

Density Ratio Estimation in Machine Learning



Masashi Sugiyama

Tokyo Institute of Technology, Japan

sugi@cs.titech.ac.jp

<http://sugiyama-www.cs.titech.ac.jp/~sugi/>

Machine Learning (ML)

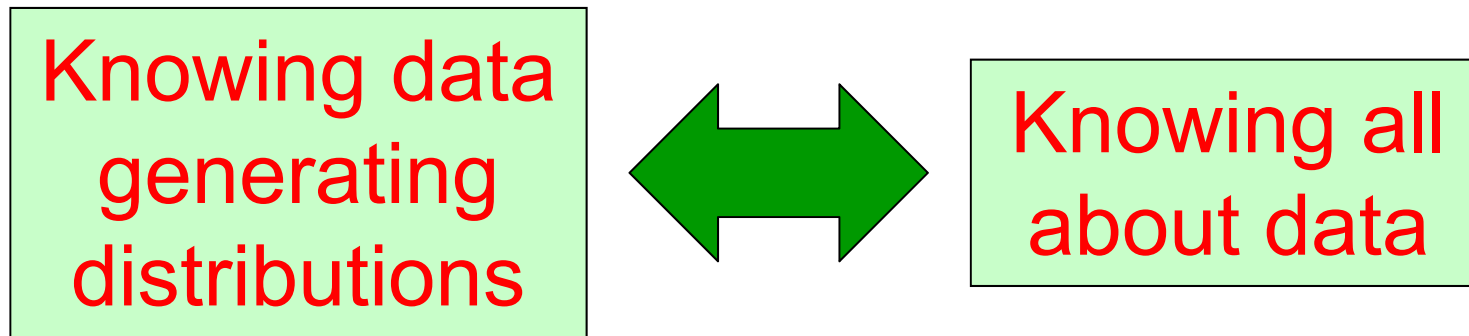
2

- **Goal:** Learn information hidden behind data
- There are **many** ML tasks:
 - Learning under non-stationarity, domain adaptation, multi-task learning, two-sample test, outlier detection, change detection in time series, independence test, feature selection, dimension reduction, independent component analysis, canonical dependency analysis, causal inference, clustering, object matching, conditional probability estimation, probabilistic classification, etc.

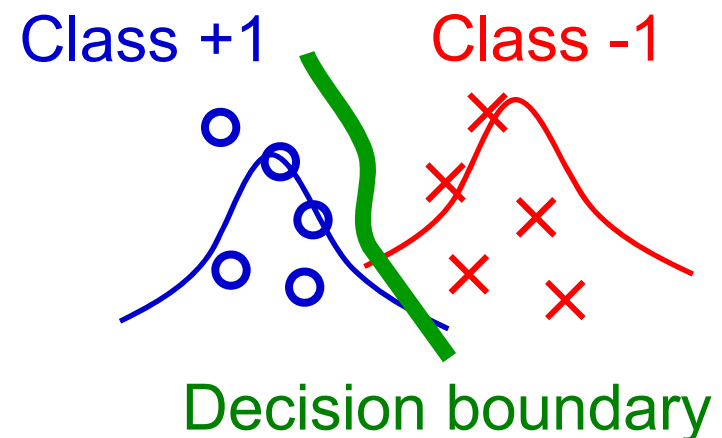
Universal Approach

3

- Learning **data-generating probability distribution** allows us to solve all ML tasks.



- Ex:** Pattern recognition is possible if data-generating probability distributions for each class are available.

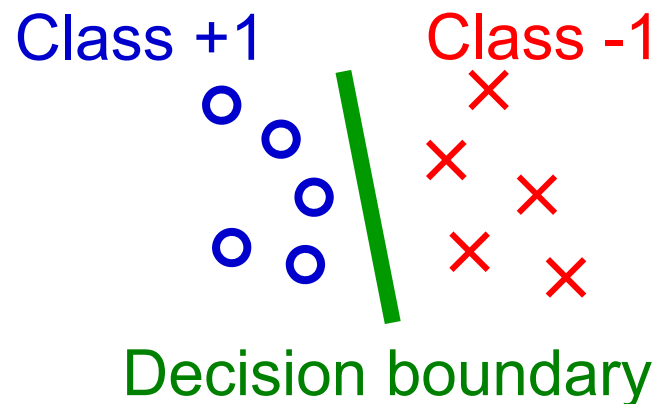


Task-Specific Approach

4

- However, estimation of probability distributions is known to be difficult.
- Avoid probability distribution estimation and solve the target task **directly**.
- **Ex: Support vector machine**
 - Directly learn a decision boundary without estimating data-generating distributions.

Cortes & Vapnik
(ML1995)



Task-Specific Approach (cont.) ⁵

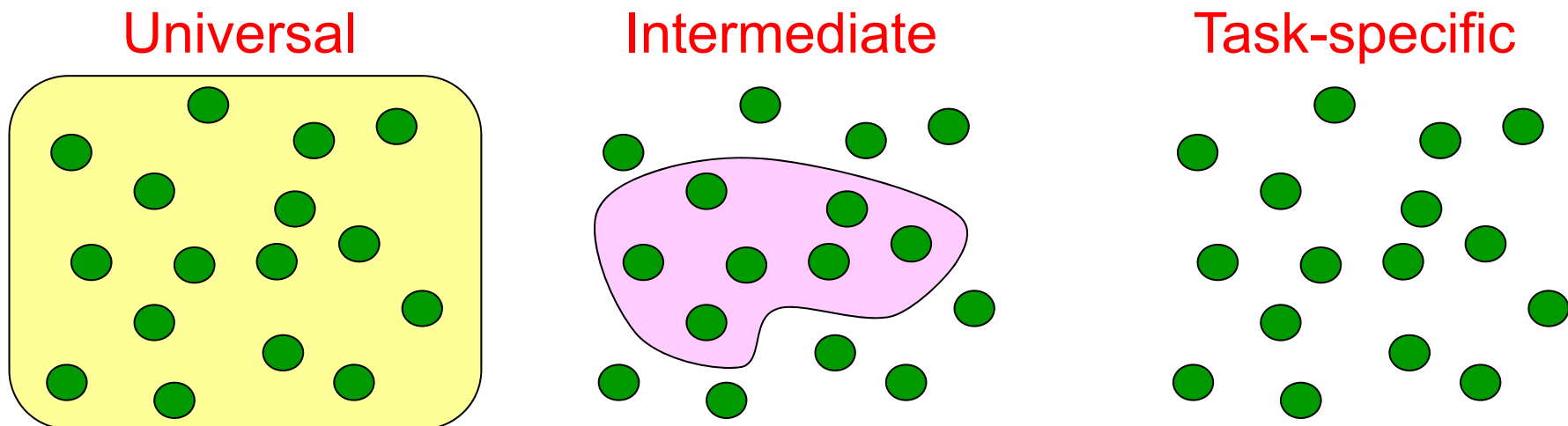
- In principle, task-specific approaches can be more accurate than the universal approach.
- However, research and development for each ML task is highly costly and cumbersome:
 - Theory, algorithms, implementation, education, etc.

Intermediate Approach

6

■ Our target: R&D for a **group** of tasks:

- Learning under non-stationarity, domain adaptation, multi-task learning, two-sample test, outlier detection, change detection in time series, independence test, feature selection, dimension reduction, independent component analysis, canonical dependency analysis, causal inference, clustering, object matching, conditional probability estimation, probabilistic classification, etc.



Density-Ratio Estimation

7

- All ML tasks listed in the previous page include **multiple** probability distributions.

$$p(\mathbf{x}), q(\mathbf{x})$$

- For solving these ML tasks, individual distributions are not necessary.
- But knowing **the density ratio** is enough.

$$r(\mathbf{x}) = \frac{p(\mathbf{x})}{q(\mathbf{x})}$$

- We **directly estimate the density ratio** without going through density estimation.

Intuitive Justification

8

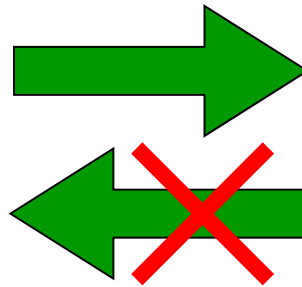
Vapnik's principle:

Vapnik (1998)

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

Knowing densities

$$p(\mathbf{x}), q(\mathbf{x})$$



Knowing ratio

$$r(\mathbf{x}) = \frac{p(\mathbf{x})}{q(\mathbf{x})}$$

- Estimating the density ratio is substantially easier than estimating densities!

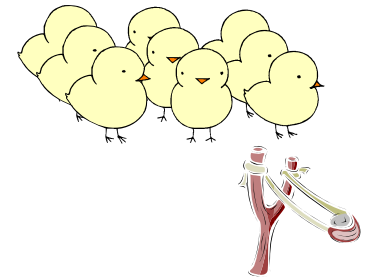
Quick Conclusions

9

- Density ratios can be accurately and efficiently estimated by simple **least-squares!**
- Many ML tasks can be solved just by LS:

- **Importance sampling:**

$$\sum_{i=1}^n \frac{p_{\text{test}}(\mathbf{x}_i)}{p_{\text{train}}(\mathbf{x}_i)} \text{loss}(\mathbf{x}_i)$$



- **Divergence estimation:**

$$\int p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x}$$

- **Mutual information estimation:**

$$\iint p(\mathbf{x}, \mathbf{y}) \log \frac{p(\mathbf{x}, \mathbf{y})}{p(\mathbf{x})p(\mathbf{y})} d\mathbf{x}d\mathbf{y}$$

- **Conditional probability estimation:**

$$p(\mathbf{y}|\mathbf{x}) = \frac{p(\mathbf{x}, \mathbf{y})}{p(\mathbf{x})}$$



Organization of This Lecture

10

1. Introduction
2. **Methods of Density Ratio Estimation**
3. Usage of Density Ratios
4. More on Density Ratio Estimation
5. Conclusions

Density Ratio Estimation: Problem Formulation

- **Goal:** Estimate the density ratio

$$r(\mathbf{x}) = \frac{p_{\text{nu}}(\mathbf{x})}{p_{\text{de}}(\mathbf{x})}$$

from data

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

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

Density Estimation Approach 12

$$r(\mathbf{x}) = \frac{p_{\text{nu}}(\mathbf{x})}{p_{\text{de}}(\mathbf{x})} \quad \begin{array}{l} \{\mathbf{x}_i^{\text{nu}}\}_{i=1}^{n_{\text{nu}}} \stackrel{i.i.d.}{\sim} p_{\text{nu}}(\mathbf{x}) \\ \{\mathbf{x}_j^{\text{de}}\}_{j=1}^{n_{\text{de}}} \stackrel{i.i.d.}{\sim} p_{\text{de}}(\mathbf{x}) \end{array}$$

■ Naïve 2-step approach:

1. Perform density estimation:

$$\hat{p}_{\text{nu}}(\mathbf{x}), \hat{p}_{\text{de}}(\mathbf{x})$$

2. Compute the ratio of estimated densities:

$$\hat{r}(\mathbf{x}) = \frac{\hat{p}_{\text{nu}}(\mathbf{x})}{\hat{p}_{\text{de}}(\mathbf{x})}$$

■ However, this works poorly because

1. is performed without regard to 2.



Organization of This Lecture

13

1. Introduction
2. Methods of Density Ratio Estimation
 - A) Density Fitting
 - B) Density-Ratio Fitting
3. Usage of Density Ratios
4. More on Density Ratio Estimation
5. Conclusions

Kullback-Leibler Importance Estimation Procedure (KLIEP)

14

Nguyen, Wainwright & Jordan (NIPS2007)
MS, Nakajima, Kashima, von Bünau & Kawanabe (NIPS2007)

- Minimize **KL divergence** from $p_{\text{nu}}(\mathbf{x})$

to $\hat{p}_{\text{nu}}(\mathbf{x}) = \hat{r}(\mathbf{x})p_{\text{de}}(\mathbf{x})$:

$$\min_{\hat{r}} \underbrace{\int p_{\text{nu}}(\mathbf{x}) \log \frac{p_{\text{nu}}(\mathbf{x})}{\hat{r}(\mathbf{x})p_{\text{de}}(\mathbf{x})} d\mathbf{x}}_{=:\text{KL}(\hat{r})}$$

$$r(\mathbf{x}) = \frac{p_{\text{nu}}(\mathbf{x})}{p_{\text{de}}(\mathbf{x})}$$

- Decomposition of KL:

$$\text{KL}(\hat{r}) = C - \int p_{\text{nu}}(\mathbf{x}) \log \hat{r}(\mathbf{x}) d\mathbf{x}$$

KLIEP: Formulation

15

■ Objective function:

$$\max_{\hat{r}} \int p_{\text{nu}}(\mathbf{x}) \log \hat{r}(\mathbf{x}) d\mathbf{x}$$

■ Constraints:

- $\hat{p}_{\text{nu}}(\mathbf{x}) = \hat{r}(\mathbf{x})p_{\text{de}}(\mathbf{x})$ is a probability density:

$$\int \hat{r}(\mathbf{x})p_{\text{de}}(\mathbf{x})d\mathbf{x} = 1 \quad \hat{r}(\mathbf{x}) \geq 0$$

■ Linear-in-parameter density-ratio model:

$$\hat{r}(\mathbf{x}) = \sum_{\ell=1}^b \alpha_{\ell} \phi_{\ell}(\mathbf{x}) = \boldsymbol{\alpha}^{\top} \boldsymbol{\phi}(\mathbf{x})$$

$\phi_{\ell}(\mathbf{x}) \geq 0$
(ex. Gauss kernel)

KLIEP: Algorithm

16

- Approximate expectations by sample averages:

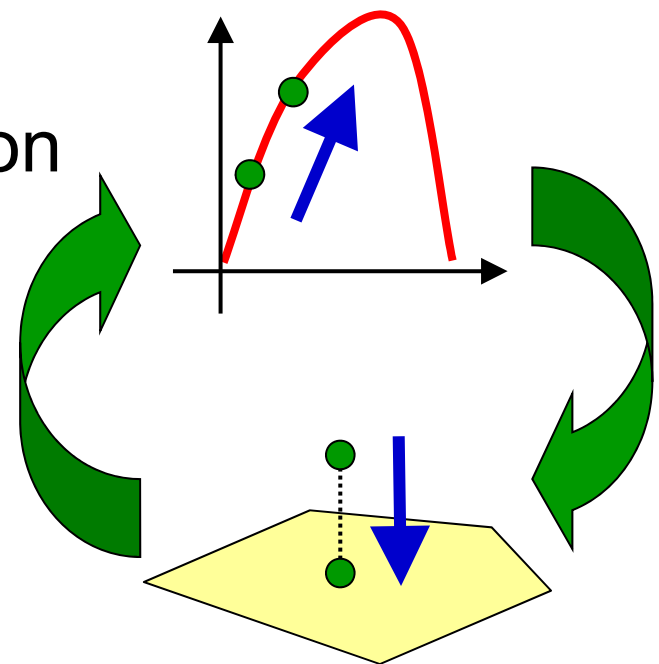
$$\max_{\alpha} \sum_{i=1}^{n_{\text{nu}}} \log(\alpha^{\top} \phi(\mathbf{x}_i^{\text{nu}})) \quad \text{subject to} \quad \frac{1}{n_{\text{de}}} \sum_{j=1}^{n_{\text{de}}} \alpha^{\top} \phi(\mathbf{x}_j^{\text{de}}) = 1 \text{ and } \alpha \geq \mathbf{0}$$

- This is **convex optimization**, so repeating

- Gradient ascent
- Projection onto the feasible region

leads to the **global solution**.

- The global solution is **sparse**!



