# Density-Difference Estimation

Masashi Sugiyama (Tokyo Institute of Technology, Japan), Takafumi Kanamori (Nagoya University, Japan), Taiji Suzuki (University of Tokyo, Japan)
Marthinus Christoffel du Plessis (Tokyo Institute of Technology, Japan), Song Liu (Tokyo Institute of Technology, Japan), Ichiro Takeuchi (Nagoya Institute of Technology, Japan)

東京工業大学 Tokyo Institute of Technology · THE UNIVERSITY OF TOKYO · NAGOYA UNIVERSITY · Nagoya Institute of Technology

## 1. This Work in A Nutshell

- **Target problem**: Estimate the difference between two densities.
- **Approach**: Avoid density estimation and directly estimate the difference in a single-shot process.
- **Theory**: Parametric and non-parametric optimality in terms of the approximation accuracy.
- **Usage**: $L^2$-distance approximation.
- **Applications**: Semi-supervised class-prior estimation and unsupervised change detection.

## 2.1 Target Problem & Motivations

- From two sets of i.i.d. samples
$$\{x_i\}_{i=1}^{n} \overset{i.i.d.}{\sim} p(x) \quad \{x'_{i'}\}_{i'=1}^{n'} \overset{i.i.d.}{\sim} p'(x)$$
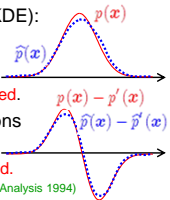estimate the difference between two densities:
$$f(x) = p(x) - p'(x)$$

- Via density-difference estimation, we want to
  - Compare probability distributions.
  - Approximate the $L^2$-distance.
$$\int \left( p(x) - p'(x) \right)^2 dx$$

## 2.2 Naïve Approach

- 2-step method of first estimating two densities separately and then computing their difference.
- Kernel density estimation (KDE):
$$\widehat{p}(x) \propto \sum_{i=1}^{n} \exp\left(-\frac{\|x - x_i\|^2}{2\sigma^2}\right)$$
  - A smoother function is obtained.
- Difference of smooth functions tends to be too smooth.
  - $L^2$-distance is under-estimated.

Cf. Anderson, Hall & Titterington (J. Multivariate Analysis 1994)

## 2.3 Vapnik's Principle

*Vladimir N. Vapnik (1998)
Statistical Learning Theory, Wiley.*

If you possess a restricted amount of information for solving some problem, try to solve the problem directly and never solve a more general problem as an intermediate step. It is possible that the available information is sufficient for a direct solution but is insufficient for solving a more general intermediate problem.

## 3.1 Least-Squares Density-Difference (LSDD) Estimation

- Density-difference model:
$$g(x) = \sum_{\ell=1}^{n+n'} \theta_\ell \exp\left(-\frac{\|x - c_\ell\|^2}{2\sigma^2}\right)$$
$$(c_1,\ldots,c_n,c_{n+1},\ldots,c_{n+n'}) = (x_1,\ldots,x_n,x'_1,\ldots,x'_{n'})$$

- Least-squares estimation:
$$\min_\theta J(\theta) \quad J(\theta) = \int \left(g(x) - f(x)\right)^2 dx$$
$$\theta = (\theta_1,\ldots,\theta_{n+n'})^\top \quad f(x) = p(x) - p'(x)$$

- Sample approximation:
$$J(\theta) = \int g(x)^2 dx - 2\int g(x)f(x)dx + \int f(x)^2 dx$$
  (Analytically computable) (Consistently estimable from samples) (Safely ignorable)
  - $\int g(x)^2 dx = \sum_{\ell,\ell'=1}^{n+n'} \theta_\ell \theta_{\ell'} (\pi\sigma^2)^{d/2} \exp\left(-\frac{\|c_\ell - c_{\ell'}\|^2}{4\sigma^2}\right) \quad x \in \mathbb{R}^d$
  - $\int g(x)f(x)dx \approx \frac{1}{n}\sum_{i=1}^{n} g(x_i) - \frac{1}{n'}\sum_{i'=1}^{n'} g(x'_{i'})$

