

# Domain Generalization via Invariant Feature Representation

Krikamol Muandet<sup>1</sup>, David Balduzzi<sup>2</sup>, Bernhard Schölkopf<sup>1</sup>

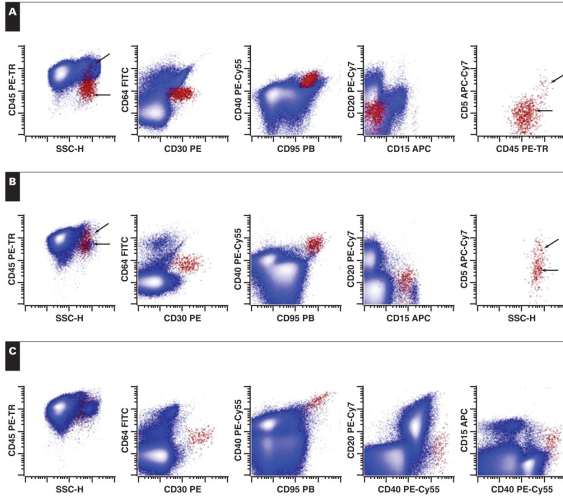
<sup>1</sup>Empirical Inference Department, MPI for Intelligent Systems

<sup>2</sup>Machine Learning Laboratory, ETH Zurich

June 18, 2013

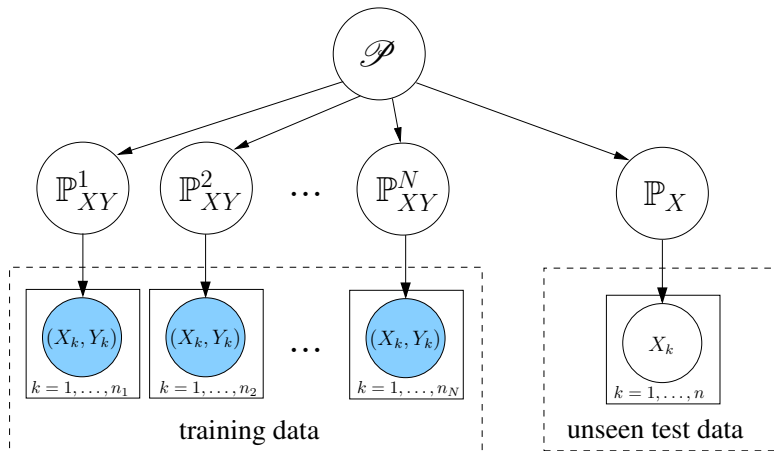
# Flow Cytometry

Image courtesy of American Journal of Clinical Pathology.



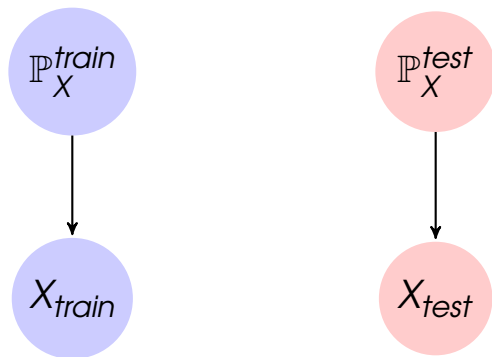
Domains = Patients ( $\mathbb{P}_{XY}$ ), Train Data  $\{X_j^{(i)}, Y_j^{(i)}\}_{j=1}^{n_i}$ .

