# Gradients Weights improve Regression and Classification-Supplementary Material

**Editor:** Somebody



(a) Wine Quality
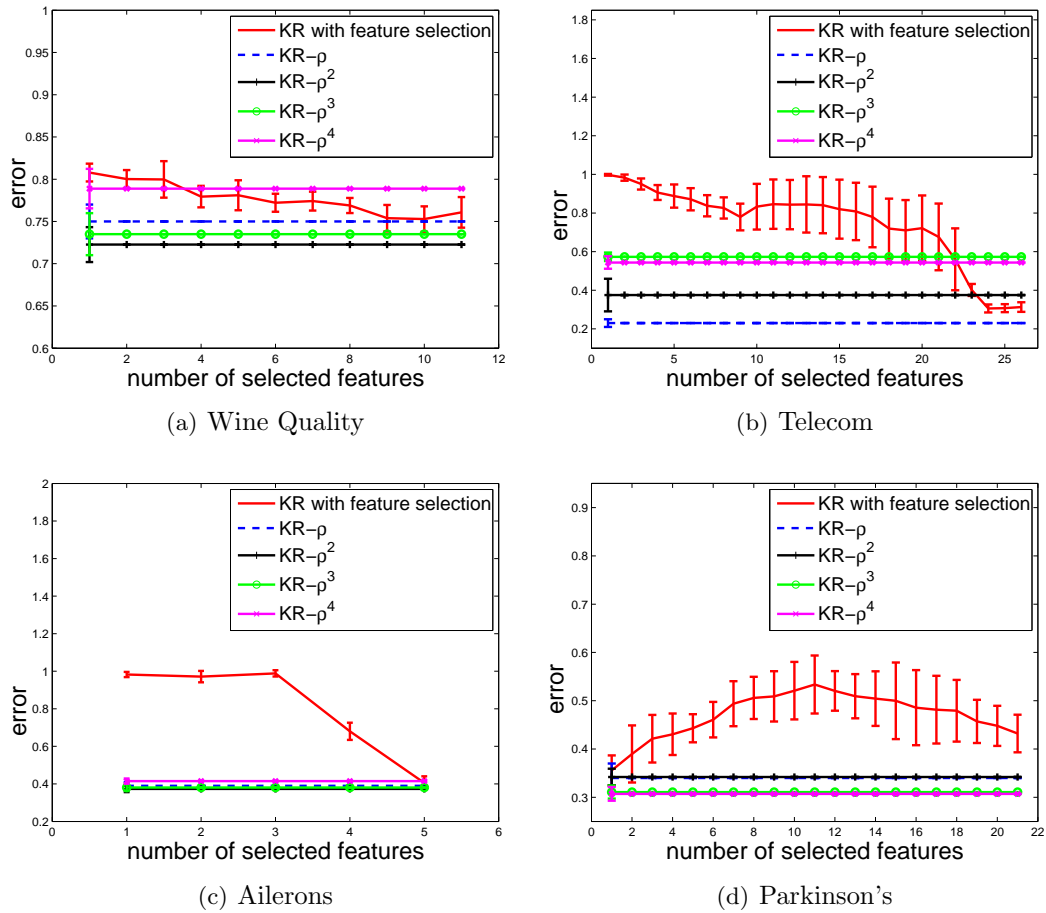
(b) Telecom

(c) Ailerons

(d) Parkinson's

Figure 1: Experiments on kernel regression with feature selection vs. gradient weights

