

Direct Density Ratio Estimation for Large-scale Covariate Shift Adaptation

Yuta Tsuboi

IBM Research, Tokyo Research Laboratory
and

Graduate School of Information Science,
Nara Institute of Science and Technology

Hisashi Kashima

IBM Research, Tokyo Research Laboratory

Shohei Hido IBM Research, Tokyo Research Laboratory

Steffen Bickel

Department of Computer Science, University of Potsdam, Germany

Masashi Sugiyama

Department of Computer Science, Tokyo Institute of Technology

Abstract

Covariate shift is a situation in supervised learning where training and test *inputs* follow different distributions even though the functional relation remains unchanged. A common approach to compensating for the bias caused by covariate shift is to reweight the loss function according to the *importance*, which is the ratio of test and training densities. We propose a novel method that allows us to directly estimate the importance from samples without going through the hard task of density estimation. An advantage of the proposed method is that the computation time is nearly independent of the number of test input samples, which is highly beneficial in recent applications with large numbers of unlabeled samples. We demonstrate through experiments that the proposed method is computationally more efficient than existing approaches with comparable accuracy. We also describe a promising result for large-scale covariate shift adaptation in a natural language processing task.

1 Introduction

An assumption that is commonly imposed—either explicitly or implicitly—in virtually all supervised learning methods is that the training and test samples follow the *same* probability distribution. However, this fundamental assumption is often violated in practice,

causing standard machine learning methods not to work as expected. In this paper, we address supervised learning problems in the absence of this fundamental assumption.

If the training and test distributions share nothing in common, we may not be able to learn anything about the test distribution from the training samples. For a meaningful discussion, the training and test distributions should be related to each other in some sense. A situation where the input distribution $p(\mathbf{x})$ is different in the training and test phases but the conditional distribution of output values, $p(y|\mathbf{x})$, remains unchanged is called *covariate shift* [21]. In many real-world applications such as robot control [27, 20, 15], bioinformatics [1, 6], spam filtering [4], natural language processing [17], brain-computer interfacing [31, 24], or econometrics [14], covariate shift is likely. Covariate shift is also naturally induced in selective sampling or active learning scenarios [10, 8, 30, 18, 23]. For this reason, learning under covariate shift is receiving a lot of attention these days in the machine learning community (such as in the NIPS2006 workshop [7] and the ECML2006 workshop [2]).

Under covariate shift, standard learning methods such as maximum likelihood estimation are no longer *consistent*, i.e., they do not produce the optimal solution even when the number of training samples tends to be infinity. Thus, there exists an estimation bias induced by covariate shift. It has been shown that the bias can be asymptotically canceled by weighting the log likelihood terms according to the *importance* [11, 21, 32]:

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

where $p_{\text{te}}(\mathbf{x})$ and $p_{\text{tr}}(\mathbf{x})$ are the test and training input densities. Since the importance is usually unknown in reality, the central issue of practical covariate shift adaptation is how to accurately estimate the importance¹.

A naive approach to importance estimation is to first estimate the training and test densities separately from the training and test input samples, and then estimate the importance by taking the ratio of the estimated densities. However, density estimation is known to be a hard problem particularly in high dimensional cases [12]. Therefore, this naive approach is usually ineffective—directly estimating the importance *without* estimating the densities is more promising. Therefore, several methods that allow us to directly obtain importance estimates without going through density estimation have been proposed recently, such as *kernel mean matching* (KMM) [16], the *logistic regression*

¹Covariate shift matters in parameter learning only when the model used for function learning is *misspecified* (i.e., the model is so simple that the true learning target function cannot be expressed) [21]. When the model is correctly (or overly) specified, the ordinary maximum likelihood estimation is still consistent. On this basis, there is a criticism that importance weighting is not needed, but just the use of a sufficiently complex model can settle the problem. However, overly complex models result in large estimation variances, and so in practice we need to choose a complex enough but not overly complex model. To choose such an appropriate model, we usually use a model selection technique such as cross-validation (CV). However, the ordinary CV score is biased due to covariate shift and we still need to importance-weight the CV score (or any other model selection criteria) for unbiasedness [21, 32, 25, 24]. For this reason, estimating the importance is indispensable when covariate shift occurs.

based method (LogReg) [3], and the *Kullback-Leibler Importance Estimation Procedure* (KLIEP) [26].

KMM is based on a special property of *universal reproducing kernel Hilbert spaces* (Gaussian reproducing kernel Hilbert spaces are typical examples) [22] and KMM allows us to directly obtain the importance estimates at the training input points. Since the KMM optimization problem is formulated as a convex quadratic programming problem, it always leads to the unique global solution. KMM has been shown to work well, as long as the kernel parameters such as the Gaussian width are chosen appropriately. However, to the best of our knowledge, there is no reliable method to determine the Gaussian width and the regularization parameter in the KMM algorithm². Therefore, the lack of model selection procedures is a critical limitation of KMM in practical applications.

LogReg builds a probabilistic classifier that separates training input samples from test input samples, and the importance can be directly estimated by LogReg. The maximum likelihood estimation of the LogReg can be formulated as a convex optimization problem, so the unique global optimal solution can be obtained. In addition, since LogReg only solves a standard supervised classification problem, the tuning parameters such as the kernel width and the regularization parameter can be optimized by the standard cross-validation (CV) procedure. This is a very useful property in practice.

KLIEP tries to match an importance-based estimation of the test input distribution to the true test input distribution in terms of the Kullback-Leibler divergence. KLIEP solves this matching problem in a non-parametric fashion. The training and test input distributions are not parameterized, but only the importance is parameterized. The KLIEP optimization problem is convex and therefore a unique global optimal solution can be obtained. Furthermore, the global solution tends to be sparse, so it is computationally efficient in the test phase. Since KLIEP is based on the minimization of the Kullback-Leibler divergence, the model selection of KLIEP, such as the choice of the kernel width and the regularization parameter, can be carried out naturally through the *likelihood CV* procedure [12], so no open tuning parameter remains.

As reviewed above, LogReg and KLIEP seem to have advantages over KMM, since they are equipped with built-in model selection procedures. On the other hand, from the viewpoint of scalability, all three of the methods have limitations—in recent applications such as spam filtering [4] and information retrieval [13], the number of test (unlabeled) samples are enormous. In these text processing applications, the distribution of training and test inputs can be changed between domains because of differences in vocabulary and writing style, and especially the large amount of unlabeled test samples is available on the Web. The purpose of this paper is to develop a computationally efficient covariate

²Intuitively, it seems possible to optimize the kernel width and the regularization parameter simply by using CV for the performance of subsequent learning algorithms. However, this is highly unreliable since the ordinary CV score is biased under covariate shift. For unbiased estimation of the prediction performance of subsequent learning algorithms, the CV procedure itself needs to be importance-weighted [32, 24]. Since the importance weight has to have been fixed when model selection is carried out using the importance weighted CV, it cannot be used for model selection of importance estimation algorithms. Note that once the importance weight has been fixed, the importance-weighted CV can be used for model selection of subsequent learning algorithms.

shift adaptation method that can deal with a large number of unlabeled data.

Our new method is primarily based on KLIEP. The key difference is that the original KLIEP uses a linearly parameterized function for modeling the importance, while we adopt a *log-linear* model. By definition, the log-linear model only takes non-negative values. This allows us to reformulate the KLIEP optimization problem as an *unconstrained* convex problem. Then we develop a new scalable estimation procedure whose computation time is nearly independent of the number of test samples. More precisely, we need to scan a large number of test samples only once to compute a summary statistic in the beginning (this precomputation can be carried out in linear time and constant storage space). The main optimization procedure does not use the test samples themselves, but only uses the summary statistic. Therefore, the computation time of the main optimization procedure is independent of the number of test samples.

The experiments show that the proposed method is computationally much more efficient than the existing approaches. Therefore the range of application of covariate shift adaptation can be greatly enlarged towards large-scale problems. As regards estimation accuracy, we experimentally show that the performance of the proposed method is comparable to the best existing methods for small and middle sized problems (since the existing methods cannot be applied to large-scale problems due to the computational costs). Thus the proposed method can be a useful alternative to the existing covariate shift adaptation methods.

The rest of this paper is organized as follows. Section 2 formulates the supervised learning problem under covariate shift and review covariate shift adaptation techniques. Section 3 reviews an existing direct importance estimation method, KLIEP. Section 4 proposes a new direct importance estimation method that is scalable to large test data sets. Section 5 illustrates how the proposed method works in covariate shift adaptation using simple regression and classification data sets. Section 6 discusses the relation of the proposed method to existing approaches. Section 7 reports experimental results comparing the computation time and estimation accuracy of the proposed and existing methods. Finally, Section 8 gives conclusions.

2 Problem Formulation

In this section, we formulate the supervised learning problem under covariate shift and briefly review existing techniques for covariate shift adaptation.

2.1 Supervised learning under covariate shift

Let $\mathbf{x} \in \mathbf{X} \subset \mathbb{R}^d$ be an input variable and $y \in Y$ be an output variable. Y is a real space in regression cases or a set of categories in classification cases. In standard supervised learning frameworks, it is assumed that \mathbf{x} is independently drawn from an input distribution and y is independently drawn from a conditional distribution both in training and test phases. In contrast, here we consider a situation called *covariate shift*

[21], i.e., the input distribution differs in the training and test phases, but the conditional distribution remains unchanged.

Suppose we have independent and identically distributed (i.i.d.) training input samples $D_{\text{tr}} = \{\mathbf{x}^{(i)}\}_{i=1}^{N_{\text{tr}}}$ from a distribution with strictly positive density $p_{\text{tr}}(\mathbf{x})$, and test input samples $D_{\text{te}} = \{\mathbf{x}^{(i)}\}_{i=1}^{N_{\text{te}}}$ from a distribution with density $p_{\text{te}}(\mathbf{x})$. In addition to the input samples, suppose we have training output samples $\{y^{(i)}\}_{i=1}^{N_{\text{tr}}}$ drawn from the conditional distribution with conditional density $p(y|\mathbf{x} = \mathbf{x}^{(i)})$, respectively. Typically, the number N_{tr} of training samples is rather small due to the high labeling cost, while the number N_{te} of test input samples is very large since they are often easily available. We denote training sample pairs of input and output as:

$$Z_{\text{tr}} = \{z^{(i)} \mid z^{(i)} = (\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^{N_{\text{tr}}}.$$

We use the following linear model:

$$f_{\boldsymbol{\theta}}(\mathbf{x}) = \langle \boldsymbol{\theta}, \boldsymbol{\phi}(\mathbf{x}) \rangle, \quad (1)$$

where $\boldsymbol{\theta}$ is the parameter vector, $\boldsymbol{\phi}(\mathbf{x}) : \mathbf{X} \rightarrow \mathfrak{R}^h$ is a basis function of \mathbf{x} , and $\langle \mathbf{u}, \mathbf{v} \rangle$ denotes the Euclidean inner product between vector \mathbf{u} and \mathbf{v} : $\langle \mathbf{u}, \mathbf{v} \rangle = \sum_{l=1}^h u_l v_l$. Note that this model can contain a bias parameter by just including a constant basis function in $\boldsymbol{\phi}(\mathbf{x})$. Throughout the paper, we suppose that this linear model is not generally specified correctly, i.e., the true input-output function is not necessarily included in the above linear model. Since we do not know the true function class in practice, dealing with misspecified models is quite realistic.

The goal of supervised learning is to learn the parameter $\boldsymbol{\theta}$ so that the output values for the test inputs can be accurately predicted. Thus our error metric (which is usually called the *generalization error*) is given by

$$\iint \text{Loss}(\mathbf{x}, y, f_{\boldsymbol{\theta}}(\mathbf{x})) p_{\text{te}}(\mathbf{x}) p(y|\mathbf{x}) d\mathbf{x} dy, \quad (2)$$

where $\text{Loss}(\mathbf{x}, y, f_{\boldsymbol{\theta}}(\mathbf{x})) : \mathbf{X} \times Y \times Y \rightarrow \mathfrak{R}$ is a loss function, such as the squared loss in a regression case or the zero-one loss in a classification case.

