

# Marginal Singularity, and the Benefits of Labels in Covariate-Shift

Samory Kpotufe\*

ORFE, Princeton University

SAMORY@PRINCETON.EDU

Guillaume Martinet†

ORFE, Princeton University

GGM2@PRINCETON.EDU

## Abstract

We present new minimax results that concisely capture the relative benefits of source and target labeled data, under covariate-shift. Namely, we show that, in general classification settings, the benefits of target labels are controlled by a *transfer-exponent*  $\gamma$  that encodes how *singular*  $Q$  is locally w.r.t.  $P$ , and interestingly allows situations where transfer did not seem possible under previous insights. In fact, our new minimax analysis – in terms of  $\gamma$  – reveals a *continuum of regimes* ranging from situations where target labels have little benefit, to regimes where target labels dramatically improve classification. We then show that a recently proposed semi-supervised procedure can be extended to adapt to unknown  $\gamma$ , and therefore requests target labels only when beneficial, while achieving nearly minimax transfer rates.

**Keywords:** Transfer learning, covariate-shift, nonparametric classification, nearest-neighbors.

## Extended Abstract

### Introduction

The goal in transfer learning is to improve prediction on a *target* distribution  $Q$  by harnessing labeled data coming from a *source* distribution  $P$ . Much of theoretical work in transfer learning concerns understanding the fundamental limitations of transfer, and in particular, proper ways of capturing *relatedness* between source  $P$  and target  $Q$ . Here we consider the common *covariate-shift* setting for classification, where it is assumed that conditional distributions  $P_{Y|X}$  and  $Q_{Y|X}$  remain the same, while marginals  $P_X$ ,  $Q_X$  are different but somewhat related.

We consider general nonparametric settings that capture a range of easy to difficult classification under  $Q$ , through standard smoothness and noise conditions (see e.g. [Audibert and Tsybakov \(2007\)](#)). Our aim is then to understand which relation between marginals  $P_X$  and  $Q_X$  control the rates of transfer, and in particular, control the relative benefits between source and target data in achieving low error under  $Q$ . A basic intuition, present in previous work, is that transfer is easiest when  $P$  assigns sufficient mass to regions of considerable  $Q$ -mass. Here, we formalize this intuition through a new *asymmetric* notion, the *transfer-exponent*  $\gamma$ , that parametrizes the behavior of ball-mass ratios  $Q_X(B(x, r))/P_X(B(x, r))$  as a function of the radius  $r$ , namely, that these ratios behave like  $r^{-\gamma}$ . The notion of  $\gamma$  can be interpreted, roughly, as capturing how close to *singular*  $Q$  is with respect to  $P$ , as it shifts mass into regions of low  $P$  mass.

---

. \* † Authors are listed in alphabetical order.

. Extended abstract. Full version appears as arXiv:1803.01833v2.

We show the pertinence of our parametrization by establishing tight minimax upper and lower bounds in terms of  $\gamma$ , under standard nonparametric conditions. The notion of  $\gamma$  is thus shown to encode a *continuum of regimes* between easy and hard transfer, and interestingly, reveals situations where transfer is possible (even at fast rates) despite  $P$  and  $Q$  seeming *unrelated* under previous notions of *relatedness*. As an example,  $\gamma$  remains well defined even when  $Q$  is singular w.r.t.  $P$  (e.g.  $Q$  puts mass on lower-dimensional structures) – in which case common notions of density-ratio and information-theoretic divergences (KL or Renyi) fail to exist, and common extensions of total-variation can be too large to characterize transfer.

Finally, we show that a recently proposed semi-supervised procedure can be extended to adapt to unknown  $\gamma$ , and therefore requests target labels only when beneficial, while achieving nearly minimax transfer rates.

## Related Work

Many insightful notions of *relatedness* are present in the literature on transfer and related problems.

