

# **Multiple Instance Learning with Bag Dissimilarities**

Veronika Cheplygina, David M.J. Tax, Marco Loog

## **Bag similarity network for deep multi-instance learning**

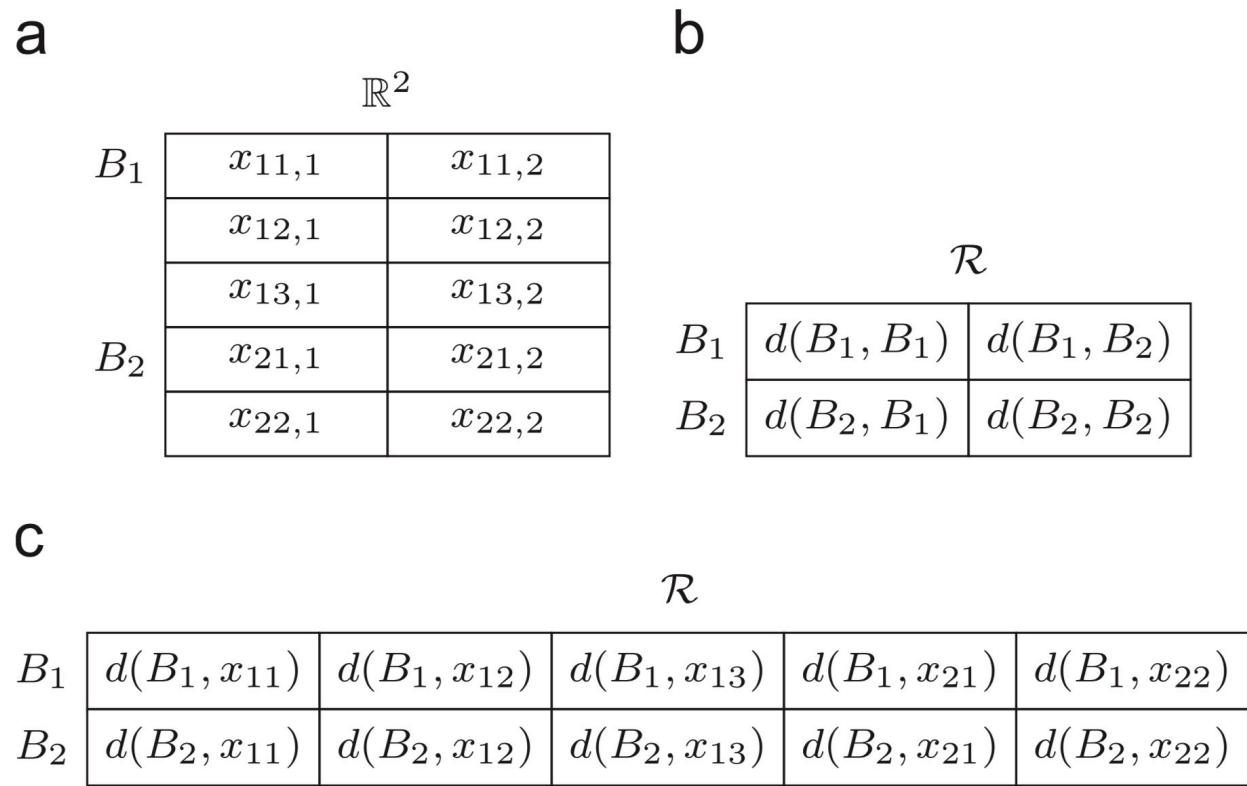
Xinggang Wang, Yongluan Yan, Peng Tang, Wenyu Liu, Xiaojie Guo

Presenter: Qingmei Wang

# Outline

- Background
- Bag Dissimilarity Representation
- Bag Similarity Network
- Summary

# Background



		$\mathbb{R}^2$	
		$B_1$	$B_2$
$B_1$	$x_{11,1}$	$x_{11,2}$	
	$x_{12,1}$	$x_{12,2}$	
	$x_{13,1}$	$x_{13,2}$	
$B_2$	$x_{21,1}$	$x_{21,2}$	
	$x_{22,1}$	$x_{22,2}$	

		$\mathcal{R}$	
		$B_1$	$B_2$
$B_1$	$d(B_1, B_1)$	$d(B_1, B_2)$	
$B_2$	$d(B_2, B_1)$	$d(B_2, B_2)$	

		$\mathcal{R}$				
		$B_1$	$B_2$	$x_{11}$	$x_{12}$	$x_{13}$
$B_1$	$d(B_1, x_{11})$	$d(B_1, x_{12})$	$d(B_1, x_{13})$	$d(B_1, x_{21})$	$d(B_1, x_{22})$	
$B_2$	$d(B_2, x_{11})$	$d(B_2, x_{12})$	$d(B_2, x_{13})$	$d(B_2, x_{21})$	$d(B_2, x_{22})$	

**Fig. 1.** Representations of a MIL problem with 2 bags and 2 features.  $B_1$  has 3 instances and  $B_2$  has 2 instances. The dimensionality of the original representation depends on the number of features, while in the dissimilarity representation, the dimensionality depends on the number of bags or instances. (a) Original MIL, (b) bag dissimilarity and (c) instance dissimilarity.

In MIL, each sample is in the form of a labeled bag that contains a set of instances associated with input features.

In MIL, the assumption is that a positive bag contains at least one positive instance, whereas a negative bag contains only negative instances.

The goal of MIL in a binary task is to train a classifier so that the labels of testing bags may be predicted.

# Background

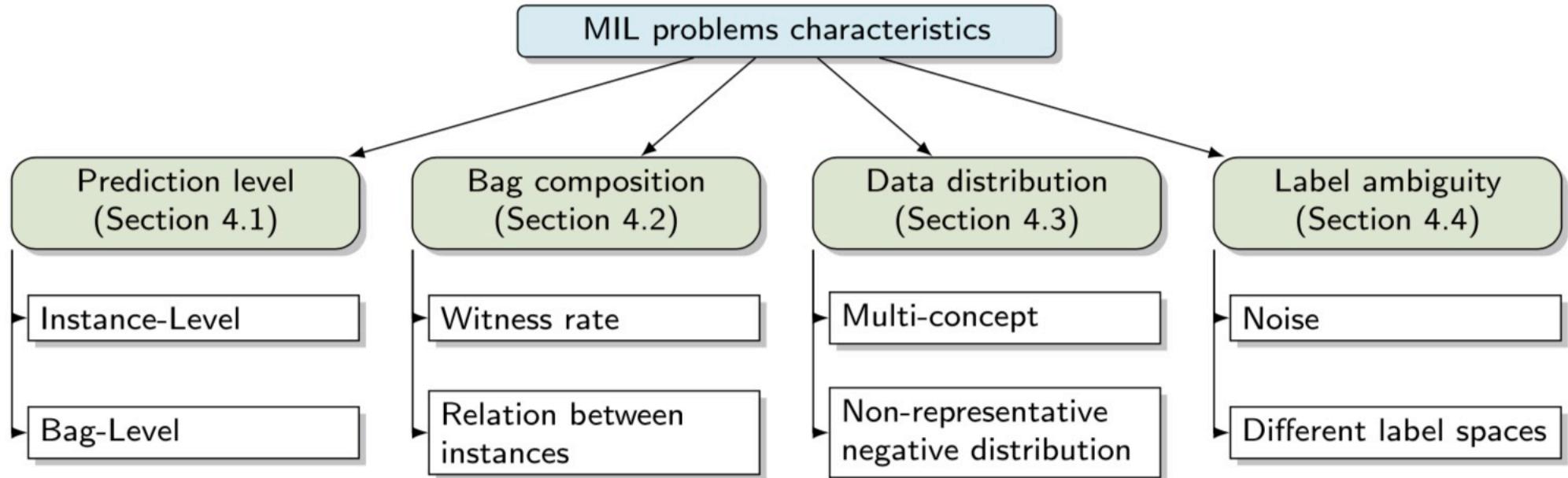


Fig. 1. Characteristics inherent to MIL problems.

