

# Robust Domain Adaptation: Representations, Weights and Inductive Bias

---

V. Bouvier, P. Very, C. Chastagnol, M. Tami and C. Hudelot

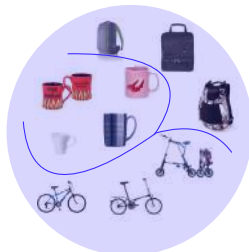


SIDETRADE

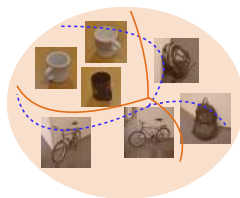
université  
PARIS-SACLAY

June 13, 2021

# Domain Adaptation



**SOURCE**  
Labelled samples



**TARGET**  
Unlabelled samples

## Setup

- **Source labelled samples:**  $(x_S^i, y_S^i)_i \sim p_S(X, Y)$
- **Target unlabelled samples:**  $(x_T^j)_j \sim p_T(X)$
- **Objective:** Learning  $h \in \mathcal{H}$  s.t.:  $h \in \arg \min_{h \in \mathcal{H}} \mathcal{E}_T(h)$

# Covariate Shift

## Covariate Shift

Labelling functions are conserved:

$$p_S(y|x) = p_T(y|x) \quad (1)$$

# Covariate Shift

## Covariate Shift

Labelling functions are conserved:

$$p_S(y|x) = p_T(y|x) \quad (1)$$

$$\epsilon_T(h) = \mathbb{E}_T[\ell(h(x), y)] = \mathbb{E}_S \left[ \frac{p_T(x)}{p_S(x)} \ell(h(x), y) \right] \quad (2)$$

# Covariate Shift

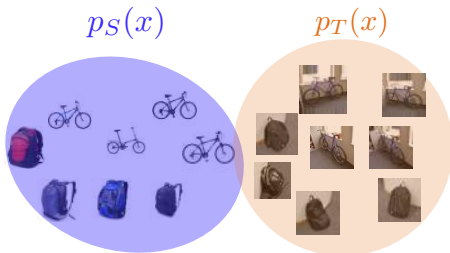
## Covariate Shift

Labelling functions are conserved:

$$p_S(y|x) = p_T(y|x) \quad (1)$$

$$\epsilon_T(h) = \mathbb{E}_T[\ell(h(x), y)] = \mathbb{E}_S \left[ \frac{p_T(x)}{p_S(x)} \ell(h(x), y) \right] \quad (2)$$

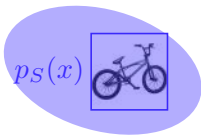
Needs overlapping supports!



# Domain Invariant Representations [Ben-David et al., Ganin et al.]



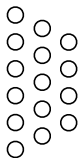
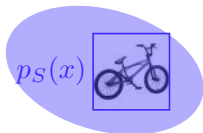
Non-overlapping  
distributions



# Domain Invariant Representations [Ben-David et al., Ganin et al.]

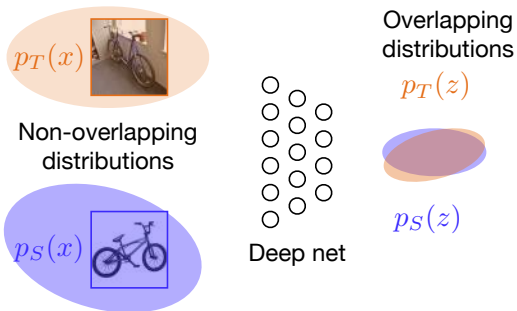


Non-overlapping  
distributions



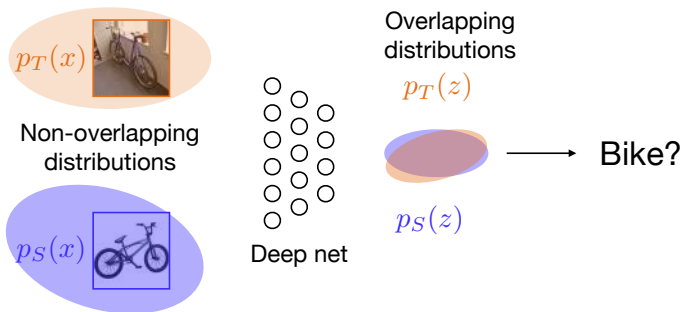
Deep net

# Domain Invariant Representations [Ben-David et al., Ganin et al.]

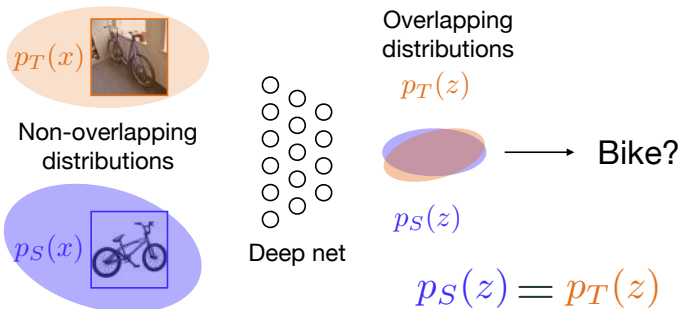




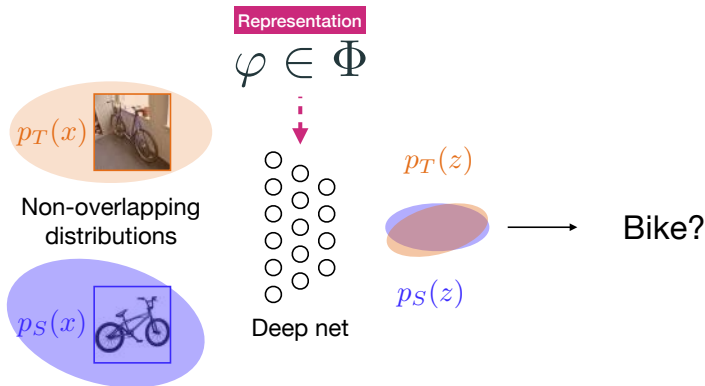
# Domain Invariant Representations [Ben-David et al., Ganin et al.]



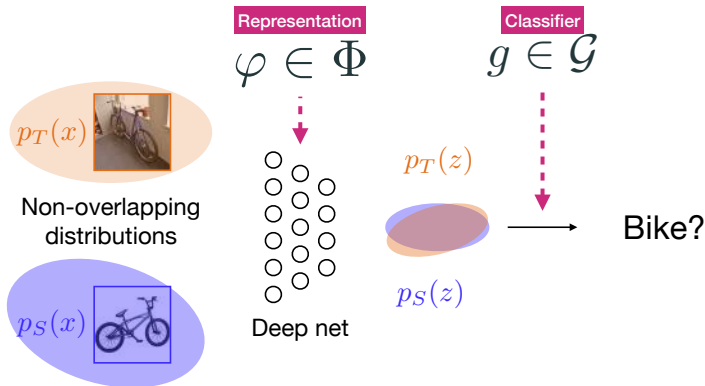
# Domain Invariant Representations [Ben-David et al., Ganin et al.]



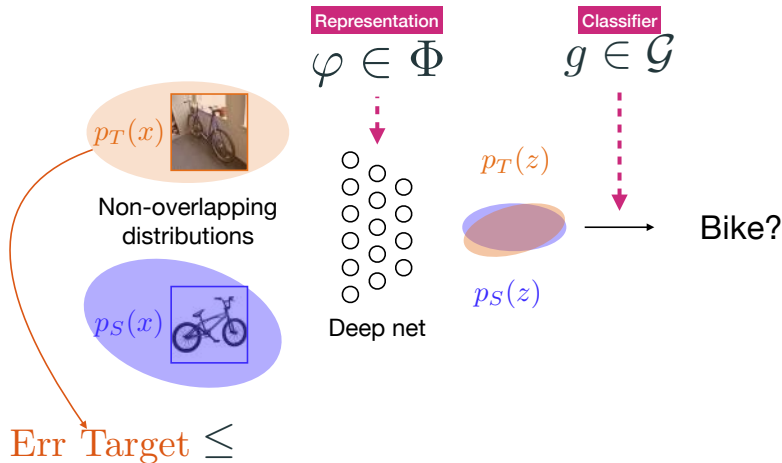
# Domain Invariant Representations [Ben-David et al., Ganin et al.]



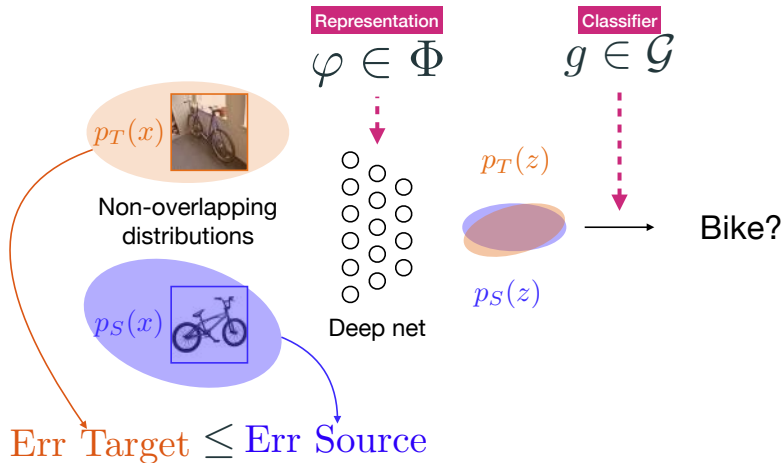
# Domain Invariant Representations [Ben-David et al., Ganin et al.]



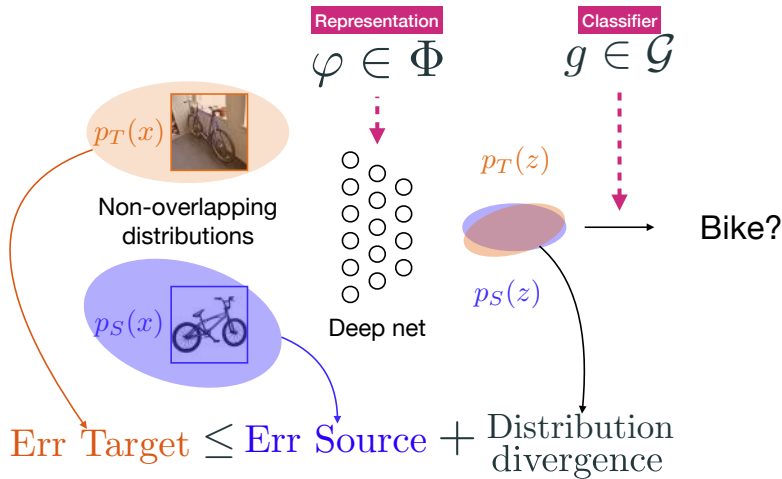
# Domain Invariant Representations [Ben-David et al., Ganin et al.]



