

---

# Appendix

## Sparse Spectrum Warped Input Measures for Nonstationary Kernel Learning

---

Anonymous Author(s)

Affiliation

Address

email

### 1 A Mathematical derivations

2 **Lemma 1.** Let  $\mathbf{v} \sim \mathcal{N}(\hat{\mathbf{v}}, \Sigma_{\mathbf{v}})$  denote a Gaussian random vector and  $\mathbf{x} \in \mathbb{R}^D$  an arbitrary point.  
3 Then  $\mathbf{z} := \mathbf{v} \odot \mathbf{x}$  is Gaussian,  $\mathbf{z} \sim \mathcal{N}(\hat{\mathbf{z}}, \Sigma_{\mathbf{z}})$ , with mean and covariance matrix given by:

$$\hat{\mathbf{z}} = \hat{\mathbf{v}} \odot \mathbf{x} \tag{1}$$

$$\Sigma_{\mathbf{z}} = \mathbf{x}\mathbf{1}^T \odot \Sigma_{\mathbf{v}} \odot \mathbf{1}\mathbf{x}^T \tag{2}$$

4 *Proof.* The element-wise vector product, which is the Hadamard product for single-column matrices,  
5 is symmetric and linear, i.e.  $\mathbf{a} \odot \mathbf{b} = \mathbf{b} \odot \mathbf{a}$  and  $\mathbf{a} \odot (\mathbf{b} + \mathbf{c}) = \mathbf{a} \odot \mathbf{b} + \mathbf{a} \odot \mathbf{c}$ , for any  $\mathbf{a}, \mathbf{b}, \mathbf{c} \in \mathbb{R}^D$ .  
6 By the linearity of the expectation, we then have  $\mathbb{E}[\mathbf{z}] = \mathbb{E}[\mathbf{v} \odot \mathbf{x}] = \mathbb{E}[\mathbf{v}] \odot \mathbf{x}$ . From the properties  
7 of the Hadamard product, one can also show that  $\mathbf{a}(\mathbf{b} \odot \mathbf{c})^T = \mathbf{a}\mathbf{b}^T \odot \mathbf{1}\mathbf{c}^T$ , where  $\mathbf{1}$  is a vector of  
8 1's. Therefore, the covariance matrix  $\Sigma_{\mathbf{z}} := \mathbb{V}[\mathbf{z}] = \mathbb{E}[(\mathbf{z} - \hat{\mathbf{z}})(\mathbf{z} - \hat{\mathbf{z}})^T]$ , is given by:

$$\begin{aligned} \mathbb{V}[\mathbf{z}] &= \mathbb{E}[(\mathbf{v} - \hat{\mathbf{v}}) \odot \mathbf{x}][(\mathbf{v} - \hat{\mathbf{v}}) \odot \mathbf{x}]^T \\ &= \mathbf{x}\mathbf{1}^T \odot \mathbb{E}[(\mathbf{v} - \hat{\mathbf{v}})(\mathbf{v} - \hat{\mathbf{v}})^T] \odot \mathbf{1}\mathbf{x}^T \\ &= \mathbf{x}\mathbf{1}^T \odot \mathbb{V}[\mathbf{v}] \odot \mathbf{1}\mathbf{x}^T, \end{aligned} \tag{3}$$

9 which concludes the proof. □

### 10 B Relationship to output warped GPs

11 A body of work has previously been developed under the title of *warped* Gaussian processes [2].  
12 As noted, this contrasts from our modelling problem because while we proceed to expand the GP's  
13 capabilities to warp the *inputs*, the WGP and extensions warps explicitly the *output* distribution of  
14 the Gaussian process. We now juxtapose the efficacy of our input warping formulation with WGP by  
15 applying SSWIM to the three challenging datasets *abalone*, *creep*, and *aileron*s experimented upon  
16 in [2, 1].

17 Our results in Table Table 1 suggest that with SSWIM we are able to improve upon both WGP  
18 and BWGP in MSE: marginal improvements for *abalone* and a significant improvement for *creep*  
19 with comparable performance in *aileron*s. Contrasting this, the output warping methods outperform  
20 unanimously on the MNLP metric. This is expected because output warping may allow one to capture  
21 non-gaussian conditional distributions which an input warping formulation cannot with a standard  
22 Gaussian process. The discussion we wish to raise here is that both aspects of manipulating the inputs  
23 and outputs of a GP can result in major improvements respectively across different metrics.

