

Gradient Weights help Nonparametric Regressors

Samory Kpotufe

Abdeslam Boularias

Toyota Technological Institute, Chicago

Max Planck Institute for Intelligent Systems, Tübingen

1 Overview

- Motivation: f does not depend equally on all dimensions, take advantage of this.
- Gradient-weighting applies to any distance-based regressor, is simple and efficient.
- Experiments on real-world problems show significant improvement in performance.

2 Preliminaries

Nonparametric regression setup:

Assume $Y = f(X) + \text{noise}$. Estimate $f(x)$ at x from sample $\{X_i, Y_i\}_1^n$.

Distance based regression:

$f_n(x) = \sum_{i=1}^n w(x, X_i) \cdot Y_i$, where $w(x, X_i)$ depends on metric ρ (usually Euclidean).

Ex: k -NN, kernel, local polynomial regressors.

3 A simple idea:

f might not depend equally on all features of X !

Let $f'_i \equiv$ derivative along coordinate i , let $\|f'_i\|_{1,\mu} \equiv \mathbb{E}_X |f'_i(X)|$.

In practice $\{\|f'_i\|_{1,\mu}\}_{i \leq d}$ has small and large elements.

Gradient weighting:

Reweight $x^i \rightarrow \sqrt{\mathbf{W}_i} \cdot x^i$, where $\mathbf{W}_i \approx \|f'_i\|_{1,\mu}$.

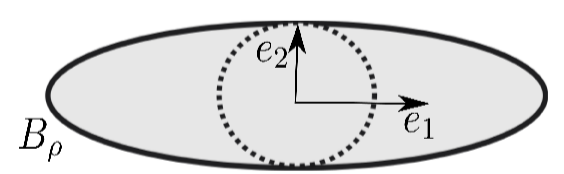
Equivalently, replace Euclidean distance with $\rho(x, x') = \sqrt{(x - x')^\top \mathbf{W} (x - x')}$.

Now run any local regressor f_n on (\mathcal{X}, ρ) .

Similar to metric learning, but cheaper: we estimate a single ρ !

More general than feature selection: f could depend on all $[d]$.

4 Intuition: suppose $\|f'_2\|_{1,\mu} \gg \|f'_1\|_{1,\mu}$.



Lower $\text{Var}(f_n(x))$: suppose $\|f'_i\|_{1,\mu}$ is large only for $i \in R \subsetneq [d]$.

- Balls B_ρ contain more points relative to Euclidean balls.
- Formally, $\mu(B_\rho(x, \rho(\mathcal{X}) \cdot \epsilon)) \approx e^{|\mathcal{R}|} \gg \epsilon^d$.

Bias of $f_n(x)$ is relatively unaffected:

- No sample is too far from x in any relevant direction.
- Formally, f has the following Lipschitz property:

$$|f(x) - f(x')| \leq \left(\sum_{i \in R} \frac{|f'_i|_\infty}{\sqrt{\mathbf{W}_i}} \right) \rho(x, x').$$

