

# Domain Generalization via Multidomain Discriminant Analysis

Shoubo Hu<sup>1</sup> Kun Zhang<sup>2</sup> Zhitang Chen<sup>3</sup> Laiwan Chan<sup>1</sup>

<sup>1</sup> The Chinese University of Hong Kong <sup>2</sup> Carnegie Mellon University <sup>3</sup> Huawei Noah's Ark Lab

## Motivation

**Background:** distribution shift, which is ubiquitous in practice, is the major source of model performance reduction when applied on previously unseen data.

### Objective (general)

Incorporate the knowledge from multiple source domains to improve the generalization ability of classifiers on unseen target domains. [1]

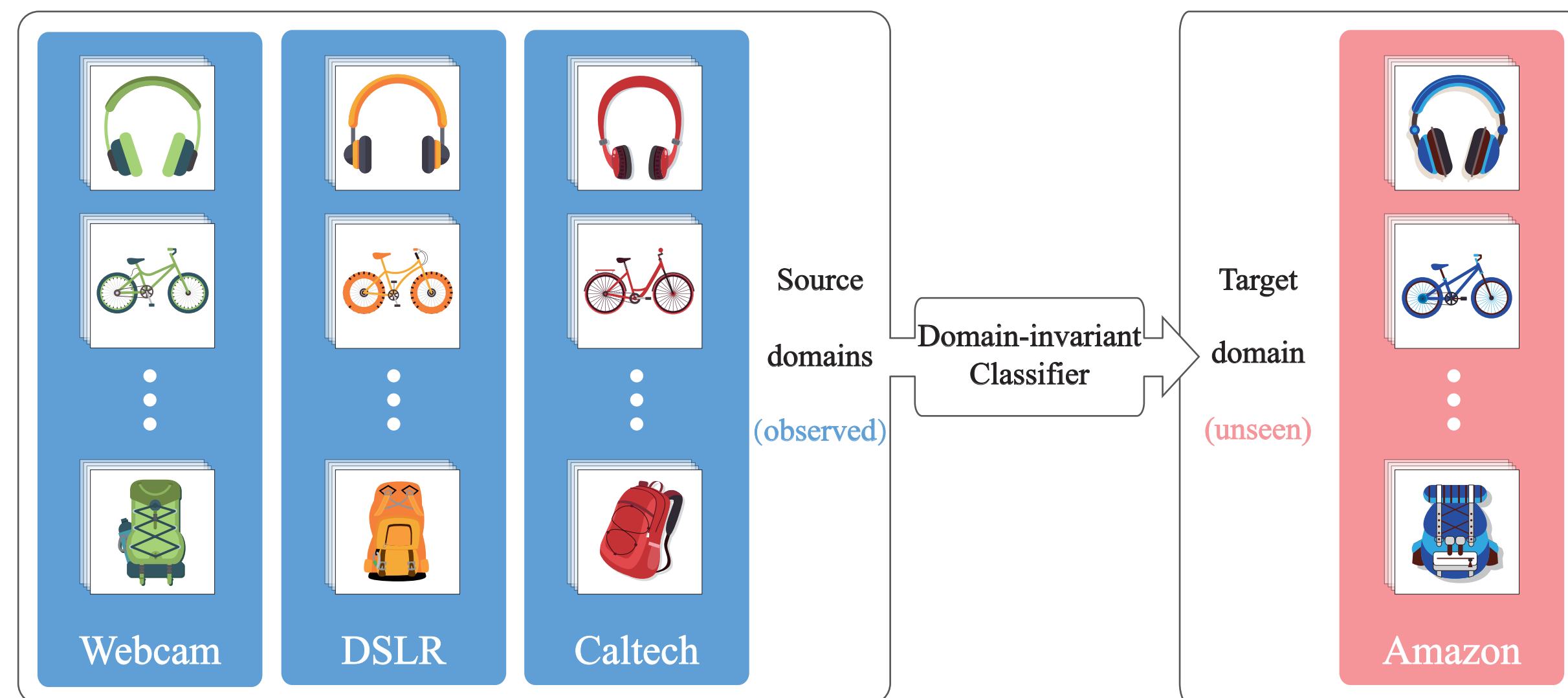


Figure: Illustration of DG on Office+Caltech Dataset. One is given source domains: Webcam, DSLR, Caltech, and aims to train a classifier generalizes well on target domain Amazon, which is unavailable in training.