Table 1: MSE and MNLP metrics for comparison with Warped and Bayesian Warped GPs [1]. MSE results for *aileron*s are  $\times 10^{-8}$ .

|      |                    | Method | abalone                           | creep                                | aileron                            |
|------|--------------------|--------|-----------------------------------|--------------------------------------|------------------------------------|
| MSE  | GP                 |        | 4.55 $\pm$ 0.14                   | 584.9 $\pm$ 71.2                     | 2.95 $\pm$ 0.16                    |
|      | BWGP               |        | 4.55 $\pm$ 0.11                   | 491.8 $\pm$ 36.2                     | 2.91 $\pm$ 0.14                    |
|      | MLWGP3             |        | 4.54 $\pm$ 0.10                   | 502.3 $\pm$ 43.3                     | <b>2.80 <math>\pm</math> 0.11</b>  |
|      | MLWGP20            |        | 4.59 $\pm$ 0.32                   | 506.3 $\pm$ 46.1                     | 3.42 $\pm$ 2.87                    |
|      | SSWIM <sub>1</sub> |        | 4.64 $\pm$ 0.13                   | 483.69 $\pm$ 64.12                   | 2.96 $\pm$ 0.08                    |
|      | SSWIM <sub>2</sub> |        | <b>4.50 <math>\pm</math> 0.11</b> | <b>279.86 <math>\pm</math> 31.88</b> | 2.83 $\pm$ 0.06                    |
| MNLP | GP                 |        | 2.17 $\pm$ 0.01                   | 4.46 $\pm$ 0.03                      | -7.30 $\pm$ 0.01                   |
|      | BWGP               |        | 1.99 $\pm$ 0.01                   | 4.31 $\pm$ 0.04                      | -7.38 $\pm$ 0.02                   |
|      | MLWGP3             |        | <b>1.97 <math>\pm</math> 0.02</b> | <b>4.21 <math>\pm</math> 0.03</b>    | -7.44 $\pm$ 0.01                   |
|      | MLWGP20            |        | 1.99 $\pm$ 0.05                   | <b>4.21 <math>\pm</math> 0.08</b>    | <b>-7.45 <math>\pm</math> 0.08</b> |
|      | SSWIM <sub>1</sub> |        | 2.18 $\pm$ 0.01                   | 4.45 $\pm$ 0.03                      | -7.24 $\pm$ 0.01                   |
|      | SSWIM <sub>2</sub> |        | 2.17 $\pm$ 0.02                   | 4.27 $\pm$ 0.03                      | -7.00 $\pm$ 0.02                   |

## 24 C Additional Experiments

### 25 C.1 Increasing number of pseudo-training points

26 For the "increasing number of pseudo-training points" experiment we used 1 layer of warping with  
 27 256 features for both the warping and top-level predictive functions.

### 28 C.2 Increasing warping depth

29 We used 256 features and 1280 pseudo-training points for all of the experiments.

### 30 C.3 Complete real-dataset experiments table

31 Table 2 contains additional real-world experiments to extend the majore experimental results from  
 32 the main paper.

Table 2: RMSE and MNLP metrics for various real world datasets.

