

# Inference and Estimation Using Nearest Neighbors

2019 The Second Korea-Japan Machine Learning Workshop  
2019. 2. 22 (Fri.)

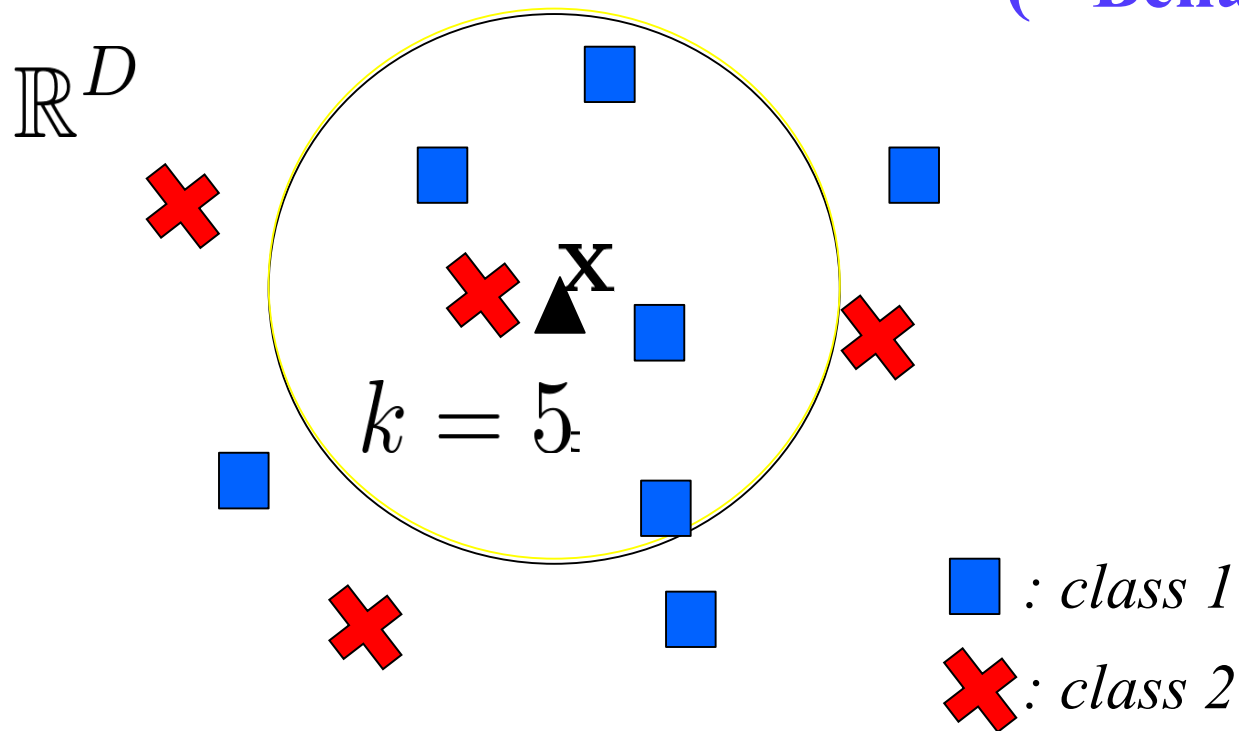
***Yung-Kyun Noh***

*Seoul National University → Hanyang University*



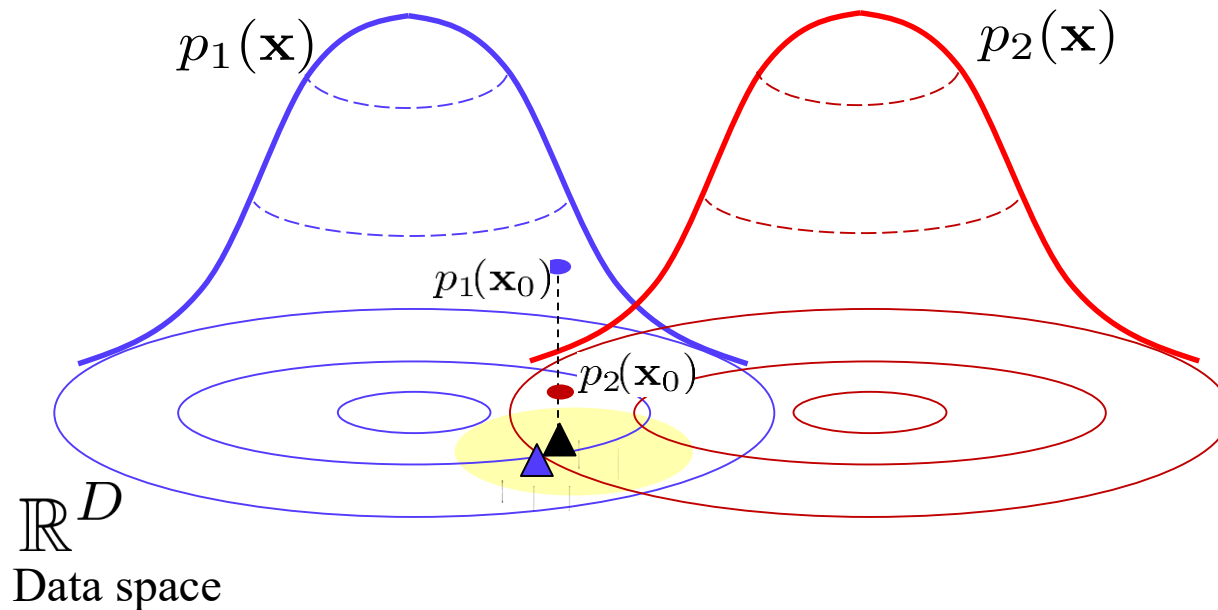
# Nearest Neighbors

- Similar data share similar properties (**= Labels?**)  
(**= Behavior**)



$\mathbf{x}_{NN} \rightarrow \mathbf{x}$ , uniformly with increasing  $N$ .

