# Opinion Mining by Transformation-Based Domain Adaptation
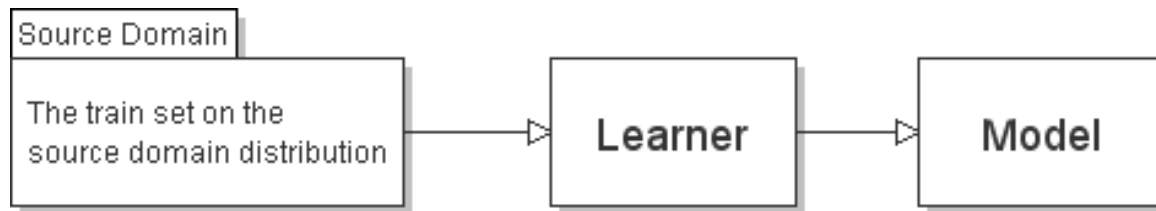
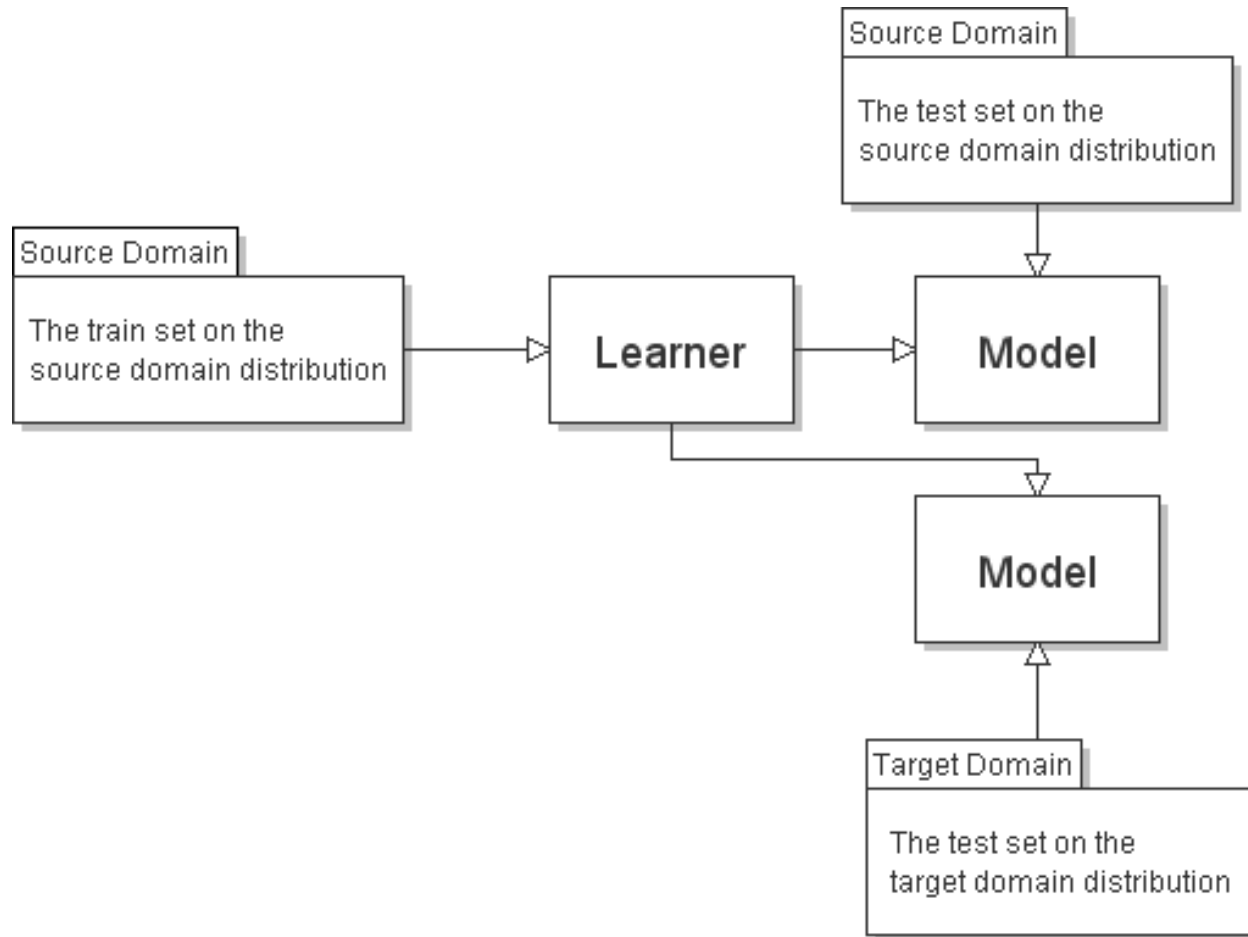Róber Ormándi, István Hegedűs, Richárd Farkas

