



Direct Density Ratio Estimation for Large-scale Covariate Shift Adaptation

Yuta Tsuboi (*), Hisashi Kashima, Shohei Hido,
IBM Research, Tokyo Research Laboratory

** Nara Institute of Science and Technology*

Steffen Bickel,

Max Planck Institute for Computer Science

Masashi Sugiyama

Tokyo Institute of Technology, Department of Computer Science

Abstract

- **Problem:** *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.
- **Contributions:**
 - **LL-KLIEP:** KLIEP (Sugiyama, *et. al.* 2007) for Log-linear models
 - Natural modeling for density ratio functions
 - **LL-KLIEP(LS):** Another optimization technique for LL-KLIEP
 - the computation time is nearly independent of the number of test input samples, and
 - the memory requirement is independent of the number of test input samples
 - which is beneficial in applications with large numbers of unlabeled samples.

Covariate shift situation

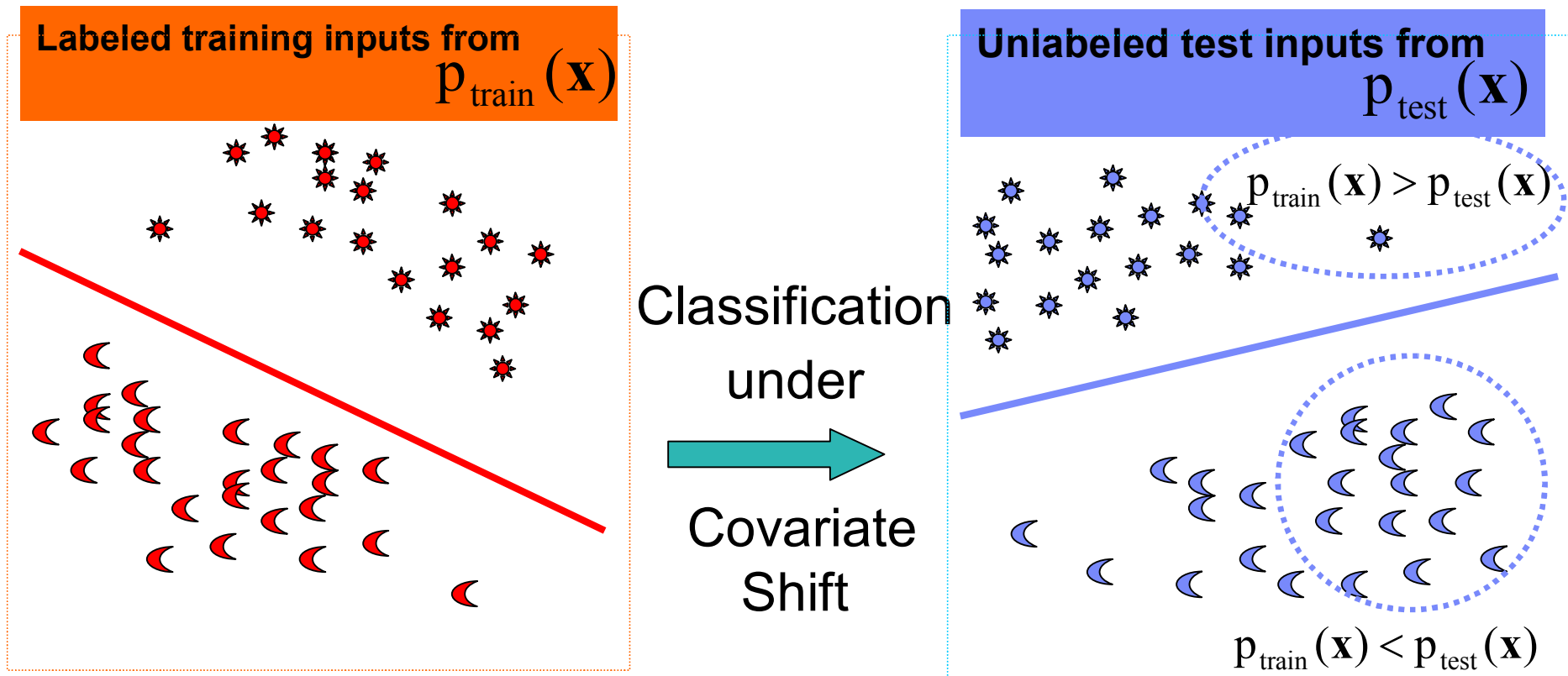
Training and test *inputs* \mathbf{x} follow different distributions

- Input distribution changes:

$$p_{\text{train}}(\mathbf{x}) \neq p_{\text{test}}(\mathbf{x})$$

- Functional relation remains unchanged:

$$p_{\text{train}}(y | \mathbf{x}) = p_{\text{test}}(y | \mathbf{x})$$

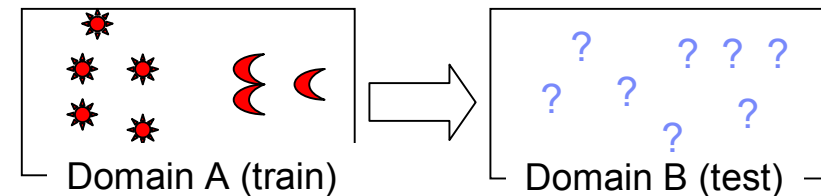


Examples of covariate shift situation

Domain Adaptation & Selective Sampling (Active Learning)

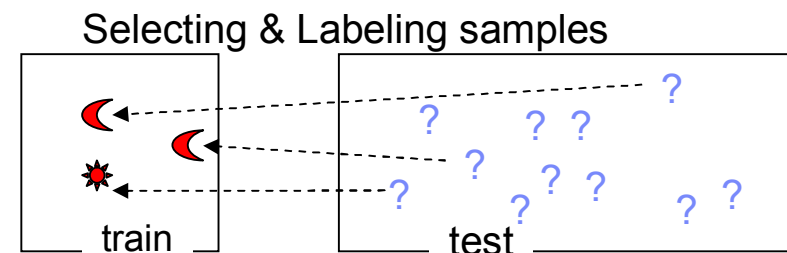
■ Domain adaptation of statistical classifiers

- The data distribution in the test domain is different from that in the training domain. (Note: the functional relation can also be changed)
 - E.g.: Spam filters can be trained on public collections of spam, but are applied to an individual person's inbox. (Personalization)



■ Selective sampling (active learning) of statistical classifiers

- The learning algorithm can actively query the teacher for labels.
- Since the learner chooses the examples by design, the data distribution of the labeled training examples is different from that of a sample pool.



