

Importance Weighting for Transfer Learning

Masashi Sugiyama



RIKEN Center for Advanced Intelligence Project/
The University of Tokyo



<http://www.ms.k.u-tokyo.ac.jp/sugi/>



Transfer Learning

Given:

- Training data $\{(\mathbf{x}_i^{\text{tr}}, y_i^{\text{tr}})\}_{i=1}^{n_{\text{tr}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{tr}}(\mathbf{x}, y)$

\mathbf{x} : Input

y : Output

Goal:

- Train a predictor $y = f(\mathbf{x})$
that works well in the test domain
(with some additional data from the test domain).

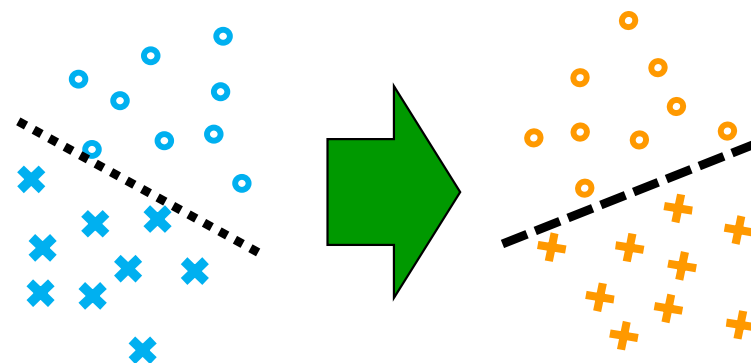
$$\min_f R(f) \quad R(f) = \mathbb{E}_{p_{\text{te}}(\mathbf{x}, y)} [\ell(f(\mathbf{x}), y)]$$

ℓ : loss function

Challenge:

- Overcome changing distributions!

$$p_{\text{tr}}(\mathbf{x}, y) \neq p_{\text{te}}(\mathbf{x}, y)$$





NIPS Workshop 2006 - Whistler

NIPS Workshop on Learning when Test and Training Inputs Have Different Distributions, Whistler 2006

Learning when test and training inputs have different distributions

Joaquin Quiñonero Candela · Masashi Sugiyama · Anton Schwaighofer · Neil D Lawrence

Workshop

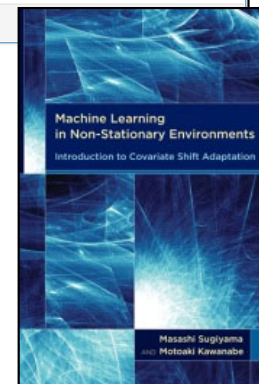
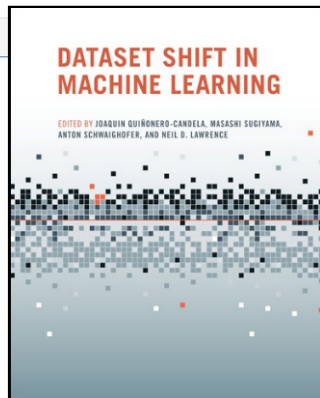
Sat Dec 09 05:00 PM -- 05:00 PM (JST) @ Nordic

Event URL: <http://ida.first.fraunhofer.de/projects/different06/> »

Many machine learning algorithms assume that the training and the test data are drawn from the same distribution. Indeed many of the proofs of statistical consistency, etc., rely on this assumption. However, in practice we are very often faced with the situation where the training and the test data both follow the same conditional distribution, $p(y|x)$, but the input distributions, $p(x)$, differ. For example, principles of experimental design dictate that training data is acquired in a specific manner that bears little resemblance to the way the test inputs may later be generated. The aim of this workshop will be to try and shed light on the kind of situations where explicitly addressing the difference in the input distributions is beneficial, and on what the most sensible ways of doing this are.

Quiñonero-Candela, Sugiyama, Schwaighofer & Lawrence (Eds.), **Dataset Shift in Machine Learning**, MIT Press, 2009.

Sugiyama & Kawanabe, **Machine Learning in Non-Stationary Environments**, MIT Press, 2012



Learning when Training and Test Inputs Have Different Distributions

Saturday December 9, 2006

Org: Joaquin Quiñonero-Candela, Anton Schwaighofer, Neil Lawrence & Masashi Sugiyama

Morning session: 7:30am–10:30am

7:30am **Opening, The organizers**

7:40am **When Training and Test Distributions are Different: Characterising Learning Transfer**, Amos Storkey, *University of Edinburgh*

8:10am **Can Adaptive Regularization Help?**, Matthias Hein, *Max Planck Institute for Biological Cybernetics*

8:40am *coffee break*

8:50am **Learning Classifiers in Distribution and Cost-sensitive Environments**, Nitesh Chawla, *University of Notre Dame*

9:20am **Optimality of Bayesian Transduction - Implications for Input Non-stationarity**, Lars Kai Hansen, *Technical University of Denmark*

9:50pm **Estimating the Joint AUC of Labelled and Unlabelled Data**, Thomas Gärtner, Gemma Garriga, Thorsten Knopp, Peter Flach and Stefan Wrobel

10:10am **A Domain Adaptation Formal Framework Addressing the Training/Test Distribution Gap**, Shai Ben-David, *University of Waterloo* and John Blitzer, *University of Pennsylvania*

Afternoon session: 3:30pm–6:30pm

3:30pm **Projection and Projectability**, David Corfield, *Max Planck Institute for Biological Cybernetics*

4:00pm **Using features of probability distributions to achieve covariate shift**, Arthur Gretton, *MPI for Biol. Cyb. and Alex Smola, National ICT Australia*

4:20pm **Active Learning, Model Selection and Covariate Shift**, Masashi Sugiyama, *Tokyo Institute of Technology*

4:50pm *coffee break*

5:00pm **Visualizing Pairwise Similarity via Semidefinite Programming**, Amir Globerson, *MIT*, and Sam Roweis, *University of Toronto*