In supervised learning under covariate shift, the following quantity called the *test domain importance* plays an important role:

$$w(\mathbf{x}) = \frac{p_{\text{te}}(\mathbf{x})}{p_{\text{tr}}(\mathbf{x})}. \quad (3)$$

The importance can be used for adjusting the difference between the training and test input distributions: for any function $A(\mathbf{x})$,

$$\int A(\mathbf{x}) p_{\text{te}}(\mathbf{x}) d\mathbf{x} = \int A(\mathbf{x}) w(\mathbf{x}) p_{\text{tr}}(\mathbf{x}) d\mathbf{x}. \quad (4)$$

2.2 Parameter learning under covariate shift

Here we review two typical parameter learning methods under covariate shift: one is *importance weighted least squares* (IWLS) for regression and the other is *importance weighted logistic regression* (IWLR) for classification.

2.2.1 IWLS

A standard learning method in regression scenarios would be ordinary least squares (LS):

$$\hat{\boldsymbol{\theta}}_{\text{LS}} \equiv \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \left[\sum_{(\mathbf{x}, y) \in Z_{\text{tr}}} (f_{\boldsymbol{\theta}}(\mathbf{x}) - y)^2 \right].$$

LS is known to be consistent under a usual setting. However, it is no longer consistent for misspecified models under covariate shift. Instead, IWLS is consistent [21]:

$$\hat{\boldsymbol{\theta}}_{\text{IWLS}} \equiv \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \left[\sum_{(\mathbf{x}, y) \in Z_{\text{tr}}} w(\mathbf{x}) (f_{\boldsymbol{\theta}}(\mathbf{x}) - y)^2 + \lambda \|\boldsymbol{\theta}\|^2 \right], \quad (5)$$

where the importance $w(\mathbf{x})$ is used as weights. Here we also added a penalty term $\lambda \|\boldsymbol{\theta}\|^2$ for regularization, where λ is a regularization parameter.

For the linear model (1), the above optimization problem is convex and the unique global solution $\hat{\boldsymbol{\theta}}_{\text{IWLS}}$ can be computed in a closed-form as

$$\hat{\boldsymbol{\theta}}_{\text{IWLS}} = (\boldsymbol{\Phi}^\top \mathbf{W} \boldsymbol{\Phi} + \lambda \mathbf{I})^{-1} \boldsymbol{\Phi}^\top \mathbf{W} \mathbf{y},$$

where \mathbf{I} is the identity matrix,

$$\begin{aligned} \Phi_{i,l} &= \phi_l(\mathbf{x}^{(i)}), \\ \mathbf{y} &= (y^{(1)}, y^{(2)}, \dots, y^{(N_{\text{tr}})})^\top, \\ \mathbf{W} &= \operatorname{diag}(\mathbf{w}^{(1)}, \mathbf{w}^{(2)}, \dots, \mathbf{w}^{(N_{\text{tr}})}). \end{aligned}$$

2.2.2 IWLR

For simplicity, we focus on the two-class case, i.e., $Y = \{-1, 1\}$; we note that it is straightforward to extend all of the discussions in this paper to multi-class cases.

Let us model the posterior probability of class y given \mathbf{x} using a parametric model $f_{\boldsymbol{\theta}}(\mathbf{x})$ as

$$p_{\boldsymbol{\theta}}(y|\mathbf{x}) = \frac{\exp(y f_{\boldsymbol{\theta}}(\mathbf{x}))}{1 + \exp(y f_{\boldsymbol{\theta}}(\mathbf{x}))}. \quad (6)$$

Then a test input sample \mathbf{x} is classified by choosing the most probable class:

$$\hat{y} = \underset{y}{\operatorname{argmax}} p_{\boldsymbol{\theta}}(y|\mathbf{x}). \quad (7)$$

A standard learning method for the above probabilistic classification scenarios would be ordinary logistic regression (LR):

$$\begin{aligned}\hat{\boldsymbol{\theta}}_{\text{LR}} &\equiv \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \left[\sum_{(\mathbf{x}, y) \in Z_{\text{tr}}} -\log p_{\boldsymbol{\theta}}(y|\mathbf{x}) \right] \\ &= \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \left[\sum_{(\mathbf{x}, y) \in Z_{\text{tr}}} (\log(1 + \exp(yf_{\boldsymbol{\theta}}(\mathbf{x}))) - yf_{\boldsymbol{\theta}}(\mathbf{x})) \right].\end{aligned}$$

Similar to the case of LS, LR is consistent under a usual setting, but is no longer consistent for misspecified models under covariate shift. Instead, IWLR is consistent:

$$\hat{\boldsymbol{\theta}}_{\text{IWLR}} \equiv \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \left[\sum_{(\mathbf{x}, y) \in Z_{\text{tr}}} w(\mathbf{x}) (\log(1 + \exp(yf_{\boldsymbol{\theta}}(\mathbf{x}))) - yf_{\boldsymbol{\theta}}(\mathbf{x})) + \lambda \|\boldsymbol{\theta}\|^2 \right]. \quad (8)$$

Here we also added a penalty term $\lambda \|\boldsymbol{\theta}\|^2$ for regularization, where λ is a regularization parameter.

This optimization problem is known to be convex and a unique optimal solution can be computed using standard non-linear optimization techniques such as a gradient descent method or some variants of the Newton method. The gradient of the above objective function is given by

$$\sum_{(\mathbf{x}, y) \in Z_{\text{tr}}} w(\mathbf{x}) (yp_{\boldsymbol{\theta}}(y|\mathbf{x})\boldsymbol{\phi}(\mathbf{x}) - y\boldsymbol{\phi}(\mathbf{x})) + 2\lambda\boldsymbol{\theta}.$$

2.3 Model selection under covariate shift

In the above learning methods, the choice of model parameters such as the basis functions $\boldsymbol{\phi}$ and the regularization parameter λ heavily affects the prediction performance. This problem is called model selection and is one of the key concerns in machine learning.

A popular model selection method in the machine learning community would be cross-validation (CV). The performance of CV is guaranteed in the sense that it gives an unbiased estimate of the generalization error. However, this useful theoretical property is no longer true under covariate shift [32]. To cope with this problem, a variant of CV called *importance weighted CV* (IWCV) has been proposed for model selection under covariate shift [24]. It has been proved that IWCV gives an unbiased estimate of the generalization error even under covariate shift.

Here, let us briefly describe the IWCV procedure. We first divide the training samples $\{z^{(i)}\}_{i=1}^{N_{\text{tr}}}$ into R disjoint subsets $\{Z^r\}_{r=1}^R$. Then we learn a function $f_{\boldsymbol{\theta}}^r(\mathbf{x})$ from $\{Z^j\}_{j \neq r}$

by IWLS/IWLR and compute its mean test error for the remaining samples Z^r :

$$\begin{aligned} \frac{1}{|Z^r|} \sum_{(\mathbf{x}, y) \in Z^r} w(\mathbf{x}) (f_{\boldsymbol{\theta}}^r(\mathbf{x}) - y)^2, \quad (\text{regression}) \\ \frac{1}{|Z^r|} \sum_{(\mathbf{x}, y) \in Z^r} w(\mathbf{x}) I(\hat{y} = y), \quad (\text{classification}) \end{aligned}$$

where $I(\cdot)$ denotes the indicator function. We repeat this procedure for $r = 1, 2, \dots, R$ and choose the model such that the average of the above mean test error is minimized.

3 Importance Estimation

As we have seen in the previous section, the importance $w(\mathbf{x})$ plays a central role in covariate shift adaptation. However, the importance is unknown in practice so we need to estimate it from samples.

Direct importance estimation methods that do not involve density estimation steps have been developed recently [16, 3, 26]. Here we review one of those direct methods called the *Kullback-Leibler Importance Estimation Procedure* (KLIEP) [26]. Other methods will be reviewed in Section 6.

3.1 KLIEP

Let us model $w(\mathbf{x})$ with the following linear model:

$$\hat{w}(\mathbf{x}) = \langle \boldsymbol{\alpha}, \boldsymbol{\psi}(\mathbf{x}) \rangle, \quad (9)$$

where $\boldsymbol{\alpha} \in \mathfrak{R}^b$ is a model parameter vector and $\boldsymbol{\psi}(\mathbf{x}) \in \mathfrak{R}^b$ is a basis function. Since the importance should be non-negative by definition, we suppose that both $\boldsymbol{\alpha}$ and $\boldsymbol{\psi}(\mathbf{x})$ are non-negative.

Using the importance estimation $\hat{w}(\mathbf{x})$, we can estimate the test input density $p_{\text{te}}(\mathbf{x})$ by

$$\hat{p}_{\text{te}}(\mathbf{x}) = p_{\text{tr}}(\mathbf{x}) \hat{w}(\mathbf{x}). \quad (10)$$

Now we learn the parameter $\boldsymbol{\alpha}$ so that the Kullback-Leibler divergence from $p_{\text{te}}(\mathbf{x})$ to $\hat{p}_{\text{te}}(\mathbf{x})$ is minimized:

$$\begin{aligned} \text{KL}[p_{\text{te}}(\mathbf{x}) || \hat{p}_{\text{te}}(\mathbf{x})] &= \int_D p_{\text{te}}(\mathbf{x}) \log \frac{p_{\text{te}}(\mathbf{x})}{p_{\text{tr}}(\mathbf{x}) \hat{w}(\mathbf{x})} d\mathbf{x} \\ &= \int_D p_{\text{te}}(\mathbf{x}) \log \frac{p_{\text{te}}(\mathbf{x})}{p_{\text{tr}}(\mathbf{x})} d\mathbf{x} - \int_D p_{\text{te}}(\mathbf{x}) \log \hat{w}(\mathbf{x}) d\mathbf{x}. \end{aligned} \quad (11)$$

Since the first term in Eq.(11) is independent of $\boldsymbol{\alpha}$, we ignore it and focus on the second term, which we denote by J_{KLIEP} :

$$J_{\text{KLIEP}} = \int_D p_{\text{te}}(\mathbf{x}) \log \hat{w}(\mathbf{x}) d\mathbf{x} \approx \frac{1}{N_{\text{te}}} \sum_{\mathbf{x} \in D_{\text{te}}} \log \hat{w}(\mathbf{x}), \quad (12)$$

where an empirical approximation based on the test input samples is used. This is the objective function to be maximized. The value of $\hat{w}(\mathbf{x})$ should be properly normalized since it is a probability density function:

$$1 = \int_D \hat{p}_{te}(\mathbf{x}) d\mathbf{x} = \int_D p_{tr}(\mathbf{x}) \hat{w}(\mathbf{x}) d\mathbf{x} \approx \frac{1}{N_{tr}} \sum_{\mathbf{x} \in D_{tr}} \hat{w}(\mathbf{x}), \quad (13)$$

where the empirical approximation based on the training samples is used.

Then the resulting optimization problem is expressed as

$$\begin{aligned} & \underset{\boldsymbol{\alpha}}{\text{maximize}} && \sum_{\mathbf{x} \in D_{te}} \log \langle \boldsymbol{\alpha}, \boldsymbol{\psi}(\mathbf{x}) \rangle \\ & \text{subject to} && \sum_{\mathbf{x} \in D_{tr}} \langle \boldsymbol{\alpha}, \boldsymbol{\psi}(\mathbf{x}) \rangle = N_{tr} \text{ and } \boldsymbol{\alpha} \geq \mathbf{0}, \end{aligned}$$

which is convex. Thus the global solution can be obtained by iteratively performing gradient ascent and feasibility satisfaction.