|                           |                    | (D, N) | (18, 8751)                        | (5, 1503)                         | (8, 1030)                         | (16, 5875)                         | (15, 17379)                         | (379, 53500)                        | (81, 21263)                       | (9, 45730)                        | (77, 583250)                       | (90, 515345)                       |
|---------------------------|--------------------|--------|-----------------------------------|-----------------------------------|-----------------------------------|------------------------------------|-------------------------------------|-------------------------------------|-----------------------------------|-----------------------------------|------------------------------------|------------------------------------|
|                           |                    | Method | elevators                         | airfoil                           | concrete                          | parkinsons                         | bikeshare                           | ct slice                            | supercond                         | protein                           | buzz                               | song                               |
| RMSE ( $\times 10^{-1}$ ) | SSWIM <sub>1</sub> |        | 2.83 $\pm$ 0.08                   | 2.38 $\pm$ 0.26                   | 3.05 $\pm$ 0.26                   | 7.63 $\pm$ 0.20                    | 0.13 $\pm$ 0.04                     | 0.46 $\pm$ 0.02                     | 3.44 $\pm$ 0.14                   | 5.91 $\pm$ 0.07                   | 2.98 $\pm$ 0.04                    | 8.12 $\pm$ 0.05                    |
|                           | SSWIM <sub>2</sub> |        | <b>2.74 <math>\pm</math> 0.08</b> | <b>2.35 <math>\pm</math> 0.22</b> | <b>3.01 <math>\pm</math> 0.31</b> | <b>7.55 <math>\pm</math> 0.15</b>  | 0.11 $\pm$ 0.03                     | <b>0.23 <math>\pm</math> 0.01</b>   | <b>3.02 <math>\pm</math> 0.04</b> | <b>5.80 <math>\pm</math> 0.08</b> | <b>2.40 <math>\pm</math> 0.01</b>  | <b>7.97 <math>\pm</math> 0.03</b>  |
|                           | DSDGP              |        | <b>2.74 <math>\pm</math> 0.06</b> | 4.30 $\pm$ 0.19                   | 5.88 $\pm$ 1.24                   | 7.94 $\pm$ 0.20                    | 0.33 $\pm$ 0.55                     | 4.81 $\pm$ 1.18                     | 5.10 $\pm$ 0.84                   | 5.96 $\pm$ 0.06                   | 3.65 $\pm$ 0.75                    | 8.46 $\pm$ 0.03                    |
|                           | DKL                |        | 3.06 $\pm$ 0.29                   | 3.19 $\pm$ 0.37                   | 3.18 $\pm$ 0.38                   | 8.84 $\pm$ 0.74                    | 0.24 $\pm$ 0.03                     | 0.52 $\pm$ 0.08                     | 3.46 $\pm$ 0.18                   | 7.15 $\pm$ 1.10                   | 4.11 $\pm$ 3.33                    | 16.66 $\pm$ 8.14                   |
|                           | RFFNS              |        | 2.83 $\pm$ 0.07                   | 3.31 $\pm$ 0.37                   | 3.46 $\pm$ 0.24                   | 8.15 $\pm$ 0.15                    | 0.05 $\pm$ 0.01                     | 4.39 $\pm$ 0.27                     | 3.85 $\pm$ 0.05                   | 6.87 $\pm$ 0.06                   | 5.70 $\pm$ 0.84                    | 8.35 $\pm$ 0.03                    |
|                           | SVGP               |        | 2.88 $\pm$ 0.10                   | 2.70 $\pm$ 0.15                   | 3.32 $\pm$ 0.26                   | 8.14 $\pm$ 0.12                    | 0.06 $\pm$ 0.03                     | 1.16 $\pm$ 0.02                     | 4.06 $\pm$ 0.05                   | 7.32 $\pm$ 0.08                   | 9.98 $\pm$ 0.02                    | 12.19 $\pm$ 0.18                   |
|                           | SGPR               |        | 4.96 $\pm$ 2.02                   | 4.24 $\pm$ 0.40                   | 5.55 $\pm$ 0.58                   | 7.86 $\pm$ 0.22                    | 0.67 $\pm$ 0.18                     | 1.79 $\pm$ 0.04                     | 4.27 $\pm$ 0.06                   | 6.45 $\pm$ 0.07                   | 2.89 $\pm$ 0.02                    | 8.40 $\pm$ 0.04                    |