5:20pm **A Divergence Prior for Adaptive Learning**, Xiao Li and Jeff Bilmes, *University of Washington*

5:40pm *discussion, everyone*

Various Scenarios

\mathbf{x} : Input

y : Output

■ Full-distribution shift:

$$p_{\text{tr}}(\mathbf{x}, y) \neq p_{\text{te}}(\mathbf{x}, y)$$

■ Covariate shift:

$$p_{\text{tr}}(\mathbf{x}) \neq p_{\text{te}}(\mathbf{x})$$

■ Class-prior shift:

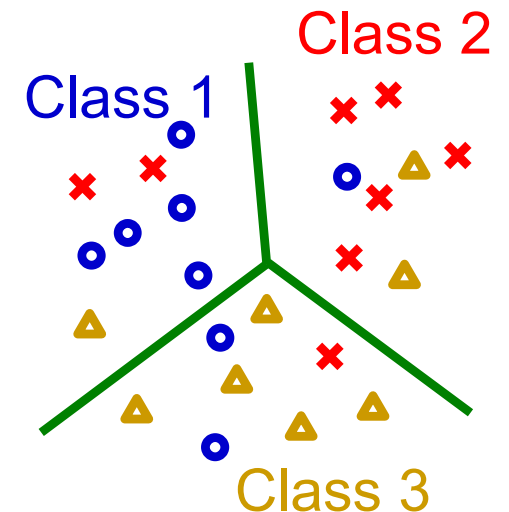
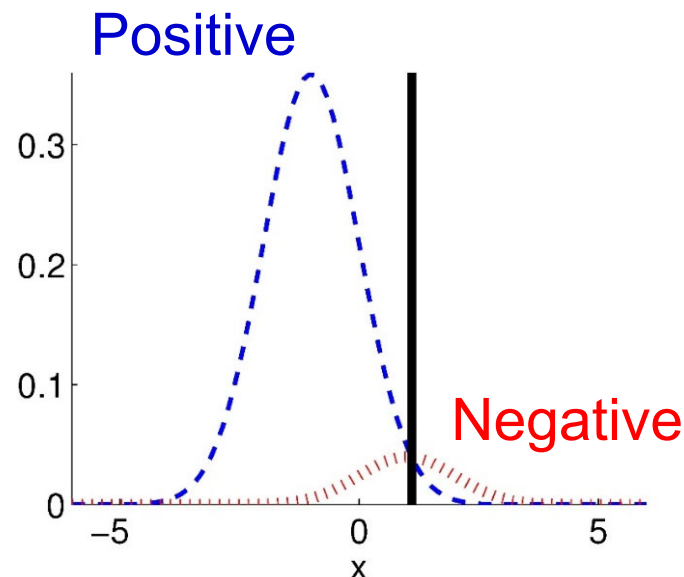
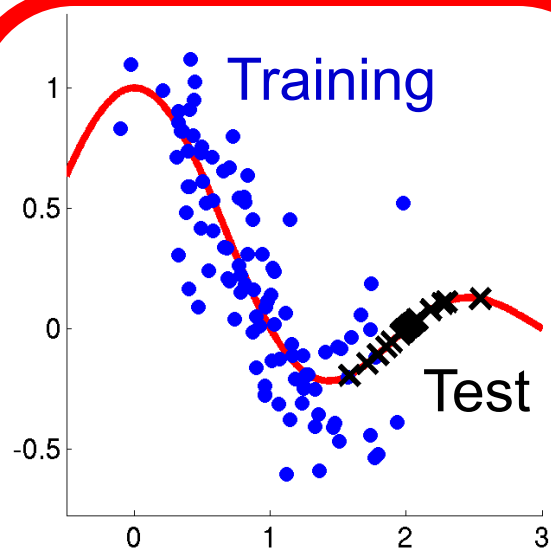
$$p_{\text{tr}}(y) \neq p_{\text{te}}(y)$$

■ Output noise:

$$p_{\text{tr}}(y|\mathbf{x}) \neq p_{\text{te}}(y|\mathbf{x})$$

■ Class-conditional shift:

$$p_{\text{tr}}(\mathbf{x}|y) \neq p_{\text{te}}(\mathbf{x}|y)$$





Contents

1. Introduction
2. **Classical results**
 - A) Importance weighting
 - B) Adaptive importance weighting
3. Recent results
 - A) Joint upper-bound minimization
 - B) Dynamic importance weighting
4. Future prospects