KLIEP: Convergence Properties¹⁷

Nguyen, Wainwright & Jordan (IEEE-IT2010)

MS, Suzuki, Nakajima, Kashima, von Bünau & Kawanabe (AISM2008)

■ **Parametric case:** $\hat{r}(\mathbf{x}) = \sum_{\ell=1}^b \alpha_{\ell} \phi_{\ell}(\mathbf{x})$

- Learned parameter converge to the optimal value with order $n^{-\frac{1}{2}}$, which is the **optimal rate**.

$$n = \min(n_{\text{nu}}, n_{\text{de}})$$

■ **Non-parametric case:** $\hat{r}(\mathbf{x}) = \sum_{\ell=1}^{n_{\text{nu}}} \alpha_{\ell} K(\mathbf{x}, \mathbf{x}_{\ell}^{\text{nu}})$

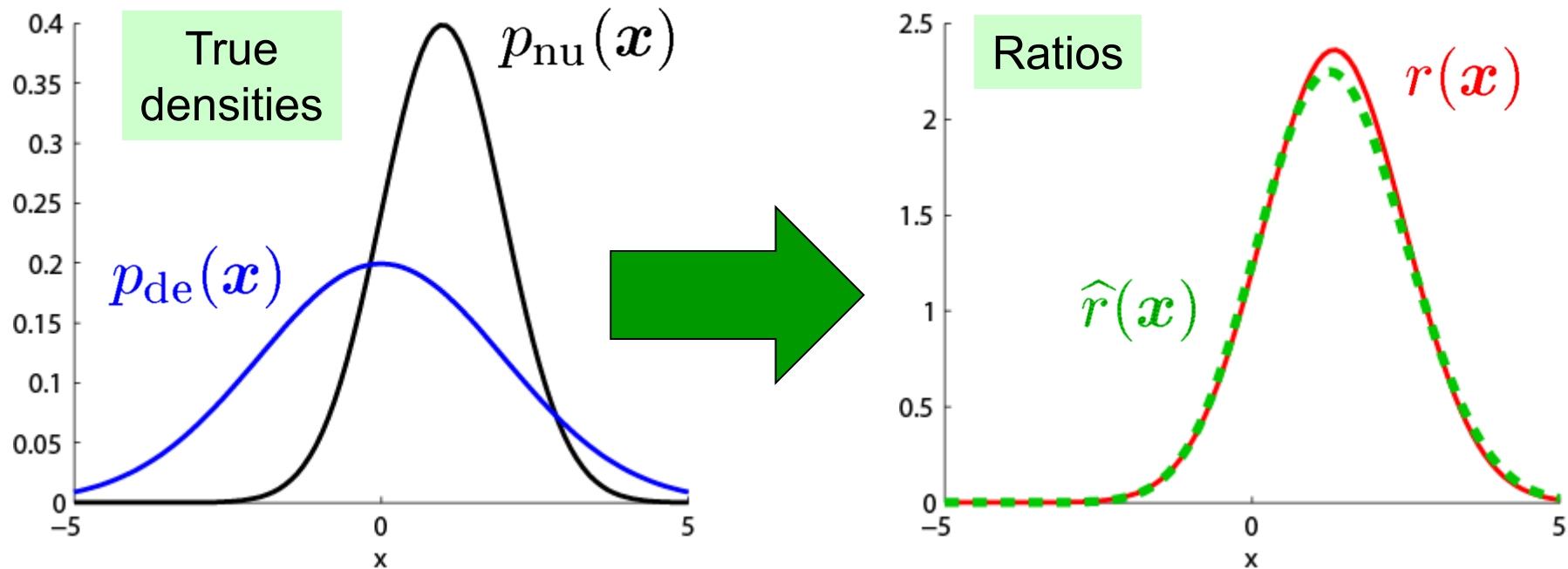
- Learned function converges to the optimal function with order $n^{-\frac{1}{2+\gamma}}$, which is the **optimal rate**.

$0 < \gamma < 2$: Complexity of the function class related to the covering number or bracketing entropy

KLIEP: Numerical Example

18

$$n_{\text{nu}} = n_{\text{de}} = 200$$



- Gaussian width can be determined by cross-validation with respect to KL.

$$\int p_{\text{nu}}(\mathbf{x}) \log \hat{r}(\mathbf{x}) d\mathbf{x}$$

KLIEP: Summary

19

- Density estimation is not involved.
- **Cross-validation** is available for kernel parameter selection.
- Variations for **various models** exist:
 - Log-linear, Gaussian mixture, PCA mixture, etc.
- An **unconstrained variant** corresponds to maximizing a lower-bound of KL divergence.

$$\int p_{\text{nu}}(\mathbf{x}) \log \frac{p_{\text{nu}}(\mathbf{x})}{p_{\text{de}}(\mathbf{x})} d\mathbf{x}$$

Nguyen, Wainwright
& Jordan (NIPS2007)



Organization of This Lecture

20

1. Introduction
2. Methods of Density Ratio Estimation
 - A) Density Fitting
 - B) Density-Ratio Fitting
3. Usage of Density Ratios
4. More on Density Ratio Estimation
5. Conclusions

Least-Squares Importance Fitting (LSIF)

21

Kanamori, Hido & MS
(NIPS2008)

- Minimize squared-loss (SQ):

$$r(\mathbf{x}) = \frac{p_{\text{nu}}(\mathbf{x})}{p_{\text{de}}(\mathbf{x})}$$

$$\min_{\hat{r}} \underbrace{\int \left(\hat{r}(\mathbf{x}) - r(\mathbf{x}) \right)^2 p_{\text{de}}(\mathbf{x}) d\mathbf{x}}_{=:\text{SQ}(\hat{r})}$$

- Decomposition and approximation of SQ:

$$\text{SQ}(\hat{r}) = \int \left(\hat{r}(\mathbf{x}) \right)^2 p_{\text{de}}(\mathbf{x}) d\mathbf{x} - 2 \int \hat{r}(\mathbf{x}) p_{\text{nu}}(\mathbf{x}) d\mathbf{x} + C$$

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

Constrained LSIF Formulation ²²

- Linear (or kernel) density-ratio model:

$$\hat{r}(\mathbf{x}) = \sum_{\ell=1}^b \alpha_{\ell} \phi_{\ell}(\mathbf{x}) = \boldsymbol{\alpha}^{\top} \boldsymbol{\phi}(\mathbf{x})$$

- **Constrained LSIF (cLSIF):**

- Non-negativity constraint with ℓ_1 -regularizer

$$\min_{\boldsymbol{\alpha}} \left[\frac{1}{2} \boldsymbol{\alpha}^{\top} \widehat{\mathbf{H}} \boldsymbol{\alpha} - \widehat{\mathbf{h}}^{\top} \boldsymbol{\alpha} + \lambda \boldsymbol{\alpha}^{\top} \mathbf{1} \right]$$

subject to $\boldsymbol{\alpha} \geq \mathbf{0}$

- A convex quadratic program with **sparse solution**.

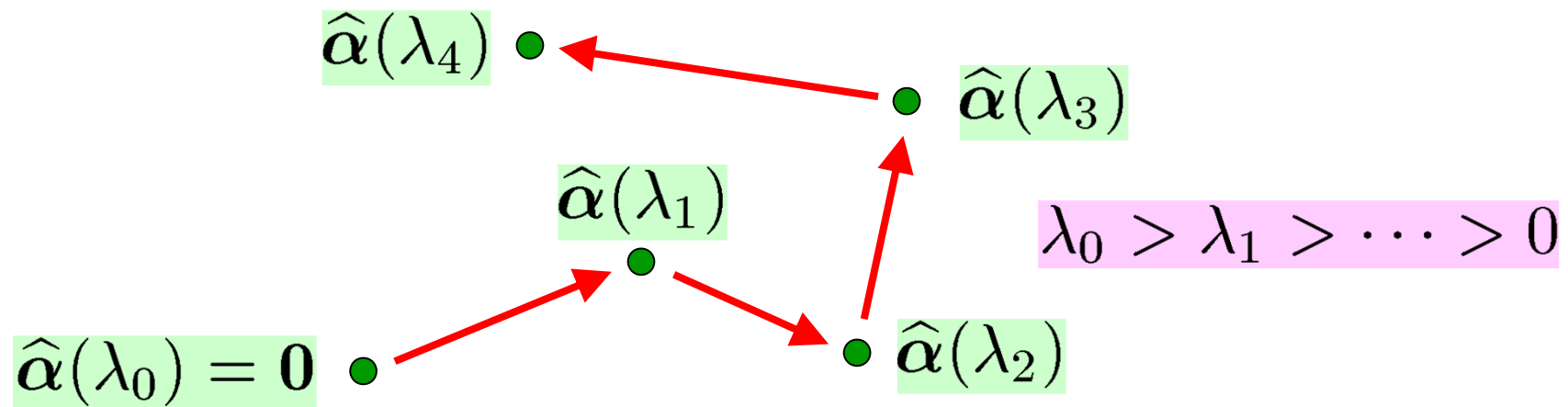
$$\widehat{\mathbf{H}} = \frac{1}{n_{\text{de}}} \sum_{j=1}^{n_{\text{de}}} \boldsymbol{\phi}(\mathbf{x}_j^{\text{de}}) \boldsymbol{\phi}(\mathbf{x}_j^{\text{de}})^{\top}$$

$$\widehat{\mathbf{h}} = \frac{1}{n_{\text{nu}}} \sum_{i=1}^{n_{\text{nu}}} \boldsymbol{\phi}(\mathbf{x}_i^{\text{nu}})$$

cLSIF: Regularization Path Tracking ²³

$$\min_{\alpha} \left[\frac{1}{2} \alpha^{\top} \widehat{H} \alpha - \widehat{h}^{\top} \alpha + \lambda \alpha^{\top} \mathbf{1} \right] \quad \text{subject to } \alpha \geq \mathbf{0}$$

- The solution path is **piece-wise linear** with respect to the regularization parameter λ .



- Solutions for all λ can be computed efficiently **without QP solvers!**

Unconstrained LSIF Formulation²⁴

$$\hat{r}(\mathbf{x}) = \sum_{\ell=1}^b \alpha_{\ell} \phi_{\ell}(\mathbf{x}) = \boldsymbol{\alpha}^{\top} \boldsymbol{\phi}(\mathbf{x})$$

■ Unconstrained LSIF (uLSIF):

- **uLSIF**: No constraint with ℓ_2 -regularizer

$$\min_{\boldsymbol{\alpha}} \left[\frac{1}{2} \boldsymbol{\alpha}^{\top} \widehat{\mathbf{H}} \boldsymbol{\alpha} - \hat{\mathbf{h}}^{\top} \boldsymbol{\alpha} + \frac{\lambda}{2} \boldsymbol{\alpha}^{\top} \boldsymbol{\alpha} \right]$$

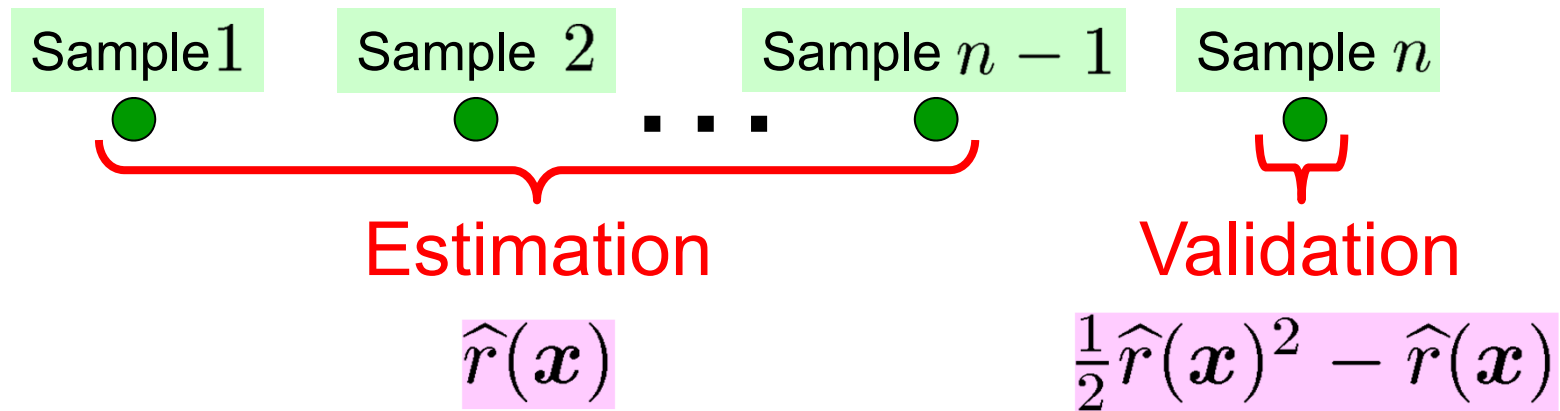
- Analytic solution is available: $(\widehat{\mathbf{H}} + \lambda \mathbf{I})^{-1} \hat{\mathbf{h}}$

$$\widehat{\mathbf{H}} = \frac{1}{n_{\text{de}}} \sum_{j=1}^{n_{\text{de}}} \boldsymbol{\phi}(\mathbf{x}_j^{\text{de}}) \boldsymbol{\phi}(\mathbf{x}_j^{\text{de}})^{\top}$$

$$\hat{\mathbf{h}} = \frac{1}{n_{\text{nu}}} \sum_{i=1}^{n_{\text{nu}}} \boldsymbol{\phi}(\mathbf{x}_i^{\text{nu}})$$

uLSIF: Analytic LOOCV Score ²⁵

■ Leave-one-out cross-validation (LOOCV):



- LOOCV generally requires n repetitions.
- However, it can be **analytically** computed for uLSIF (Sherman-Woodbury-Morrison formula).
- Computation time including model selection is **significantly reduced**.

uLSIF: Theoretical Properties 26

■ Parametric convergence:

- Learned parameter converge to the optimal value with order $n^{-\frac{1}{2}}$, which is the **optimal rate**.

$$n = \min(n_{\text{nu}}, n_{\text{de}}) \quad \text{Kanamori, Hido \& MS (JMLR2009)}$$

■ Non-parametric convergence:

- Learned function converges to the optimal function with order $n^{-\frac{1}{2+\gamma}}$ (depending on the bracketing entropy), which is the **optimal rate**.

$$0 < \gamma < 2 \quad \text{Kanamori, Suzuki \& MS (MLJ2012)}$$

■ Non-parametric numerical stability:

- uLSIF has the **smallest condition number** among a class of density ratio estimators.

$$\text{Kanamori, Suzuki \& MS (MLJ2012)}$$

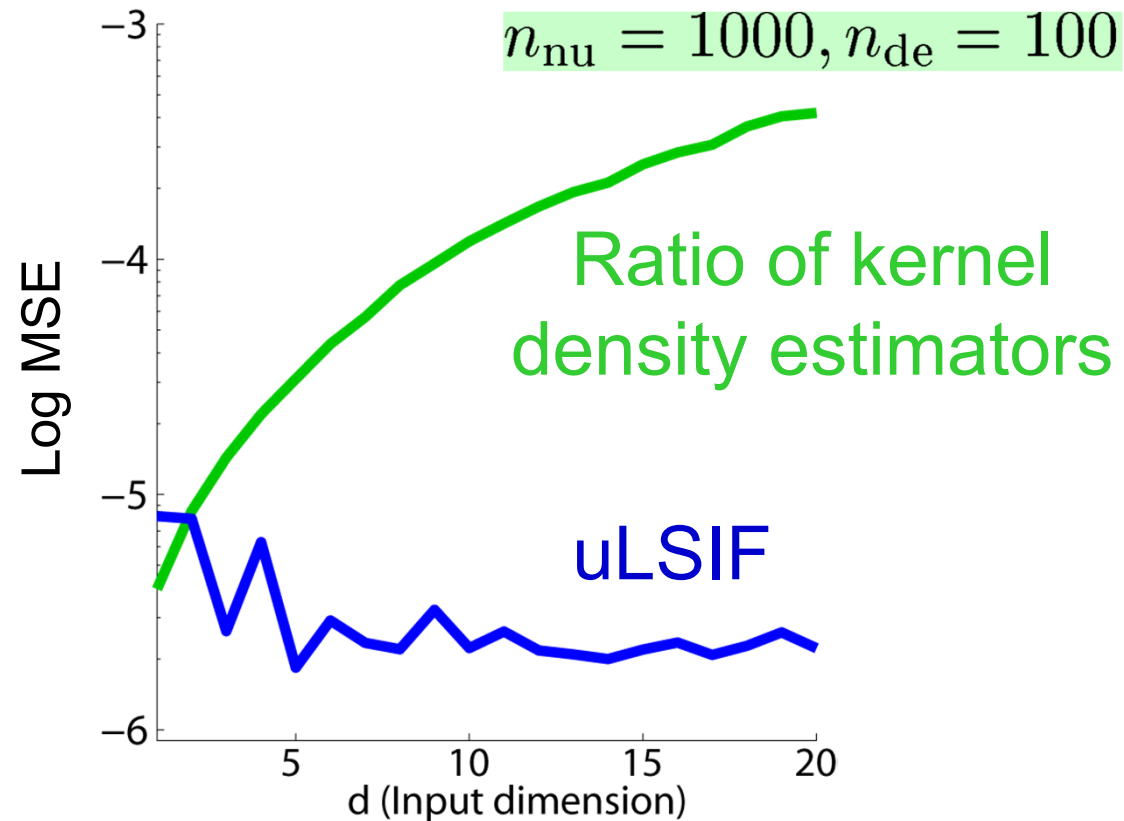
uLSIF: Numerical Example

27

$$p_{\text{nu}}(\mathbf{x}) = N(\mathbf{x}; (0, 0, \dots, 0)^\top, \mathbf{I}_d)$$

$$p_{\text{de}}(\mathbf{x}) = N(\mathbf{x}; (1, 0, \dots, 0)^\top, \mathbf{I}_d)$$

$$r(\mathbf{x}) = \frac{p_{\text{nu}}(\mathbf{x})}{p_{\text{de}}(\mathbf{x})}$$



LSIF: Summary

28

- LS formulation is computationally efficient:
 - **cLSIF**: Regularization path tracking
 - **uLSIF**: Analytic solution and LOOCV
- Gives an accurate approximator of **Pearson (PE) divergence** (an f -divergence):

$$\int p_{\text{de}}(\mathbf{x}) \left(\frac{p_{\text{nu}}(\mathbf{x})}{q_{\text{de}}(\mathbf{x})} - 1 \right)^2 d\mathbf{x}$$

- Analytic solution of uLSIF allows us to compute the **derivative** of PE divergence approximator:
 - Useful in dimension reduction, independent component analysis, causal inference etc.