In the limit,  $p_1(\mathbf{x}) = p_1(\mathbf{x}_{NN})$  &  $p_2(\mathbf{x}) = p_2(\mathbf{x}_{NN})$



For classification as an example:

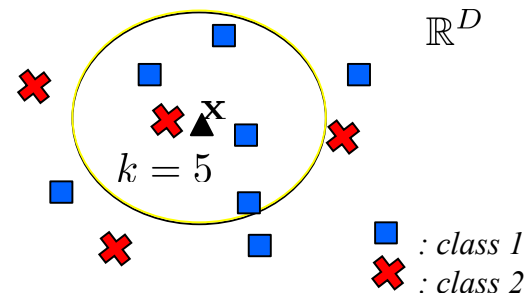
$$\epsilon_{(k=1)} = \int \frac{p_1(\mathbf{x})p_2(\mathbf{x})}{p_1(\mathbf{x}) + p_2(\mathbf{x})} d\mathbf{x} \leq 2E_{Bayes}(1 - E_{Bayes})$$

[T. Cover and P. Hart, *IEEE TIT*, 1967]

# Applications of Using Nearest Neighbors

- Prediction using  $k$ -Nearest Neighbor Information

- $k$ -Nearest Neighbor Classification
- $k$ -Nearest Neighbor Regression



- Estimation using  $k$ -Nearest Neighbor Information

[Leonenko, N., Pronzato, L., & Savani, V., 2008]

$$KL(p_1(\mathbf{x})||p_2(\mathbf{x})) = - \int p_1(\mathbf{x}) \log \frac{p_2(\mathbf{x})}{p_1(\mathbf{x})} d\mathbf{x}$$



$$\widehat{KL}(p_1(\mathbf{x})||p_2(\mathbf{x})) = \frac{1}{N_1} \sum_{\mathbf{x} \sim p_1} \log \frac{u_2(\mathbf{x})}{u_1(\mathbf{x})}$$

$$u_c = N_c d_c^D$$

$d_c$  is a distance to the nearest neighbor in class  $c$  from  $\mathbf{x} \in \mathbb{R}^D$ .

# Similar Formulations

- Nadaraya-Watson estimator for kernel classification/regression

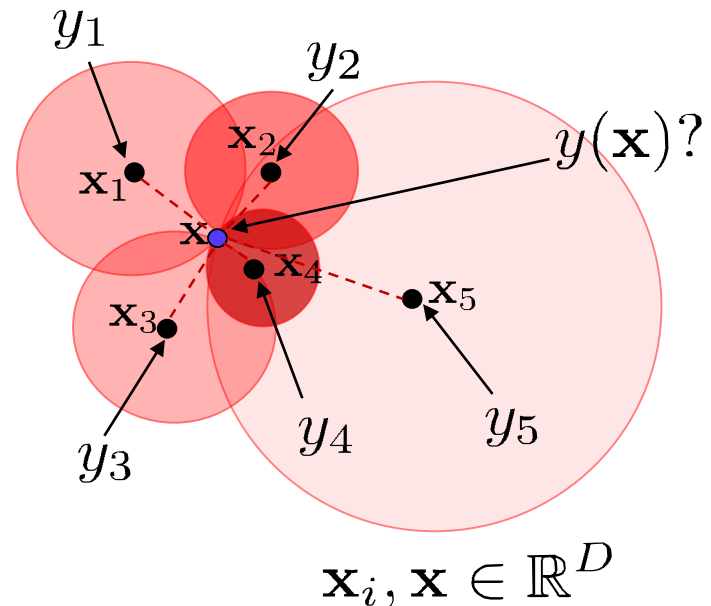
$$\hat{y}_N(\mathbf{x}) = \frac{\sum_{i=1}^N K(\mathbf{x}_i, \mathbf{x}) y_i}{\sum_{i=1}^N K(\mathbf{x}_i, \mathbf{x})}$$

$$\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^N$$

Kernel weight with respect to the distance

$$K(\mathbf{x}_i, \mathbf{x}) = K\left(\frac{\|\mathbf{x}_i - \mathbf{x}\|}{h}\right)$$

↖  
bandwidth



# Bias Analysis

- $k$ -Nearest Neighbor Classification

[R. R. Snapp et al. *The Annals of Statistics*, 1998]

[Y.-K. Noh et al. *IEEE TPAMI*, 2018]

$$E_{NN} \cong \int \frac{p_1(\mathbf{x})p_2(\mathbf{x})}{p_1(\mathbf{x}) + p_2(\mathbf{x})} d\mathbf{x} \quad \dots \textcircled{1}$$

$$+ \frac{1}{4D} \int \frac{\mathbb{E}_{d_N} [d_N^2 | \mathbf{x}]}{(p_1 + p_2)^2} [p_1^2 \nabla^2 p_2 + p_2^2 \nabla^2 p_1 - p_1 p_2 (\nabla^2 p_1 + \nabla^2 p_2)] d\mathbf{x} \quad \dots \textcircled{2}$$

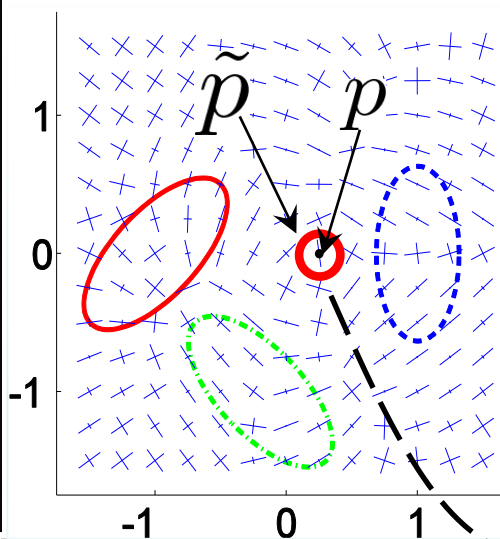
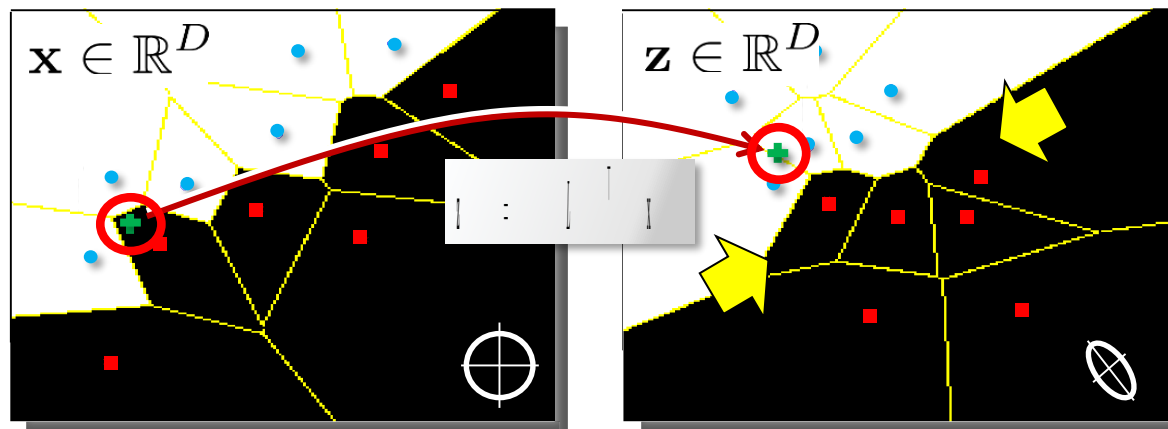
①: Asymptotic  $NN$  Error