(a) Barrett, joint 1

(b) Barrett, joint 5

(c) Sarcos, joint 1

(d) Sarcos, joint 5
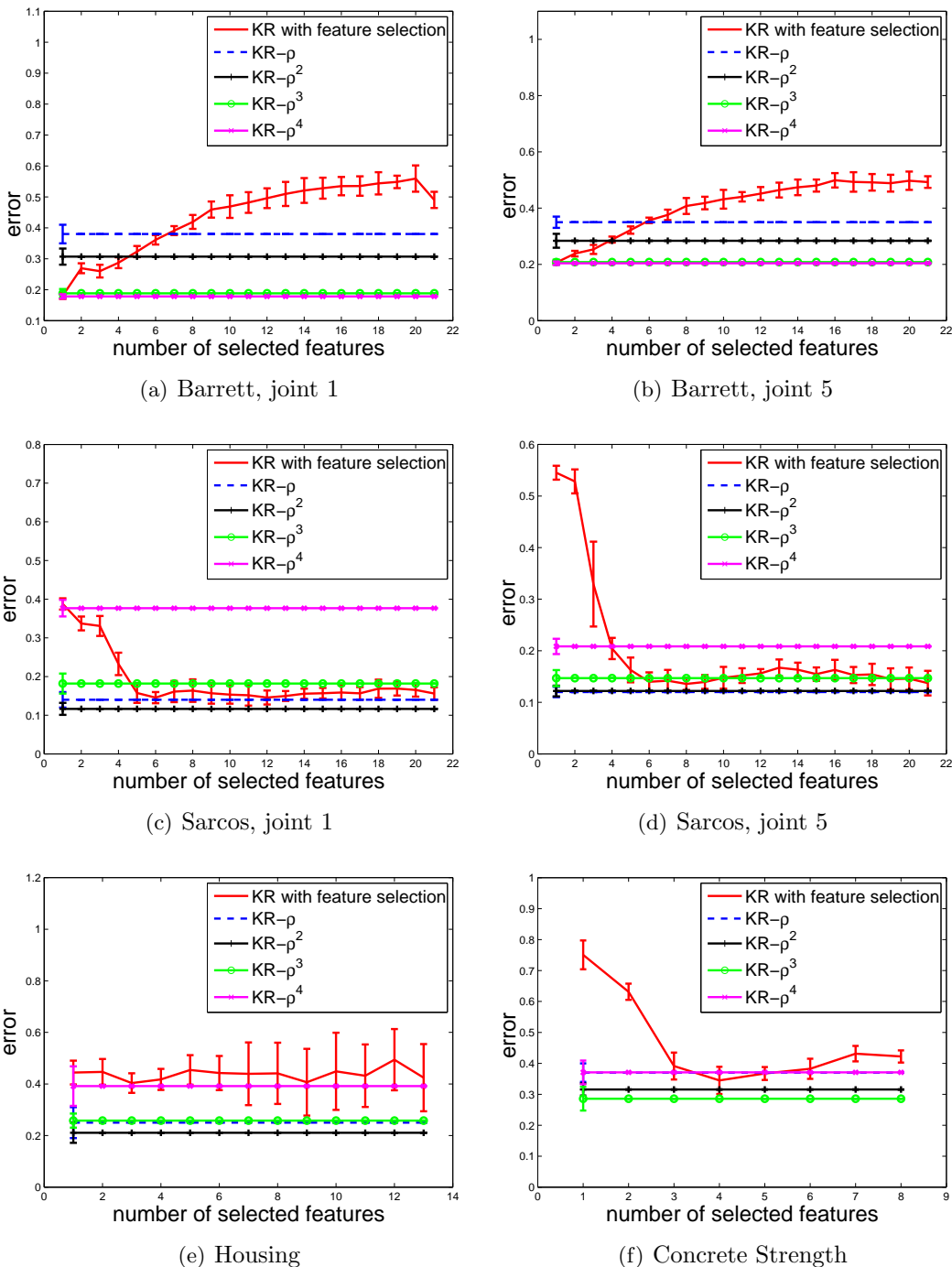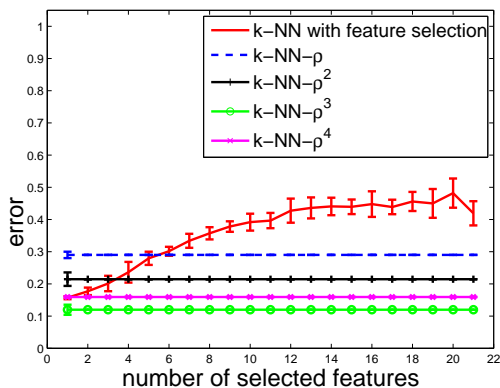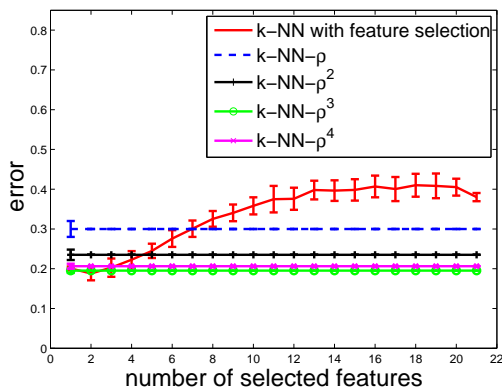
(e) Housing

(f) Concrete Strength

Figure 2: Experiments on kernel regression with feature selection vs. gradient weights