A first line of work considers refinements of total-variation which encode changes in classification error from  $P$  to  $Q$  (restricted to a hypothesis class  $\mathcal{H}$ ). The most common such measures are the so-called  $d_{\mathcal{A}}$ -divergence (Ben-David et al., 2010a,b; Germain et al., 2013) and  $\mathcal{Y}$ -discrepancy (Mansour et al., 2009a; Mohri and Medina, 2012; Cortes et al.). These notions are the first to capture – through *differences* in mass over space – the intuition that transfer is easiest when  $P$  has sufficient mass in regions of substantial  $Q$ -mass. Typical excess-error bounds on classifiers learned from source (and some or no target) data are of the form  $o_p(1) + C \cdot \text{divergence}(P, Q)$ . In other words, transfer seems impossible when these divergences are large; however, we show that there are ranges of reasonable situations ( $0 \leq \gamma < \infty$ ) where fast transfer is possible even when such divergences are large. Furthermore, while such divergences are symmetric, the notion of  $\gamma$  is not, thus capturing the fact that transfer might be easy from  $P$  to  $Q$  but not from  $Q$  to  $P$ .

Another prominent line of work, which has led to many practical procedures, considers so-called density-ratios  $f_Q/f_P$  or more generally, Radon-Nikodym derivatives  $dQ/dP$ , as a way to capture the similarity between  $P$  and  $Q$  (Quionero-Candela et al., 2009; Sugiyama et al., 2012). It is often assumed in such work that  $dQ/dP$  is bounded, which corresponds to assuming  $\gamma = 0$ . Typical excess-error bounds are dominated by the estimation rates for  $dQ/dP$  (see e.g. rates for  $\alpha$ -Hölder  $dQ/dP$ ,  $\alpha \rightarrow 0$ , in Kpotufe (2017)), which unfortunately could be arbitrarily higher than the minimax rates we establish for the boundary case with  $\gamma = 0$ .

Finally, another line of work instead considers information-theoretic measures such as KL-divergence or Renyi divergence (Sugiyama et al., 2008; Mansour et al., 2009b). In particular, such divergences are closer in spirit to our notion of transfer-exponent  $\gamma$  (viewing  $\gamma$  as roughly characterizing the log of mass-ratios between), but are also undefined in typical scenarios with structured data where  $Q_X$  might be singular w.r.t.  $P_X$ .

## Result Overview

Our first results consider transfer settings where the learner has access to  $n_P$  labeled samples drawn from  $P$  and  $n_Q$  labeled samples drawn from  $Q$ , where typically  $n_P \gg n_Q$ . The label  $Y$  is assumed in  $\{0, 1\}$ , while the input  $X$  belongs to a compact metric space  $\mathcal{X}$ .

We work under common smoothness and low noise conditions, namely, we assume the *regression function*  $\eta(x) \doteq E[Y|X = x]$  to be  $\alpha$ -Hölder, and also that  $Q_X(0 < |\eta(X) - 1/2| \leq t) \lesssim t^\beta$

(see e.g. [Audibert and Tsybakov \(2007\)](#)). A *transfer exponent* is then defined, roughly, as any quantity  $\gamma$  that satisfies:

$$\forall x, \forall \text{ small } r, \quad P_X(B(x, r)) \gtrsim Q_X(B(x, r)) \cdot r^\gamma.$$

Two main distributional regimes are considered, which capture the difficulty of vanilla classification under  $Q_X$ . The first regime, (DM) (for *doubling measure*), roughly assumes that  $Q_X$  behaves like a uniform measure on its support (this is the most common assumption in nonparametric classification, and is sometimes termed the *strong-density assumption*). The second regime, (BCN) (for *bounded covering number*), allows for general  $Q_X$  and is most difficult with slower rates. Both regimes introduce a parameter  $d$  that might be viewed as a notion of *dimension* of the marginal  $Q_X$ .

For exact definitions we refer the reader to the archived version of this work ([Kpotufe and Martinet, 2018](#)).

Our minimax rates are then of the following form.

