
Supplementary Material for Support and Invertibility in Unsupervised Domain Adaptation

Fredrik D. Johansson
MIT

David Sontag
MIT

Rajesh Ranganath
NYU

A Proofs

A.1 Proof of bounds for support sufficiency divergence

Lemma 1. *The support sufficiency divergence is bounded with $0 \leq d_{\text{supp}}^\epsilon(p \parallel q) \leq 1$, and the bounds are tight.*

Proof. The lower bound holds and is tight because

$$\begin{aligned} d_{\text{supp}}^\epsilon(p \parallel q) &= \int_x (q(x) - p(x)) \delta_{p,q}(x) \\ &= \int_x \max(q(x) - p(x), 0) \mathbf{1}[p(x) \leq \epsilon] \mathbf{1}[p(x) \leq \epsilon]. \end{aligned}$$

which is clearly non-negative. Moreover, for $\epsilon \leq \inf_x p(x)$, $d_{\text{supp}}^\epsilon(p \parallel q) = 0$. The upper bound holds trivially as $\delta_{p,q}(x) \leq 1$. For tightness, let q, p be discrete densities over two states, $q = [1., 0.]$ and $p = [0., 1.]$. Then with $\epsilon > 0$, $d_{\text{supp}}^\epsilon(p \parallel q) = 1$. \square

Recall that

$$w_{p,q}^\epsilon(x) = \begin{cases} q(x)/p(x) & \text{if } p(x) \geq \epsilon \\ 1 & \text{otherwise} \end{cases} \quad (1)$$

A.2 Proof of Lemma 1

Lemma 2. *Let p, q be densities over \mathcal{X} . Further, define $\delta_{p,q}(x) = \mathbf{1}[p(x) \leq \epsilon \text{ and } p(x) \leq q(x)]$. Then,*

$$\mathbb{E}_q[f(x)] \leq \mathbb{E}_p[w_{p,q}^\epsilon(x)f(x)] + M \cdot (\mathbb{E}_q[\delta_{p,q}^\epsilon(x)] - \mathbb{E}_p[\delta_{p,q}^\epsilon(x)])$$

Proof. We have,

$$\begin{aligned} \mathbb{E}_q[f(x)] &= \int_x q(x)f(x)dx && \text{(By def.)} \\ &= \int_{x:p(x)>\epsilon} q(x)f(x)dx + \int_{x:p(x)\leq\epsilon} q(x)f(x)dx && \text{(Dividing domain)} \\ &= \int_{x:p(x)>\epsilon} \frac{q(x)}{p(x)}p(x)f(x)dx + \int_{x:p(x)\leq\epsilon} q(x)f(x)dx && \text{(Multiply and divide)} \\ &\leq \mathbb{E}_p[w_{p,q}^\epsilon(x)f(x)] + \int_{x:p(x)\leq\epsilon} (q(x) - p(x))f(x)dx && \text{(Add and subtract } p(x)) \\ &\leq \mathbb{E}_p[w_{p,q}^\epsilon(x)f(x)] + M \int_{\substack{x:p(x)\leq\epsilon \\ p(x)\leq q(x)}} (q(x) - p(x))dx && \text{(By assmp. \& } f, p, q \geq 0) \end{aligned}$$

$$\begin{aligned}
&= \mathbb{E}_p [w_{p,q}^\epsilon(x) f(x)] \\
&+ M \int_x (q(x) - p(x)) \underbrace{\mathbb{1}[p(x) \leq \epsilon \wedge p(x) \leq q(x)]}_{\delta_{p,q}^\epsilon(x)} dx && \text{(By domain of int.)} \\
&= \mathbb{E}_p [w_{p,q}^\epsilon(x) f(x)] + M \cdot \left(\mathbb{E}_q[\delta_{p,q}^\epsilon(x)] - \mathbb{E}_p[\delta_{p,q}^\epsilon(x)] \right) && \text{(By def.)}
\end{aligned}$$