Organization of This Lecture

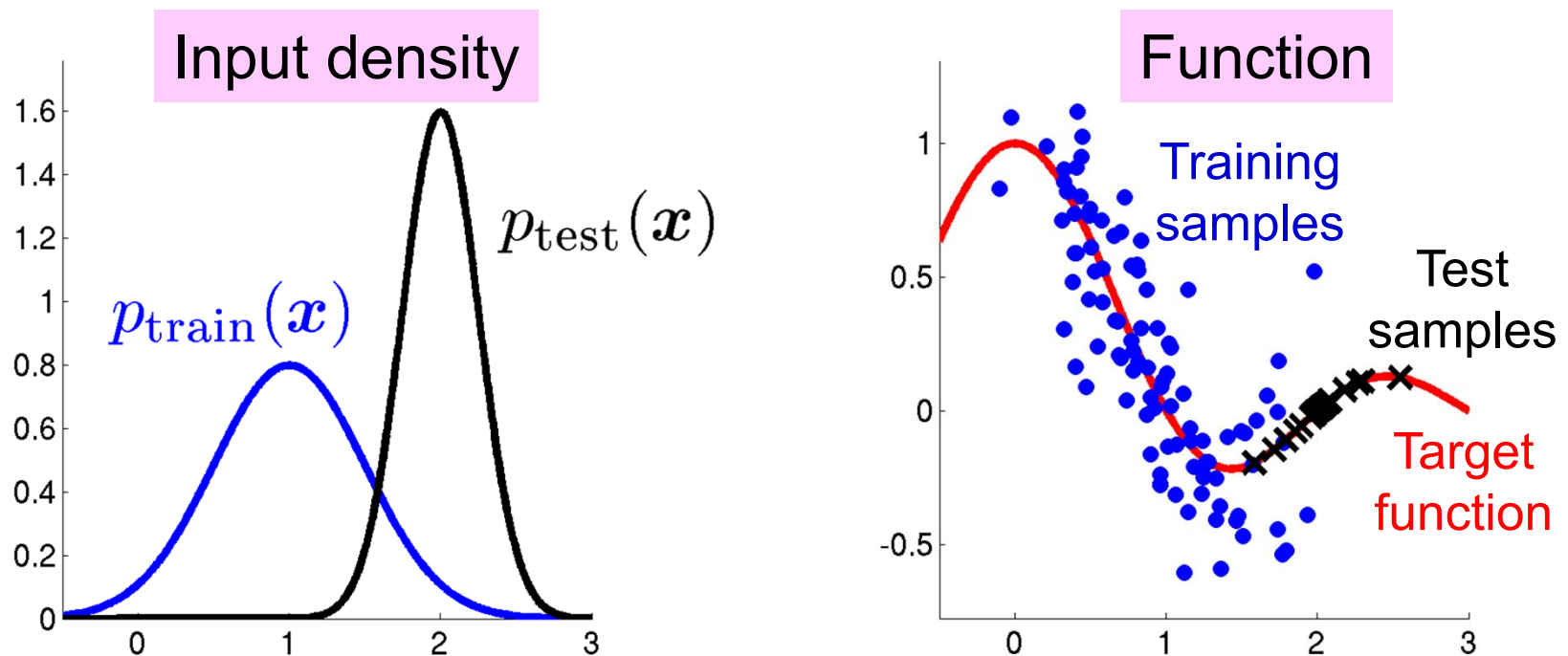
29

1. Introduction
2. Methods of Density Ratio Estimation
3. Usage of Density Ratios
 - A) Importance sampling
 - B) Distribution comparison
 - C) Mutual information estimation
 - D) Conditional probability estimation
4. More on Density Ratio Estimation
5. Conclusions

Learning under Covariate Shift³⁰

■ Covariate shift: Shimodaira (JSPI2000)

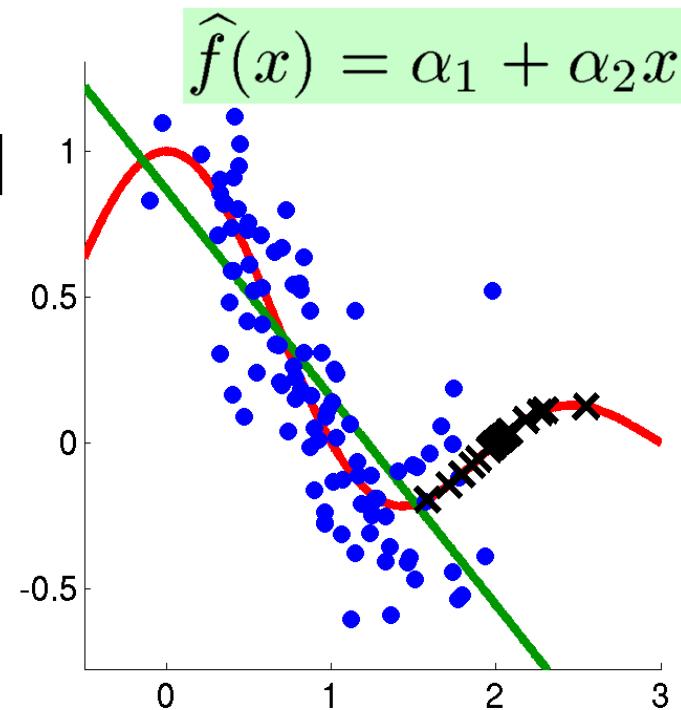
- Training/test input distributions are different, but target function remains unchanged.
- (Weak) extrapolation.



Ordinary Least-Squares (OLS)³¹

$$\min_{\alpha} \left[\sum_{i=1}^n \left(\hat{f}(x_i) - y_i \right)^2 \right]$$

- In standard setting, OLS is **consistent**, i.e., the learned function converges to the best solution when $n \rightarrow \infty$.
- Under covariate shift, OLS is **no longer consistent**.



Law of Large Numbers

- Sample average converges to the population mean:

$$\frac{1}{n} \sum_{i=1}^n \text{loss}(\mathbf{x}_i) \longrightarrow \int \text{loss}(\mathbf{x}) p_{\text{train}}(\mathbf{x}) d\mathbf{x}$$

$$\mathbf{x}_i \stackrel{i.i.d.}{\sim} p_{\text{train}}(\mathbf{x})$$

- We want to estimate the expectation over **test input points** only using **training input points** $\{\mathbf{x}_i\}_{i=1}^n$.

$$\int \text{loss}(\mathbf{x}) p_{\text{test}}(\mathbf{x}) d\mathbf{x}$$

Importance Weighting

- **Importance**: Ratio of test and training input densities

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

- **Importance-weighted average**:

$$\frac{1}{n} \sum_{i=1}^n \frac{p_{test}(\mathbf{x}_i)}{p_{train}(\mathbf{x}_i)} \text{loss}(\mathbf{x}_i)$$

$$\mathbf{x}_i \stackrel{i.i.d.}{\sim} p_{train}(\mathbf{x})$$

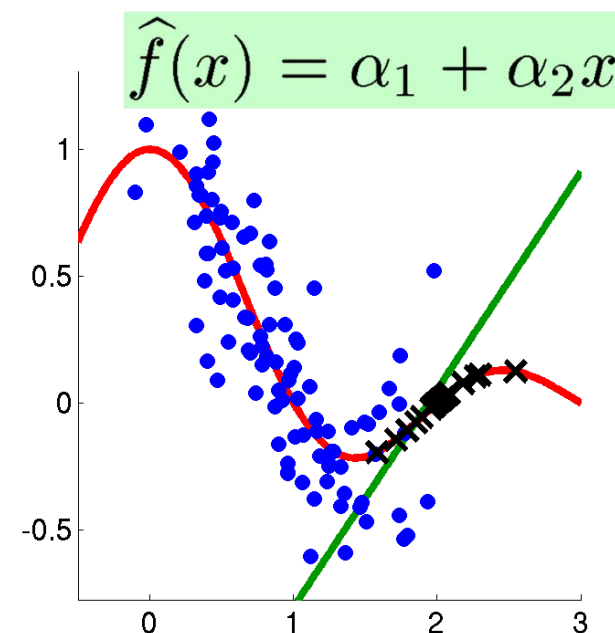
$$\rightarrow \int \frac{p_{test}(\mathbf{x})}{p_{train}(\mathbf{x})} \text{loss}(\mathbf{x}) p_{train}(\mathbf{x}) d\mathbf{x}$$

$$= \int \text{loss}(\mathbf{x}) p_{test}(\mathbf{x}) d\mathbf{x}$$

Importance-Weighted Least-Squares

$$\min_{\alpha} \left[\sum_{i=1}^n \frac{p_{test}(\mathbf{x}_i)}{p_{train}(\mathbf{x}_i)} \left(\hat{f}(\mathbf{x}_i) - y_i \right)^2 \right]$$

- IWLS is **consistent even under covariate shift**.
- The idea is applicable to **any likelihood-based methods!**
 - Support vector machine, logistic regression, conditional random field, etc.



Model Selection

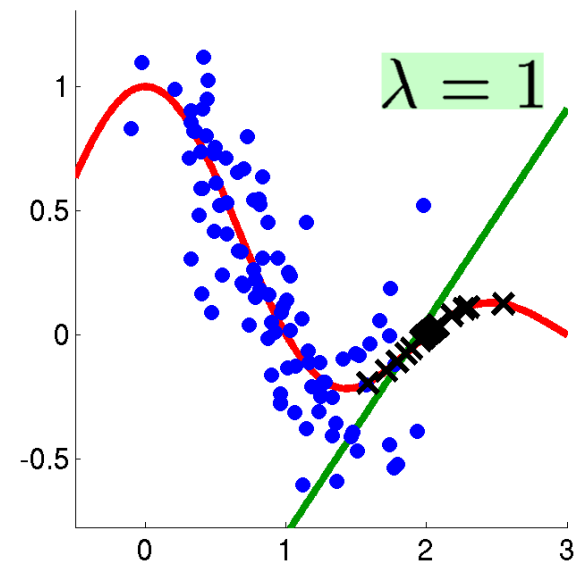
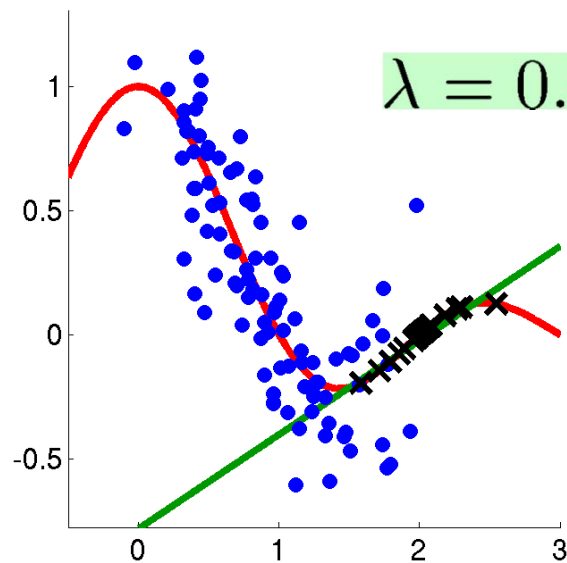
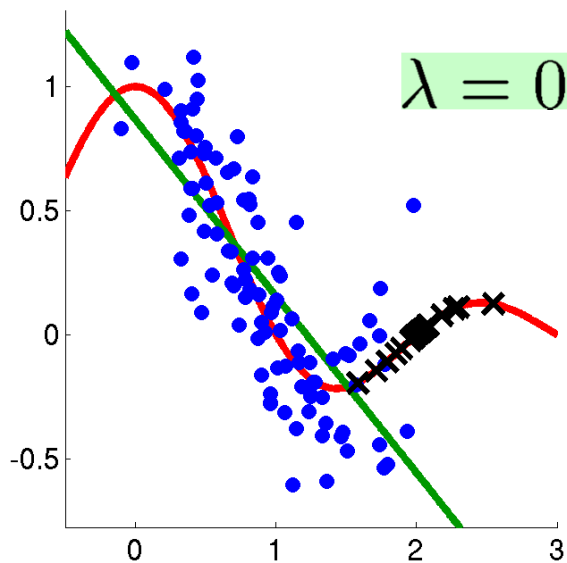
35

- Controlling **bias-variance trade-off** is important.
 - **No weighting**: low-variance but high-bias
 - **Importance weighting**: low-bias but high-variance

■ “Flattened”-IWLS:

Shimodaira (JSPI2000)

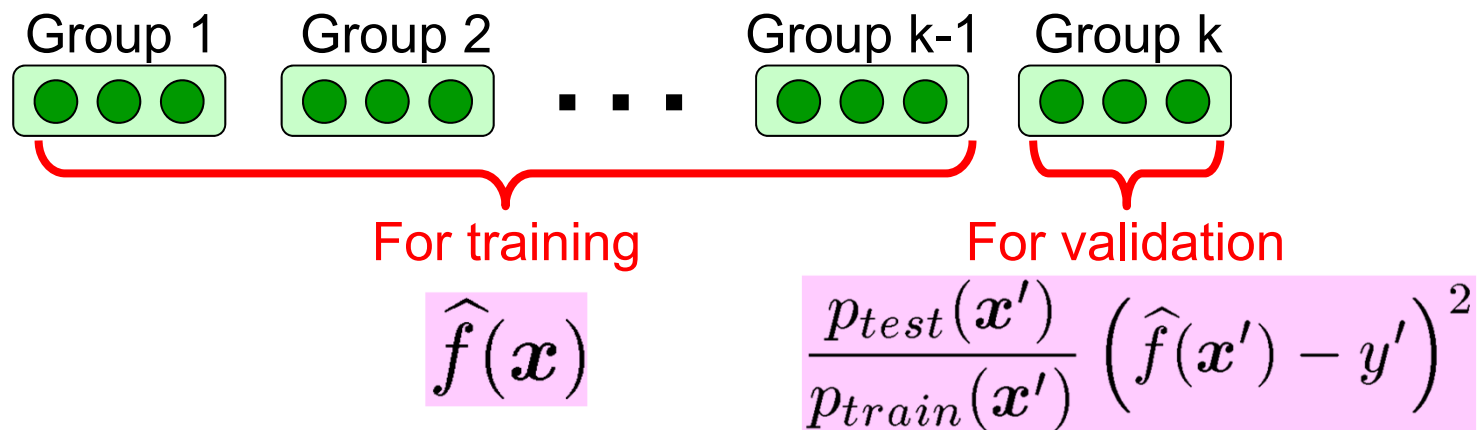
$$\min_{\alpha} \left[\sum_{i=1}^n \left(\frac{p_{test}(\mathbf{x}_i)}{p_{train}(\mathbf{x}_i)} \right)^{\lambda} \left(\hat{f}(\mathbf{x}_i) - y_i \right)^2 \right]$$



Model Selection

■ Importance weighting also plays a central role for unbiased model selection:

- Akaike information criterion (regular models)
Shimodaira (JSPI2000)
- Subspace information criterion (linear models)
MS & Müller (Stat&Dec.2005)
- Cross-validation (arbitrary models)
MS, Krauledat & Müller (JMLR2007)



Experiments: Speaker Identification³⁷

Yamada, MS & Matsui (SigPro2010)

- NTT Japanese speech dataset. Matsui & Furui (ICASSP1993)
- Text-independent speaker identification accuracy for 10 male speakers.
- Kernel logistic regression (KLR) with sequence kernel.