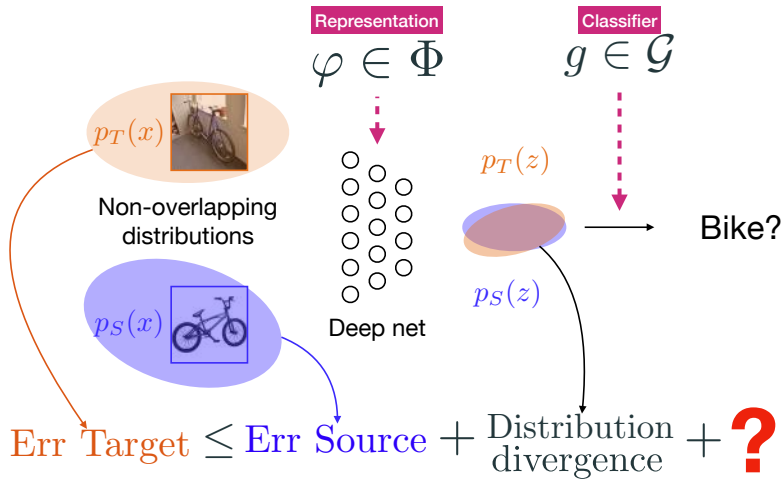
# Domain Invariant Representations [Ben-David et al., Ganin et al.]



# Domain Invariant Representations [Ben-David et al., Ganin et al.]

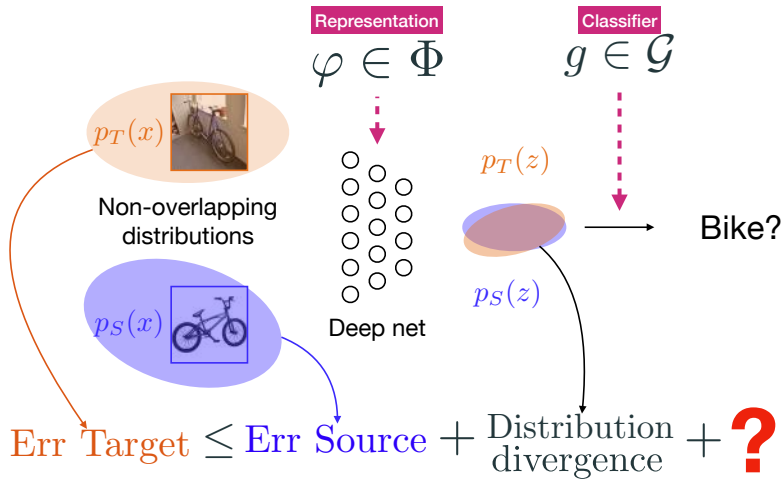


# Domain Invariant Representations [Ben-David et al., Ganin et al.]

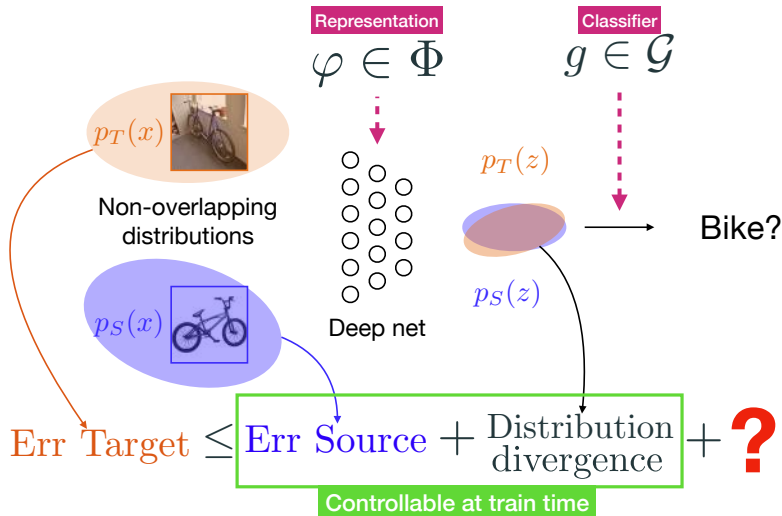




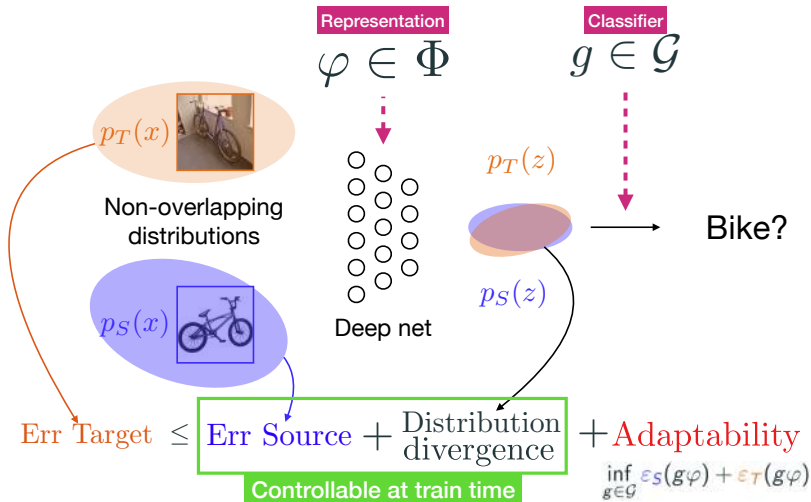
# Domain Invariant Representations [Ben-David et al., Ganin et al.]



# Domain Invariant Representations [Ben-David et al., Ganin et al.]



# Domain Invariant Representations [Ben-David et al., Ganin et al.]



# Domain Invariant Representations: Limits

$$\varepsilon_T(g\varphi) \leq \underbrace{\varepsilon_S(g\varphi) + d_G(\varphi)}_{\text{Controllable}} + \underbrace{\lambda_G(\varphi)}_{\text{Not controllable}} \quad (3)$$

## An unexpected trade-off

Let  $\psi$  be a representation which is a richer feature extractor than  $\varphi$ :

$$\mathcal{G} \circ \varphi \subset \mathcal{G} \circ \psi$$

# Domain Invariant Representations: Limits

$$\varepsilon_T(g\varphi) \leq \underbrace{\varepsilon_S(g\varphi) + d_G(\varphi)}_{\text{Controllable}} + \underbrace{\lambda_G(\varphi)}_{\text{Not controllable}} \quad (3)$$

## An unexpected trade-off

Let  $\psi$  be a representation which is a richer feature extractor than  $\varphi$ :

$$\mathcal{G} \circ \varphi \subset \mathcal{G} \circ \psi$$

Then,

$$d_G(\varphi) \leq d_G(\psi) \text{ while } \lambda_G(\psi) \leq \lambda_G(\varphi) \quad (4)$$

# Domain Invariant Representations: Limits

$$\varepsilon_T(g\varphi) \leq \underbrace{\varepsilon_S(g\varphi) + d_G(\varphi)}_{\text{Controllable}} + \underbrace{\lambda_G(\varphi)}_{\text{Not controllable}} \quad (3)$$

## An unexpected trade-off

Let  $\psi$  be a representation which is a richer feature extractor than  $\varphi$ :

$$\mathcal{G} \circ \varphi \subset \mathcal{G} \circ \psi$$

Then,

$$d_G(\varphi) \leq d_G(\psi) \text{ while } \lambda_G(\psi) \leq \lambda_G(\varphi) \quad (4)$$

▷ The benefit of representation invariance must be higher than the loss of adaptability, which is impossible to guarantee in practice.

# Domain Invariant Representations: Limits

$$\varepsilon_T(g\varphi) \leq \underbrace{\varepsilon_S(g\varphi) + d_G(\varphi)}_{\text{Controllable}} + \underbrace{\lambda_G(\varphi)}_{\text{Not controllable}} \quad (3)$$

## An unexpected trade-off

Let  $\psi$  be a representation which is a richer feature extractor than  $\varphi$ :

$$\mathcal{G} \circ \varphi \subset \mathcal{G} \circ \psi$$

Then,

$$d_G(\varphi) \leq d_G(\psi) \text{ while } \lambda_G(\psi) \leq \lambda_G(\varphi) \quad (4)$$

▷ The benefit of representation invariance must be higher than the loss of adaptability, which is impossible to guarantee in practice.

**Invariance is conflicting with label shift [Zhao et al., ICML2019]**

$$\lambda_G(\varphi) \geq \frac{1}{2} (\text{JS}(Y) - \text{JS}(Z))^2 \quad (5)$$

If  $\text{JS}(Z) \rightarrow 0$ ,  $\lambda_G(\varphi)$  can not be small if  $\text{JS}(Y)$  is high...

# Outline

---



1. A new trade-off between **invariance** and **transferability**

# Outline

1. A new trade-off between **invariance** and **transferability**
  - Introduce a new error term named *transferability error*
  - Reconcile Invariant Representations and Weights

1. A new trade-off between **invariance** and **transferability**
  - Introduce a new error term named *transferability error*
  - Reconcile Invariant Representations and Weights
2. Role of **Inductive Bias**:
  - Weights  $\triangleright$  *induce new property of invariance of representations and the labelling function*
  - Classifier  $\triangleright$  Feedback for better representations invariance.

1. A new trade-off between **invariance** and **transferability**
  - Introduce a new error term named *transferability error*
  - Reconcile Invariant Representations and Weights
2. Role of **Inductive Bias**:
  - Weights ▷ *induce new property of invariance of representations and the labelling function*
  - Classifier ▷ Feedback for better representations invariance.
3. A new algorithm for Robust Unsupervised Domain Adaptation ▷ RUDA
  - Robust to strong label shift
  - Evaluation on two benchmarks

## **A new trade-off between invariance and transferability**

---

# Three ingredients

1. **INV** ▷ captures the difference between source and target distribution of representations.