TSD-2010

# The Domain Adaptation Task

- Given two datasets from different domains, namely the source (S) and the target (T) domains
- The source has a huge, the target often has just a few amount of labeled samples ($|S|>>|T|$)
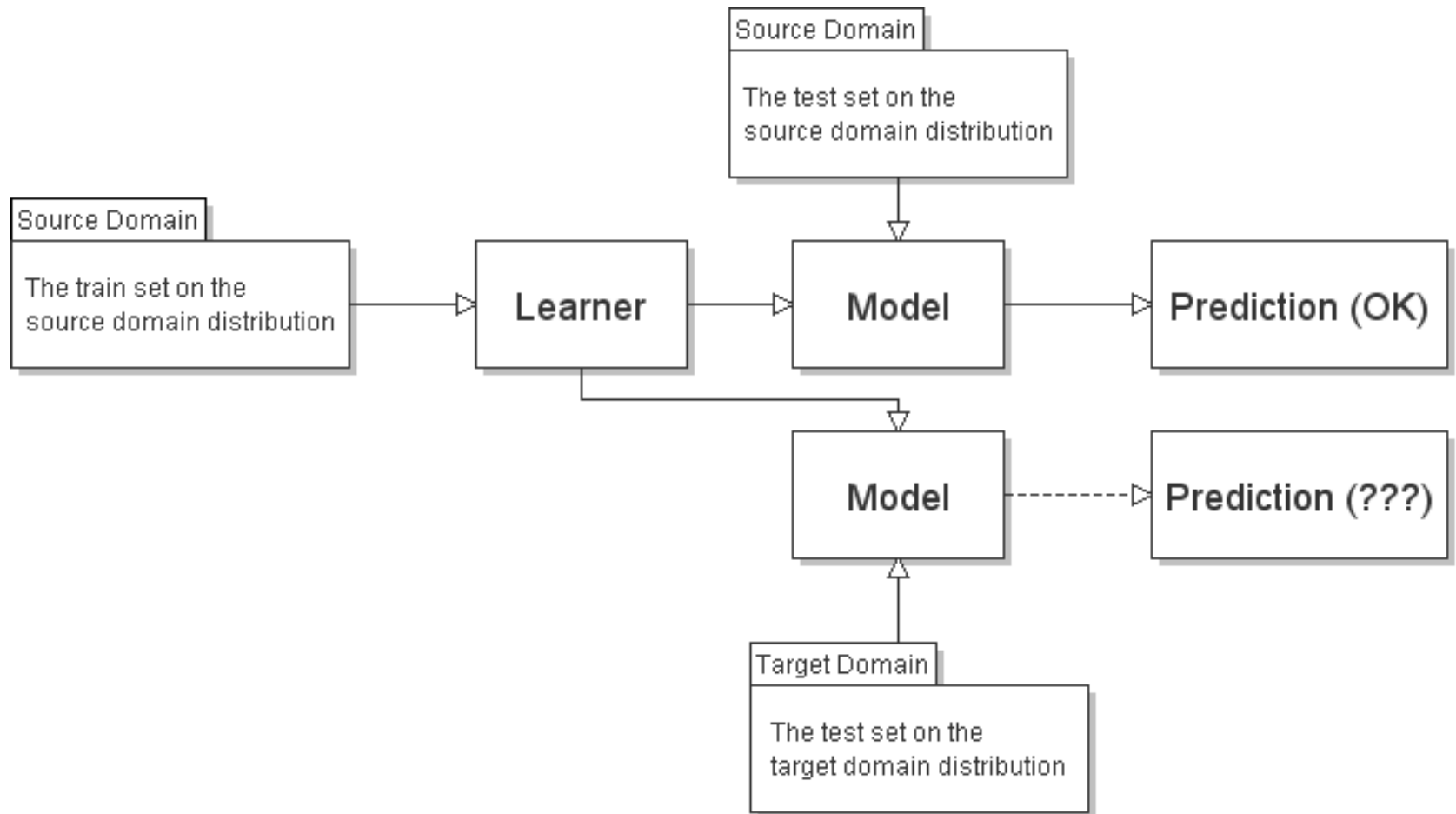- The task is to set the labels of the target domain as precisely as possible, using the labeled source samples
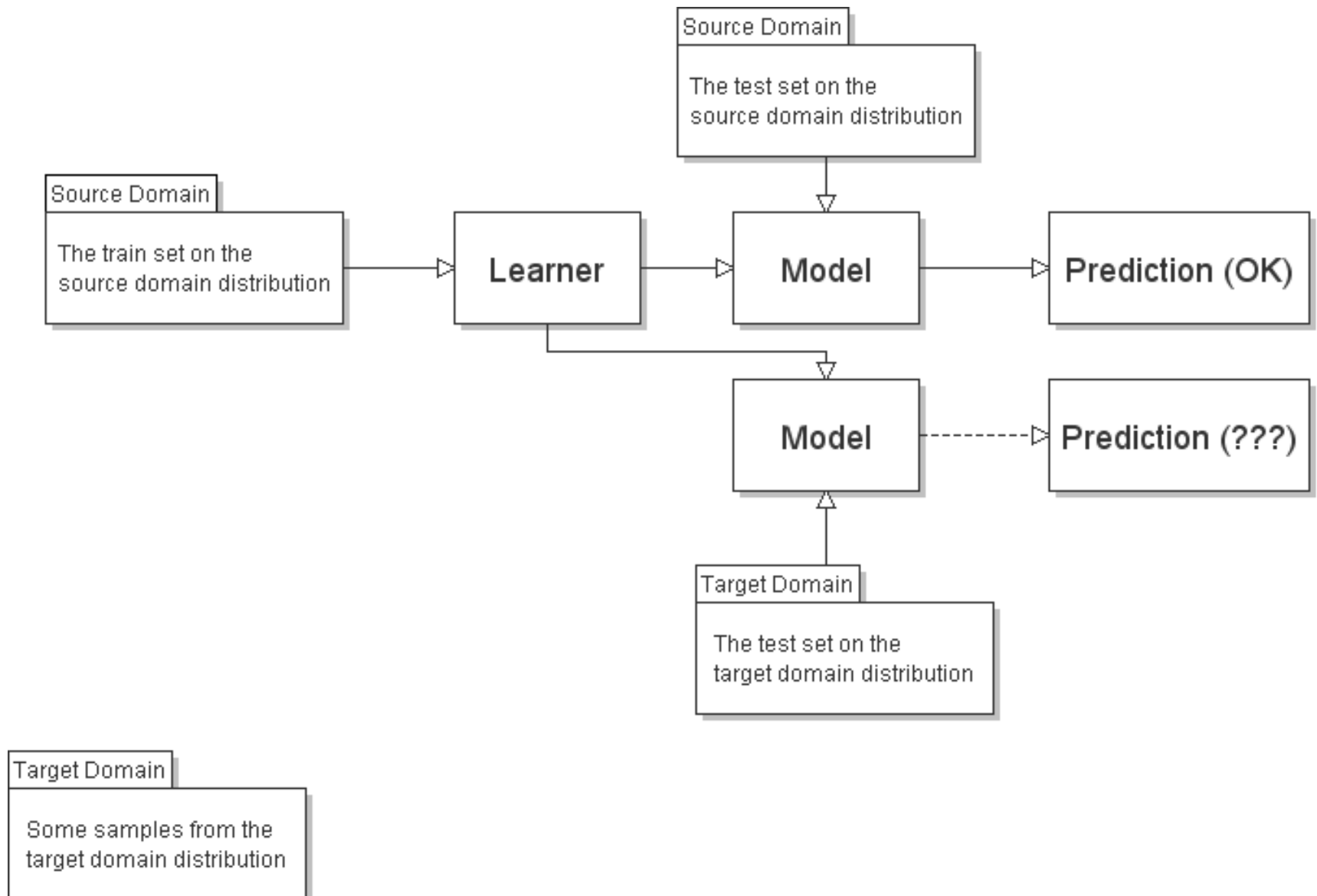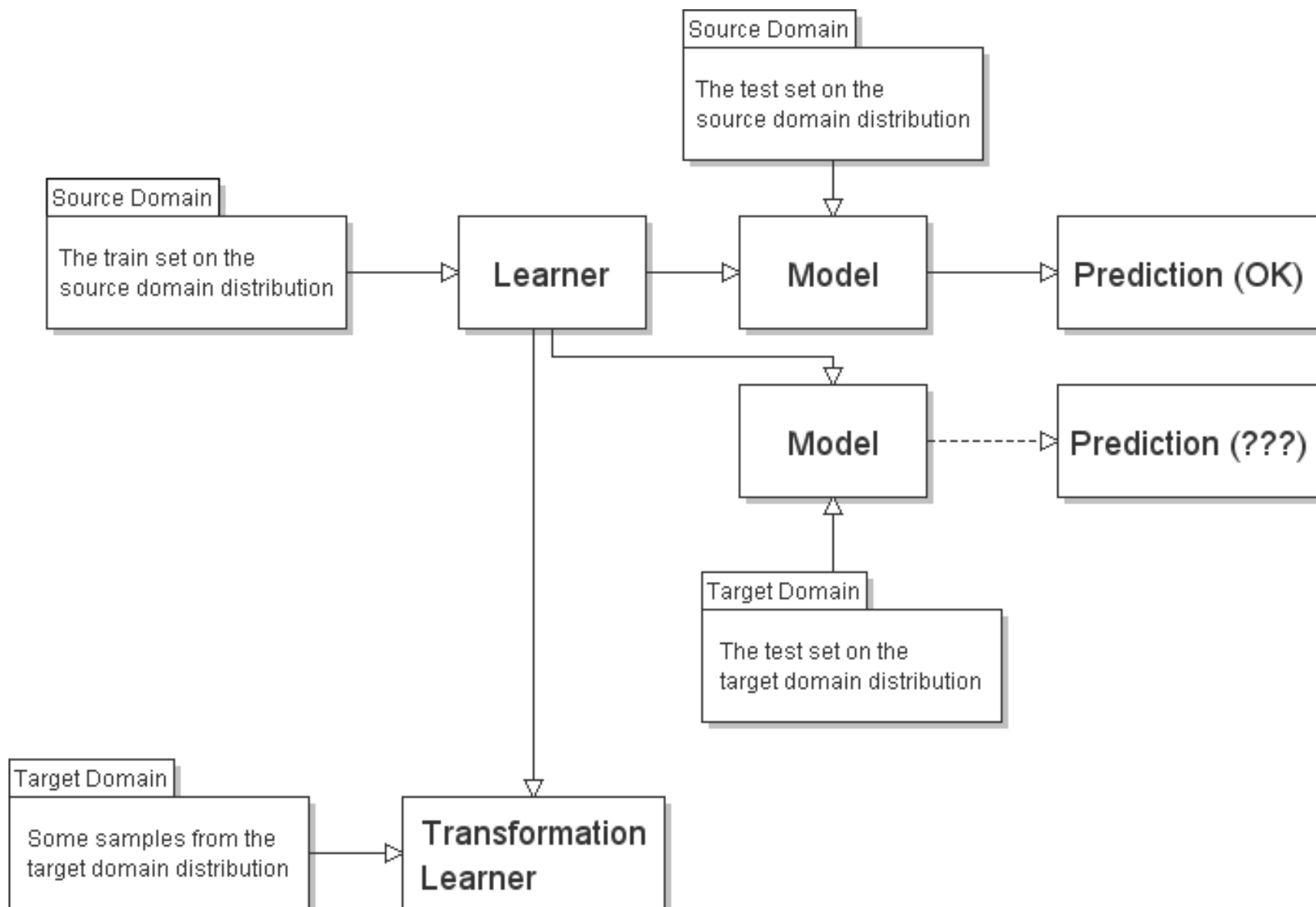
# Basic Machine Learning

Source Domain

The train set on the source domain distribution → Learner → Model

# Basic Machine Learning



Source Domain
The test set on the source domain distribution

Source Domain
The train set on the source domain distribution

Learner

Model

Model

Target Domain
The test set on the target domain distribution

# Basic Machine Learning

# Domain Adaptation

# Our Approach

# Our Approach

# Our Approach



Source Domain
The test set on the source domain distribution

Source Domain
The train set on the source domain distribution

Learner

Model

Prediction (OK)

Model

Prediction (???)

Target Domain
The test set on the target domain distribution

Target Domain
Some samples from the target domain distribution

Transformation Learner

Transformed Model

Prediction (OK)
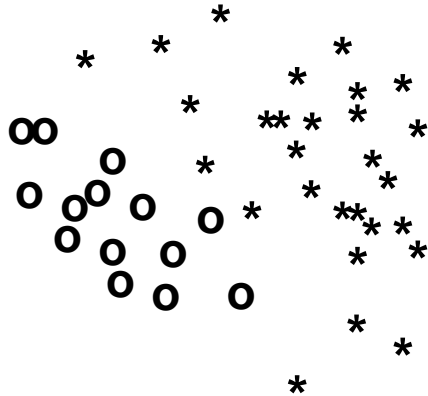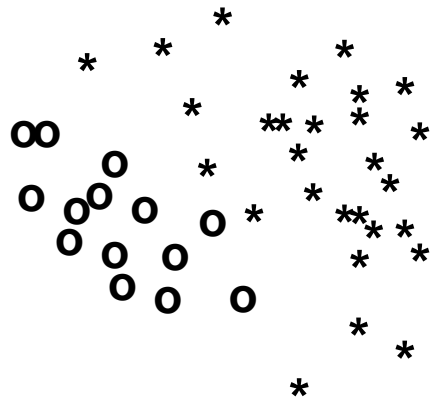
# Transformation Based Domain Adaptation