A common approach for covariate shift situation

Weighting the training examples by importance.

- **Density ratio (importance):** $w(\mathbf{x}) = \frac{p_{\text{test}}(\mathbf{x})}{p_{\text{train}}(\mathbf{x})}$.

- **Example: Importance Weighted Logistic Regression (IWLR)**
- Weighted Log-likelihood for Logistic Regression (LR)

$$\iint p_{\text{test}}(\mathbf{x}) p(y | \mathbf{x}) p_{\theta}(y | \mathbf{x}) d\mathbf{x} dy$$

LR model

Labeled training inputs

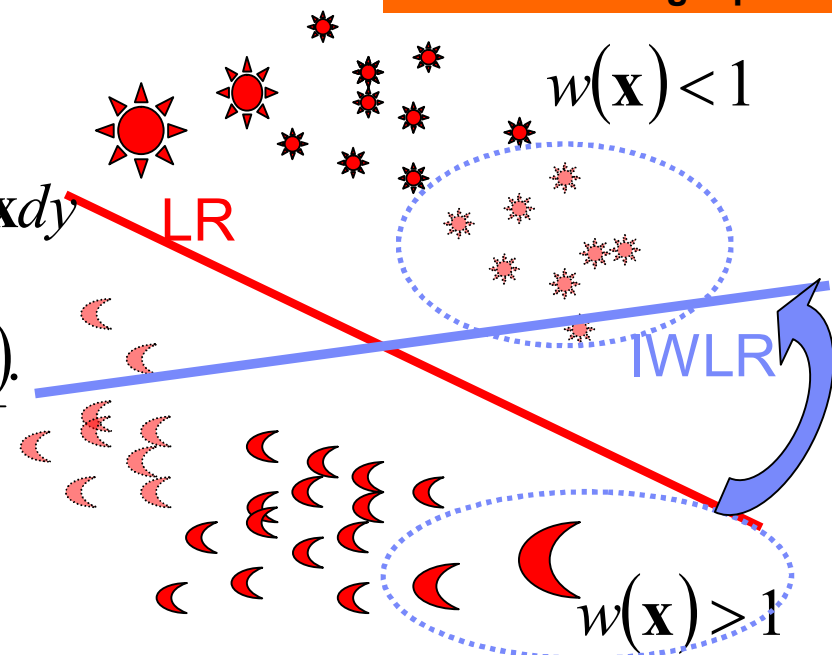
$$= \iint p_{\text{train}}(\mathbf{x}) \frac{p_{\text{test}}(\mathbf{x})}{p_{\text{train}}(\mathbf{x})} p(y | \mathbf{x}) \log p_{\theta}(y | \mathbf{x}) d\mathbf{x} dy$$

$$\approx \sum_{(\mathbf{x}, y) \in Z_{\text{train}}} w(\mathbf{x}) (y f_{\theta}(\mathbf{x}) - \log(1 + \exp(y f_{\theta}(\mathbf{x}))))$$

Training data

Density ratio
(Importance)

Log-likelihood



We need to estimate the density ratio from samples. Importance Estimation

- **Problem setting:** i.i.d. training and test samples are given

Training inputs: $D_{\text{tr}} = \{x_i\}_{i=1}^{N_{\text{tr}}}$ from $P_{\text{train}}(\mathbf{x})$

Test inputs: $D_{\text{te}} = \{x_i\}_{i=1}^{N_{\text{te}}}$ from $P_{\text{test}}(\mathbf{x})$

- Naïve approach: estimate $P_{\text{train}}(\mathbf{x})$ and $P_{\text{test}}(\mathbf{x})$ separately, and take the ratio of the density estimates
- However, density $P(\mathbf{x})$ estimation is the hard problem particularly in high dimensional cases.

Modeling Density Ratio by Log-linear Model

- We use a **log-linear model**:

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

$\boldsymbol{\alpha}$: model parameter

$\boldsymbol{\psi}(\mathbf{x})$: basis function

- Log-linear model
 - $\hat{w}(\mathbf{x})$ takes only non-negative values.
 - **natural modeling for ratio** ($\boldsymbol{\alpha}$ and $\boldsymbol{\psi}(\mathbf{x})$ can be an arbitrary value)

- Test density** is approximated by

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

$$p_{test}(\mathbf{x}) = p_{train}(\mathbf{x}) \frac{p_{test}(\mathbf{x})}{p_{train}(\mathbf{x})}$$

- Learn $\boldsymbol{\alpha}$** so that $\hat{p}_{test}(\mathbf{x})$ approximates $p_{test}(\mathbf{x})$ well.