## Problem Setup

Notation	Description	Notation	Description
$X, Y$	feature/label variable	$\mathbf{x}, y$	feature/label instance
$m, n$	# domains/instances	$Z$	domain-invariant latent variable
$\mathbb{P}_j^s$	class-conditional distribution	$\mu_j^s$	kernel mean embedding of $\mathbb{P}_j^s$
$u_j$	mean representation of class $j$	$\bar{u}$	mean representation of $\mathcal{D}$

### Model assumptions

A domain is defined to be a joint distribution  $\mathbb{P}(X, Y)$ .

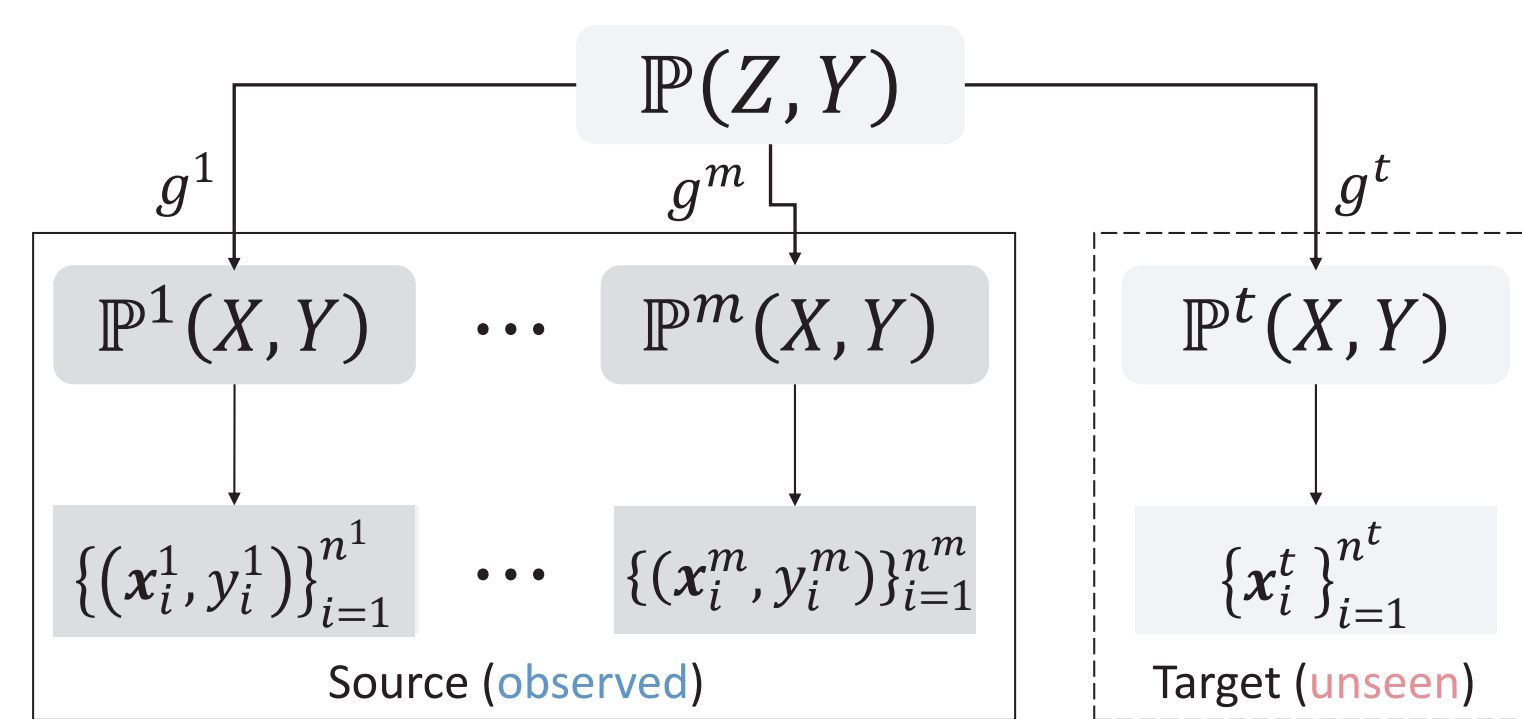


Figure: Domain generalization model assumption.  $m$  domains are uniformly sampled from a set of domains and are called the **source domains**. A model trained on  $m$  source domains is expected to generalize well on an unseen domain  $\mathbb{P}^t(X, Y)$ , which is called the **target domain**.