# Three ingredients

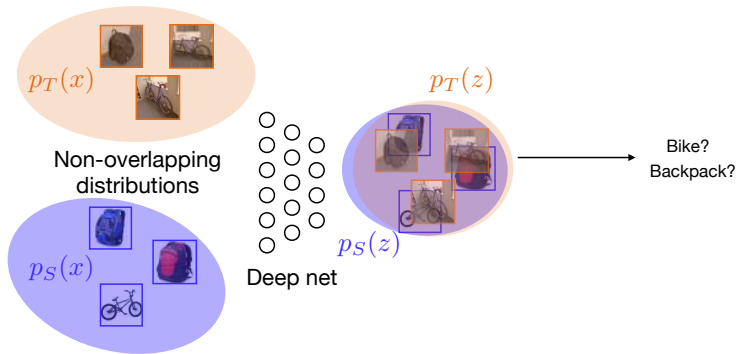
1. **INV** ▷ captures the difference between source and target distribution of representations.
2. **TSE** ▷ catches if the coupling between representations and labels shifts across domains.

# Three ingredients

1. **INV** ▷ captures the difference between source and target distribution of representations.
2. **TSF** ▷ catches if the coupling between representations and labels shifts across domains.
3. Reconcile Weights and Invariant Representations.



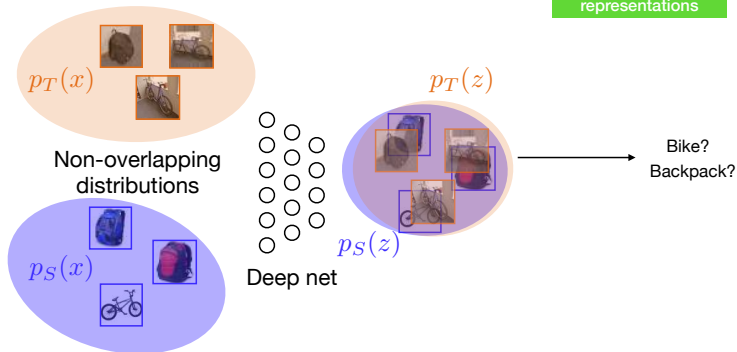
# A new trade-off between invariance and transferability



# A new trade-off between invariance and transferability

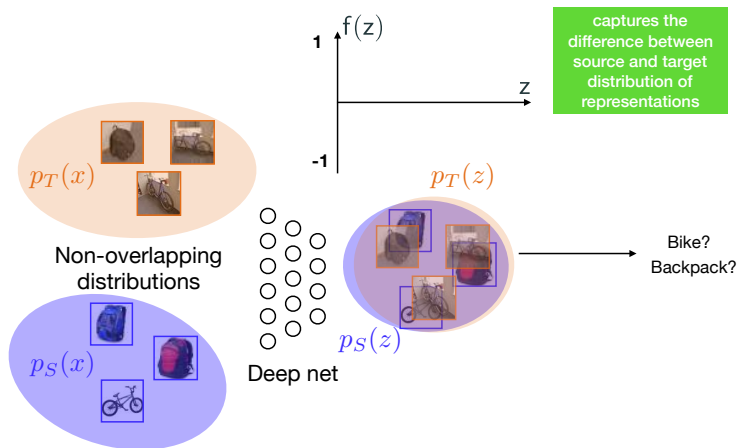
$$\text{INV}(\varphi) := \sup_{f \in \mathcal{F}} \mathbb{E}_S[f(Z)] - \mathbb{E}_T[f(Z)]$$

captures the  
difference between  
source and target  
distribution of  
representations



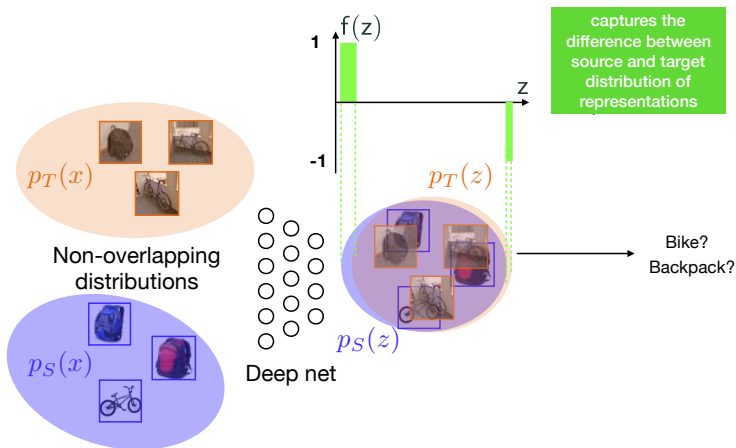
# A new trade-off between invariance and transferability

$$\text{INV}(\varphi) := \sup_{f \in \mathcal{F}} \mathbb{E}_S[f(Z)] - \mathbb{E}_T[f(Z)]$$



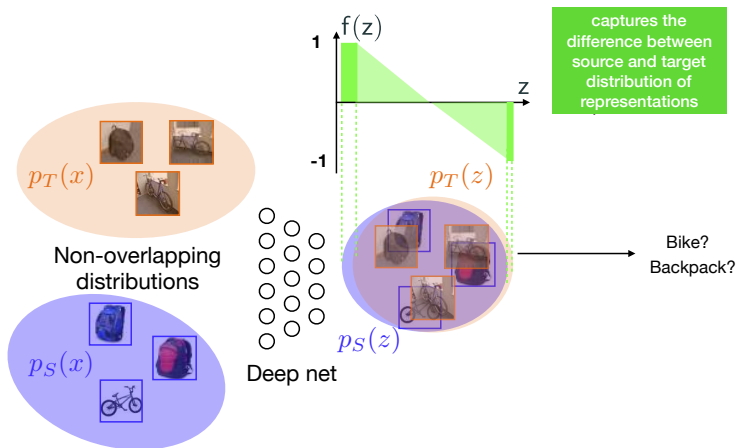
# A new trade-off between invariance and transferability

$$\text{INV}(\varphi) := \sup_{f \in \mathcal{F}} \mathbb{E}_S[f(Z)] - \mathbb{E}_T[f(Z)]$$



# A new trade-off between invariance and transferability

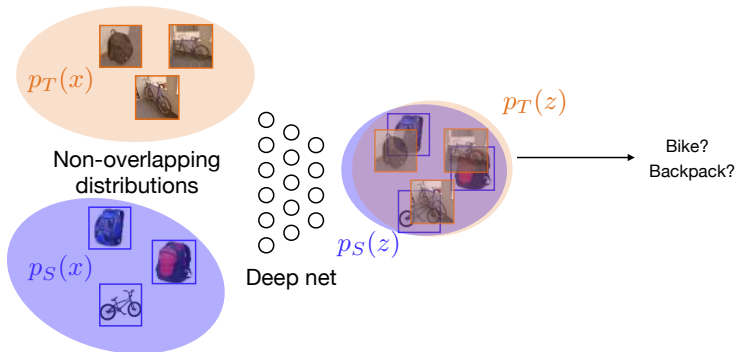
$$\text{INV}(\varphi) := \sup_{f \in \mathcal{F}} \mathbb{E}_S[f(Z)] - \mathbb{E}_T[f(Z)]$$



# A new trade-off between invariance and transferability

$$\text{TSF}(\varphi) := \sup_{f \in \mathcal{F}_c} \mathbb{E}_S[Y \cdot f(Z)] - \mathbb{E}_T[Y \cdot f(Z)]$$

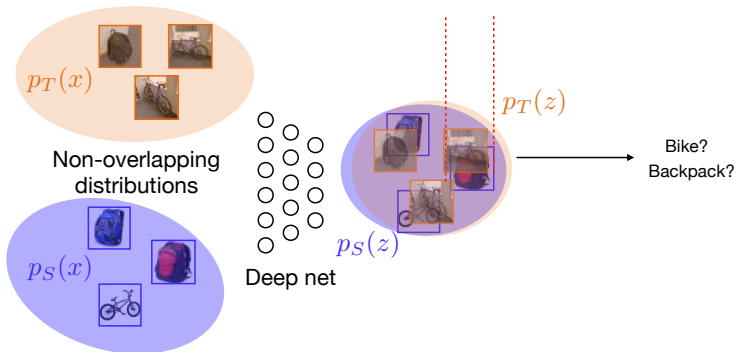
catches if the coupling between representations and labels shifts across domains



# A new trade-off between invariance and transferability

$$\text{TSF}(\varphi) := \sup_{f \in \mathcal{F}_c} \mathbb{E}_S[Y \cdot f(Z)] - \mathbb{E}_T[Y \cdot f(Z)]$$

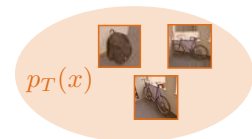
catches if the coupling between representations and labels shifts across domains



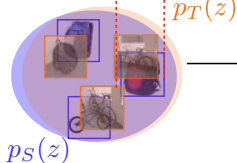
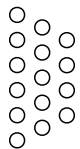
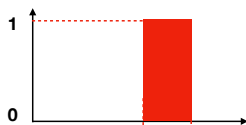
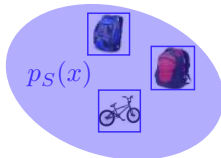
# A new trade-off between invariance and transferability

$$\text{TSF}(\varphi) := \sup_{f \in \mathcal{F}_c} \mathbb{E}_S[Y \cdot f(Z)] - \mathbb{E}_T[Y \cdot f(Z)]$$

catches if the coupling between representations and labels shifts across domains



Non-overlapping distributions



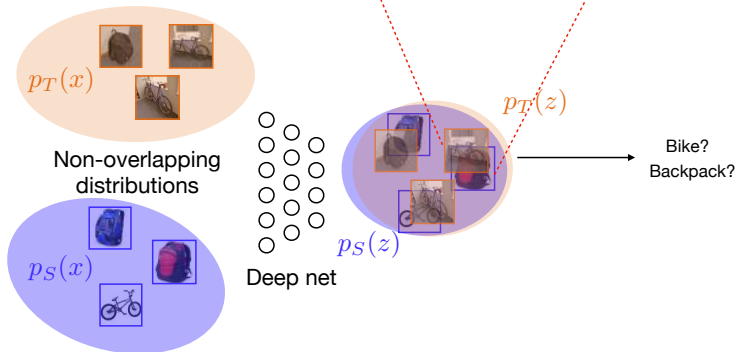
Bike?  
Backpack?