**Theorem 1 (Sketch)** *Call  $\mathcal{T}_{(DM)}$  (resp.  $\mathcal{T}_{(BCN)}$ ) the class of all the tuples  $(P, Q)$  under (DM) (resp. (BCN)) regime. Let  $\mathcal{T} \in \{\mathcal{T}_{(DM)}, \mathcal{T}_{(BCN)}\}$ , we have then:*

$$\inf_{\hat{h}} \sup_{(P, Q) \in \mathcal{T}} \mathbb{E}[\mathcal{E}_Q(\hat{h})] \asymp \left( n_P^{d_0/(d_0 + \gamma/\alpha)} + n_Q \right)^{-(\beta+1)/d_0},$$

where  $\mathcal{E}_Q$  represents the excess error, the infimum is taken over all classifiers  $\hat{h}$  learned on the data, the expectation is taken w.r.t. the data,  $d_0 = 2 + d/\alpha$  when  $\mathcal{T} = \mathcal{T}_{(DM)}$ , and  $d_0 = 2 + \beta + d/\alpha$  when  $\mathcal{T} = \mathcal{T}_{(BCN)}$ .

Our upper-bounds are established with a generic  $k$ -NN classifier defined over the combined source and target sample. In particular, our results imply new convergence rates of independent interest for vanilla  $k$ -NN under the BCN regime, which complements recent developments on vanilla  $k$ -NN ([Samworth et al., 2012](#); [Chaudhuri and Dasgupta, 2014](#); [Shalev-Shwartz and Ben-David, 2014](#); [Gadat et al., 2014](#)). On the other hand, our lower-bounds are established over any learner with access to both source and target samples, and interestingly, which is also allowed access to infinite unlabeled source and target data (i.e., full knowledge of  $P_X$  and  $Q_X$ ). In other words, the above rates cannot be improved (beyond constants) with access to unlabeled data, which is often an important consideration in practice given the cost of target labels ([Huang et al., 2007](#); [Ben-David and Uner, 2012](#)).

Finally, we address semi-supervised situations where the learner has access to  $n_Q$  unlabeled target data, along with  $n_P$  labeled source data, and is allowed to request (as few as possible) target labels in order to improve classification ([Saha et al., 2011](#); [Chen et al., 2011](#); [Chattopadhyay et al., 2013](#)). An early theoretical treatment of this can be found in ([Yang et al., 2013](#)), but which however considers a transfer setting with fixed marginal but varying conditionals (labeling functions). For our setting of covariate-shift, we build on a recent approach of [Berlind and Uner \(2015\)](#) which constructs so-called  $k$ - $2k$  covers, to help limit label requests to regions of low  $P$  mass. In this work, we show a strategy for choosing  $k$  from data (building on so-called *Lepski's method* ([Lepski and Spokoiny, 1997](#))), so as to nearly attain the above minimax rates with no a priori knowledge of distributional parameters, nor of  $\gamma$ . Furthermore, labeling complexity is shown to be controlled by unknown  $\gamma$ , i.e. the resulting approach requests labels only when *useful*, as controlled by  $\gamma$  and relative sample sizes  $n_P, n_Q$ .