Source Domain

# Transformation Based Domain Adaptation
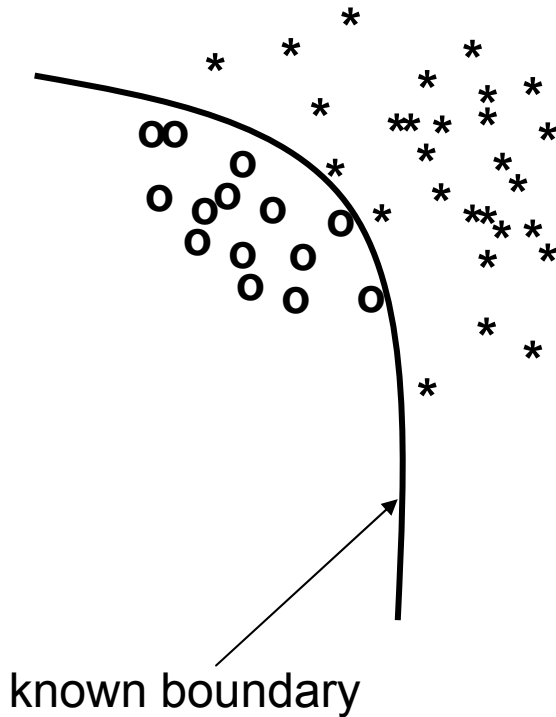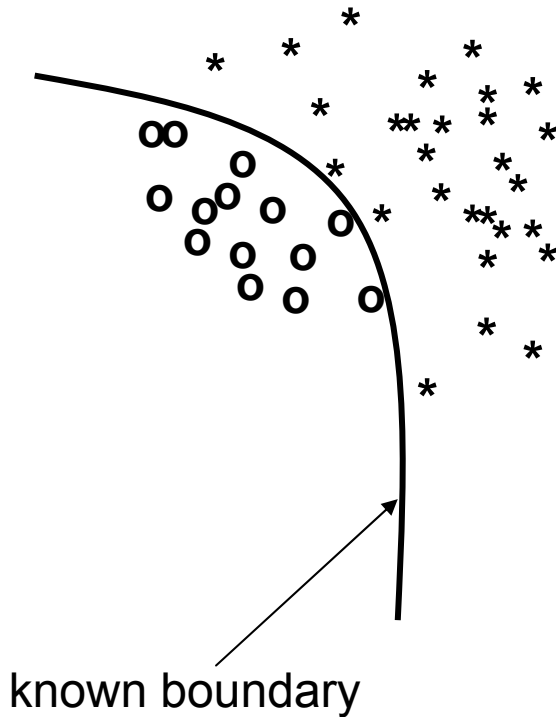
## Source Domain

use a machine learning method

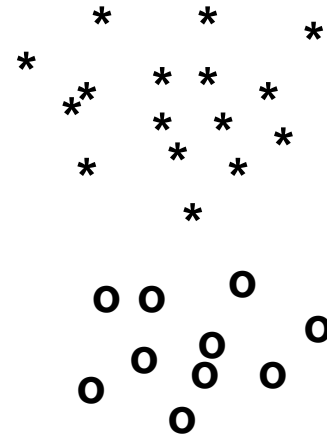# Transformation Based Domain Adaptation

<u>Source Domain</u>



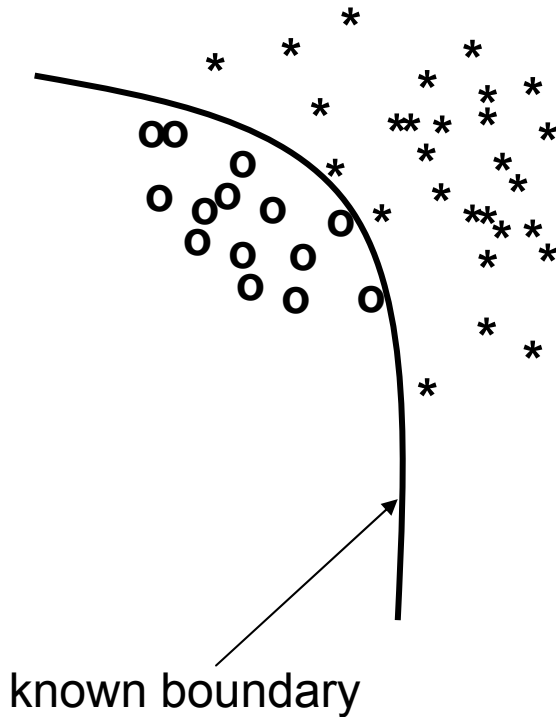known boundary

# Transformation Based Domain Adaptation

## Source Domain

## Target Domain

known boundary

# Transformation Based Domain Adaptation

## Source Domain



known boundary

## Target Domain



unknown boundary

# Transformation Based Domain Adaptation

**Source Domain**

$\Phi$ transformation

**Target Domain**

known boundary

unknown boundary

# Transformation Based Domain Adaptation



Source Domain

$\Phi$ transformation

Target Domain

known boundary

unknown boundary

Task: to minimize this error function

$$E_{T,p_{\mathcal{D}_S}}(W) = \frac{1}{2} \sum_{x \in T} (t(x) - p_{\mathcal{D}_S}(Wx))^2$$

# Transformation Based Domain Adaptation

**Source Domain**

**Φ** transformation

**Target Domain**

* * *
* * *
* * *
* * * * *
* * * * * *
∞
o o o * * * * *
o o o o o *
o o o o *
o o o o *
* * *
* * *
*
*

unknown boundary

known boundary

Task: to minimize this error function
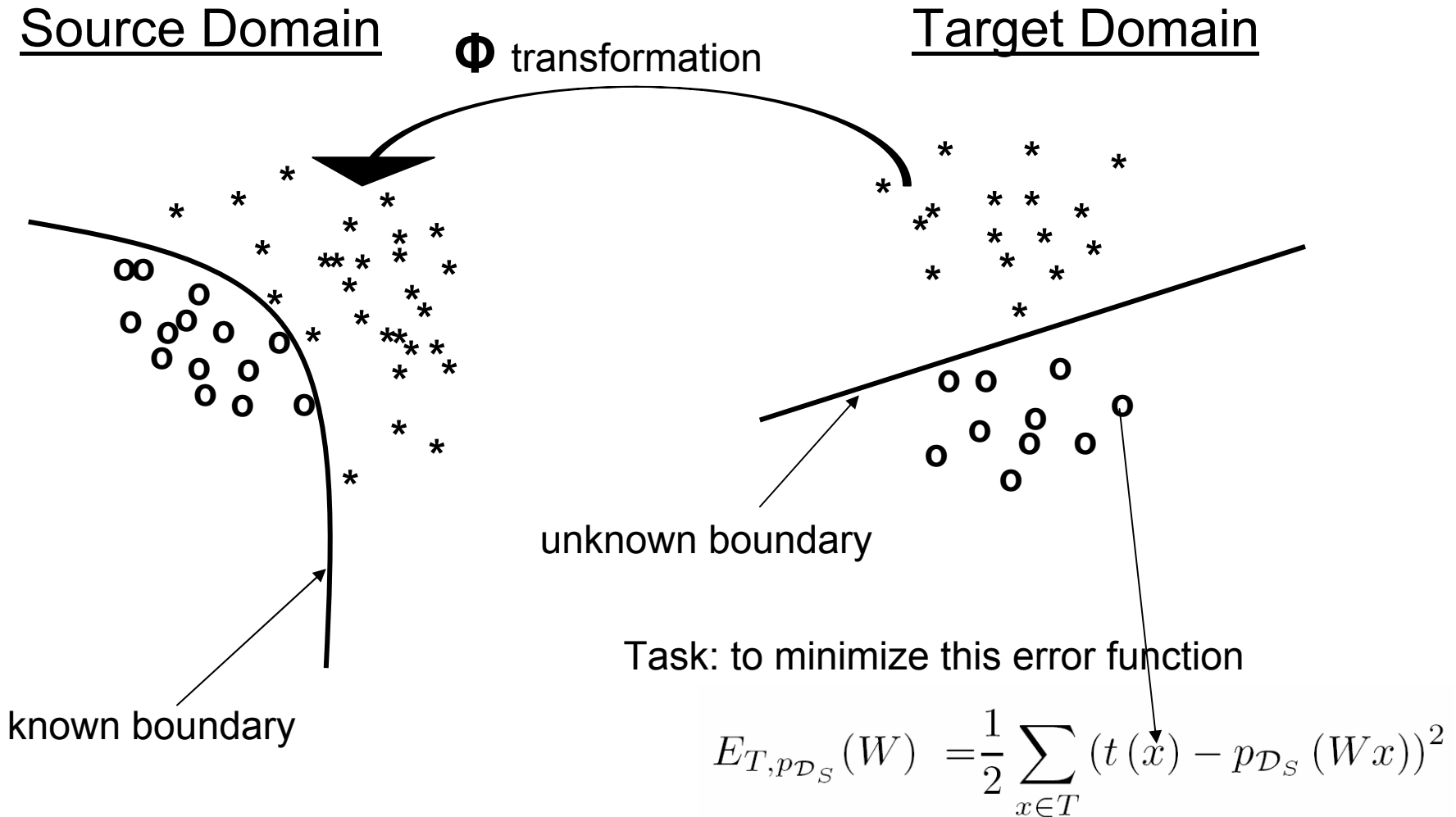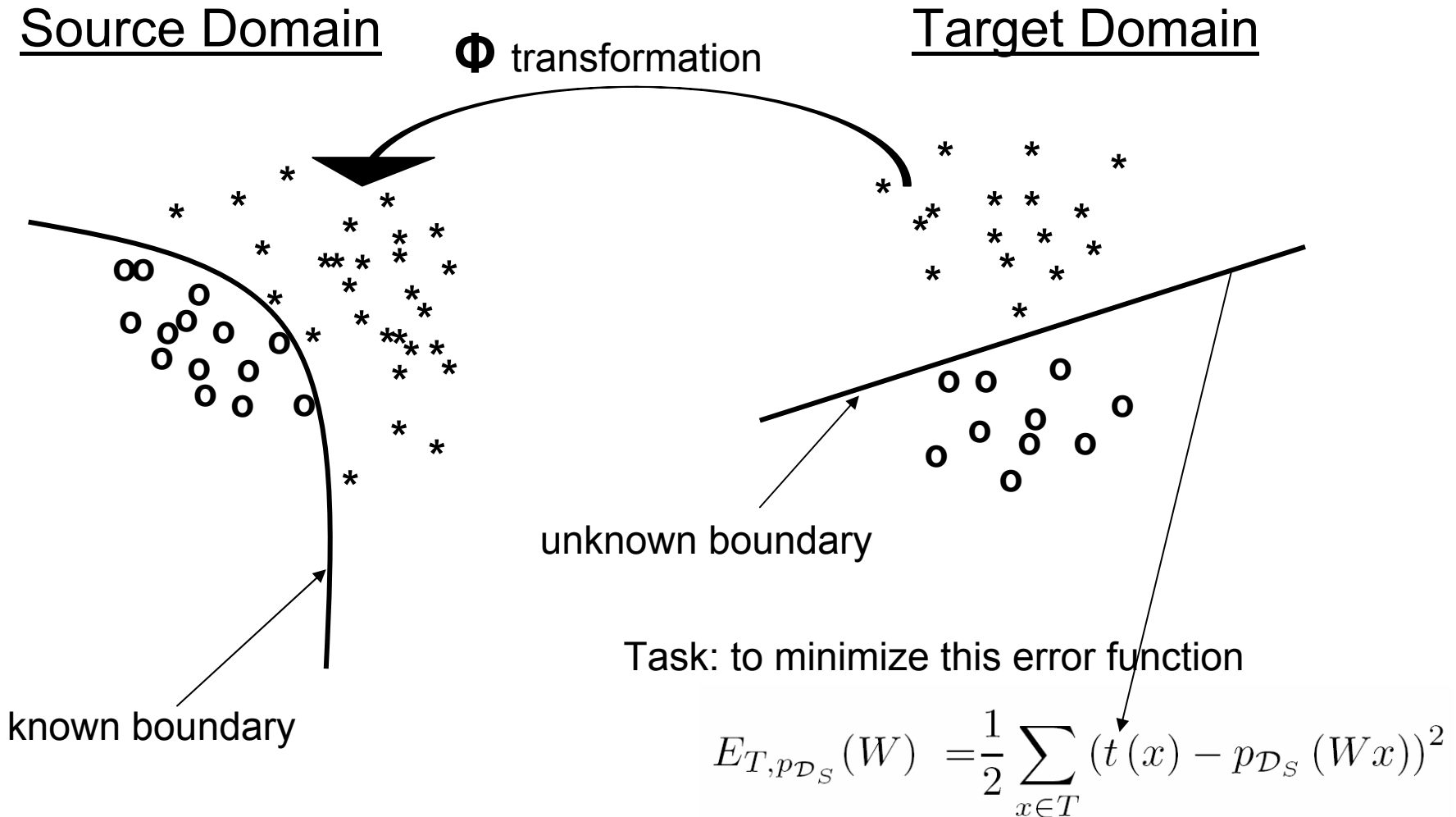
$$E_{T, p_{\mathcal{D}_S}}(W) = \frac{1}{2} \sum_{x \in T} (t(x) - p_{\mathcal{D}_S}(Wx))^2$$