|                           | RFFS               |        | 2.87 $\pm$ 0.10                   | 3.28 $\pm$ 0.24                   | 3.33 $\pm$ 0.30                   | 8.24 $\pm$ 0.17                    | <b>0.03 <math>\pm</math> 0.00</b>   | 2.34 $\pm$ 0.05                     | 3.89 $\pm$ 0.06                   | 6.91 $\pm$ 0.07                   | 3.78 $\pm$ 0.14                    | 8.36 $\pm$ 0.04                    |
| MNLP ( $\times 10^{-1}$ ) | SSWIM <sub>1</sub> |        | 1.81 $\pm$ 0.55                   | 1.25 $\pm$ 2.50                   | 10.22 $\pm$ 4.15                  | 11.95 $\pm$ 0.47                   | -11.89 $\pm$ 0.15                   | -11.24 $\pm$ 0.05                   | 3.55 $\pm$ 0.32                   | 8.95 $\pm$ 0.12                   | 2.03 $\pm$ 0.13                    | 12.08 $\pm$ 0.05                   |
|                           | SSWIM <sub>2</sub> |        | 1.65 $\pm$ 0.63                   | <b>1.05 <math>\pm</math> 1.74</b> | 5.19 $\pm$ 2.59                   | 12.50 $\pm$ 0.44                   | -11.78 $\pm$ 0.07                   | <b>-11.79 <math>\pm</math> 0.02</b> | <b>2.82 <math>\pm</math> 0.29</b> | <b>8.82 <math>\pm</math> 0.15</b> | <b>-0.09 <math>\pm</math> 0.04</b> | <b>11.93 <math>\pm</math> 0.04</b> |
|                           | DSDGP              |        | <b>1.16 <math>\pm</math> 0.21</b> | 9.19 $\pm$ 0.20                   | 11.02 $\pm$ 1.06                  | <b>11.91 <math>\pm</math> 0.24</b> | -23.28 $\pm$ 8.29                   | 6.62 $\pm$ 2.61                     | 7.36 $\pm$ 1.62                   | 9.04 $\pm$ 0.10                   | 3.80 $\pm$ 2.02                    | 12.52 $\pm$ 0.04                   |
|                           | DKL                |        | 7.57 $\pm$ 0.15                   | 7.70 $\pm$ 0.17                   | 7.69 $\pm$ 0.20                   | 13.17 $\pm$ 1.12                   | 6.82 $\pm$ 0.01                     | 6.83 $\pm$ 0.01                     | 7.76 $\pm$ 0.10                   | 11.02 $\pm$ 1.53                  | 9.01 $\pm$ 4.65                    | 42.64 $\pm$ 44.77                  |
|                           | RFFNS              |        | 1.44 $\pm$ 0.24                   | 2.74 $\pm$ 1.20                   | 3.31 $\pm$ 0.45                   | 12.18 $\pm$ 0.18                   | -11.97 $\pm$ 0.00                   | 5.95 $\pm$ 0.66                     | 4.66 $\pm$ 0.12                   | 10.39 $\pm$ 0.08                  | 8.78 $\pm$ 1.87                    | 12.39 $\pm$ 0.04                   |
|                           | SVGP               |        | 1.78 $\pm$ 0.30                   | 1.19 $\pm$ 0.33                   | <b>2.83 <math>\pm</math> 0.56</b> | 12.21 $\pm$ 0.14                   | <b>-27.70 <math>\pm</math> 1.24</b> | -5.98 $\pm$ 0.13                    | 5.32 $\pm$ 0.12                   | 11.10 $\pm$ 0.09                  | 63.31 $\pm$ 3.44                   | 18.02 $\pm$ 0.09                   |
|                           | SGPR               |        | 11.59 $\pm$ 22.38                 | 6.45 $\pm$ 0.47                   | 8.48 $\pm$ 2.10                   | 12.39 $\pm$ 0.30                   | -13.67 $\pm$ 0.98                   | -3.14 $\pm$ 0.26                    | 5.58 $\pm$ 0.10                   | 10.05 $\pm$ 0.13                  | 1.11 $\pm$ 0.12                    | 11.97 $\pm$ 0.07                   |
|                           | RFFS               |        | 1.59 $\pm$ 0.26                   | 2.72 $\pm$ 0.69                   | 3.05 $\pm$ 0.96                   | 12.29 $\pm$ 0.20                   | -11.98 $\pm$ 0.00                   | -0.33 $\pm$ 0.22                    | 4.79 $\pm$ 0.13                   | 10.45 $\pm$ 0.09                  | 4.41 $\pm$ 0.37                    | 12.40 $\pm$ 0.05                   |