Witness rate : the scenario where the number of positive instances in a bag

# Bag Dissimilarity Representation

Represent a MIL bag by its dissimilarities to prototype objects in a representation set R .  
R is taken to be a subset of size M of the training set of size N (typically M<N) .

$$\mathbf{d}(B_i, \mathcal{T}) = [d(B_i, B_1), \dots d(B_i, B_M)]$$

Each bag is represented by a single feature vector.

MIL problem can be viewed as a regular supervised learning problem.

# Bag Dissimilarity Representation

## Point set

Hausdorff distance applied to bags uses the maximum mismatch between the instances of the respective bags

$$d_h(B_i, B_j) = \max_k \min_l d(\mathbf{x}_{ik}, \mathbf{x}_{jl}) \quad (1)$$

## Modified versions of the Hausdorff distance

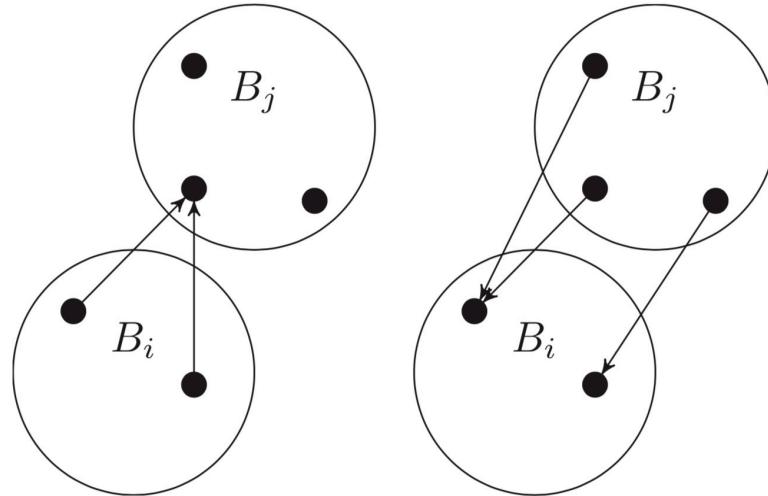
$$d_{minmin}(B_i, B_j) = \min_k \min_l d(\mathbf{x}_{ik}, \mathbf{x}_{jl}) \text{ and} \quad (2)$$

$$d_{meanmin}(B_i, B_j) = \frac{1}{n_i} \sum_{k=1}^{n_i} \min_l d(\mathbf{x}_{ik}, \mathbf{x}_{jl}). \quad (3)$$

$$d_{meanmean}(B_i, B_j) = \frac{1}{n_i n_j} \sum_{k=1}^{n_i} \sum_{l=1}^{n_j} d(\mathbf{x}_{ik}, \mathbf{x}_{jl}). \quad (4)$$

# Bag Dissimilarity Representation

## Point set

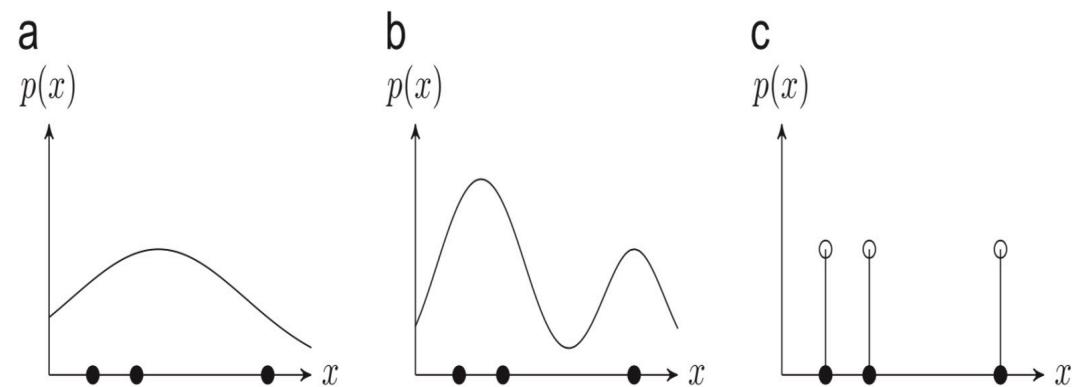


**Fig. 2.** Minimum instance distances between two bags. The bag dissimilarity is defined as the minimum, maximum, average or other statistic of these distances. The directions of the arrows show that there are two, possibly asymmetric, dissimilarities: that of  $B_i$  to  $B_j$ , and that of  $B_j$  to  $B_i$ .

shows the first step in computing dissimilarities between two bags. The arrows in each diagram are the minimum distances between instances of two bags

# Bag Dissimilarity Representation

## Bags as instance distributions



**Fig. 3.** Different ways to represent a 1-D bag with instances at  $x=0.5$ ,  $x=1$  and  $x=2.5$  as a distribution. From left to right: normal distribution, Parzen density,  $\delta$ -peaks. The type of approximation influences the choice of distribution distance that can be applied.

view each bag as a probability distribution in instance space, and define a bag dissimilarity as a distribution distance; approximate the instance distributions and provide distances between approximated distributions.

a number of approaches to approximate a 1D distribution

# Bag Dissimilarity Representation

Bags as instance distributions

Bag dissimilarity measure

$$d_{\text{Maha}}(B_i, B_j) = (\mu_i - \mu_j)^\top \left( \frac{1}{2}\Sigma_i + \frac{1}{2}\Sigma_j \right)^{-1} (\mu_i - \mu_j). \quad (5)$$

$$d_{\text{CS}}(B_i, B_j) = -\log \left( \frac{K_{\sigma_i + \sigma_j}(B_i, B_j)}{(K_{2\sigma_i}(B_i, B_i)K_{2\sigma_j}(B_j, B_j))^{1/2}} \right)$$

where

$$K_\sigma(B_i, B_j) = \sum_{\substack{\mathbf{x}_k \in B_i \\ \mathbf{x}_l \in B_j}} \frac{\exp\left(-\frac{1}{2\sigma^2} (\mathbf{x}_k - \mathbf{x}_l)^\top (\mathbf{x}_k - \mathbf{x}_l)\right)}{(2\pi\sigma^2)^{d/2}}. \quad (6)$$

# Bag Dissimilarity Representation

## Bags as instance distributions

Bag dissimilarity measure :earth movers distance (EMD)

EMD measures the minimum amount of work to transform one probability into another probability distribution

$$d_{\text{EMD}}(B_i, B_j) = \sum_{\mathbf{x}_k \in B_i, \mathbf{x}_l \in B_j} f(\mathbf{x}_k, \mathbf{x}_l) d(\mathbf{x}_k, \mathbf{x}_l) \quad (7)$$

where  $f(\mathbf{x}_k, \mathbf{x}_l)$  is the flow that minimizes the overall distance, and that is subject to constraints that ensure that the only available amounts of earth are transported into available holes, and that all of the earth is indeed transported:  $f(\mathbf{x}_k, \mathbf{x}_l) \geq 0$ ,  $\sum_{\mathbf{x}_k \in B_i} f(\mathbf{x}_k, \mathbf{x}_l) \leq 1/n_j$ ,  $\sum_{\mathbf{x}_l \in B_j} f(\mathbf{x}_k, \mathbf{x}_l) \leq 1/n_i$  and  $\sum_{\mathbf{x}_k \in B_i, \mathbf{x}_l \in B_j} f(\mathbf{x}_k, \mathbf{x}_l) = 1$ .