3.2 Model selection by likelihood CV

The performance of KLIEP depends on the choice of the basis functions $\boldsymbol{\psi}(\mathbf{x})$ (and possibly an additional regularization parameter). Since KLIEP is based on the maximization of the score J_{KLIEP} , it would be natural to select the model such that J_{KLIEP} is maximized. The expectation over $p_{te}(\mathbf{x})$ involved in J_{KLIEP} can be numerically approximated by *likelihood CV* (LCV) [12] as follows: First, divide the test samples D_{te} into R disjoint subsets $\{D_{te}^r\}_{r=1}^R$. Then, obtain an importance estimate $\hat{w}^r(\mathbf{x})$ from $\{D_{te}^t\}_{t \neq r}^R$ and approximate the score J_{KLIEP} using D_{te}^r as

$$\hat{J}_{\text{KLIEP}}^r = \frac{1}{|D_{te}^r|} \sum_{\mathbf{x} \in D_{te}^r} \hat{w}^r(\mathbf{x}). \quad (14)$$

This procedure is repeated for $r = 1, 2, \dots, R$ and choose the model such that the average of \hat{J}_{KLIEP}^r for all r is maximized.

One of the potential general limitations of CV is that it is not reliable in small sample cases, since data splitting by CV further reduces the sample size. A key advantage of the LCV procedure is that, not the training samples, but the test input samples are cross-validated. This contributes greatly to improving the model selection accuracy, since the number of training samples is typically limited while there are lots of test input samples available.

As basis functions, it is suggested to use Gaussian kernels centered at a subset of the test input points D_{te} [26]:

$$K_s(\mathbf{x}, \mathbf{x}_l) = \exp \left\{ -\frac{\|\mathbf{x} - \mathbf{x}_l\|^2}{2s^2} \right\}, \quad (15)$$

where $\mathbf{x}_l \in D_{te}$ is a template test sample and s is the kernel width. This is a heuristic to allocate many kernels at high test input density regions since many kernels may be needed in the region where the output of the target function is large. In the original paper, the number of Gaussian centers was fixed at $N_{te}/10$ for computational efficiency and the kernel width s was chosen by LCV.

4 KLIEP for Log-linear Models

As shown above, KLIEP has its own model selection procedure and has been shown to work well in importance estimation [26]. However, it has a weakness in computation time. In each step of gradient ascent, the summation over all test input samples needs to be computed, which is prohibitively slow in large-scale problems. The main contribution of this paper is to extend KLIEP so that it can deal with large sets of test input data.

4.1 LL-KLIEP

In the original KLIEP, a linearly parameterized model (9) is used for modeling the importance function. Here, we propose using a (normalized) log-linear model for modeling the importance $w(\mathbf{x})$ as

$$\hat{w}(\mathbf{x}) = \frac{\exp(\langle \boldsymbol{\alpha}, \boldsymbol{\psi}(\mathbf{x}) \rangle)}{\frac{1}{N_{tr}} \sum_{\mathbf{x}' \in D_{tr}} \exp(\langle \boldsymbol{\alpha}, \boldsymbol{\psi}(\mathbf{x}') \rangle)}, \quad (16)$$

where the denominator guarantees the normalization constraint (13)³. By definition, the log-linear model takes only non-negative values. Therefore, we no longer need the non-negative constraint for the parameter (and the basis functions).

Then the optimization problem becomes *unconstrained*:

$$\underset{\boldsymbol{\alpha}}{\text{maximize}} J_{LL-KLIEP}(\boldsymbol{\alpha}),$$

where

$$\begin{aligned} J_{LL-KLIEP}(\boldsymbol{\alpha}) &= \frac{1}{N_{te}} \sum_{\mathbf{x} \in D_{te}} \log \hat{w}(\mathbf{x}) \\ &= \frac{1}{N_{te}} \sum_{\mathbf{x} \in D_{te}} \langle \boldsymbol{\alpha}, \boldsymbol{\psi}(\mathbf{x}) \rangle - \log \frac{1}{N_{tr}} \sum_{\mathbf{x} \in D_{tr}} \exp(\langle \boldsymbol{\alpha}, \boldsymbol{\psi}(\mathbf{x}) \rangle). \end{aligned} \quad (17)$$

³The log-linear model can have numerical problems since it contains an exponential function. To cope with this problem, we do not directly compute the value of $\hat{w}(\mathbf{x})$, but we compute it in the exponential of the logarithmic domain, i.e.,

$$\exp(\log \hat{w}(\mathbf{x})) = \exp(\langle \boldsymbol{\alpha}, \boldsymbol{\psi}(\mathbf{x}) \rangle - \log \frac{1}{N_{tr}} \sum_{\mathbf{x}' \in D_{tr}} \exp(\langle \boldsymbol{\alpha}, \boldsymbol{\psi}(\mathbf{x}') \rangle)).$$

To further stabilize the computation, we compute the logarithmic sum of the exponential functions as $\log(\exp(a) + \exp(b)) = \log(1 + \exp(b - a))$, where we pick the smaller exponent as b .

Below, we refer to this method as *LL-KLIEP* (log-linear KLIEP). In practice, we may add a penalty term for regularization:

$$j(\boldsymbol{\alpha}) = J_{\text{LL-KLIEP}}(\boldsymbol{\alpha}) - \frac{\|\boldsymbol{\alpha}\|^2}{2\sigma^2}, \quad (18)$$

where σ^2 is a regularization parameter.

An advantage of LL-KLIEP over the original KLIEP is its computational efficiency. The gradient of $j(\boldsymbol{\alpha})$ can be computed as

$$\begin{aligned} \frac{\partial j(\boldsymbol{\alpha})}{\partial \boldsymbol{\alpha}} &= \frac{1}{N_{\text{te}}} \sum_{\mathbf{x} \in D_{\text{te}}} \boldsymbol{\psi}(\mathbf{x}) - \sum_{\mathbf{x} \in D_{\text{tr}}} \frac{\exp(\langle \boldsymbol{\alpha}, \boldsymbol{\psi}(\mathbf{x}) \rangle)}{\sum_{\mathbf{x}' \in D_{\text{tr}}} \exp(\langle \boldsymbol{\alpha}, \boldsymbol{\psi}(\mathbf{x}') \rangle)} \boldsymbol{\psi}(\mathbf{x}) - \frac{\boldsymbol{\alpha}}{\sigma^2} \\ &= F - \frac{1}{N_{\text{tr}}} \sum_{\mathbf{x} \in D_{\text{tr}}} \hat{w}(\mathbf{x}) \boldsymbol{\psi}(\mathbf{x}) - \frac{\boldsymbol{\alpha}}{\sigma^2}, \end{aligned} \quad (19)$$

where

$$F = \frac{1}{N_{\text{te}}} \sum_{\mathbf{x} \in D_{\text{te}}} \boldsymbol{\psi}(\mathbf{x}).$$

This means that once we pre-compute the value of F , we do not need to use the test samples when we compute the gradient. This contributes greatly to reducing the computation time when the number of test samples is large. In addition, we do not need to store all of the test samples in memory since we only need the value of F . The required storage capacity is only $\Omega(cN_{\text{tr}})$, where c is the average number of non-zero basis entries.

As the model selection of KLIEP, LCV can be used to find the optimal hyper-parameters of LL-KLIEP. Since $J_{\text{LL-KLIEP}}(\boldsymbol{\alpha})$ is evaluated using both training and test samples, both test and training samples can be divided into R disjoint subset $\{D_{\text{te}}^r\}_{r=1}^R$ and $\{D_{\text{tr}}^r\}_{r=1}^R$ in the LCV procedure. After the estimation of $\hat{w}^r(\mathbf{x})$, $J_{\text{LL-KLIEP}}(\boldsymbol{\alpha})$ can be approximated as

$$\hat{J}_{\text{LL-KLIEP}}^r(\boldsymbol{\alpha}) = \frac{1}{|D_{\text{te}}^r|} \sum_{\mathbf{x} \in D_{\text{te}}^r} \langle \boldsymbol{\alpha}, \boldsymbol{\psi}(\mathbf{x}) \rangle - \log \frac{1}{|D_{\text{tr}}^r|} \sum_{\mathbf{x} \in D_{\text{tr}}^r} \exp(\langle \boldsymbol{\alpha}, \boldsymbol{\psi}(\mathbf{x}) \rangle).$$

Although the proposed optimization procedure may be more efficient than original KLIEP, there still exists a potential weakness: we still need to use all the test samples when computing the values of $J_{\text{LL-KLIEP}}(\boldsymbol{\alpha})$ or $j(\boldsymbol{\alpha})$. The value of $J_{\text{LL-KLIEP}}(\boldsymbol{\alpha})$ is needed when we choose a model by LCV, and the value of $j(\boldsymbol{\alpha})$ is often utilized in line search or in the stopping criterion.

4.2 LL-KLIEP(LS)

Here, we introduce another optimization technique for LL-KLIEP that enables us to overcome the above weakness. Our basic idea is to encourage the derivative of the convex

objective function to be zero. We use a squared norm to measure the ‘magnitude’ of the derivative (19):

$$J_{\text{LS}}(\boldsymbol{\alpha}) = \frac{1}{2} \left\| \frac{\partial j(\boldsymbol{\alpha})}{\partial \boldsymbol{\alpha}} \right\|^2. \quad (20)$$

The partial derivative of Eq.(20) with respect to $\boldsymbol{\alpha}$ is expressed as

$$\frac{\partial J_{\text{LS}}(\boldsymbol{\alpha})}{\partial \boldsymbol{\alpha}} = \frac{\partial^2 j(\boldsymbol{\alpha})}{\partial^2 \boldsymbol{\alpha}} \frac{\partial j(\boldsymbol{\alpha})}{\partial \boldsymbol{\alpha}}. \quad (21)$$

Note that the first component of Eq. (21) is the Hessian matrix of $j(\boldsymbol{\alpha})$:

$$\frac{\partial^2 j(\boldsymbol{\alpha})}{\partial^2 \boldsymbol{\alpha}} = \left(\sum_{\mathbf{x} \in D_{\text{tr}}} \frac{1}{N_{\text{tr}}} w(\mathbf{x}) (\bar{\boldsymbol{\psi}}(\mathbf{x}) - \boldsymbol{\psi}(\mathbf{x})) \boldsymbol{\psi}(\mathbf{x})^T - \frac{I}{\sigma^2} \right),$$

where

$$\bar{\boldsymbol{\psi}}(\mathbf{x}) = \sum_{\mathbf{x} \in D_{\text{tr}}} \frac{1}{N_{\text{tr}}} w(\mathbf{x}) \boldsymbol{\psi}(\mathbf{x}),$$

$\boldsymbol{\psi}(\mathbf{x})^T$ is the transpose of $\boldsymbol{\psi}(\mathbf{x})$, and I is the identity matrix. When we explicitly compute the Hessian matrix of $j(\boldsymbol{\alpha})$, the computational complexity of the derivative is $O(b^2 N_{\text{tr}})$, which is independent of N_{te} . Also, the required storage space is independent of N_{te} : $\Omega(b^2 + cN_{\text{tr}})$. We refer to this approach as *LL-KLIEP(LS1-a)* below.

The computation time and storage space of LL-KLIEP(LS1-a) are quadratic functions of the number of parameters b , which could be a bottleneck in high dimensional problems. To cope with this problem, we propose two approaches.

One approach for high dimensional problems is directly computing the product between the Hessian matrix and the gradient vector of $j(\boldsymbol{\alpha})$ without storing the Hessian matrix:

$$\frac{\partial J_{\text{LS}}(\boldsymbol{\alpha})}{\partial \boldsymbol{\alpha}} = \left(\sum_{\mathbf{x} \in D_{\text{tr}}} \frac{1}{N_{\text{tr}}} w(\mathbf{x}) (\bar{\boldsymbol{\psi}}(\mathbf{x}) - \boldsymbol{\psi}(\mathbf{x})) \langle \boldsymbol{\psi}(\mathbf{x}), G \rangle - \frac{G}{\sigma^2} \right), \quad (22)$$

where $G = \frac{\partial j(\boldsymbol{\alpha})}{\partial \boldsymbol{\alpha}}$. Since the inner product $\langle \boldsymbol{\psi}(\mathbf{x}), G \rangle$ requires $O(b)$ time, we can compute Eq.(22) in total with $O(bN_{\text{tr}})$ computation time and $\Omega(cN_{\text{tr}})$ space. We call this approach as *LL-KLIEP(LS1-b)*, which is still independent of N_{te} and suitable for high dimensional problems comparing with LL-KLIEP(LS1-a).