# Transformation Based Domain Adaptation

**Source Domain**     **Φ** transformation     **Target Domain**

unknown boundary

known boundary

Task: to minimize this error function

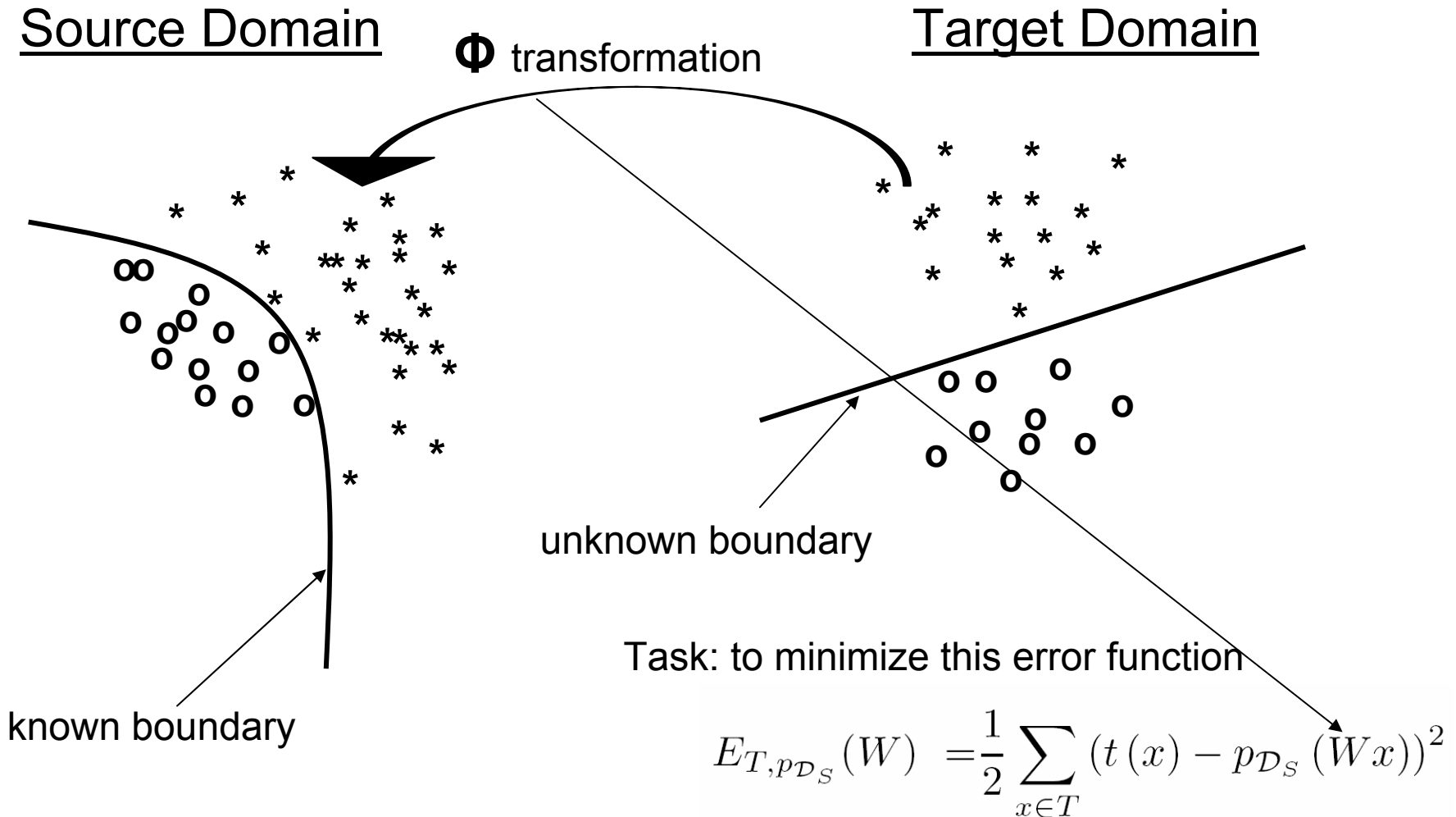$$E_{T, p_{\mathcal{D}_S}}(W) = \frac{1}{2} \sum_{x \in T} \left( t(x) - p_{\mathcal{D}_S}(Wx) \right)^2$$