# A new trade-off between invariance and transferability

$$\text{TSF}(\varphi) := \sup_{f \in \mathcal{F}_c} \mathbb{E}_S[Y \cdot f(Z)] - \mathbb{E}_T[Y \cdot f(Z)]$$

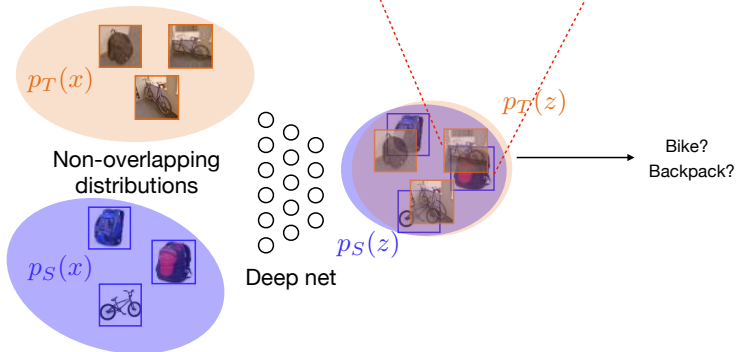
catches if the coupling between representations and labels shifts across domains



# A new trade-off between invariance and transferability

$$\text{TSF}(\varphi) := \sup_{f \in \mathcal{F}_c} \mathbb{E}_S[Y \cdot f(Z)] - \mathbb{E}_T[Y \cdot f(Z)]$$

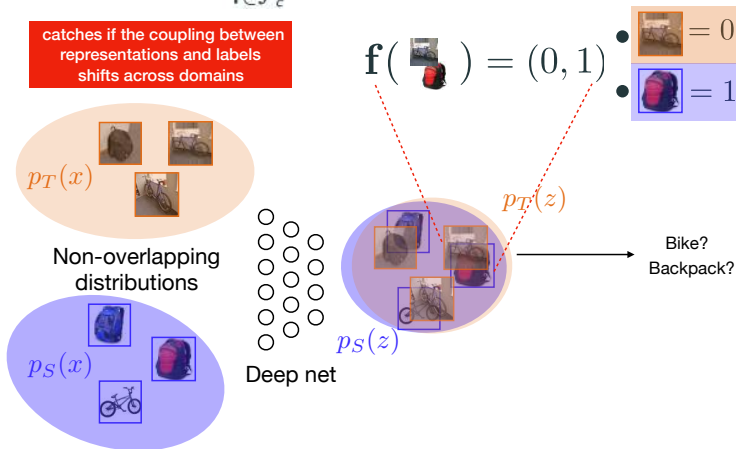
catches if the coupling between representations and labels shifts across domains



# A new trade-off between invariance and transferability

$$\text{TSF}(\varphi) := \sup_{f \in \mathcal{F}_c} \mathbb{E}_S[Y \cdot f(Z)] - \mathbb{E}_T[Y \cdot f(Z)]$$

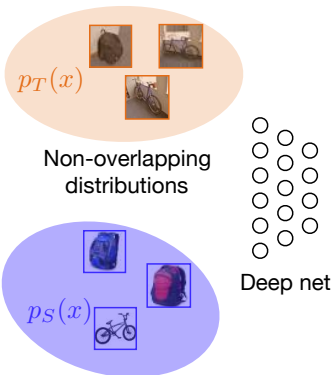
catches if the coupling between representations and labels shifts across domains



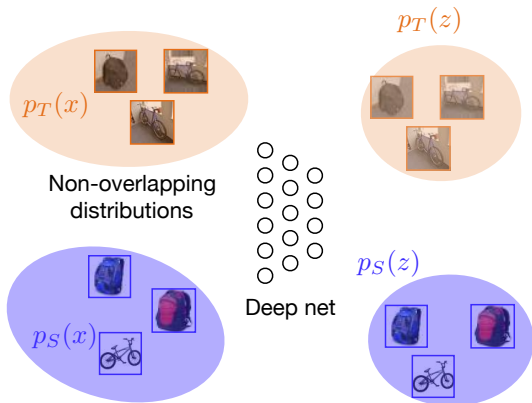
# Reconciling Weights and Invariant Representations

---

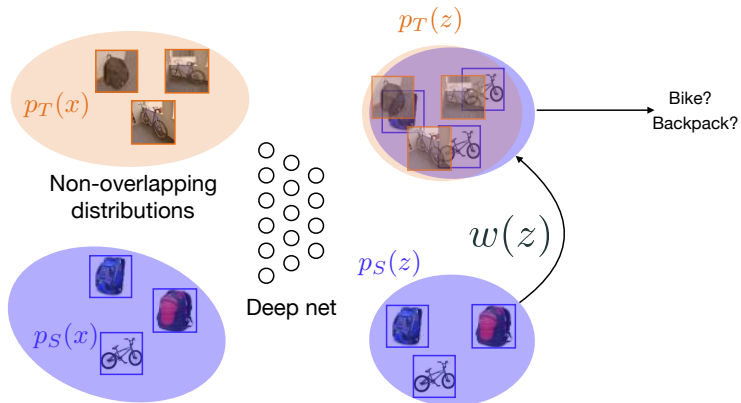
# Reconciling Weights and Invariant Representations



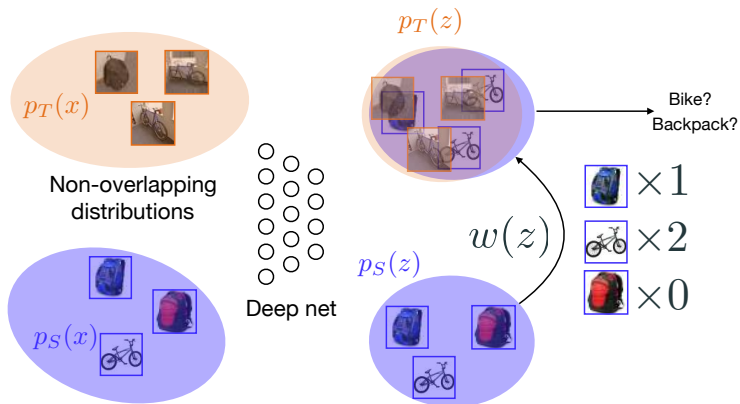
# Reconciling Weights and Invariant Representations



# Reconciling Weights and Invariant Representations

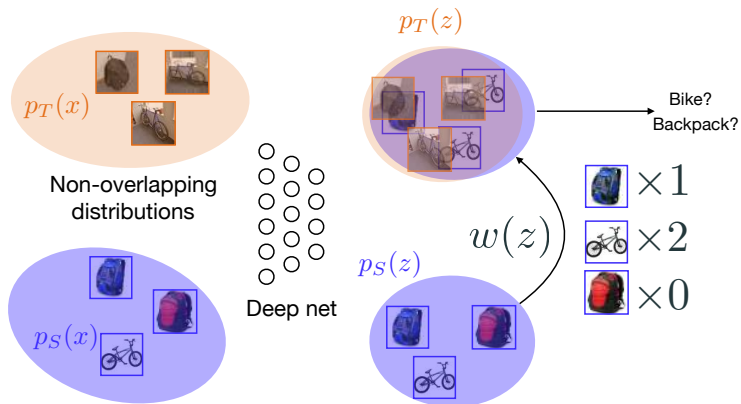


# Reconciling Weights and Invariant Representations



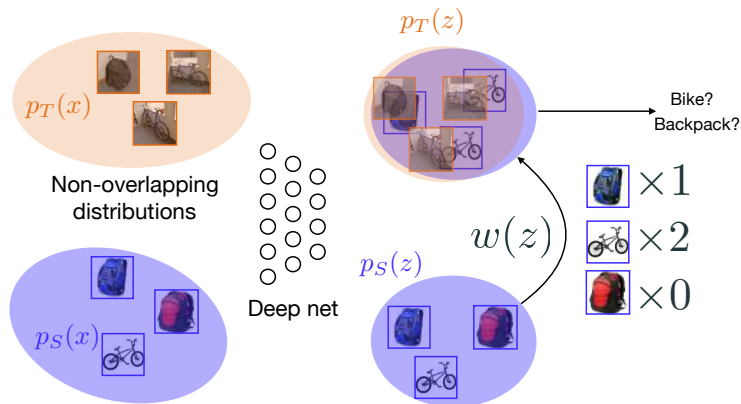


# Reconciling Weights and Invariant Representations



$$\text{INV}(w, \varphi) := \sup_{f \in \mathcal{F}} \mathbb{E}_S[w(Z)f(Z)] - \mathbb{E}_T[f(Z)]$$

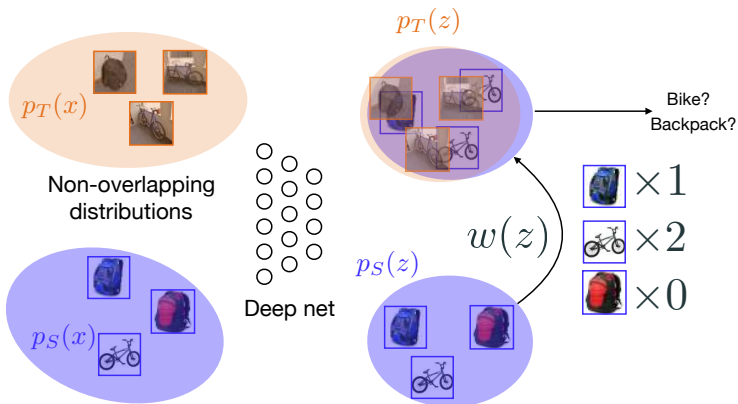
# Reconciling Weights and Invariant Representations



$$\text{INV}(w, \varphi) := \sup_{f \in \mathcal{F}} \mathbb{E}_S[w(Z)f(Z)] - \mathbb{E}_T[f(Z)]$$

$$\text{TSF}(w, \varphi) := \sup_{f \in \mathcal{F}_c} \mathbb{E}_S[w(Z)Y \cdot f(Z)] - \mathbb{E}_T[Y \cdot f(Z)]$$