# Experiments

**Table 2**

Point set, symmetrized dissimilarity, logistic and SVM classifiers. AUC and standard error ( $\times 100$ ),  $5 \times 10$ -fold cross-validation. Bold=best (or not significantly worse) result per dataset.

Classifier	Data	$D_{minmin}$	$D_{meanmin}$	$D_{maxmin}$	$D_{meanmean}$
Logistic	C25	<b>61.4 (2.3)</b>	54.8 (2.2)	47.5 (2.2)	53.4 (2.1)
	C50	<b>98.6 (0.4)</b>	79.6 (1.9)	50.3 (3.3)	65.6 (3.0)
	D25	86.2 (1.8)	<b>97.8 (0.5)</b>	96.6 (0.6)	69.4 (2.4)
	D50	91.7 (1.1)	<b>100.0 (0.0)</b>	99.6 (0.2)	<b>100.0 (0.0)</b>
	M25	54.4 (2.2)	50.8 (1.9)	60.5 (2.3)	<b>71.4 (1.9)</b>
	M50	71.5 (2.4)	78.4 (2.1)	<b>84.3 (1.5)</b>	69.4 (2.4)
	Musk 1	88.2 (1.8)	<b>90.2 (1.7)</b>	<b>91.8 (1.6)</b>	84.3 (1.8)
	Musk 2	92.0 (1.2)	<b>92.6 (1.3)</b>	<b>93.3 (1.2)</b>	82.8 (1.7)
	African	<b>96.3 (0.4)</b>	94.4 (0.6)	94.2 (0.6)	91.1 (0.7)
	Ajax	68.4 (1.6)	<b>98.1 (0.5)</b>	<b>97.2 (0.7)</b>	87.8 (1.1)
	Alt.ath	49.2 (0.8)	<b>88.5 (1.7)</b>	83.7 (1.7)	85.2 (1.8)
	BrCr	89.6 (0.6)	<b>93.6 (0.4)</b>	91.1 (0.5)	82.3 (0.7)
	Web	69.7 (4.0)	<b>77.0 (3.2)</b>	66.8 (3.7)	69.9 (3.3)
LibSVM	C25	<b>57.7 (2.4)</b>	52.2 (2.6)	45.9 (2.1)	40.8 (1.5)
	C50	<b>98.6 (0.4)</b>	83.9 (2.0)	46.1 (2.4)	66.4 (3.0)
	D25	72.0 (2.4)	<b>78.9 (2.0)</b>	<b>82.7 (1.6)</b>	68.2 (2.5)
	D50	92.9 (1.1)	<b>100.0 (0.0)</b>	99.5 (0.2)	<b>100.0 (0.0)</b>
	M25	47.0 (2.2)	50.6 (2.3)	66.2 (2.3)	<b>71.6 (2.0)</b>
	M50	72.0 (2.4)	<b>78.9 (2.0)</b>	<b>82.7 (1.6)</b>	68.2 (2.5)
	Musk 1	92.0 (1.2)	<b>93.4 (1.2)</b>	<b>93.4 (1.3)</b>	88.2 (1.5)
	Musk 2	94.0 (1.3)	<b>95.4 (1.4)</b>	<b>95.3 (1.2)</b>	90.3 (1.5)
	African	<b>96.6 (0.3)</b>	<b>96.7 (0.4)</b>	95.5 (0.5)	90.1 (0.7)
	Ajax	71.1 (1.4)	<b>98.6 (0.4)</b>	<b>97.8 (0.5)</b>	84.0 (1.2)
	Alt.ath	50.0 (0.0)	<b>94.9 (1.0)</b>	91.4 (1.1)	94.2 (1.1)
	BrCr	87.8 (0.6)	<b>95.5 (0.3)</b>	92.6 (0.5)	54.4 (2.4)
	Web	53.2 (4.8)	<b>76.0 (2.7)</b>	43.3 (3.6)	<b>77.6 (3.3)</b>

# Experiments

**Table 3**

Distribution dissimilarities. AUC and standard error ( $\times 100$ ),  $5 \times 10$ -fold cross-validation. Bold=best (or not significantly worse) result per dataset.

Classifier	Data	$D_{\text{Maha}}$	$D_{\text{EMD}}$	$D_{\text{CS}}$
Logistic	Musk 1	70.1 (2.3)	<b>88.7 (1.9)</b>	84.6 (2.0)
	Musk 2	<b>91.5 (1.3)</b>	–	88.8 (1.5)
	African	59.9 (1.5)	<b>93.3 (0.6)</b>	85.9 (1.0)
	Ajax	86.9 (1.8)	<b>97.9 (0.4)</b>	95.5 (0.7)
	Alt.ath	49.9 (2.5)	<b>84.0 (1.9)</b>	59.2 (2.9)
	BrCr	63.3 (1.2)	<b>94.5 (0.4)</b>	89.4 (0.8)
	Web	–	69.4 (4.1)	<b>75.7 (3.4)</b>
LibSVM	Musk 1	76.5 (2.9)	<b>89.8 (1.6)</b>	<b>88.2 (1.7)</b>
	Musk 2	<b>96.0 (0.9)</b>	–	87.4 (1.8)
	African	64.8 (1.5)	<b>94.7 (0.4)</b>	<b>94.3 (0.5)</b>
	Ajax	87.3 (1.7)	<b>98.9 (0.3)</b>	98.1 (0.3)
	Alt.ath	47.0 (2.3)	<b>87.4 (1.7)</b>	41.9 (2.5)
	BrCr	59.7 (1.1)	<b>95.5 (0.3)</b>	93.9 (0.4)
	Web	–	<b>77.7 (2.7)</b>	69.5 (3.8)