②: Residual due to *Finite Sampling* .

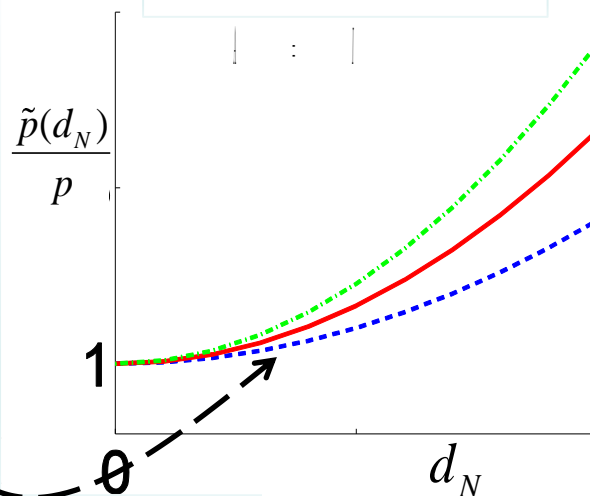
# Change of Metric

[Y.-K. Noh et al. *IEEE TPAMI*, 2018]

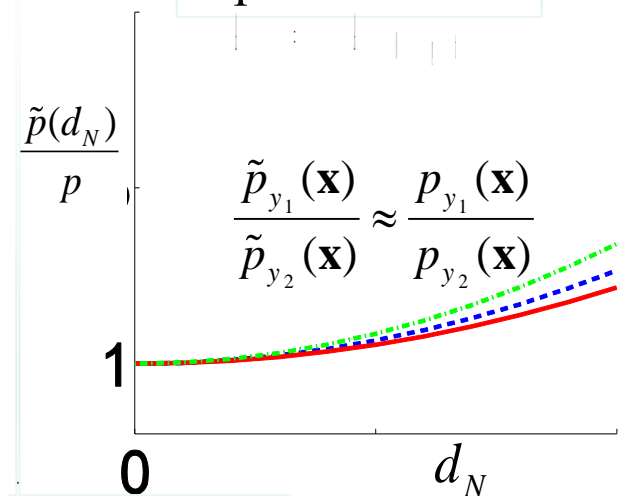
$$d(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{(\mathbf{x}_i - \mathbf{x}_j)^T A (\mathbf{x}_i - \mathbf{x}_j)}, \quad A = LL^T \succ 0$$