## References

- Jean-Yves Audibert and Alexandre B Tsybakov. Fast learning rates for plug-in classifiers. *The Annals of Statistics*, 35(2):608–633, 2007.
- Shai Ben-David and Ruth Uner. On the hardness of domain adaptation and the utility of unlabeled target samples. In *International Conference on Algorithmic Learning Theory*, pages 139–153. Springer, 2012.
- Shai Ben-David, John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jennifer Wortman Vaughan. A theory of learning from different domains. *Machine learning*, 79(1-2):151–175, 2010a.
- Shai Ben-David, Tyler Lu, Teresa Luu, and Dávid Pál. Impossibility theorems for domain adaptation. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, pages 129–136, 2010b.
- Christopher Berlind and Ruth Uner. Active nearest neighbors in changing environments. In *International Conference on Machine Learning*, pages 1870–1879, 2015.
- Rita Chattopadhyay, Wei Fan, Ian Davidson, Sethuraman Panchanathan, and Jieping Ye. Joint transfer and batch-mode active learning. In Sanjoy Dasgupta and David McAllester, editors, *Proceedings of the 30th International Conference on Machine Learning*, volume 28 of *Proceedings of Machine Learning Research*, pages 253–261, 2013.
- Kamalika Chaudhuri and Sanjoy Dasgupta. Rates of convergence for nearest neighbor classification. In *Advances in Neural Information Processing Systems*, pages 3437–3445, 2014.
- Minmin Chen, Kilian Q Weinberger, and John Blitzer. Co-training for domain adaptation. In *Advances in neural information processing systems*, pages 2456–2464, 2011.
- Corinna Cortes, Mehryar Mohri, and Andrés Muñoz Medina. Adaptation based on generalized discrepancy. *Machine Learning Research*, forthcoming. URL <http://www.cs.nyu.edu/~mohri/pub/daj.pdf>.
- Sébastien Gadat, Thierry Klein, and Clément Marteau. Classification with the nearest neighbor rule in general finite dimensional spaces: necessary and sufficient conditions. *arXiv preprint arXiv:1411.0894*, 2014.
- Pascal Germain, Amaury Habrard, François Laviolette, and Emilie Morvant. A pac-bayesian approach for domain adaptation with specialization to linear classifiers. In *International Conference on Machine Learning*, pages 738–746, 2013.
- Jiayuan Huang, Arthur Gretton, Karsten M Borgwardt, Bernhard Schölkopf, and Alex J Smola. Correcting sample selection bias by unlabeled data. In *Advances in neural information processing systems*, pages 601–608, 2007.
- Samory Kpotufe. Lipschitz density-ratios, structured data, and data-driven tuning. In *Artificial Intelligence and Statistics*, pages 1320–1328, 2017.

- Samory Kpotufe and Guillaume Martinet. Marginal singularity, and the benefits of labels in covariate-shift. *arXiv preprint arXiv:1803.01833*, 2018.
- Oleg V Lepski and VG Spokoiny. Optimal pointwise adaptive methods in nonparametric estimation. *The Annals of Statistics*, pages 2512–2546, 1997.
- Yishay Mansour, Mehryar Mohri, and Afshin Rostamizadeh. Domain adaptation: Learning bounds and algorithms. *arXiv preprint arXiv:0902.3430*, 2009a.
- Yishay Mansour, Mehryar Mohri, and Afshin Rostamizadeh. Multiple source adaptation and the rényi divergence. In *Proceedings of the Twenty-Fifth Conference on Uncertainty in Artificial Intelligence*, pages 367–374. AUAI Press, 2009b.
- Mehryar Mohri and Andres Munoz Medina. New analysis and algorithm for learning with drifting distributions. In *International Conference on Algorithmic Learning Theory*, pages 124–138. Springer, 2012.
- Joaquin Quionero-Candela, Masashi Sugiyama, Anton Schwaighofer, and Neil D Lawrence. *Dataset shift in machine learning*. The MIT Press, 2009.
- Avishek Saha, Piyush Rai, Hal Daumé, Suresh Venkatasubramanian, and Scott L DuVall. Active supervised domain adaptation. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 97–112. Springer, 2011.
- Richard J Samworth et al. Optimal weighted nearest neighbour classifiers. *The Annals of Statistics*, 40(5):2733–2763, 2012.
- Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: From theory to algorithms*. Cambridge university press, 2014.
- Masashi Sugiyama, Shinichi Nakajima, Hisashi Kashima, Paul V Buenau, and Motoaki Kawanabe. Direct importance estimation with model selection and its application to covariate shift adaptation. In *Advances in neural information processing systems*, pages 1433–1440, 2008.
- Masashi Sugiyama, Taiji Suzuki, and Takafumi Kanamori. *Density ratio estimation in machine learning*. Cambridge University Press, 2012.
- Liu Yang, Steve Hanneke, and Jaime Carbonell. A theory of transfer learning with applications to active learning. *Machine learning*, 90(2):161–189, 2013.