| | Barrett joint 1 | Barrett joint 5 | SARCOS joint 1 | SARCOS joint 5 | Housing |
|---|---|---|---|---|---|
| KR-unnormalized | $0.98 \pm 0.03$ | $0.90 \pm 0.03$ | $0.16 \pm 0.02$ | $0.32 \pm 0.03$ | $0.73 \pm 0.09$ |
| KR-r-normalized | $0.97 \pm 0.01$ | $0.89 \pm 0.02$ | $0.19 \pm 0.01$ | $0.33 \pm 0.04$ | $0.75 \pm 0.07$ |
| KR-normalized | $0.50 \pm 0.02$ | $0.50 \pm 0.03$ | $0.16 \pm 0.02$ | $0.14 \pm 0.02$ | $0.37 \pm 0.08$ |
| KR-normalized-$\rho$ | $0.38 \pm 0.03$ | $0.35 \pm 0.02$ | $0.14 \pm 0.02$ | $0.12 \pm 0.01$ | $0.25 \pm 0.06$ |
| KR-normalized-$\rho^2$ | $\mathbf{0.30} \pm 0.03$ | $\mathbf{0.28} \pm 0.03$ | $\mathbf{0.11} \pm 0.02$ | $\mathbf{0.12} \pm 0.01$ | $\mathbf{0.21} \pm 0.04$ |
| KR-normalized-$\rho^3$ | $0.18 \pm 0.02$ | $0.20 \pm 0.01$ | $0.18 \pm 0.03$ | $0.14 \pm 0.02$ | $0.25 \pm 0.03$ |
| KR-normalized-$\rho^4$ | $0.17 \pm 0.01$ | $0.20 \pm 0.01$ | $0.37 \pm 0.02$ | $0.20 \pm 0.01$ | $0.39 \pm 0.08$ |
| 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-$\rho$ 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-unnormalized | $0.45 \pm 0.03$ | $0.92 \pm 0.01$ | $0.23 \pm 0.02$ | $0.43 \pm 0.02$ | $0.75 \pm 0.09$ |
| KR-r-normalized | $0.43 \pm 0.04$ | $0.86 \pm 0.02$ | $0.23 \pm 0.02$ | $0.45 \pm 0.02$ | $0.75 \pm 0.06$ |
| KR-normalized | $0.42 \pm 0.05$ | $0.75 \pm 0.03$ | $0.30 \pm 0.02$ | $0.40 \pm 0.02$ | $0.38 \pm 0.03$ |
| KR-normalized-$\rho$ | $0.37 \pm 0.03$ | $0.75 \pm 0.02$ | $\mathbf{0.23} \pm 0.02$ | $0.39 \pm 0.02$ | $\mathbf{0.34} \pm 0.03$ |
| KR-normalized-$\rho^2$ | $\mathbf{0.31} \pm 0.02$ | $\mathbf{0.72} \pm 0.02$ | $0.37 \pm 0.08$ | $\mathbf{0.37} \pm 0.02$ | $\mathbf{0.34} \pm 0.02$ |
| KR-normalized-$\rho^3$ | $0.28 \pm 0.04$ | $0.73 \pm 0.03$ | $0.57 \pm 0.02$ | $0.38 \pm 0.02$ | $0.31 \pm 0.02$ |
| KR-normalized-$\rho^4$ | $0.37 \pm 0.04$ | $0.78 \pm 0.02$ | $0.54 \pm 0.03$ | $0.41 \pm 0.01$ | $0.30 \pm 0.01$ |
| 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-$\rho$ 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-unnormalized | $0.96 \pm 0.01$ | $0.80 \pm 0.03$ | $0.11 \pm 0.01$ | $0.19 \pm 0.01$ | $0.53 \pm 0.08$ |
| $k$-NN-r-normalized | $0.96 \pm 0.01$ | $0.78 \pm 0.04$ | $0.12 \pm 0.01$ | $0.22 \pm 0.02$ | $0.53 \pm 0.07$ |
| $k$-NN-normalized | $0.41 \pm 0.02$ | $0.40 \pm 0.02$ | $0.08 \pm 0.01$ | $0.08 \pm 0.01$ | $0.28 \pm 0.09$ |
| $k$-NN-normalized-$\rho$ | $0.29 \pm 0.01$ | $0.30 \pm 0.02$ | $0.07 \pm 0.01$ | $0.07 \pm 0.01$ | $0.22 \pm 0.06$ |
| $k$-NN-normalized-$\rho^2$ | $\mathbf{0.21} \pm 0.02$ | $\mathbf{0.23} \pm 0.01$ | $\mathbf{0.06} \pm 0.01$ | $\mathbf{0.06} \pm 0.01$ | $\mathbf{0.18} \pm 0.03$ |
| $k$-NN-normalized-$\rho^3$ | $0.11 \pm 0.02$ | $0.19 \pm 0.01$ | $0.08 \pm 0.01$ | $0.05 \pm 0.01$ | $0.23 \pm 0.03$ |
| $k$-NN-normalized-$\rho^4$ | $0.15 \pm 0.01$ | $0.20 \pm 0.01$ | $0.37 \pm 0.01$ | $0.16 \pm 0.01$ | $0.31 \pm 0.07$ |
| $k$-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-$\rho$ time | $0.13 \pm 0.04$ | $0.16 \pm 0.03$ | $0.14 \pm 0.01$ | $0.13 \pm 0.01$ | $0.08 \pm 0.01$ |