Further, $p = q$ implies equality when $\epsilon \geq \sup_x q(x)$. □

A.3 Proof of Theorem 2

Lemma 3. Assume that $p_t(Y | X) = p_s(Y | X)$. Define $Z = \phi(X)$ and let $h(x) = f(\phi(x))$. Then,

$$\mathbb{E}_{x,y \sim q(x,y)} [\ell(h(x), y)] = \mathbb{E}_{z,y \sim q(z,y)} [\ell(f(z), y)]$$

Proof.

$$\begin{aligned}
\mathbb{E}_{z,y \sim q(z,y)} [\ell(f(z), y)] &= \int_{z,y} q(z,y) \ell(f(z), y) dz dy \\
&= \int_{z,y} \ell(f(z), y) \int_{x \in \phi^{-1}(z)} q(x,y) dx dz dy \\
&= \int_{x,y} q(x,y) \int_z \mathbb{1}[z = \phi(x)] \ell(f(z), y) dz dx dy \\
&= \int_{x,y} q(x,y) \ell(h(x), y) dx dy \\
&= \mathbb{E}_{x,y \sim q(x,y)} [\ell(h(x), y)]
\end{aligned}$$

□

Lemma 4.

$$R_t(h) = E_{q(z)p(y|z)}[\ell(f(z), y)] + \eta_\phi^\ell(f, y)$$

Proof. By Lemma A.3

$$R_t(h) = E_{q(x,y)}[\ell(h(x), y)] = E_{q(z,y)}[\ell(f(z), y)]$$

We have that

$$\begin{aligned}
E_{q(z)q(y|z)}[\ell(f(z), y)] &= \int_z \int_y q(z)q(y|z) \ell(f(z), y) dy dz \\
&= \int_z \int_y \left(\int_{x: \phi(x)=z} q(x) dx \right) q(y|z) \ell(f(z), y) dy dz \\
&= \int_z \int_y \int_{x \in \phi^{-1}(z)} q(x)q(y|\phi(x)) \ell(h(\phi(x)), y) dx dy dz \\
&= \int_x \int_y q(x)q(y|\phi(x)) \ell(h(\phi(x)), y) dx dy \\
&= \int_x \int_y q(x) \ell(h(\phi(x)), y) \left(\underbrace{q(y|x)}_{=p(y|x) \text{ (assmp.)}} + q(y|\phi(x)) - q(y|x) \right) dx dy
\end{aligned}$$

and by the same argument,

$$E_{q(z)p(y|z)}[\ell(f(z), y)] = \int_x \int_y q(x) \ell(h(\phi(x)), y) (p(y | x) + p(y | \phi(x)) - p(y | x)) dx dy$$

and as a result,

$$\begin{aligned} & E_{q(z)q(y|z)}[\ell(f(z), y)] - E_{q(z)p(y|z)}[\ell(f(z), y)] \\ &= \mathbb{E}_q(x) \left[\mathbb{E}_{q(y|\phi(x))} \ell(f(\phi(x)), y) \right] - \mathbb{E}_{q(y|x)} \ell(f(\phi(x)), y) \\ &+ \mathbb{E}_q(x) \left[\mathbb{E}_{p(y|\phi(x))} \ell(f(\phi(x)), y) \right] - \mathbb{E}_{p(y|x)} \ell(f(\phi(x)), y) \\ &= \eta_\phi^\ell(f, y) \end{aligned}$$

□

Theorem 2 (Restated). Consider any feature representation $z = \phi(x)$ with $\phi : \mathcal{X} \rightarrow \mathcal{Z}$ and prediction function $f : \mathcal{Z} \rightarrow \mathcal{Y}$, and define $h = f \circ \phi$. Further, let $p(Z)$ and $q(Z)$ be the two distributions induced by the representation ϕ applied to X distributed according to $p(X), q(X)$. Further, assume that for any hypothesis h and a loss function ℓ , $\sup_{x \in \mathcal{X}, h \in \mathcal{H}} \ell(h(x), y) \leq M$. Now, with $\epsilon > 0$, we have the following result.

$$R_q(h) \leq \mathbb{E}_p [w_{p_s, p_t}^\epsilon(z) \ell(f(z), y)] + M d_{\text{supp}}^\epsilon(p(z) \parallel q(z)) + \eta_\phi^\ell(f, y) .$$