Euclidean metric



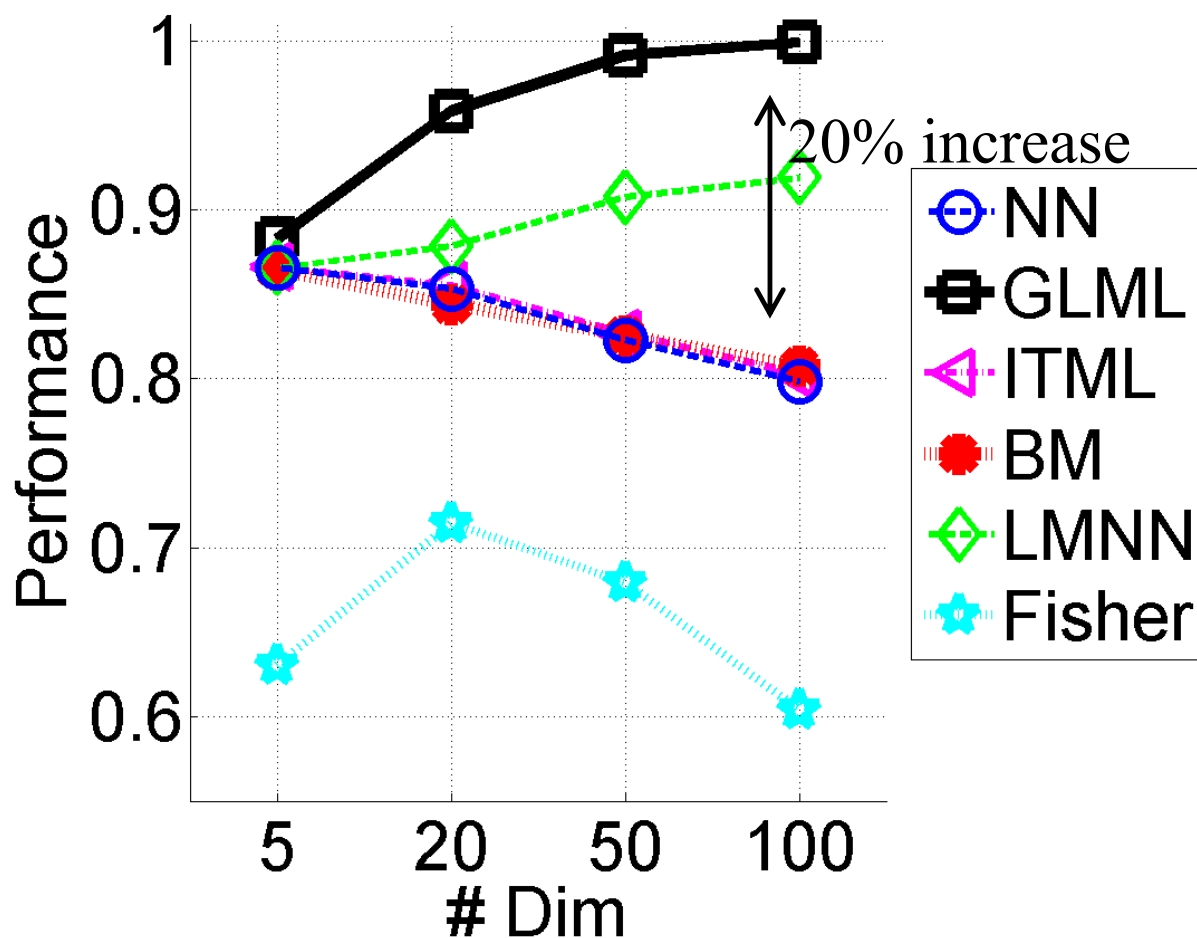
Optimal metric



# Nearest Neighbor Classification with Metric

$\nabla^2 p_1, \nabla^2 p_2, p_1, p_2$   $\leftarrow$  Obtain from generative models

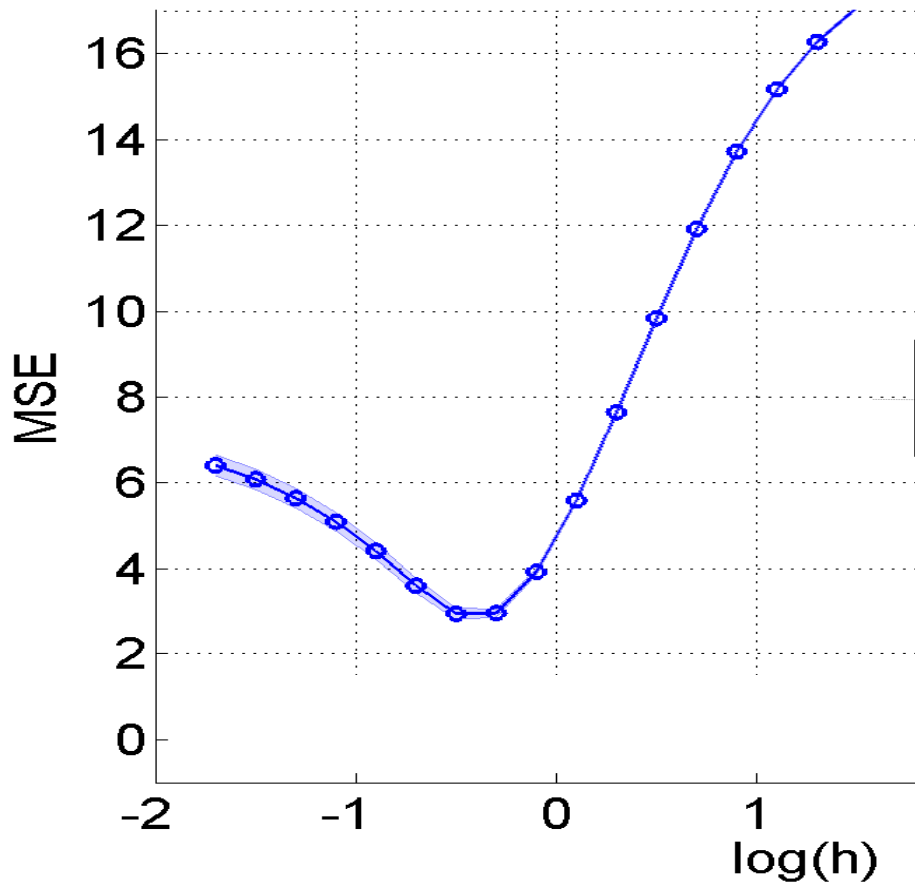
[Y.-K. Noh et al. *IEEE TPAMI*, 2018]





# Bandwidth and Nadaraya-Watson Regression

$$\hat{y}_N(\mathbf{x}) = \frac{\sum_{i=1}^N K_h(\mathbf{x}_i, \mathbf{x}) y_i}{\sum_{j=1}^N K_h(\mathbf{x}_j, \mathbf{x})}$$



# Bias Analysis

- $k$ -Nearest Neighbor Classification

$$\lim_{N \rightarrow \infty} \hat{y}_N(\mathbf{x}) = \mathbb{E}_{p(y|\mathbf{x})}[y] \quad (h \rightarrow 0)$$

→ Minimizes mean square error (MSE)

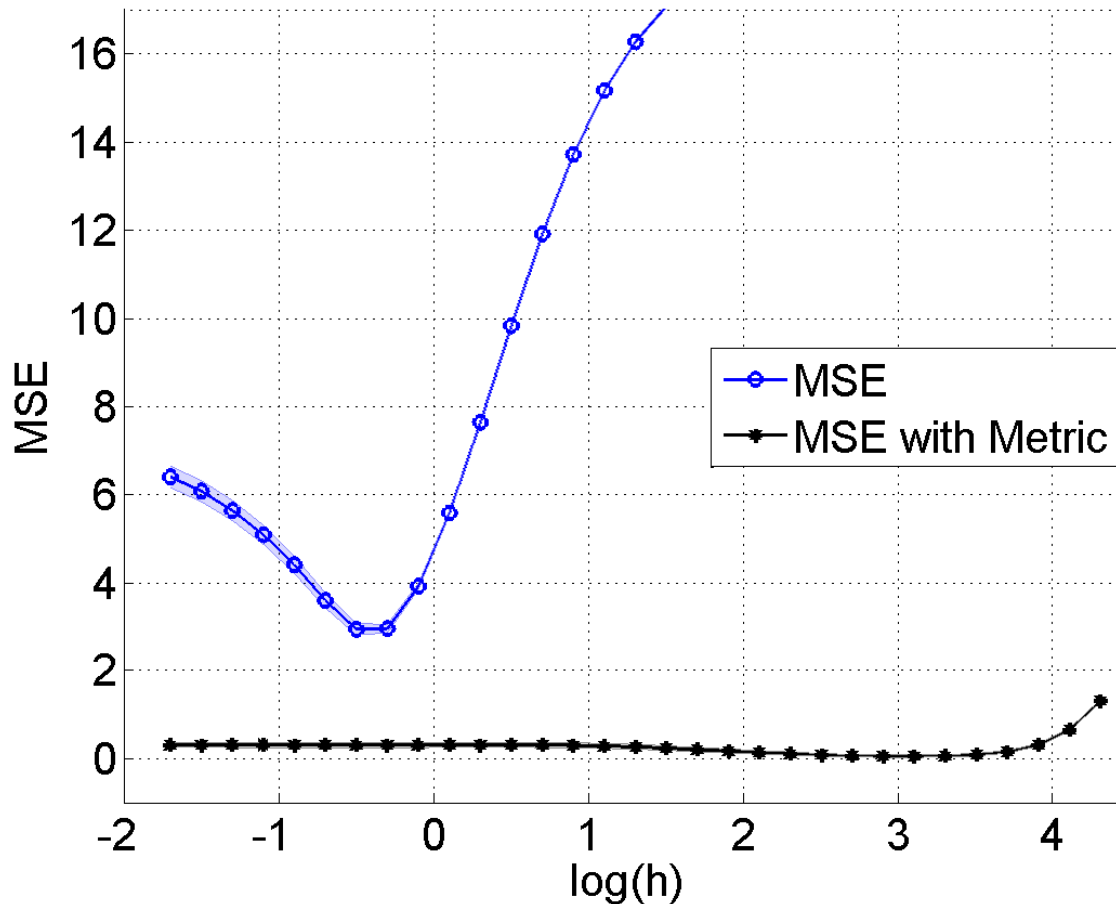
→ Metric independent asymptotic property