→ The denominator guarantees $\hat{p}_{test}(\mathbf{x})$ be a probability density function

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

Training data set

Kullback—Leibler (KL) Divergence

- Minimize KL divergence between $p_{\text{test}}(\mathbf{x})$ and $\hat{p}_{\text{test}}(\mathbf{x})$:

$$\arg \min_{\alpha} \text{KL}[p_{\text{test}}(\mathbf{x}) \parallel \hat{p}_{\text{test}}(\mathbf{x})]$$

$$\hat{p}_{\text{test}}(\mathbf{x}) = p_{\text{train}}(\mathbf{x})\hat{w}(\mathbf{x})$$

$$\begin{aligned} & \text{KL}[p_{\text{test}}(\mathbf{x}) \parallel \hat{p}_{\text{test}}(\mathbf{x})] \\ &= \int p_{\text{test}}(\mathbf{x}) \log \frac{p_{\text{test}}(\mathbf{x})}{\hat{p}_{\text{train}}(\mathbf{x})\hat{w}(\mathbf{x})} d\mathbf{x} \\ &= \underbrace{\int p_{\text{test}}(\mathbf{x}) \log \frac{p_{\text{test}}(\mathbf{x})}{p_{\text{train}}(\mathbf{x})} d\mathbf{x}}_{\text{constant}} - \underbrace{\int p_{\text{test}}(\mathbf{x}) \log \hat{w}(\mathbf{x}) d\mathbf{x}}_{\text{relevant}} \end{aligned}$$

Learning Importance Function

- Thus,

$$\arg \min_{\alpha} \text{KL}[p_{\text{test}}(\mathbf{x}) \parallel \hat{p}_{\text{test}}(\mathbf{x})]$$

$$\Leftrightarrow \arg \max_{\alpha} \int p_{\text{test}}(\mathbf{x}) \log \hat{w}(\mathbf{x}) d\mathbf{x} \quad \text{Objective function}$$

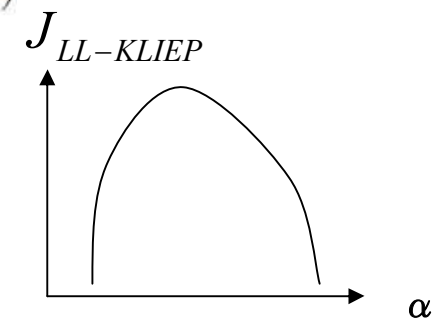
- Empirical approximation of objective function (*LL-KLIEP*)

$$J_{LL-KLIEP}(\alpha) = \frac{1}{N_{te}} \sum_{\mathbf{x} \in D_{te}} \log \hat{w}(\mathbf{x})$$

Test data set

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

- Unconstraint convex optimization:
 - standard gradient ascent method can be used.
 - unique global solution is available.



Kullback-Leibler Importance Estimation Procedure (KLIEP) for Log-linear Models: LL-KLIEP

- Regularized version of LL-KLIEP

$$\begin{aligned}
 j(\alpha) &= \frac{1}{N_{te}} \sum_{\mathbf{x} \in D_{te}} \langle \alpha, \psi(\mathbf{x}) \rangle \\
 &\quad - \log \frac{1}{N_{tr}} \sum_{\mathbf{x} \in D_{tr}} \exp(\langle \alpha, \psi(\mathbf{x}) \rangle) - \frac{\|\alpha\|^2}{2\sigma^2}
 \end{aligned}$$

regularizer

- Gradient of the objective function

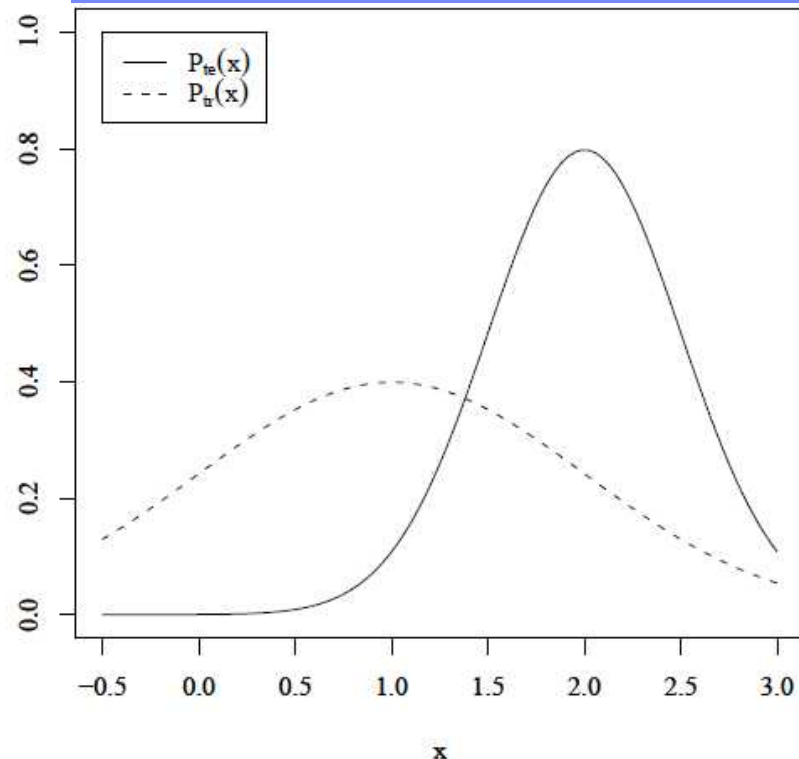
$$\begin{aligned}
 \frac{\partial j(\alpha)}{\partial \alpha} &= \frac{1}{N_{te}} \sum_{\mathbf{x} \in D_{te}} \psi(\mathbf{x}) \\
 &\quad - \sum_{\mathbf{x} \in D_{tr}} \frac{\exp(\langle \alpha, \psi(\mathbf{x}) \rangle)}{\sum_{\mathbf{x}' \in D_{te}} \exp(\langle \alpha, \psi(\mathbf{x}') \rangle)}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

Samples were generated from two Gaussian distributions. We used 100 Gaussian basis functions (Gaussian kernels) centered at randomly chosen test input samples.

