

# Optimal nonparametric classification via radial distance

A. Okuno<sup>1,2</sup>, R. Cao<sup>3,2</sup>, K. Nakagawa<sup>4</sup>, H. Shimodaira<sup>3,2</sup>

<sup>1</sup>Institute of Statistical Mathematics    <sup>2</sup>RIKEN AIP    <sup>3</sup>Kyoto University

<sup>4</sup>Nomura Asset Management Co., Ltd.

## Related Publications of Ours

*In English:*

- (1) Akifumi Okuno and Hidetoshi Shimodaira. “Extrapolation Towards Imaginary 0-Nearest Neighbour and Its Improved Convergence Rate”, Advances in Neural Information Processing Systems 33 (**NeurIPS 2020**), pages 21889-21899.
- (2) Ruixing Cao\*, Akifumi Okuno\*, Kei Nakagawa and Hidetoshi Shimodaira. “Improving Nonparametric Classification via Local Radial Regression with an Application to Stock Prediction”, 23 pages. **arXiv:2112.13951**. (\*First co-author)

*In Japanese (JSAI Annual Proceedings):*

- (3) Ruixing Cao, Takuma Tanaka, Akifumi Okuno, Hidetoshi Shimodaira. “A Study on Regression and Loss Functions for Multiscale  $k$ -Nearest Neighbour”. 2021. 4pages.
- (4) Takuma Tanaka, Akifumi Okuno, Hidetoshi Shimodaira. “Extreme Multi-Label Classification of Images via Multiscale  $k$ -Nearest Neighbour”. 2021. 4pages.

We had no opportunity to make a presentation for the study conducted 2 years ago...

# Background

## Regression/classification

Let  $(X, Y) \in \mathcal{X} \times \mathbb{R}$  be a pair of covariate and response variables, defined with

$$f(x) = \mathbb{E}(Y \mid X = x).$$

In the regression problem, we estimate the function  $f$  from i.i.d. observations  $\{(x_i, y_i)\}_{i=1}^n$ .

- ▶  $f_{\theta}(x) = \langle \theta_1, x \rangle + \theta_2$  (linear regression),
- ▶  $f_{\theta}(x) = \langle \theta_3, \sigma(\langle \theta_1, x \rangle + \theta_2) \rangle + \theta_4$  (neural network), and so forth.

### Problem in this study:

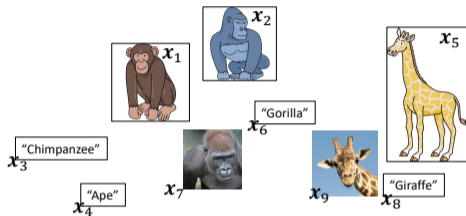
In real-world situations, several different forms of covariates are mixed. For instance,

$$\mathcal{X} = \bigcup_{q \in \mathbb{N}} \mathbb{R}^q;$$

typical regression functions cannot be applied to both  $x_1 \in \mathbb{R}^{q_1}$  and  $x_2 \in \mathbb{R}^{q_2}$  simultaneously.

# Examples

## (Example 1) Mixture of pictures/drawings/texts

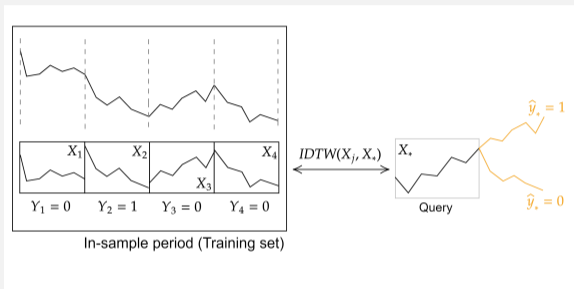


$$d(\text{ape photo}, \text{gorilla cartoon}) = 0.1,$$

$$d(\text{ape photo}, \text{"Giraffe"}) = 1.3, \dots$$

- ▶  $x_j$ : an image or a word,
- ▶  $y_i$ : whether the object represents primates (霊長類 in Japanese) or not.

## (Example 2) Time series of different lengths



- ▶  $\mathbf{x}_i$ : stock price record in  $i$ th month ( $\in \mathbb{R}^{\text{days in } i\text{th month}}$ ),
- ▶  $y_i$ : whether the price increases ( $y_i = 1$ ) or decreases ( $y_i = 0$ ).