Regression under Covariate Shift

6

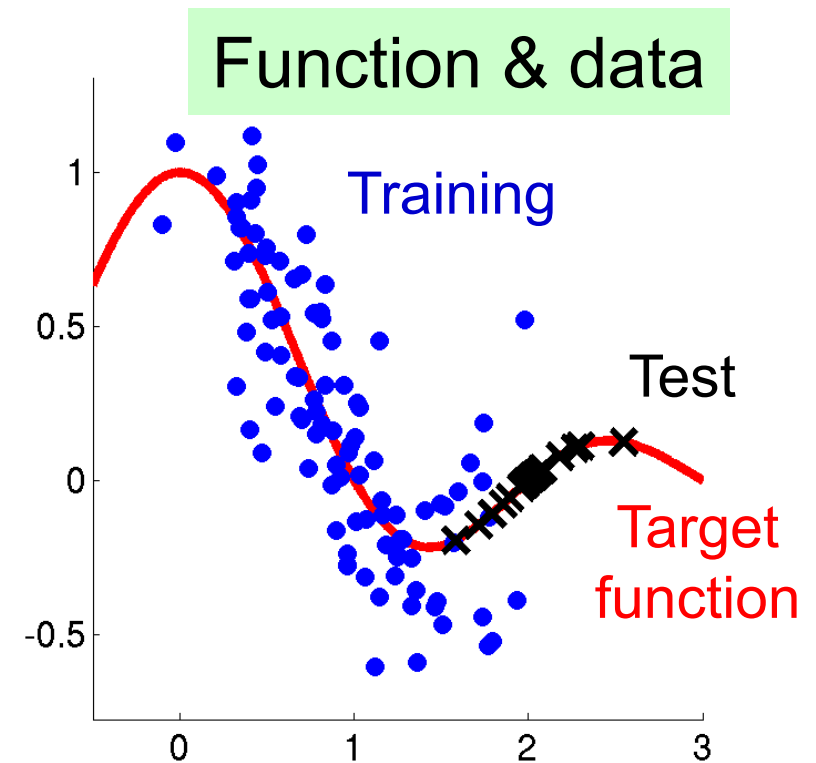
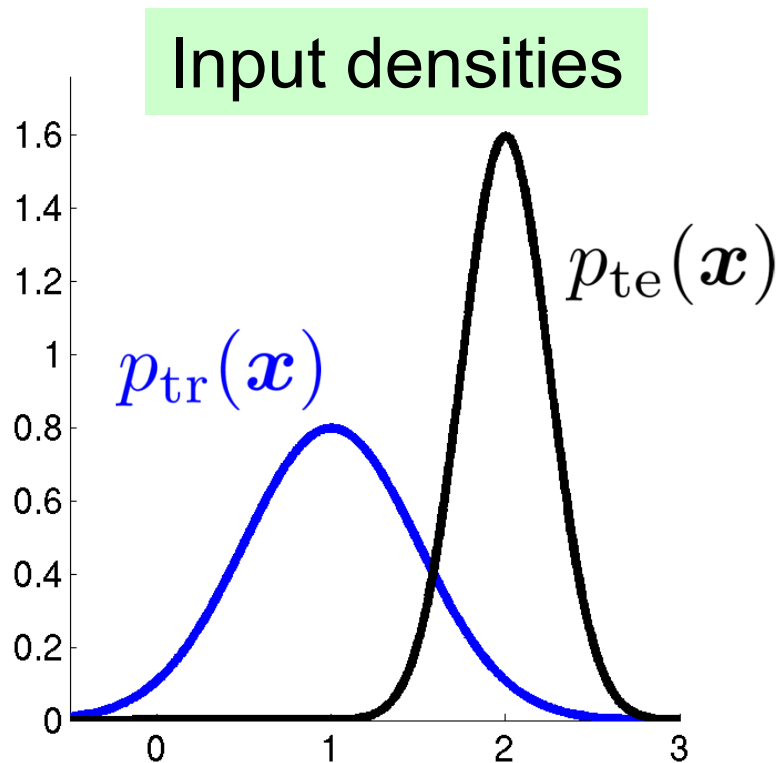
■ Covariate shift: Shimodaira (JSPI2000)

- Training and test input distributions are different:

$$p_{\text{tr}}(\mathbf{x}) \neq p_{\text{te}}(\mathbf{x})$$

- But the output-given-input distribution remains unchanged:

$$p_{\text{tr}}(y|\mathbf{x}) = p_{\text{te}}(y|\mathbf{x}) = p(y|\mathbf{x})$$



Empirical Risk Minimization (ERM)

7

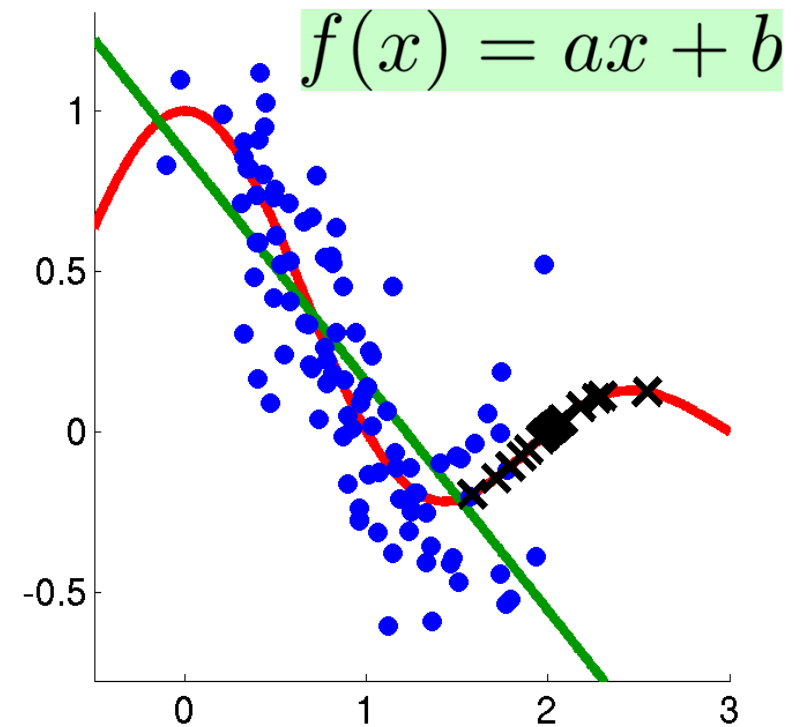
$$\operatorname{argmin}_{f \in \mathcal{F}} \left[\sum_{i=1}^{n_{\text{tr}}} \ell(f(\mathbf{x}_i^{\text{tr}}), y_i^{\text{tr}}) \right]$$

$$\{(\mathbf{x}_i^{\text{tr}}, y_i^{\text{tr}})\}_{i=1}^{n_{\text{tr}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{tr}}(\mathbf{x}, y)$$

■ Generally, ERM is **consistent**:

- Learned function converges to the optimal solution when $n_{\text{tr}} \rightarrow \infty$.