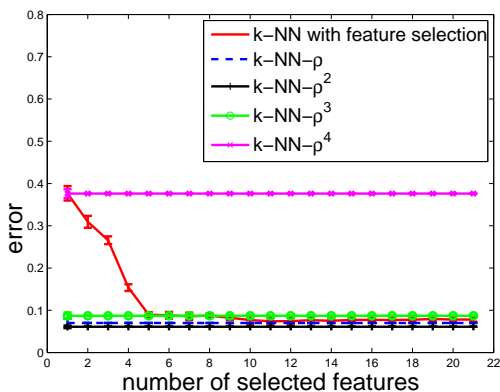| | Concrete Strength | Wine Quality | Telecom | Ailerons | Parkinson's |
|---|---|---|---|---|---|
| $k$-NN-unnormalized | $0.40 \pm 0.07$ | $0.88 \pm 0.01$ | $0.15 \pm 0.02$ | $0.42 \pm 0.02$ | $0.63 \pm 0.04$ |
| $k$-NN-r-normalized | $0.37 \pm 0.08$ | $0.85 \pm 0.02$ | $0.17 \pm 0.02$ | $0.44 \pm 0.02$ | $0.58 \pm 0.02$ |
| $k$-NN-normalized | $0.40 \pm 0.04$ | $0.73 \pm 0.04$ | $\mathbf{0.13} \pm 0.02$ | $0.37 \pm 0.01$ | $0.22 \pm 0.01$ |
| $k$-NN-normalized-$\rho$ | $0.38 \pm 0.03$ | $0.72 \pm 0.03$ | $0.17 \pm 0.02$ | $\mathbf{0.34} \pm 0.01$ | $\mathbf{0.20} \pm 0.01$ |
| $k$-NN-normalized-$\rho^2$ | $\mathbf{0.31} \pm 0.06$ | $\mathbf{0.70} \pm 0.01$ | $0.34 \pm 0.05$ | $0.34 \pm 0.01$ | $\mathbf{0.20} \pm 0.01$ |
| $k$-NN-normalized-$\rho^3$ | $0.26 \pm 0.02$ | $0.71 \pm 0.01$ | $0.55 \pm 0.03$ | $0.36 \pm 0.01$ | $0.22 \pm 0.01$ |
| $k$-NN-normalized-$\rho^4$ | $0.38 \pm 0.05$ | $0.78 \pm 0.01$ | $0.52 \pm 0.02$ | $0.45 \pm 0.01$ | $0.25 \pm 0.01$ |
| $k$-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-$\rho$ time | $0.11 \pm 0.01$ | $0.15 \pm 0.01$ | $0.15 \pm 0.01$ | $0.11 \pm 0.01$ | $0.15 \pm 0.01$ |

Table 1: Normalized mean square prediction errors and average prediction time per point (in milliseconds). The top five tables are for KR vs KR-$\rho$, KR-$\rho^2$,KR-$\rho^3$ and KR-$\rho^4$, the bottom five for $k$-NN vs $k$-NN-$\rho$, $k$-NN-$\rho^2$, $k$-NN-$\rho^3$ and $k$-NN-$\rho^4$. In the methods labeled as normalized, each feature was divided by its empirical standard deviation in the training data. In the methods labeled as unnormalized, we used the raw data sets. In the methods labeled as r-normalized, each data point (row) was normalized so that its $l_2$ norm is equal to 1. This latter type of normalization does not seem to have a significant benefit. In fact, the prediction errors using unnormalized and r-normalized data are not significantly different.
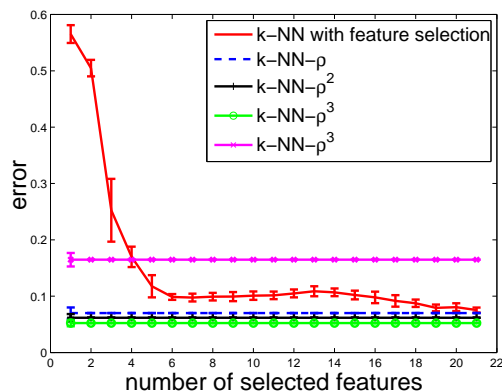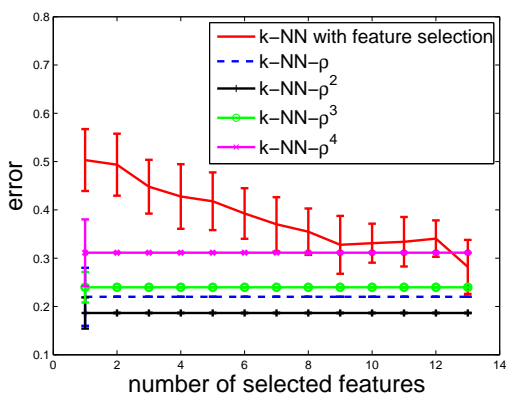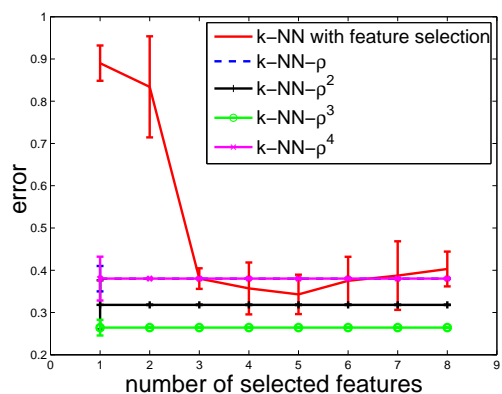
(a) Barrett, joint 1

