

ON THE HARDNESS OF DOMAIN ADAPTATION

Tamas Madarasz & Michael Rabadi

April 15, 2015

QUESTIONS

- DA under the covariate shift assumption can succeed if we are given/can estimate the weight ratios. But what happens if we are not given these ratios and some (all?) of the points in the **source** sample are not found in the **target** sample?
- Can we find satisfying error guarantees without prior knowledge about the **source** but strong prior knowledge about the **target** task?
- Can we trade the 'expensive' labeled examples from the **source** distribution for unlabeled **target**-generated data.

ASSUMPTIONS

- Covariate shift
- Pointwise density ratio between source and target distributions is bounded below by 0.5
- Realizable hypothesis class with smallest possible VC dimension

Solvability:

Let \mathcal{W} be a class of triples (P_S, P_T, l) of source and target distributions over some domain \mathcal{X} and a labeling function l .

The DA learner \mathcal{A} **(ϵ, δ, m, n) -solves** DA for the class \mathcal{W} , if for \forall triples $(P_S, P_T, l) \in \mathcal{W}$,

labeled sample S of size m i.i.d. from P_S ,

unlabeled sample T of size n from P_T ,

with probability at least $1 - \delta$ \mathcal{A} outputs a function h with $Err_l^T(h) \leq \epsilon$.

d_A distance

$d_{H\Delta H}$ distance: $d_{H\Delta H}(P_S, P_T) = \sup_{A \in H\Delta H} |P_T(A) - P_S(A)|$

where $H\Delta H = \{h_1\Delta h_2 \mid h_1, h_2 \in H\}$,

$h_1\Delta h_2 = \{x \in \mathcal{X} \mid h_1(x) \neq h_2(x)\}$

Weight Ratio: For $\mathcal{B} \subseteq 2^{\mathcal{X}}$, source and target distributions P_S and P_T

$$C_{\mathcal{B}}(P_S, P_T) = \inf_{b \in \mathcal{B}, P_T(b) \neq 0} \frac{P_S(b)}{P_T(b)}$$

We use $C(P_S, P_T)$ if \mathcal{B} is the collection of all sets that are P_S and P_T -measurable.

We will bound the weight ratio by $C(P_S, P_T) \geq 1/2$.

THEOREM 1

Theorem 1: Let \mathcal{X} be a finite domain, and $\epsilon + \delta < \frac{1}{2}$.

No algorithm can (ϵ, δ, s, t) -solve the DA problem for the class \mathcal{W} of triples (P_S, P_T, l) with

- $C(P_S, P_T) \geq 1/2$
- $d_{H_{1,0} \Delta H_{1,0}}(P_S, P_T) = 0$
- $opt_T^l(H_{1,0}) = 0$

if $s + t < \sqrt{(1 - 2(\epsilon + \delta))|\mathcal{X}|} - 2$.

We prove this theorem by reducing the related **Left/Right Problem** (Kelly et al. 2010) to Domain Adaptation under specific assumptions.

LEFT/RIGHT PROBLEM

The Left/Right problem: Given three finite samples, L, R and M of points from some domain set \mathcal{X} , with

- L i.i.d. sample from some distribution P,
- R an i.i.d. sample from some distribution Q over \mathcal{X} ,

If M is an i.i.d. sample generated by one of these two probability distributions, can we tell which one it is?

LEFT/RIGHT PROBLEM II

An algorithm (δ, l, r, m) -**solves** the L/R problem if, given samples L , R and M of sizes l , r and m respectively, it gives the correct answer with probability at least $1 - \delta$.

LEFT/RIGHT PROBLEM III

More formally

Let $\mathcal{W}_n^{uni} = \{(U_A, U_B, U_C) : A \cup B = \{1, \dots, n\}, A \cap B = \emptyset, |A| = |B|, \text{ and } C = A \text{ or } C = B\}$, where U_Y denotes the uniform distribution over the set Y .

Lemma 1: For any given sample sizes l for L , r for R and m for M and any $0 < \gamma < 1/2$, if $k = \max\{l, r\} + m$, then for $n > (k + 1)^2 / (1 - 2\gamma)$ no algorithm has probability of success greater than $1 - \gamma$ over the class \mathcal{W}_n^{uni} .

LEFT/RIGHT PROBLEM IV

Idea for Proof (Batu et al. 2010): The L/R problem is permutation invariant (doesn't depend on permutations of X). Find large enough n that with high probability no element is repeated more than once across the three samples.

Proposition 1: Let X be a finite domain of size n . For every $0 < \delta < 1$, with probability $> (1 - \delta)$, an i.i.d. sample of size at most $\sqrt{\delta n} - \delta$ uniformly drawn over X , contains no repeated elements.

Permutation invariance

For multi-sets L, R, M of size $\leq n$, from distributions P or Q over some domain X the **fingerprint F** is defined as $\{C_{i,j,k} | 1 \leq i, j, k \leq n\}$, where $C_{i,j,k}$ is the number of elements of X , that appear exactly i times in L , j times in R and k times in M .

Proposition 2: If there exists an algorithm A for testing some permutation-invariant property of distributions, then there exists an algorithm for that same task that gets as input only the fingerprints of the samples that A takes and enjoys the same guarantee on its probability of success.