### Objective (our method)

We aim to learn a feature transformation,  $h(X) : \mathcal{X} \mapsto \mathbb{R}^q$ , from the input space to a  $q$ -dimensional transformed space  $\mathbb{R}^q$  such that

- source instances of the same class are close in  $\mathbb{R}^q$ ;
- source instances of different classes are distant in  $\mathbb{R}^q$ .

#### Postulate 1. Independence of cause and mechanism [2]

If  $Y$  causes  $X$  ( $Y \rightarrow X$ ), then the marginal distribution of the cause,  $\mathbb{P}(Y)$ , and the conditional distribution of the effect given the cause,  $\mathbb{P}(X|Y)$ , are "independent" in the sense that  $\mathbb{P}(X|Y)$  contains no information about  $\mathbb{P}(Y)$ .

According to the postulate above, we factorize the joint distributions in the causal direction

$$\mathbb{P}(X, Y) = \mathbb{P}(Y)\mathbb{P}(X|Y), \quad (1)$$

and manipulate the class-conditional distributions  $\mathbb{P}^s(X|Y = j)$  for  $s = 1, \dots, m$  and  $j = 1, \dots, c$  instead of marginal distributions in most previous works [3].

## Regularization Measures

### Within-class measures (objective 1)

$$\text{Average Domain Discrepancy } \Psi^{add} := \frac{1}{c \binom{m}{2}} \sum_{j=1}^c \sum_{1 \leq s < s' \leq m} \|\mu_j^s - \mu_j^{s'}\|_{\mathcal{H}}^2$$

$$\text{Multidomain within-class scatter } \Psi^{mws} := \frac{1}{n} \sum_{j=1}^c \sum_{s=1}^m \sum_{i=1}^{n_j^s} \|\phi(\mathbf{x}_{i \in j}^s) - u_j\|_{\mathcal{H}}^2$$