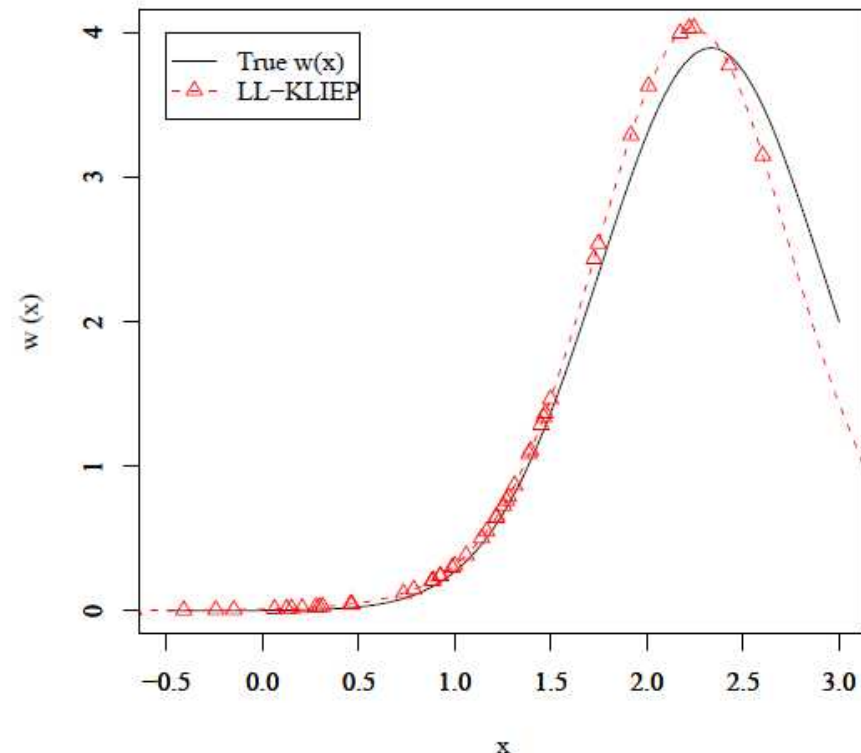
$$\hat{w}(\mathbf{x}) = \frac{\exp(\langle \boldsymbol{\alpha}, \boldsymbol{\psi}(\mathbf{x}) \rangle)}{\frac{1}{N_{tr}} \sum_{\mathbf{x}' \in D_{tr}} \exp(\langle \boldsymbol{\alpha}, \boldsymbol{\psi}(\mathbf{x}') \rangle)}$$

$$\psi_l(\mathbf{x}) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}_l^{\text{test}}\|^2}{2s^2}\right)$$

Training and Test Densities



Estimated Importance



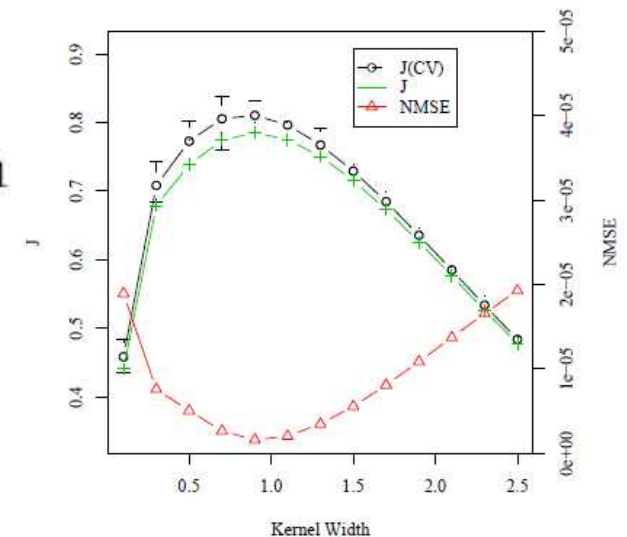
Model selection of KLIEP/LL-KLIEP Likelihood Cross Validation (LCV)

- The performance of KLIEP depends on the choice of the basis functions $\psi(\mathbf{x})$
 → How to choose hyper parameters, e.g., the kernel width s for Gaussian kernels:

$$K_s(x, x_l) = \exp \left\{ -\frac{\|x - x_l\|^2}{2s^2} \right\},$$

- However, the correct value of importance for each \mathbf{x} is not available for unknown distributions $p_{\text{train}}(\mathbf{x})$ and $p_{\text{test}}(\mathbf{x})$
 → unsupervised learning setting

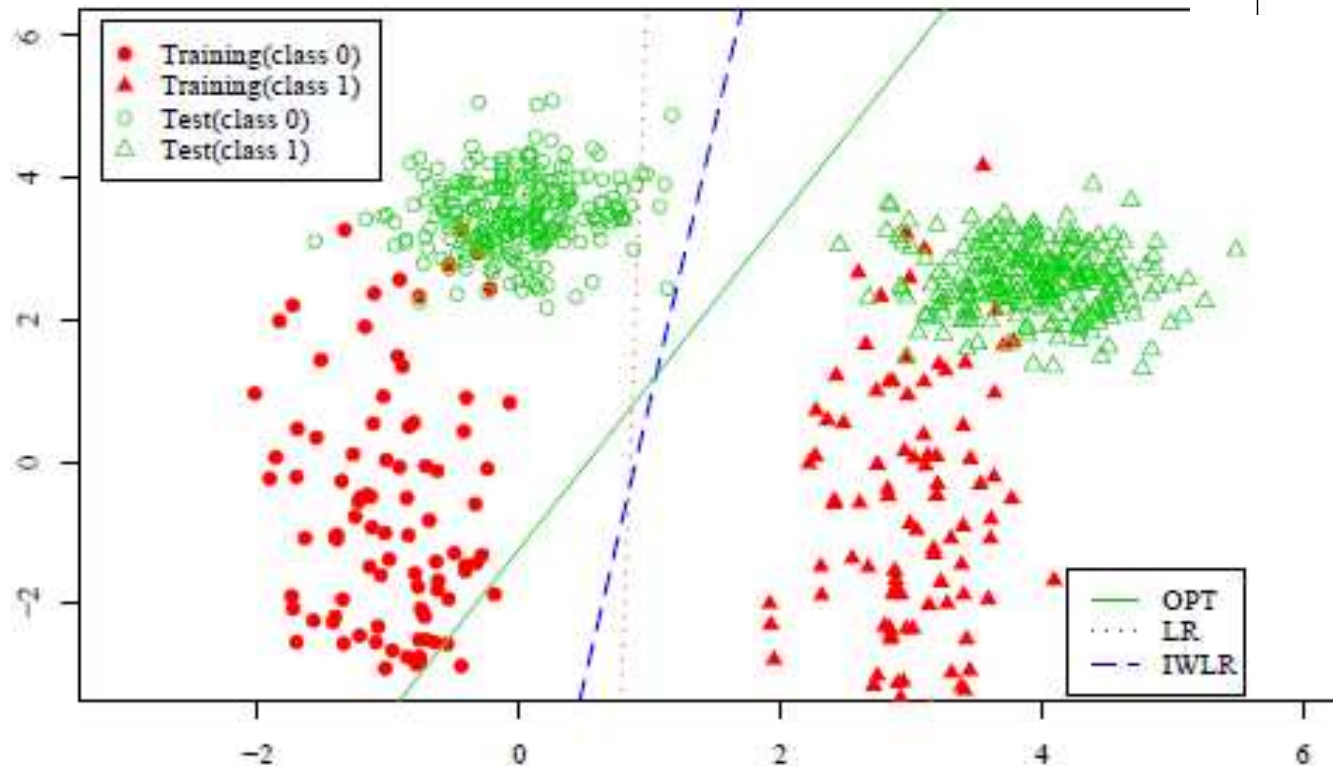
- LCV: Select the model such that maximized $J(\alpha)$
 - Divide test samples into R disjoint subsets: $\{D_{\text{te}}^r\}_{r=1}^R$
 - Learn importance: $\hat{w}^r(\mathbf{x})$ from $\{D_{\text{te}}^t\}_{t \neq r}^R$
 - Evaluate a model by likelihood: $\frac{1}{|D_{\text{te}}^r|} \sum_{\mathbf{x} \in D_{\text{te}}^r} \hat{w}^r(\mathbf{x})$