## Distance-based approaches

Let  $\mathcal{X} = \mathbb{R}^q$ ,  $h > 0$  and consider a kernel smoother

$$\hat{f}_h^{(\text{KS})}(\mathbf{x}) = \frac{1}{|\mathcal{N}_h(\mathbf{x})|} \sum_{j \in \mathcal{N}_h(\mathbf{x})} y_j, \quad \mathcal{N}_h(\mathbf{x}) := \{i \mid d(\mathbf{x}, \mathbf{x}_i) \leq h\}.$$

Then, the distance  $d$  can be replaced with other discrepancy functions formally. For instance, if  $\{\mathbf{x}_i\}$  represents the time-series of different lengths, we may employ

$$d(\mathbf{x}, \mathbf{x}') := \text{DynamicTimeWarping}(\mathbf{x}, \mathbf{x}').$$

 Distance-based approaches are (formally) widely-applicable.

Kernel smoother with  $k = |\mathcal{N}_h(\mathbf{x})|$  is called *k-nearest neighbour (k-NN)* estimator.

## (Higher-order) asymptotic bias and its correction

Consider the simple case  $\mathcal{X} = \mathbb{R}^q$ . While the kernel smoother and  $k$ -NN estimators are *consistent*, i.e.,

$$\hat{f} \rightarrow^P f,$$

they are *not minimax optimal* if  $f$  is highly-smooth. Conventional local polynomial (LPoR) estimator corrects the asymptotic bias.

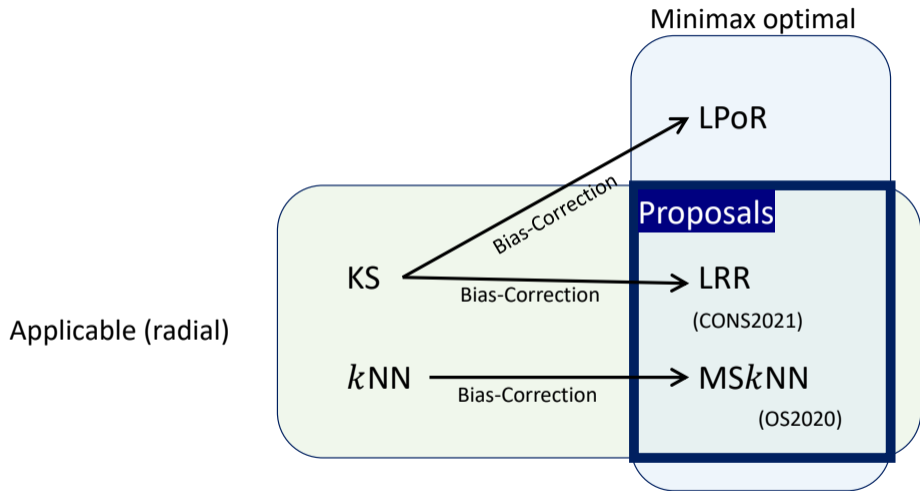
- ▶ KS and  $k$ NN: 👍 *widely-applicable*, 👎 *not optimal*.
- ▶ LPoR: 👎 *not widely-applicable*, 👍 *optimal*.

**Problem:** can we correct the asymptotic bias while holding the applicability?



# Proposal

# Overview

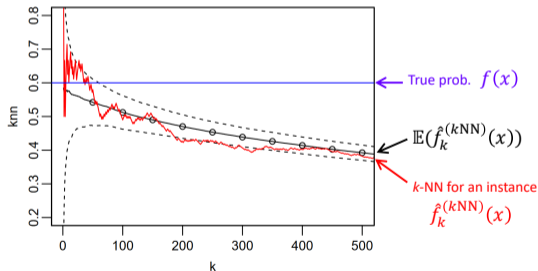


# Multiscale $k$ -NN (Okuno and Shimodaira, NeurIPS2020)

$k$ -NN estimator has larger bias as  $k$  increases.

## Idea

Extrapolating  $k$ -NN estimators from  $k = k_1, k_2, \dots, k_V$  to  $r = 0$  (via  $r_k := d(\mathbf{x}, \mathbf{x}_{(k)})$ ) yields imaginary 0-NN estimator, which is also called *multiscale  $k$ -NN (MSkNN) estimator*.



👍 *minimax optimal* and *widely-applicable*.

# Local Radial Regression (Cao, Okuno, Nakagawa and Shimodaira)

We define a *local radial regression (LRR)*:

$$\hat{f}^{(\text{LRR})}(\mathbf{x}) = \hat{\tau}(0), \quad \hat{\tau} := \arg \min_{\tau \in \mathcal{P}(1,q)} \sum_{i=1}^n w(r_i) \{Y_i - \tau(r_i)\}^2,$$

equipped with the radial distance  $r_i = d(\mathbf{x}, \mathbf{x}_i)$ , decreasing non-negative function  $w$  and a polynomial function  $\tau$  to be trained.