# Domain Generalization



## Related Works

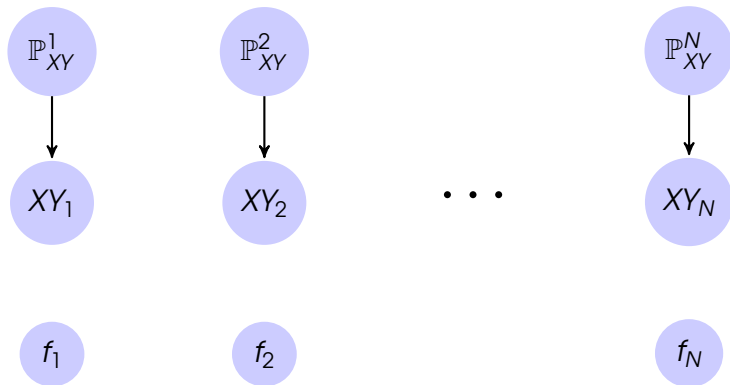
**Domain Adaptation (Bickel, Brückner, and Scheffer 2009)**



Deal with a mismatch between training and test distributions.

# Related Works

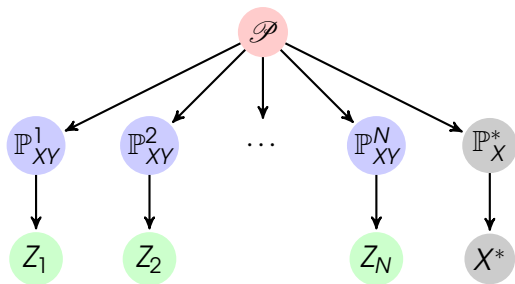
## Multitask Learning (Caruana 1997)



Learn multiple tasks simultaneously.

# Domain Generalization

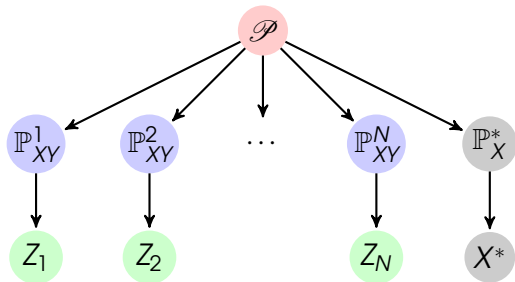
Blanchard, Lee, and Scott 2011



Generalize from multiple source domains to previously unseen domains.

# Domain Generalization

## Problem Setting



**Train:** The joint distributions  $\mathbb{P}^1_{XY}, \mathbb{P}^2_{XY}, \dots, \mathbb{P}^N_{XY} \sim \mathcal{P}$ .

**Prediction:** An unseen distribution  $\mathbb{P}^*_X \sim \mathcal{P}$ .

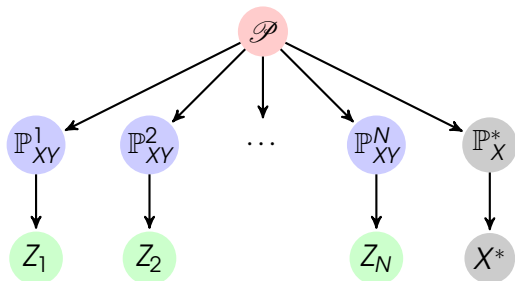
**Goal:** Learn  $f: \mathfrak{X} \times \mathcal{X} \rightarrow \mathcal{Y}$ .

**Assume:**  $\mathbb{P}^1_{Y|X} \approx \mathbb{P}^2_{Y|X} \approx \dots \approx \mathbb{P}^N_{Y|X}$ .

i.e. functional relationship is stable

# Domain Generalization

## Problem Setting



**Train:** The joint distributions  $\mathbb{P}^1_{XY}, \mathbb{P}^2_{XY}, \dots, \mathbb{P}^N_{XY} \sim \mathcal{P}$ .

**Prediction:** An unseen distribution  $\mathbb{P}^*_X \sim \mathcal{P}$ .

**Goal:** Learn  $f: \mathfrak{X} \times \mathcal{X} \rightarrow \mathcal{Y}$ .

**Assume:**  $\mathbb{P}^1_{Y|X} \approx \mathbb{P}^2_{Y|X} \approx \dots \approx \mathbb{P}^N_{Y|X}$ .

i.e. functional relationship is stable

## Domain Adaptation under Target and Conditional Shift

K. Zhang, B. Schölkopf, K. Muandet, and Z. Wang (ICML2013)



# Objective

Find feature representation,  $\mathcal{B}(X)$  that is *invariant* across domains.