- Regularized training criterion:
$$\widehat{\theta} = \underset{\theta}{\operatorname{argmin}} \left[\theta^\top H \theta - 2\widehat{h}^\top \theta + \lambda \theta^\top \theta\right]$$
$$H_{\ell,\ell'} = (\pi\sigma^2)^{d/2} \exp\left(-\frac{\|c_\ell - c_{\ell'}\|^2}{4\sigma^2}\right)$$
$$\widehat{h}_\ell = \frac{1}{n}\sum_{i=1}^{n}\exp\left(-\frac{\|x_i - c_\ell\|^2}{2\sigma^2}\right) - \frac{1}{n'}\sum_{i'=1}^{n'}\exp\left(-\frac{\|x'_{i'} - c_\ell\|^2}{2\sigma^2}\right)$$

- Solution can be computed analytically:
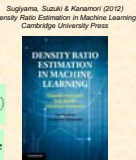$$\widehat{\theta} = (H + \lambda I_b)^{-1} \widehat{h}$$

- Cross-validation is available for objectively tuning Gaussian width $\sigma^2$ and regularization parameter $\lambda$.
- MATLAB code is available:
http://sugiyama-www.cs.titech.ac.jp/~sugi/software/LSDD/

### Related Work: Density Ratio Estimation
$$r(x) = \frac{p(x)}{p'(x)}$$
*Sugiyama, Suzuki & Kanamori (2012) Density Ratio Estimation in Machine Learning Cambridge University Press*

- **Density ratio is a versatile tool:**
  - Importance sampling: $\sum_i \frac{p_{test}(x_i)}{p_{train}(x_i)} \mathrm{loss}(x_i)$
  - Divergence estimation: $\int p'(x) f\left(\frac{p(x)}{p'(x)}\right) dx$
  - Mutual information estimation: $\iint p(x,y) \log\frac{p(x,y)}{p(x)p(y)} dx dy$
  - Conditional probability estimation: $p(y|x) = \frac{p(x,y)}{p(x)}$

## 3.2 Theoretical Analyses

- **Parametric case:** $g(x) = \sum_{\ell=1}^{b} \theta_\ell \upsilon_\ell(x) = \theta^\top \psi(x)$

For $\frac{n}{n+n'} \overset{n,n'\to\infty}{\longrightarrow} \eta \in [0,1]$, we have
$$\left(\sqrt{\frac{1}{n}+\frac{1}{n'}}\right)^{-1} (\widehat{\theta} - \theta^*) \rightsquigarrow N\left(0, H^{-1}((1-\eta)V_p + \eta V_{p'})H^{-1}\right)$$
$\widehat{\theta}$ : LSDD estimator  $V_p = \int (\psi(x) - v_p)(\psi(x) - v_p)^\top p(x)dx$
$\theta^*$ : Optimal parameter  $H = \int \psi(x)\psi(x)^\top dx$  $v_p = \int \psi(x)p(x)dx$

Optimal convergence rate is achieved.

- **Non-parametric case:** Gaussian RKHS $\mathcal{H}$
$$\widehat{f} = \underset{g\in\mathcal{H}}{\arg\min}\left[\|g\|_{L^2}^2 - 2\left(\frac{1}{n}\sum_{i=1}^{n}g(x_i) + \frac{1}{n'}\sum_{i'=1}^{n'}g(x'_{i'})\right) + \lambda\|g\|_{\mathcal{H}}^2\right]$$
$$n = n' \quad k_\gamma(x,x') = \exp\left(-\frac{\|x-x'\|^2}{\gamma^2}\right)$$

For all $\rho,\rho' > 0$, there exists a constant $K$ such that, for appropriately chosen $\lambda, \gamma$, the following inequality holds with probability $1 - 4e^{-\tau}$ for all $\tau, n \geq 1$:
$$\|\widehat{f} - f\|_{L^2}^2 + \lambda\|\widehat{f}\|_{\mathcal{H}}^2 \leq K\left(n^{-\frac{2\alpha}{1+\frac{2\alpha}{d}+\rho}} + \tau n^{-1+\rho'}\right)$$
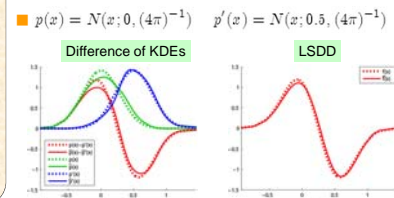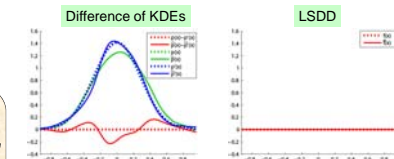$\alpha \geq 0$ : Regularity of Besov space that contains true $f$  $x \in \mathbb{R}^d$

Optimal convergence rate is achieved.