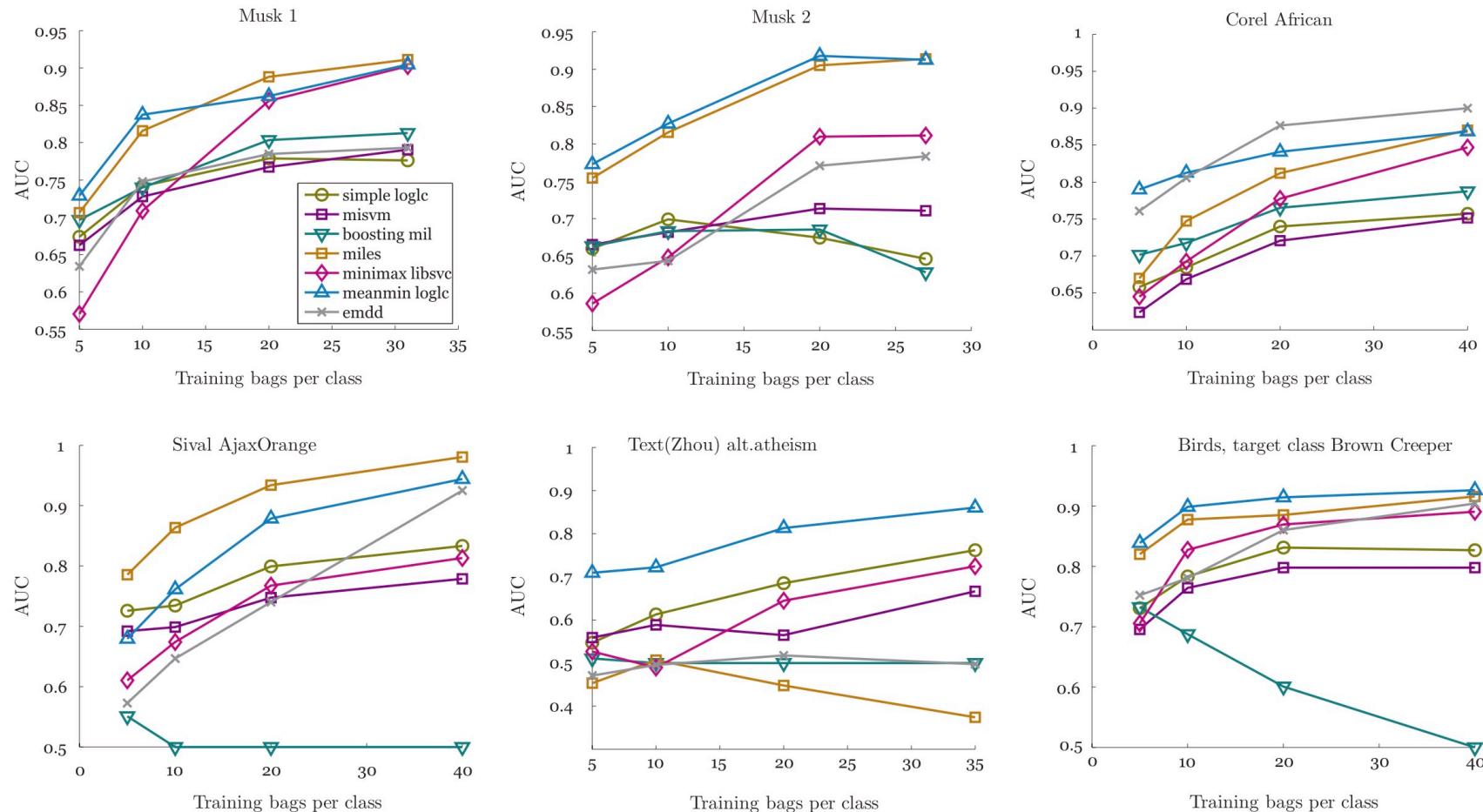
# Experiments

**Table 4**

Properties of dissimilarity matrices: negative eigenfraction (in %), negative eigenratio and non-metricity fraction (in %).

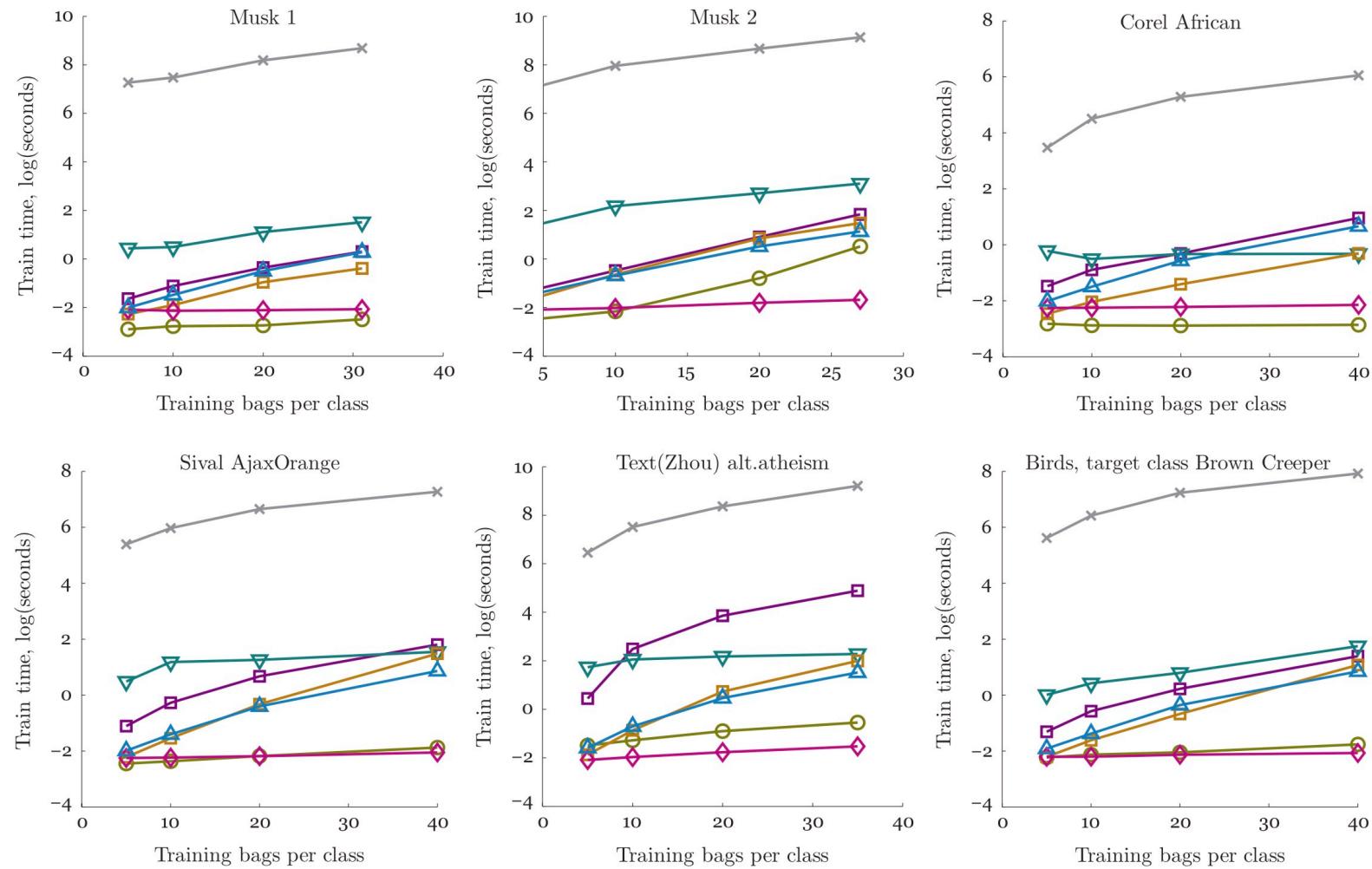
Measure	Data	$D_{minmin}$	$D_{meanmin}$	$D_{maxmin}$	$D_{mean}$	$D_{Maha}$	$D_{EMD}$	$D_{CS}$
NEF	Musk 1	31.2	25.5	22.6	21.0	48.9	27.4	27.1
	Musk 2	31.1	25.3	21.7	20.5	28.4	—	96.3
	Afr	47.2	37.8	33.1	29.7	50.0	49.5	35.4
	Ajax	47.9	32.6	41.3	40.0	49.6	29.9	32.5
	Alt	49.5	21.0	22.0	93.1	36.9	37.4	9.2
	BrCr	45.5	30.7	33.6	51.1	49.9	37.3	27.3
	Web	18.6	4.8	3.7	25.9	—	10.3	100.0
NER	Musk 1	0.4	0.2	0.2	0.2	1.0	0.2	0.2
	Musk 2	0.4	0.2	0.2	0.1	0.3	—	28.1
	Afr	0.7	0.4	0.3	0.3	1.0	1.0	0.3
	Ajax	0.5	0.2	0.3	1.1	1.0	0.2	0.3
	Alt	1.0	0.5	0.6	24.1	0.8	0.7	0.8
	BrCr	0.5	0.4	0.3	0.8	1.0	0.6	0.2
	Web	1.0	1.1	0.9	1.1	—	0.5	0.0
NMF	Musk 1	5.5	2.5	1.6	0.7	27.1	3.8	5.5
	Musk 2	7.2	3.2	1.9	0.8	7.9	—	93.3
	Afr	20.1	8.1	8.8	2.3	27.8	11.4	12.4
	Ajax	17.5	0.1	0.5	0.0	18.9	0.3	4.2
	Alt	0.0	0.3	0.9	0.0	1.7	9.7	0.0
	BrCr	7.8	0.8	2.0	4.2	26.7	3.5	13.1
	Web	2.4	0.0	0.0	0.0	—	0.0	97.3