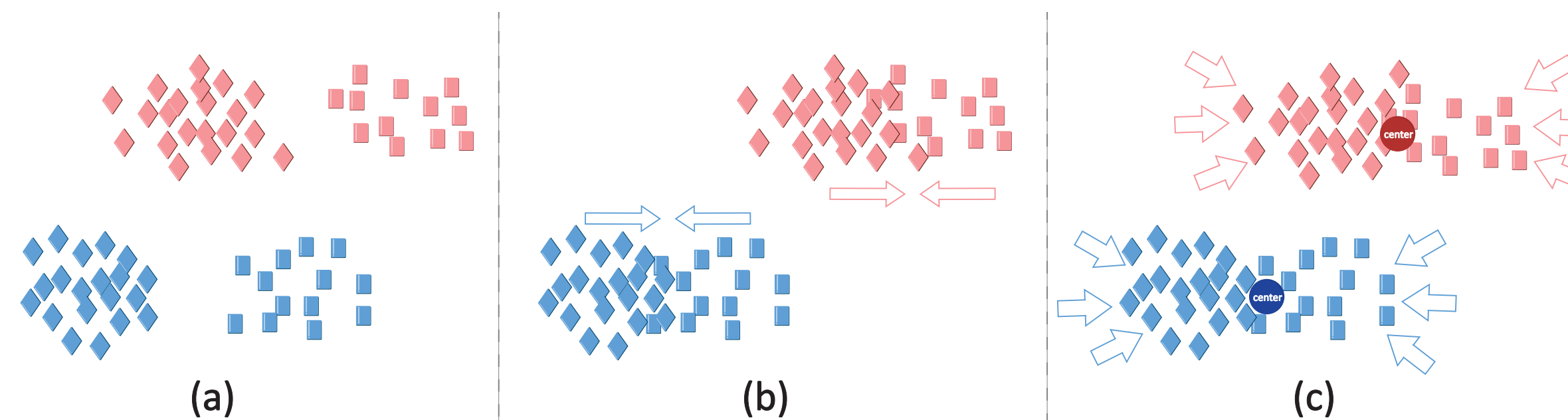


Figure: Illustration of  $\Psi^{add}$  and  $\Psi^{mws}$  (colors - classes; markers - domains). (a) The distribution of data in the subspace  $\mathbb{R}^q$  transformed by some  $\mathbf{W}^0$ . (b) Minimizing  $\Psi^{add}$  makes the means within each class closer. (c) Minimizing  $\Psi^{mws}$  makes the distribution of each class more compact towards its mean representation.

### Between-class measures (objective 2)

$$\text{Average class discrepancy } \Psi^{acd} := \frac{1}{\binom{c}{2}} \sum_{1 \leq j < j' \leq c} \|u_j - u_{j'}\|_{\mathcal{H}}^2$$

$$\text{Multidomain between-class scatter } \Psi^{mbs} := \frac{1}{n} \sum_{j=1}^c n_j \|u_j - \bar{u}\|_{\mathcal{H}}^2$$

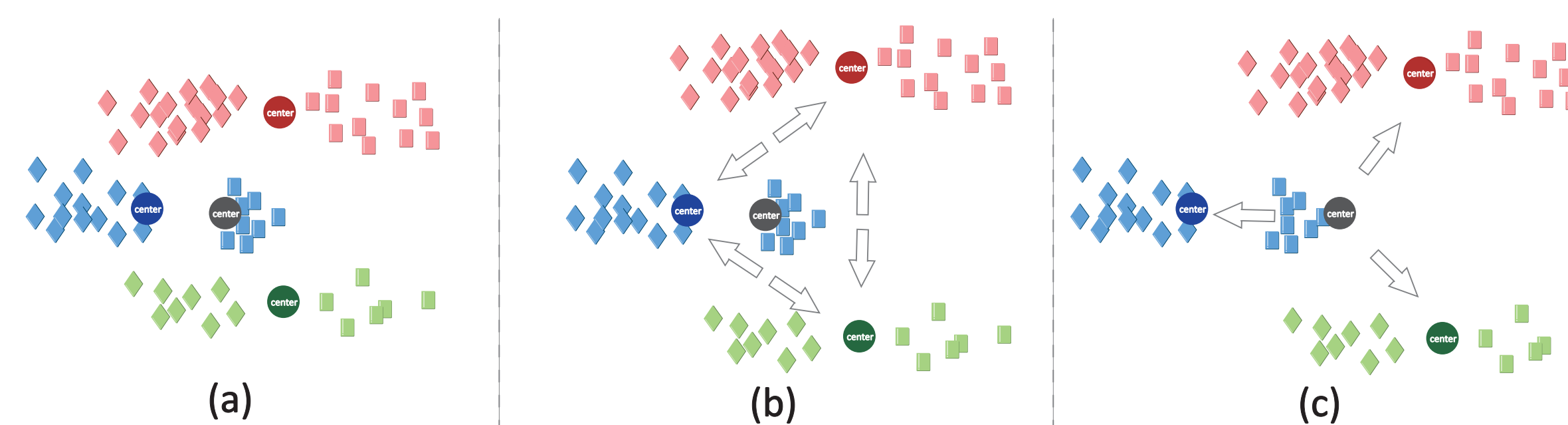


Figure: Illustration of  $\Psi^{acd}$  and  $\Psi^{mbs}$  (colors - classes; markers - domains). (a) The distribution of data in the subspace  $\mathbb{R}^q$  transformed by some  $\mathbf{W}^0$ . (b) Maximizing  $\Psi^{acd}$  makes the distances between each pair of mean representations larger. (c) Maximizing  $\Psi^{mbs}$  makes the average distance between the overall mean and the mean representation of different classes larger.