(b) Barrett, joint 5

(c) Sarcos, joint 1

(d) Sarcos, joint 5

(e) Housing

(f) Concrete Strength

Figure 3: Experiments on $k$-NN regression with feature selection

|                  | Covertype          | IJCNN                 | MAGIC Gamma         | Shuttle               | Page Blocks           |
|------------------|--------------------|-----------------------|---------------------|-----------------------|-----------------------|
| SVM error        | $0.25\pm0.01$      | $0.0576\pm0.0067$     | $\mathbf{0.1491}\pm0.0100$ | $0.0058\pm0.0027$ | $0.0345\pm0.0048$ |
| SVM-$\rho$ error | $\mathbf{0.24}\pm0.01$ | $0.0531\pm0.0059$ | $0.1507\pm0.0107$   | $\mathbf{0.0034}\pm0.0025$ | $0.0342\pm0.0051$ |
| SVM-$\rho^2$ error | $\mathbf{0.24}\pm0.01$ | $\mathbf{0.0521}\pm0.0065$ | $0.1540\pm0.0113$ | $\mathbf{0.0034}\pm0.0018$ | $\mathbf{0.0336}\pm0.0042$ |

Table 2: Classification error rates of a support vector machine suggest that pre-multiplying features by their gradient weight, which corresponds to the $\rho^2$ metric, also improves performance of that classifier in several cases.
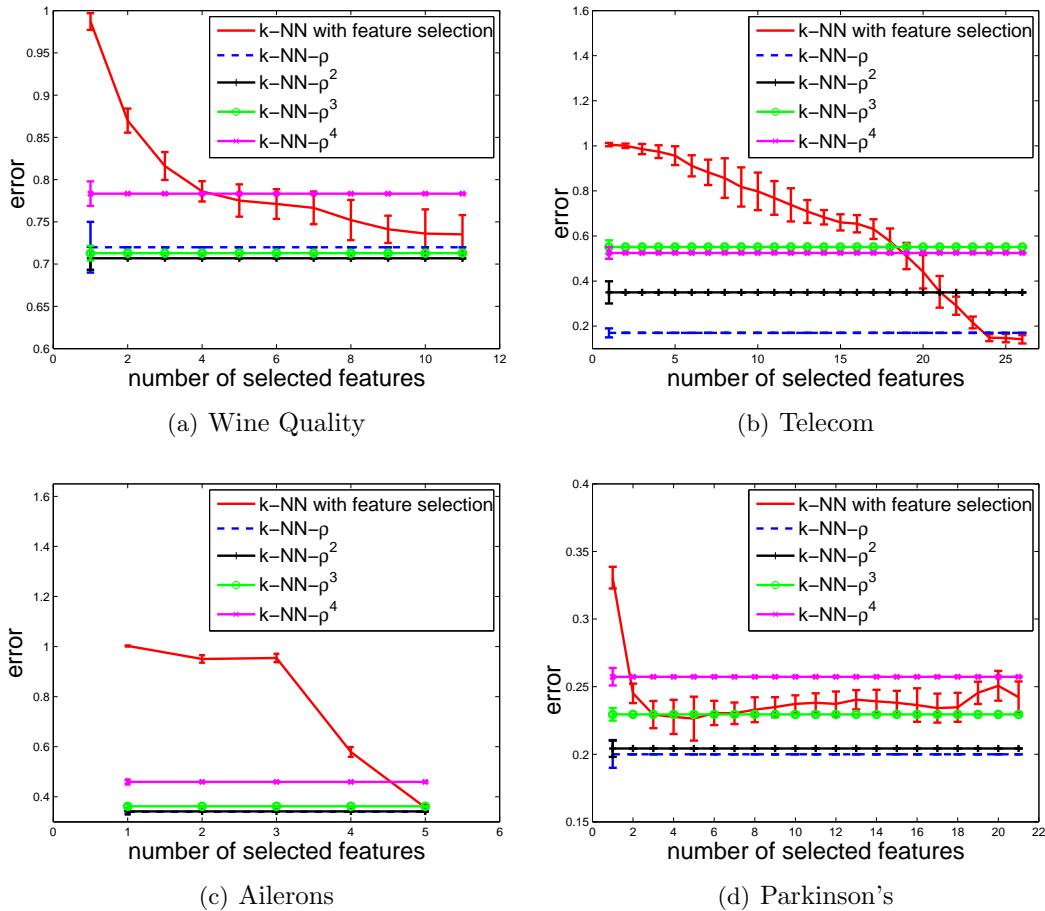