5 Gradient norm estimator: estimate $\|f'_i\|_{1,\mu}$ and not f'_i

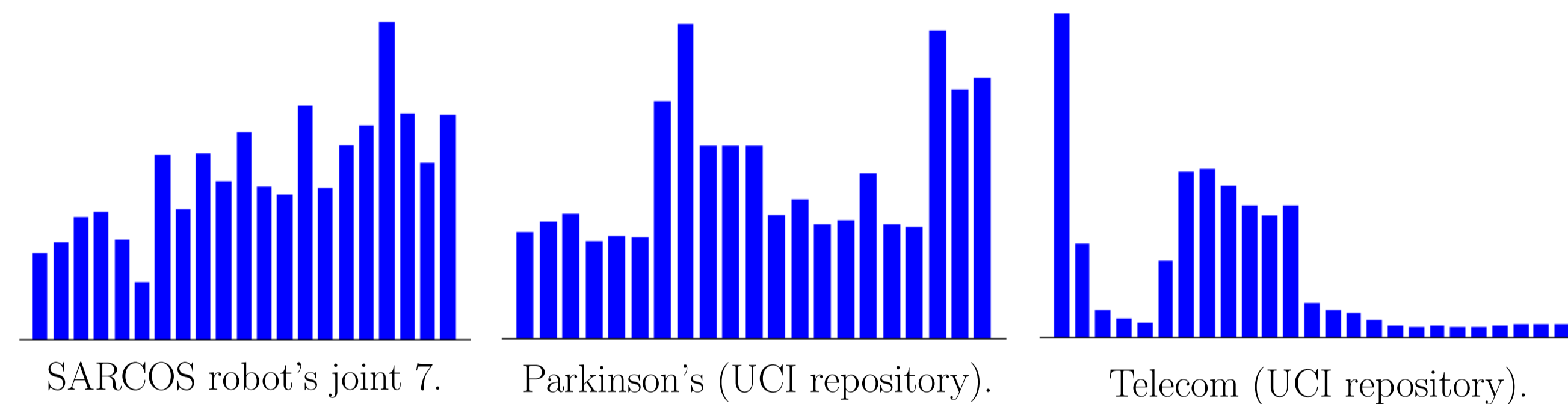
$$\mathbf{W}_i \triangleq \mathbb{E}_n \frac{|f_{n,h}(X + te_i) - f_{n,h}(X - te_i)|}{2t} \cdot 1_{\{A_{n,i}(X)\}},$$

where $A_{n,i}(X) \equiv$ we are confident in both estimates $f_{n,h}(X \pm te_i)$.

- Fast preprocessing, and online: just 2 estimates of $f_{n,h}$ at X . Metric learning optimizes over a space of possible metrics.
- Only $2 \times 2 \times 1$ param. search grid to adapt to d dimensions. KR with bandwidths $\{h_i\}_1^d$ needs $d \times d$ param. search grid.
- General: preprocessing for any distance-based regressor. Other methods apply to particular regressors, e.g. Rodeo for local linear regression, metric learning for KR.

6 Practicality of the approach

On typical real-world data, $\|f'_i\|_{1,\mu}$ varies a lot!



7 Consistency of the gradient weights estimator

Distributional assumptions:

- \mathcal{X} has bounded diameter 1 and μ has mass everywhere on \mathcal{X} : $\forall x \in \mathcal{X}, \forall h > 0, \mu(B(x, h)) \geq C_\mu h^d$.
- f is continuously differentiable on the τ -envelope $\mathcal{X} + B(0, \tau)$, f'_i is uniformly continuous on $\mathcal{X} + B(0, \tau)$ and $\sup_{x \in \mathcal{X} + B(0, \tau)} |f'_i(x)| \leq |f'_i|_{\text{sup}}$.

Theorem 1. Under general regularity conditions on μ , and smoothness conditions on ∇f , we have with probability $\geq 1 - \delta$:

$$\begin{aligned} |\mathbf{W}_i - \|f'_i\|_{1,\mu}| &\leq \frac{1}{t} \left(\sqrt{\frac{A(n)}{nh^d}} + h \cdot \sum_{i \in [d]} |f'_i|_{\text{sup}} \right) \\ &\quad + 2 |f'_i|_{\text{sup}} \left(\sqrt{\frac{\ln 2d/\delta}{n}} + \mu(\partial_{t,i}(\mathcal{X})) \right) + \epsilon_{t,i}. \end{aligned}$$

Boundary parameter: $\partial_{t,i}(\mathcal{X}) \triangleq \{x : \{x + te_i, x - te_i\} \not\subseteq \mathcal{X}\}$.

Smoothness parameter: $\epsilon_{t,i} \triangleq \sup_{x \in \mathcal{X}, s \in [-t, t]} |f'_i(x) - f'_i(x + se_i)|$.

