





Centre d'Enseignement et de Recherche en Technologies de l'Information et Systèmes

## Manifold-Adaptive Dimension Estimation

<u>Amir massoud Farahmand<sup>(1)</sup></u>, Csaba Szepesvári<sup>(1)</sup>, Jean-Yves Audibert<sup>(2)</sup>

Department of Computing Science, University of Alberta, Canada
(2) CERTIS, Ecole Nationale des Ponts, France







# High-Dimensional Data Everywhere

- Vision
- Sensor Fusion
- Feature Expansion
  - Kernel



## Curse of Dimensionality



# Practical Implications

- Thou shall reduce the dimension of the data before working with it
- Thou shall not ac lite a tres unnecessarily!
- Thou shall no accept projects with highdímensional data!
- ... !



0.5

0

-0.5 1.5

0.5

0

-0.5

0.5

-0.5

- Smoothness
- Sparsity
- Low noise at boundary

#### $\sqrt{Lower dimensional submanifold}$

- LLE, IsoMap, Laplacian Eigenmap, Hessian Eigenmap, ...
- Semi-supervised Learning, Reinforcement Learning, ...

## Goal

- Manifold-adaptive machine learning methods
- Convergence rate independent of the dimension of the input space

#### Many open questions! Here: dimension estimation (:

# Why?

- Needed in various learning methods
- Not known a priori

## New?

- Many existing methods [Pettis et al. (1979), Kegl (2002), Costa & Hero (2004), Levina & Bickel (2005), Hein & Audibert (2005)]
- No rigorous analysis
  - Asymptotic result [Levina & Bickel (2005)]

## Our Contribution

- New algorithm
  - K-NN
- Manifold-adaptive convergence rate

#### General Idea

 $P(X_i \in B(x,r)) = \eta(x,r)r^d$ 





 $\ln\left(P(X_i \in B(x,r))\right) = \ln(\eta(x,r)) + d\ln(r)$ 



#### $\ln\left(P(X_i \in B(x,r))\right) = \ln(\eta(x,r)) + d\ln(r)$





#### $\ln\left(P(X_i \in B(x,r))\right) = \ln(\eta(x,r)) + d\ln(r)$









 $\left| \left\langle \hat{r}(\left\lceil k/2 \right\rceil) \right\rangle$ 

 $\hat{d}(x) = \frac{\ln 2}{\ln(\hat{r}^k(x)/\hat{r}^{\lceil k/2 \rceil}(x))}$ 

$$\hat{d}(X_i) = \frac{\ln 2}{\ln(\hat{r}^{(k)}(X_i)/\hat{r}^{(\lceil k/2\rceil)}(X_i))}$$

**Theorem:** Under some regularity assumptions on  $\eta$ , provided that  $\frac{n}{k} > \Omega(2^d)$ , with probability at least  $1 - \delta$ ,

٠

$$|\hat{d}(X_i) - d| \le O\left(d\left[\left(\frac{k}{n}\right)^{\frac{1}{d}} + \sqrt{\frac{\ln(4/\delta)}{k}}\right]\right)$$

#### Issues

$$\hat{d}(X_i) = \frac{\ln 2}{\ln(\hat{r}^{(k)}(X_i)/\hat{r}^{(\lceil k/2 \rceil)}(X_i))}$$

#### High variance of $\hat{d}(X_i)$ Inefficient use of data $r \ll 1 \Longrightarrow k \ll n$

• Averaging 
$$\hat{d}_{avg} = \left[\frac{1}{n}\sum_{i=1}^{n}\hat{d}(X_i)\right]$$
  
• Voting  $\hat{d}_{vote} = \arg\max_{d'}\sum_{i=1}^{n}I\{[\hat{d}(X_i)] = d'\}$ 

**Theorem:** 

$$\mathbb{P}\left(\hat{d}_{\text{vote}} \neq d\right) \leq e^{-\frac{c'n}{(c^d k)^2}},$$
$$\mathbb{P}\left(\hat{d}_{\text{avg}} \neq d\right) \leq e^{-\frac{c''n}{(Dc^d k)^2}}$$

## Experiments

## Varying the Manifold Dimension



#### Varying Embedding Space Dimension



## Other Datasets

DATA SET	N = 50	N=100	N=500	N=1000	N = 5000
$S^1$	98 (99)	100(100)	100(100)	100(100)	100 (100)
$S^3$	75(19)	95(20)	100(15)	100(19)	100(62)
$S^5$	33~(5)	50(10)	100 (9)	98(2)	100(0)
$S^7$	18(2)	17(3)	57(1)	54(1)	100(0)
Sinusoid	92 (98)	100(100)	100(100)	100(100)	100(100)
10-Möbius	69 (47)	13 (74)	100 (98)	100 (99)	100 (100)
SWISS ROLL	62 (71)	49(91)	88 (96)	100 (100	100(100)

## Conclusions and Future Work

- New algorithm
- Competitive results
- Manifold-adaptive convergence rate
- Other ML methods?
- K-NN regression can!
- Penalized least squares in the works
- Dimension Reduction?





## Lower Bound

Assume that  $m_n$  is a regression function that estimate random variable Y based on X and  $D_n = \{(X_1, Y_1), ..., (X_n, Y_n)\}$ , and m(X) = E[Y|X]. What is the best possible performance of  $m_n$  in  $L_2$  sense, i.e.  $E\{||m_n(X) - m(X)||^2\}$ ?

For the class of  $D^{(p,C)}$  of (X,Y) distributions, when  $X \in \mathbb{R}^D$ , we have the the following behavior:

$$E\{\|m_n(X) - m(X)\|^2\} > O\left(n^{-\frac{2p}{2p+D}}\right)$$

Two sources of error:

- Approximation Error: assuming fixed  $\eta(x, r)$
- Estimation Error: estimating  $P(X \in B(x, r))$  with the empirical estimate k/n.

Both of them can be controlled by changing the size of neighborhood r (which is related to k/n).

#### Effect of k and n









## Experiments Noise Effect











## **Exponential Rate**