# Experiments



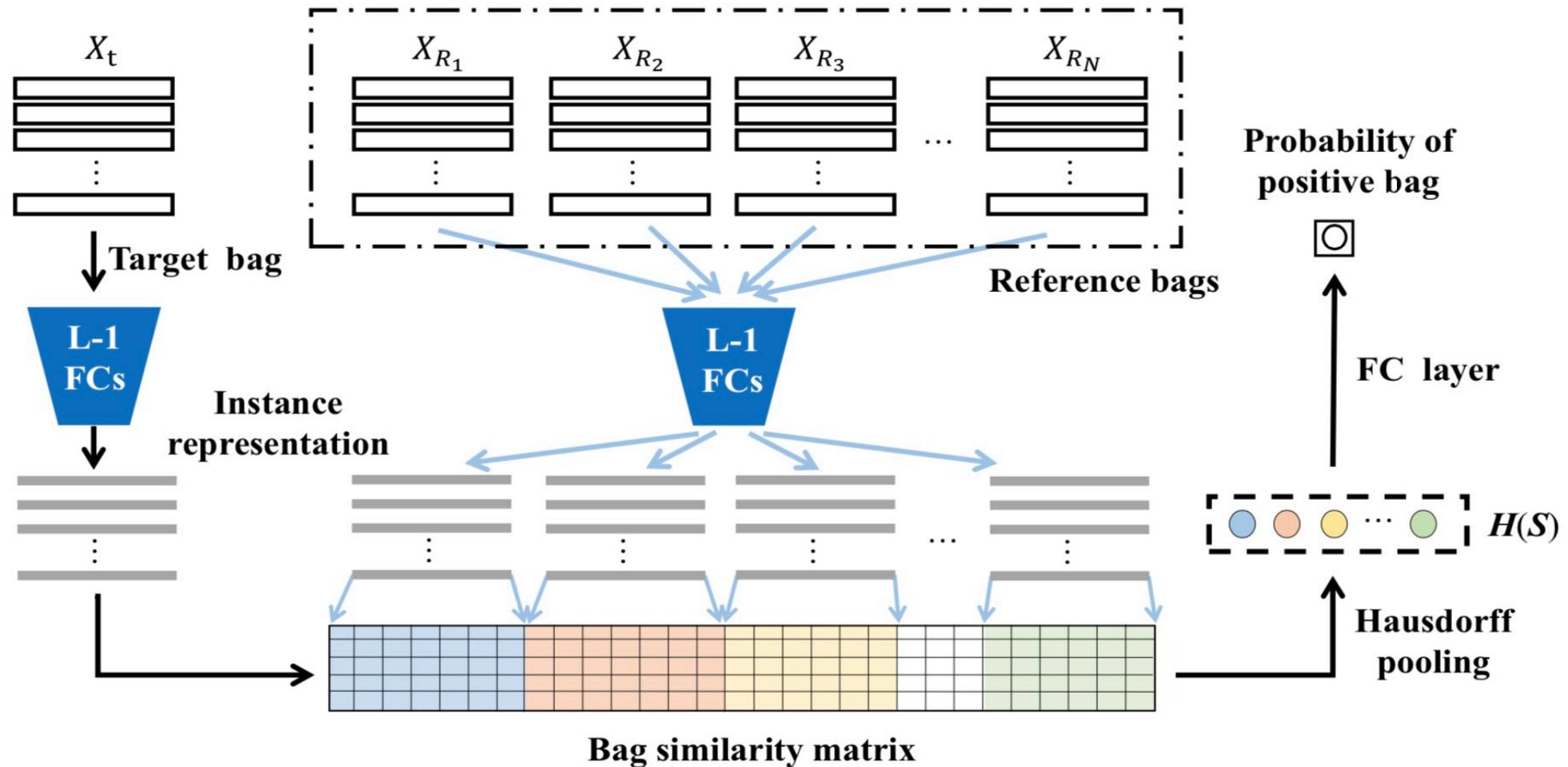
**Fig. 7.** Learning curves (AUC performance) for Musk 1, Musk 2, African, AjaxOrange, alt.atheism and Brown Creeper datasets. The standard deviations are generally around 0.1, but are lower for MILES, Minimax and MInD at larger training sizes.

# Experiments



**Fig. 8.** Training time curves for Musk 1, Musk 2, African, AjaxOrange, alt.atheism and Brown Creeper datasets. The standard deviations in time are all quite low except for EM-DD.

# Bag similarity network for MIL



**Fig. 2.** Architecture of bag similarity network.

# Bag similarity network for MIL

## Instance-level representations

all training bags  $X = \{X_1, X_2, \dots, X_N\}$  as reference bags

$$X_R^{L-1} = \{X_{R_1}^{L-1}, X_{R_2}^{L-1}, \dots, X_{R_N}^{L-1}\}, \quad (1)$$

where  $X_{R_i}^{L-1} = \{x_{R_i1}^{L-1}, x_{R_i2}^{L-1}, \dots, x_{R_im_i}^{L-1}\}$  denotes the learned instance-level representations of reference bag  $X_{R_i}$ .

# Bag similarity network for MIL

## Bag similarity matrix

$$S = [f(X_t^{L-1}, X_{R_1}^{L-1}), \dots, f(X_t^{L-1}, X_{R_N}^{L-1})], \quad (2)$$

where the similarity between  $X_t$  and  $X_{R_i}$  is represented by  $f(X_t^{L-1}, X_{R_i}^{L-1})$ , which is defined as

$$\begin{pmatrix} f(x_{t1}^{L-1}, x_{R_i1}^{L-1}) & f(x_{t1}^{L-1}, x_{R_i2}^{L-1}) & \dots & f(x_{t1}^{L-1}, x_{R_im_i}^{L-1}) \\ f(x_{t2}^{L-1}, x_{R_i1}^{L-1}) & f(x_{t2}^{L-1}, x_{R_i2}^{L-1}) & \dots & f(x_{t2}^{L-1}, x_{R_im_i}^{L-1}) \\ \dots & \dots & \dots & \dots \\ f(x_{tm_t}^{L-1}, x_{R_i1}^{L-1}) & f(x_{tm_t}^{L-1}, x_{R_i2}^{L-1}) & \dots & f(x_{tm_t}^{L-1}, x_{R_im_i}^{L-1}) \end{pmatrix}, \quad (3)$$

where  $m_t$  and  $m_i$  denote the number of instances in  $X_t$  and  $X_{R_i}$ , respectively. Therefore,  $X_t$  can be represented by a bag similarity matrix, which describes the similarity of  $X_t$  to reference bags  $X_R = \{X_{R_1}, X_{R_2}, \dots, X_{R_N}\}$ .