### 33 C.4 Extended discussion

34 It is imperative to note here our aim is not to demand any algorithmic dominance when comparing  
 35 methods. Firstly, this is a fruitless pursuit due to the diversity and dependence of all advanced  
 36 GP methods on hyperparameter optimization, and secondly it is not a constructive approach to the  
 37 communal development of new methodologies to claim some benchmark task superiority. Rather,  
 38 encouraged by the results of this paper, we invite the discussion to move beyond stationary kernels  
 39 and inquire upon and interpret the effectiveness of new Gaussian process methodologies through the

40 *more general perspective* of input space warping for unearthing hidden nonstationary patterns within  
 41 data.

## 42 C.5 Overfitting analysis

43 We ran an overfitting analysis of  $SSWIM_1$  to observe the effect of over-optimising with respect  
 44 to the marginal likelihood. We ran with 256 features, 1280 pseudo-training points, for 150 steps,  
 45 with 10 repeats, and evaluated the test RMSE and test MNLP on the test set for every single epoch of  
 46 optimisation. This test bench is provided in the supplementary code. The results are averaged with  
 47 mean and standard deviations of the training curves displayed in Figure 1. We can see that we are  
 48 quite resistant to overfitting except for RMSE in the *elevators* dataset and the MNLP in the *concrete*  
 49 and *parkinsons* datasets. The causes of this could be explained by the underlying flexibility of the  
 50 proposed method which allows the model to become overconfident in what it has learned with respect  
 51 to the data it has observed. In fact, we observed similar overfitting behaviour for similar optimisation  
 52 periods with DKL and DSDGP.

53 This analysis leads to some interesting observations and recommendations for future algorithm  
 54 development in more expressive GP methodologies: 1. the marginal likelihood is no panacea  
 55 although it is easy to think it is, and 2. other loss functions and training schemes, such as leave-  
 56 one-out cross validation. Actually, these results corroborate long known discussions from [3] about  
 57 the risk of overfitting from trusting the marginal likelihood with standard optimisation procedures,  
 58 however their importance seems to have been largely ignored in evaluation of recent methodology  
 59 innovations in the GP literature. We believe that a more open discussion should be on the table for  
 60 analysing the interplay between model expressiveness and the effect this has on overfitting; this is  
 61 especially pertinent to the GP literature which has placed a large emphasis on the importance of the  
 62 marginal likelihood has a valid hyperparameter optimisation loss.

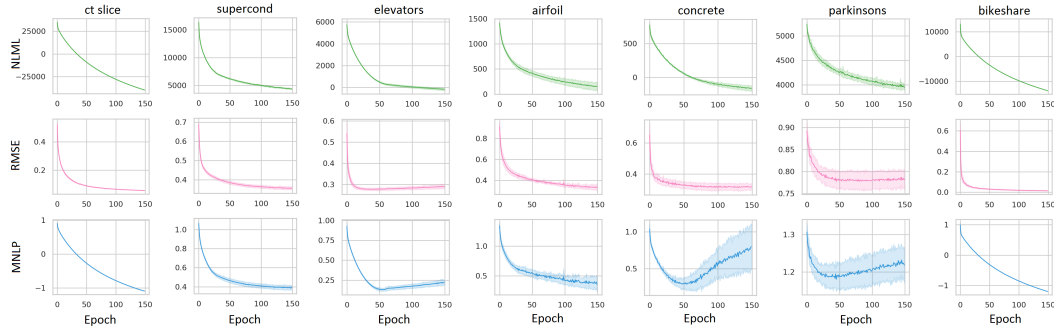


Figure 1: Empirical analysis of overfitting behaviour in SSWIM

## 63 C.6 Pseudo-training points in 2D

64 Figure 2 contains an interpretation of input warping and the pseudo-training points in higher di-  
 65 mensions. We use the exponential 2D function from [4] as a case where it is intuitive to observe  
 66 how SSWIM reacts to topological differences in the underlying function. The function consists of  
 67 a mostly flat surface with abrupt spiking occurring at a corner of the domain as seen in Figure 2  
 68 (b). If we were to assume a homogeneous domain, and model our data with a standard SSGP with  
 69 stationary RBF kernel, the kernel’s hyperparameters would be optimised to provide a homogeneous  
 70 representation resulting in conflict between the spiked area in the corner and flat areas elsewhere. By  
 71 directly manipulating the input domain with our warping, we are able to transform the input domain  
 72 with a continuous warping that consequently allows accurate representation as seen in Figure 2 (c).