(a) Wine Quality

(b) Telecom

(c) Ailerons

(d) Parkinson's

Figure 4: Experiments on $k$-NN regression with feature selection

(a) Covertype with SVM

(b) IJCNN with SVM

(c) Thyroid with SVM
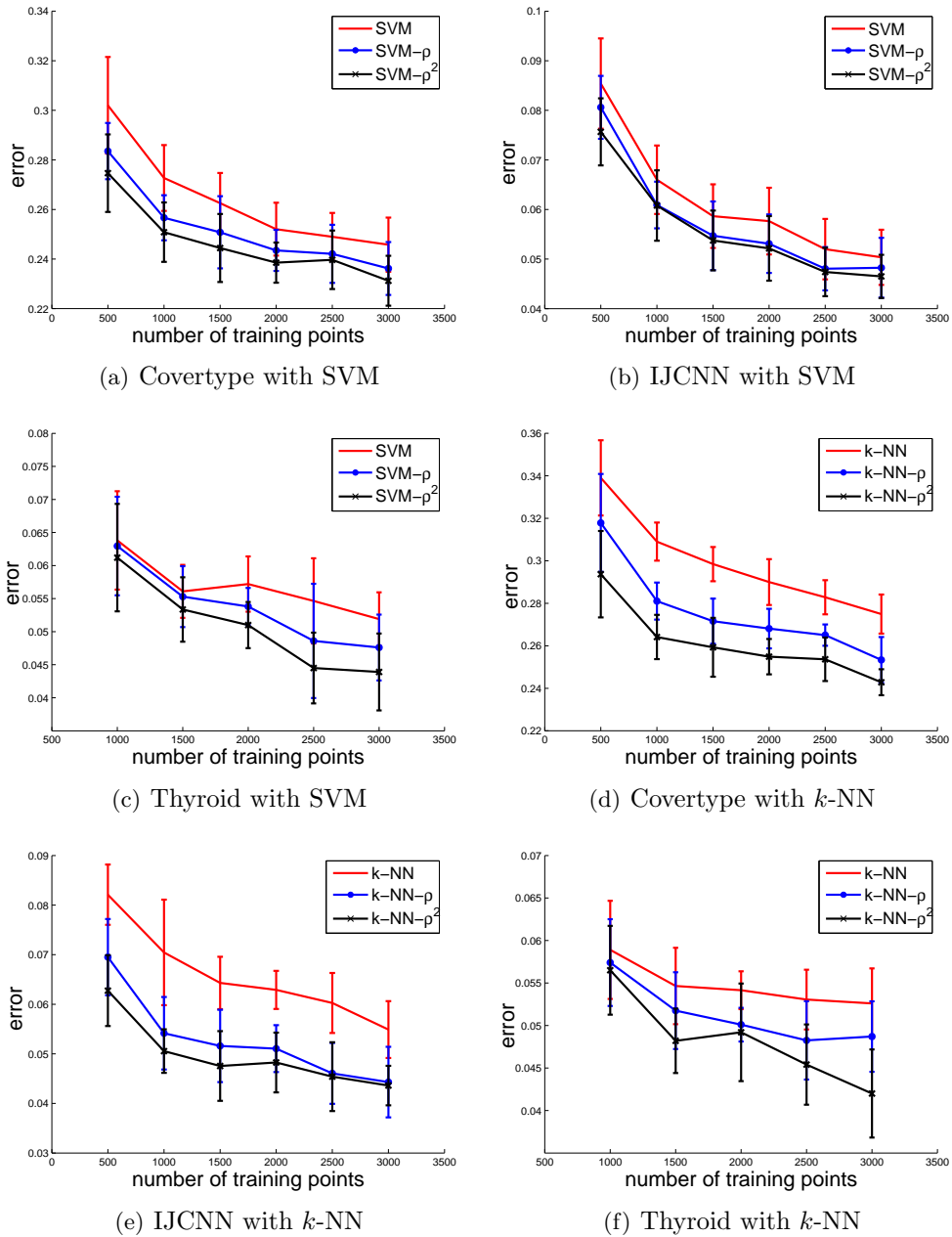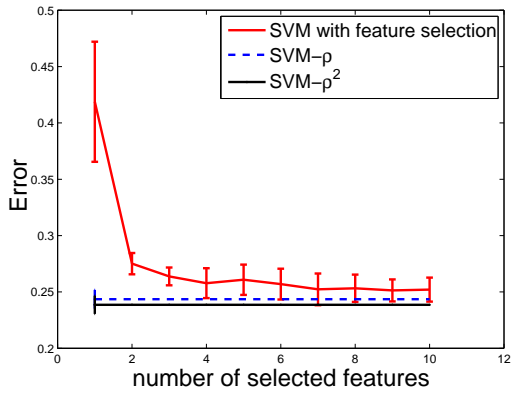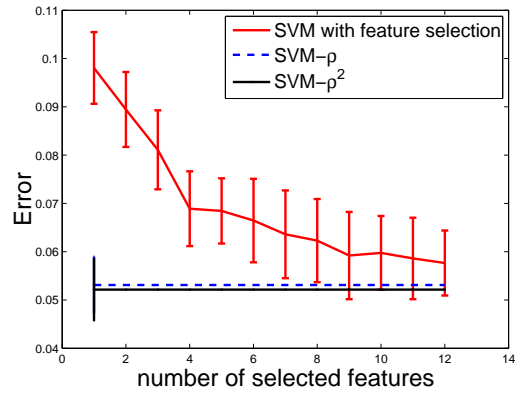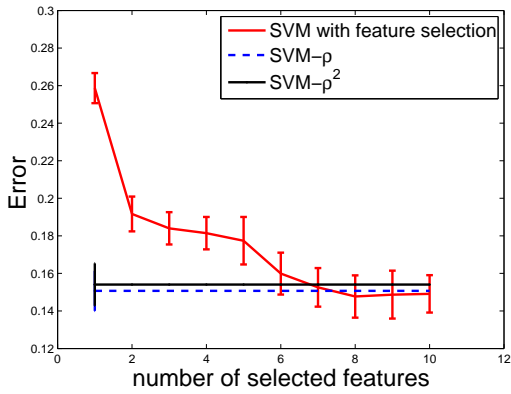
(d) Covertype with $k$-NN