# Bag similarity network for MIL

## Hausdorff pooling

$$H(S) = [h(f(X_t^{L-1}, X_{R_1}^{L-1})), h(f(X_t^{L-1}, X_{R_2}^{L-1})), \dots, h(f(X_t^{L-1}, X_{R_N}^{L-1}))], \quad (4)$$

where  $h(\cdot)$  can be one of the following:

$$\begin{cases} \text{max-max pooling : } h(A) = \max_i \max_j A_{ij}; \\ \text{mean-max pooling : } h(A) = \frac{1}{N} \sum_{i=1}^N \max_j A_{ij}; \\ \text{min-max pooling : } h(A) = \min_i \max_j A_{ij}; \\ \text{mean-mean pooling : } h(A) = \frac{1}{N} \frac{1}{M} \sum_{i=1}^N \sum_{j=1}^M A_{ij}, \end{cases} \quad (5)$$

# Bag similarity network for MIL

## Training loss and optimization

To predict bag labels, it is natural to choose the standard cross entropy loss function, which is the same as the loss function in MI-Net

$$L(\hat{Y}_t, Y_t) = -((1 - Y_t) \log(1 - \hat{Y}_t) + Y_t \log \hat{Y}_t), \quad (6)$$

where  $\hat{Y}_t$  is the probability that bag  $B_t$  is positive, and  $Y_t$  is the label of  $X_t$ .

# Datasets

**Table 1**

Parameter details for training BSN with MI-Net, including initial learning rate (LR), weight decay (WD), and the standard deviation of initial weights (W-Std). The parameters for Colon Cancer dataset are the same as in Attention-based MI-Net [13].

Dataset	LR	WD	W-Std
MUSK1	0.001	0.01	0.05
MUSK2	0.001	0.01	0.05
Fox	0.0005	0.01	0.05
Tiger	0.001	0.01	0.05
Elephant	0.0005	0.01	0.05
20 Newsgroups	0.001	0.001	0.1
Messidor	0.0005	0.001	0.05

**MUSK1** and **MUSK2** are used to predict the molecular activity of drugs, where bags are molecules and instances are different conformations of these molecules.

**Fox**, **Tiger**, and **Elephant** are widely adopted MIL datasets for solving localized content-based image retrieval problems. The bags are images and instances are image segments.

**20 Newsgroups** contains posts from newsgroups on 20 subjects for text categorization.

**Messidor** is a public diabetic retinopathy screening dataset that contains 1200 eye fundus images from 654 diabetic and 546 healthy patients.

**Colon Cancer** contains 100 H&E images generated from normal or malignant tissue appearance. The majority of the nuclei of each cell are marked in each image.

# Experiments

**Table 2**

Bag classification results ( $mean \pm std$ ) of different methods on MUSK1 and MUSK2 (task: drag activation prediction), as well as Fox, Tiger, and Elephant (task: content-based image retrieval).

Dataset	MUSK1	MUSK2	Fox	Tiger	Elephant
mi-SVM	0.874	0.836	0.582	0.784	0.822
MI-SVM	0.779	0.843	0.578	0.840	0.814
MI-Kernel	0.880	0.893	0.603	0.842	0.843
EM-DD	$0.849 \pm 0.098$	$0.869 \pm 0.108$	$0.609 \pm 0.101$	$0.730 \pm 0.096$	$0.771 \pm 0.098$
mi-Graph	$0.889 \pm 0.073$	$0.903 \pm 0.086$	$0.616 \pm 0.079$	$0.860 \pm 0.083$	$0.869 \pm 0.078$
miVLAD	$0.871 \pm 0.097$	$0.872 \pm 0.095$	$0.620 \pm 0.098$	$0.811 \pm 0.087$	$0.850 \pm 0.080$
miFV	$0.909 \pm 0.089$	$0.884 \pm 0.094$	$0.621 \pm 0.109$	$0.813 \pm 0.083$	$0.852 \pm 0.081$
MInD	$0.893 \pm 0.019$	$0.888 \pm 0.034$	<b><math>0.651 \pm 0.011</math></b>	$0.819 \pm 0.021$	$0.857 \pm 0.018$
MI-Net	$0.893 \pm 0.099$	$0.872 \pm 0.096$	$0.627 \pm 0.080$	$0.832 \pm 0.087$	$0.891 \pm 0.074$
Att. Net	$0.892 \pm 0.040$	$0.858 \pm 0.048$	$0.615 \pm 0.043$	$0.839 \pm 0.022$	$0.868 \pm 0.022$
Gated Att. Net	$0.900 \pm 0.050$	$0.863 \pm 0.042$	$0.603 \pm 0.029$	$0.845 \pm 0.018$	$0.857 \pm 0.027$
BSN	<b><math>0.931 \pm 0.094</math></b>	<b><math>0.906 \pm 0.109</math></b>	$0.640 \pm 0.111$	<b><math>0.878 \pm 0.092</math></b>	<b><math>0.907 \pm 0.073</math></b>

Table 2 provides the average accuracy and the corresponding standard deviation of methods under comparison.BSN achieves superior performance on MUSK1,MUSK2 ,Tiger and Elephant, and competitive performance on Fox .

# Experiments

**Table 3**

Bag classification results of different methods on 20 Newsgroups (task: text categorization).

Dataset	MI-Kernel	mi-Graph	miFV	MInD	MI-Net	Att. Net	Gated Att. Net	BSN
alt.atheism	0.602 ± 0.039	0.655 ± 0.040	0.848 ± 0.119	0.861 ± 0.089	0.846 ± 0.101	0.784 ± 0.084	0.780 ± 0.074	<b>0.903 ± 0.101</b>
comp.graphics	0.470 ± 0.033	0.778 ± 0.016	0.594 ± 0.120	0.825 ± 0.118	0.831 ± 0.123	0.774 ± 0.081	0.764 ± 0.073	<b>0.861 ± 0.134</b>
comp.windows.misc	0.510 ± 0.052	0.631 ± 0.015	0.615 ± 0.172	0.730 ± 0.094	0.730 ± 0.112	0.686 ± 0.088	0.700 ± 0.080	<b>0.769 ± 0.101</b>
comp.ibm.pc.hardware	0.469 ± 0.036	0.595 ± 0.027	0.665 ± 0.147	0.780 ± 0.127	0.803 ± 0.155	0.632 ± 0.087	0.640 ± 0.080	<b>0.813 ± 0.29</b>
comp.sys.mac.hardware	0.445 ± 0.032	0.617 ± 0.048	0.660 ± 0.157	0.835 ± 0.098	0.811 ± 0.138	0.744 ± 0.084	0.754 ± 0.082	<b>0.865 ± 0.113</b>
comp.window.x	0.508 ± 0.043	0.698 ± 0.021	0.768 ± 0.155	0.785 ± 0.111	0.836 ± 0.135	0.766 ± 0.093	0.780 ± 0.075	<b>0.869 ± 0.111</b>
misc.forsale	0.518 ± 0.025	0.552 ± 0.027	0.565 ± 0.146	0.729 ± 0.102	0.707 ± 0.119	0.706 ± 0.076	0.674 ± 0.072	<b>0.768 ± 0.128</b>
rec.autos	0.529 ± 0.033	0.720 ± 0.037	0.667 ± 0.166	0.775 ± 0.088	0.812 ± 0.129	0.762 ± 0.081	0.724 ± 0.091	<b>0.830 ± 0.121</b>
rec.motorcycles	0.506 ± 0.035	0.640 ± 0.028	0.802 ± 0.144	0.577 ± 0.102	0.853 ± 0.119	0.750 ± 0.097	0.814 ± 0.066	<b>0.868 ± 0.116</b>
rec.sport.baseball	0.517 ± 0.028	0.647 ± 0.031	0.779 ± 0.148	0.837 ± 0.078	0.871 ± 0.113	0.774 ± 0.080	0.790 ± 0.078	<b>0.887 ± 0.113</b>
rec.sport.hockey	0.513 ± 0.034	0.850 ± 0.025	0.823 ± 0.137	0.833 ± 0.096	0.918 ± 0.111	0.936 ± 0.041	0.932 ± 0.045	<b>0.947 ± 0.085</b>
sci.crypt	0.563 ± 0.036	0.696 ± 0.021	0.760 ± 0.146	0.768 ± 0.122	0.808 ± 0.154	0.804 ± 0.063	0.748 ± 0.088	<b>0.826 ± 0.142</b>
sci.electronics	0.506 ± 0.019	0.871 ± 0.017	0.555 ± 0.156	<b>0.940 ± 0.078</b>	0.928 ± 0.090	0.854 ± 0.053	0.828 ± 0.064	0.927 ± 0.088
sci.med	0.506 ± 0.019	0.621 ± 0.039	0.783 ± 0.125	0.832 ± 0.091	0.862 ± 0.111	0.772 ± 0.090	0.742 ± 0.101	<b>0.885 ± 0.107</b>
sci.space	0.547 ± 0.025	0.757 ± 0.034	0.818 ± 0.131	0.796 ± 0.110	0.820 ± 0.087	0.888 ± 0.062	<b>0.894 ± 0.063</b>	0.869 ± 0.112
soc.religion.christian	0.492 ± 0.034	0.590 ± 0.047	0.814 ± 0.138	0.841 ± 0.140	0.829 ± 0.122	0.726 ± 0.088	0.708 ± 0.100	<b>0.878 ± 0.113</b>
talk.politics.guns	0.477 ± 0.038	0.585 ± 0.060	0.747 ± 0.150	0.806 ± 0.094	0.782 ± 0.095	0.714 ± 0.074	0.708 ± 0.100	<b>0.814 ± 0.117</b>
talk.politics.mideast	0.559 ± 0.028	0.736 ± 0.026	0.793 ± 0.135	0.830 ± 0.108	0.825 ± 0.119	0.750 ± 0.084	0.784 ± 0.064	<b>0.867 ± 0.127</b>
talk.politics.misc	0.515 ± 0.037	0.704 ± 0.036	0.697 ± 0.152	0.720 ± 0.119	0.748 ± 0.140	0.788 ± 0.091	0.806 ± 0.078	<b>0.829 ± 0.135</b>
talk.religion.misc	0.554 ± 0.043	0.633 ± 0.035	0.739 ± 0.151	0.725 ± 0.104	0.778 ± 0.114	0.738 ± 0.074	0.746 ± 0.082	<b>0.821 ± 0.115</b>
average	0.515	0.679	0.726	0.791	0.820	0.767	0.766	<b>0.855</b>