■ However, covariate shift makes ERM **inconsistent**:



$$\operatorname{argmin}_{f \in \mathcal{F}} \left[\frac{1}{n_{\text{tr}}} \sum_{i=1}^{n_{\text{tr}}} \ell(f(\mathbf{x}_i^{\text{tr}}), y_i^{\text{tr}}) \right] \xrightarrow{n_{\text{tr}} \rightarrow \infty} \operatorname{argmin}_{f \in \mathcal{F}} \underbrace{\mathbb{E}_{p(y|\mathbf{x}) p_{\text{tr}}(\mathbf{x})} [\ell(f(\mathbf{x}), y)]}_{\neq R(f)}$$



Contents

1. Introduction
2. Classical results
 - A) Importance weighting
 - B) Adaptive importance weighting
3. Recent results
 - A) Joint upper-bound minimization
 - B) Dynamic importance weighting
4. Future prospects

Importance-Weighted ERM (IWERM) 9

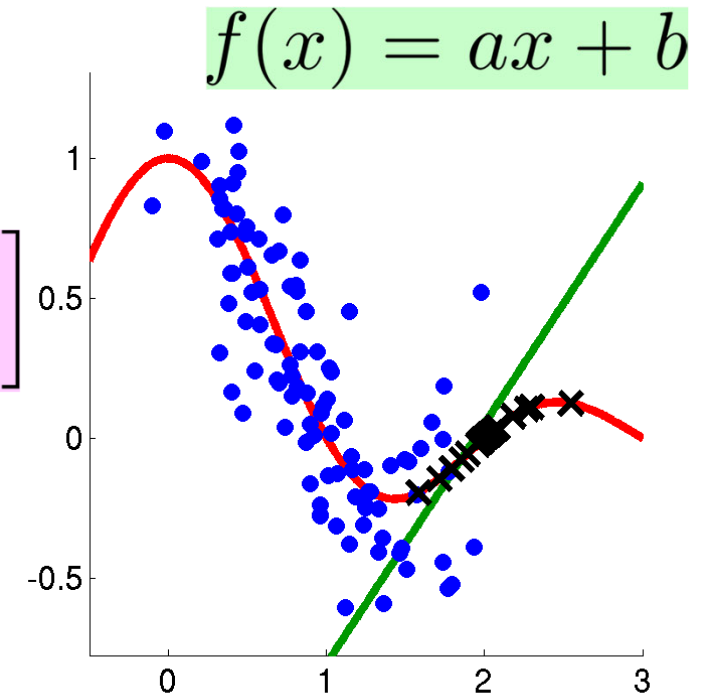
$$\operatorname{argmin}_{f \in \mathcal{F}} \left[\sum_{i=1}^{n_{\text{tr}}} \underbrace{\frac{p_{\text{te}}(\mathbf{x}_i^{\text{tr}})}{p_{\text{tr}}(\mathbf{x}_i^{\text{tr}})}}_{\text{Importance}} \ell(f(\mathbf{x}_i^{\text{tr}}), y_i^{\text{tr}}) \right]$$

- IWERM is consistent even under covariate shift:

$$\operatorname{argmin}_{f \in \mathcal{F}} \left[\frac{1}{n_{\text{tr}}} \sum_{i=1}^{n_{\text{tr}}} \frac{p_{\text{te}}(y_i^{\text{tr}})}{p_{\text{tr}}(y_i^{\text{tr}})} \ell(f(\mathbf{x}_i^{\text{tr}}), y_i^{\text{tr}}) \right]$$

$$\xrightarrow{n_{\text{tr}} \rightarrow \infty} \operatorname{argmin}_{f \in \mathcal{F}} \left[\mathbb{E}_{p(y|\mathbf{x}) p_{\text{tr}}(\mathbf{x})} \left[\frac{p_{\text{te}}(\mathbf{x})}{p_{\text{tr}}(\mathbf{x})} \ell(f(\mathbf{x}), y) \right] \right]$$

$$= \operatorname{argmin}_{f \in \mathcal{F}} \left[\underbrace{\mathbb{E}_{p(y|\mathbf{x}) p_{\text{te}}(\mathbf{x})} [\ell(f(\mathbf{x}), y)]}_{= R(f)} \right]$$



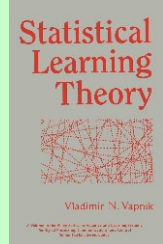
- How can we know the importance weight?



Vapnik's principle:

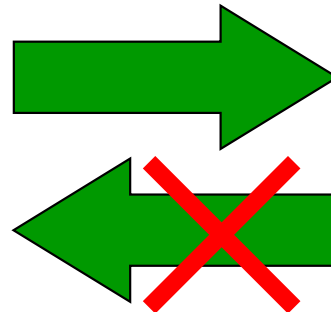
Vapnik (Wiley, 1998)

When solving a problem of interest,
one should not solve a more general problem
as an intermediate step



Knowing densities

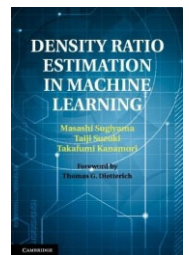
$$p_{\text{te}}(\mathbf{x}), p_{\text{tr}}(\mathbf{x})$$



Knowing ratio

$$r^*(\mathbf{x}) = \frac{p_{\text{te}}(\mathbf{x})}{p_{\text{tr}}(\mathbf{x})}$$

- Estimating the density ratio is substantially easier than estimating both the densities!
- Various **direct density-ratio estimators** were developed.