Training data	Speech length	IWKLR+IWCV+KLIEP	KLR+CV
9 months before	1.5 [sec]	91.0 %	88.2 %
	3.0 [sec]	95.0 %	92.9 %
	4.5 [sec]	97.7 %	96.1 %
6 months before	1.5 [sec]	91.0 %	87.7 %
	3.0 [sec]	95.3 %	91.1 %
	4.5 [sec]	97.4 %	93.4 %
3 months before	1.5 [sec]	94.8 %	91.7 %
	3.0 [sec]	97.9 %	96.3 %
	4.5 [sec]	98.8 %	98.3 %



Experiments: Text Segmentation³⁸

Tsuboi, Kashima, Hido, Bickel & MS (JIP2009)

こんな失敗はご愛敬だよ。
→ こんな／失敗／は／ご／愛敬／だ／よ／。

- Japanese word segmentation dataset.

Tsuboi, Kashima, Mori, Oda & Matsumoto (COLING2008)

- Adaptation from daily conversation to medical domain.
- Segmentation by conditional random field (CRF).

	IWCRF+IWCV +KLIEP	CRF+CV	CRF+CV (use additional test labels)
F-measure (larger is better)	94.46	92.30	94.43

Semi-supervised adaptation with importance weighting is comparable to supervised adaptation!

Other Applications

39

■ Age prediction from faces:

- Illumination change

Ueki, MS & Ihara (ICPR2010)

■ Brain-computer interface:

- Mental condition change

MS, Krauledat & Müller (JMLR2007)

Li, Kambara, Koike & MS (IEEE-TBME2010)

■ Robot control:

- Efficient sample reuse

Hachiya, Akiyama, MS & Peters (NN2009)

Hachiya, Peters & MS (NeCo2011)



Organization of This Lecture

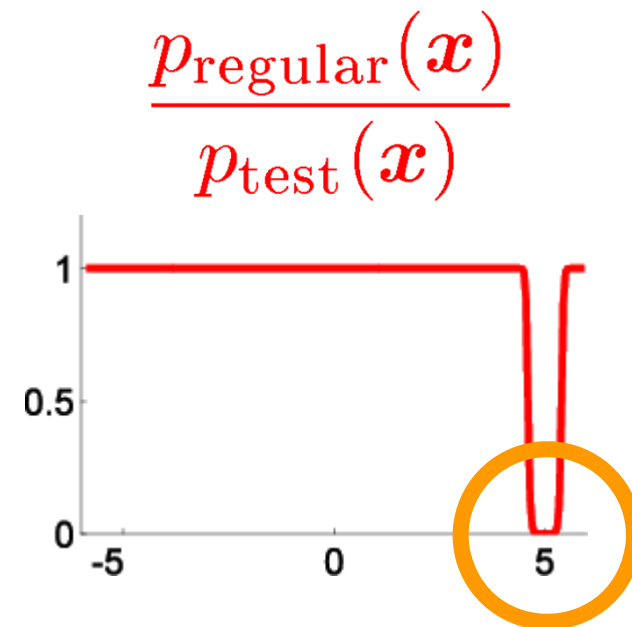
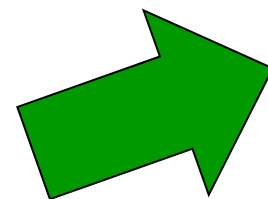
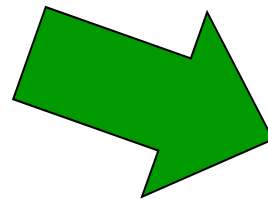
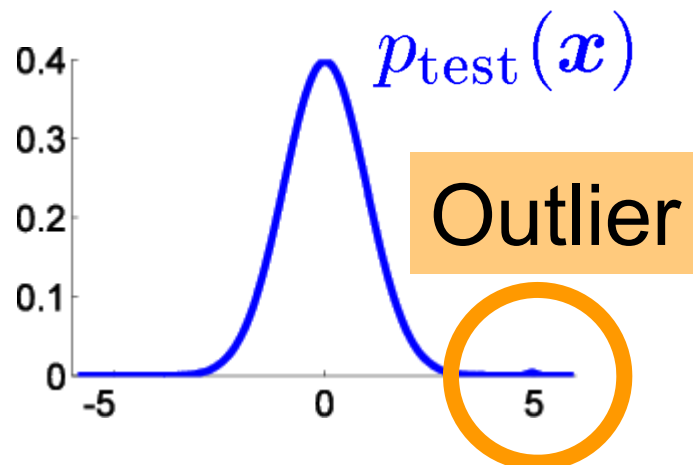
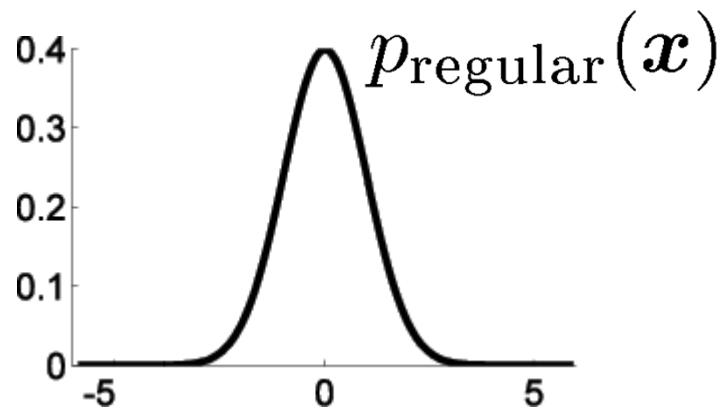
40

1. Introduction
2. Methods of Density Ratio Estimation
3. Usage of Density Ratios
 - A) Importance sampling
 - B) Distribution comparison
 - C) Mutual information estimation
 - D) Conditional probability estimation
4. More on Density Ratio Estimation
5. Conclusions

Inlier-Based Outlier Detection ⁴¹

Hido, Tsuboi, Kashima, MS & Kanamori (ICDM2008, KAIS2011)
Smola, Song & Teo (AISTATS2009)

- **Goal:** Given a set of inlier samples, find outliers in a test set (if exist)



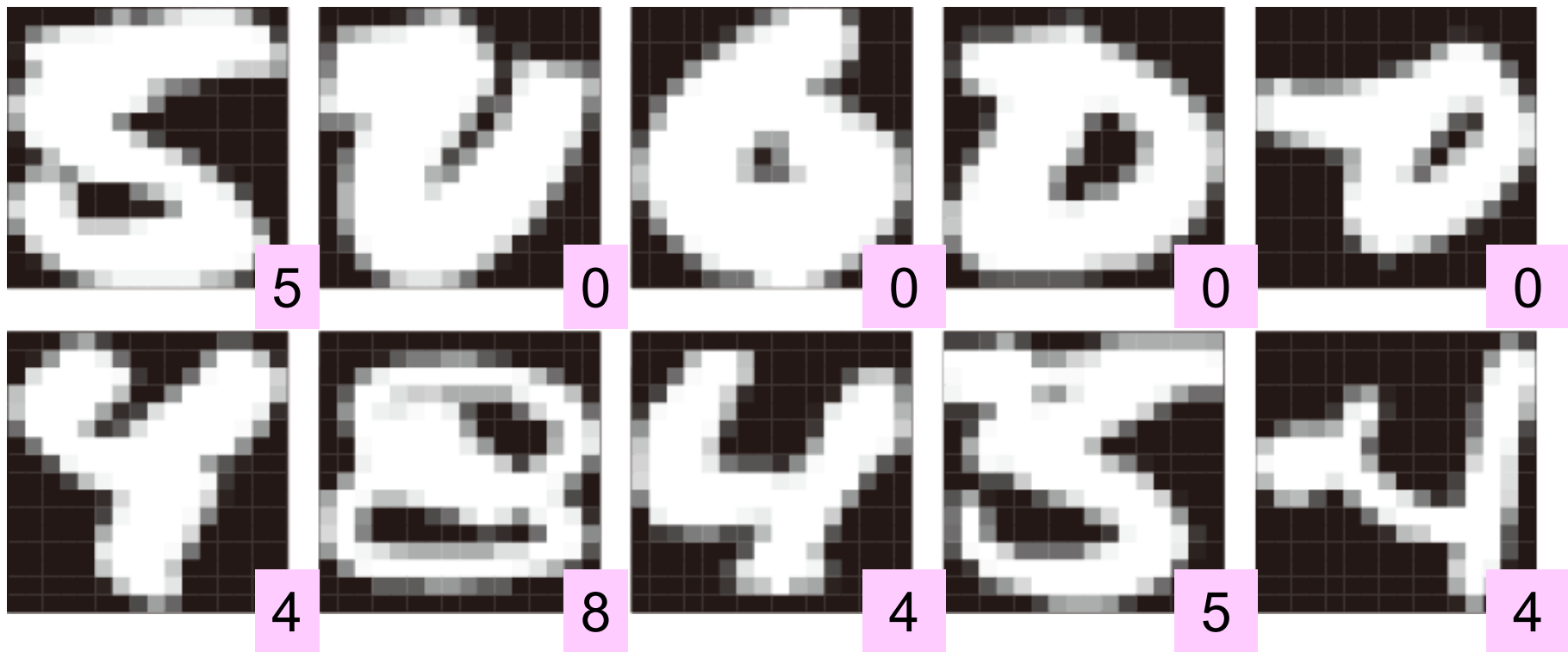
Tuning parameters can be optimized in terms of ratio approximation error

Experiments

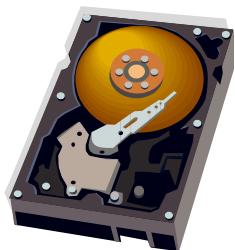
42

Hido, Tsuboi, Kashima, MS & Kanamori (ICDM2008, KAIS2011)

- Top10 outliers in the USPS test dataset found based on the USPS training dataset.



Most of them are not readable even by human.



Failure Prediction in Hard-Disk Drives

- Self-Monitoring And Reporting Technology (SMART):
Murray, Hughes & Kreutz-Delgado (JMLR 2005)

	Least-squares density ratio	One-class SVM	Local outlier factor	
			#NN=5	#NN=30
AUC (larger is better)	0.881	0.843	0.847	0.924
Comp. time	1	26.98	65.31	

- LOF works well, given #NN is set appropriately. But there is no objective model selection method.
- Density ratio method can use cross-validation for model selection, and is computationally efficient.

OSVM: Schölkopf, Platt, Shawe-Taylor, Smola & Williamson (NeCo2001)

LOF: Breunig, Kriegel, Ng & Sander (SIGMOD2000)

Other Applications

44

- Steel plant diagnosis Hirata, Kawahara & MS (Patent2011)
- Printer roller quality control Takimoto, Matsugu & MS (DMSS2009)
- Loan customer inspection Hido, Tsuboi, Kashima, MS & Kanamori (KAIS2011)
- Sleep therapy Kawahara & MS (SADM2012)

Divergence Estimation

45

Nguyen, Wainwright & Jordan (IEEE-IT2010)
MS, Suzuki, Ito, Kanamori & Kimura (NN2011)

- **Goal:** Estimate a divergence functional from

$$\{\mathbf{x}_i\}_{i=1}^n \stackrel{i.i.d.}{\sim} p(\mathbf{x}) \quad \{\mathbf{x}'_j\}_{j=1}^{n'} \stackrel{i.i.d.}{\sim} p'(\mathbf{x})$$

- **Kullback-Leibler divergence:** $\int p(\mathbf{x}) \log \frac{p(\mathbf{x})}{p'(\mathbf{x})} d\mathbf{x}$

- **Pearson divergence:** $\int p'(\mathbf{x}) \left(\frac{p(\mathbf{x})}{p'(\mathbf{x})} - 1 \right)^2 d\mathbf{x}$
(an f -divergence)

- Use density ratio estimation: $r(\mathbf{x}) = \frac{p(\mathbf{x})}{p'(\mathbf{x})}$

Real-World Applications

46

■ Regions-of-interest detection in images:

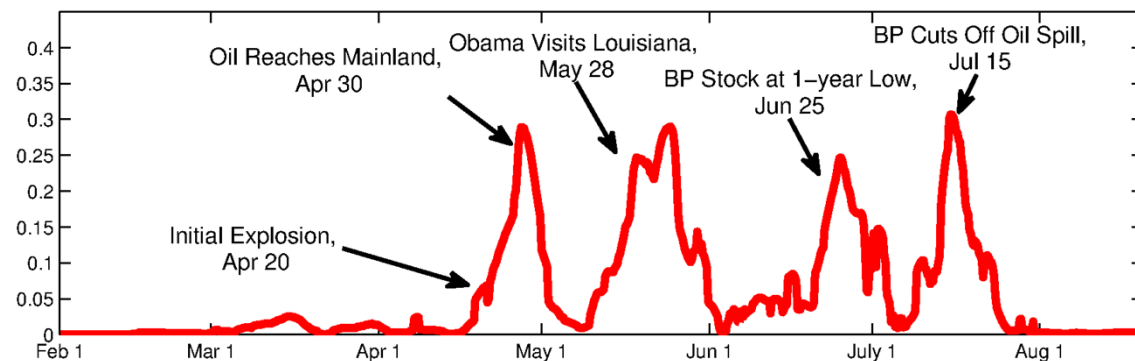
Yamanaka, Matsugu & MS
(IEEJ2011)

■ Event detection in movies:

Matsugu, Yamanaka & MS
(VECTaR2011)

■ Event detection from Twitter data:

Liu, Yamada, Collier
& MS (arXiv2012)





Organization of This Lecture

47

1. Introduction
2. Methods of Density Ratio Estimation
3. Usage of Density Ratios
 - A) Importance sampling
 - B) Distribution comparison
 - C) Mutual information estimation
 - D) Conditional probability estimation
4. More on Density Ratio Estimation
5. Conclusions

Mutual Information Estimation ⁴⁸

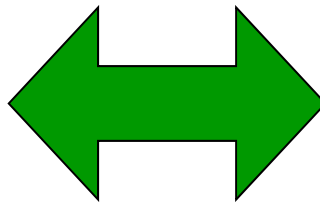
Suzuki, MS, Sese & Kanamori (FSDM2008)

- **Mutual information (MI):** Shannon (1948)

$$\text{MI} = \iint p(\mathbf{x}, \mathbf{y}) \log \frac{p(\mathbf{x}, \mathbf{y})}{p(\mathbf{x})p(\mathbf{y})} d\mathbf{x}d\mathbf{y}$$

- MI works as an **independence measure**:

$$\text{MI} = 0$$



x and y are statistically independent

- Use KL-based density ratio estimation (KLIEP):