It can be seen that BSN outperforms all competitors in all cases except for *sci.electronics*, where the accuracy of MInD (best) and MI-Net (second best) is higher by 1.3% and 0.1%, respectively. The average accuracy on all 20 datasets indicates that both MI-Net and BSN outperform the others.

# Experiments

**Table 4**

Comparison of different methods ( $mean \pm std$ ) for bag classification on Messidor dataset.

mi-SVM	MI-SVM	miVLAD	miFV	MInD	MI-Net	Att.Net	Gated Att. Net	BSN
0.620 ± 0.039	0.640 ± 0.050	0.691 ± 0.037	0.715 ± 0.047	0.665 ± 0.071	0.730 ± 0.051	0.703 ± 0.041	0.698 ± 0.048	<b>0.737</b> ± 0.050

Table 4 shows that BSN achieves the best result on Medical image diagnosis

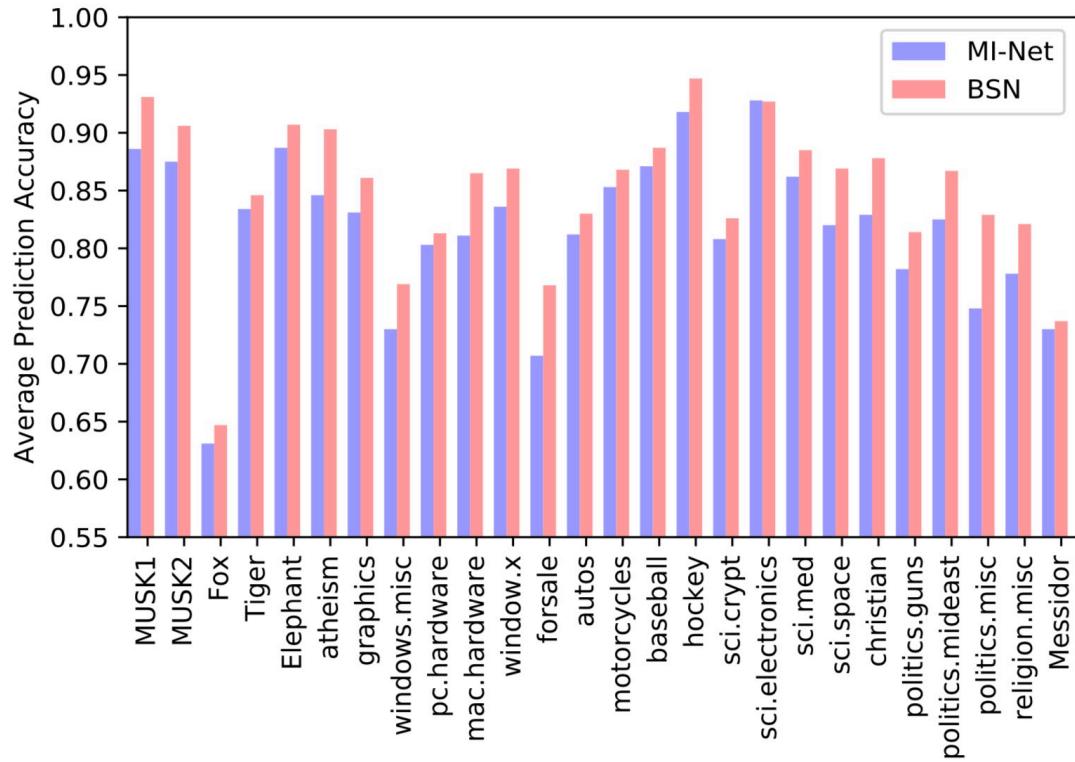
**Table 5**

Comparison of different methods ( $mean \pm std$ ) for bag classification on Colon Cancer dataset.

Method	Accuracy	precision	Recall	F-score
Instance+max	$0.842 \pm 0.021$	$0.866 \pm 0.017$	$0.816 \pm 0.031$	$0.839 \pm 0.023$
Embedding+max	$0.824 \pm 0.015$	$0.884 \pm 0.014$	$0.753 \pm 0.020$	$0.813 \pm 0.017$
Att. Net	<b><math>0.904 \pm 0.011</math></b>	<b><math>0.953 \pm 0.014</math></b>	$0.855 \pm 0.017$	<b><math>0.901 \pm 0.011</math></b>
Gated Att. Net	$0.898 \pm 0.020$	$0.944 \pm 0.016$	$0.851 \pm 0.035$	$0.893 \pm 0.022$
BSN	$0.869 \pm 0.008$	$0.820 \pm 0.019$	<b><math>0.983 \pm 0.019</math></b>	$0.886 \pm 0.030$

Table 5 presents Att. Net obtains the best precision ( $95.3 \pm 1.4\%$ ), whereas BSN obtains the best recall ( $98.3 \pm 1.9\%$ ).