Note that, under distributional assumptions, $\mu(\partial_{t,i}(\mathcal{X})) \xrightarrow{t \rightarrow 0} 0$ and $\epsilon_{t,i} \xrightarrow{t \rightarrow 0} 0$. Thus, if for example $h = C/\log^2 n$, $t = \sqrt{h}$, we have $\mathbf{W}_i \xrightarrow{P} \|f'_i\|_{1,\mu}$ for all $i \in [d]$.

8 Experiments

We tested the performance of both Kernel Regression (KR) and k -NN without using the gradient weights, and with the gradient weights (KR- ρ and k -NN- ρ).

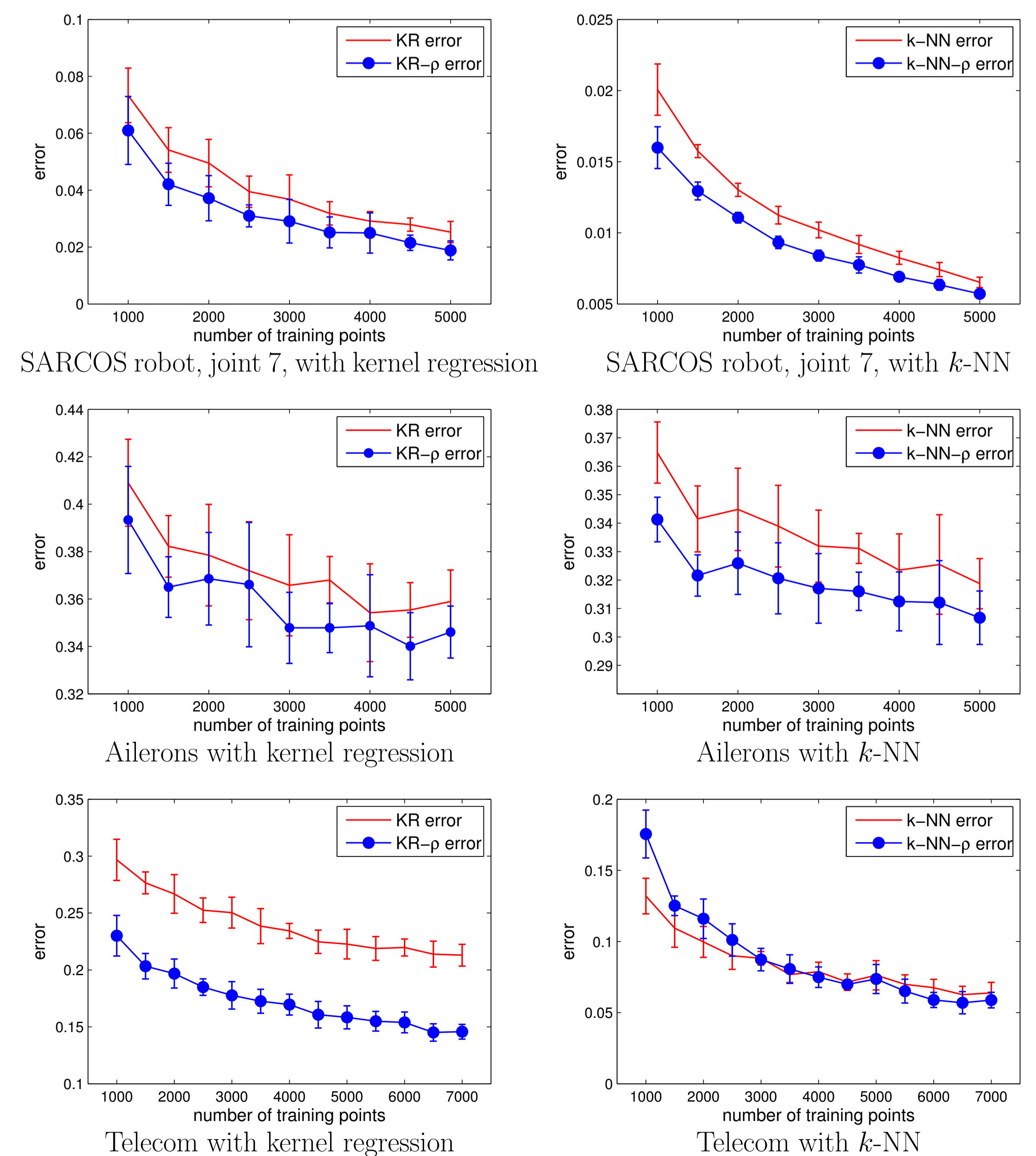
	Barrett joint 1	Barrett joint 5	SARCOS joint 1	SARCOS joint 5	Housing
KR error	0.50 \pm 0.02	0.50 \pm 0.03	0.16 \pm 0.02	0.14 \pm 0.02	0.37 \pm 0.08
KR- ρ error	0.38 \pm 0.03	0.35 \pm 0.02	0.14 \pm 0.02	0.12 \pm 0.01	0.25 \pm 0.06
KR time	0.39 \pm 0.02	0.37 \pm 0.01	0.28 \pm 0.05	0.23 \pm 0.03	0.10 \pm 0.01
KR- ρ time	0.41 \pm 0.03	0.38 \pm 0.02	0.32 \pm 0.05	0.23 \pm 0.02	0.11 \pm 0.01