LEFT/RIGHT PROBLEM VI

Proof of Lemma 1 Let $\delta = 1 - 2\gamma$. Then with probability $> (1 - \delta)$ the input to the Left/Right problem over \mathcal{W}_n^{uni} has no repeated elements, and the fingerprint F has

$C_{1,0,0} = l, C_{0,1,0} = r, C_{0,0,1} = m$, and $C_{i,j,k} = 0$ for all other i, j and k , *independently* of M coming from U_A or U_B .

Let A be some algorithm, $p = P(A \text{ outputs } U_A | F)$

If $p \geq 0.5$, then $P(A \text{ errs}) > (1 - \delta)/2$ for all triples where C is equal to B .

If $p \leq 0.5$, then $P(A \text{ errs}) > (1 - \delta)/2$ on all triples where C is equal to A .

But $(1 - \delta)/2 = \gamma$, thus, no algorithm can (γ, l, r, m) -solve the Left/Right problem for the class \mathcal{W}_n^{uni} .

LEFT/RIGHT PROBLEM TO DA

How do we get from DA to the L/R problem?

For n , let \mathcal{W}_n^{DA} be the class of triples (P_S, P_T, l) , where

- \mathcal{X} is a finite set of size n ,
- $P_S \sim U(\mathcal{X})$
- P_T is uniform over some subset U of X of size $n/2$
- $l(x) = 1$ for $x \in U$ and $l(x) = 0$ for $x \in \mathcal{X} \setminus U$.

Notably $C(P_S, P_T) = 1/2$.

Lemma 2. For $n \in \mathcal{N}$ and an algorithm A that can (ϵ, δ, s, t) -solve DA for \mathcal{W}_n^{DA} given that the Target is realizable by $H_{1,0}$, we can construct an algorithm that $(\epsilon + \delta, s, s, t + 1)$ -solves the Left/Right problem on \mathcal{W}_n^{uni} .

LEFT/RIGHT PROBLEM TO DA II

Proof

Given sample $L = \{l_1, l_2, \dots, l_s\}$ and $R = \{r_1, r_2, \dots, r_s\}$, sample M of size $t + 1$ for L/R in \mathcal{W}_n^{uni}

Let $T = M \setminus \{p\}$, $p \in M$ chosen uniformly at random.

Construct S from $L \times \{0\} \cup R \times \{1\}$ by successively flipping an unbiased coin, and choosing the next element from $L \times \{0\} \cup R \times \{1\}$ accordingly.

Then P_T of this DA problem has marginal either U_A or U_B and the labeling function of this Domain Adaptation instance is $l(x) = 0$ if $x \in A$ and $l(x) = 1$ if $x \in B$.

If \mathcal{A} outputs h on input S and T , then Left/Right problem outputs U_A if $h(p) = 0$, U_B if $h(p) = 1$, and $(\epsilon + \delta, s, s, t + 1)$ -solves the Left/Right problem.

THEOREM 1 CONCLUDED

Thus if an algorithm A that can (ϵ, δ, s, t) -solve DA by Lemma 1 we have $|X| \leq (s + (t + 1) + 1)^2 / (1 - 2(\epsilon + \delta))$, thus

$s + t \geq \sqrt{(1 - 2(\epsilon + \delta))|X|} - 2$ completing the proof of Theorem 1.

Corollary:

Theorem 2:

Let $\mathcal{X} = [0, 1]^d$, $\epsilon, \delta > 0$ s.t $\epsilon + \delta < 1/2$.

Let $\lambda > 1$ and let \mathcal{W} be the set of triples (P_S, P_T, l) of distributions over \mathcal{X} , with covariance shift and realizability as before, and λ -Lipschitz labeling function l .

Then no DA-learner can (s, t, ϵ, δ) -solve the DA-problem for the class \mathcal{W} unless $s + t \geq \sqrt{(\lambda + 1)^d (1 - 2(\epsilon + \delta))} - 2$.

THEOREM 2

Proof:

Let $G \subseteq \mathcal{X}$ be the points of a grid in $[0, 1]^d$ with distance $1/\lambda$. Then $|G| = (\lambda + 1)^d$, and any labeling function $l : G \rightarrow \{0, 1\}$ is λ -Lipschitz.

As G is finite, the bound follows from Theorem 1.

CITATIONS

Ben-David, S., and Urner, R. (2012). On the hardness of domain adaptation and the utility of unlabeled target samples. In Algorithmic Learning Theory (pp. 139-153). Springer Berlin Heidelberg.

Bernstein, S. (1946). The theory of Probabilities.

Cortes, C., Mansour, Y., and Mohri, M. (2010). Learning bounds for importance weighting. In Advances in neural information processing systems (pp. 442-450).

Jin, J., Wilson, J., Nobel, A. (2014). Jensen's and Holder's Inequalities.

Kelly, B.G., Tularak, T., Wagner, A.B., Viswanath, P.: Universal hypothesis testing in the learning-limited regime.(2010) IEEE, ISIT.

ALGORITHM

Algorithm A Input An i.i.d. sample S from P_S labeled by l , an unlabeled i.i.d. sample T from P_T , and a margin parameter γ .

Step 1 Partition the domain $[0, 1]^d$ into a collection B of boxes (axis-aligned rectangles) with sidelength (γ/\sqrt{d}) .

Step 2 Obtain sample S' by removing every point in S , which is sitting in a box that is not hit by T .

Step 3 Output an ERM classifier from H for the sample S'