# Ablation study

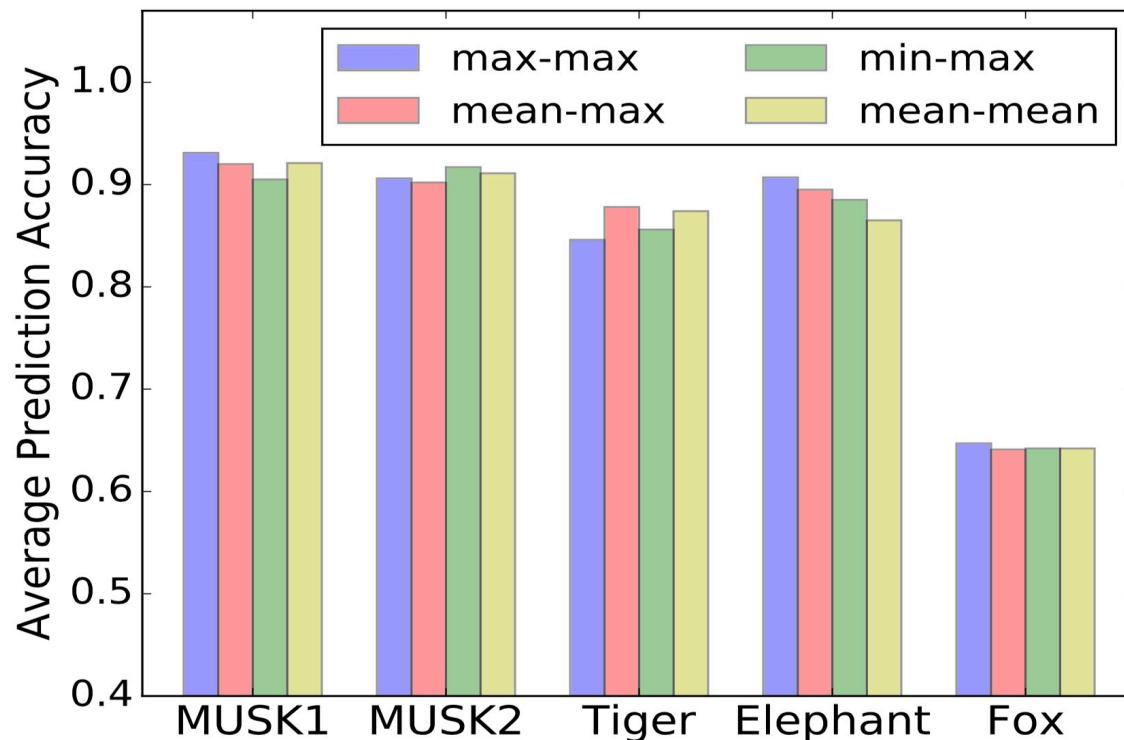


**Fig. 3.** Effectiveness of the proposed bag similarity representation on MUSK1, MUSK2, Fox, Tiger, Elephant, 20 Newsgroups, and Messidor datasets.

The proposed bag similarity representation are superior to other methods such as MI-Net. To justify the argument and indicate the improvement, Fig. 3 offer a detailed comparison between the proposed BSN and MI-Net.

The average accuracy of BSN was improved. As both BSN and MI-Net are based on the same instance feature learning network, the results clearly demonstrate that bag similarity representation is useful for MIL problems.

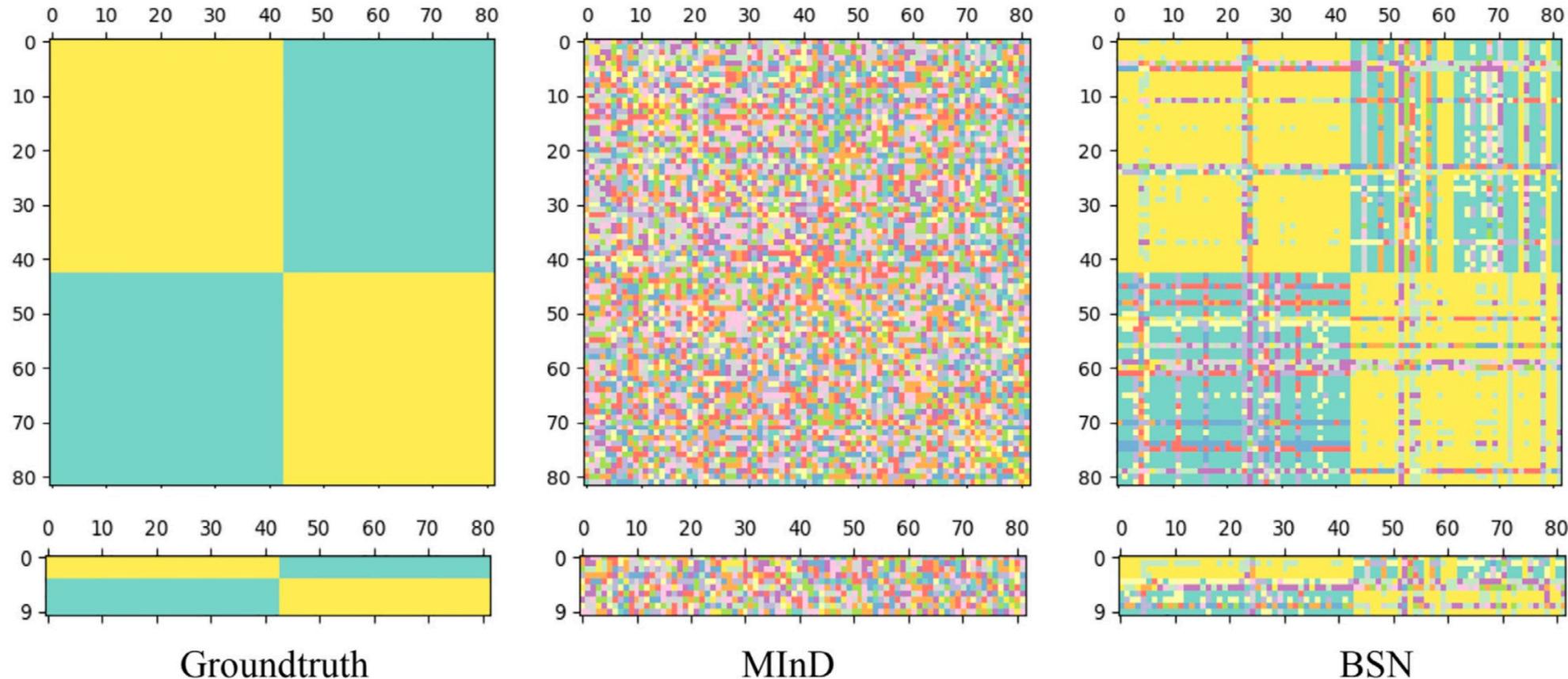
# Ablation study



**Fig. 4.** Comparison of different MIL pooling methods on MUSK1, MUSK2, Fox, Tiger, and Elephant datasets.

The paper introduces four Hausdorff pooling methods to convert the bag similarity matrix into a bag-level feature vector: max–max, mean–max, min–max, and mean– mean pooling functions, which are all differentiable. Usually, the second operator in Hausdorff polling is preferably a max or mean function.

# Bag similarity & Dissimilarity Representation for MIL



**Fig. 1.** Comparison between the fixed similarity metric (MInd) [4] and the learned similarity matrix by BSN. The **top** row shows the similarity matrices of training bags, and the **bottom** row shows the similarity matrices of testing bags. Even though the diagonal-blockness of the similarity matrix by BSN is not perfect, it is quite satisfactory, whereas MInd lacks this desired property.

# Summary

- Proposed a learnable bag similarity and bag dissimilarities framework representation for MIL.;
- Point set distances and distribution distances are considered;
- Characteristics of MIL problems affect the effect of different algorithms.