(e) IJCNN with $k$-NN

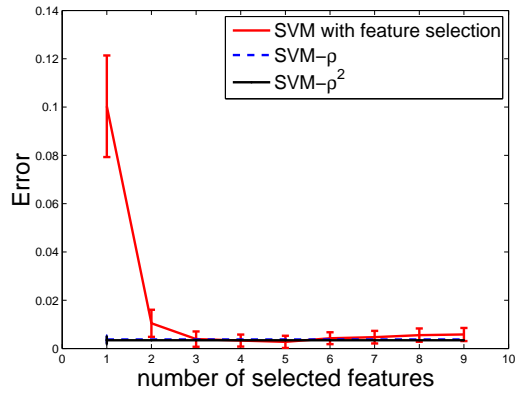(f) Thyroid with $k$-NN

Figure 5: Classification results with gradient weights

Figure 6: Experiments on SVM with feature selection

(a) covertype

(b) IJCNN
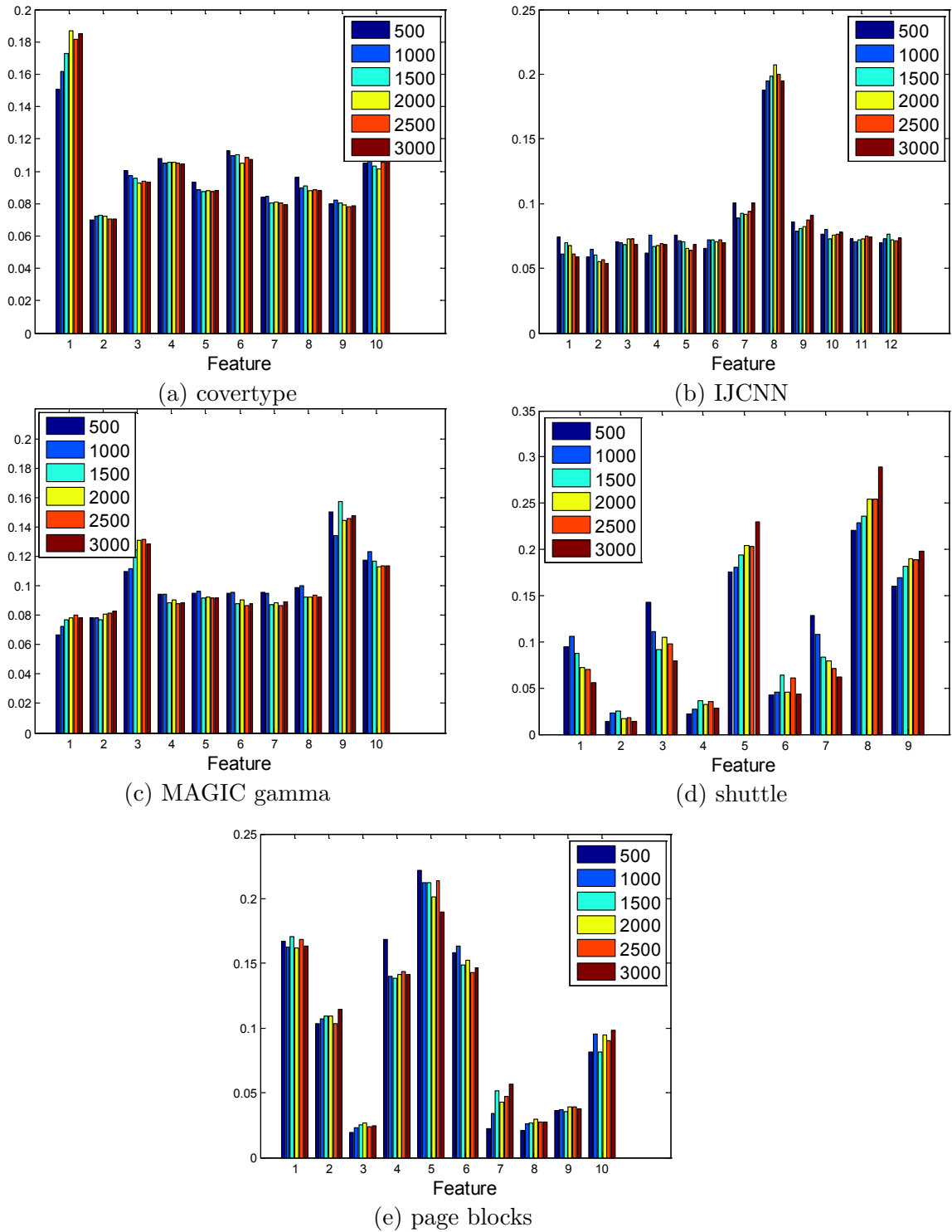
(c) MAGIC gamma

(d) shuttle

(e) page blocks

Figure 7: Normalized gradient weights obtained by weighted $k$-NN with a Gaussian kernel using different training sizes.

| Dataset | $k$-NN | $k$-NN-$\rho$ |
|---------|--------|---------------|
| Ailerons | $0.3364 \pm 0.0087$ | $0.3161 \pm 0.0058$ |
| Concrete | $0.2884 \pm 0.0311$ | $\mathbf{0.2040} \pm 0.0234$ |
| Housing | $0.2897 \pm 0.0632$ | $\mathbf{0.2389} \pm 0.0604$ |
| Wine | $0.6633 \pm 0.0119$ | $0.6615 \pm 0.0134$ |
| Barrett1 | $0.1051 \pm 0.0150$ | $\mathbf{0.0843} \pm 0.0229$ |
| Barrett5 | $0.1095 \pm 0.0096$ | $\mathbf{0.0984} \pm 0.0244$ |
| Sarcos1 | $0.1222 \pm 0.0074$ | $\mathbf{0.0769} \pm 0.0037$ |
| Sarcos5 | $0.0870 \pm 0.0051$ | $0.0779 \pm 0.0026$ |
| Parkinson | $0.3638 \pm 0.0443$ | $\mathbf{0.3181} \pm 0.0477$ |
