# Unsupervised Domain Adaptation: From Practice to Theory



#### John Blitzer



## **Unsupervised Domain Adaptation**



#### Avante Deep Fryer; Black

Title: lid does not work well...

I love the way the Tefal deep fryer cooks, however, I am returning my second one due to a defective lid closure. The lid may close initially, but after a few uses it no longer stays closed. I won't be buying this one again.





#### Source





## **Target-Specific Features**



This book was horrible. I read half, suffering from a headache the entire time, and eventually i lit it on fire. 1 less copy in the world. Don't waste your money. I wish i had the time spent reading this book back. It wasted my life

Avante Deep Fryer; Black

amazon.com<sup>®</sup>

Title: lid does not work well...

I love the way the Tefal deep fryer cooks, however, I am returning my second one due to a defective lid closure. The lid may close initially, but after a few uses it no longer stays closed. I won't be buying this one again.

Target



#### Source



#### Learning Shared Representations





# Shared Representations: A Quick Review

Blitzer et al. (2006, 2007). <u>Shared CCA</u>.

Tasks: Part of speech tagging, sentiment.

Xue et al. (2008). Probabilistic LSA

Task: Cross-lingual document classification.

Guo et al. (2009). Latent Dirichlet Allocation

Task: Named entity recognition

Huang et al. (2009). Hidden Markov Models

Task: Part of Speech Tagging





























Adaptation Learning Theory:  $\epsilon_{ ext{tar}}$ 

$$\epsilon_{\text{target}} \leq ??$$



1. A computable (source) sample bound on target error

- 2. A formal description of empirical phenomena
  - Why do shared representations algorithms work?
- 3. Suggestions for future research



- Target Generalization Bounds using Discrepancy Distance
  [BBCKPW 2009]
  [Mansour et al. 2009]
  - h
- Coupled Subspace Learning [BFK 2010]





# Formalizing Domain Adaptation

# Source distribution

 $(x,y) \sim \Pr_S[x,y]$ 

#### Target distribution

$$(x,y) \sim \Pr_T[x,y]$$

Source labeled data

 $x \sim \Pr_S[x]$ 



Target unlabeled data

$$x \sim \Pr_T[x]$$

 $y \sim \Pr_S[y|x]$ 



# Formalizing Domain Adaptation

# Source distribution

$$(x,y) \sim \Pr_S[x,y]$$

#### Target distribution

$$(x,y) \sim \Pr_T[x,y]$$

Source labeled data

 $y \sim \Pr_S[y|x]$ 

 $x \sim \Pr_S[x]$ 



Target unlabeled data

$$x \sim \Pr_T[x]$$

$$y \sim \Pr_T[y|x]$$

Semi-supervised adaptation

Some target labels



# Formalizing Domain Adaptation

# Source distribution

 $(x,y) \sim \Pr_S[x,y]$ 

#### Target distribution

$$(x,y) \sim \Pr_T[x,y]$$

Source labeled data



Target unlabeled data

$$x \sim \Pr_T[x]$$

 $x \sim \Pr_S[x]$  $y \sim \Pr_S[y|x]$ 

Semi-supervised adaptation

Not in this talk



# A Generalization Bound

S, T: Source and target  $\mathcal{H}$ : Hypothesis class n: Sample size  $\hat{S}$ : Labeled S sample  $\hat{T}$ : Unlabeled T sample  $h^*$ : best  $h \in \mathcal{H}$ 

With probability  $1 - \delta$ , for h the ERM of  $\hat{S}$ :  $\epsilon_T(h) - \epsilon_T(h^*) \leq$ 



S, T: Source and target  $\mathcal{H}$ : Hypothesis class n: Sample size  $\hat{S}$ : Labeled S sample  $\hat{T}$ : Unlabeled T sample  $h^*$ : best  $h \in \mathcal{H}$ Bound from [MMR09]

With probability  $1 - \delta$ , for h the ERM of  $\hat{S}$ :





When good source models go bad

 $\operatorname{disc}_{\mathcal{H}}(S,T) = \max_{\substack{h,h^* \in \mathcal{H}}} |E_S[h(x) \neq h^*(x)] - E_T[h(x) \neq h^*(x)]|$ 



















#### When good source models go bad

 $\operatorname{disc}_{\mathcal{H}}(S,T) = \max_{h,h^* \in \mathcal{H}} |E_S[h(x) \neq h^*(x)] - E_T[h(x) \neq h^*(x)]|$ 





Learn pairs of hypotheses to discriminate source from target





Learn pairs of hypotheses to discriminate source from target





Learn pairs of hypotheses to discriminate source from target





Linear Hypothesis Class:  $h(x) = \operatorname{sgn} \left(\beta \cdot x\right)$ 

Induced classes from projections  $\beta \cdot \Pi x$   $\Pi = \Pi \Pi$ 





Linear Hypothesis Class:  $h(x) = \operatorname{sgn} \left(\beta \cdot x\right)$ 

Induced classes from projections  $\beta \cdot \Pi x$   $\Pi = \Pi \Pi$ 





# Problems with the Proxy



 $\Pi$  ignores target-unique features!