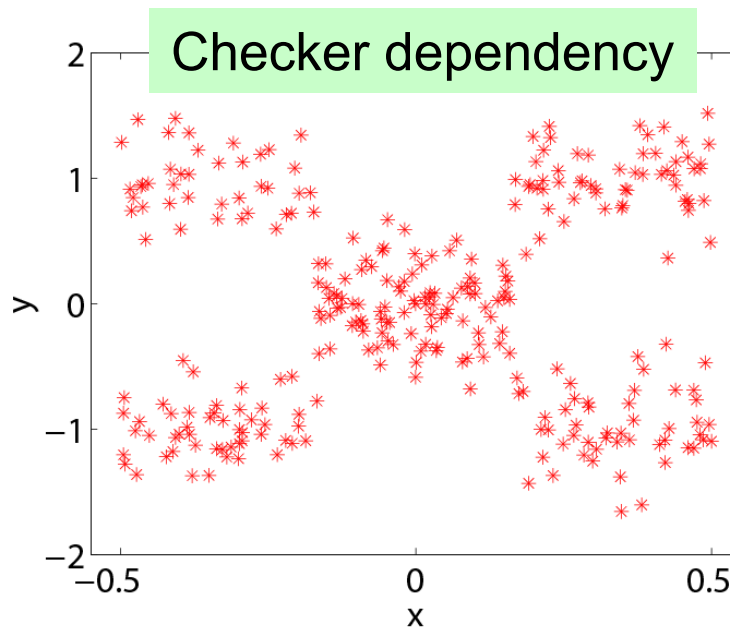
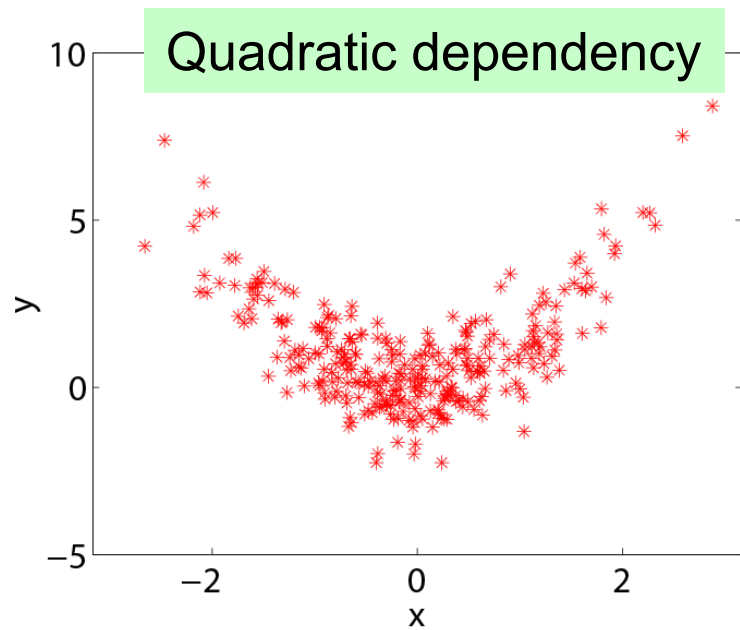
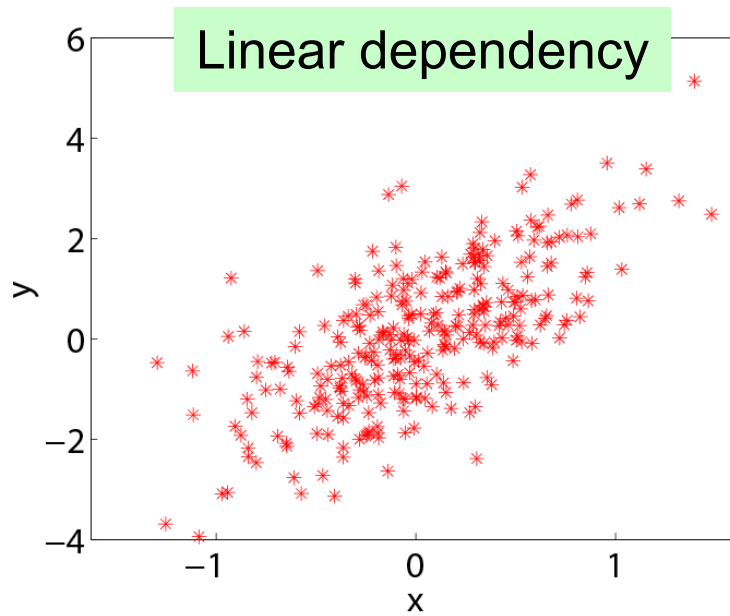
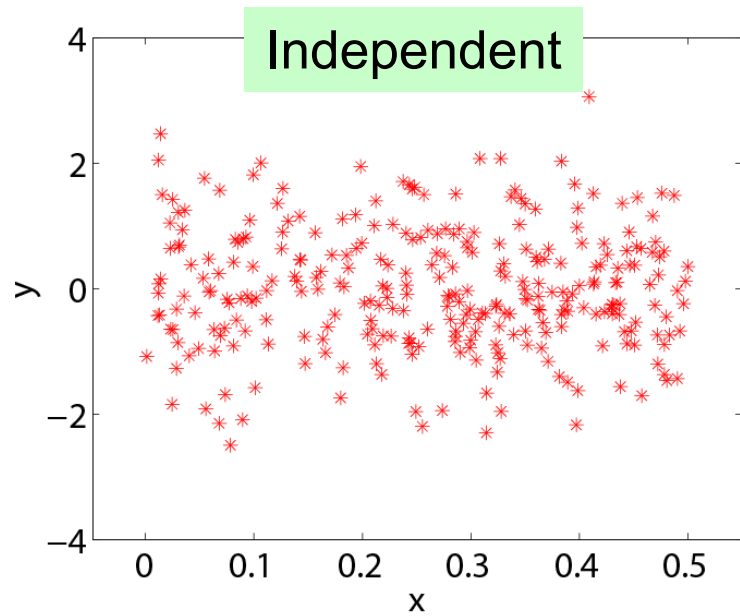
$$r(\mathbf{x}, \mathbf{y}) = \frac{p(\mathbf{x}, \mathbf{y})}{p(\mathbf{x})p(\mathbf{y})}$$

Experiments: Methods Compared⁴⁹

- KL-based density ratio method.
- Kernel density estimation (KDE).
- K-nearest neighbor density estimation (KNN). Kraskov, Stögbauer & Grassberger (PRE2004)
 - The number of NNs is a tuning parameter.
- Edgeworth expansion density estimation (EDGE). van Hulle (NeCo2005)

Datasets for Evaluation

50

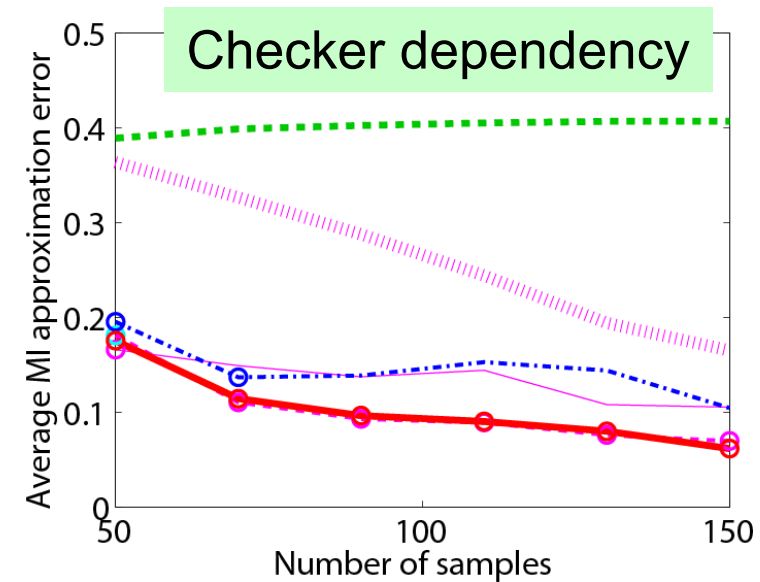
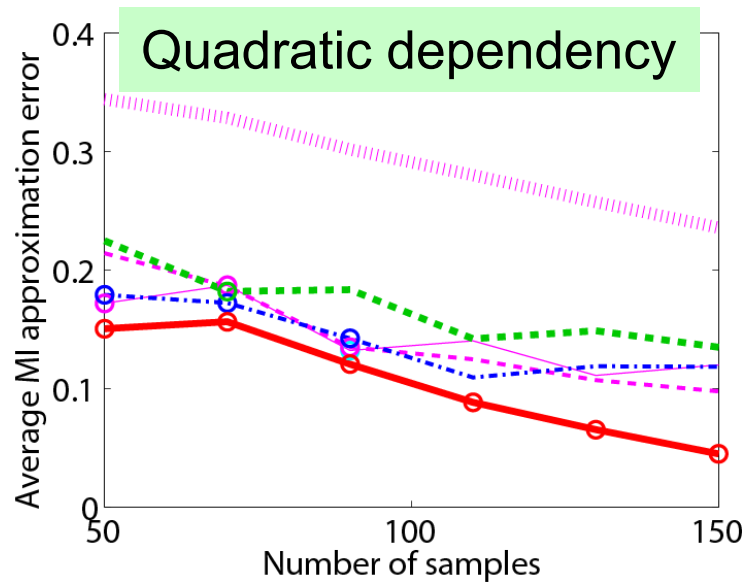
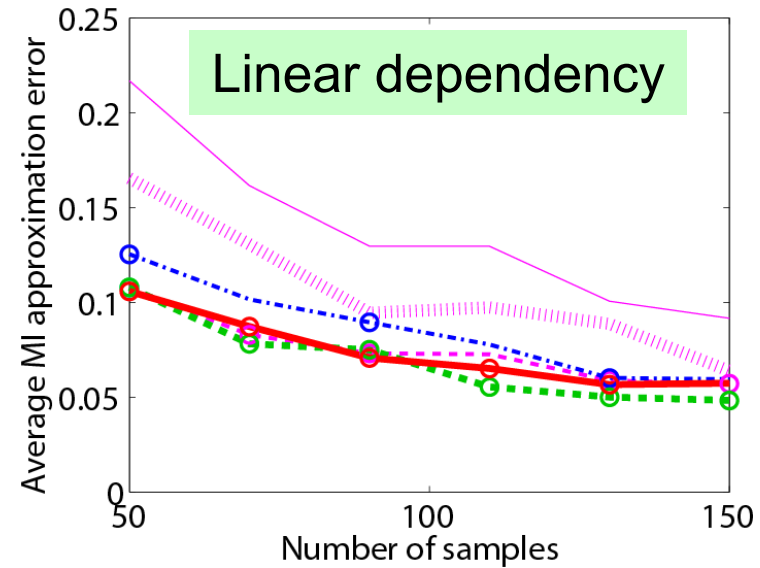
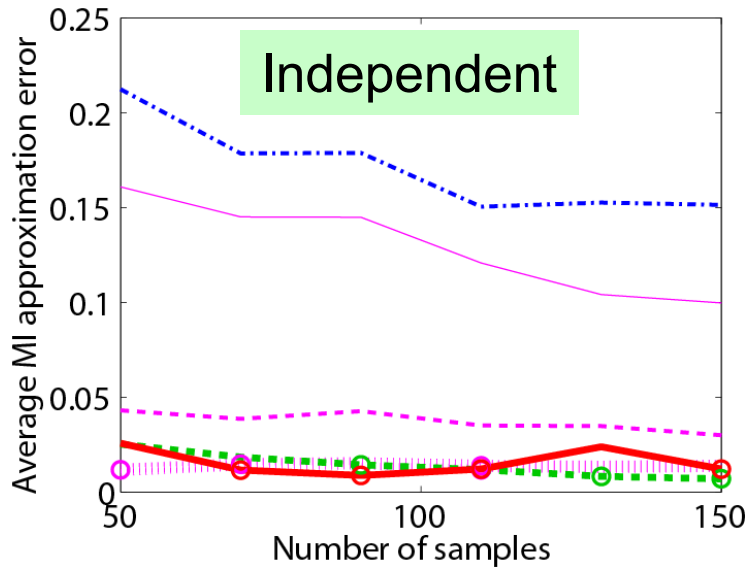


MI Approximation Error

51



$$\text{Error} = |\widehat{\text{MI}} - \text{MI}|$$



Estimation of Squared-Loss Mutual Information (SMI)

52

Suzuki, MS, Sese & Kanamori (BMC Bioinfo. 2009)

- Ordinary MI is based on the KL-divergence.
- SMI is based on the Pearson divergence:

$$\text{SMI} = \iint p(\mathbf{x})p(\mathbf{y}) \left(\frac{p(\mathbf{x}, \mathbf{y})}{p(\mathbf{x})p(\mathbf{y})} - 1 \right)^2 d\mathbf{x}d\mathbf{y}$$

- Can also be used as an independence measure.
- Can be approximated **analytically** and efficiently by least-squares density ratio estimation (uLSIF).

Usage of SMI Estimator

53

■ Between input and output:

- Feature ranking

Suzuki, MS, Sese & Kanamori
(BMCBioinfo 2009)

- Sufficient dimension reduction

Suzuki & MS (NeCo2012)

- Clustering

MS, Yamada, Kimura & Hachiya (ICML2011)
Kimura & MS (JACIII2011)

■ Between inputs:

- Independent component analysis

Suzuki & MS
(NeCo2010)

- Object matching

Yamada & MS (AISTATS2011)

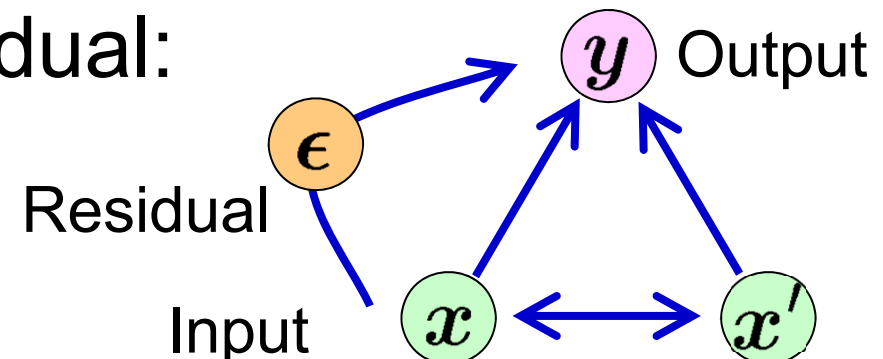
- Canonical dependency analysis

Karasuyama
& MS (NN2012)

■ Between input and residual:

- Causal inference

Yamada & MS (AAAI2010)



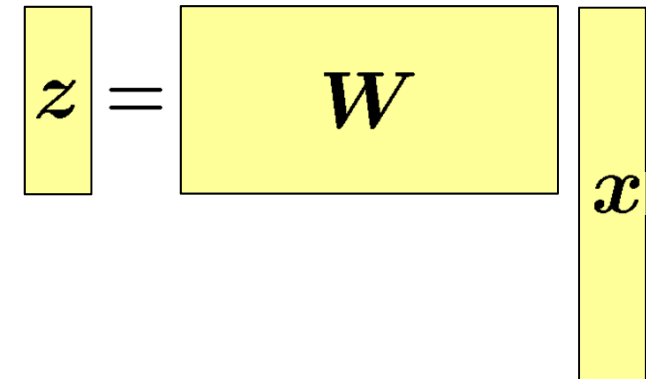
Sufficient Dimension Reduction⁵⁴

Li (JASA1991)

■ Input: x

■ Output: y

■ Projected input: $z = Wx$



$$WW^T = I$$

■ **Goal:** Find W so that z contains all information on y , i.e., $y \perp\!\!\!\perp x \mid z$

● In terms of SMI: Suzuki & MS (NeCo2012)

$$y \perp\!\!\!\perp x \mid z \iff \max_W \text{SMI}(Wx, y)$$

Sufficient Dimension Reduction⁵⁵ via SMI Maximization

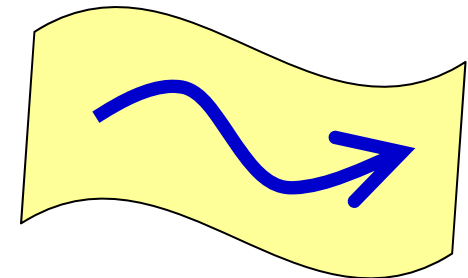
- Let's solve $\max_{\mathbf{W}} \widehat{\text{SMI}}(\mathbf{W})$ subject to $\mathbf{W}\mathbf{W}^\top = \mathbf{I}$.

$$\widehat{\text{SMI}}(\mathbf{W}) = 2\hat{\mathbf{h}}^\top \hat{\boldsymbol{\alpha}} - \hat{\boldsymbol{\alpha}}^\top \widehat{\mathbf{H}} \hat{\boldsymbol{\alpha}} - 1 \quad \hat{\boldsymbol{\alpha}} : \text{uLSIF solution}$$

- Since \mathbf{W} is on a Grassmann manifold,
natural gradient gives the steepest direction:

Amari (NeCo1998)

$$\mathbf{W} \leftarrow \mathbf{W} + \epsilon \frac{\partial \widehat{\text{SMI}}}{\partial \mathbf{W}} \left(\mathbf{I} - \mathbf{W}^\top \mathbf{W} \right)$$



- A computationally efficient heuristic update is also available.

