75.5±11.2	41.3±5.6	39.4±6.9
	SSKLR	25.4±19.3	12.1±2.5	24.2±17.2	18.1±10.1	24.6±17.0	17.3±9.9
	GLSVM	89.7±2.1	44.5±1.2	75.7±5.8	75.7±5.7	40.4±3.8	33.9±4.0
Distribution-based supervised	SMM _{AR}	93.2±0.9	93.1±1.0	82.2±13.4	78.9±18.4	20.5±3.3	11.7±3.9
Distribution-based semi-supervised	DSSL	98.8±0.5	98.8±0.5	99.9±0.2	99.9±0.2	56.5±5.1	55.6±5.0

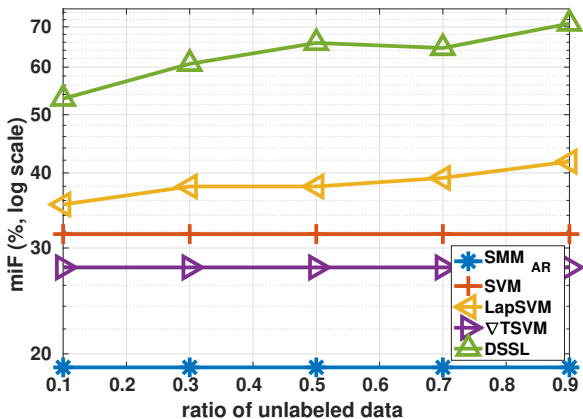
Experiments Analysis (1/3)

Varying ratios of labeled data



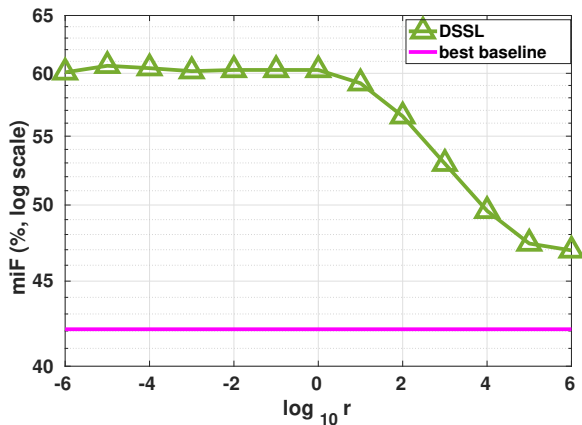
Experiments Analysis (2/3)

Varying ratios of unlabeled data



Experiments Analysis (3/3)

Impact of parameter r to the performance



Outline

- 1 Problem Overview
- 2 Kernel Mean Embeddings for Feature Extraction
- 3 The Proposed DSSL for Semi-Supervised Learning
- 4 Experiments
- 5 Conclusion**

Conclusion

We propose a novel method, i.e., Distribution-based Semi-Supervised Learning (DSSL) for human activity recognition

- 1 All orders of statistical moments features are extracted implicitly and automatically
- 2 DSSL relaxes SMM_{AR} 's full supervision assumption, and exploit unlabeled instances to learn an underlying data structure
- 3 DSSL is the first attempt on semi-supervised learning with distributions, with rigorous theoretical proofs provided.
- 4 Extensive experiments to show the efficacy of DSSL.

Questions?