In the other approach for high dimensional problems, we make use of the *representer theorem* [29]. Our idea is to represent the parameter $\boldsymbol{\alpha}$ as a linear combination of the input samples:

$$\boldsymbol{\alpha} = \sum_{\mathbf{x} \in D_{\text{tr}}} \boldsymbol{\psi}(\mathbf{x}) \beta_{\mathbf{x}},$$

where $\{\beta_{\mathbf{x}}\}_{\mathbf{x} \in D_{\text{tr}}}$ is a data-wise parameter. Then Eq.(20) can be rewritten as

$$J_{\text{LS}}(\{\beta_{\mathbf{x}}\}_{\mathbf{x} \in D_{\text{tr}}}) = \frac{1}{2} \left\| F - \sum_{\mathbf{x} \in D_{\text{tr}}} \boldsymbol{\psi}(\mathbf{x}) \omega(\mathbf{x}) - \sum_{\mathbf{x} \in D_{\text{tr}}} \frac{\boldsymbol{\psi}(\mathbf{x}) \beta_{\mathbf{x}}}{\sigma^2} \right\|^2, \quad (23)$$

where

$$\omega(\mathbf{x}) = \frac{\exp(\sum_{\mathbf{x}' \in D_{\text{tr}}} K(\mathbf{x}, \mathbf{x}') \beta_{\mathbf{x}'})}{\sum_{\mathbf{x}'' \in D_{\text{tr}}} \exp(\sum_{\mathbf{x}' \in D_{\text{tr}}} K(\mathbf{x}'', \mathbf{x}') \beta_{\mathbf{x}'})}, \quad (24)$$

$$K(\mathbf{x}, \mathbf{x}') = \langle \boldsymbol{\psi}(\mathbf{x}), \boldsymbol{\psi}(\mathbf{x}') \rangle.$$

The partial derivative of Eq.(23) with respect to $\beta_{\mathbf{x}}$ is:

$$\frac{\partial j_{\text{LS}}(\{\beta_{\mathbf{x}}\}_{\mathbf{x} \in D_{\text{tr}}})}{\partial \beta_{\mathbf{x}}} = \left\langle F - \sum_{\mathbf{x}' \in D_{\text{tr}}} \boldsymbol{\psi}(\mathbf{x}') \left(\omega(\mathbf{x}') - \frac{\beta_{\mathbf{x}'}}{\sigma^2} \right), \sum_{\mathbf{x}' \in D_{\text{tr}}} \omega(\mathbf{x}') \boldsymbol{\psi}(\mathbf{x}') \langle \boldsymbol{\varphi}(\mathbf{x}'), \boldsymbol{\psi}(\mathbf{x}) \rangle - \frac{\boldsymbol{\psi}(\mathbf{x})}{\sigma^2} \right\rangle, \quad (25)$$

where $\boldsymbol{\varphi}(\mathbf{x}) = \sum_{\mathbf{x}' \in D_{\text{tr}}} \omega(\mathbf{x}') \boldsymbol{\psi}(\mathbf{x}') - \boldsymbol{\psi}(\mathbf{x})$. By the change of variables, it is not required to calculate the partial derivative with respect to $\boldsymbol{\alpha}$ so that we can avoid the computation of the Hessian matrix of $j(\boldsymbol{\alpha})$. We refer to this approach as *LL-KLIEP(LS2)*.

The computation of LL-KLIEP(LS2) requires $O(bN_{\text{tr}}^2)$ time and $\Omega(N_{\text{tr}}^2 + cN_{\text{tr}})$ space. The computation time is linear with respect to the number of parameters b and the storage space is independent of b . This is also a improvement over the direct computation of the partial derivative in Eq.(21).

For LL-KLIEP(LS), LCV can also be computed very efficiently. In each validation set using D_{te}^r and D_{tr}^r , we can compute the validation error as

$$\hat{J}_{\text{LL-KLIEP(LS)}}^r = \left\| F^r - \sum_{\mathbf{x} \in D_{\text{tr}}^r} \hat{w}^r(\mathbf{x}) \boldsymbol{\psi}(\mathbf{x}) \right\|^2,$$

where

$$F^r = \frac{1}{|D_{\text{te}}^r|} \sum_{\mathbf{x} \in D_{\text{te}}^r} \boldsymbol{\psi}(\mathbf{x}).$$

Note that, once the mean basis vectors F^r are calculated for all R disjoint subsets of D_{te} , $\hat{J}_{\text{LL-KLIEP(LS)}}^r$ can be evaluated independently of the size of the test data D_{te}^r .

The computational complexity and storage space of each method are summarized in Table 1. In term of the complexity analysis, LL-KLIEP(LS1-b) is the best solution for the large amount of test inputs. We verified the analysis by the computational experiments in Section 7.1.

5 Illustrative Examples

In this section, we illustrate the behavior of the proposed LL-KLIEP and show how it can be applied in covariate shift adaptation.

Table 1: Computational complexity and space requirements. N_{tr} is the number of training samples, N_{te} is the number of test samples, b is the number of parameters, and c is the average number of non-zero basis entries. ‘‘Precomp.’’ denotes the computational complexity of once-off precomputation.

	Computational complexity			Space requirement	
	Precomp.	Objective	Derivative	Objective	Derivative
KLIEP	0	$bN_{\text{tr}} + bN_{\text{te}}$	$bN_{\text{tr}} + bN_{\text{te}}$	$cN_{\text{tr}} + cN_{\text{te}}$	$cN_{\text{tr}} + cN_{\text{te}}$
LL-KLIEP	bN_{te}	$bN_{\text{tr}} + bN_{\text{te}}$	bN_{tr}	$cN_{\text{tr}} + cN_{\text{te}}$	cN_{tr}
LL-KLIEP(LS1-a)	bN_{te}	bN_{tr}	b^2N_{tr}	cN_{tr}	$b^2 + cN_{\text{tr}}$
LL-KLIEP(LS1-b)	bN_{te}	bN_{tr}	bN_{tr}	cN_{tr}	cN_{tr}
LL-KLIEP(LS2)	bN_{te}	bN_{tr}^2	bN_{tr}^2	cN_{tr}	$N_{\text{tr}}^2 + cN_{\text{tr}}$

5.1 Regression under covariate shift

Let us consider an illustrative regression problem of learning

$$f(x) = \text{sinc}(x).$$

Let the training and test input densities be $p_{\text{tr}}(x) = \mathcal{N}(x; 1, 1^2)$ and $p_{\text{te}}(x) = \mathcal{N}(x; 2, 0.5^2)$, where $\mathcal{N}(x; \mu, \sigma^2)$ denotes the Gaussian density with mean μ and variance σ^2 . We create the training output value $\{y^{(i)}\}_{i=1}^{N_{\text{tr}}}$ as $y^{(i)} = f(x^{(i)}) + \epsilon^{(i)}$, where the noise $\{\epsilon^{(i)}\}_{i=1}^{N_{\text{tr}}}$ has density $\mathcal{N}(\epsilon; 0, 0.25^2)$. Let the number of training samples be $N_{\text{tr}} = 200$ and the number of test samples be $N_{\text{te}} = 1000$. These settings imply that we are considering an extrapolation problem (see Figure 1(a)).

We used 100 Gaussian basis functions centered at randomly chosen test input samples. Figure 1(b) shows the actual importance $w(x)$ and an estimated importance $\hat{w}(x)$ by using LL-KLIEP, where the hyper-parameters such as the Gaussian width and the regularization parameter are selected by LCV. We also tested LL-KLIEP(LS1-a), LL-KLIEP(LS1-b), and LL-KLIEP(LS2), but we omit their graphs since their solutions are almost identical to the solution of LL-KLIEP.

Figure 2 depicts the values of the true $J_{\text{LL-KLIEP}}$ (see Eq.(17)) and its estimate by 5-fold LCV. The means, the 25 percentiles, and the 75 percentiles over each validation are plotted as functions of the kernel width s for the different $\sigma = 0.5, 1, 2$. We also plot the normalized mean squared error of the estimated importance:

$$\text{NMSE} = \frac{1}{N_{\text{tr}}} \sum_{\mathbf{x} \in D_{\text{tr}}} \left(\frac{\hat{w}(\mathbf{x})}{\sum_{\mathbf{x}' \in D_{\text{tr}}} \hat{w}(\mathbf{x}')} - \frac{w(\mathbf{x})}{\sum_{\mathbf{x}' \in D_{\text{tr}}} w(\mathbf{x}')} \right)^2. \quad (26)$$

The graph shows that LCV gives a very good estimate of $J_{\text{LL-KLIEP}}$ and also NMSE. Figure 2 also shows that σ value affects the importance estimation in terms of NMSE.

Figure 3 shows the true learning target function and functions learned by ordinary LS and IWLS with a linear basis function (Figure 3(a)), i.e., $\phi(x) = (1, x)^\top$, and a quadratic

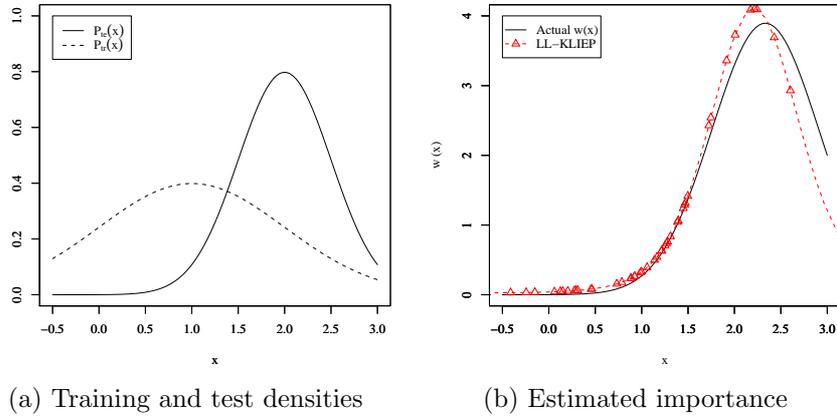


Figure 1: Importance estimation.

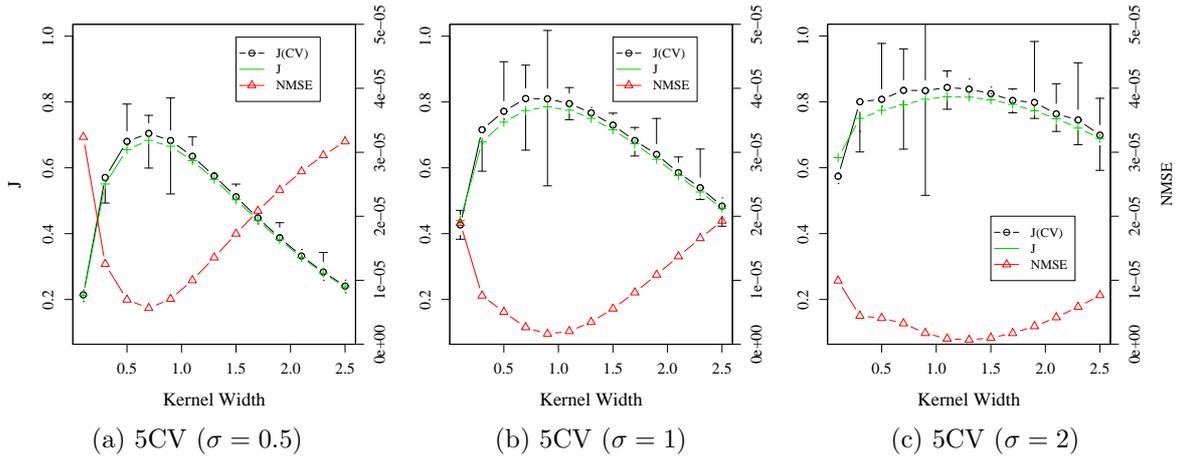


Figure 2: Model selection curve.