Proof. By Lemma 4, we have that

$$R_q(h) \leq E_{q(z)p(y|z)}[\ell(f(z), y)] + \eta_\phi^\ell(f, y) .$$

Further,

$$\begin{aligned} E_{q(z)p(y|z)}[\ell(f(z), y)] &= \iint_{z \in \mathcal{Z}, y \in \mathcal{Y}} q(z) p(y | z) \ell(f(z), y) dy dz \\ &= \iint_{z: p(z) \geq \epsilon, y \in \mathcal{Y}} q(z) p(y | z) \ell(f(z), y) dy dz \\ &+ \iint_{z: p(z) < \epsilon, y \in \mathcal{Y}} q(z) p(y | z) \ell(f(z), y) dy dz \\ &= \iint_{z \in \mathcal{Z}, y \in \mathcal{Y}} w_{p_s, p_t}^\epsilon(z) p(z) p(y | z) \ell(f(z), y) dy dz \\ &+ \int_{z: \substack{p(z) < \epsilon \\ p(z) \leq q(z)}} \underbrace{(q(z) - p(z))}_{\geq 0} \underbrace{\int_{y \in \mathcal{Y}} p(y | z) \ell(f(z), y) dy}_{\in [0, M]} dz \\ &+ \underbrace{\int_{z: \substack{p(z) < \epsilon \\ p(z) > q(z)}} (q(z) - p(z)) \int_{y \in \mathcal{Y}} p(y | z) \ell(f(z), y) dy dz}_{\leq 0} . \end{aligned}$$

□

A.4 Kernel support divergence

Theorem 3 may be viewed as a measuring differences in density only where supports differ significantly. In the case where \mathcal{L} is a Hilbert space, similar to the maximum mean discrepancy (Gretton *et al.*, 2012), we may decompose $d_{\text{supp}}^{\mathcal{L}, \epsilon}(p \parallel q)$ using reproducing kernels.

Lemma 5. Let \mathcal{H} be the reproducing-kernel Hilbert space with kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$. Then,

$$\begin{aligned} d_{\text{supp}}^\epsilon(p \parallel q)_{\mathcal{G}} &= \mathbb{E}_{x, x' \sim p} [\delta_p^\epsilon(x, x') k(x, x')] \\ &- 2 \mathbb{E}_{x \sim p, x' \sim q} [\delta_p^\epsilon(x, x') k(x, x')] + \mathbb{E}_{x, x' \sim q} [\delta_p^\epsilon(x, x') k(x, x')] \end{aligned} \quad (2)$$

Proof. The proof follows from Gretton *et al.* (2012) and can be found in Appendix A.4. \square

Let $\delta_p^\epsilon(x) = \mathbb{1}[p(x) \leq \epsilon]$ and $\delta_p^\epsilon(x, x') = \delta_p^\epsilon(x)\delta_p^\epsilon(x')$

$$\begin{aligned} d_{\text{supp}}^\epsilon(p \parallel q)_{\mathcal{G}} &:= \sup_{g \in \mathcal{G}} \left| \mathbb{E}_q[\delta_p^\epsilon(x)g(x)] - \mathbb{E}_p[\delta_p^\epsilon(x)g(x)] \right| \\ &= \mathbb{E}_{x, x' \sim p} [\delta_p^\epsilon(x, x')k(x, x')] \\ &\quad - 2 \mathbb{E}_{x \sim p, x' \sim q} [\delta_p^\epsilon(x, x')k(x, x')] \\ &\quad + \mathbb{E}_{x, x' \sim q} [\delta_p^\epsilon(x, x')k(x, x')] \end{aligned} \quad (3)$$

Proof. Follow http://alex.smola.org/teaching/iconip2006/iconip_3.pdf page 18–20. \square