Cf. Eberts & Steinwart (NIPS2011)

## 3.3 Numerical Examples

- $p(x) = p'(x) = N(x; 0, (4\pi)^{-1})$  $n = n' = 200$


Difference of KDEs / LSDD

- $p(x) = N(x; 0, (4\pi)^{-1})$  $p'(x) = N(x; 0.5, (4\pi)^{-1})$


Difference of KDEs / LSDD

## 4.1 $L^2$-Distance Approximation

$$L^2(p,p') = \int \left(p(x) - p'(x)\right)^2 dx$$
$$f(x) = p(x) - p'(x)$$

- **Naïve approximators via LSDD:**
  - $L^2(p,p') = \int f(x)^2 dx \approx \widehat{\theta}^\top H \widehat{\theta}$
  - $L^2(p,p') = \int f(x)\left(p(x) - p'(x)\right)dx \approx \widehat{h}^\top \widehat{\theta}$

- Consider their linear combination:
$$\beta\widehat{h}^\top\widehat{\theta} + (1-\beta)\widehat{\theta}^\top H\widehat{\theta}$$

- For small $\lambda$, Taylor expansion gives
$$\beta\widehat{h}^\top\widehat{\theta} + (1-\beta)\widehat{\theta}^\top H\widehat{\theta}$$
$$= \widehat{h}^\top H^{-1}\widehat{h} - \lambda(2-\beta)\widehat{h}^\top H^{-2}\widehat{h} + o_p(\lambda)$$

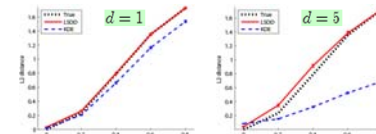- Bias caused by regularization can be eliminated by $\beta = 2$.
$$\widehat{L}^2(\mathcal{X},\mathcal{X}') = 2\widehat{h}^\top\widehat{\theta} - \widehat{\theta}^\top H\widehat{\theta}$$

## 4.2 Numerical Examples

- $L^2$-Distance Approximation  $n = n' = 100$
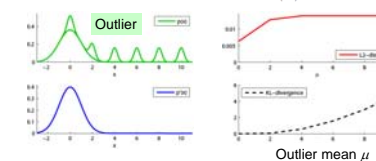$$p(x) = N(x; \mu, 0, \ldots, 0)^\top, (4\pi)^{-1}I_d)$$
$$p'(x) = N(x; 0, 0, \ldots, 0)^\top, (4\pi)^{-1}I_d)$$


$d = 1$ / $d = 5$

- $L^2$-Distance vs. KL-Divergence
$$L^2(p,p') = \int \left(p(x) - p'(x)\right)^2 dx$$
$$\mathrm{KL}(p\|p') = \int p(x)\log\frac{p(x)}{p'(x)}dx$$


Outlier

Outlier mean $\mu$

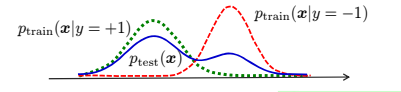$L^2$-distance is more robust than KL-divergence.

Cf. Basu, Harris, Hjort & Jones (Biometrika1998)
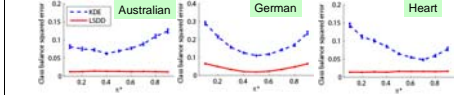
## 5.1 Class-Balance Estimation

- Pattern recognition when class balances are different in training and test phases.
- If test class-balance is known, weighted learning eliminates estimation bias.
- When test class-balance is unknown, fit mixture of training class-wise input densities to test input density:  Du Plessis & Sugiyama (ICML2012)
$$\min_{\pi\in[0,1]} L^2\left(p_{test}(x), q_\pi(x)\right)$$
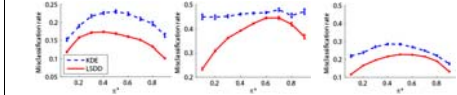$$q_\pi(x) = \pi p_{train}(x|y=+1) + (1-\pi)p_{train}(x|y=-1)$$


$p_{train}(x|y=+1)$, $p_{train}(x|y=-1)$, $p_{test}(x)$

- Class-balance estimation error  $\pi^*$: True test class-balance


Australian / German / Heart

- Classification accuracy:



Regularized kernel least-squares classifier with class-balance weighting is used.

## 5.2 Change Detection

- **Goal**: Find change points in time-series.
- Use $L^2$-distance between past and current data as change score:  Kawahara & Sugiyama (SADM2012)




CENSREC Speech Data / HASC Accelerometer Data