 *minimax optimal* and *widely-applicable*.

# Theory

We prove the convergence rate for the plug-in type classifier

$$\hat{g}(x) := \mathbb{1}(\hat{f}(x) \geq 1/2) \in \{0, 1\}.$$

Going through a very bothersome calculation to prove the optimality, we have:

	Abbrev.	Application	Higher-order opt.
Multi-layer perceptron	MLP/NN	 limited	 No
Kernel smoother	KS	 wide	 No
$k$ -Nearest neighbour	$k$ NN	 wide	 No
Local polynomial regression	LPoR	 limited	 Optimal
Multiscale $k$ NN (OS2020)	MS $k$ NN	 wide	 Optimal
Local radial regression (CONS2021)	LRR	 wide	 Optimal

cf. optimal rate is  $\mathcal{E}(\hat{g}_n) \asymp n^{-2\beta/(2\beta+d)}$  when assuming  $\beta$ -Hölder condition.

## Proof overview

The optimality is shown with  $d(\mathbf{x}, \mathbf{x}') = \|\mathbf{x} - \mathbf{x}'\|_2$ . For multiscale  $k$ -NN, the  $k$ -NN estimator satisfies

$$\begin{aligned}\hat{f}_k(\mathbf{x}_*) &\approx \mathbb{P}(Y = 1 \mid \|X - \mathbf{x}_*\|_2 \leq r) \\ &= f(X_*) + \sum_{c=1}^{\lfloor \beta/2 \rfloor} \underbrace{b_c^*(\mathbf{x}_*) r^{2c}}_{(*)} + \delta_{\beta,r}(\mathbf{x}_*), \quad |\delta_{\beta,r}(\mathbf{x}_*)| \lesssim r^\beta,\end{aligned}$$

for  $r = r(k) := \|\mathbf{x}_* - \mathbf{x}_{(k)}\|_2$  and large  $k \in \mathbb{N}$ . Therefore, multiscale  $k$ -NN estimator removes the higher-order bias term  $(*)$  by the extrapolation (via regression).

□

# Experiments

**Table:**  $n$ : sample size,  $d$ : dimension,  $m$ : #categories.

Sample average and the standard deviation for the prediction accuracy are computed on 10 times experiments. Best scores are **bolded**, and second best scores are underlined.