1. A computable bound



2. Description of shared representations X

3. Suggestions for future research





Target Generalization Bounds using Discrepancy Distance
 [BBCKPW 2009]
 [Mansour et al. 2009]

 Coupled Subspace Learning [BFK 2010]





Assumption 1:  $\mathbb{E}_S[Y|x] = \mathbb{E}_T[Y|x] = \beta \cdot x$ 

## $\beta \cdot x$ can be decomposed as

target-specific can't be estimated from source alone ... yet











Assumption 2: 
$$\mathbb{E}_{S}[Y|x] = \beta_{S} \cdot (\Pi_{S}x)$$
  
 $\mathbb{E}_{T}[Y|x] = \beta_{T} \cdot (\Pi_{T}x)$ 

- Projections  $\Pi_S = \Pi_S \Pi_S \quad \Pi_T = \Pi_T \Pi_T$
- $\Pi_T$  couples (works well) and -(don't buy)

 $\Pi_S \ \& \ \Pi_T$  learned from unlabeled data

# Visualizing Dimensionality Reduction



# Visualizing Dimensionality Reduction



# **Representation Soundness**



# **Representation Soundness**



# **Representation Soundness**



![](_page_42_Picture_0.jpeg)

![](_page_42_Figure_1.jpeg)

![](_page_43_Picture_0.jpeg)

Input: Labeled source instances  $(x_i, y_i)_{i=1}^n$ Unlabeled target instances  $x_T$ 

1) Compute  $\Pi_S$  and  $\Pi_T$  (LDA, HMM, CCA) 2)  $\left( [\hat{\beta}]_{S,T} \right) = \underset{[\beta]_{S,T}}{\operatorname{argmin}} \sum_i \left( [\beta]_{S,T} \Pi_T [x_i]_{S,T} - y_i \right)^2$ 

3) For target instance x, predict  $[\beta]_{S,T}\Pi_T x$ 

![](_page_44_Picture_0.jpeg)

# Let $\Sigma_T = I$ n = num source instances $\Sigma_{S \to T} = \sum_i (\Pi_T[x_i]_{S,T}) (\Pi_T[x_i]_{S,T})^\top$ $\lambda_j = \text{eigenvalues of } \Sigma_{S \to T}$

Under perfect adaptation, we have

$$\ell_T([\hat{\beta}]_{S,T}) - \ell_T(\beta_T^*) \le$$

![](_page_45_Picture_0.jpeg)

# $\sum \frac{d}{dt}$ when $S = T^{[T_n]_{S,T}}$ Under perfect adaptation, where $\ell_T([\hat{\beta}]_{S,T}) - \ell_T(\beta_T^*) \leq \left(\frac{\sum j \ \overline{\lambda_j}}{m}\right)$

Computing  $\Pi_S$  and  $\Pi_T$ 

Canonical Correlation Analysis (CCA) [Hotelling 1935]

1) Divide feature space into disjoint views

Do **not buy** the Shark portable steamer. The trigger mechanism is **defective**.

![](_page_46_Figure_4.jpeg)

2) Find maximally correlating projections  $\Pi_{T} = \begin{bmatrix} \Pi_{T}^{1} & 0 \\ 0 & \Pi_{T}^{2} \end{bmatrix}$ 

![](_page_47_Picture_0.jpeg)

Canonical Correlation Analysis (CCA) [Hotelling 1935]

Ando and Zhang (ACL 2005)

Kakade and Foster (COLT 2006)

2) Find maximally correlating projections  $\Pi_{T} = \begin{bmatrix} \Pi_{T}^{1} & 0 \\ 0 & \Pi_{T}^{2} \end{bmatrix}$ 

# Square Loss: Kitchen Appliances

![](_page_48_Figure_1.jpeg)

# Square Loss: Kitchen Appliances

![](_page_49_Figure_1.jpeg)

![](_page_50_Picture_0.jpeg)

![](_page_50_Figure_1.jpeg)

![](_page_51_Figure_0.jpeg)

![](_page_51_Figure_1.jpeg)

![](_page_52_Figure_0.jpeg)

![](_page_52_Figure_1.jpeg)

![](_page_53_Figure_0.jpeg)

![](_page_53_Figure_1.jpeg)

![](_page_54_Picture_0.jpeg)

### Only label really new target instances

$$\Sigma_{T \to S} = \sum_{x_i \in T} \left( \prod_{S} [x_i]_{S,T} \right) \left( \prod_{S} [x_i]_{S,T} \right)^{\top}$$

Order  $x \in \mathcal{X}_T$  by

Piyush Rai et al. (2010)

$$\frac{\|\Pi_T x\|_{\Sigma_T^{-1}}^2}{\|\Pi_T x\|_{\Sigma_T^{-1}}^2}$$

Ratio is

- 1 when S = T
- $\infty$  when  $\Pi_T x$  has no shared part

![](_page_55_Picture_0.jpeg)

1. A computable bound

![](_page_55_Picture_2.jpeg)

2. Description of shared representations

3. Suggestions for future research

![](_page_55_Picture_5.jpeg)

 $\checkmark$ 

![](_page_56_Picture_0.jpeg)

1. Theory can help us understand domain adaptation better

2. Good theory suggests new directions for future research

- 3. There's still a lot left to do
  - Connecting supervised and unsupervised adaptation
  - Unsupervised adaptation for problems with structure

![](_page_57_Picture_0.jpeg)

#### **Collaborators**

Shai Ben-David Koby Crammer Dean Foster Sham Kakade

Alex Kulesza Fernando Pereira Jenn Wortman

#### <u>References</u>

Ben-David et al. <u>A Theory of Learning from Different Domains</u>. Machine Learning 2009. Mansour et al. <u>Domain Adaptation: Learning Bounds and Algorithms</u>. COLT 2009.