| TeleComm | $0.0864 \pm 0.0094$ | $0.0688 \pm 0.0074$ |

Table 3: Regression results, with ten random runs per data set. For each method, the values of $k$ as well as $t$ (the bandwidth used to estimate finite differences for calculating the gradients) were set by two fold cross-validation on the training set.

| Dataset | $k$-NN | $k$-NN$-\rho$ |
|---------|--------|---------------|
| Cover Type | $0.2279 \pm 0.0091$ | $0.2135 \pm 0.0064$ |
| Gamma | $0.1775 \pm 0.0070$ | $0.1680 \pm 0.0075$ |
| Page Blocks | $0.0349 \pm 0.0042$ | $0.0361 \pm 0.0048$ |
| Shuttle | $0.0037 \pm 0.0025$ | $0.0024 \pm 0.0016$ |
| Musk | $0.2279 \pm 0.0091$ | $0.2135 \pm 0.0064$ |
| IJCNN | $0.0540 \pm 0.0061$ | $0.0459 \pm 0.0058$ |
| RNA | $0.1042 \pm 0.0063$ | $0.0673 \pm 0.0062$ |

Table 4: Classification results, with ten random runs per data set. For each method, the values of $k$ as well as $t$ (the bandwidth used to estimate finite differences for calculating the gradients) were set by two fold cross-validation on the training set.