- Bias

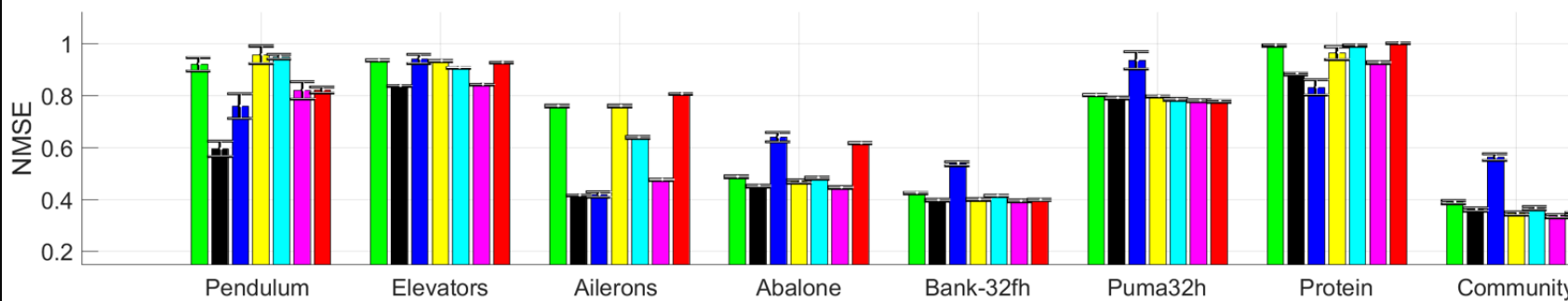
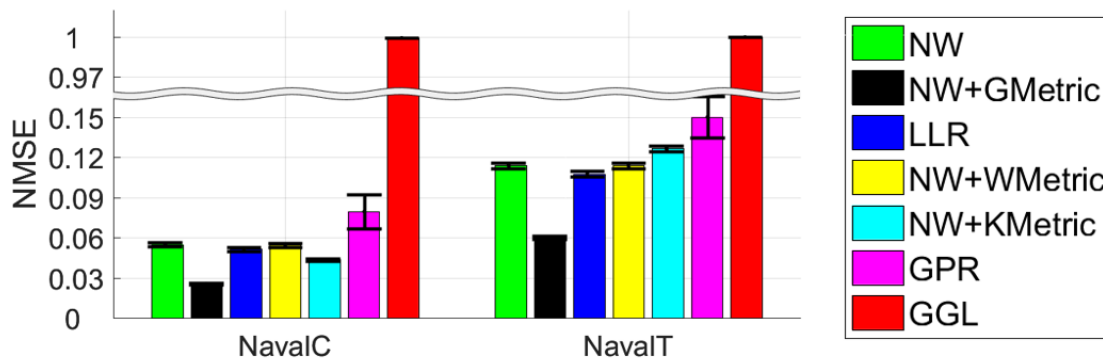
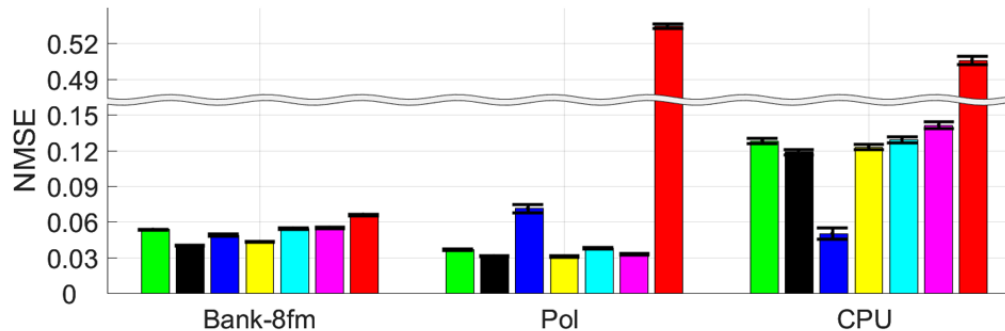
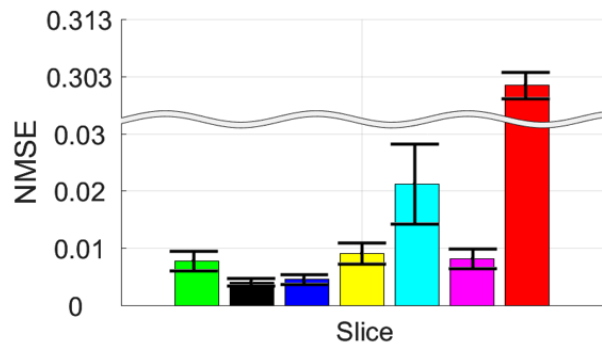
$$\mathbb{E} [\hat{y}(\mathbf{x}) - y(\mathbf{x})] = h^2 \left( \frac{\nabla^\top p(\mathbf{x}) \nabla y(\mathbf{x})}{p(\mathbf{x})} + \frac{\nabla^2 y(\mathbf{x})}{2} \right) + o(h^4)$$