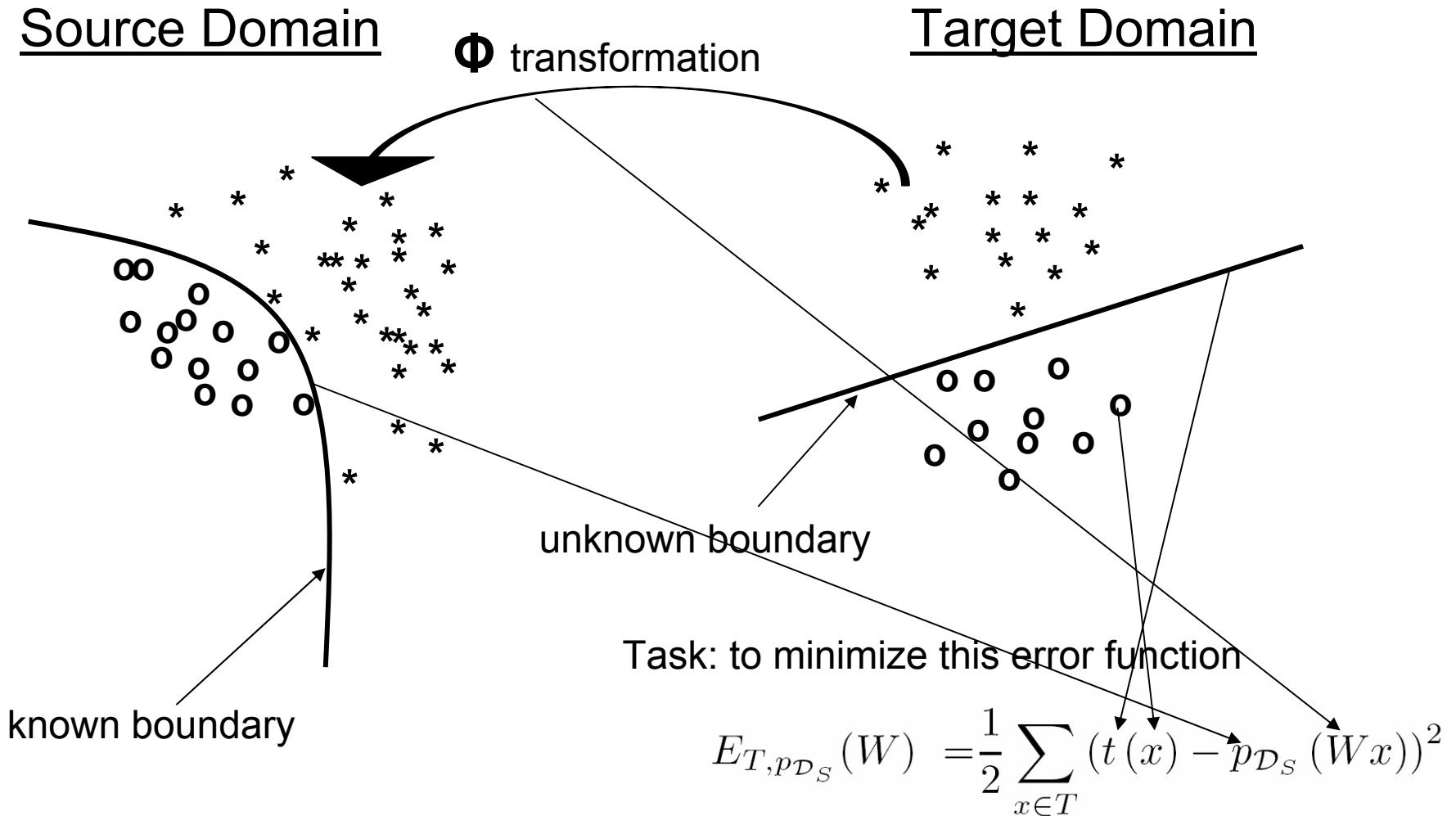# Transformation Based Domain Adaptation

Source Domain

Target Domain

$\Phi$ transformation

unknown boundary

known boundary

Task: to minimize this error function

$$E_{T, p_{\mathcal{D}_S}}(W) = \frac{1}{2} \sum_{x \in T} \left( t(x) - p_{\mathcal{D}_S}(Wx) \right)^2$$

# Transformation Based Domain Adaptation



Source Domain

Φ transformation

Target Domain

unknown boundary

known boundary

Task: to minimize this error function

$$E_{T,p_{\mathcal{D}_S}}(W) = \frac{1}{2}\sum_{x\in T}(t(x) - p_{\mathcal{D}_S}(Wx))^2$$

# Transformation Based Domain Adaptation

Source Domain

$\Phi$ transformation

Target Domain

unknown boundary

known boundary

Task: to minimize this error function

$$E_{T, p_{\mathcal{D}_S}}(W) = \frac{1}{2} \sum_{x \in T} (t(x) - p_{\mathcal{D}_S}(Wx))^2$$
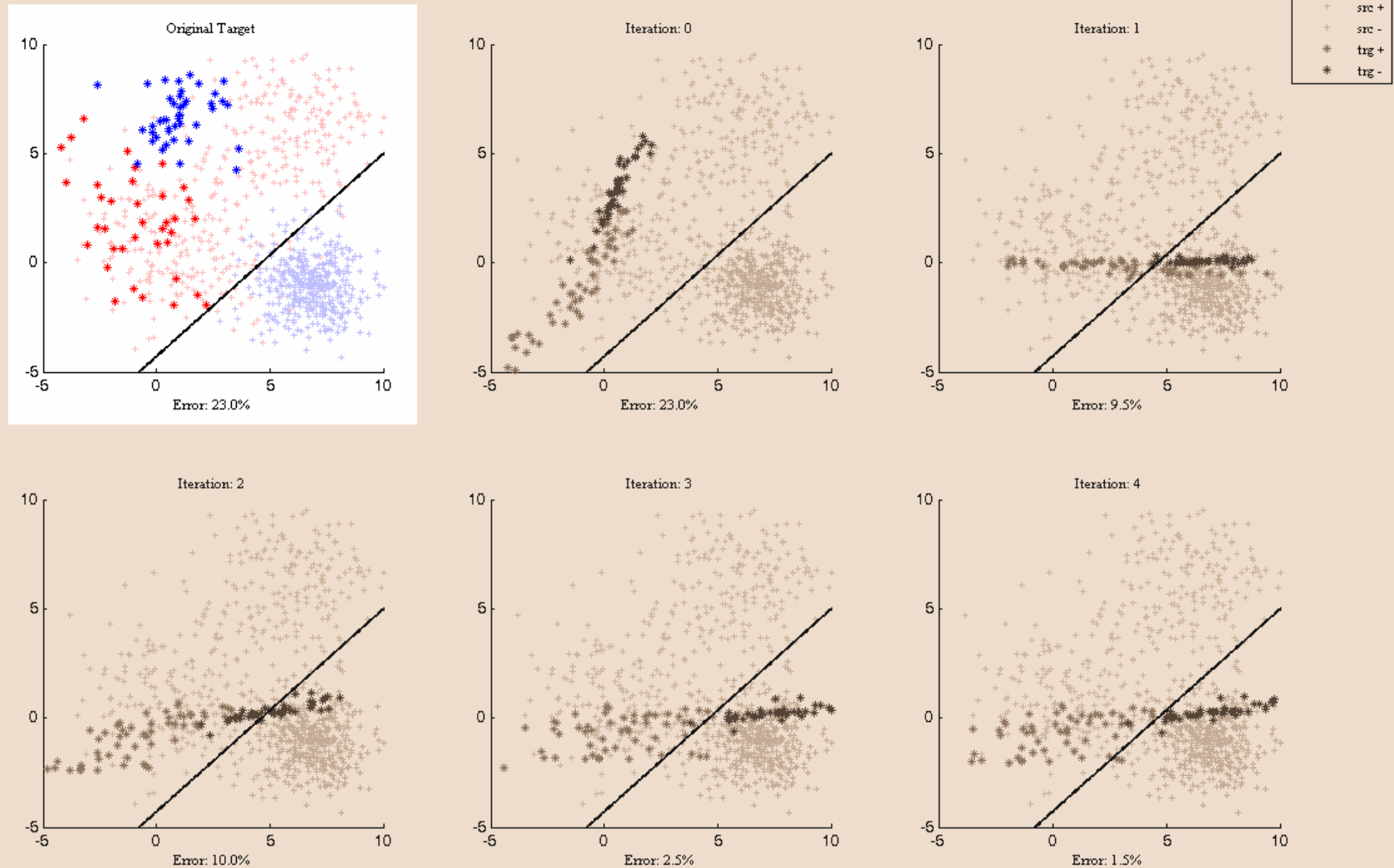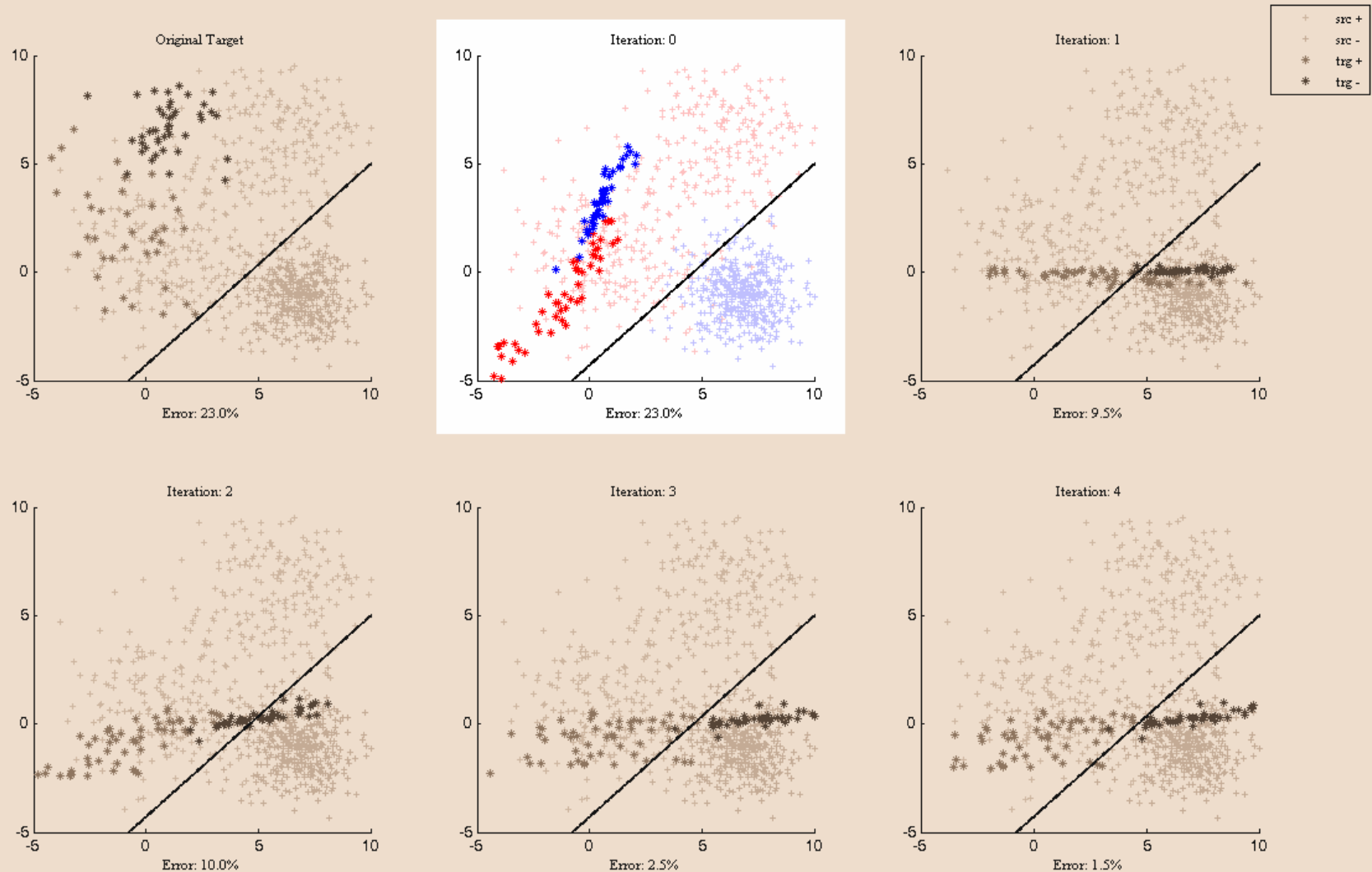
# Synthetic Database