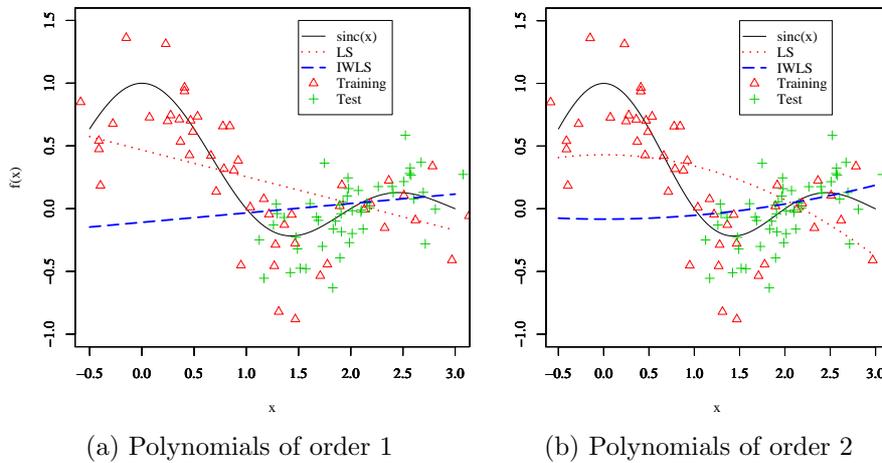


Figure 3: Regression under covariate shift (True and learned functions).

basis function (Figure 3(b)), i.e., $\phi(x) = (1, x, x^2)^\top$ (Section 2.2). The regularization parameter λ was selected by CV for LS and IWCV for IWLS (Sections 2.3). The results show that the learned function using IWLS goes reasonably well through the test samples, while that of ordinary LS overfits the training samples. Note that the output of the test samples are not used to obtain the learned functions.

5.2 Classification under covariate shift

Next, let us consider two illustrative binary classification problems, where two-dimensional samples were generated from Gaussian distributions (see Table 2 and Figure 4). These data sets correspond to a ‘linear shift’ and a ‘non-linear shift’ (rotation).

Let the number of the training samples be $N_{\text{tr}} = 200$ and that of the test samples be $N_{\text{te}} = 1000$ (only 500 test samples are plotted for clarity). We used LR/IWLR for the training classifiers (see Section 2.2), and employed CV/IWCV for the regularization parameter tuning (see Section 2.3). We used a linear basis function for LR/IWLR: $\phi(\mathbf{x}) = (1, \mathbf{x}^\top)^\top$.

Figure 4 shows the decision boundaries obtained by LR+CV and IWLR+IWCV. For references, we also show ‘OPT’, which is the optimal decision boundary obtained using the test input-output samples. For the data set depicted in Figure 4(a), the correct classification rate of LR+CV is 99.1% while that of IWLR+IWCV is 100%. For the data set depicted in Figure 4(b), the correct classification rate of LR+CV is 97.2% while that of IWLR+IWCV is 99.1%. Thus, for both cases, the prediction performance is improved by importance weighting.

6 Related Work and Discussion

In this section, we compare the proposed LL-KLIEP with existing importance estimation approaches.

6.1 Kernel density estimator

The kernel density estimator (KDE) is a non-parametric technique to estimate a density $p(\mathbf{x})$ from samples $\{\mathbf{x}_l\}_{l=1}^N$. For the Gaussian kernel, KDE is expressed as

$$\hat{p}(\mathbf{x}) = \frac{1}{(2\pi s^2)^{d/2} N} \sum_{l=1}^N K_s(\mathbf{x}, \mathbf{x}_l), \quad (27)$$

where $K_s(\mathbf{x}, \mathbf{x}')$ is the Gaussian kernel (15). The performance of KDE depends on the choice of the kernel width s , which can be optimized by LCV [12]. Note that LCV corresponds to choosing s such that the Kullback-Leibler divergence from $p(\mathbf{x})$ to $\hat{p}(\mathbf{x})$ is minimized.

KDE can be used for importance estimation by first obtaining $\hat{p}_{\text{tr}}(\mathbf{x})$ and $\hat{p}_{\text{te}}(\mathbf{x})$ separately from $\{\mathbf{x}^{(i)}\}_{i=1}^{N_{\text{tr}}}$ and $\{\mathbf{x}^{(i)}\}_{i=1}^{N_{\text{te}}}$ and then estimating the importance as $\hat{w}(\mathbf{x}) =$

Table 2: Specifications of illustrative classification data.

		Training $p_{\text{tr}}(\mathbf{x}, y)$		Test $p_{\text{te}}(\mathbf{x}, y)$	
		$y = 0$	$y = 1$	$y = 0$	$y = 1$
Figure 4(a)	μ	$(-1, -1)$	$(3, -1)$	$(0, 3.5)$	$(4, 2.5)$
	Σ	$\begin{pmatrix} 0.25 & 0 \\ 0 & 4 \end{pmatrix}$	$\begin{pmatrix} 0.25 & 0 \\ 0 & 4 \end{pmatrix}$	$\begin{pmatrix} 0.25 & 0 \\ 0 & 0.25 \end{pmatrix}$	$\begin{pmatrix} 0.25 & 0 \\ 0 & 0.25 \end{pmatrix}$
Figure 4(b)	μ	$(-1, 0)$	$(4, 2)$	$(0, 2)$	$(3, 1)$
	Σ	$\begin{pmatrix} 0.75 & 0 \\ 0 & 1.5 \end{pmatrix}$	$\begin{pmatrix} 0.75 & 0 \\ 0 & 1.5 \end{pmatrix}$	$\begin{pmatrix} 0.75 & 0.5 \\ 0.01 & 0.1 \end{pmatrix}$	$\begin{pmatrix} 0.75 & 0.5 \\ 0.01 & 0.1 \end{pmatrix}$

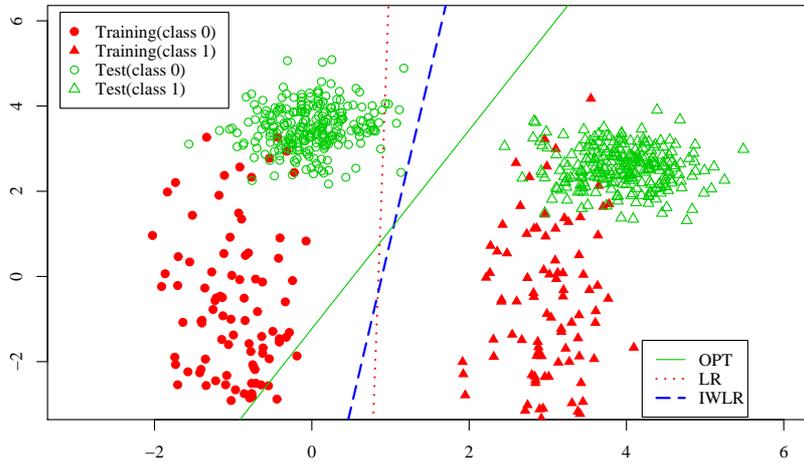
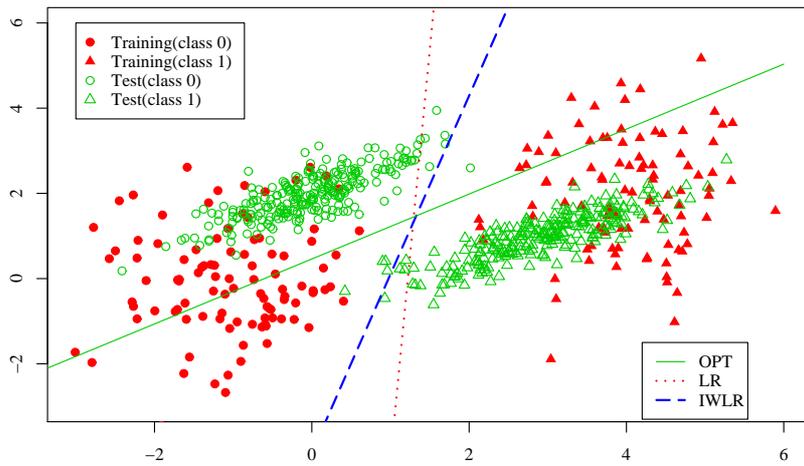
(a) $p_{\text{te}}(\mathbf{x})$ is linearly shifted from $p_{\text{tr}}(\mathbf{x})$.(b) $p_{\text{te}}(\mathbf{x})$ is rotated from $p_{\text{tr}}(\mathbf{x})$.

Figure 4: Classification examples under covariate shift.

$\hat{p}_{\text{te}}(\mathbf{x})/\hat{p}_{\text{tr}}(\mathbf{x})$. A potential limitation of this approach is that KDE suffers from the *curse of dimensionality* [12], since the number of samples needed to maintain the same approximation quality grows exponentially as the dimension of the input space increases. This is particularly critical when estimating $p_{\text{tr}}(\mathbf{x})$ since the number of training input samples is typically limited. In addition, model selection by LCV is unreliable in such cases, since data splitting in the CV procedure further reduces the sample size. Therefore, in high-dimensional cases LL-KLIEP may be more reliable than the KDE-based approach.

6.2 Kernel mean matching

The kernel mean matching (KMM) method avoids density estimation and directly gives an estimate of the importance at the training input points [16].

The basic idea of KMM is to find $w(\mathbf{x})$ such that the mean discrepancy between nonlinearly transformed samples drawn from $p_{\text{te}}(\mathbf{x})$ and $p_{\text{tr}}(\mathbf{x})$ is minimized in a *universal reproducing kernel Hilbert space* [22]. The Gaussian kernel (15) is an example of kernels that induce universal reproducing kernel Hilbert spaces and it has been shown that the solution of the following optimization problem agrees with the true importance:

$$\begin{aligned} \min_{w(\mathbf{x})} \left\| \int K_s(\mathbf{x}, \cdot) p_{\text{te}}(\mathbf{x}) d\mathbf{x} - \int K_s(\mathbf{x}, \cdot) w(\mathbf{x}) p_{\text{tr}}(\mathbf{x}) d\mathbf{x} \right\|_{\mathcal{F}}^2 \\ \text{subject to } \int w(\mathbf{x}) p_{\text{tr}}(\mathbf{x}) d\mathbf{x} = 1 \quad \text{and} \quad w(\mathbf{x}) \geq 0, \end{aligned}$$

where $\|\cdot\|_{\mathcal{F}}$ denotes the norm in the Gaussian reproducing kernel Hilbert space and $K_s(\mathbf{x}, \mathbf{x}')$ is the Gaussian kernel (15).

An empirical version of the above problem is reduced to the following quadratic program:

$$\begin{aligned} \min_{\{w(\mathbf{x})\}_{\mathbf{x} \in D_{\text{tr}}}} \left[\frac{1}{2} \sum_{\mathbf{x}, \mathbf{x}' \in D_{\text{tr}}} w(\mathbf{x}) w(\mathbf{x}') K_s(\mathbf{x}, \mathbf{x}') - \sum_{\mathbf{x} \in D_{\text{tr}}} w(\mathbf{x}) \kappa(\mathbf{x}) \right] \\ \text{subject to } \left| \sum_{\mathbf{x} \in D_{\text{tr}}} w(\mathbf{x}) - N_{\text{tr}} \right| \leq N_{\text{tr}} \epsilon, \text{ and} \\ 0 \leq w(\mathbf{x}) \leq B \text{ for all } \mathbf{x} \in D_{\text{tr}}, \end{aligned}$$

where

$$\kappa(\mathbf{x}) = \frac{N_{\text{tr}}}{N_{\text{te}}} \sum_{\mathbf{x}' \in D_{\text{te}}} K_s(\mathbf{x}, \mathbf{x}').$$

B (≥ 0), ϵ (≥ 0), and s (≥ 0) are tuning parameters. The solution $\{w(\mathbf{x})\}_{\mathbf{x} \in D_{\text{tr}}}$ is an estimate of the importance at the training input points.