Yamada, Niu, Takagi & MS (ACML2011)

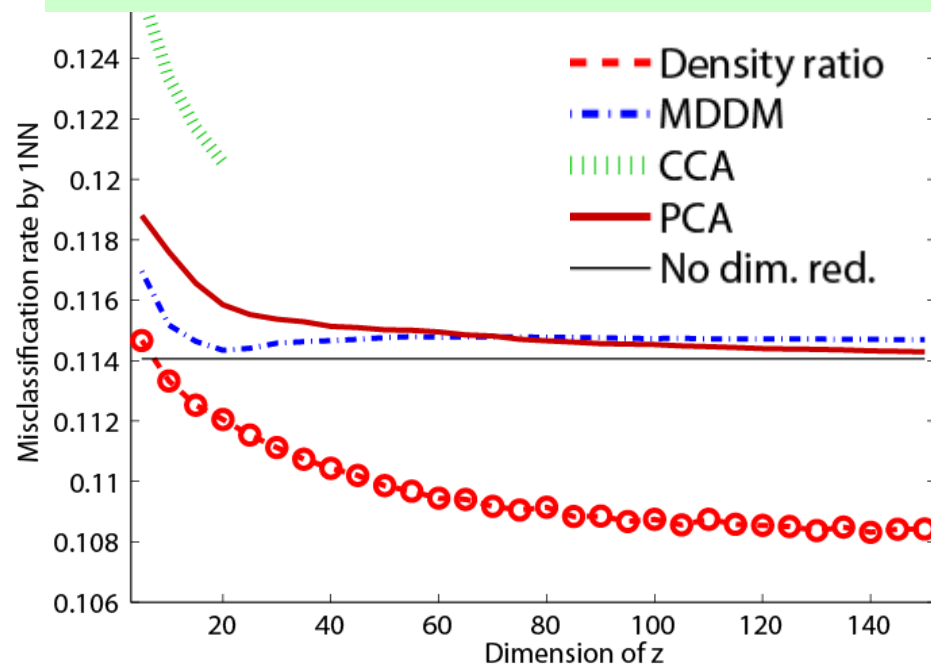
Experiments

56

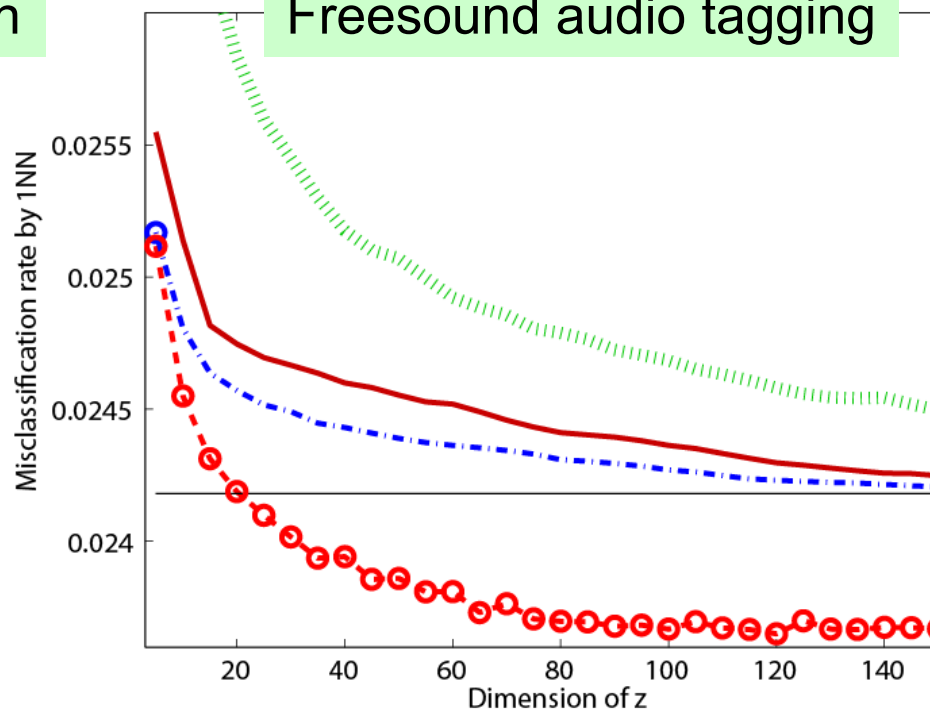
Yamada, Niu, Takagi & MS (ACML2011)

Dimension reduction for multi-label data:

Pascal VOC 2010 image classification



Freesound audio tagging



- **MDDM**: Multi-label dimensionality reduction via dependence maximization (MDDM) Zhang & Zhou (ACM-TKDD2010)
- **CCA**: Canonical correlation analysis
- **PCA**: Principal component analysis



Organization of This Lecture

57

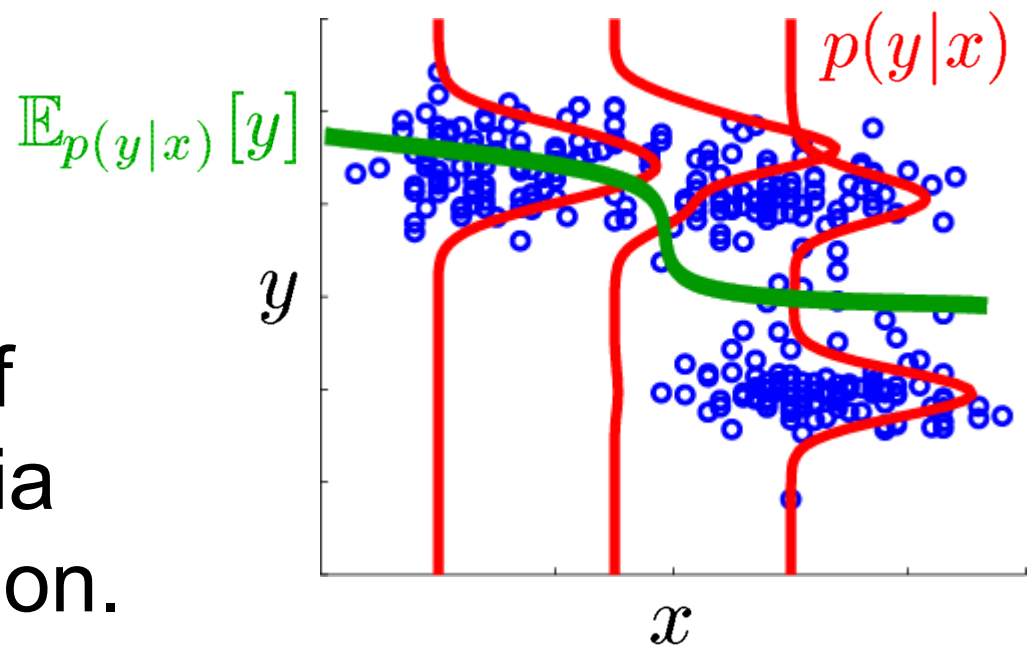
1. Introduction
2. Methods of Density Ratio Estimation
3. Usage of Density Ratios
 - A) Importance sampling
 - B) Distribution comparison
 - C) Mutual information estimation
 - D) Conditional probability estimation
4. More on Density Ratio Estimation
5. Conclusions

Conditional Density Estimation⁵⁸

$$p(\mathbf{y}|\mathbf{x}) = \frac{p(\mathbf{x}, \mathbf{y})}{p(\mathbf{x})}$$

MS, Takeuchi, Suzuki, Kanamori,
Hachiya & Okanohara (IEICE-ED2010)

- Regression = Conditional **mean** estimation
- However, regression is not informative enough for **complex** data analysis:
 - **Multi-modality**
 - **Asymmetry**
 - **Hetero-scedasticity**
- Directly estimation of conditional density via density-ratio estimation.

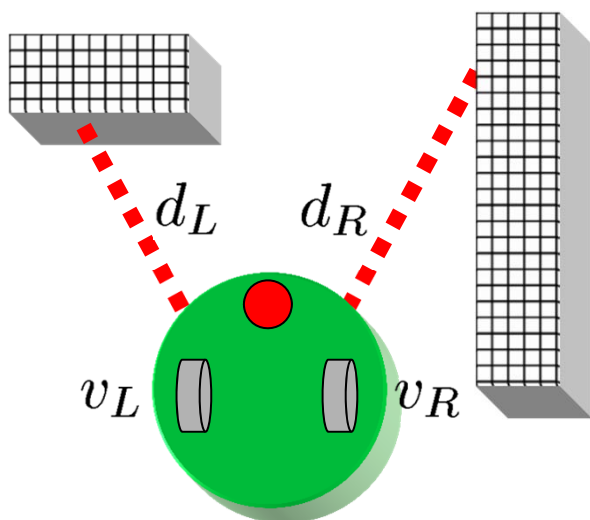


Experiments: Transition Estimation for Mobile Robot

- **Transition probability** $p(s'|s, a)$: Probability of being at state s' when action a is taken at s .

Khepera robot

- **State:** Infrared sensors
- **Action:** Wheel speed



Mean (std.) test negative log-likelihood over 10 runs (smaller is better)
(red: comparable by 5% t-test)

Data	uLSIF	ϵ -KDE	MDN
Khepera1	1.69(0.01)	2.07(0.02)	1.90(0.36)
Khepera2	1.86(0,01)	2.10(0.01)	1.92(0.26)
Pendulum1	1.27(0.05)	2.04(0.10)	1.44(0.67)
Pendulum2	1.38(0.05)	2.07(0.10)	1.43(0.58)
Comp. Time	1	0.164	1134

ϵ -KDE: ϵ -neighbor kernel density estimation

MDN: Mixture density network Bishop (Book2006)

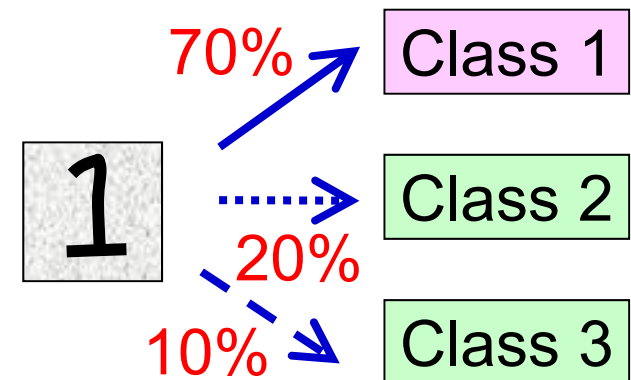
Probabilistic Classification

60

$$p(\mathbf{y}|\mathbf{x}) = \frac{p(\mathbf{x}, \mathbf{y})}{p(\mathbf{x})}$$

MS (IEICE-ED2010)

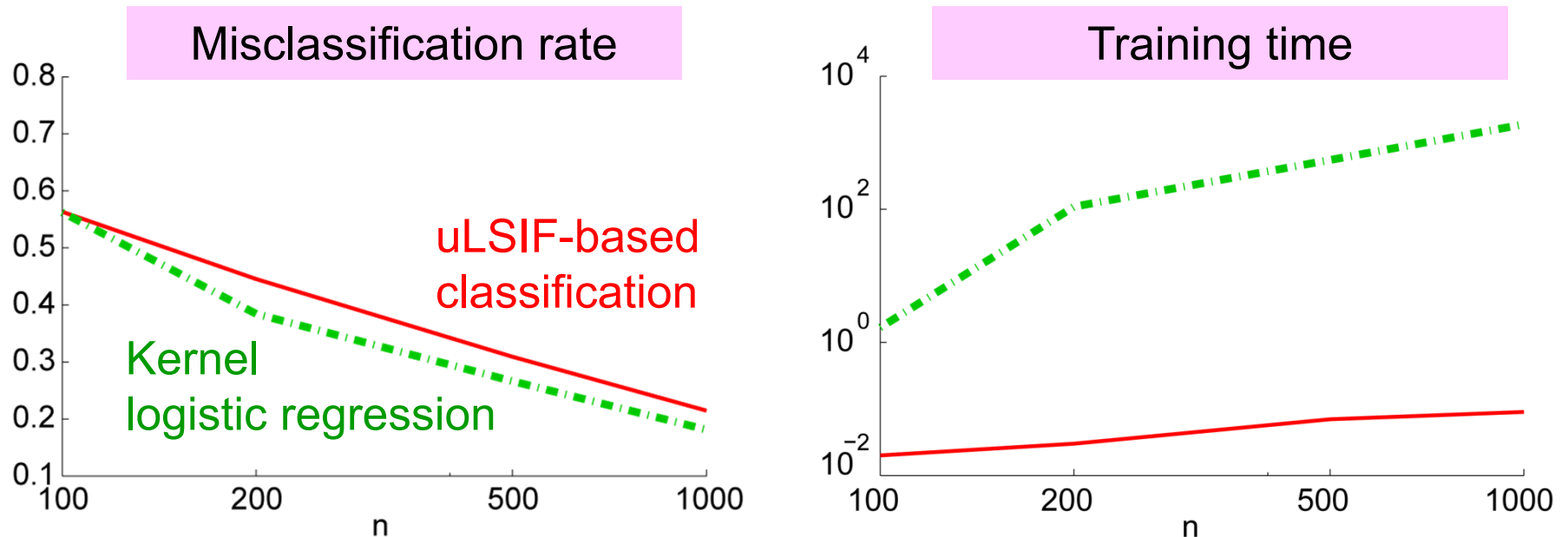
- If y is **categorical**, conditional probability estimation corresponds to learning **class-posterior probability**.
- Least-squares density ratio estimation (uLSIF) provides **an analytic estimator**:
 - Computationally efficient alternative to kernel logistic regression.
 - No normalization term included.
 - Classwise training is possible.



Numerical Example

61

Letter dataset (26 classes):



uLSIF-based classification method:

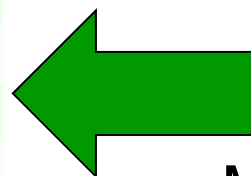
- Comparable accuracy with KLR.
- **Training is 1000 times faster!**

More Experiments

62

Yamada, MS, Wichern & Simm (IEICE2011)

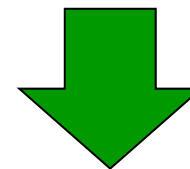
Dataset	uLSIF	KLR
Aeroplane	82.6(1.0)	83.0(1.3)
Bicycle	77.7(1.7)	76.6(3.4)
Bird	68.7(2.0)	70.8(2.2)
Boat	74.4(2.0)	72.8(2.6)
Bottle	65.4(1.8)	62.1(4.3)
Bus	85.4(1.4)	85.6(1.4)
Car	73.0(0.8)	72.1(1.2)
Cat	73.6(1.4)	74.1(1.7)
Chair	71.0(1.0)	70.5(1.0)
Cow	71.7(3.2)	69.3(3.6)
Diningtable	75.0(1.6)	71.4(2.7)
Dog	69.6(1.0)	69.4(1.8)
Horse	64.4(2.5)	61.2(3.2)
Motorbike	77.0(1.7)	75.9(3.3)
Person	67.6(0.9)	67.0(0.8)
Pottedplant	66.2(2.6)	61.9(3.2)
Sheep	77.8(1.6)	74.0(3.8)
Sofa	67.4(2.7)	65.4(4.6)
Train	79.2(1.3)	78.4(3.0)
Tvmonitor	76.7(2.2)	76.6(2.3)
Training time [sec]	0.7	24.6



Pascal VOC 2010
image classification:

Mean AUC (std) over 50 runs
(red: comparable by 5% t-test)

Freesound audio tagging:
Mean AUC (std) over 50 runs



	uLSIF	KLR
AUC	70.1(9.6)	66.7(10.3)
Training time [sec]	0.005	0.612

Other Applications

63

- Action recognition from accelerometer

Hachiya, MS & Ueda (Neurocomputing2011)

- Age prediction from faces

Ueki, MS, Ihara & Fujita (ACPR2011)



Organization of This Lecture

64

1. Introduction
2. Methods of Density Ratio Estimation
3. Usage of Density Ratios
4. **More on Density Ratio Estimation**
 - A) **Unified Framework**
 - B) Dimensionality Reduction
 - C) Relative Density Ratios
5. Conclusions