# Reconciling Weights and Invariant Representations

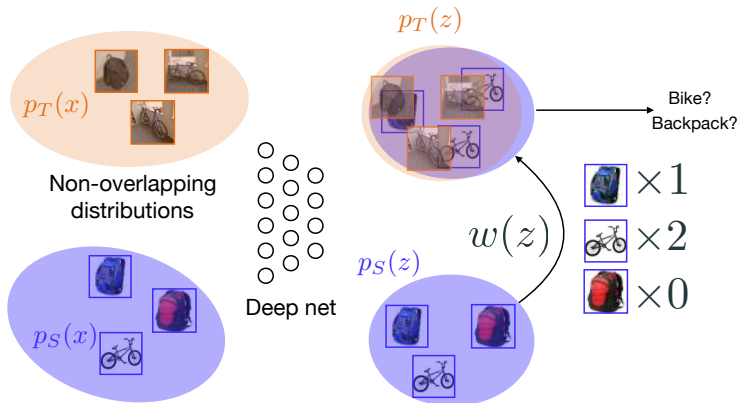


$$\text{INV}(w, \varphi) := \sup_{f \in \mathcal{F}} \mathbb{E}_S[w(Z)f(Z)] - \mathbb{E}_T[f(Z)]$$

$$\text{TFS}(w, \varphi) := \sup_{f \in \mathcal{F}_c} \mathbb{E}_S[w(Z)Y \cdot f(Z)] - \mathbb{E}_T[Y \cdot f(Z)]$$

$$\boxed{\varepsilon_T(g\varphi) \leq \varepsilon_{w \cdot S}(g\varphi)}$$

# Reconciling Weights and Invariant Representations

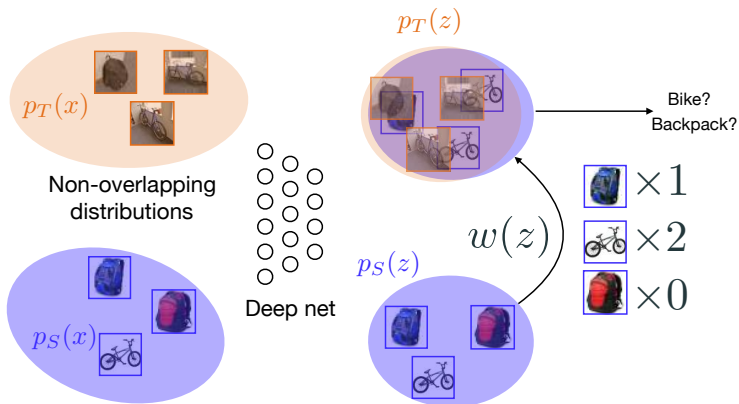


$$\text{INV}(w, \varphi) := \sup_{f \in \mathcal{F}} \mathbb{E}_S[w(Z)f(Z)] - \mathbb{E}_T[f(Z)]$$

$$\text{TSF}(w, \varphi) := \sup_{f \in \mathcal{F}_c} \mathbb{E}_S[w(Z)Y \cdot f(Z)] - \mathbb{E}_T[Y \cdot f(Z)]$$

$$\varepsilon_T(g\varphi) \leq \varepsilon_{w,S}(g\varphi) + 6 \cdot \text{INV}(w, \varphi)$$

# Reconciling Weights and Invariant Representations

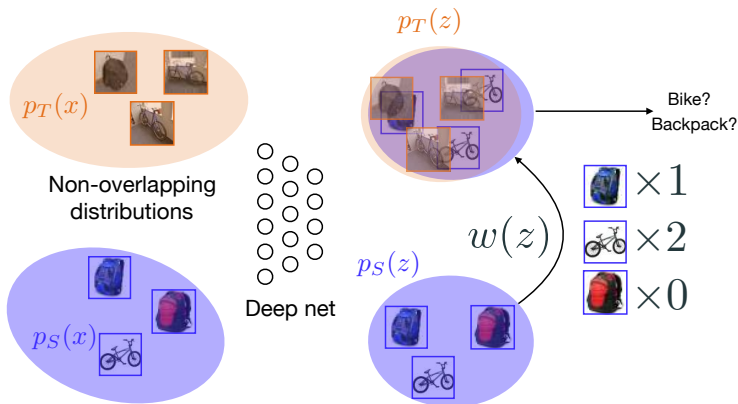


$$\text{INV}(w, \varphi) := \sup_{f \in \mathcal{F}} \mathbb{E}_S[w(Z)f(Z)] - \mathbb{E}_T[f(Z)]$$

$$\text{TSF}(w, \varphi) := \sup_{f \in \mathcal{F}_c} \mathbb{E}_S[w(Z)Y \cdot f(Z)] - \mathbb{E}_T[Y \cdot f(Z)]$$

$$\mathcal{E}_T(g\varphi) \leq \mathcal{E}_{w \cdot S}(g\varphi) + 6 \cdot \text{INV}(w, \varphi) + 2 \cdot \text{TSF}(w, \varphi)$$

# Reconciling Weights and Invariant Representations



$$\text{INV}(w, \varphi) := \sup_{f \in \mathcal{F}} \mathbb{E}_S[w(Z)f(Z)] - \mathbb{E}_T[f(Z)]$$

$$\text{TSF}(w, \varphi) := \sup_{f \in \mathcal{F}_c} \mathbb{E}_S[w(Z)Y \cdot f(Z)] - \mathbb{E}_T[Y \cdot f(Z)]$$

$$\epsilon_T(g\varphi) \leq \epsilon_{w \cdot S}(g\varphi) + 6 \cdot \text{INV}(w, \varphi) + 2 \cdot \text{TSF}(w, \varphi) + \epsilon_T(\mathbf{f}_T\varphi)$$

# Reconciling Weights and Invariant Representations

$$\varepsilon_T(g\varphi) \leq \varepsilon_{w.S}(g\varphi) + 6 \cdot \text{INV}(w, \varphi) + 2 \cdot \text{TSF}(w, \varphi) + \varepsilon_T(\mathbf{f}_T\varphi) \quad (6)$$

# Reconciling Weights and Invariant Representations

$$\varepsilon_T(g\varphi) \leq \varepsilon_{w \cdot S}(g\varphi) + 6 \cdot \text{INV}(w, \varphi) + 2 \cdot \text{TSF}(w, \varphi) + \varepsilon_T(\mathbf{f}_T\varphi) \quad (6)$$

## Weighting the source domain

$$\varepsilon_{w \cdot S}(g\varphi) := \mathbb{E}_S[w(Z)\ell(g(Z), Y)] \quad (7)$$

## A new trade-off

$$\varepsilon_T(\mathbf{f}_T\varphi) := \inf_{\mathbf{f} \in \mathcal{F}_c} \varepsilon_T(\mathbf{f}\varphi) \quad (8)$$



# Reconciling Weights and Invariant Representations

$$\varepsilon_T(g\varphi) \leq \varepsilon_{w \cdot S}(g\varphi) + 6 \cdot \text{INV}(w, \varphi) + 2 \cdot \text{TSF}(w, \varphi) + \varepsilon_T(\mathbf{f}_T\varphi) \quad (6)$$

## Weighting the source domain

$$\varepsilon_{w \cdot S}(g\varphi) := \mathbb{E}_S[w(Z)\ell(g(Z), Y)] \quad (7)$$

## A new trade-off

$$\varepsilon_T(\mathbf{f}_T\varphi) := \inf_{\mathbf{f} \in \mathcal{F}_c} \varepsilon_T(\mathbf{f}\varphi) \quad (8)$$

## Tightness

$\text{INV}(w, \varphi) = 0$  and  $\text{TSF}(w, \varphi) = 0$ , then,

$$\implies p_T(y|z) = p_S(y|z) \quad \text{and} \quad w(z) = \frac{p_T(z)}{p_S(z)}$$

▷ Much smaller than the **adaptability**.

## Remaining challenges

$$\varepsilon_T(g\varphi) \leq \varepsilon_{w.S}(g\varphi) + 6 \cdot \text{INV}(w, \varphi) + 2 \cdot \text{TSF}(w, \varphi) + \varepsilon_T(\mathbf{f}_T\varphi) \quad (9)$$

# Remaining challenges

$$\varepsilon_T(g\varphi) \leq \varepsilon_{w \cdot S}(g\varphi) + 6 \cdot \text{INV}(w, \varphi) + 2 \cdot \text{TSF}(w, \varphi) + \varepsilon_T(\mathbf{f}_T\varphi) \quad (9)$$

1. Classifier  $\triangleright$  *address the lack of labelled data in the target domain.*

$$\text{TSF}(w, \varphi) := \sup_{\mathbf{f} \in \mathcal{F}_c} \mathbb{E}_S[w(Z)Y \cdot \mathbf{f}(Z)] - \mathbb{E}_T[\underbrace{Y}_{??} \cdot \mathbf{f}(Z)]$$

2. Weights  $\triangleright$  induce invariance property on representations (see the paper for more details).

## Remaining challenges

$$\varepsilon_T(g\varphi) \leq \varepsilon_{w \cdot S}(g\varphi) + 6 \cdot \text{INV}(w, \varphi) + 2 \cdot \text{TSF}(w, \varphi) + \varepsilon_T(\mathbf{f}_T\varphi) \quad (9)$$

1. *Classifier* ▷ address the lack of labelled data in the target domain.

$$\text{TSF}(w, \varphi) := \sup_{\mathbf{f} \in \mathcal{F}_c} \mathbb{E}_S[w(Z)Y \cdot \mathbf{f}(Z)] - \mathbb{E}_T[\underbrace{Y}_{??} \cdot \mathbf{f}(Z)]$$

2. Weights ▷ induce invariance property on representations (see the paper for more details).

# Inductive design of the classifier

---

$$\widehat{\text{TSF}}(\varphi, \tilde{g}) := \sup_{f \in F_C} \mathbb{E}_S[Y \cdot f(Z)] - \mathbb{E}_T[\tilde{g}(Z) \cdot f(Z)] \quad (10)$$

# Inductive design of the classifier

$$\widehat{\text{TSF}}(\varphi, \tilde{g}) := \sup_{f \in F_C} \mathbb{E}_S[Y \cdot f(Z)] - \mathbb{E}_T[\tilde{g}(Z) \cdot f(Z)] \quad (10)$$

## Insight

The **best source classifier** is not the **best target classifier**, and, it is possible to improve the **best source classifier**, e.g., **specific architecture** or a **well-suited regularization**.