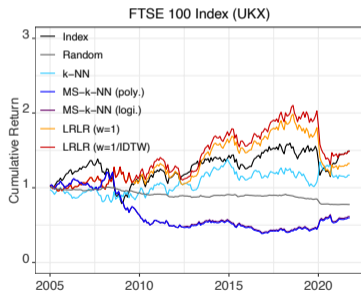
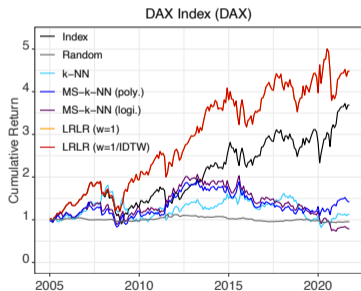
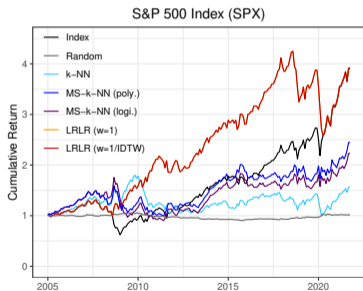
Dataset	$n$	$d$	$m$	kNN			MSkNN	
				$w_i = 1/k$	$w_i \geq 0$	$w_i \in \mathbb{R}$	via $r(k)$	via $\log k$
Iris	150	4	3	0.83 $\pm$ 0.04	0.92 $\pm$ 0.05	0.92 $\pm$ 0.04	<u>0.93</u> $\pm$ 0.04	<b>0.96</b> $\pm$ 0.04
Glass identification	213	9	6	0.58 $\pm$ 0.06	<u>0.64</u> $\pm$ 0.06	<b>0.67</b> $\pm$ 0.05	<u>0.64</u> $\pm$ 0.05	<u>0.64</u> $\pm$ 0.05
Ecoli	335	7	8	0.80 $\pm$ 0.03	<b>0.85</b> $\pm$ 0.03	<u>0.84</u> $\pm$ 0.02	<b>0.85</b> $\pm$ 0.02	<u>0.84</u> $\pm$ 0.02
Diabetes	768	8	2	<b>0.75</b> $\pm$ 0.03	<u>0.74</u> $\pm$ 0.03	0.70 $\pm$ 0.04	<b>0.75</b> $\pm$ 0.03	0.71 $\pm$ 0.03
Biodegradation	1054	41	2	<u>0.84</u> $\pm$ 0.02	<b>0.86</b> $\pm$ 0.03	0.79 $\pm$ 0.02	<b>0.86</b> $\pm$ 0.02	0.80 $\pm$ 0.02
Banknote	1371	4	2	0.95 $\pm$ 0.01	<u>0.98</u> $\pm$ 0.01	0.97 $\pm$ 0.01	<u>0.98</u> $\pm$ 0.01	<b>0.99</b> $\pm$ 0.00
Yeast	1484	8	10	<u>0.57</u> $\pm$ 0.02	<b>0.58</b> $\pm$ 0.02	0.54 $\pm$ 0.03	<b>0.58</b> $\pm$ 0.02	0.54 $\pm$ 0.02
Wireless localization	2000	7	4	<u>0.97</u> $\pm$ 0.00	<b>0.98</b> $\pm$ 0.00	<b>0.98</b> $\pm$ 0.01	<b>0.98</b> $\pm$ 0.00	<b>0.98</b> $\pm$ 0.01
Spambase	4600	57	2	<u>0.90</u> $\pm$ 0.01	<b>0.91</b> $\pm$ 0.00	0.86 $\pm$ 0.01	<b>0.91</b> $\pm$ 0.00	0.87 $\pm$ 0.01
Robot navigation	5455	24	4	0.81 $\pm$ 0.01	<b>0.86</b> $\pm$ 0.01	0.81 $\pm$ 0.01	<u>0.84</u> $\pm$ 0.01	<u>0.84</u> $\pm$ 0.01
Page blocks	5473	10	5	<u>0.95</u> $\pm$ 0.01	<u>0.95</u> $\pm$ 0.01	<b>0.96</b> $\pm$ 0.01	<b>0.96</b> $\pm$ 0.01	<b>0.96</b> $\pm$ 0.01
MAGIC	19020	10	2	0.82 $\pm$ 0.00	0.82 $\pm$ 0.00	<b>0.84</b> $\pm$ 0.01	<u>0.83</u> $\pm$ 0.00	<u>0.83</u> $\pm$ 0.00
Avila	20867	10	12	0.63 $\pm$ 0.01	0.68 $\pm$ 0.01	<b>0.70</b> $\pm$ 0.01	<u>0.69</u> $\pm$ 0.00	<b>0.70</b> $\pm$ 0.01



# Application to Stock Prediction

MSkNN and LRR are applied to the stock prediction problem of S&P500, S&P500/TSX, EURO, ...

- ▶  $\mathbf{x}_i$ : stock price record in  $i$ th month ( $\in \mathbb{R}^{\text{“days in } i\text{th month”}}$ ),
- ▶  $y_i$ : whether the price increases ( $y_i = 1$ ) or decreases ( $y_i = 0$ ).















# Application to Stock Prediction

**Table:** Predictive classification accuracy. A higher score is better: the best and the second-best are bolded and underlined, respectively.

	S&P 500	S&P/TSX	EURO.	FTSE.	DAX	CAC.	TOPIX	Hang Seng
random	0.492	0.495	0.498	0.482	0.492	0.490	0.493	0.486
$k$ -NN	0.574	<u>0.594</u>	<u>0.510</u>	<u>0.500</u>	0.530	<u>0.525</u>	<u>0.500</u>	<u>0.564</u>
MSkNN	<u>0.604</u>	0.559	<b>0.525</b>	0.485	<u>0.545</u>	0.495	<b>0.515</b>	0.530
LRLR	<b>0.649</b>	<b>0.609</b>	0.505	<b>0.574</b>	<b>0.609</b>	<b>0.550</b>	0.465	<b>0.574</b>

# Conclusion

# Conclusion

	Abbrev.	Applicability	Optimality
Multi-layer perceptron	MLP/NN	 limited	 No
Kernel smoother	KS	 Good	 No
$k$ -Nearest neighbour	$k$ NN	 Good	 No
Local polynomial regression	LPoR	 Less	 Optimal
Multiscale $k$ NN (OS2020)	MSkNN	 Good	 Optimal
Local radial regression (CONS2021)	LRR	 Good	 Optimal

- ▶ We proposed a *widely-applicable* and *optimal* MSkNN and LRR estimators.
- ▶ This study was mainly based on:
  - (1) Okuno and Shimodaira (NeurIPS2020)
  - (2) Cao, Okuno, Nakagawa, and Shimodaira (arXiv:2112.13951)
- ▶ Contact Info.: [okuno@ism.ac.jp](mailto:okuno@ism.ac.jp)