Classification example under Covariate shift 2-dimensional samples were generated from Gaussian distributions

- We used Importance Weighted Logistic Regression (IWLR)

	Training $p_{tr}(\mathbf{x}, y)$		Test $p_{te}(\mathbf{x}, y)$	
	$y = 0$	$y = 1$	$y = 0$	$y = 1$
μ	$(-1, -1)$	$(3, -1)$	$(0, 3.5)$	$(4, 2.5)$
Σ	$\begin{pmatrix} 0.25 & 0 \\ 0 & 4 \end{pmatrix}$		$\begin{pmatrix} 0.25 & 0 \\ 0 & 0.25 \end{pmatrix}$	



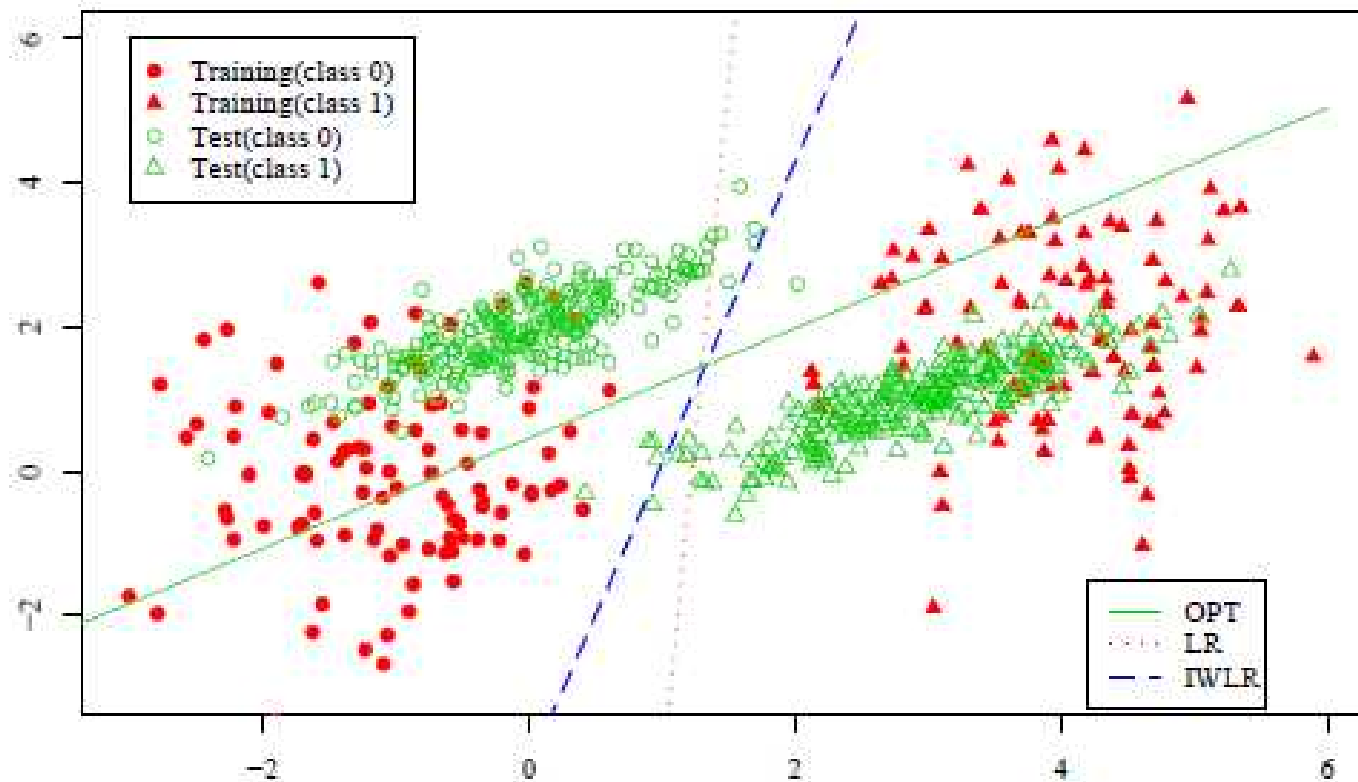
Correct classification rate of LR is 99.1% while that of IWLR is 100%.