# Inductive design of the classifier

$$\widehat{\text{TSF}}(\varphi, \tilde{g}) := \sup_{f \in F_c} \mathbb{E}_S[Y \cdot f(Z)] - \mathbb{E}_T[\tilde{g}(Z) \cdot f(Z)] \quad (10)$$

## Insight

The **best source classifier** is not the **best target classifier**, and, it is possible to improve the **best source classifier**, e.g., **specific architecture** or a **well-suited regularization**.

## Inductive design

We say that there is an inductive design of a classifier at level  $0 < \beta \leq 1$  if for any representations  $\varphi$ , noting  $g_S = \arg \min_{g \in \mathcal{G}} \varepsilon_S(g\varphi)$ , we can determine  $\tilde{g}$  such that:

$$\varepsilon_T(\tilde{g}\varphi) \leq \beta \varepsilon_T(g_S\varphi) \quad (11)$$

- $\beta$ -strong when  $\beta < 1$ ,
- weak when  $\beta = 1$ .



## Generalization guarantees with inductive bias

A revisited version of the bound:

$$\varepsilon_T(g_S\varphi) \leq \varepsilon_S(g_S\varphi) + 6 \cdot \text{INV}(\varphi) + 2 \cdot \widehat{\text{TSF}}(\varphi, \tilde{g}) + \varepsilon_T(\tilde{g}\varphi) + \varepsilon_T(\mathbf{f}_T\varphi) \quad (12)$$

# Generalization guarantees with inductive bias

A revisited version of the bound:

$$\varepsilon_T(g_S\varphi) \leq \varepsilon_S(g_S\varphi) + 6 \cdot \text{INV}(\varphi) + 2 \cdot \widehat{\text{TSSF}}(\varphi, \tilde{g}) + \varepsilon_T(\tilde{g}\varphi) + \varepsilon_T(\mathbf{f}_T\varphi) \quad (12)$$

- $\widehat{\text{TSSF}}(\varphi, \tilde{g}) \triangleright$  transferability of representations with respect to the inductive design:

$$\widehat{\text{TSSF}}(\varphi, \tilde{g}) := \sup_{\mathbf{f} \in F_C} \mathbb{E}_S[Y \cdot \mathbf{f}(Z)] - \mathbb{E}_T[\tilde{g}(Z) \cdot \mathbf{f}(Z)] \quad (13)$$

# Generalization guarantees with inductive bias

A revisited version of the bound:

$$\varepsilon_T(g_S\varphi) \leq \varepsilon_S(g_S\varphi) + 6 \cdot \text{INV}(\varphi) + 2 \cdot \widehat{\text{TSSF}}(\varphi, \tilde{g}) + \varepsilon_T(\tilde{g}\varphi) + \varepsilon_T(\mathbf{f}_T\varphi) \quad (12)$$

- $\widehat{\text{TSSF}}(\varphi, \tilde{g}) \triangleright$  transferability of representations with respect to the inductive design:

$$\widehat{\text{TSSF}}(\varphi, \tilde{g}) := \sup_{\mathbf{f} \in F_C} \mathbb{E}_S[Y \cdot \mathbf{f}(Z)] - \mathbb{E}_T[\tilde{g}(Z) \cdot \mathbf{f}(Z)] \quad (13)$$

- $\varepsilon_T(\tilde{g}\varphi) \triangleright$  How good is the inductive design.

# Generalization guarantees with inductive bias

A revisited version of the bound:

$$\varepsilon_T(g_S\varphi) \leq \varepsilon_S(g_S\varphi) + 6 \cdot \text{INV}(\varphi) + 2 \cdot \widehat{\text{TSF}}(\varphi, \tilde{g}) + \varepsilon_T(\tilde{g}\varphi) + \varepsilon_T(\mathbf{f}_T\varphi) \quad (12)$$

- $\widehat{\text{TSF}}(\varphi, \tilde{g}) \triangleright$  transferability of representations with respect to the inductive design:

$$\widehat{\text{TSF}}(\varphi, \tilde{g}) := \sup_{\mathbf{f} \in F_C} \mathbb{E}_S[Y \cdot \mathbf{f}(Z)] - \mathbb{E}_T[\tilde{g}(Z) \cdot \mathbf{f}(Z)] \quad (13)$$

- $\varepsilon_T(\tilde{g}\varphi) \triangleright$  How good is the inductive design.

Here comes the guarantees

$$\boxed{\varepsilon_T(g_S\varphi)} \leq \varepsilon_S(g_S\varphi) + 6 \cdot \text{INV}(\varphi) + 2 \cdot \widehat{\text{TSF}}(\varphi, \tilde{g}) + \beta \boxed{\varepsilon_T(g_S\varphi)} + \varepsilon_T(\mathbf{f}_T\varphi)$$

# Generalization guarantees with inductive bias

A revisited version of the bound:

$$\varepsilon_T(g_S\varphi) \leq \varepsilon_S(g_S\varphi) + 6 \cdot \text{INV}(\varphi) + 2 \cdot \widehat{\text{TSF}}(\varphi, \tilde{g}) + \varepsilon_T(\tilde{g}\varphi) + \varepsilon_T(\mathbf{f}_T\varphi) \quad (12)$$

- $\widehat{\text{TSF}}(\varphi, \tilde{g}) \triangleright$  transferability of representations with respect to the inductive design:

$$\widehat{\text{TSF}}(\varphi, \tilde{g}) := \sup_{\mathbf{f} \in F_C} \mathbb{E}_S[Y \cdot \mathbf{f}(Z)] - \mathbb{E}_T[\tilde{g}(Z) \cdot \mathbf{f}(Z)] \quad (13)$$

- $\varepsilon_T(\tilde{g}\varphi) \triangleright$  How good is the inductive design.

Here comes the guarantees

$$\boxed{\varepsilon_T(g_S\varphi)} \leq \frac{1}{1-\beta} \left( \varepsilon_S(g_S\varphi) + 6 \cdot \text{INV}(\varphi) + 2 \cdot \widehat{\text{TSF}}(\varphi, \tilde{g}) + \varepsilon_T(\mathbf{f}_T\varphi) \right)$$

# Generalization guarantees with inductive bias

A revisited version of the bound:

$$\varepsilon_T(g_S\varphi) \leq \varepsilon_S(g_S\varphi) + 6 \cdot \text{INV}(\varphi) + 2 \cdot \widehat{\text{TSF}}(\varphi, \tilde{g}) + \varepsilon_T(\tilde{g}\varphi) + \varepsilon_T(\mathbf{f}_T\varphi) \quad (12)$$

- $\widehat{\text{TSF}}(\varphi, \tilde{g}) \triangleright$  transferability of representations with respect to the inductive design:

$$\widehat{\text{TSF}}(\varphi, \tilde{g}) := \sup_{\mathbf{f} \in F_C} \mathbb{E}_S[Y \cdot \mathbf{f}(Z)] - \mathbb{E}_T[\tilde{g}(Z) \cdot \mathbf{f}(Z)] \quad (13)$$

- $\varepsilon_T(\tilde{g}\varphi) \triangleright$  How good is the inductive design.

Here comes the guarantees

$$\varepsilon_T(\tilde{g}\varphi) \leq \frac{\beta}{1 - \beta} \left( \varepsilon_S(g_S\varphi) + 6 \cdot \text{INV}(\varphi) + 2 \cdot \widehat{\text{TSF}}(\varphi, \tilde{g}) + \varepsilon_T(\mathbf{f}_T\varphi) \right)$$

## Generalization guarantees with inductive bias

$$\epsilon_T(\tilde{g}\varphi) \leq \frac{\beta}{1-\beta} \left( \epsilon_S(g_S\varphi) + 6 \cdot \text{INV}(\varphi) + 2 \cdot \widehat{\text{TSF}}(\varphi, \tilde{g}) + \epsilon_T(\mathbf{f}_T\varphi) \right)$$

## Generalization guarantees with inductive bias

$$\epsilon_T(\tilde{g}\varphi) \leq \frac{\beta}{1-\beta} \left( \epsilon_S(g_S\varphi) + 6 \cdot \text{INV}(\varphi) + 2 \cdot \widehat{\text{TSF}}(\varphi, \tilde{g}) + \epsilon_T(\mathbf{f}_T\varphi) \right)$$



## Generalization guarantees with inductive bias

$$\epsilon_T(\tilde{g}\varphi) \leq \frac{\beta}{1-\beta} \left( \epsilon_S(g_S\varphi) + 6 \cdot \text{INV}(\varphi) + 2 \cdot \widehat{\text{TSF}}(\varphi, \tilde{g}) + \epsilon_T(\mathbf{f}_T\varphi) \right)$$

- Target labels are only involved in  $\epsilon_T(\mathbf{f}_T\varphi)$  which reflects the level of noise when fitting labels from representations  $\triangleright$  *transferability is now free of target labels.*

## Generalization guarantees with inductive bias

$$\varepsilon_T(\tilde{g}\varphi) \leq \frac{\beta}{1-\beta} \left( \varepsilon_S(g_S\varphi) + 6 \cdot \text{INV}(\varphi) + 2 \cdot \widehat{\text{TSF}}(\varphi, \tilde{g}) + \varepsilon_T(\mathbf{f}_T\varphi) \right)$$

- Target labels are only involved in  $\varepsilon_T(\mathbf{f}_T\varphi)$  which reflects the level of noise when fitting labels from representations  $\triangleright$  *transferability is now free of target labels*.
- the weaker the inductive bias ( $\beta \rightarrow 1$ ), the higher the bound and vice versa.