Bregman (BR) Divergence

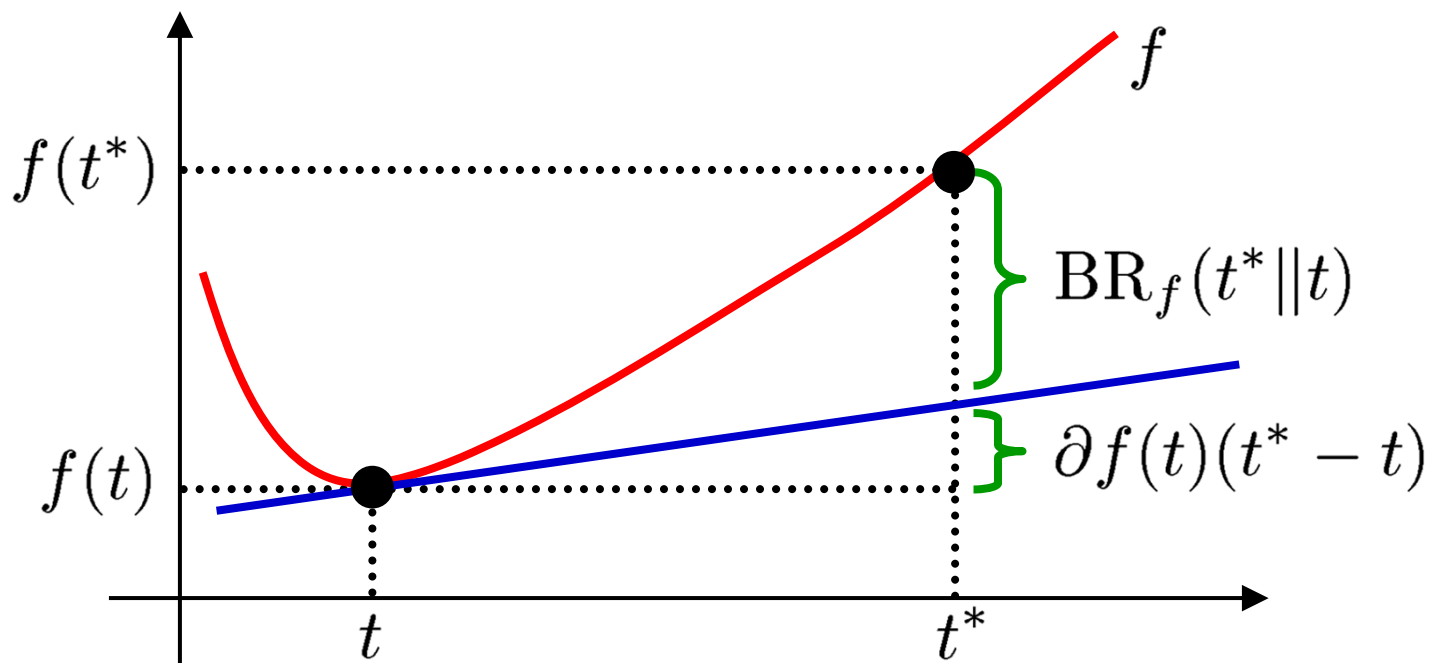
65

Bregman (1967)

- f : Differentiable convex function
- BR divergence with function f :

$$\text{BR}_f(t^* || t) := f(t^*) - f(t) - \underbrace{\partial f(t)(t^* - t)}$$

Linear prediction from $f(t)$ to $f(t^*)$



Density-Ratio Fitting under BR Divergence

MS, Suzuki & Kanamori (AISM2012)

- Fit a ratio model $\hat{r}(\mathbf{x})$ to true ratio $r(\mathbf{x})$ under the BR divergence:

$$\min_{\hat{r}} \text{BR}_f(\hat{r})$$

$$\text{BR}_f(\hat{r}) = \int p_{\text{de}}(\mathbf{x}) \nabla f(\hat{r}(\mathbf{x})) \hat{r}(\mathbf{x}) d\mathbf{x} - \int p_{\text{de}}(\mathbf{x}) f(\hat{r}(\mathbf{x})) d\mathbf{x}$$

$$- \int p_{\text{nu}}(\mathbf{x}) \nabla f(\hat{r}(\mathbf{x})) d\mathbf{x} + C$$

$$\approx \frac{1}{n_{\text{de}}} \sum_{j=1}^{n_{\text{de}}} \nabla f(\hat{r}(\mathbf{x}_j^{\text{de}})) \hat{r}(\mathbf{x}_j^{\text{de}}) - \frac{1}{n_{\text{de}}} \sum_{j=1}^{n_{\text{de}}} f(\hat{r}(\mathbf{x}_j^{\text{de}}))$$

$$- \frac{1}{n_{\text{nu}}} \sum_{i=1}^{n_{\text{nu}}} \nabla f(\hat{r}(\mathbf{x}_i^{\text{nu}})) + C$$

$$r(\mathbf{x}) = \frac{p_{\text{nu}}(\mathbf{x})}{p_{\text{de}}(\mathbf{x})}$$

Unified View

67

- Logistic regression:

$$f(t) = t \log t - (1 + t) \log(1 + t)$$

- (Extended) kernel mean matching:

$$f(t) = (t - 1)^2 / 2$$

$$\min_{\hat{r}} \|\nabla J(\hat{r})\|^2$$

- KL-based method:

$$f(t) = t \log t - t$$

- uLSIF:

$$f(t) = (t - 1)^2 / 2$$

$$\min_{\hat{r}} J(\hat{r})$$

- Robust estimator (power divergence):

$$f(t) = \alpha^{-1} (t^{1+\alpha} - t) \quad \alpha > 0$$



Organization of This Lecture

68

1. Introduction
2. Methods of Density Ratio Estimation
3. Usage of Density Ratios
4. **More on Density Ratio Estimation**
 - A) Unified Framework
 - B) **Dimensionality Reduction**
 - C) Relative Density Ratios
5. Conclusions

Direct Density-Ratio Estimation⁶⁹ with Dimensionality Reduction (D^3)

- Directly density-ratio estimation without density estimation is promising.
- However, for **high-dimensional data**, density-ratio estimation is still challenging.
- We combine direct density-ratio estimation with **dimensionality reduction!**

Hetero-distributional Subspace (HS)⁷⁰

MS, Kawanabe & Chui (NN2010)

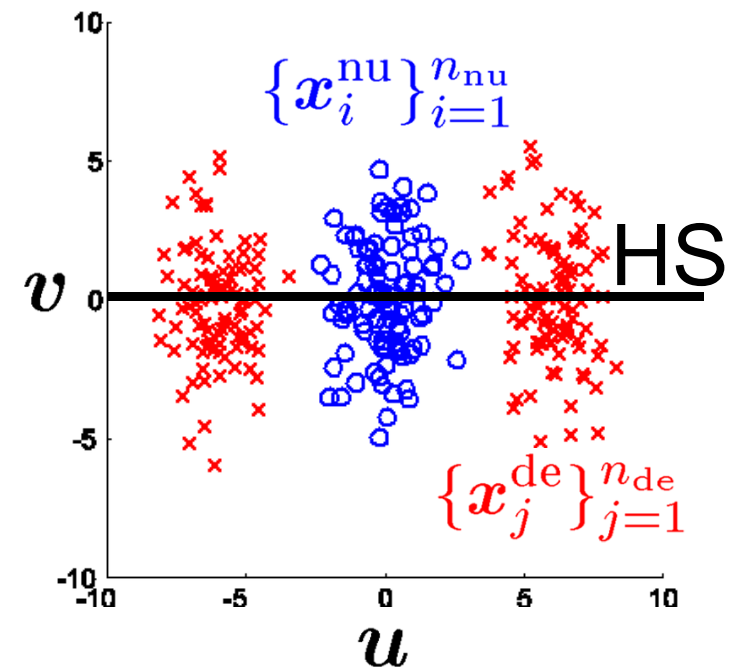
- **Key assumption:** $p_{\text{nu}}(\mathbf{x})$ and $p_{\text{de}}(\mathbf{x})$ are different only in a subspace (called HS).

$$r(\mathbf{x}) = \frac{p_{\text{nu}}(\mathbf{x})}{p_{\text{de}}(\mathbf{x})}$$

$$\begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix} = \begin{bmatrix} \mathbf{U} \\ \mathbf{V} \end{bmatrix} \mathbf{x}$$

$$= \frac{p(\mathbf{v}|\mathbf{u})p_{\text{nu}}(\mathbf{u})}{p(\mathbf{v}|\mathbf{u})p_{\text{de}}(\mathbf{u})} = \frac{p_{\text{nu}}(\mathbf{u})}{p_{\text{de}}(\mathbf{u})}$$

$\begin{bmatrix} \mathbf{U} \\ \mathbf{V} \end{bmatrix}$: Full-rank and orthogonal



- This allows us to estimate the density ratio **only within the low-dimensional HS!**

Characterization of HS

71

MS, Yamada, von Büнау, Suzuki, Kanamori & Kawanabe (NN2011)

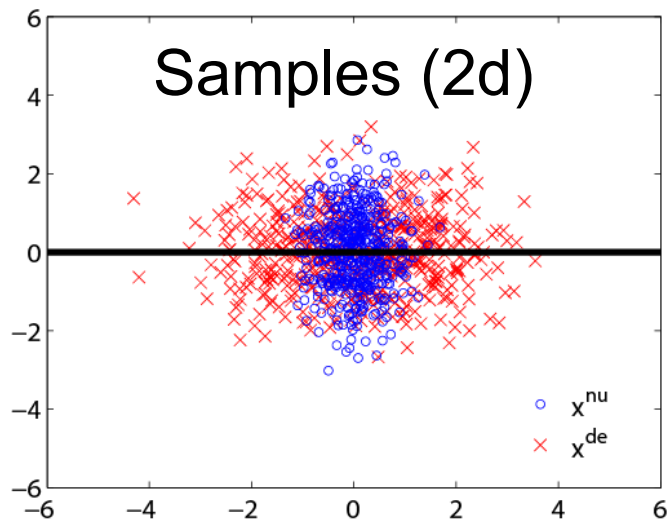
- HS is given as the **maximizer of the Pearson divergence** with respect to U :

$$\text{PE}[p_{\text{nu}}(\mathbf{u}), p_{\text{de}}(\mathbf{u})] = \int \left(\frac{p_{\text{nu}}(\mathbf{u})}{p_{\text{de}}(\mathbf{u})} - 1 \right)^2 p_{\text{de}}(\mathbf{u}) d\mathbf{u}$$

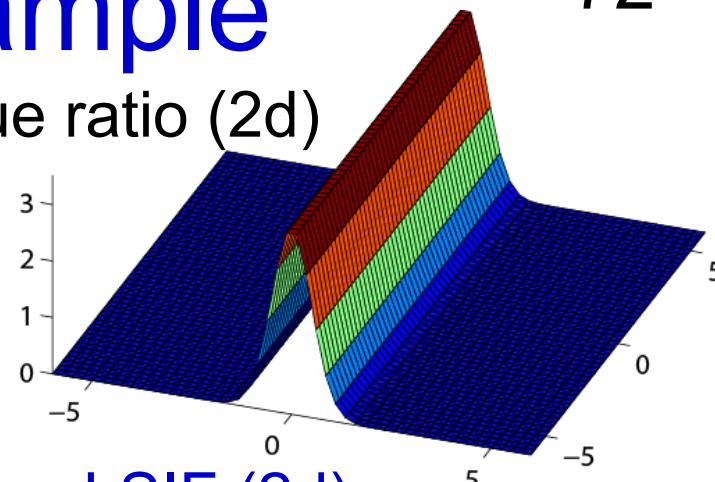
- PE can be **analytically** approximated by uLSIF (with good convergence property).
- HS search by
 - Natural gradient
 - A heuristic update Yamada & MS (AAAI2011)

Numerical Example

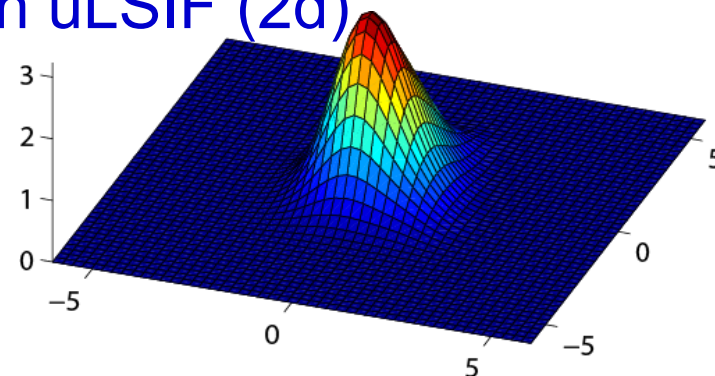
72



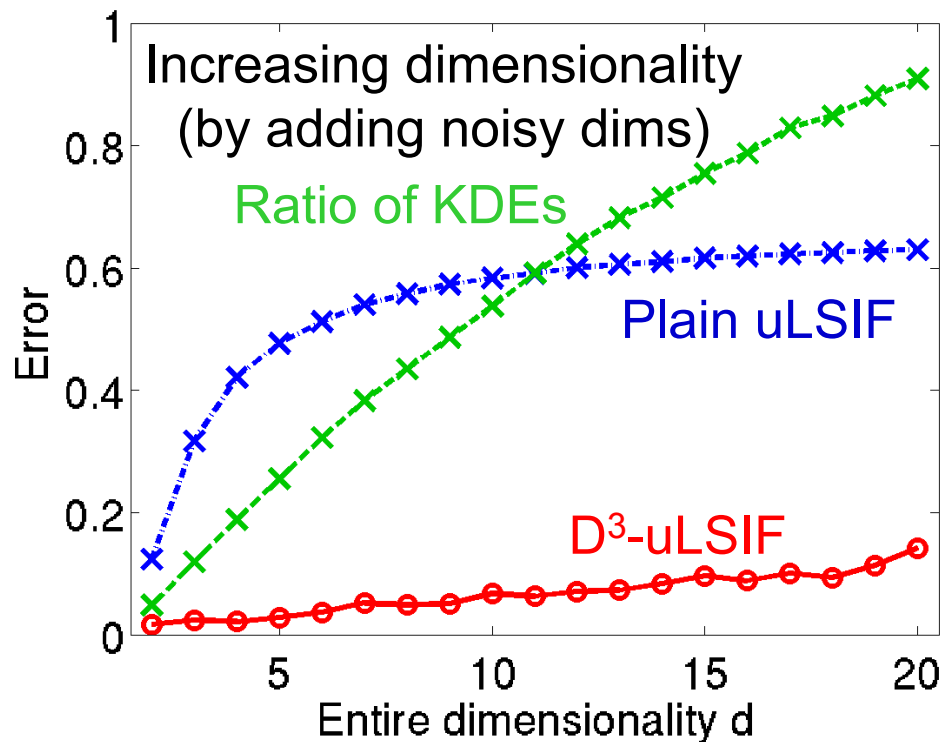
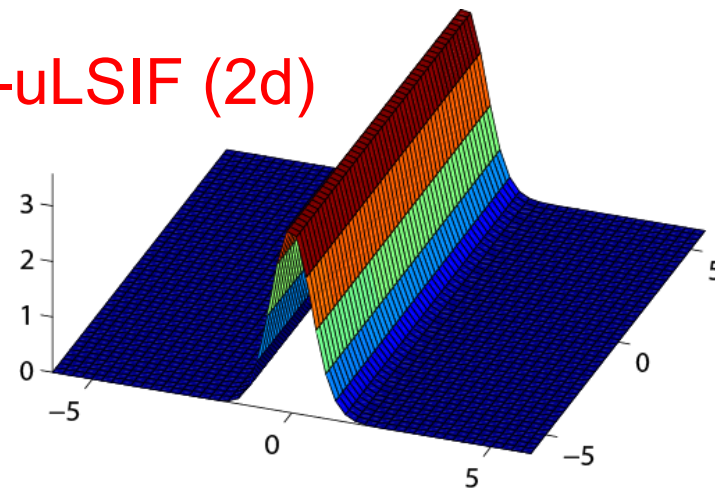
True ratio (2d)



Plain uLSIF (2d)



D³-uLSIF (2d)





Organization of This Lecture

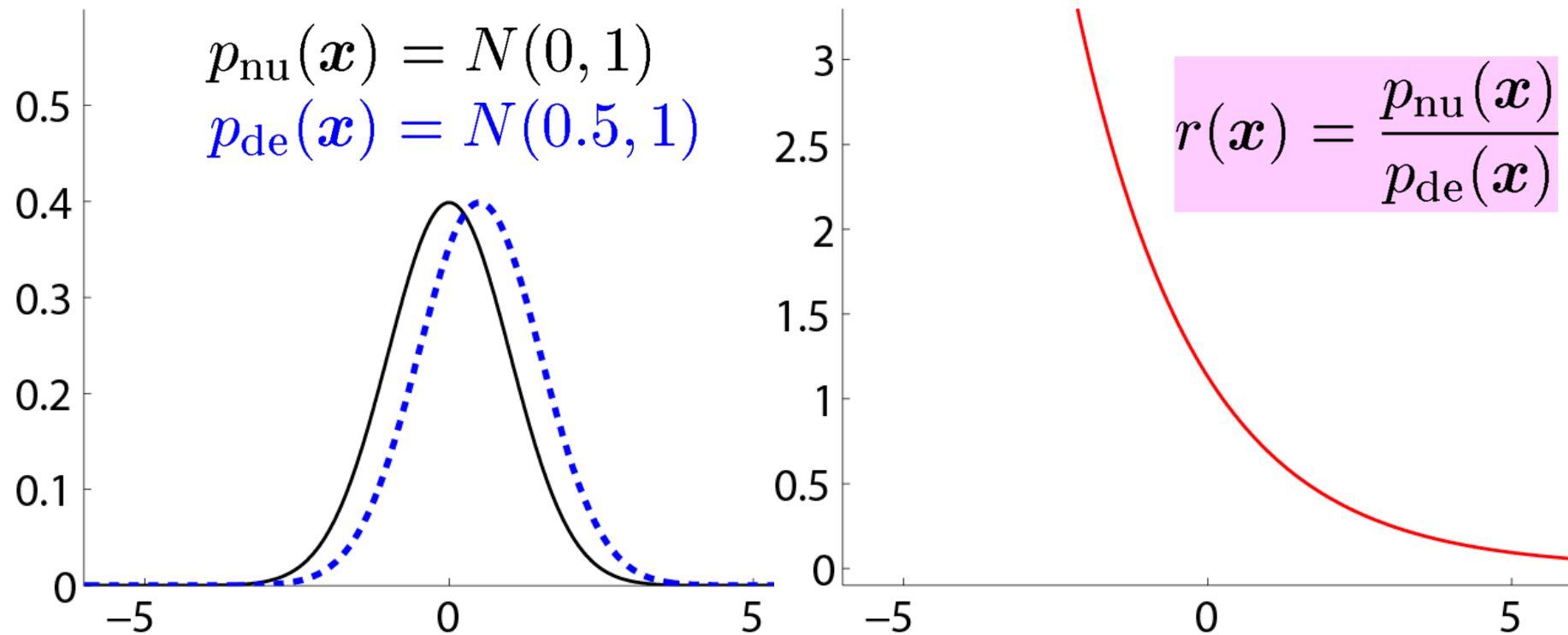
73

1. Introduction
2. Methods of Density Ratio Estimation
3. Usage of Density Ratios
4. **More on Density Ratio Estimation**
 - A) Unified Framework
 - B) Dimensionality Reduction
 - C) **Relative Density Ratios**
5. Conclusions

Weakness of Density Ratios

74

- Density ratio can diverge to **infinity**:



- Estimation becomes unreliable!