More info in <http://hangwei12358.github.io/>

References

- [1] Mikhail Belkin, Partha Niyogi, and Vikas Sindhwani. “Manifold Regularization: A Geometric Framework for Learning from Labeled and Unlabeled Examples”. In: *Journal of Machine Learning Research* 7 (2006), pp. 2399–2434. URL: <http://www.jmlr.org/papers/v7/belkin06a.html>.
- [2] Olivier Chapelle and Alexander Zien. “Semi-Supervised Classification by Low Density Separation”. In: *AISTATS*. 2005.
- [3] Nils Y. Hammerla et al. “On preserving statistical characteristics of accelerometry data using their empirical cumulative distribution”. In: *ISWC*. 2013, pp. 65–68.
- [4] Jessica Lin et al. “Experiencing SAX: a novel symbolic representation of time series”. In: *Data Min. Knowl. Discov.* 15.2 (2007), pp. 107–144.
- [5] Ryunosuke Matsushige, Koh Kakusho, and Takeshi Okadome. “Semi-supervised learning based activity recognition from sensor data”. In: *GCCE*. 2015, pp. 106–107.
- [6] Thomas Plötz, Nils Y. Hammerla, and Patrick Olivier. “Feature Learning for Activity Recognition in Ubiquitous Computing”. In: *IJCAI*. 2011, pp. 1729–1734.
- [7] Hangwei Qian, Sinno Jialin Pan, and Chunyan Miao. “Sensor-Based Activity Recognition via Learning From Distributions”. In: *AAAI*. AAAI Press, 2018.
- [8] Maja Stikic, Diane Larlus, and Bernt Schiele. “Multi-graph Based Semi-supervised Learning for Activity Recognition”. In: *ISWC*. 2009, pp. 85–92.

Kernel Mean Embeddings of Distributions

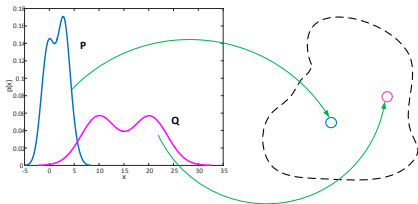


Figure 2: Illustration of the kernel mean embedding of two different distributions

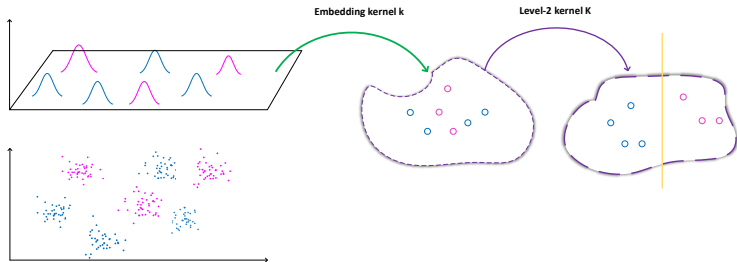
Injectivity[smola2007hilbert]

A universal kernel k can promise an injective mean map
 $\mu : P_X \rightarrow \mu[P_X]$.

SMM_{AR} Framework

$$\langle \hat{\mu}_{\mathbb{P}_x}, \hat{\mu}_{\mathbb{P}_z} \rangle = \tilde{k}(\hat{\mu}_{\mathbb{P}_x}, \hat{\mu}_{\mathbb{P}_z}) = \frac{1}{n_x \times n_z} \sum_{i=1}^{n_x} \sum_{j=1}^{n_z} k(\mathbf{x}_i, \mathbf{z}_j), \quad (8)$$

$$\tilde{k}(\mu_{\mathbb{P}_x}, \mu_{\mathbb{P}_z}) = \langle \psi(\mu_{\mathbb{P}_x}), \psi(\mu_{\mathbb{P}_z}) \rangle \quad (9)$$



Problem Formulation of SMM_{AR}

- Training set: $\{(P_i, y_i)\}, i \in \{1, \dots, N\}, x_i \sim P_i, x_i = \{x_{i1}, \dots, x_{im_i}\}, y_i \in \{1, \dots, L\}$
- Multi-class classifier $\rightarrow C_L^2$ binary classifiers
 $f, y = f(\phi(\mu_x)) + b$
- Primal Optimization problem:

$$\begin{aligned}
 \underset{f, b}{\operatorname{argmin}} \quad & \frac{1}{2} \|f\|_{\mathcal{H}}^2 + C \sum_{i=1}^N \xi_i \\
 \text{s.t.} \quad & y_i f(\phi(\mu_{x_i})) + b \\
 & y_i f(\phi(\mu_i)) \geq 1 - \xi_i, \forall i \\
 & \xi_i \geq 0, \forall i
 \end{aligned} \tag{10}$$