Since KMM does not require density estimates, it is expected to work well even in high dimensional cases. However, the performance is dependent on the tuning parameters B , ϵ , and s and they cannot be optimized easily, e.g., by CV, since estimates of the importance

are available only at the training input points. Thus, an out-of-sample extension is needed to apply KMM in the CV framework, but this currently seems to be an open research issue.

Here, we show that LL-KLIEP(LS2) (see Eq.(23)) has a tight connection to KMM. Up to irrelevant constants, Eq.(23) without a regularizer can be expressed as

$$\frac{1}{2} \sum_{\mathbf{x}, \mathbf{x}' \in D_{\text{tr}}} w(\mathbf{x})w(\mathbf{x}')K_s(\mathbf{x}, \mathbf{x}') - \sum_{\mathbf{x} \in D_{\text{tr}}} w(\mathbf{x})\kappa(\mathbf{x}),$$

which is exactly the same form as the objective function of KMM. Thus, KMM and LL-KLIEP(LS2) share a common objective function, although they are derived from very different frameworks.

However, KMM and LL-KLIEP(LS2) still have a significant difference—KMM directly optimizes the importance values $\{w(\mathbf{x})\}_{\mathbf{x} \in D_{\text{tr}}}$, while LL-KLIEP(LS2) optimizes the parameter $\{\beta_{\mathbf{x}}\}_{\mathbf{x} \in D_{\text{tr}}}$ in the importance model (24). Thus, LL-KLIEP(LS2) learns the entire importance function and therefore it allows us to *interpolate* the value of the importance function at any input point. This interpolation property is a significant advantage over KMM since it allows us to use LCV for model selection. Therefore, LL-KLIEP(LS2) may be regarded as an extension of KMM.

6.3 Logistic regression discriminating training and test input data

Another method to directly estimate the importance weights is to use a probabilistic classifier. Let us assign a selector variable $\delta = -1$ to the training inputs and $\delta = 1$ to the test inputs. This means that the training and test input densities are written as

$$p_{\text{tr}}(\mathbf{x}) = p(\mathbf{x}|\delta = -1), \quad p_{\text{te}}(\mathbf{x}) = p(\mathbf{x}|\delta = 1).$$

A simple calculation shows that the importance can be expressed in terms of δ as [3]:

$$w(\mathbf{x}) = \frac{p(\delta = -1)}{p(\delta = 1)} \frac{p(\delta = 1|\mathbf{x})}{p(\delta = -1|\mathbf{x})}. \quad (28)$$

The probability ratio $p(\delta = -1)/p(\delta = 1)$ may be simply estimated using the ratio of the numbers of training and test input samples. The conditional probability $p(\delta|\mathbf{x})$ may be learned by discriminating between the test input samples and the training input samples using LR, where δ plays the role of a class variable (cf. Eq.(6)). Let us train the LR model by regularized maximum likelihood estimation. The objective function to be maximized is given by

$$\text{LR}(\boldsymbol{\alpha}) = \sum_{\mathbf{x} \in D_{\text{te}} \cup D_{\text{tr}}} \delta_{\mathbf{x}} \langle \boldsymbol{\alpha}, \boldsymbol{\psi}(\mathbf{x}) \rangle - \sum_{\mathbf{x} \in D_{\text{te}} \cup D_{\text{tr}}} \log(1 + \exp(\delta_{\mathbf{x}} \langle \boldsymbol{\alpha}, \boldsymbol{\psi}(\mathbf{x}) \rangle)) - \frac{\|\boldsymbol{\alpha}\|^2}{2\sigma^2}, \quad (29)$$

where the first term is the main likelihood term, the second term is a normalizer, and the third term is a regularizer. Since this is a convex optimization problem, the global

Table 3: Relation between the proposed and related methods.

	Model selection	Direct importance model	Optimization
KDE	Available	Not Available	Analytic
KMM	Not Available	Non-parametric	Constraint quadratic program
LogReg	Available	Log-linear	Unconstraint non-linear
KLIEP	Available	Linear	Constraint non-linear
LL-KLIEP	Available	Log-linear	Unconstraint non-linear

solution can be obtained by standard non-linear optimization methods. The gradient of the objective function is given as

$$\frac{\partial \text{LR}(\boldsymbol{\alpha})}{\partial \boldsymbol{\alpha}} = \sum_{\mathbf{x} \in D_{\text{te}} \cup D_{\text{tr}}} \delta_{\mathbf{x}} \boldsymbol{\psi}(\mathbf{x}) - \sum_{\mathbf{x} \in D_{\text{te}} \cup D_{\text{tr}}} \delta_{\mathbf{x}} p_{\boldsymbol{\alpha}}(\delta_{\mathbf{x}} | \mathbf{x}) \boldsymbol{\psi}(\mathbf{x}) - \frac{\|\boldsymbol{\alpha}\|^2}{2\sigma^2}. \quad (30)$$

Then the importance estimate is given by

$$\hat{w}(\mathbf{x}) = \frac{N_{\text{tr}}}{N_{\text{te}}} \exp(\langle \boldsymbol{\alpha}, \boldsymbol{\psi}(\mathbf{x}) \rangle). \quad (31)$$

We refer to this approach as *LogReg*.

Eq.(31) shows that the function model of the importance in LogReg is actually the same as that of LL-KLIEP except for a scaling factor (cf. Eq.(16)). However, the optimization criteria of LL-KLIEP and LogReg are different—in LL-KLIEP, the summation is taken only over the training or test input samples but not both, while the summation in LogReg is over both the training and test input samples. This difference is significant since LogReg does not allow us to use the computational trick we proposed in Section 4.2. Thus LL-KLIEP has the advantage in computation time and storage space consumption over LogReg.

Bickel et al. [3] proposed simultaneous optimization of both importance estimator and classifier. Although their method can perform better than our two stage method which solves importance estimation and classifier’s parameter estimation separately, it has a weakness in model selection. Since the hyper-parameter of their method supposed to be tuned for test samples, CV is not applicable if no labeled test sample is available. On the other hand, our method can select hyper-parameters for both importance estimation by LCV and classification by IWCV.

The characteristics of the proposed and related methods are summarized in Table 3.

7 Experiments

In this section, we experimentally compare the performance of LL-KLIEP with existing methods.

7.1 Toy experiments

Let $p_{\text{tr}}(\mathbf{x}) = \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d)$ and $p_{\text{te}}(\mathbf{x}) = \mathcal{N}((1, 0, \dots, 0)^\top, 0.75^2 \mathbf{I}_d)$. The task is to estimate the importance at the training input points:

$$w(\mathbf{x}) = \frac{p_{\text{te}}(\mathbf{x})}{p_{\text{tr}}(\mathbf{x})} \text{ for } \mathbf{x} \in D_{\text{tr}}.$$

We compared KLIEP, KDE, KMM, LogReg, LL-KLIEP, LL-KLIEP(LS1-a), LL-KLIEP(LS1-b), and LL-KLIEP(LS2). For LL-KLIEP, LL-KLIEP(LS1-a), LL-KLIEP(LS1-b), and LL-KLIEP(LS2), we used 5-fold LCV to choose the regularization parameter σ and the kernel width s . For KLIEP, we use 5-fold LCV to choose the kernel width s . For KDE, we used 5-fold LCV to choose the kernel widths for the training and test densities. For KMM, we used $B = 1000$ and $\epsilon = (\sqrt{N_{\text{tr}}} - 1)/\sqrt{N_{\text{tr}}}$ following the suggestion in the original KMM paper [16]. We tested two different values of the kernel width ($s = 0.1$ and $s = 1.0$) for KMM since there is no reliable method to determine the kernel width. For LogReg, we used 5-fold CV to choose the regularization parameter σ and the kernel width s .

We fixed the number of test input samples at $N_{\text{te}} = 1000$ and considered the following setting for the number of training input samples N_{tr} and the input dimension d :

1. $N_{\text{tr}} = 100$ and $d = 2, 4, \dots, 20$.
2. $d = 10$ and $N_{\text{tr}} = 50, \dots, 150$.

We ran the simulation 100 times for each d and N_{tr} , and evaluated the estimation accuracy of $\{w(\mathbf{x})\}_{\mathbf{x} \in D_{\text{tr}}}$ by the mean NMSE (see Eq.(26)).

The mean NMSE over 100 trials is plotted in Figure 5. The filled plot markers indicate the best method and comparable ones based on the *Wilcoxon* signed rank test at the significance level 1% in terms of the NMSE. We omitted the graphs of LL-KLIEP(LS1-a), LL-KLIEP(LS1-b) and LL-KLIEP(LS2) since they are almost identical to the result of LL-KLIEP. Figure 5(a) shows that the error of KDE sharply increases as the input dimension grows, while LL-KLIEP, KLIEP, and LogReg tend to give much smaller errors than KDE. Figure 5(b) shows that the errors of all methods tend to decrease as the number of training samples grows. Again, LL-KLIEP and LogReg are shown to work well. These would be the fruit of directly estimating the importance without going through density estimation. The results of LL-KLIEP and LogReg are slightly better than KLIEP, perhaps because the original KLIEP does not contain a regularizer; we believe that the performance of KLIEP could be improved by adding a regularizer as used in LL-KLIEP and LogReg. KMM also works reasonably well, as long as that the kernel width s is chosen appropriately. However, the performance of KMM is highly dependent on s and determining its appropriate value may be difficult. Overall, the accuracy of LL-KLIEP is comparable to the best existing approaches.

Next, we compared the computational cost of LL-KLIEP, LL-KLIEP(LS1-a), LL-KLIEP(LS1-b), LL-KLIEP(LS2), and LogReg, which have good accuracy in the previous

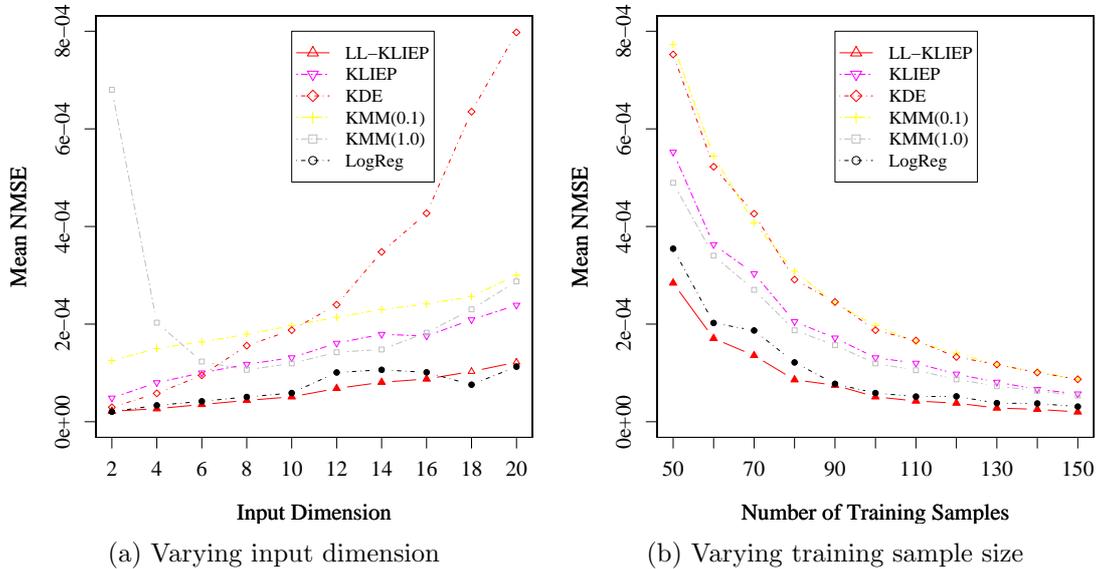


Figure 5: Mean NMSE over 100 trials. The filled plot markers indicate the best method and comparable ones based on the *Wilcoxon* signed rank test at the significance level 1% in terms of the NMSE. ‘KMM(s)’ denotes KMM with kernel width s .

experiments. We investigated the entire computation time of all of them including cross-validation and the precomputation times for the test samples. Note that the Gaussian width s and the regularization parameter σ are chosen over the 5×5 equidistant grid in this experiment for all the methods. We fixed the input dimension at $d = 10$ and changed the number of training input points $N_{\text{tr}} = 10^2, 10^3$ and the number of test samples $N_{\text{te}} = 10^2, 10^3, \dots, 10^6$. We repeated the experiments 100 times for each N_{tr} and N_{te} on the PC server with an Intel[®] Xeon[®] 2.66GHz. All of them are implemented on R (<http://www.r-project.org>) and *conjugate gradient* method was used to optimize their objective functions.