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

Classification example under Covariate shift 2-dimensional samples were generated from Gaussian distributions

- We used Importance Weighted Logistic Regression (IWLR).

$$\begin{array}{c|cc|cc} \mu & (-1,0) & (4,2) & (0,2) & (3,1) \\ \Sigma & \begin{pmatrix} 0.75 & 0 \\ 0 & 1.5 \end{pmatrix} & & \begin{pmatrix} 0.75 & 0.5 \\ 0.01 & 0.1 \end{pmatrix} & \end{array}$$



Correct classification rate of LR is 97.2% while that of IWLR is 99.1%.

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

Related Work of Density Ratio Estimation

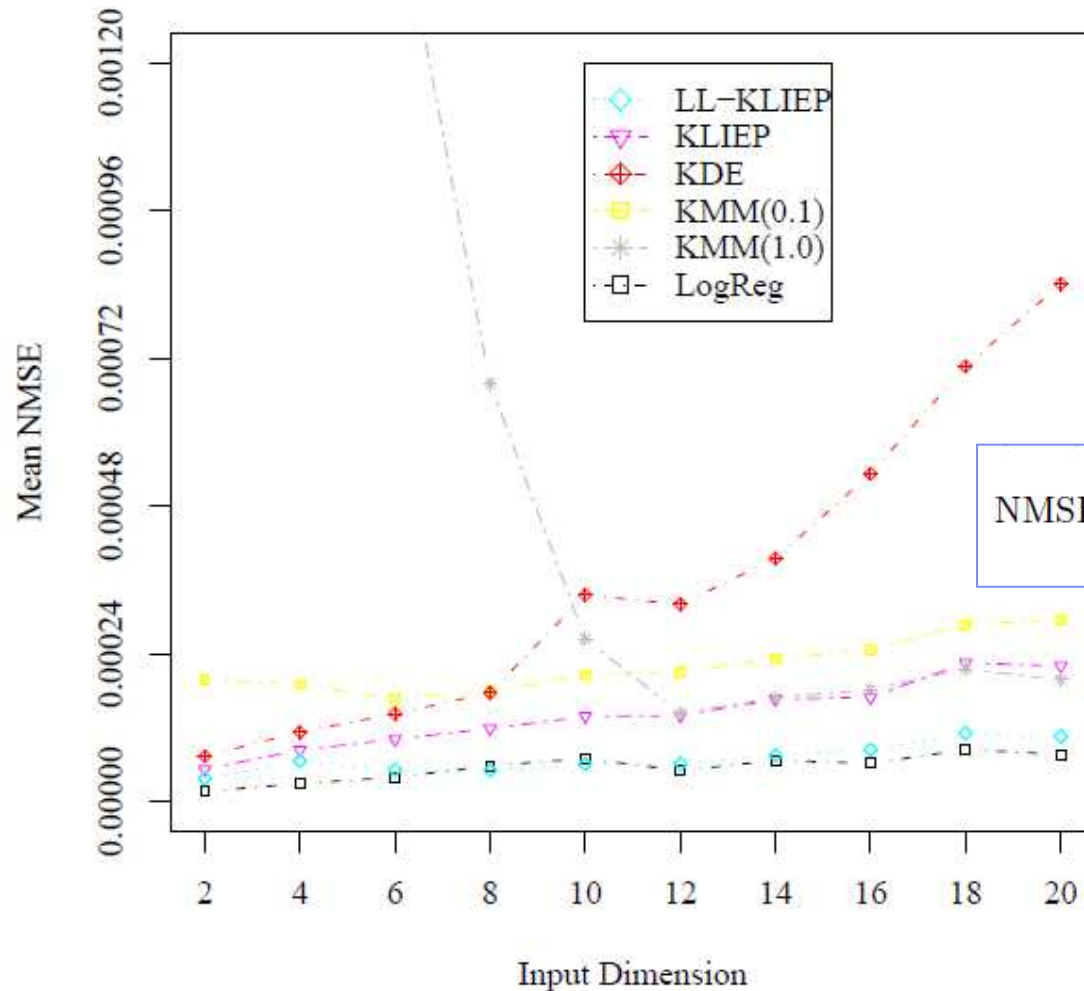
$$w(\mathbf{x}) = \frac{p_{\text{test}}(\mathbf{x})}{p_{\text{train}}(\mathbf{x})}.$$

- **Kernel density estimator (KDE)**
 - Separately estimate training and test input densities.
 - Gaussian kernel width is chosen by likelihood cross-validation.
- **Kernel Mean Matching (KMM)** (Huang *et al.*, NIPS2006)
 - Direct importance estimation method in universal reproducing kernel Hilbert spaces (RKHS)
 - There is no model selection method for kernel width.
- **Logistic regression (LogReg)** (Beckel *et al.*, ICML2007)
 - Classifier discriminating training and test input data.
 - Gaussian kernel width is chosen by likelihood cross-validation.
- **Kullback-Leibler Importance Estimation Procedure (KLIEP)** (Sugiyama *et al.*, NIPS2007)
 - Direct importance estimation method using KL Divergence.
 - Gaussian kernel width is chosen by likelihood cross-validation.

Experiments varying input dimension

$$p_{\text{tr}}(\mathbf{x}) = \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d)$$

$$p_{\text{te}}(\mathbf{x}) = \mathcal{N}((1, 0, \dots, 0)^\top, 0.75^2 \mathbf{I}_d)$$



Mean NMSE over 100 trials.

KMM (s) denotes KMM with kernel width s

NMSE:

Normalized Mean Squared Error

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

KDE: Suffers from the curse of dimensionality

KMM: Performance depends on kernel width

KLIEP, LogReg, and LL-KLIEP: Model selection by LCV works well

Disadvantage: LL-KLIEP and previous methods require to use all test inputs in their optimization procedure.