Sugiyama, Suzuki & Kanamori,
Density Ratio Estimation
in Machine Learning
(Cambridge University Press, 2012)

Least-Squares Importance Fitting (LSIF)

11

Kanamori, Hido & Sugiyama (JMLR2009)

- Given training and test input data:

$$\{\mathbf{x}_i^{\text{tr}}\}_{i=1}^{n_{\text{tr}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{tr}}(\mathbf{x}) \quad \{\mathbf{x}_j^{\text{te}}\}_{j=1}^{n_{\text{te}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{te}}(\mathbf{x})$$

- Directly fit a model r to $r^*(\mathbf{x}) = \frac{p_{\text{te}}(\mathbf{x})}{p_{\text{tr}}(\mathbf{x})}$ by LS:

$$\min_r Q(r) \quad Q(r) = \int \left(r(\mathbf{x}) - r^*(\mathbf{x}) \right)^2 p_{\text{tr}}(\mathbf{x}) d\mathbf{x}$$

- Empirical approximation:

$$Q(r) = \int r(\mathbf{x})^2 p_{\text{tr}}(\mathbf{x}) d\mathbf{x} - 2 \int r(\mathbf{x}) p_{\text{te}}(\mathbf{x}) d\mathbf{x} + C$$

$$\approx \frac{1}{n_{\text{tr}}} \sum_{i=1}^{n_{\text{tr}}} r(\mathbf{x}_i^{\text{tr}})^2 - \frac{2}{n_{\text{te}}} \sum_{j=1}^{n_{\text{te}}} r(\mathbf{x}_j^{\text{te}}) + C$$



Contents

1. Introduction
2. **Classical results**
 - A) Importance weighting
 - B) **Adaptive importance weighting**
3. Recent results
 - A) Joint upper-bound minimization
 - B) Dynamic importance weighting
4. Future prospects

- Importance-weighted empirical risk minimizer

$$\operatorname{argmin}_{f \in \mathcal{F}} \left[\sum_{i=1}^{n_{\text{tr}}} \frac{p_{\text{te}}(\mathbf{x}_i^{\text{tr}})}{p_{\text{tr}}(\mathbf{x}_i^{\text{tr}})} \ell(f(\mathbf{x}_i^{\text{tr}}), y_i^{\text{tr}}) \right]$$

has **no bias**, but has **large variance**.

- The ordinary empirical risk minimizer

$$\operatorname{argmin}_{f \in \mathcal{F}} \left[\sum_{i=1}^{n_{\text{tr}}} \ell(f(\mathbf{x}_i^{\text{tr}}), y_i^{\text{tr}}) \right]$$

has **small variance** (statistically efficient),
but has **large bias**.

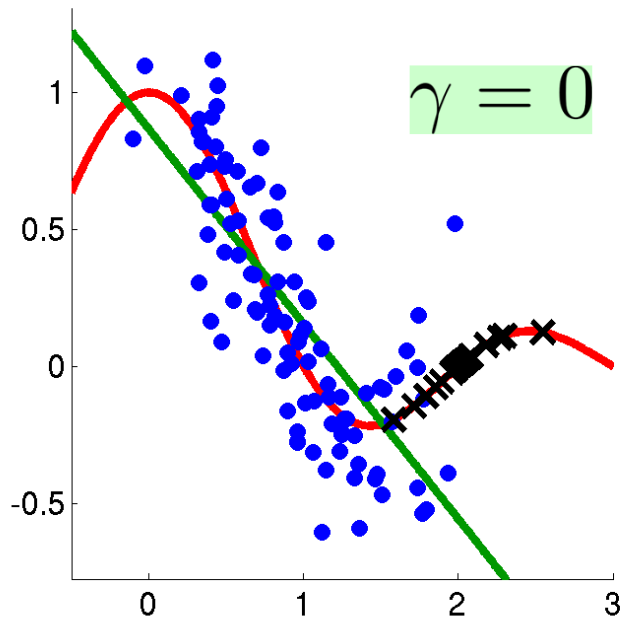
- How can we control the **bias-variance trade-off**?

Flattened Importance Weighting

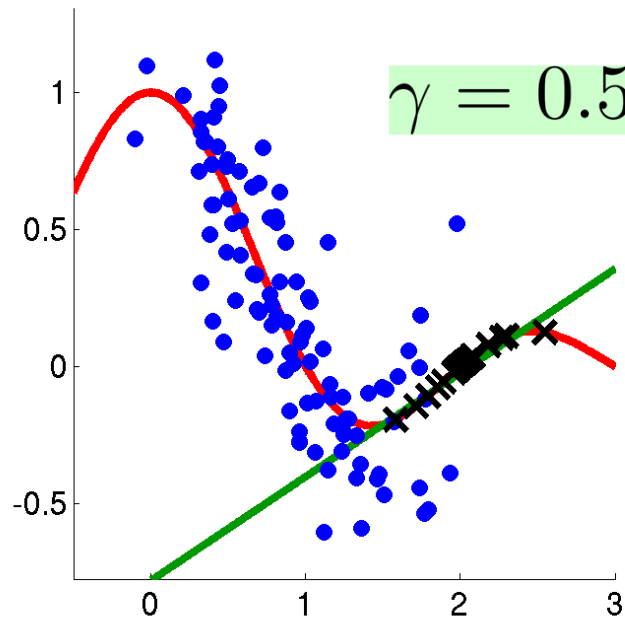
$$\min_f \left[\sum_{i=1}^{n_{\text{tr}}} \left(\frac{p_{\text{te}}(\mathbf{x}_i^{\text{tr}})}{p_{\text{tr}}(\mathbf{x}_i^{\text{tr}})} \right)^\gamma \ell(f(\mathbf{x}_i^{\text{tr}}), y_i^{\text{tr}}) \right]$$

Shimodaira (JSPI2000)