Figure 6 shows the average elapsed times for LL-KLIEP, LL-KLIEP(LS1-a), LL-KLIEP(LS1-b), LL-KLIEP(LS2), and LogReg. The results show that the computational cost of LL-KLIEP and LogReg increases as the amount of test data N_{te} grows, but the computational cost of LL-KLIEP(LS) is nearly independent of the number of test data N_{te} . This is in good agreement with our theoretical analysis in Section 4.2. Thus the cost of dealing with a large amount of test data in each optimization step is much higher than that at one time precomputation.

We also compared the memory usage of LL-KLIEP, LL-KLIEP(LS1-a), LL-KLIEP(LS1-b), LL-KLIEP(LS2), and LogReg. We used the same implementation and computational environment as the previous experiments. Figure 7 shows the memory usage of each method. The results show that the space requirement of LL-KLIEP and LogReg increases as the amount of test data N_{te} grows, but that of LL-KLIEP(LS) is

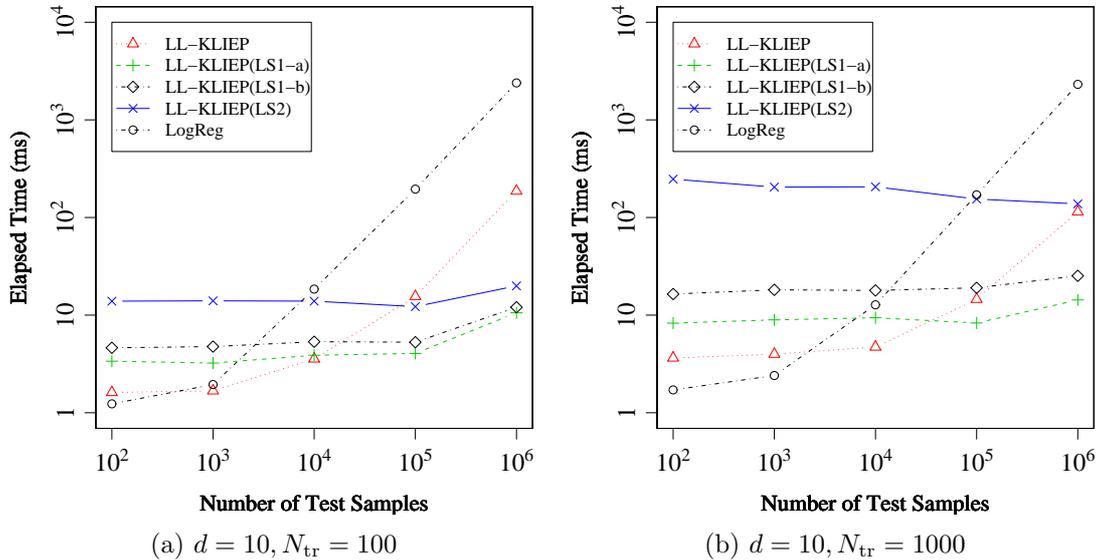


Figure 6: Average computation time over 100 trials. The horizontal axis represents the number of test samples (N_{te}), and the vertical axis represents the elapsed time (millisecond), respectively.

independent of the number of test data N_{te} .

In addition, we compared the computational cost of LL-KLIEP, LL-KLIEP(LS1-a), LL-KLIEP(LS1-b), and LL-KLIEP(LS2) in detail. We examined the combination of the following setting for the number of test samples N_{te} , the number of training inputs N_{tr} , and the input dimension d :

- $N_{te} = 10^2, 10^3, \dots, 10^6$
- $N_{tr} = 10^2, 10^3$
- $d = 10^2, 10^3, 10^4$.

In this experiment, we used linear basis function so that the number of bases is equivalent to the input dimension. Since the computational time of cross-validation is conceptually a scalar multiple of that of each optimization step, we compared the computational time including the precomputation times for the test inputs after the model parameters are fixed. We repeated the experiments 100 times for each N_{te} , N_{tr} , and d using the same implementation and computational environment as the previous experiments.

Figure 8 shows the average elapsed times for LL-KLIEP, LL-KLIEP(LS1-a), LL-KLIEP(LS1-b), and LL-KLIEP(LS2). When $d = 10^3$, the result of $N_{te} = 10^6$ was excluded because of the large memory requirements. As we expected,

1. LL-KLIEP is faster than LL-KLIEP(LS) when the number of test samples is small,

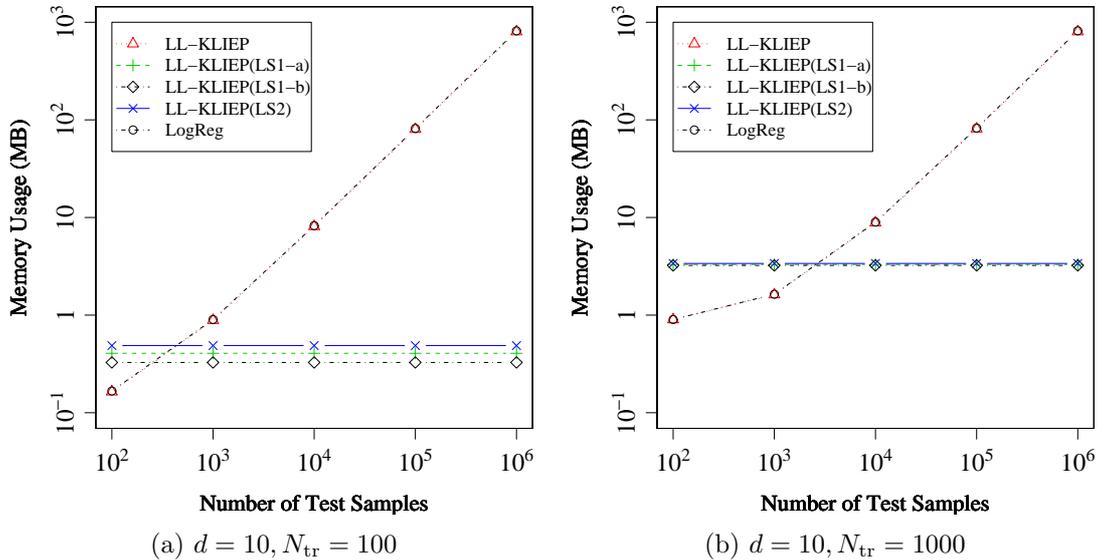


Figure 7: Memory usage. The horizontal axis represents the number of test samples (N_{te}), and the vertical axis represents the memory usage (MB), respectively.

2. LL-KLIEP(LS1-a) is faster than LL-KLIEP(LS2) for lower dimensional data,
3. both LL-KLIEP(LS1-b) and LL-KLIEP(LS2) are advantageous for high dimensional problems, and
4. the computational cost of LL-KLIEP(LS1-b) increases for the larger amount of training inputs.

Since LL-KLIEP(LS1-b) outperformed LL-KLIEP(LS2) in all the settings, we conclude LL-KLIEP(LS1-b) is more suitable for high dimensional problems. However, LL-KLIEP(LS1-a) is faster than LL-KLIEP(LS1-b) for lower dimensional data against the complexity analysis. One reason of this result might be that the computation of LL-KLIEP(LS1-b) is relatively complex in the iteration of the training inputs comparing with LL-KLIEP(LS1-a). Therefore, LL-KLIEP(LS1-a) run faster than LL-KLIEP(LS1-b) if the number of dimension is small enough not to ignore this overhead.

7.2 Natural language processing task

In this section, we show the experimental result that LL-KLIEP is applied to the domain adaptation task of a natural language processing (NLP).

Since statistical NLP systems tend to perform worse in different domains because of differences in vocabulary and writing style, the domain adaptation of NLP system is one of the important issues in NLP [17, 5, 9].

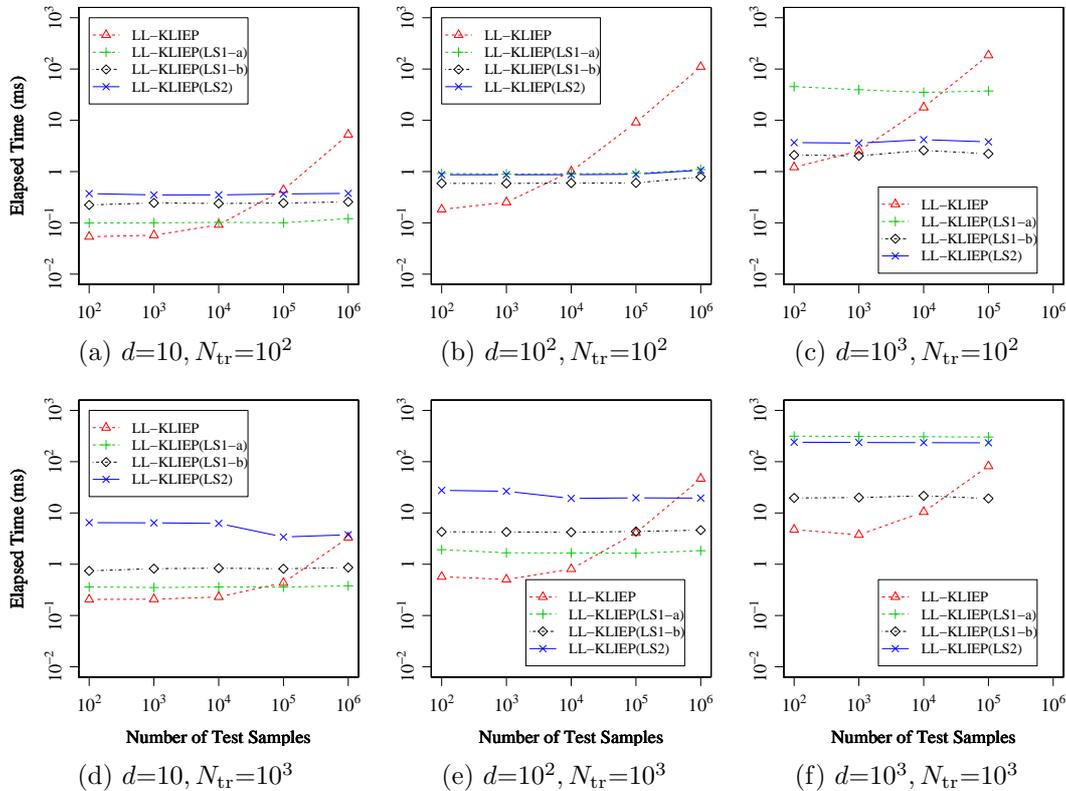


Figure 8: Average computation time over 100 trials. The horizontal axis represents the number of test samples (N_{te}), and the vertical axis represents the elapsed time (millisecond), respectively.

In this experiments, we conducted domain adaptation experiments for the Japanese word segmentation task. It is not trivial to detect word boundaries for non-segmented languages such as Japanese or Chinese. In the word segmentation task, $\mathbf{x}=(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T) \in \mathbf{X}$ represents a given sequence of character boundaries of a sentence and $\mathbf{y}=(y_1, y_2, \dots, y_T) \in \mathbf{Y}$ is a sequence of the corresponding labels, which specify whether the current position is a word boundary. It is reasonable to consider the domain adaptation task of word segmentation system as covariate shift adaptation problem since word segmentation policy ($p(\mathbf{y}|\mathbf{x})$) is rarely changed between domains in the same language, but the distribution of characters ($p(\mathbf{x})$) tends to be changed between domains.