## Multidomain Discriminant Analysis

### The optimization problem

We unify regularization measures and solve the following optimization problem:

$$\arg \max \frac{\Psi^{acd} + \Psi^{mbs}}{\Psi^{add} + \Psi^{mws}}. \quad (2)$$

We term the proposed method Multidomain Discriminant Analysis (MDA) and summarize the algorithm below

**Input:**  $\mathcal{D} = \{\mathcal{D}^s\}_{s=1}^m$  - the set of instances from  $m$  domains;  
 $\alpha, \beta, \gamma$  - trade-off parameters;

Feature transformation learning	Target feature transformation
<b>Output:</b> Optimal projection $\mathbf{B}_{n \times q}$ ; corresponding eigenvalues $\Gamma$ .	<b>Output:</b> the transformed target features $\mathbf{X}^t$
<ul style="list-style-type: none"> <li>Construct kernel matrix <math>\mathbf{K}</math>, whose entry on <math>i</math>th row and <math>i'</math>th column <math>[\mathbf{K}]_{ii'} = k(\mathbf{x}_i, \mathbf{x}_{i'})</math>;</li> <li>Compute matrices corresponding to regularization measures;</li> <li>Center the kernel matrix as <math>\mathbf{K} \leftarrow \mathbf{K} - \mathbf{1}_n \mathbf{K} - \mathbf{K} \mathbf{1}_n + \mathbf{1}_n \mathbf{K} \mathbf{1}_n</math>, where <math>\mathbf{1}_n \in \mathbb{R}^{n \times n}</math> denotes a matrix with all entries equal to <math>\frac{1}{n}</math>;</li> <li>Solve for the projection <math>\mathbf{B}</math> and corresponding eigenvalues <math>\Gamma</math>, then select <math>q</math> leading components.</li> </ul>	<ul style="list-style-type: none"> <li>Denote the set of instances from the target domain by <math>\mathcal{D}^t</math>, one first constructs the kernel matrix <math>\mathbf{K}^t</math>, where <math>[\mathbf{K}^t]_{ii'} = k(\mathbf{x}_i^t, \mathbf{x}_{i'}^t), \forall \mathbf{x}_i^t \in \mathcal{D}^t, \forall \mathbf{x}_{i'}^t \in \mathcal{D}^t</math>;</li> <li>Center the kernel matrix as <math>\mathbf{K}^t \leftarrow \mathbf{K}^t - \mathbf{1}_{n^t} \mathbf{K}^t - \mathbf{K}^t \mathbf{1}_{n^t} + \mathbf{1}_{n^t} \mathbf{K}^t \mathbf{1}_{n^t}</math>, where <math>n^t</math> is the number of instances in <math>\mathcal{D}^t</math>;</li> <li>Then the transformed features of the target domain are given by <math>\mathbf{X}^t = \mathbf{K}^t \mathbf{B} \Gamma^{-\frac{1}{2}}</math>.</li> </ul>

## Learning Theory Analysis

**Theorem 3.** Under assumptions 2 - 4, and assuming that all source sample sets are of the same size, i.e.  $n^s = \bar{n}$  for  $s = 1, \dots, m$ , then with probability at least  $1 - \delta$  there is

$$\sup_{\|f\|_{\mathcal{H}_k} \leq 1} \left| \frac{1}{m} \sum_{s=1}^m \frac{1}{n^s} \sum_{i=1}^{n^s} \ell(f(\hat{X}_i^s; \mathbf{W}), y_i^s) - \mathcal{E}(f, \infty) \right| \leq U_\ell \left( \sqrt{\frac{\log 2\delta^{-1}}{2m\bar{n}}} + \sqrt{\frac{\log \delta^{-1}}{2m}} \right) + \sqrt{\text{tr}(\mathbf{B}^T \mathbf{K} \mathbf{B})} \left( c_1 \sqrt{\frac{\log 2\delta^{-1} m}{\bar{n}}} + c_2 \left( \sqrt{\frac{1}{m\bar{n}}} + \sqrt{\frac{1}{m}} \right) \right). \quad (3)$$

The first term is of order  $O(m^{-1/2})$  and converges to zero as  $m \rightarrow \infty$ . The second term involves the size of the distortion  $\text{tr}(\mathbf{B}^T \mathbf{K} \mathbf{B})$  introduced by  $\mathbf{B}$ . Therefore, a poor choice of  $\mathbf{B}$  would loose the guarantee.