# For $x$ & $y$ Jointly Gaussian

- Learned metric is not sensitive to the bandwidth



[Y.-K. Noh, et al., *NeurIPS*, 2017]



[Y.-K. Noh, et al., *NeurIPS*, 2017]

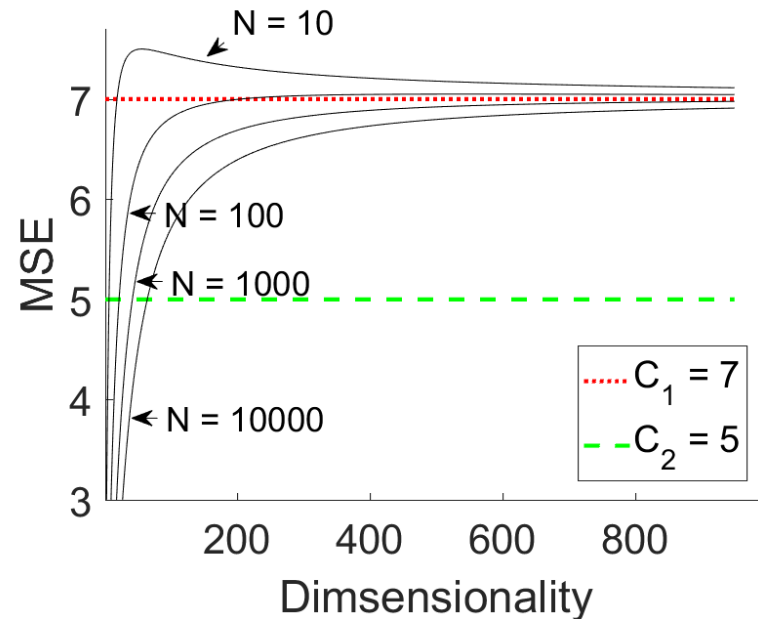
# Variance Reduction is Not Critical in High-Dimensions

[Y.-K. Noh, et al., *NeurIPS*, 2017]

$$f(h) = h^4 C_1 + \frac{1}{N h^D} C_2.$$

$$h_* = N^{-\frac{1}{D+4}} \left( \frac{D \cdot C_2}{4 \cdot C_1} \right)^{\frac{1}{D+4}}$$

$$\lim_{D \rightarrow \infty} f(h_*) = C_1$$



## Proposition

Reducing the variance is not important in a high dimensional space once the bias is minimized and the bandwidth selection is followed.

# Information-theoretic Measure Estimation

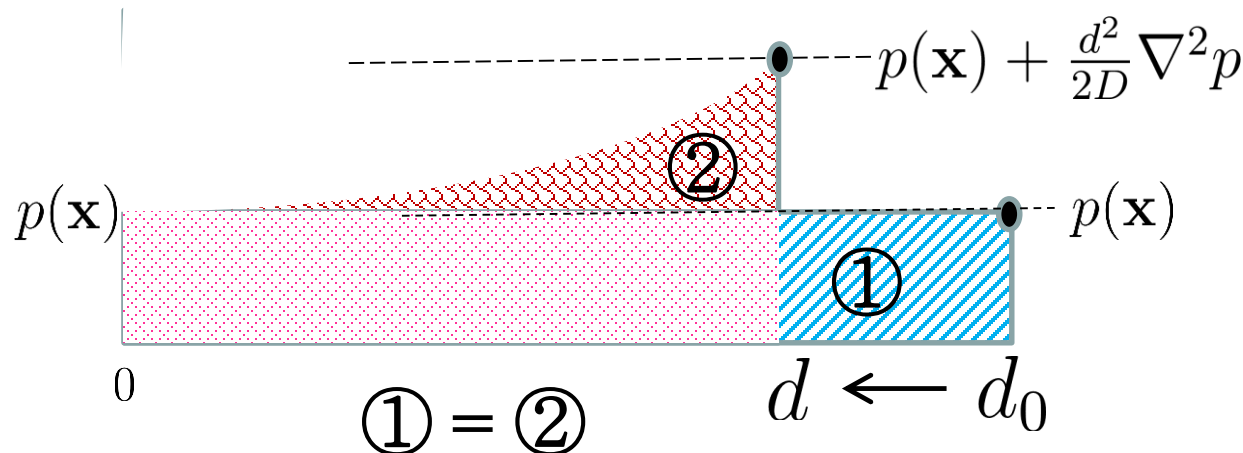
$u_c = N_c d_c^D$   $d_c$  is a distance to the nearest neighbor in class  $c$  from  $\mathbf{x}$ .

**Metric invariant**

$$\frac{1}{N_1} \sum_{\mathbf{x} \sim p_1} \log \frac{u_2}{u_1} \rightarrow - \int p_1 \log \frac{p_2}{p_1} d\mathbf{x}$$