	Concrete Strength	Wine Quality	Telecom	Ailerons	Parkinson's
KR error	0.42 \pm 0.05	0.75 \pm 0.03	0.30 \pm 0.02	0.40 \pm 0.02	0.38 \pm 0.03
KR- ρ error	0.37 \pm 0.03	0.75 \pm 0.02	0.23 \pm 0.02	0.39 \pm 0.02	0.34 \pm 0.03
KR time	0.14 \pm 0.02	0.19 \pm 0.02	0.15 \pm 0.01	0.20 \pm 0.01	0.30 \pm 0.03
KR- ρ time	0.14 \pm 0.01	0.19 \pm 0.02	0.16 \pm 0.01	0.21 \pm 0.01	0.30 \pm 0.03

	Barrett joint 1	Barrett joint 5	SARCOS joint 1	SARCOS joint 5	Housing
k -NN error	0.41 \pm 0.02	0.40 \pm 0.02	0.08 \pm 0.01	0.08 \pm 0.01	0.28 \pm 0.09
k -NN- ρ error	0.29 \pm 0.01	0.30 \pm 0.02	0.07 \pm 0.01	0.07 \pm 0.01	0.22 \pm 0.06
k -NN time	0.21 \pm 0.04	0.16 \pm 0.03	0.13 \pm 0.01	0.13 \pm 0.01	0.08 \pm 0.01
k -NN- ρ time	0.13 \pm 0.04	0.16 \pm 0.03	0.14 \pm 0.01	0.13 \pm 0.01	0.08 \pm 0.01

	Concrete Strength	Wine Quality	Telecom	Ailerons	Parkinson's
k -NN error	0.40 \pm 0.04	0.73 \pm 0.04	0.13 \pm 0.02	0.37 \pm 0.01	0.22 \pm 0.01
k -NN- ρ error	0.38 \pm 0.03	0.72 \pm 0.03	0.17 \pm 0.02	0.34 \pm 0.01	0.20 \pm 0.01
k -NN time	0.10 \pm 0.01	0.15 \pm 0.01	0.16 \pm 0.02	0.12 \pm 0.01	0.14 \pm 0.01
k -NN- ρ time	0.11 \pm 0.01	0.15 \pm 0.01	0.15 \pm 0.01	0.11 \pm 0.01	0.15 \pm 0.01

Normalized mean square prediction errors and average prediction time per point (in milliseconds).



Normalized mean square prediction error over 2000 points for varying training sizes.