B Model

We may bound $d_{\text{supp}}^\epsilon(p \parallel q)$ using the hinge loss as follows,

$$\begin{aligned} d_{\text{supp}}^\epsilon(p \parallel q) &\leq \mathbb{E}_{x \sim q} \left[\max \left(0, 2 - \frac{p(x)}{\epsilon} \right) \max \left(0, 2 - \frac{p(x)}{q(x)} \right) \right] \\ &\quad - \mathbb{E}_{x \sim p} \left[\max \left(0, 1 - \frac{p(x)}{\epsilon} \right) \max \left(0, 1 - \frac{p(x)}{q(x)} \right) \right] \\ &=: \tilde{d}_{\text{supp}}^\epsilon(p \parallel q). \end{aligned}$$

C Experiments

In Figure 1, we see that the embeddings learned using DANN models under label marginal shift show worse separation between classes, than the embeddings learned under equal label marginal distributions.

D Consistent domain-invariant variable selection

Consider a matrix $A \in \{0, 1\}^{k \times d}$ with $k < d$ such that $\forall j : \sum_{i=1}^k a_{ij} \leq 1$ and $\forall i : \sum_{j=1}^d a_{ij} \leq 1$. In other words, A is a variable selection operator on X . Now, assume that $Z := \Phi(X) := AX$ is sufficient for Y on p_s and p_t and that $p_s(AX) = p_t(AX)$. Further, assume labeled data is observed under p and unlabeled data observed under p_t . Then, is Y identifiable based on domain-invariance and source predictive loss?

Condition 1 (Smoothness). With $\Sigma_L = \{f : \sum_{k \in \mathbb{Z}^d} k_2^2 \langle f, \varphi_k \rangle^2 \leq L; \forall j \in \{1, \dots, d\} \text{ for } L > 0, f \text{ is } L\text{-smooth if } f \in \Sigma_L, \text{ with } \varphi_k \text{ the trigonometric fourier basis.}$

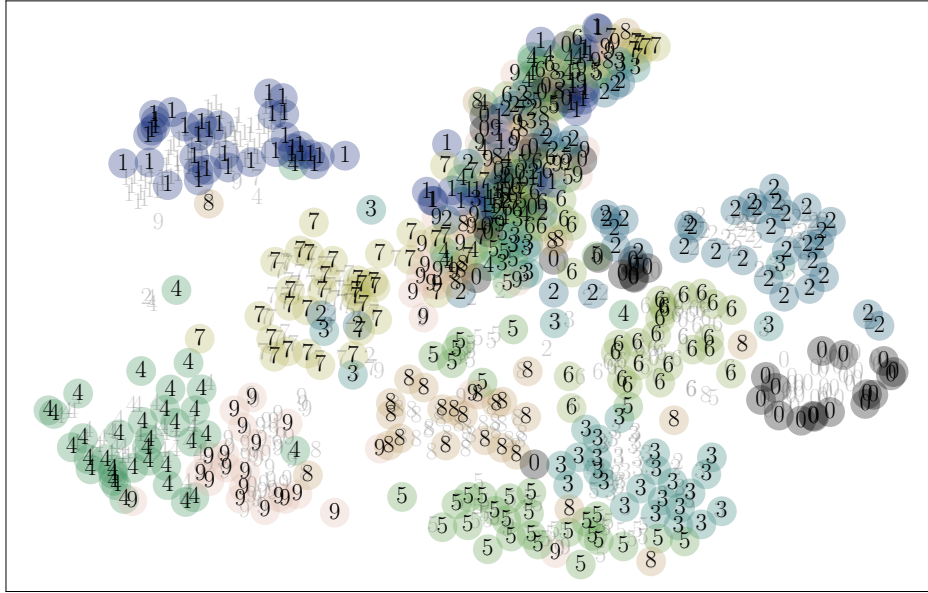
Condition 2 (Identifiability). A sufficient set of variables J , such that $\exists \bar{f} : f(x) = \bar{f}(x_J)$ for all $x \in \mathbb{R}^d$, is κ -identifiable if for all $j \in J$,