## Generalization guarantees with inductive bias

$$\varepsilon_T(\tilde{g}\varphi) \leq \frac{\beta}{1-\beta} \left( \varepsilon_S(g_S\varphi) + 6 \cdot \text{INV}(\varphi) + 2 \cdot \widehat{\text{TSF}}(\varphi, \tilde{g}) + \varepsilon_T(\mathbf{f}_T\varphi) \right)$$

- Target labels are only involved in  $\varepsilon_T(\mathbf{f}_T\varphi)$  which reflects the level of noise when fitting labels from representations  $\triangleright$  *transferability is now free of target labels*.
- the weaker the inductive bias ( $\beta \rightarrow 1$ ), the higher the bound and vice versa.
- **Takeaways**  $\triangleright$  if a regularization is available, it will interact with the transferability error.

# **Robust Unsupervised Domain Adaptation (RUDA)**

---

# Assumptions

- Weak inductive design of the classifier:
  - **Classifier:**  $\tilde{g} \leftarrow g_S (\beta = 1)$
  - *No theoretical guarantees from inductive classifier.*

# Assumptions

- Weak inductive design of the classifier:
  - **Classifier:**  $\tilde{g} \leftarrow g_S (\beta = 1)$
  - *No theoretical guarantees from inductive classifier.*
- Weights controls the invariance error:  $w(z) = \frac{p_T(z)}{p_S(z)}$

# Assumptions

- Weak inductive design of the classifier:
  - **Classifier:**  $\tilde{g} \leftarrow g_S (\beta = 1)$
  - *No theoretical guarantees from inductive classifier.*
- Weights controls the invariance error:  $w(z) = \frac{p_T(z)}{p_S(z)}$
- Bring strong robustness to the adaptation procedure  $\triangleright$  *stress-test on dataset with label strong label shift.*

# Assumptions

- Weak inductive design of the classifier:
  - **Classifier:**  $\tilde{g} \leftarrow g_S$  ( $\beta = 1$ )
  - *No theoretical guarantees from inductive classifier.*
- Weights controls the invariance error:  $w(z) = \frac{p_T(z)}{p_S(z)}$
- Bring strong robustness to the adaptation procedure  $\triangleright$  *stress-test on dataset with label strong label shift.*

$$\left\{ \begin{array}{l} \varphi^* = \arg \min_{\varphi \in \Phi} \epsilon_{w(\varphi) \cdot S}(g_{w \cdot S} \varphi) + \lambda \cdot \widehat{\text{TSF}}(w, \varphi, g_{w \cdot S}) \\ \text{such that } w(\varphi) = \arg \min_w \text{INV}(w, \varphi) \end{array} \right.$$

$\triangleright$  More details about the procedure in the paper.



# Experiments

---

# Experiments

Table 1: Accuracy (%) on the Office-31 dataset.

Method		A→W	W→A	A→D	D→A	D→W	W→D	Avg
Standard	ResNet-50	68.4 ± 0.2	60.7 ± 0.3	68.9 ± 0.2	62.5 ± 0.3	96.7 ± 0.1	99.3 ± 0.1	76.1
	DANN	82.0 ± 0.4	67.4 ± 0.5	79.7 ± 0.4	68.2 ± 0.4	96.9 ± 0.2	99.1 ± 0.1	82.2
	CDAN	93.1 ± 0.2	68.0 ± 0.4	89.8 ± 0.3	70.1 ± 0.4	98.2 ± 0.2	100. ± 0.0	86.6
	CDAN+E	94.1 ± 0.1	69.3 ± 0.4	<b>92.9 ± 0.2</b>	<b>71.0 ± 0.3</b>	<b>98.6 ± 0.1</b>	<b>100. ± 0.0</b>	<b>87.7</b>
5 × [16 ~ 31]	RUDA	<b>94.3 ± 0.3</b>	<b>70.7 ± 0.3</b>	92.1 ± 0.3	70.7 ± 0.1	98.5 ± 0.1	100. ± 0.0	87.6
	RUDA <sub>w</sub>	92.0 ± 0.3	67.9 ± 0.3	91.1 ± 0.3	70.2 ± 0.2	<b>98.6 ± 0.1</b>	100. ± 0.0	86.6
	ResNet-50	72.4 ± 0.7	59.5 ± 0.1	79.0 ± 0.1	61.6 ± 0.3	97.8 ± 0.1	99.3 ± 0.1	78.3
	DANN	67.5 ± 0.1	52.1 ± 0.8	69.7 ± 0.0	51.5 ± 0.1	89.9 ± 0.1	75.9 ± 0.2	67.8
5 × [16 ~ 31]	CDAN	82.5 ± 0.4	62.9 ± 0.6	81.4 ± 0.5	65.5 ± 0.5	98.5 ± 0.3	99.8 ± 0.0	81.6
	RUDA	85.4 ± 0.8	66.7 ± 0.5	81.3 ± 0.3	64.0 ± 0.5	98.4 ± 0.2	99.5 ± 0.1	82.1
	IWAN	72.4 ± 0.4	54.8 ± 0.8	75.0 ± 0.3	54.8 ± 1.3	97.0 ± 0.0	95.8 ± 0.6	75.0
	CDAN <sub>w</sub>	81.5 ± 0.5	64.5 ± 0.4	80.7 ± 1.0	65 ± 0.8	<b>98.7 ± 0.2</b>	99.9 ± 0.1	81.8
	RUDA <sub>w</sub>	<b>87.4 ± 0.2</b>	<b>68.3 ± 0.3</b>	<b>82.9 ± 0.4</b>	<b>68.8 ± 0.2</b>	<b>98.7 ± 0.1</b>	100. ± 0.0	<b>83.8</b>

RUDA performs similarly than SOTA approaches

Table 2: Accuracy (%) on the Digits dataset.

Method	Shift of [0 ~ 5]	U→M					Avg	M→U					Avg	Avg
		5%	10%	15%	20%	100%		5%	10%	15%	20%	100%		
DANN		41.7	51.0	59.6	69.0	94.5	63.2	34.5	51.0	59.6	63.6	90.7	59.9	63.2
CDAN		<u>50.7</u>	62.2	82.9	82.8	<b>96.9</b>	75.1	32.0	69.7	78.9	81.3	<b>93.9</b>	71.2	73.2
RUDA		44.4	58.4	80.0	<u>84.0</u>	95.5	72.5	34.9	59.0	76.1	78.8	93.3	68.4	70.5
IWAN		73.7	74.4	78.4	77.5	95.7	79.9	72.2	82.0	84.3	86.0	92.0	83.3	81.6
CDAN <sub>w</sub>		68.3	78.8	84.9	<b>88.4</b>	96.6	83.4	69.4	80.0	83.5	87.8	93.7	82.9	83.2
RUDA <sub>w</sub>		<b>78.7</b>	<b>82.8</b>	<b>86.0</b>	86.9	93.9	<b>85.7</b>	<b>78.7</b>	<b>87.9</b>	<b>88.2</b>	<b>89.3</b>	<b>92.5</b>	<b>87.3</b>	<b>86.5</b>

# Experiments

Table 1: Accuracy (%) on the Office-31 dataset.

Method		A→W	W→A	A→D	D→A	D→W	W→D	Avg
Standard	ResNet-50	68.4 ± 0.2	60.7 ± 0.3	68.9 ± 0.2	62.5 ± 0.3	96.7 ± 0.1	99.3 ± 0.1	76.1
	DANN	82.0 ± 0.4	67.4 ± 0.5	79.7 ± 0.4	68.2 ± 0.4	96.9 ± 0.2	99.1 ± 0.1	82.2
	CDAN	93.1 ± 0.2	68.0 ± 0.4	89.8 ± 0.3	70.1 ± 0.4	98.2 ± 0.2	100. ± 0.0	86.6
	CDAN+E	94.1 ± 0.1	69.3 ± 0.4	<b>92.9 ± 0.2</b>	<b>71.0 ± 0.3</b>	<b>98.6 ± 0.1</b>	<b>100. ± 0.0</b>	<b>87.7</b>
	RUDA	<b>94.3 ± 0.3</b>	<b>70.7 ± 0.3</b>	92.1 ± 0.3	70.7 ± 0.1	98.5 ± 0.1	100. ± 0.0	87.6
	RUDA <sub>w</sub>	92.0 ± 0.3	67.9 ± 0.3	91.1 ± 0.3	70.2 ± 0.2	98.6 ± 0.1	100. ± 0.0	86.6
6 ~ 31	ResNet-50	72.4 ± 0.7	59.5 ± 0.1	79.0 ± 0.1	61.6 ± 0.3	97.8 ± 0.1	99.3 ± 0.1	78.3
	DANN	67.5 ± 0.1	52.1 ± 0.8	69.7 ± 0.0	51.5 ± 0.1	89.9 ± 0.1	75.9 ± 0.2	67.8
	CDAN	82.5 ± 0.4	62.9 ± 0.6	81.4 ± 0.5	65.5 ± 0.5	98.5 ± 0.3	99.8 ± 0.0	81.6
	RUDA	85.4 ± 0.8	66.7 ± 0.5	81.3 ± 0.3	64.0 ± 0.5	98.4 ± 0.2	99.5 ± 0.1	82.1
	IWAN	72.4 ± 0.4	54.8 ± 0.8	75.0 ± 0.3	54.8 ± 1.3	97.0 ± 0.0	95.8 ± 0.6	75.0
	CDAN <sub>w</sub>	81.5 ± 0.5	64.5 ± 0.4	80.7 ± 1.0	65 ± 0.8	<b>98.7 ± 0.2</b>	99.9 ± 0.1	81.8
	RUDA <sub>w</sub>	<b>87.4 ± 0.2</b>	<b>68.3 ± 0.3</b>	<b>82.9 ± 0.4</b>	<b>68.8 ± 0.2</b>	<b>98.7 ± 0.1</b>	<b>100. ± 0.0</b>	<b>83.8</b>

RUDA still performs well even when stress with strong label shift

Table 2: Accuracy (%) on the Digits dataset.