- 2D points and 2 classes
- Source train/test 800/200 samples
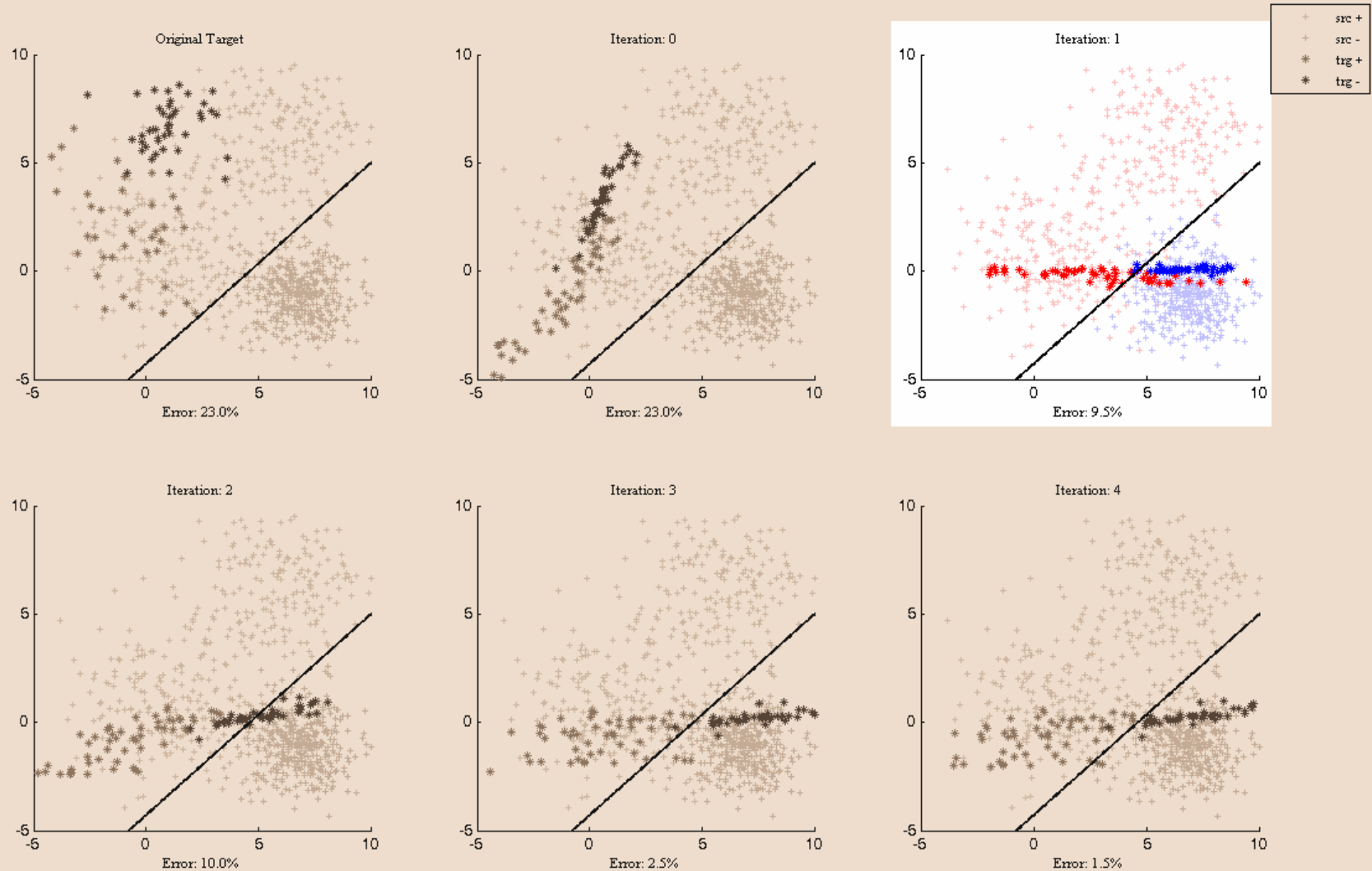- Target distribution is the „rotated" source distribution (by 90°)