① minimize the distance between empirical distributions  $\hat{\mathbb{P}}_X^1, \hat{\mathbb{P}}_X^2, \dots, \hat{\mathbb{P}}_X^N$  of the transformed samples  $\mathcal{B}(X)$ .

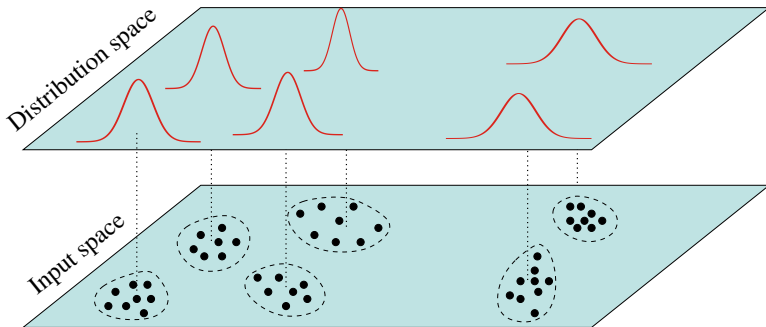
$$\mathbb{P}_{Y|X}^1 \cdot \mathbb{P}_X^1 \quad \mathbb{P}_{Y|X}^2 \cdot \mathbb{P}_X^2 \quad \dots \quad \mathbb{P}_{Y|X}^N \cdot \mathbb{P}_X^N$$

② preserve functional relationship between  $X$  and  $Y$ .

# Minimizing Distributional Variance

## Hilbert space embedding

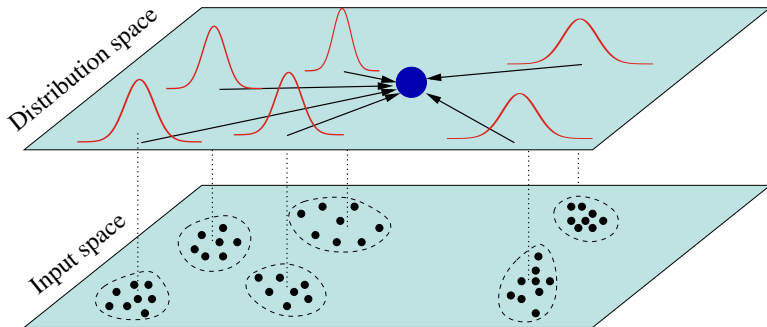
$$\mu : \mathfrak{P}_{\mathcal{X}} \rightarrow \mathcal{H}, \quad \mathbb{P} \mapsto \int_{\mathcal{X}} k(x, \cdot) d\mathbb{P}(x) =: \mu_{\mathbb{P}}.$$



# Minimizing Distributional Variance

Find transformation  $\mathcal{B}$  that minimizes

$$\mathbb{V}_{\mathcal{H}}(\mathcal{P}) = \frac{1}{N} \sum_{i=1}^N \|\mu_i \mathcal{B} - \bar{\mu} \mathcal{B}\|_{\mathcal{H}}^2$$



# Minimizing Distributional Variance

- ▶ Minimizing distributional variance **alone** does not necessarily help with generalization!
  - ▶ Setting  $\mathcal{B} = \mathbf{0}$  gives zero distributional variance!
- ▶ We **also** need to preserve the functional relationship between  $X$  and  $Y$  encoded in  $\mathbb{P}_{Y|X}$ .

# Preserving Functional Relationship

## Central Subspace

The central subspace  $C$  is the minimal subspace that captures the functional relationship between  $X$  and  $Y$ , i.e.  $Y \perp\!\!\!\perp X|C^\top X$ .