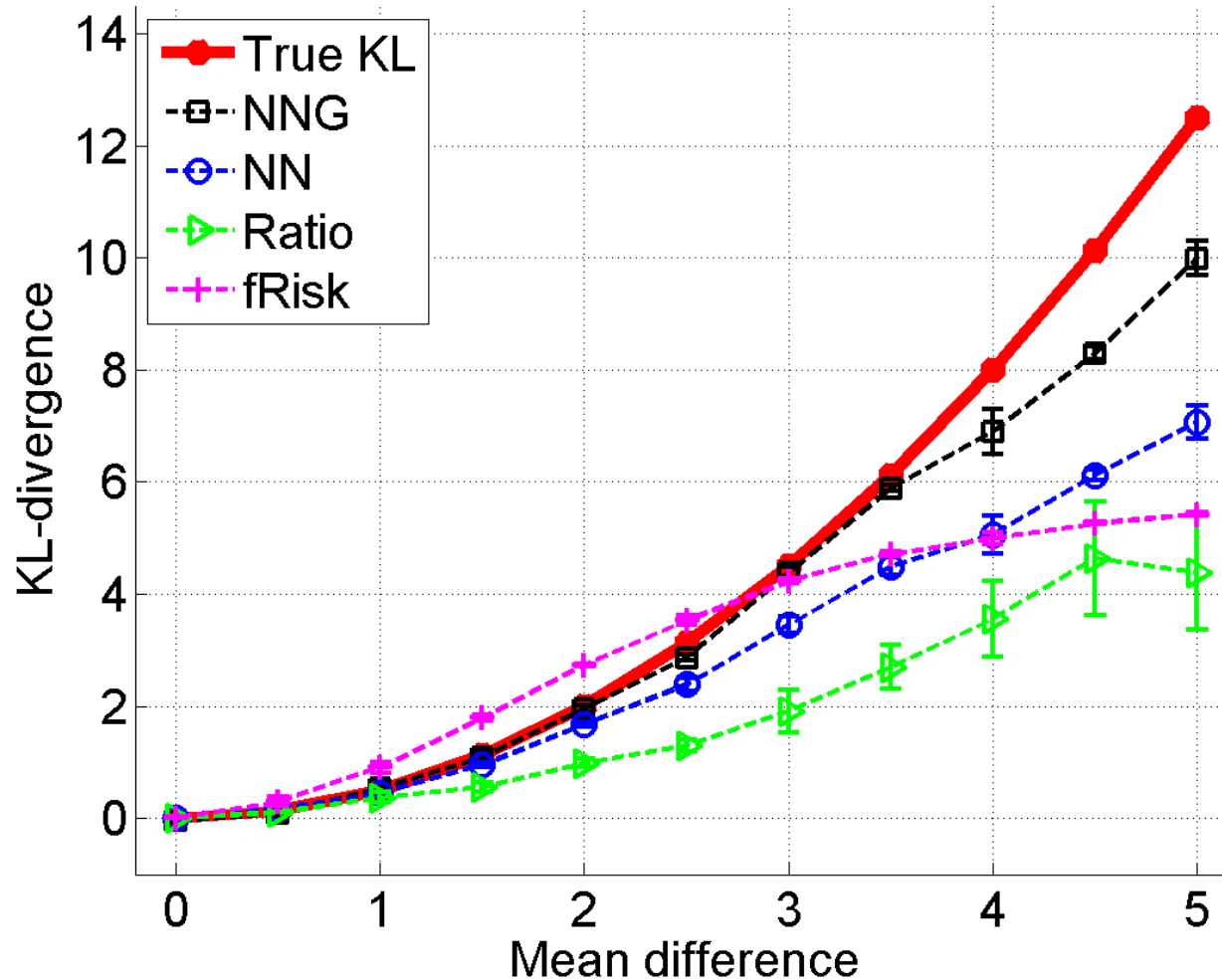
**Metric dependent**

$$+ \frac{1}{2D\gamma^{\frac{2}{D}}} \int \left[ (N_1 p_1)^{-\frac{2}{D}} \frac{\nabla^2 p_1}{p_1} - (N_2 p_2)^{-\frac{2}{D}} \frac{\nabla^2 p_2}{p_2} \right] d\mathbf{x}$$



# Increase the KL-Divergence of Two Gaussians and its Estimation

[Y.-K. Noh, et al., *NeCo*, 2018]



# MAKING GENERAL ESTIMATORS FOR F-DIVERGENCES



# Estimation of the General $f$ -Divergences

- Shannon Entropy Estimation

[D. Lombardi and S. Pant, *Phys. Rev. E*, 2016]

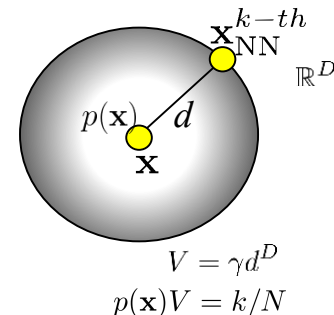
[A. Kraskov, H. Stögbauer, and P. Grassberger, *Phys. Rev. E*, 2004]

$$\hat{H}(X) = \psi(N) - \psi(k) + \log(\gamma) + \frac{1}{N} \sum_{i=1}^N \log d(\mathbf{x}_i)^D$$

$$\gamma = \frac{\pi^{\frac{D}{2}}}{\Gamma(1 + D/2)}, \quad \psi(t) = \Gamma(t)^{-1} \frac{d\Gamma(t)}{dt}$$

Note that  $p(\mathbf{x})V = k/N$

$$\text{In this case, } p(\mathbf{x}) = \frac{k}{VN} = \frac{k}{(\gamma d(\mathbf{x})^D) \cdot N}$$



# Density Estimator and Entropy Estimator

- Loftsgaarden and Quesenberry (1965)

$$\hat{p}(\mathbf{x}) = \frac{k}{\gamma N d(\mathbf{x})^D}$$

- Shannon Entropy Estimator

$$\begin{aligned}\hat{H}(X) &= \psi(N) - \psi(k) + \log(\gamma) + \frac{1}{N} \sum_{i=1}^N \log d(\mathbf{x}_i)^D \\ &= -\frac{1}{N} \sum_{i=1}^N \log \hat{p}(\mathbf{x}_i) + (\psi(N) - \log(N)) - (\psi(k) - \log(k))\end{aligned}$$

# Historical Remarks of Making Plug-in Estimators

[N. Leonenko, L. Pronzato, & V. Savani, *Annals of Statistics*, 2008]

[B. Póczos and J. Schneider, *AISTATS*, 2011]

Shannon entropy

$$\hat{H}(X) = -\frac{1}{N} \sum_{i=1}^N \log \hat{p}(\mathbf{x}_i) - \psi(k) \quad \text{Plug-in and correction}$$
$$\rightarrow - \int p(\mathbf{x}) \log p(\mathbf{x}) d\mathbf{x}$$

Rényi and Tsallis entropies

$$\hat{H}_q(X) = \frac{1}{N} \frac{\Gamma(k)}{\Gamma(k+1-q)} \sum_{i=1}^N (\hat{p}(\mathbf{x}_i))^{1-q} \quad \text{Plug-in and correction}$$
$$\rightarrow - \int p(\mathbf{x})^q d\mathbf{x}$$



[Moon, K. & Hero, A., 2014] considers the general  $f$ -divergence plug-in estimator