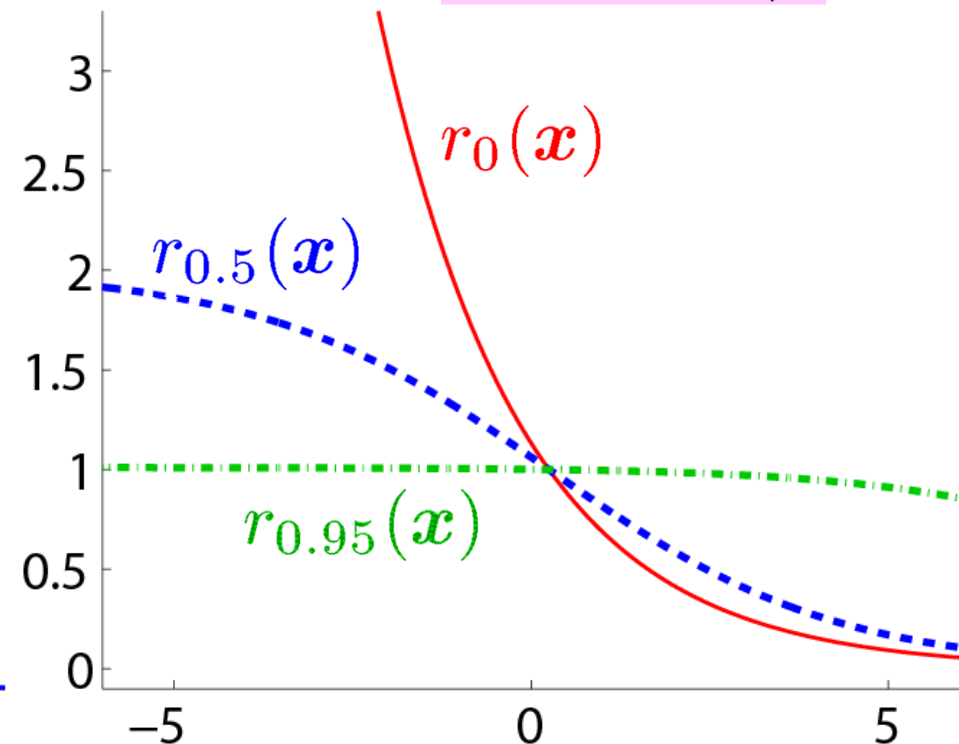
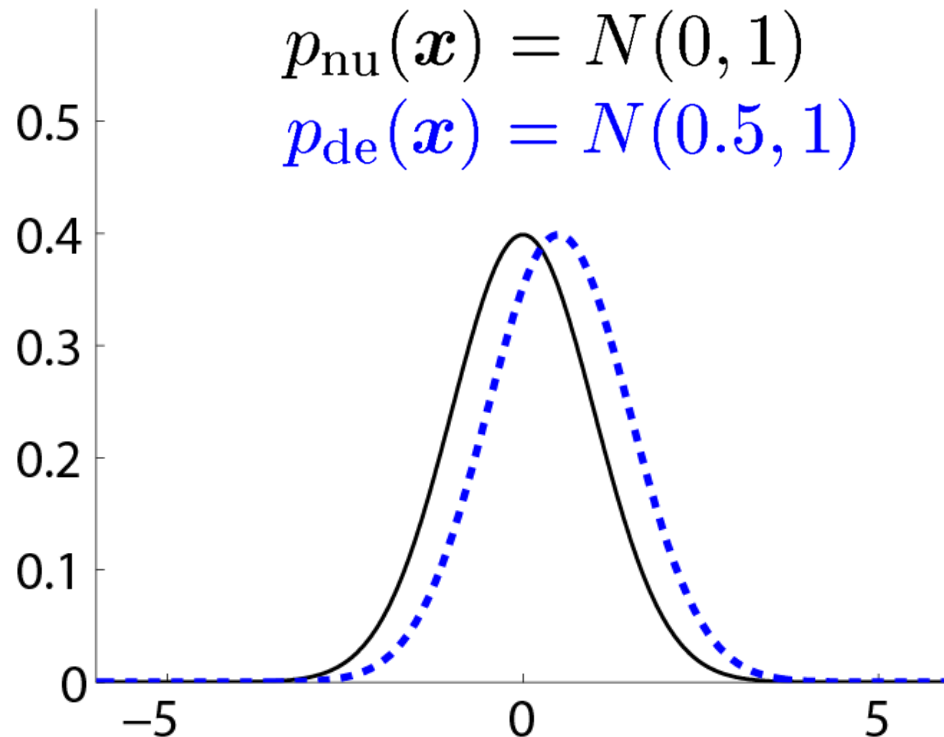
Relative Density Ratios

75

Yamada, Suzuki, Kanamori, Hachiya & MS (NIPS2011)

$$r_{\beta}(\mathbf{x}) = \frac{p_{\text{nu}}(\mathbf{x})}{\beta p_{\text{nu}}(\mathbf{x}) + (1 - \beta)p_{\text{de}}(\mathbf{x})} \quad 0 \leq \beta < 1$$

■ Bounded for any $p_{\text{nu}}(\mathbf{x}), p_{\text{de}}(\mathbf{x})$: $r_{\beta}(\mathbf{x}) < \frac{1}{\beta}$



Estimation of Relative Ratios

76

■ Linear model: $\hat{r}(\mathbf{x}) = \sum_{\ell=1}^b \alpha_{\ell} \phi_{\ell}(\mathbf{x}) = \boldsymbol{\alpha}^{\top} \boldsymbol{\phi}(\mathbf{x})$

- Relative unconstrained least-squares importance fitting (RuLSIF):

$$\min_{\hat{r}} \int \left(\hat{r}(\mathbf{x}) - r_{\beta}(\mathbf{x}) \right)^2 q_{\beta}(\mathbf{x}) d\mathbf{x} \quad r_{\beta}(\mathbf{x}) = \frac{p_{\text{nu}}(\mathbf{x})}{q_{\beta}(\mathbf{x})}$$

$$q_{\beta}(\mathbf{x}) = \beta p_{\text{nu}}(\mathbf{x}) + (1 - \beta) p_{\text{de}}(\mathbf{x})$$

- The solution can be computed analytically:

$$\operatorname{argmin}_{\boldsymbol{\alpha}} \left[\frac{1}{2} \boldsymbol{\alpha}^{\top} \widehat{\mathbf{H}} \boldsymbol{\alpha} - \widehat{\mathbf{h}}^{\top} \boldsymbol{\alpha} + \frac{\lambda}{2} \boldsymbol{\alpha}^{\top} \boldsymbol{\alpha} \right] = (\widehat{\mathbf{H}} + \lambda \mathbf{I})^{-1} \widehat{\mathbf{h}}$$

$$\widehat{\mathbf{H}} = \frac{\beta}{n_{\text{de}}} \sum_{j=1}^{n_{\text{de}}} \boldsymbol{\phi}(\mathbf{x}_j^{\text{de}}) \boldsymbol{\phi}(\mathbf{x}_j^{\text{de}})^{\top} + \frac{1 - \beta}{n_{\text{nu}}} \sum_{i=1}^{n_{\text{nu}}} \boldsymbol{\phi}(\mathbf{x}_i^{\text{nu}}) \boldsymbol{\phi}(\mathbf{x}_i^{\text{nu}})^{\top} \quad \widehat{\mathbf{h}} = \frac{1}{n_{\text{nu}}} \sum_{i=1}^{n_{\text{nu}}} \boldsymbol{\phi}(\mathbf{x}_i^{\text{nu}})$$

Relative Pearson Divergence 77

$$\text{PE}_\beta[p_{\text{nu}}(\mathbf{x}), p_{\text{de}}(\mathbf{x})] = \frac{1}{2} \int \left(r_\beta(\mathbf{x}) - 1 \right)^2 q_\beta(\mathbf{x}) d\mathbf{x}$$

$$r_\beta(\mathbf{x}) = \frac{p_{\text{nu}}(\mathbf{x})}{q_\beta(\mathbf{x})} \quad q_\beta(\mathbf{x}) = \beta p_{\text{nu}}(\mathbf{x}) + (1 - \beta) p_{\text{de}}(\mathbf{x})$$

- Relative Pearson divergence can be more reliably approximated:

$$\widehat{\text{PE}}_\beta - \text{PE}_\beta = \mathcal{O}_p(n^{-1/2} c \|r_\beta\|_\infty + \lambda_n \max(1, R(r_\beta)^2))$$

$$n = \min(n_{\text{nu}}, n_{\text{de}}) \quad \lambda_n \rightarrow o(1) \quad \text{and} \quad \lambda_n^{-1} = o(n^{2/(2+\gamma)}), \quad 0 < \gamma < 2$$

$$\|r_\beta\|_\infty = \max_{\mathbf{x}} r_\beta(\mathbf{x}) = \left\| \left(\beta + (1 - \beta)/r(\mathbf{x}) \right)^{-1} \right\|_\infty < \frac{1}{\beta} \quad r(\mathbf{x}) = \frac{p_{\text{nu}}(\mathbf{x})}{p_{\text{de}}(\mathbf{x})}$$



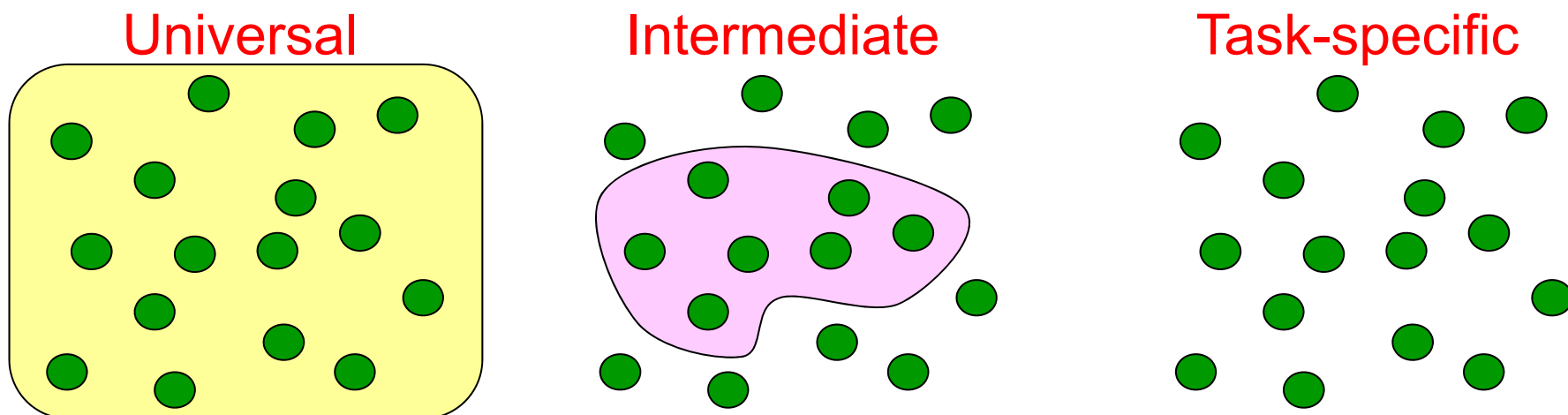
Organization of This Lecture

78

1. Introduction
2. Methods of Density Ratio Estimation
3. Usage of Density Ratios
4. More on Density Ratio Estimation
5. **Conclusions**

Conclusions

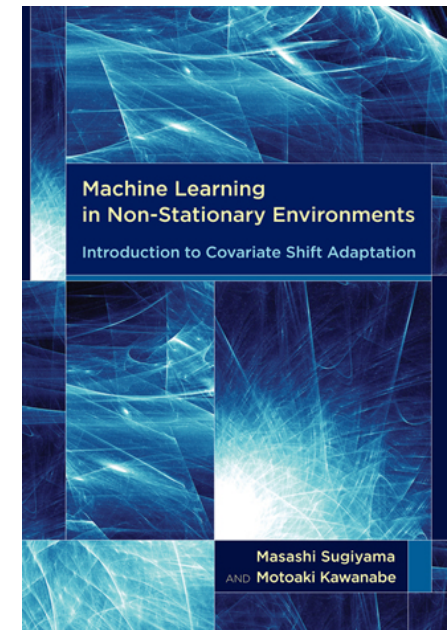
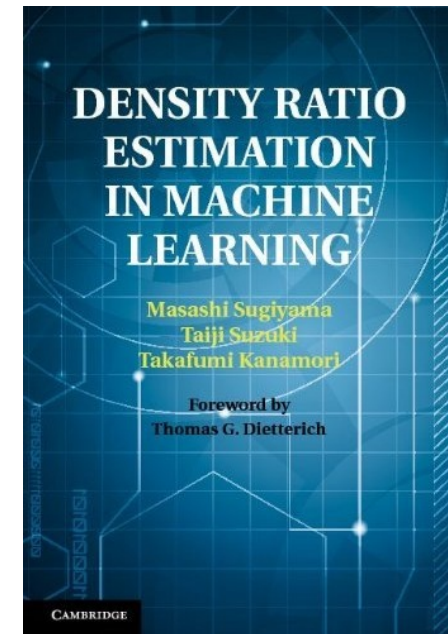
- Estimating data-generating probability distributions is universal, but inaccurate.
- Solving each task directly is ideal, but costly.
- **Density ratio estimation** is realistic compromise:
 - Systematically avoiding density estimation.
 - Applicable to a reasonably rich class of tasks.
 - Useful in many real-world problems.



Books on Density Ratios

80

- Sugiyama, Suzuki & Kanamori,
**Density Ratio Estimation
in Machine Learning**,
Cambridge University Press, 2012
- Sugiyama & Kawanabe
**Machine Learning
in Non-Stationary Environments**,
MIT Press, 2012



Acknowledgements

81

- **Colleagues:** Hirotaka Hachiya, Shohei Hido, Yasuyuki Ihara, Hisashi Kashima, Motoaki Kawanabe, Manabu Kimura, Masakazu Matsugu, Shin-ichi Nakajima, Klaus-Robert Müller, Jun Sese, Jaak Simm, Ichiro Takeuchi, Masafumi Takimoto, Yuta Tsuboi, Kazuya Ueki, Paul von Büнау, Gordon Wichern, Makoto Yamada.
- **Funding Agencies:** Ministry of Education, Culture, Sports, Science and Technology, Alexander von Humboldt Foundation, Okawa Foundation, Microsoft Institute for Japanese Academic Research Collaboration Collaborative Research Project, IBM Faculty Award, Mathematisches Forschungsinstitut Oberwolfach Research-in-Pairs Program, Asian Office of Aerospace Research and Development, Support Center for Advanced Telecommunications Technology Research Foundation, Japan Science and Technology Agency
- Papers, articles, and software of density ratio estimation are available from

<http://sugiyama-www.cs.titech.ac.jp/~sugi/>