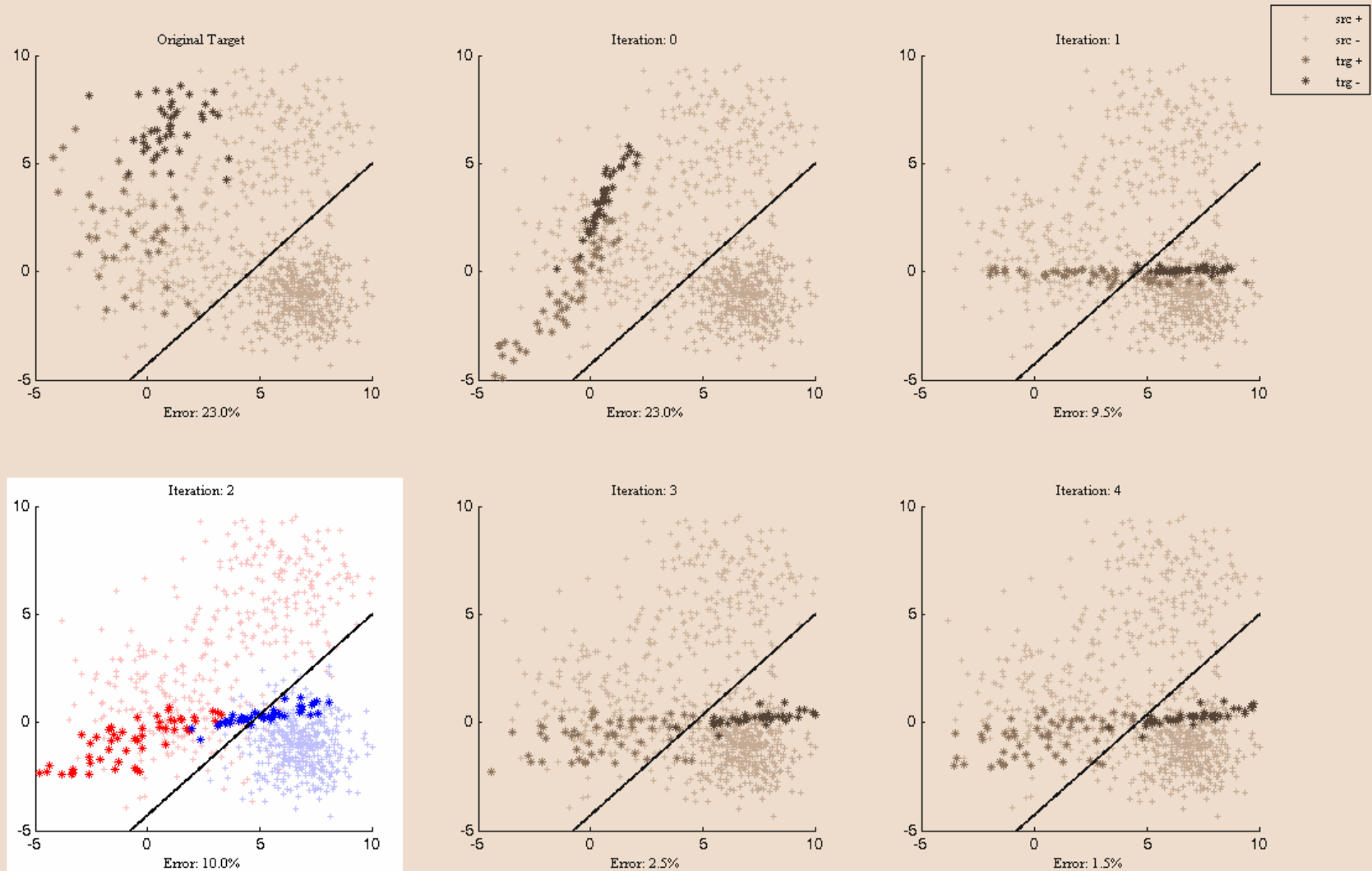
# Results on Synthetic Dataset
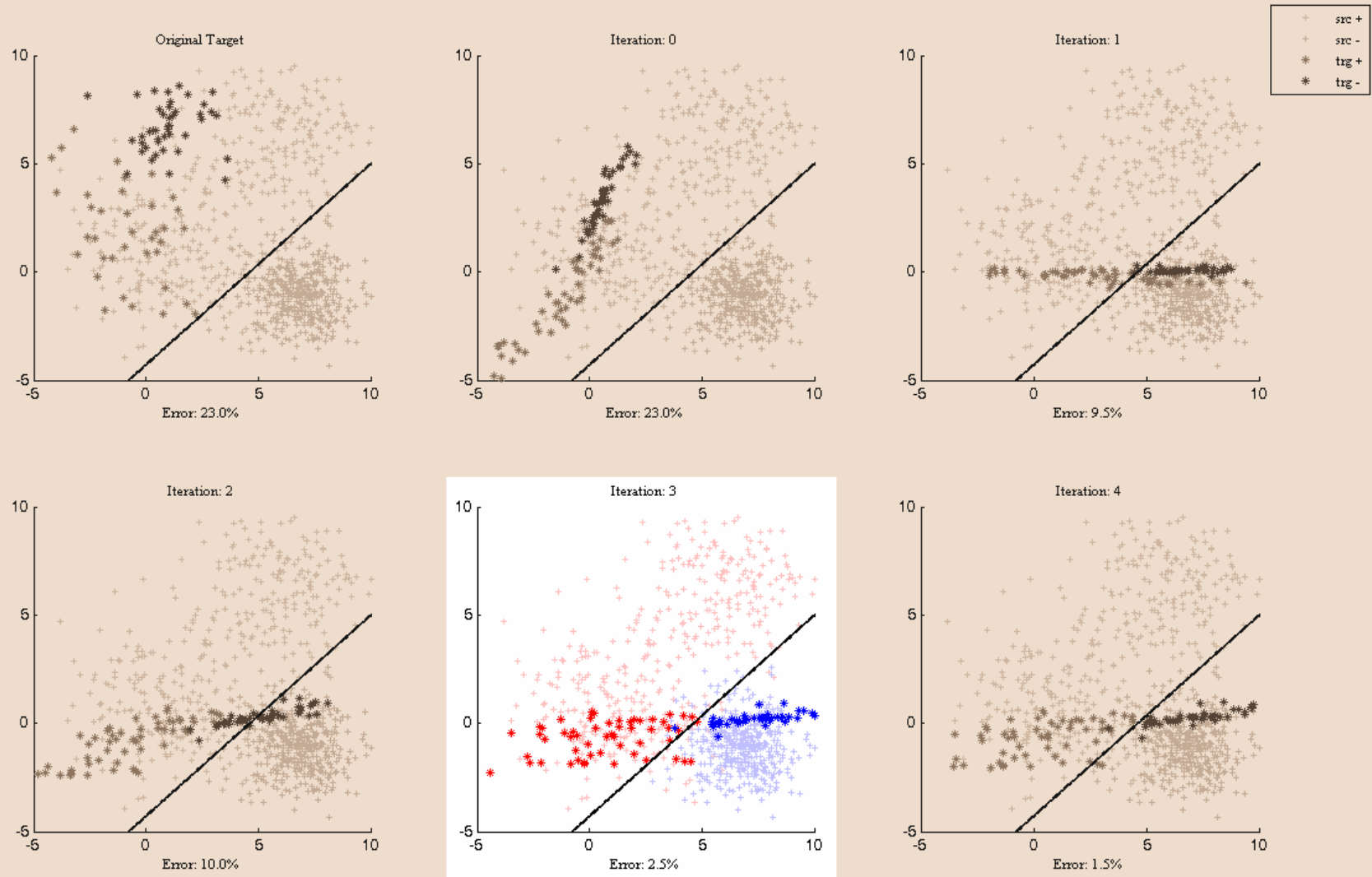
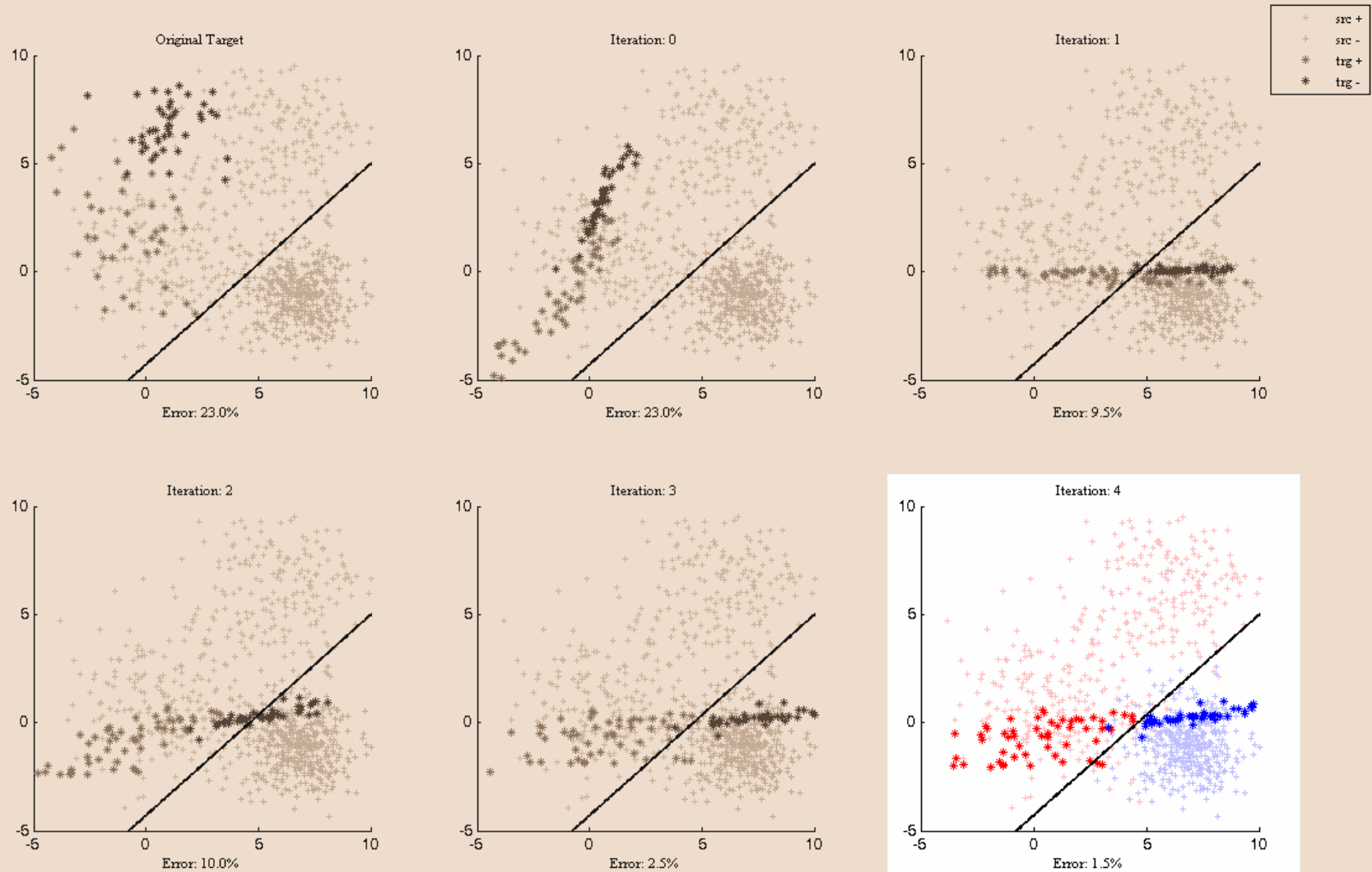# Results on Synthetic Dataset