- We need to iterate over all test inputs when computing the values of the objective function:

$$j(\alpha) = \frac{1}{N_{te}} \sum_{x \in D_{te}} \langle \alpha, \psi(x) \rangle$$

Evaluation over Test data set

$$- \log \frac{1}{N_{tr}} \sum_{x \in D_{tr}} \exp(\langle \alpha, \psi(x) \rangle) - \frac{\|\alpha\|^2}{2\sigma^2}$$

Evaluation over Training data set

- However, the derivative of the objective function requires the evaluation of all test samples once.

$$\frac{\partial j(\alpha)}{\partial \alpha} = F - \frac{1}{N_{tr}} \sum_{x \in D_{tr}} \hat{w}(x) \psi(x) - \frac{\alpha}{\sigma^2}$$

$$F = \frac{1}{N_{te}} \sum_{x \in D_{te}} \psi(x): \text{Independent from } \alpha \rightarrow \text{Pre-computing the value}$$

An optimization technique w/o the objective function evaluation LL-KLIEP(LS1)

- Idea: the derivative of the convex objective function to be zero at the optimum point.
→ Minimizing a squared norm to measure the ‘magnitude’ of the derivative:

Objective function for LL-KLIEP(LS1) $J_{LS}(\alpha) = \frac{1}{2} \left\| \frac{\partial j(\alpha)}{\partial \alpha} \right\|^2 .$

- The partial derivative of LL-KLIEP(LS1):

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

- Computation time & memory size are independent of N_{te} .
 - However, the derivative is a quadratic function of the number of parameters, which could be a bottleneck in high dimensional problems.

LL-KLIEP(LS) for the high-dimensional data LL-KLIEP(LS2)

- Idea: representing the parameter α as a linear combination of the training inputs (representer theorem (Wahba 1990)):

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

where $\{\beta_{\mathbf{x}}\}_{\mathbf{x} \in D_{\text{tr}}}$ is a data-wise parameter.

Objective function for LL-KLIEP(LS2)

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

where

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

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

- Now, the computation time is linear w.r.t. the number of parameters (quadratic w.r.t. the number of the training inputs).

LL-KLIEP (LS): No iteration and no storage for N_{te} in optimization → Well-suited to the applications with the large amount of test samples

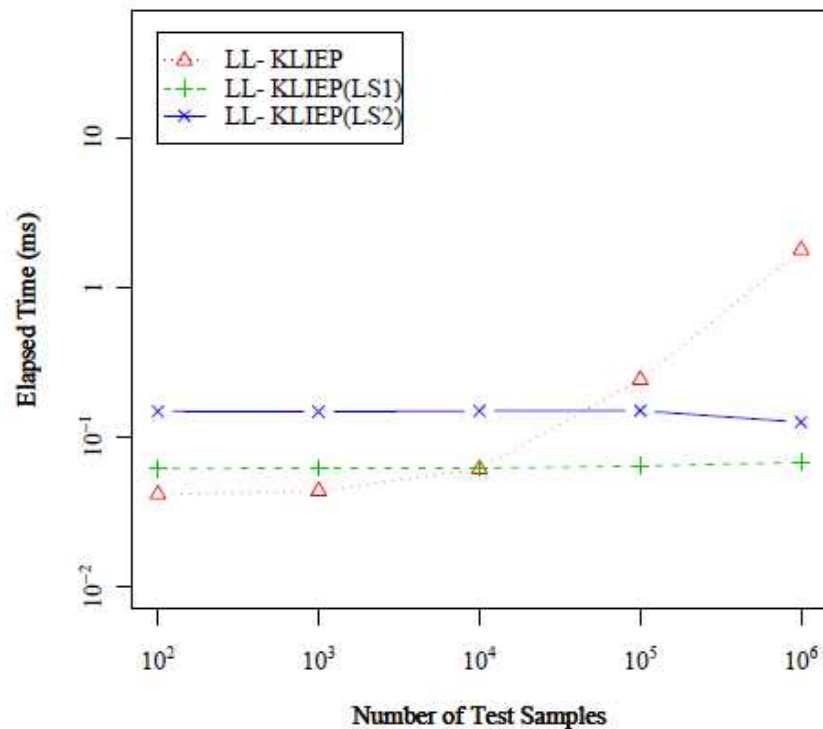
Computational complexity and space requirements. N_{tr} is the number of training samples, N_{te} is the number of test samples, b is the number of parameters, and c is the average number of non-zero basis entries.

	Computational complexity			Space requirement	
	Pre. Comp. (once)	Objective	Derivative	Objective	Derivative
KLIEP	0	$bN_{tr} + bN_{te}$	$bN_{tr} + bN_{te}$	$cN_{tr} + cN_{te}$	$cN_{tr} + cN_{te}$
LL-KLIEP	bN_{te}	$bN_{tr} + bN_{te}$	bN_{tr}	$cN_{tr} + cN_{te}$	cN_{tr}
LL-KLIEP(LS1)	bN_{te}	bN_{tr}	b^2N_{tr}	cN_{tr}	$b^2 + cN_{tr}$
LL-KLIEP(LS2)	bN_{te}	bN_{tr}^2	bN_{tr}^2	cN_{tr}	$N_{tr}^2 + cN_{tr}$

- LL-KLIEP (LS1) : For lower-dimensional and large-scale training data.
- LL-KLIEP (LS2): For higher-dimensional and moderate-size training data.

Average computation time (including Pre-comp.) over 100 trials We varied the number of test inputs, and fixed the number of training inputs.