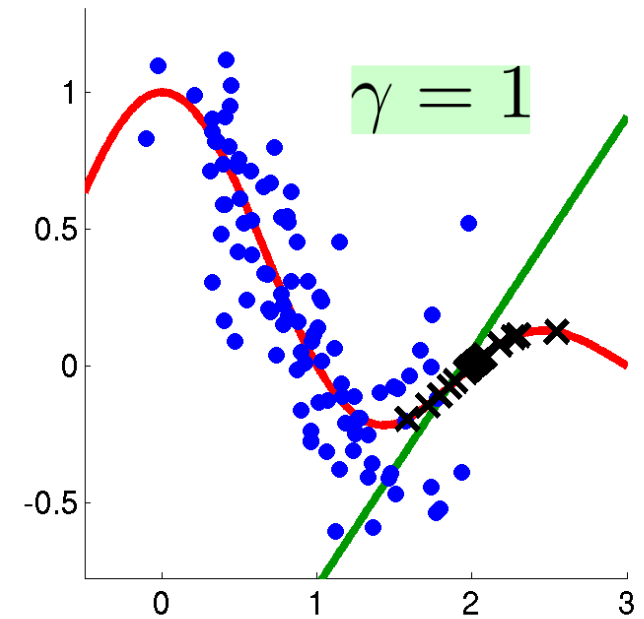
$$0 \leq \gamma \leq 1$$



Large bias, small variance



(Intermediate)



Small bias, large variance

■ **Flattening factor** γ may be chosen by

- Importance-weighted Akaike information criterion
- Importance-weighted cross-validation

Shimodaira
(JSPI2000)

Sugiyama, Krauledat
& Müller (JMLR2007)

- Even with direct methods, reliably estimating the importance weight is hard:

- $r^*(\mathbf{x})$ could be highly fluctuated.

$$r^*(\mathbf{x}) = \frac{p_{\text{te}}(\mathbf{x})}{p_{\text{tr}}(\mathbf{x})}$$

- Then, flattening unreliable importance estimator $\hat{r}(\mathbf{x})$ by power factor γ is also unreliable.

$$\min_f \left[\sum_{i=1}^{n_{\text{tr}}} \hat{r}(\mathbf{x}_i^{\text{tr}})^\gamma \ell(f(\mathbf{x}_i^{\text{tr}}), y_i^{\text{tr}}) \right]$$

- Let's use **relative importance weight**:

Yamada, Suzuki, Kanamori, Hachiya & Sugiyama (NIPS2011, NeCo2013)

$$r_\beta(\mathbf{x}) = \frac{p_{\text{te}}(\mathbf{x})}{\beta p_{\text{tr}}(\mathbf{x}) + (1 - \beta) p_{\text{te}}(\mathbf{x})}$$

$$0 \leq \beta \leq 1$$

- Directly estimable for each β by relative LSIF.



Contents

1. Introduction
2. Classical results
 - A) Importance weighting
 - B) Adaptive importance weighting
3. **Recent results**
 - A) Joint upper-bound minimization
 - B) Dynamic importance weighting
4. Future prospects

From Two-Step Adaptation to One-Step Adaptation

- The classical approaches are **two steps**:

1. Weight estimation (e.g., LSIF):

$$\hat{r} = \operatorname{argmin}_r \mathbb{E}_{p_{\text{tr}}(\mathbf{x})} [(r(\mathbf{x}) - r^*(\mathbf{x}))^2]$$

2. Weighted predictor training (e.g., IWERM):

$$\hat{f} = \operatorname{argmin}_f \mathbb{E}_{p_{\text{tr}}(\mathbf{x}, y)} [\hat{r}(\mathbf{x}) \ell(f(\mathbf{x}), y)]$$

- Can we integrate these two steps?



Contents

1. Introduction
2. Classical results
 - A) Importance weighting
 - B) Adaptive importance weighting
3. **Recent results**
 - A) **Joint upper-bound minimization**
 - B) Dynamic importance weighting
4. Future prospects

Joint Upper-Bound Minimization

19

Zhang et al. (ACML2020, SNCS2021)

■ Suppose we are given

- Labeled training data: $\{(\mathbf{x}_i^{\text{tr}}, y_i^{\text{tr}})\}_{i=1}^{n_{\text{tr}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{tr}}(\mathbf{x}, y)$
- Unlabeled test data: $\{\mathbf{x}_j^{\text{te}}\}_{j=1}^{n_{\text{te}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{te}}(\mathbf{x})$

■ **Goal:** We want to minimize the test risk.

$$R_\ell(f) = \mathbb{E}_{p_{\text{te}}(\mathbf{x}, y)}[\ell(f(\mathbf{x}), y)]$$

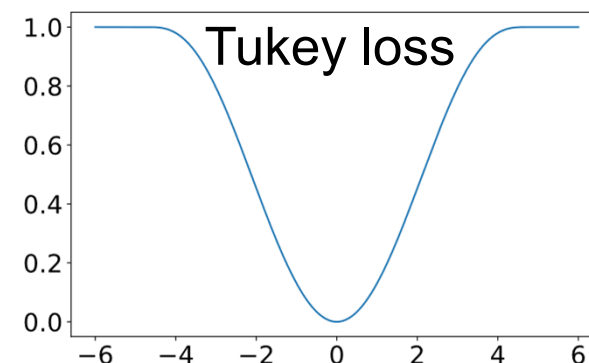
 ℓ : evaluation loss

■ We use **two losses** $\ell(\leq 1)$, $\ell'(\geq \ell)$.

 ℓ' : surrogate loss

For example:

- ℓ : 0/1, ℓ' : hinge or softmax cross-entropy (classification)
- ℓ : Tukey, ℓ' : squared (regression)



- For $\ell \leq 1, \ell' \geq \ell, r \geq 0$,
the test risk is upper-bounded as

$$\frac{1}{2} R_\ell(f)^2 \leq J_{\ell'}(r, f)$$

$$R_\ell(f) = \mathbb{E}_{p_{\text{te}}(\mathbf{x}, y)} [\ell(f(\mathbf{x}), y)]$$

$$J_{\ell'}(r, f) = \left(\mathbb{E}_{p_{\text{tr}}(\mathbf{x}, y)} [r(\mathbf{x}) \ell'(f(\mathbf{x}), y)] \right)^2 \leftarrow \text{IWERM}$$
$$+ \mathbb{E}_{p_{\text{tr}}(\mathbf{x})} [(r(\mathbf{x}) - r^*(\mathbf{x}))^2] \leftarrow \text{LSIF}$$

- In terms of this upper-bound minimization,
2-step (LSIF followed by IWERM) is not optimal:
- Let's directly minimize the upper bound w.r.t. r, f !