## Experiments

**Synthetic data.** Data: two-dimensional Gaussian. Domains: two source domains and one target domain. Classes: three classes in each domain.



Figure: Class Prior Distributions  $\mathbb{P}(Y)$  in Synthetic Experiments.

Table: Accuracy (%) of Synthetic Experiments (**bold red** and **bold** indicate the best and second best).

$\mathbb{P}^1(Y)$	(a)	(b)	(c)	(d)	(e)	(a)	(a)	(a)	(a)
$\mathbb{P}^2(Y)$	(a)	(a)	(a)	(a)	(a)	(b)	(c)	(d)	(e)
SVM	56.00	34.00	33.33	33.33	33.33	33.33	40.00	36.00	60.00
KPCA	66.00	62.00	66.67	33.33	33.33	65.33	36.00	40.00	14.00
KFD	78.67	38.67	46.00	74.67	47.33	49.33	34.00	19.33	76.00
L-SVM	56.00	60.00	64.00	62.00	60.67	64.67	45.33	46.00	59.33
DICA	<b>93.33</b>	84.67	76.00	<b>84.00</b>	84.67	54.00	<b>95.33</b>	71.33	<b>88.67</b>
SCA	79.33	72.00	<b>84.67</b>	57.33	76.00	59.33	84.67	61.33	81.33
CIDG	90.67	<b>87.33</b>	74.67	77.33	<b>86.67</b>	<b>83.33</b>	<b>92.00</b>	<b>82.00</b>	86.00
MDA	<b>96.67</b>	<b>96.00</b>	<b>97.33</b>	<b>94.00</b>	<b>94.00</b>	<b>91.33</b>	<b>95.33</b>	<b>94.00</b>	<b>94.00</b>

**VLCS Datasets.** Data: DeCAF<sub>6</sub> features of 4096 dimensions. Domains: V(VOC2007), L(La-belMe), C(Caltech), and S(SUN09). Classes: five classes (bird, car, chair, dog, and person).

Table: Accuracy (%) of VLCS Dataset

Target	V	L	C	S	V,L	V,C	V,S	L,C	L,S	C,S
1NN	60.19	53.57	89.94	55.74	57.26	58.54	50.59	66.06	58.13	66.25
SVM	<b>68.57</b>	59.26	<b>93.99</b>	<b>65.27</b>	<b>61.80</b>	<b>64.39</b>	<b>55.89</b>	70.08	<b>64.10</b>	<b>71.09</b>
KPCA	60.69	54.86	83.89	55.61	57.54	57.50	49.46	67.48	56.05	66.15
KFD	61.64	<b>60.54</b>	86.78	58.75	57.33	46.84	53.20	70.03	61.64	67.87
L-SVM	58.14	39.87	75.56	52.92	52.25	56.64	48.27	61.24	56.65	66.27
CCSA	60.39	58.80	86.88	59.87	59.27	55.02	51.56	69.94	61.41	68.49
DICA	62.71	59.38	86.15	57.28	58.11	55.08	55.17	70.01	61.44	70.30
SCA	62.13	58.24	88.48	<b>60.66</b>	<b>60.66</b>	57.59	54.66	<b>71.90</b>	61.57	70.71
CIDG	64.16	57.91	90.11	59.48	60.54	54.56	55.77	70.74	62.48	69.83
MDA	<b>66.86</b>	<b>61.78</b>	<b>92.64</b>	59.58	59.60	<b>63.72</b>	<b>55.98</b>	<b>72.88</b>	<b>62.83</b>	<b>72.00</b>

## References

- [1] Gilles Blanchard, Gyemin Lee, and Clayton Scott. Generalizing from several related classification tasks to a new unlabeled sample. In *Advances in Neural Information Processing Systems (NIPS)*, pages 2178–2186, 2011.
- [2] Dominik Janzing and Bernhard Schölkopf. Causal inference using the algorithmic markov condition. *IEEE Transactions on Information Theory*, 56(10):5168–5194, 2010.
- [3] Krikamol Muandet, David Balduzzi, and Bernhard Schölkopf. Domain generalization via invariant feature representation. In *Proceedings of the 30th International Conference on Machine Learning (ICML 2013)*, pages 10–18, 2013.