- we used linear basis function so that the number of bases is equivalent to the input dimension.
 - > d : input dimension = #parameter,
 N_{tr} : The number of training inputs, N_{te} : The number of test inputs



Lower-dimensional data

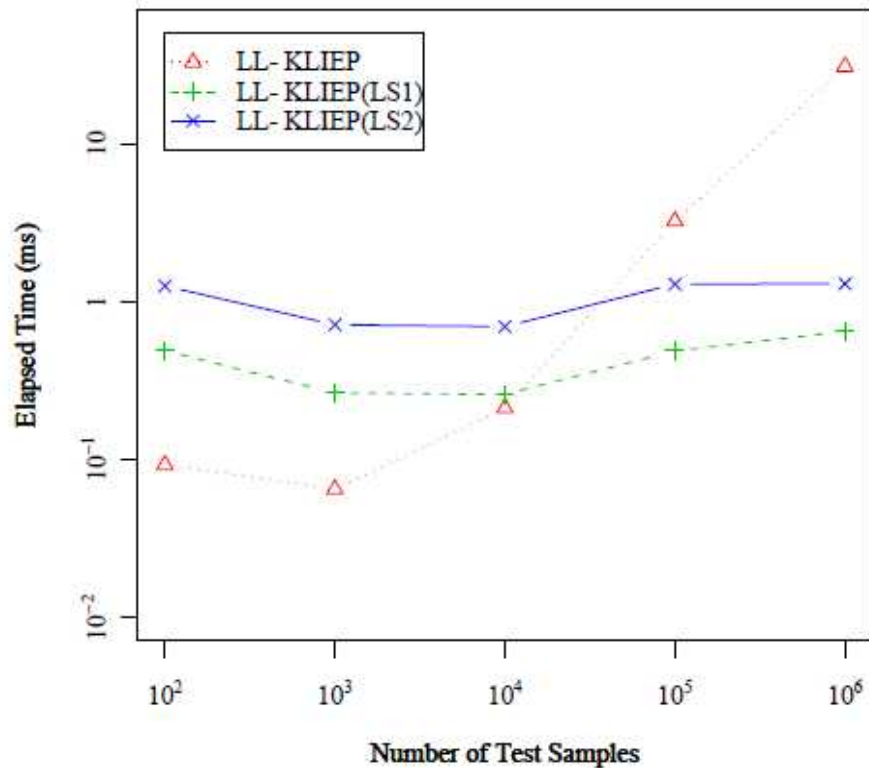
The computation time of LL-KLIEP(LS) is independent from the number of test inputs.

(a) $d = 10, N_{tr} = 100$

Average computation time (including Pre-comp.) over 100 trials We varied the number of test inputs, and fixed the number of training inputs.

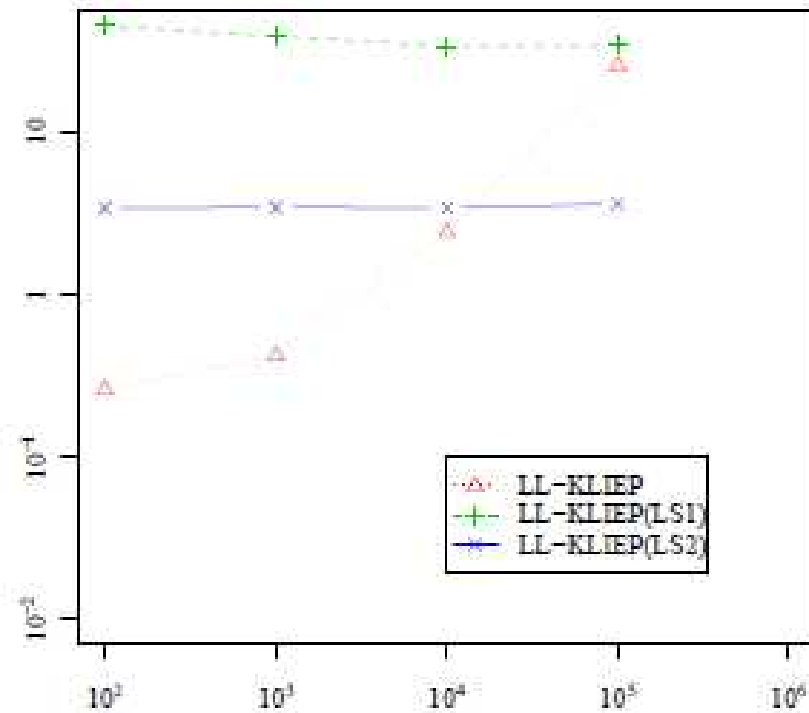
- d : input dimension = #parameter,
 N_{tr} : The number of training inputs, N_{te} : The number of test inputs

Moderate-dimensional data



(b) $d = 100, N_{tr} = 100$

Higher-dimensional data



(c) $d = 1000, N_{tr} = 100$

Related work: Kernel Mean Matching (KMM)
LL-KLIEP (LS2) without a regularizer has the same form as the objective function of KMM.

- Moment matching method:

Objective function for KMM

$$\min_{\{w(\mathbf{x})\}_{\mathbf{x} \in D_{\text{tr}}}} \left[\frac{1}{2} \sum_{\mathbf{x}, \mathbf{x}' \in D_{\text{tr}}} w(\mathbf{x})w(\mathbf{x}')K_s(\mathbf{x}, \mathbf{x}') - \sum_{\mathbf{x} \in D_{\text{tr}}} w(\mathbf{x})\kappa(\mathbf{x}) \right]$$

subject to $\left| \sum_{\mathbf{x} \in D_{\text{tr}}} w(\mathbf{x}) - N_{\text{tr}} \right| \leq N_{\text{tr}}\epsilon$, and

$$0 \leq w(\mathbf{x}) \leq B \text{ for all } \mathbf{x} \in D_{\text{tr}},$$

where

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

Disadvantage of KMM.

The estimates of $w(\mathbf{x})$ are only available for training samples \rightarrow Cannot optimize hyper parameters by CV

- The objective function of LL-KLIEP (LS2):

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

Related work: Logistic regression (LogReg) Classifier discriminating training and test input data

- Selector variable $\delta = -1$ to the training input samples and $\delta = 1$ to the test input samples:

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

- Importance can be $w(\mathbf{x}) = \frac{p(\delta = -1)}{p(\delta = 1)} \frac{p(\delta = 1|\mathbf{x})}{p(\delta = -1|\mathbf{x})}$.
- 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.

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

Empirical estimation

- Objective function: regularized maximum likelihood estimation
- Disadvantage: summation over both training and test samples in CV.