$$\int_{[0,1]^d} (f(x) - \int_0^1 f(x) dx_j)^2 dx \geq \kappa.$$

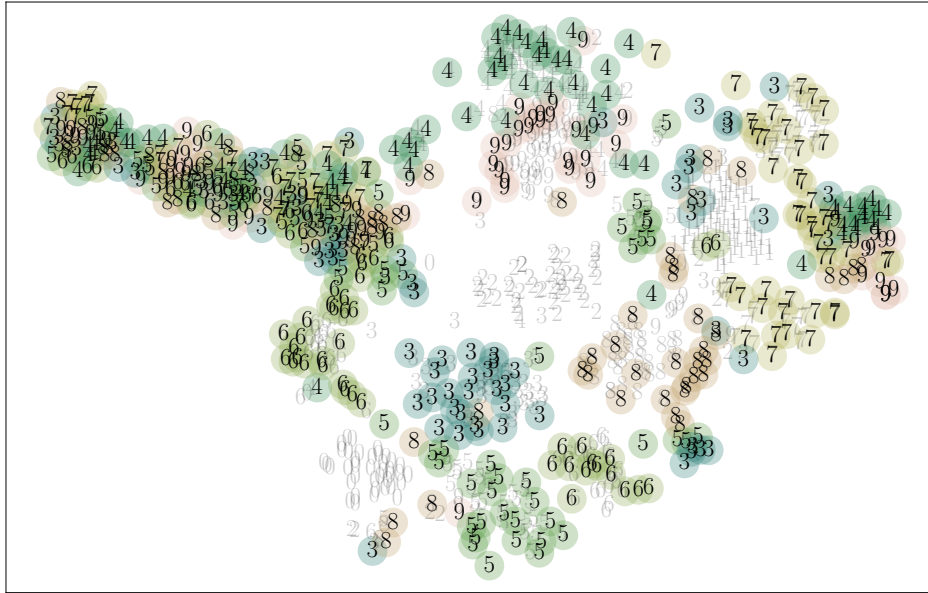
Condition 3 (Positive bounded support). The density $p_s(x)$ has positive bounded support over $[0, 1]^d$ if with $p_{s\text{min}} > 0, \forall x \in [0, 1]^d : p_s(x) \geq p_{s\text{min}}$ and $\forall x \notin [0, 1]^d : p_s(x) = 0$.

Condition 4 (Bounded ∞ -norm and 2-norm). A function f has bounded ∞ -norm and 2-norm with respect to p if $\Pr_{X \sim p_s}(|f(X)| \leq L_\infty) = 1$ and $\mathbb{E}_{X \sim p_s}[f(X)^2] \leq L_2^2$.

Condition 5 (Sub-gaussian additive noise). The observed outcome may be written as $Y_i = f(x_i) + \sigma \epsilon_i$ with $\mathbb{E}[e^{t\epsilon_i} | X_i] \leq e^{t^2/2}$ for all $t > 0$.



(a) MNIST \rightarrow MNIST-M



(b) MNIST \rightarrow MNIST-M $\setminus \{0, 1, 2\}$

Figure 1: Embeddings learned by DANN with equal (top) and unequal (bottom) label marginal distributions. In MNIST-M $\setminus \{0, 1, 2\}$, all images of digits 0,1,2 have been removed. Grey digits are embeddings and labels of the source domain. Black digits against colored background are from the target domain.

Theorem 1 (Variable selection in non-parametric regression (Comminges *et al.*, 2012)). *Assume that Conditions 1–5 hold, with known parameters $p_{s_{\min}}$, $\theta = 2L/\kappa$ and L_2 . Then, there is an estimator \hat{J} that satisfies $\Pr(\hat{J} \neq J) \leq (8d/d^*)^{-d^*}$.*

Comminges *et al.* (2012) give a constructive proof of Theorem 1 in which the chosen estimator is allowed to depend on the density $p_s(x)$.

References

Comminges, L., Dalalyan, A.S. *et al.* (2012). Tight conditions for consistency of variable selection in the context of high dimensionality. *The Annals of Statistics*, **40**, 2667–2696.

Gretton, A., Borgwardt, K.M., Rasch, M.J., Schölkopf, B. & Smola, A. (2012). A kernel two-sample test. *Journal of Machine Learning Research*, **13**, 723–773.