Method	Shift of [0 ~ 5]	U→M					M→U					Avg		
		5%	10%	15%	20%	100%	5%	10%	15%	20%	100%			
DANN		41.7	51.0	59.6	69.0	94.5	63.2	34.5	51.0	59.6	63.6	90.7	59.9	63.2
CDAN		50.7	62.2	62.9	82.8	<b>96.9</b>	75.1	32.0	69.7	<b>75.9</b>	81.3	<b>93.9</b>	71.2	73.2
RUDA		44.4	58.4	80.0	84.0	95.5	72.5	34.9	59.0	70.1	78.8	93.3	68.4	70.5
IWAN		73.7	74.4	78.4	77.5	95.7	79.9	72.2	82.0	84.3	86.0	92.0	83.3	81.6
CDAN <sub>w</sub>		68.3	78.8	84.9	<b>88.4</b>	96.6	83.4	69.4	80.0	83.5	87.8	93.7	82.9	83.2
RUDA <sub>w</sub>		<b>78.7</b>	<b>82.8</b>	<b>86.0</b>	86.3	93.9	<b>85.7</b>	<b>78.7</b>	<b>87.9</b>	<b>88.2</b>	<b>89.3</b>	92.5	<b>87.3</b>	<b>86.5</b>

## Conclusion

---

# Conclusion

1. New bound of the target risk which unifies weights and representations in UDA.

$$\epsilon_T(g\varphi) \leq \epsilon_{w,S}(g\varphi) + 6 \cdot \text{INV}(w, \varphi) + 2 \cdot \text{TSE}(w, \varphi) + \epsilon_T(\mathbf{f}_T\varphi)$$

# Conclusion

1. New bound of the target risk which unifies weights and representations in UDA.
2. Theoretical analysis of the role of inductive bias when designing both weights and the classifier.

$$\epsilon_T(\tilde{g}\varphi) \leq \frac{\beta}{1-\beta} \left( \epsilon_S(g_S\varphi) + 6 \cdot \text{INV}(\varphi) + 2 \cdot \widehat{\text{TSF}}(\varphi, \tilde{g}) + \epsilon_T(\mathbf{f}_T\varphi) \right)$$

# Conclusion

1. New bound of the target risk which unifies weights and representations in UDA.
2. Theoretical analysis of the role of inductive bias when designing both weights and the classifier.
3. New learning procedure  $\triangleright$  *weak inductive bias can make adaptation more robust even when stressed by strong label shift between source and target domains.*

$$\left\{ \begin{array}{l} \varphi^* = \arg \min_{\varphi \in \Phi} \varepsilon_{w(\varphi) \cdot S}(\mathbf{g}_{w \cdot S} \varphi) + \lambda \cdot \widehat{\text{TSF}}(w, \varphi, \mathbf{g}_{w \cdot S}) \\ \text{such that } w(\varphi) = \arg \min_w \text{INV}(w, \varphi) \end{array} \right.$$

# Conclusion

1. New bound of the target risk which unifies weights and representations in UDA.
2. Theoretical analysis of the role of inductive bias when designing both weights and the classifier.
3. New learning procedure  $\triangleright$  *weak inductive bias can make adaptation more robust even when stressed by strong label shift between source and target domains.*

This work leaves room for in-depth study of stronger inductive bias by providing both theoretical and empirical foundations.



**Thank you!**

$$\varepsilon_T(g\varphi) \leq \underbrace{\varepsilon_S(g\varphi) + d_G(\varphi)}_{\text{Controllable}} + \underbrace{\lambda_G(\varphi)}_{\text{Not controllable}} \quad (14)$$

## An unexpected trade-off

Let  $\psi$  be a representation which is a richer feature extractor than  $\varphi$ :

$\mathcal{G} \circ \varphi \subset \mathcal{G} \circ \psi$ . Then,

$$d_G(\varphi) \leq d_G(\psi) \text{ while } \lambda_G(\psi) \leq \lambda_G(\varphi) \quad (15)$$

## Invariance is conflicting with label shift [Zhao et al., ICML2019]

$$\lambda_{\mathcal{G}}(\varphi) \geq \frac{1}{2} (\text{JS}(Y) - \text{JS}(Z))^2 \quad (16)$$

If  $\text{JS}(Z) \rightarrow 0$ ,  $\lambda_{\mathcal{G}}(\varphi)$  can not be small if  $\text{JS}(Y)$  is high...  $\triangleright$  *We should weight the source distributions! But how...*

- how to design weights?
- how weights interact with invariance?
- why predictions are important in UDA?

# Our strategy

▷ emerges from the sup / inf duality computed on a small hypothesis class  $\mathcal{G} \circ \varphi$ .

## Our strategy

Express both **invariance** and **transferability** of representations as supremum over a large space of critic functions

▷ See the paper for details about property of the critic functions

## 2 critic functions space

- $\mathcal{F}$  ▷ 'large' function space from  $\mathcal{Z}$  to  $[-1, 1]$
  - $\mathcal{F}_C$  ▷ 'large' function space from  $\mathcal{Z}$  to  $[-1, 1]^C$
- ▷ Typically continuous functions.

## 2 critic functions space

- $\mathcal{F}$  ▷ 'large' function space from  $\mathcal{Z}$  to  $[-1, 1]$
- $\mathcal{F}_C$  ▷ 'large' function space from  $\mathcal{Z}$  to  $[-1, 1]^C$

▷ Typically continuous functions.

## 2 errors

- captures the difference between source and target distribution of representations:

$$\text{INV}(\varphi) := \sup_{f \in \mathcal{F}} \mathbb{E}_S[f(Z)] - \mathbb{E}_T[f(Z)] \quad (17)$$

- catches if the coupling between  $Z$  and  $Y$  shifts across domains:

$$\text{TSF}(\varphi) := \sup_{\mathbf{f} \in \mathcal{F}_c} \mathbb{E}_S[Y \cdot \mathbf{f}(Z)] - \mathbb{E}_T[Y \cdot \mathbf{f}(Z)] \quad (18)$$

## Why designing weights?

- Prediction weighting [Partial Adversarial Domain Adaptation, Cao et al. 2018] ▷ *Estimated labels are used to re-weight the source domain:*

$$w(x) = \frac{p_T(g(z))}{p_S(g(z))}$$

- Entropy conditioning [Conditional Adversarial Domain Adaptation, Long et al. 2018] ▷ *Transfer only confident samples:*

$$w(x) \propto 1 + e^{-H(g(z))} \text{ where } H \text{ is the entropy.}$$

## Inductive design of weights

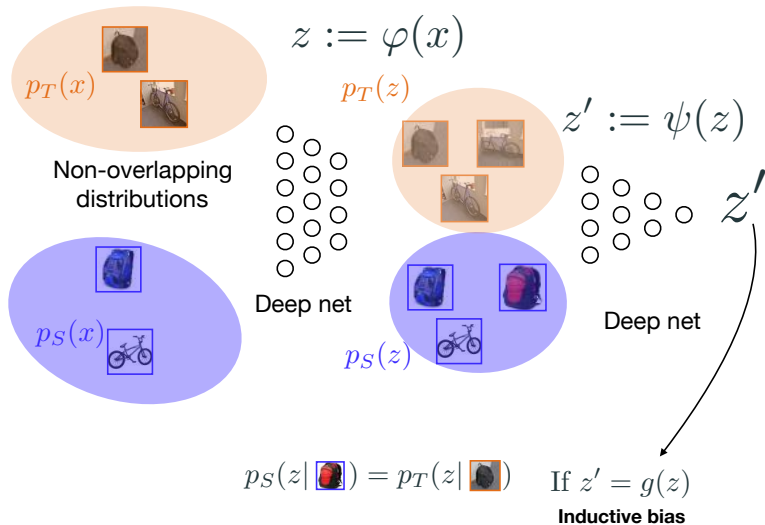
It exists a function  $\psi : \mathcal{Z} \rightarrow \mathcal{Z}'$ ,  $\psi(z) =: z'$  s.t.  $w$  is a function of  $Z'$ .

## Weights enforces new invariance

If the bound is tight, then,

$$w(z') = \frac{p_T(z')}{p_S(z')} \quad \text{and} \quad p_S(z|z') = p_T(z|z') \quad (19)$$

# Role of weights



## Detailed view of RUDA

$$\begin{aligned}\epsilon_T(\tilde{g}\varphi) &\leq \frac{\beta}{1-\beta} \{ \epsilon_{w \cdot S}(g_{w \cdot S}\varphi) \\ &\quad + 6 \cdot \text{INV}(w, \varphi) \triangleright \text{Set weight s.t. } \text{INV}(w, \varphi) = 0 \\ &\quad + 2 \cdot \widehat{\text{TSF}}(w, \varphi, \tilde{g}) \triangleright \text{Set inductive classifier s.t. } \tilde{g} = g \\ &\quad + \epsilon_T(\mathbf{f}_T\varphi) \}\end{aligned}$$

- Set weight s.t.  $\text{INV}(w, \varphi) = 0$ :

$$w(z) := \frac{p_T(z)}{p_S(z)} = \frac{1 - d(z)}{d(z)} \quad (20)$$

where  $d$  is a domain classifier (trained to map 1 in the source domain and 0 in the target domain.)

- Set inductive classifier s.t.  $\tilde{g} = g$ :  $\beta = 1$   
*"This is a weak inductive design ( $\beta = 1$ ), thus, theoretical guarantee from bound 4 is not applicable. However, there is empirical evidence that showed that predicted labels help in UDA"*



## Detailed view of RUDA

$$\left\{ \begin{array}{l} \theta_{\varphi}^* = \arg \min_{\theta_{\varphi}} \mathcal{L}_c(\theta_g, \theta_{\varphi} | \theta_d) + \lambda \cdot \widehat{\mathcal{L}}_{\text{TSE}}(\theta_{\varphi}, \theta_d | \theta_d, \theta_g) \\ \theta_g = \arg \min_{\theta_g} \mathcal{L}_c(\theta_g, \theta_{\varphi} | \theta_d) \\ \theta_d = \arg \min_{\theta_d} \mathcal{L}_{\text{INV}}(\theta_d | \theta_{\varphi}) \end{array} \right. \quad (21)$$

$$\left\{ \begin{array}{l} \varphi^* = \arg \min_{\varphi \in \Phi} \varepsilon_{w(\varphi) \cdot S}(g_{w \cdot S} \varphi) + \lambda \cdot \widehat{\text{TSE}}(w, \varphi, g_{w \cdot S}) \\ \text{such that } w(\varphi) = \arg \min_w \text{INV}(w, \varphi) \end{array} \right. \quad (\text{RUDA})$$

# Detailed view of RUDA