In the experiments, we used the same data and feature set as Tsuboi et al. [28] which aims to adapt a word segmentation system from a daily conversation domain to a medical domain. One of the characteristics of NLP tasks is high dimensionality. The total number of distinct features was about $d = 300K$ in this data. In this experiment, we used labeled (i.e. word segmented) $N_{tr} = 13,000$ sentences for the source domain and unlabeled (i.e. unsegmented) $N_{te} = 53,834$ sentences for the target domain. In addition, we have 1,000 labeled target domain sentences for the evaluation of the domain adaptation task.

As a word segmentation model, we employed *Conditional Random Fields* (CRFs), which are the generalization of the LR (6) for structured prediction [19]. CRFs model the conditional probability $p_{\theta}(\mathbf{y}|\mathbf{x})$ of output structure \mathbf{y} given an input \mathbf{x} :

$$p_{\theta}(\mathbf{y}|\mathbf{x}) = \frac{\exp(\langle \theta, \Phi(\mathbf{x}, \mathbf{y}) \rangle)}{\sum_{\mathbf{y} \in \mathcal{Y}} \exp(\langle \theta, \Phi(\mathbf{x}, \mathbf{y}) \rangle)}, \quad (32)$$

where $\Phi(\mathbf{x}, \mathbf{y}) : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^h$ denote a basis function mapping from a pair of \mathbf{x} and \mathbf{y} to a h dimensional vector. Although conventional CRF learning algorithms minimize the regularized negative log-likelihood, we implemented a importance weighted CRF training algorithm (IWCRF) in the same way as IWLR:

$$\hat{\theta}_{\text{IWCRF}} \equiv \underset{\theta}{\operatorname{argmin}} \left[\sum_{(\mathbf{x}, \mathbf{y}) \in Z_{\text{tr}}} -w(\mathbf{x}) \log p_{\theta}(\mathbf{y}|\mathbf{x}) + \lambda \|\theta\|^2 \right].$$

To estimate the importance of each sentence, we used the source domain data as training data D_{tr} and the unlabeled target domain data as test data D_{te} . We defined the basis function of the importance model $\hat{w}(\mathbf{x})$ as the average value of features for CRFs in a sentence:

$$\psi(\mathbf{x}) = \frac{1}{T} \sum_{t=1}^T \mathbf{x}_t. \quad (33)$$

The performance measure in the experiments is the standard F measure score, $F = 2RP/(R + P)$ where

$$R = \frac{\# \text{ of correct words}}{\# \text{ of words in test data}} \times 100$$

$$P = \frac{\# \text{ of correct words}}{\# \text{ of words in system output}} \times 100.$$

In addition, we tuned the hyper-parameter of IWCRF based on an *importance weighted F measure score*, IWF, which the number of correct words is weighted by the importance of the sentence that these words belong to:

$$\text{IWF}(D) = \frac{2 \times \text{IWR}(D) \times \text{IWP}(D)}{\text{IWR}(D) + \text{IWP}(D)}$$

for the validation set D where

$$\text{IWR}(D) = \frac{\sum_{(\mathbf{x}, \mathbf{y}) \in D} w(\mathbf{x}) \sum_{v_t \in \mathbf{y}} [\hat{v}_t = v_t]}{\sum_{(\mathbf{x}, \mathbf{y}) \in D} \sum_{v_t \in \mathbf{y}} w(\mathbf{x})} \times 100$$

$$\text{IWP}(D) = \frac{\sum_{(\mathbf{x}, \mathbf{y}) \in D} w(\mathbf{x}) \sum_{v_t \in \mathbf{y}} [\hat{v}_t = v_t]}{\sum_{(\mathbf{x}, \mathbf{y}) \in D} \sum_{\hat{v}_t \in \hat{\mathbf{y}}} w(\mathbf{x})} \times 100, \text{ and}$$

Table 4: Word segmentation performance in the target domain. “CRF + 1,000” stands for the performance of a CRF additionally using 1,000 manual word segmentations of the target domain.

	F	R	P
CRF	92.30	90.58	94.08
IWCRF (LL-KLIEP)	94.46	94.32	94.59
IWCRF (LogReg)	93.68	94.30	93.07
CRF + 1,000	94.43	93.49	95.39

v_t denotes a t -th word in a sentence (\mathbf{x}, \mathbf{y}) and \hat{v}_t denotes a t -th word of a system prediction $\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y}} p_{\theta}(\mathbf{y}|\mathbf{x})$. We used 1/10 of training data D_{tr} as the validation set.

Then, we compared the target domain performance between CRF, IWCRF, and CRF which was trained additionally using 1,000 manual word segmentations, described in Tsuboi et al. [28], denoted as “CRF + 1,000”. For importance estimation, we compared LL-KLIEP and LogReg for which we employed 5-fold CV to find the optimal hyperparameter σ .

Table 4 shows the result of the performance of each method. Surprisingly, the F score of IWCRF (LL-KLIEP) outperformed not only that of CRF, but also that of “CRF + 1,000”, so the benefit of the importance weighting is worth the manual labeling of 1,000 words. Especially, the empirical result shows that the covariate shift adaptation technique improves the coverage (R) in the target domain. Comparing with LogReg, LL-KLIEP results in higher final system performance. Since it is easy to obtain large amount of unlabeled text data in NLP tasks, we believe the domain adaptation of NLP tasks is one of the promising application of the proposed method.

8 Conclusions

In this paper, we addressed the problem of estimating the importance for covariate shift adaptation. We proposed a scalable direct importance estimation method called *LL-KLIEP*. The computation time of LL-KLIEP is nearly independent of the amount of test data, which is a significant advantage over existing approaches when we deal with a large number of test samples. Our experiments highlighted this advantage and we experimentally confirmed that the accuracy of the proposed method is comparable to the best existing methods. Finally, a natural language processing experiment shows a promising result of the proposed method for large-scale covariate shift adaptation.

References

- [1] P. Baldi and S. Brunak. *Bioinformatics: The Machine Learning Approach*. MIT Press, Cambridge, 1998.

- [2] S. Bickel. ECML 2006 discovery challenge workshop, 2006.
- [3] S. Bickel, M. Brückner, and T. Scheffer. Discriminative learning for differing training and test distributions. In *Proceedings of the 24th international conference on Machine learning*, pages 81 – 88. ACM Press, 2007.
- [4] S. Bickel and T. Scheffer. Dirichlet-enhanced spam filtering based on biased samples. In *Advances in Neural Information Processing Systems*, pages 161–168. MIT Press, Cambridge, MA, 2007.
- [5] J. Blitzer, R. McDonald, and F. Pereira. Domain adaptation with structural correspondence learning. In *In Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 120–128, 2006.
- [6] K. M. Borgwardt, A. Gretton, M. J. Rasch, H.-P. Kriegel, B. Schölkopf, and A. J. Smola. Integrating structured biological data by kernel maximum mean discrepancy. *Bioinformatics*, 22(14):e49–e57, 2006.
- [7] J. Q. Candela, N. Lawrence, A. Schwaighofer, and M. Sugiyama. NIPS 2006 workshop on learning when test and training inputs have different distributions, 2006.
- [8] D. A. Cohn, Z. Ghahramani, and M. I. Jordan. Active learning with statistical models. *Journal of Artificial Intelligence Research*, 4:129–145, 1996.
- [9] H. Daumé III. Frustratingly easy domain adaptation. In *In Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics*, pages 256–263, 2007.
- [10] V. V. Fedorov. *Theory of Optimal Experiments*. Academic Press, New York, 1972.
- [11] G. S. Fishman. *Monte Carlo: Concepts, Algorithms, and Applications*. Springer-Verlag, Berlin, 1996.
- [12] W. Härdle, M. Müller, S. Sperlich, and A. Werwatz. *Nonparametric and Semiparametric Models*. Springer, Berlin, 2004.
- [13] D. Hawking, E. Voorhees, N. Craswell, and P. Bailey. Overview of the trec-8 web track. In *Proceedings of the Eighth Text REtrieval Conference*, pages 131–150, 1999.
- [14] J. J. Heckman. Sample selection bias as a specification error. *Econometrica*, 47(1):153–162, 1979.
- [15] H. Hirotaka, T. Akiyama, M. Sugiyama, and J. Peters. Adaptive importance sampling with automatic model selection in value function approximation. In *Proceedings of the Twenty-Third AAAI Conference on Artificial Intelligence*, pages 1351–1356, Chicago, USA, 2008.

- [16] J. Huang, A. J. Smola, A. Gretton, K. M. Borgwardt, and B. Schölkopf. Correcting sample selection bias by unlabeled data. In *Advances in Neural Information Processing Systems*, pages 601–608, Cambridge, MA, 2007. MIT Press.
- [17] J. Jiang and C. Zhai. Instance weighting for domain adaptation in nlp. In *Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics*, pages 264–271, 2007.
- [18] T. Kanamori and H. Shimodaira. Active learning algorithm using the maximum weighted log-likelihood estimator. *Journal of Statistical Planning and Inference*, 116(1):149–162, 2003.
- [19] J. Lafferty, A. McCallum, and F. Pereira. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *Proceedings of the 18th International Conference on Machine Learning*, pages 282–289, 2001.
- [20] C. R. Shelton. *Importance Sampling for Reinforcement Learning with Multiple Objectives*. PhD thesis, Massachusetts Institute of Technology, 2001.
- [21] H. Shimodaira. Improving predictive inference under covariate shift by weighting the log-likelihood function. *Journal of Statistical Planning and Inference*, 90(2):227–244, 2000.
- [22] I. Steinwart. On the influence of the kernel on the consistency of support vector machines. *Journal of Machine Learning Research*, 2:67–93, 2001.
- [23] M. Sugiyama. Active learning in approximately linear regression based on conditional expectation of generalization error. *Journal of Machine Learning Research*, 7:141–166, Jan. 2006.
- [24] M. Sugiyama, M. Krauledat, and K.-R. Müller. Covariate shift adaptation by importance weighted cross validation. *Journal of Machine Learning Research*, 8:985–1005, May 2007.
- [25] M. Sugiyama and K.-R. Müller. Input-dependent estimation of generalization error under covariate shift. *Statistics & Decisions*, 23(4):249–279, 2005.
- [26] M. Sugiyama, S. Nakajima, H. Kashima, P. von Bünau, and M. Kawanabe. Direct importance estimation with model selection and its application to covariate shift adaptation. In *Advances in Neural Information Processing Systems*, Cambridge, MA, 2008. MIT Press.
- [27] R. S. Sutton and G. A. Barto. *Reinforcement Learning: An Introduction*. MIT Press, Cambridge, MA, 1998.
- [28] Y. Tsuboi, H. Kashima, S. Mori, H. Oda, and Y. Matsumoto. Training conditional random fields using incomplete annotations. In *Proceedings of 22nd International Conference on Computational Linguistics*, pages 897–904, 2008.

- [29] G. Wahba. *Spline Model for Observational Data*. Society for Industrial and Applied Mathematics, Philadelphia and Pennsylvania, 1990.
- [30] D. P. Wiens. Robust weights and designs for biased regression models: Least squares and generalized M-estimation. *Journal of Statistical Planning and Inference*, 83(2):395–412, 2000.
- [31] J. R. Wolpaw, N. Birbaumer, D. J. McFarland, G. Pfurtscheller, and T. M. Vaughan. Brain-computer interfaces for communication and control. *Clinical Neurophysiology*, 113(6):767–791, 2002.
- [32] B. Zadrozny. Learning and evaluating classifiers under sample selection bias. In *Proceedings of the 21st International Conference on Machine Learning*, New York, NY, 2004. ACM Press.