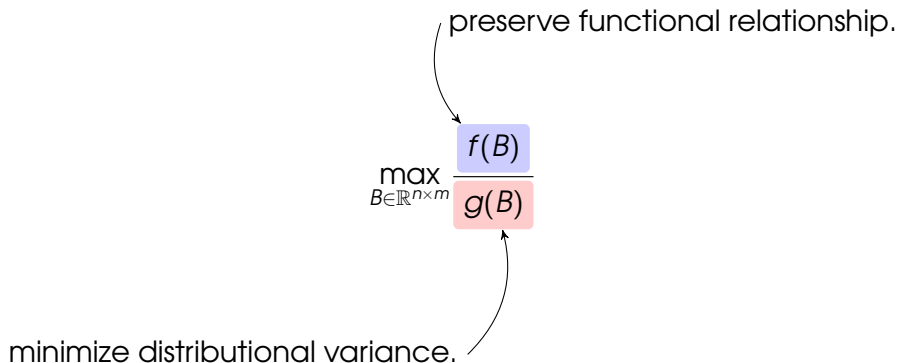
**Theorem** (Li 1991; Kim and Pavlovic 2011; Muandet 2013)

*If  $B$  maximizes*

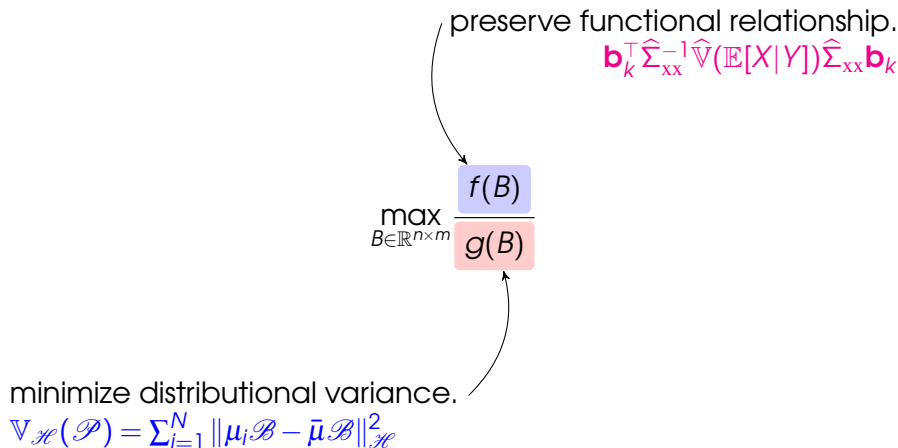
$$\mathbf{b}_k^\top \Sigma_{xx}^{-1} \mathbb{V}(\mathbb{E}[X|Y]) \Sigma_{xx} \mathbf{b}_k$$

*then  $Y \perp\!\!\!\perp X|B^\top X$ .*

# Domain-Invariant Component Analysis



# Domain-Invariant Component Analysis



# Domain-Invariant Component Analysis

preserve functional relationship.

$$\mathbf{b}_k^\top \hat{\Sigma}_{xx}^{-1} \hat{V}(\mathbb{E}[X|Y]) \hat{\Sigma}_{xx} \mathbf{b}_k$$

$$\max_{B \in \mathbb{R}^{n \times m}} \frac{\frac{1}{n} \text{tr}(B^\top L(L + n\epsilon I_n)^{-1} K^2 B)}{\text{tr}(B^\top K Q K B + B K B)}$$

minimize distributional variance.

$$\mathbb{V}_{\mathcal{H}}(\mathcal{P}) = \sum_{i=1}^N \|\mu_i \mathcal{B} - \bar{\mu} \mathcal{B}\|_{\mathcal{H}}^2$$



# Domain-Invariant Component Analysis

## Maximization Problem

$$\max_{B \in \mathbb{R}^{n \times m}} \frac{\frac{1}{n} \text{tr}(B^T L (L + n\epsilon I_n)^{-1} K^2 B)}{\text{tr}(B^T K Q K B + B K B)}$$