- Under some mild conditions, the test risk of the empirical solution $\hat{f} = \operatorname{argmin}_{f \in \mathcal{F}} \min_r \hat{J}_{\ell'}(r, f)$ is upper-bounded as

$$R_{\ell}(\hat{f}) \leq \sqrt{2} \min_{f \in \mathcal{F}} R_{\ell'}(f) + \mathcal{O}_p(n_{\text{tr}}^{-1/4} + n_{\text{te}}^{-1/4})$$

$$\hat{J}_{\ell'}(r, f) = \left(\frac{1}{n_{\text{tr}}} \sum_{i=1}^{n_{\text{tr}}} r(\mathbf{x}_i^{\text{tr}}) \ell'(f(\mathbf{x}_i^{\text{tr}}), y_i^{\text{tr}}) \right)^2 + \left(\frac{1}{n_{\text{tr}}} \sum_{i=1}^{n_{\text{tr}}} r(\mathbf{x}_i^{\text{tr}})^2 - \frac{2}{n_{\text{te}}} \sum_{j=1}^{n_{\text{te}}} r(\mathbf{x}_j^{\text{tr}}) + C \right)$$

$$\{(\mathbf{x}_i^{\text{tr}}, y_i^{\text{tr}})\}_{i=1}^{n_{\text{tr}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{tr}}(\mathbf{x}, y)$$

$$\{\mathbf{x}_j^{\text{te}}\}_{j=1}^{n_{\text{te}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{te}}(\mathbf{x})$$

$$R_{\ell}(\hat{f}) = \mathbb{E}_{p_{\text{te}}(\mathbf{x}, y)}[\ell(\hat{f}(\mathbf{x}), y)]$$

$$R_{\ell'}(f) = \mathbb{E}_{p_{\text{te}}(\mathbf{x}, y)}[\ell'(f(\mathbf{x}), y)]$$



Contents

1. Introduction
2. Classical results
 - A) Importance weighting
 - B) Adaptive importance weighting
3. **Recent results**
 - A) Joint upper-bound minimization
 - B) **Dynamic importance weighting**
4. Future prospects

Dynamic Importance Weighting

23

Fang et al. (NeurIPS2020)

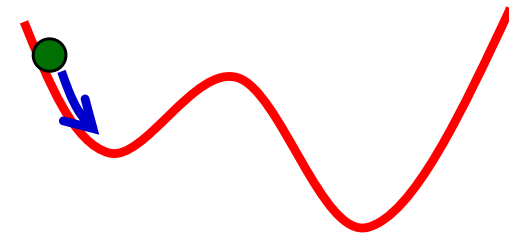
- Deep learning adopts **stochastic optimization**:

$$f \leftarrow f - \eta \nabla \hat{R}(f) \quad \eta > 0: \text{Learning rate}$$

- Let's learn

- Importance weight r
- predictor f

dynamically in the **mini-batch-wise** manner.



■ Suppose we are given

- (Large) labeled training data: $\{(\mathbf{x}_i^{\text{tr}}, y_i^{\text{tr}})\}_{i=1}^{n_{\text{tr}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{tr}}(\mathbf{x}, y)$
- (Small) labeled test data: $\{(\mathbf{x}_j^{\text{te}}, y_j^{\text{te}})\}_{j=1}^{n_{\text{te}}} \stackrel{\text{i.i.d.}}{\sim} p_{\text{te}}(\mathbf{x}, y)$

■ For **each mini-batch** $\{(\bar{\mathbf{x}}_i^{\text{tr}}, \bar{y}_i^{\text{tr}})\}_{i=1}^{\bar{n}_{\text{tr}}}, \{(\bar{\mathbf{x}}_j^{\text{te}}, \bar{y}_j^{\text{te}})\}_{j=1}^{\bar{n}_{\text{te}}}$
importance weights are estimated by
kernel mean matching for loss values:

Huang, et al. (NeurIPS2007)

$$\frac{1}{\bar{n}_{\text{tr}}} \sum_{i=1}^{\bar{n}_{\text{tr}}} r_i \ell(f(\bar{\mathbf{x}}_i^{\text{tr}}), \bar{y}_i^{\text{tr}}) \approx \frac{1}{\bar{n}_{\text{te}}} \sum_{j=1}^{\bar{n}_{\text{te}}} \ell(f(\bar{\mathbf{x}}_j^{\text{te}}), \bar{y}_j^{\text{te}})$$