73 The intuition and utility of our proposed *pseudo-training* as virtual training points is visualised from  
 74 a birds-eye view in Figure 2 (d). Before optimisation, we initialise these points uniformly across  
 75 the domain of the training space. It is clear that the learned positions of the pseudo-training data  
 76 have transformed their spatial locations away from uniform. At the bottom left they appear to have  
 77 clustered near the discontinuity of the test function while in the remaining corners of the space the

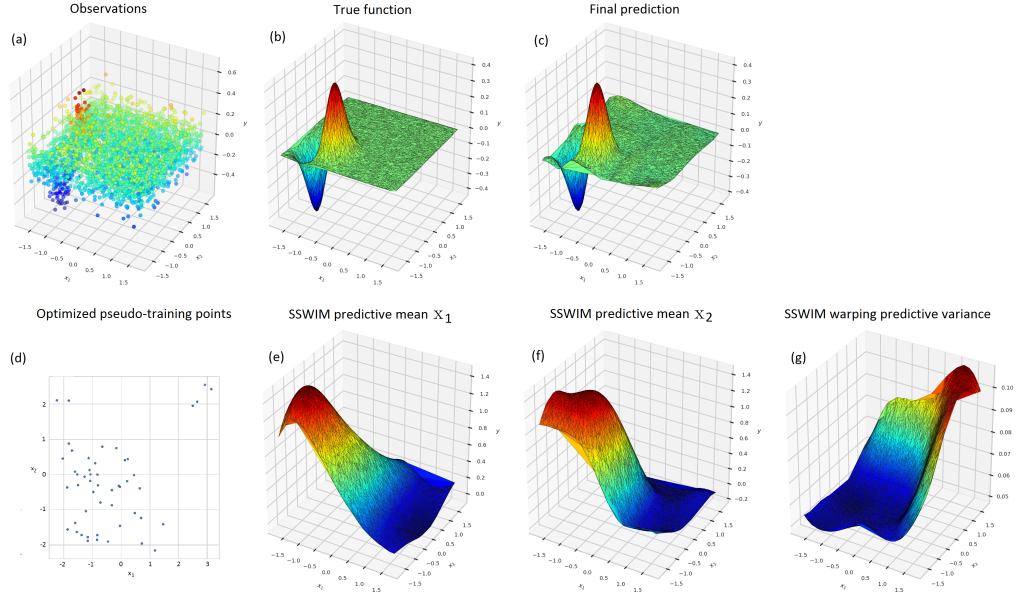


Figure 2: Visualisation of learning spatial positioning of pseudo training points in 2D. We demonstrate spatial nonstationarity with the exponential 2D function from [4]. (a) Noisy training data, (b) True function surface, (c) SSWIM prediction conditioned on training data, (d) Learned pseudo-training point positions, (e) Learned warping predictive mean for  $x_1$ , (f) Learned warping predictive mean for  $x_2$ , (g) Learned warping predictive variance.

points have spread away. The predictive mean of the learned warping function is thus visualised across the domain, for both  $x_1$  and  $x_2$ , in Figure 2 (e) and (f). Furthermore, Figure 2 (g) shows the predictive variance of the warping function and we can see how a lower amount of spatial uncertainty arises both from the noisy pseudo-targets as well as the pseudo-inputs from (d).

## D Datasets and Experimental Conditions

For completeness, we specify for each dataset used, the data dimensionality and sample size, the raw data source, the modelling objective (i.e. the target) as defined for the original problem, and any target variable pre-processing excluding standardisation. The dimensionality  $D$  is reported for the inputs  $X$  (i.e. excluding the target variable  $y$ ). All problems are single output regression tasks. Number of samples  $N$  is reported for the entire dataset before train/test splitting is applied. Note that we do not alter the raw data files and pre-processing is applied through code exactly as in the provided supplementary code. *dataloader.py*.