## Generalized Eigenvalue Problem

$$\frac{1}{n} L (L + n\epsilon I)^{-1} K^2 B = (K Q K + K + \lambda I) B \Gamma$$

# Learning guarantee

## Theorem

*Under reasonable assumptions, it holds with probability at least  $1 - \delta$  that,*

$$\mathbb{E}[\text{error}] \leq c_1 \mathbb{V}_{\mathcal{H}}(\mathcal{P} \cdot \mathcal{B}) + L(n, N) .$$

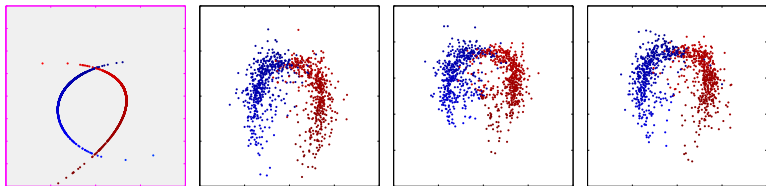
- ▶ Bound depends on the distributional variance.
- ▶  $L(n, N) \rightarrow 0$  as samples  $n$  and domains  $N$  go to infinity.

# Experimental Results

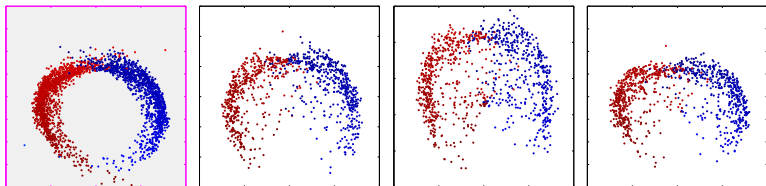
## Synthetic Data

- ▶ Generate 10 collections of  $n_i \sim \text{Poisson}(200)$  data points.
- ▶ For each collection,  $x \sim \mathcal{N}(\mathbf{0}, \Sigma_i)$  where  $\Sigma_i \sim \mathcal{W}(0.2 \times I_5, 10)$ .
- ▶ The output value is  $y = \text{sign}(b_1^\top x + \varepsilon_1) \cdot \log(|b_2^\top x + c + \varepsilon_2|)$ , where  $\varepsilon_1, \varepsilon_2 \sim \mathcal{N}(0, 1)$ .

## Experimental Results: synthetic data



COIR



DICA

# Experimental Results

## Real-world Data

- ▶ Flow cytometry dataset (classification).
- ▶ Parkinson's telemonitoring dataset (regression).

## Learning algorithms

- ▶ **Pooling SVM** : pool data from all domains and apply standard SVM.
- ▶ **Distributional SVM** : apply the kernel

$$\kappa((\mathbb{P}^i, x_k^i), (\mathbb{P}^j, x_l^j)) = K(\mathbb{P}^i, \mathbb{P}^j) \cdot k(x_k^i, x_l^j)$$

(Blanchard, Lee, and Scott 2011).

## Experimental Results: Flow cytometry

Methods	Pooling SVM	Distributional SVM
Input	92.03 $\pm$ 8.21	93.19 $\pm$ 7.20
KPCA	91.99 $\pm$ 9.02	93.11 $\pm$ 6.83
COIR	92.40 $\pm$ 8.63	92.92 $\pm$ 8.20
UDICA	92.51 $\pm$ 5.09	92.74 $\pm$ 5.01
DICA	<b>92.72<math>\pm</math>6.41</b>	<b>94.80<math>\pm</math>3.81</b>

Similar results for Parkinson's telemonitoring dataset.

# Conclusion

Domain-Invariance Component Analysis (DICA) finds an **invariant representation** that

- ▶ minimizes “differences” between domains
- ▶ while preserving discriminative information.

To learn more, please come to our poster!

Thank you!