■ **No covariate shift assumption is needed!**



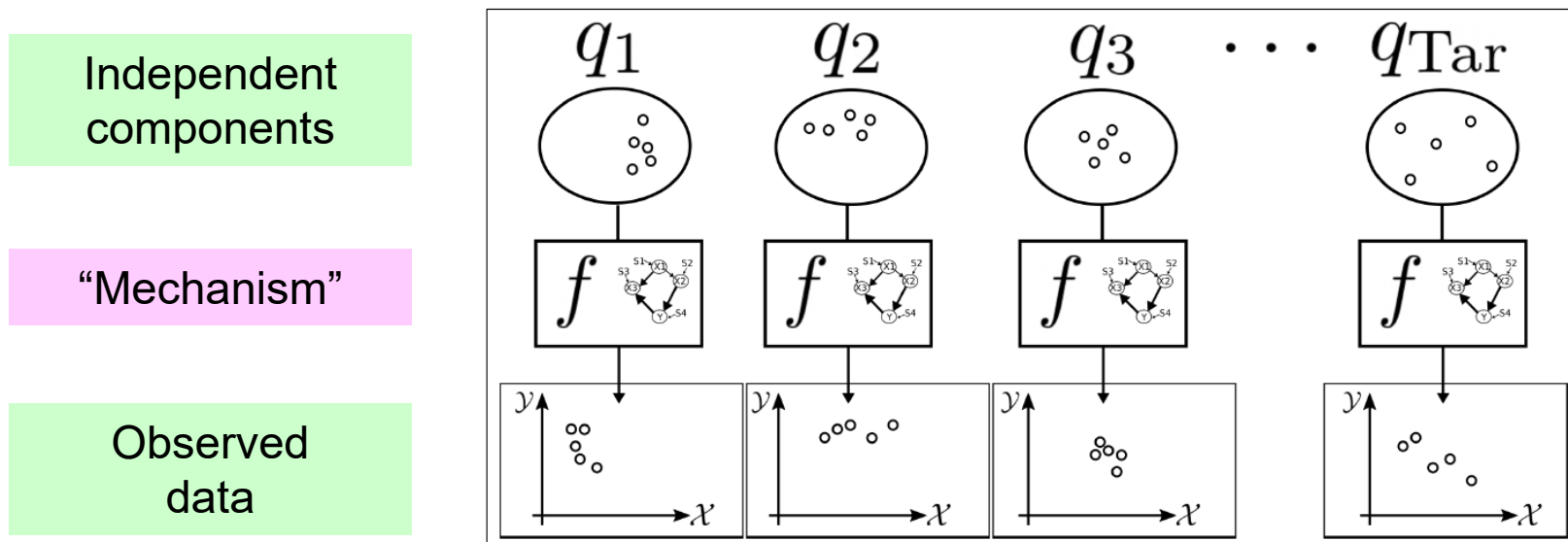
Contents

1. Introduction
2. Classical results
 - A) Importance weighting
 - B) Adaptive importance weighting
3. Recent results
 - A) Joint upper-bound minimization
 - B) Dynamic importance weighting
4. **Future prospects**

Conclusions

- In transfer learning with importance weighting, simultaneously performing **importance estimation** and **predictor training** is promising.
- What should we do if training and test distributions look very different?
 - **Mechanism transfer!**

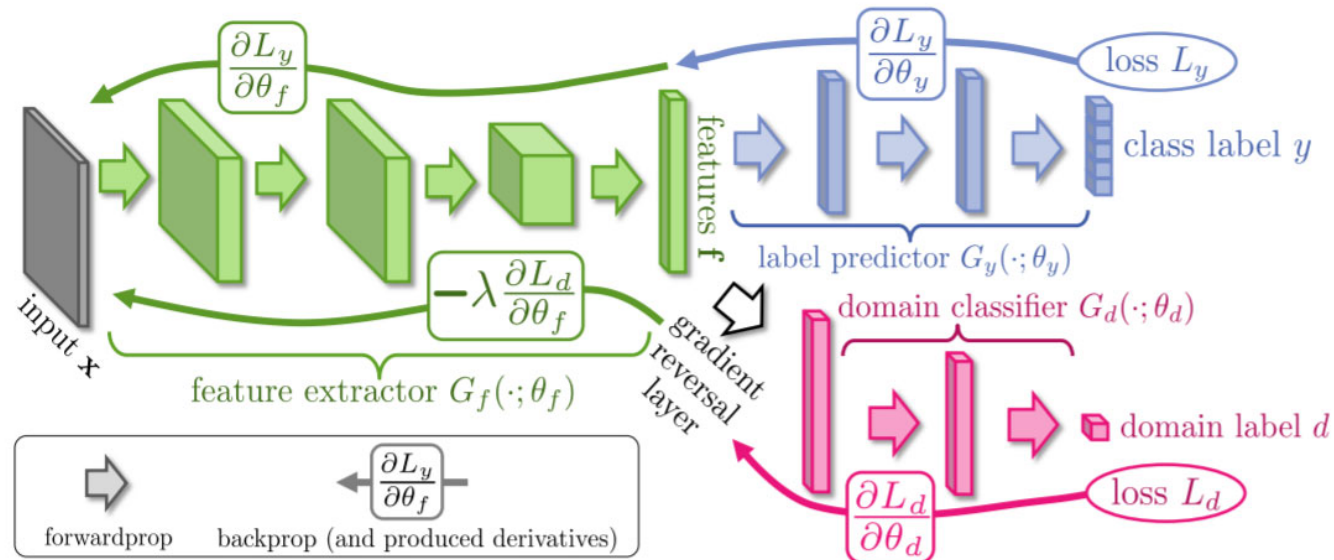
Teshima, Sato & Sugiyama (ICML2020)



Future Prospects: Domain Matching

- **Domain matching** would be another popular approach for transfer learning in deep learning:

Ben-David, Blitzer, Crammer & Pereira (NIPS2006)
Ganin & Lempitsky (ICML2015)



- Can we **combine** domain matching and importance weighting for better performance?

Future Prospects: Classification with Noisy Labels

■ **Output shift:** $p_{\text{tr}}(y|\mathbf{x}) \neq p_{\text{te}}(y|\mathbf{x})$

\mathbf{x} : Input pattern

y : Class label

● **Noise transition** connects two distributions:

$$p_{\text{tr}}(\bar{y}|\mathbf{x}) = \sum_y p(\bar{y}|y)p_{\text{te}}(y|\mathbf{x})$$

\bar{y} : Noisy class label

■ **Back/forward loss correction** yields consistency.

Patrini, Rozza, Menon, Nock & Qu (CVPR2017)

■ Estimation of noise transition **only from noisy training data** is the current challenge.

Xia et al. (NeurIPS2019), Yao et al. (NeurIPS2020), Xia et al. (NeurIPS2020),
Zhang et al. (ICML2021), Li et al. (ICML2021), Berthon et al. (ICML2021)

■ **Can we use transfer learning techniques to better solve noisy label classification?**