# Plug-in Nearest Neighbor $f$ -divergence Estimator

- Kullback-Leibler Divergence

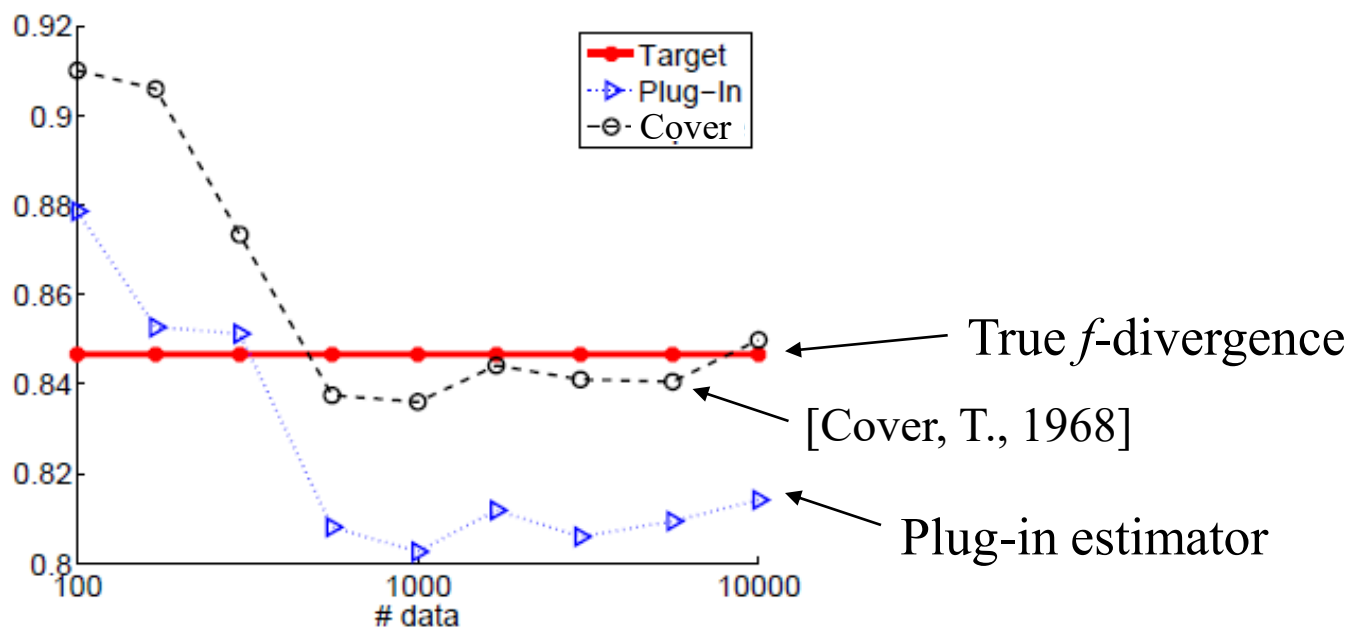
$$\frac{1}{N} \sum_{i=1}^N \log \hat{p}(\mathbf{x}_i) - \log \hat{q}(\mathbf{x}_i) \Rightarrow - \int p(\mathbf{x}) \log \left( \frac{q(\mathbf{x})}{p(\mathbf{x})} \right) d\mathbf{x}$$

- Tsallis-alpha Divergence

$$\frac{1}{\alpha - 1} \left( \frac{1}{N} \frac{\Gamma(k)^2}{\Gamma(k - \alpha + 1)\Gamma(k + \alpha - 1)} \sum_{i=1}^N \left( \frac{\hat{q}(\mathbf{x}_i)}{\hat{p}(\mathbf{x}_i)} \right)^{1-\alpha} - 1 \right) \\ \Rightarrow \frac{1}{\alpha - 1} \left( \int \left( \frac{p(\mathbf{x})}{q(\mathbf{x})} \right)^{\alpha} p(\mathbf{x}) d\mathbf{x} - 1 \right)$$

# Plug-in methods do not work for general $f$ -divergences

$$\frac{1}{N_p} \left( \sum_{i=1}^{N_p} \mathbf{1}(d_p > d_q) \right) \xrightarrow{[\text{Cover, T., 1968}]} \int \frac{p(\mathbf{x})q(\mathbf{x})}{p(\mathbf{x}) + q(\mathbf{x})} d\mathbf{x}$$



(b) Nonparametric estimator performance

[Noh, Y.-K. Ph.D. thesis, 2011]

# Obtaining the General $f$ -Divergence Estimator

$$X_{1:m} \sim p(\mathbf{x}), \quad Y_{1:n} \sim q(\mathbf{x})$$

$$\hat{T}(X_{1:m}, Y_{1:n}) = \frac{1}{m} \sum_{i=1}^m \phi(u(\mathbf{x}_i), v(\mathbf{x}_i)) \quad \begin{array}{l} u = m\gamma d_p^D \\ v = n\gamma d_q^D \end{array}$$

$$\rightarrow \int p(\mathbf{x}) f\left(\frac{p(\mathbf{x})}{q(\mathbf{x})}\right) d\mathbf{x}$$

$$\phi(u(\mathbf{x}_i), v(\mathbf{x}_i)) = \frac{(k-1)!(l-1)!}{u^{k-1}v^{l-1}} \mathcal{L}_{(u,v)}^{-1} \left[ \frac{f(s,t)}{s^k t^l} \right]$$

↑  
Inverse Laplace Transform

arXiv:1805.08342

---

## Nearest neighbor density functional estimation based on inverse Laplace transform

---

**Shouvik Ganguly<sup>†\*</sup> Jongha Ryu<sup>†\*</sup> Young-Han Kim<sup>†</sup> Yung-Kyun Noh<sup>‡</sup> Daniel D. Lee<sup>§</sup>**  
<sup>†</sup>University of California, San Diego   <sup>‡</sup>Seoul National University   <sup>§</sup>University of Pennsylvania  
<sup>†</sup>{shgangul, jongharyu, yhk}@ucsd.edu, <sup>‡</sup>nohyung@snu.ac.kr, <sup>§</sup>ddlee@seas.upenn.edu

# Summary

- Asymptotically, nearest neighbor methods are very nice. (*In terms of Theory!!*)
- With finite samples, bias treatment using geometry change can improve the conventional nonparametric methods significantly (in high-dimensional space).
- General and systematic way of obtaining  $f$ -divergence using nearest neighbor information.





Yung-Kyun Noh  
nohyung@snu.ac.kr