# Results on Synthetic Dataset

# Results on Synthetic Dataset

# Results on Synthetic Dataset

# Results on Synthetic Dataset
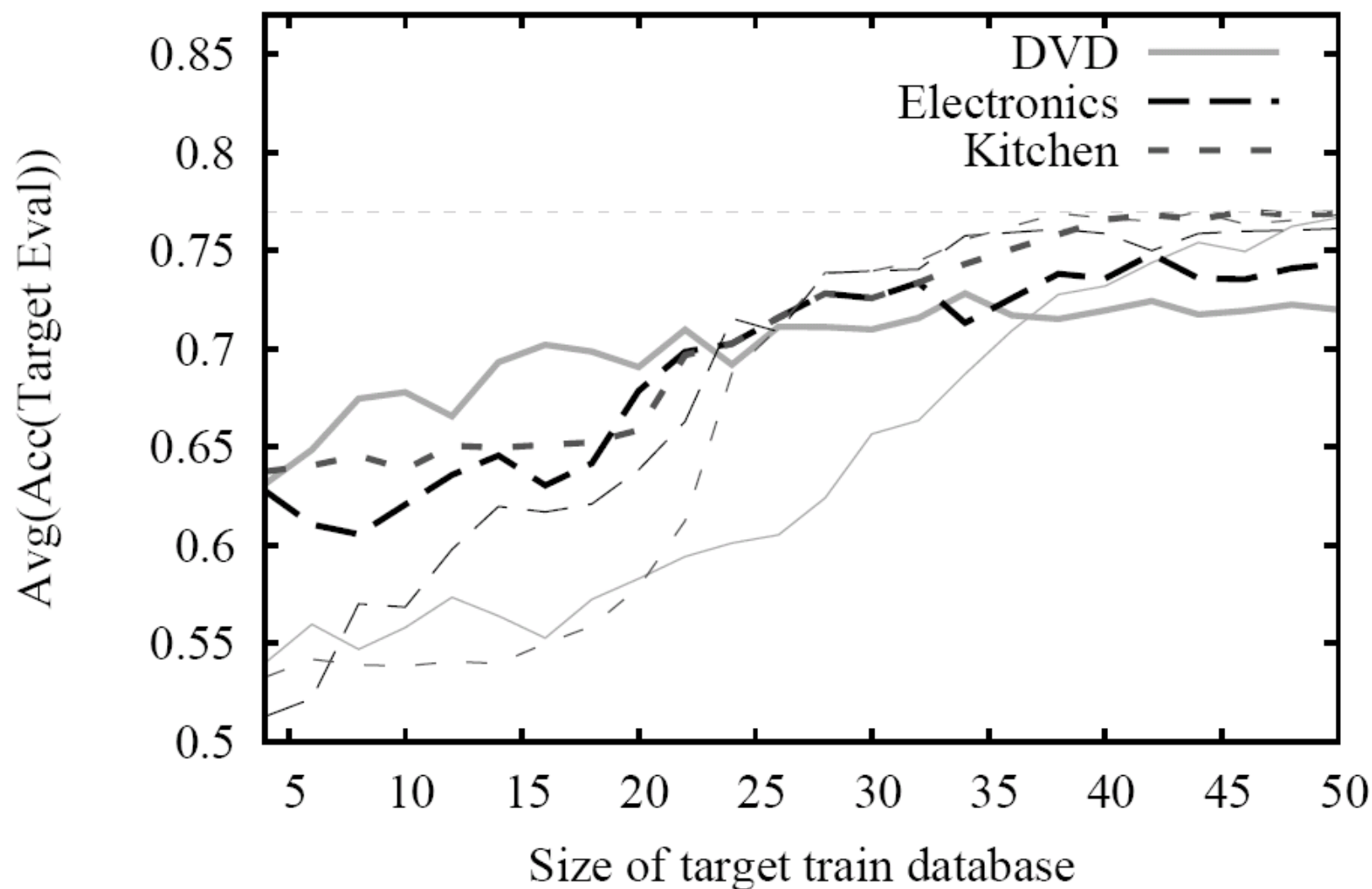
# Multi-Domain Sentiment Database

- Reviews for 4 types of product/domain from the Amazon.com
  - Books
  - DVDs
  - Electronics
  - Kitchen appliances
- Balanced 1000-1000 positive and negative samples in ~ 5000 dimensions

# PreProcessing steps

- Train/test cut, 80/20% randomly
- Feature reduction by InfoGain score and PCA
  - we held the feature with > 0 score
  - compressed the dataset by PCA
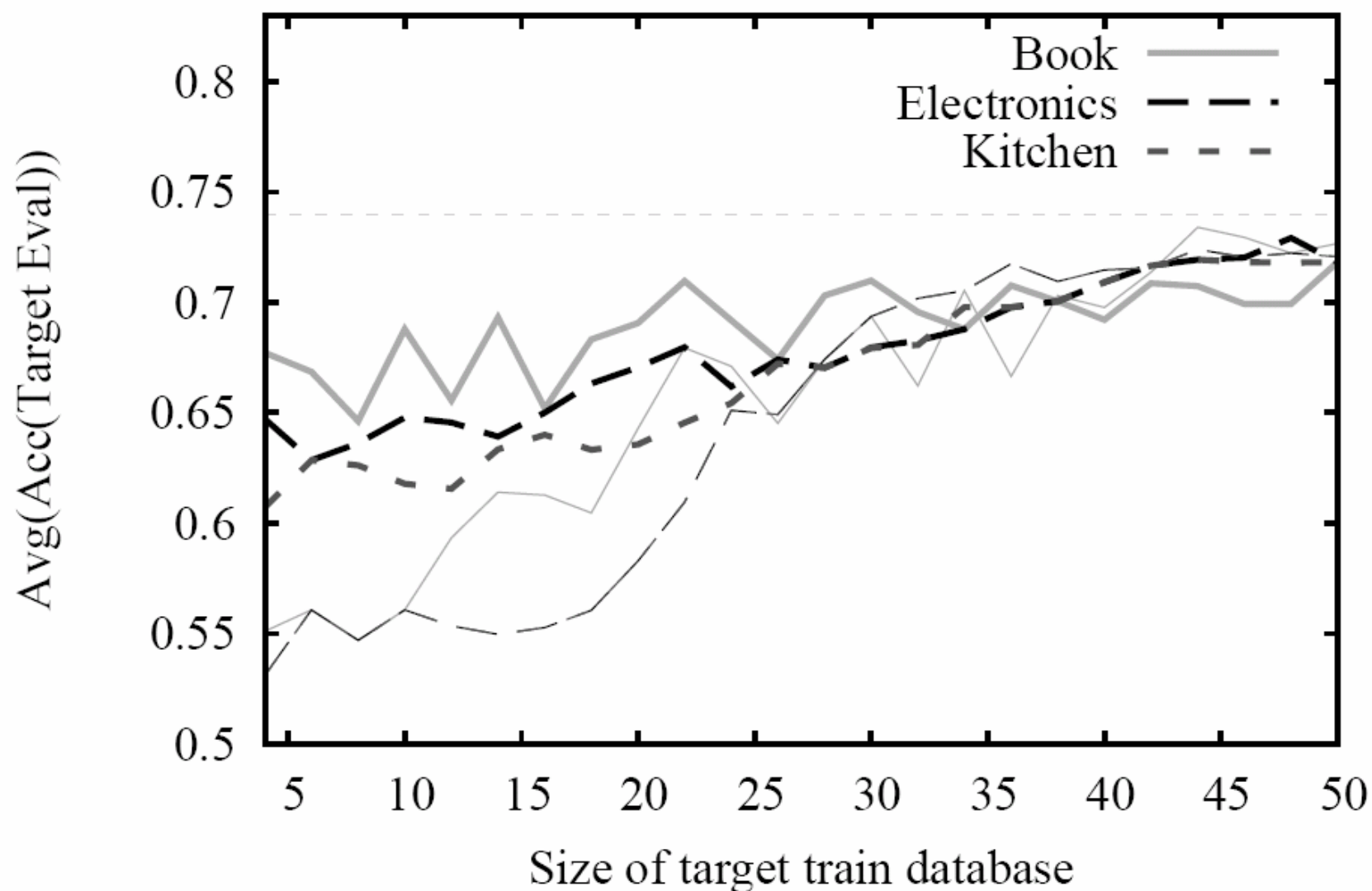- The results shows the average accuracy score of 10 runs
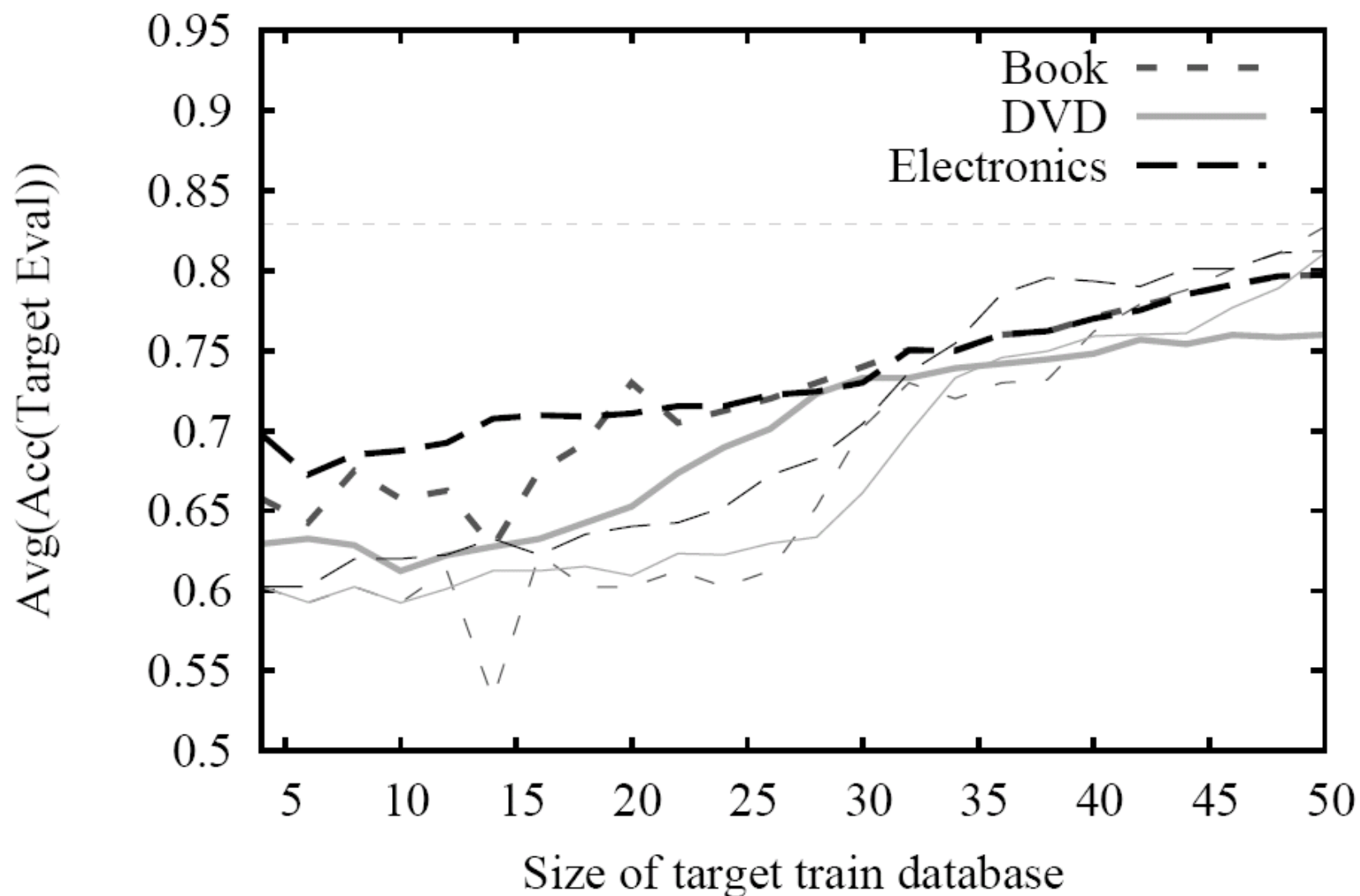
# Results



Book database as target domain

# Results



DVD database as target domain

# Results



Kitchen database as target domain

# Results



Electronics database as target domain