### elevators

$D$  : 18

$N$  : 8751

Source: [https://web.archive.org/web/\\*/http://www.liacc.up.pt/~ltorgo/](https://web.archive.org/web/*/http://www.liacc.up.pt/~ltorgo/)

Regression/\*

Preprocessing: None

Target: goal

### airfoil

$D$  : 5

$N$  : 1503

Source: <https://archive.ics.uci.edu/ml/datasets/Airfoil+Self-Noise>

101 *Preprocessing:* None  
 102 *Target:* Sound pressure in decibels

103 **concrete**

104 *D* : 8  
 105 *N* : 1030  
 106 *Source:* <https://archive.ics.uci.edu/ml/datasets/concrete+compressive+strength>  
 107 *Preprocessing:* None  
 108 *Target:* Compressive Strength

109 **parkinsons**

110 *D* : 16  
 111 *N* : 5875  
 112 *Source:* <https://archive.ics.uci.edu/ml/datasets/Parkinsons+Telemonitoring>  
 113 *Preprocessing:* Drop the first 5 columns as they are not used in the original problem  
 114 *Target:* Total udpr

115 **bikeshare**

116 *D* : 15  
 117 *N* : 17379  
 118 *Source:* <https://archive.ics.uci.edu/ml/datasets/Bike+Sharing+Dataset>  
 119 *Preprocessing:* None  
 120 *Target:* Number of bike shares per hour

121 **ct slice**

122 *D* : 379  
 123 *N* : 53500  
 124 *Source:* <https://archive.ics.uci.edu/ml/datasets/Relative+location+of+CT+slices+on+axial+axis>  
 125 *Preprocessing:* Drop Patient ID. Drop columns which have constant value throughout entire dataset.  
 126 *Target:* Reference (relative location)

128 **supercond**

129 *D* : 81  
 130 *N* : 21263  
 131 *Source:* <https://archive.ics.uci.edu/ml/datasets/Superconductivity+Data>  
 132 *Preprocessing:* None  
 133 *Target:* Critical Temperature

134 **protein**

135 *D* : 9  
 136 *N* : 45730  
 137 *Source:* <https://archive.ics.uci.edu/ml/datasets/Physicochemical+Properties+of+Protein+Tertiary+Structure>  
 138 *Preprocessing:*  $\log(1 + y)$  transform for target  $y$   
 139 *Target:* RMSD

141 **buzz**

142 *D* : 77  
 143 *N* : 583250  
 144 *Source:* <https://archive.ics.uci.edu/ml/datasets/Buzz+in+social+media+#>  
 145 *Preprocessing:*  $\log(1 + y)$  transform of target  $y$   
 146 *Target:* Mean Number of Active Discussion (NAD)

147 **song**

148 *D* : 90

149 *N* : 515345

150 *Source*: <https://archive.ics.uci.edu/ml/datasets/Buzz+in+social+media+#>

151 *Preprocessing*: None

152 *Target*: Year of song release

153 **abalone**

154 *D* : 9

155 *N* : 4177

156 *Source*: [https://web.archive.org/web/\\*/http://www.liacc.up.pt/~ltorgo/](https://web.archive.org/web/*/http://www.liacc.up.pt/~ltorgo/)

157 *Regression*/\*

158 *Preprocessing*: None

159 *Target*: Number of Rings

160 **creep**

161 *D* : 30

162 *N* : 2066

163 *Source*: <http://www.phase-trans.msm.cam.ac.uk/map/data/materials/creeprupt-b.html#down>

164 *Preprocessing*: None

165 *Target*: Rupture stress

167 **aileron**

168 *D* : 40

169 *N* : 7154

170 *Source*: [https://web.archive.org/web/\\*/http://www.liacc.up.pt/~ltorgo/](https://web.archive.org/web/*/http://www.liacc.up.pt/~ltorgo/)

171 *Regression*/\*

172 *Preprocessing*: None

173 *Target*: goal

174 **References**

- 175 [1] M. Lázaro-Gredilla, “Bayesian warped Gaussian processes,” in *Advances in Neural Information Processing*  
176 *Systems (NIPS)*, 2012.
- 177 [2] E. Snelson, Z. Ghahramani, and C. E. Rasmussen, “Warped Gaussian processes,” in *Advances in Neural*  
178 *Information Processing Systems (NIPS)*, 2004.
- 179 [3] C. E. Rasmussen, “Gaussian processes in machine learning,” in *Advanced lectures on machine learning*,  
180 Springer, 2004.
- 181 [4] R. B. Gramacy, *Bayesian treed Gaussian process models*. PhD dissertation, 2005.