

Proceedings

A Progressive Framework for Sparse Gaussian Process Regression

V.R. Lalchand ^{1,2}, A.C. Faul ^{1,2}

¹ University of Cambridge

² The Alan Turing Institute

* Correspondence: vr308@cam.ac.uk

Version July 27, 2018 submitted to Proceedings

Abstract: In their standard form Gaussian processes (GPs) provide a powerful non-parametric framework for regression and classification tasks. Their one limiting property is their $\mathcal{O}(N^3)$ scaling where N is the number of training data points. In this paper we present a framework for GP training with sequential selection of training data points using an intuitive selection metric. The greedy forward selection strategy is devised to target two factors - regions of high predictive uncertainty and underfit. Under this technique the complexity of GP training is reduced to $\mathcal{O}(M^3)$ where ($M \ll N$) if M data points (out of N) are eventually selected. The sequential nature of the algorithm circumvents the need to invert the covariance matrix of dimension $N \times N$ and enables the use of favourable matrix inverse update identities. We outline the algorithm and sequential updates to the posterior mean and variance at each stage of training. We demonstrate our method on selected one dimensional functions to aid visualisation and show that the loss in accuracy due to using a subset of data points is marginal compared to the computational gains.

Keywords: gaussian processes; bayesian; on-line training; greedy approach; sequential training

1. Introduction

A GP model is nonparametric, this means that the complexity of the model grows as more data is observed. Another attractive feature of GP models is the behaviour of the predictive variance (also referred to in this context as uncertainty), which is naturally higher in regions away from the training data. This is intuitive as in regions where there is no training data there is higher uncertainty about the interpolating function. The application of GPs to the regression task involves the computation of a matrix inverse, this leads to the $\mathcal{O}(N^3)$ scaling. Further, $\mathcal{O}(N^2)$ space is required to store a dense covariance matrix in memory.

There has been significant interest in finding sparse approximations to the full Gaussian process in order to speed up training and prediction times to $\mathcal{O}(NM^2)$ where M is the size of an auxiliary set, typically a subset of the training data. [1] provides an excellent overview of sparse approximations. A common theme of the sparse methods involves developing a low-rank approximation to the covariance matrix. We describe the most influential approaches in section 3 but to aid in understanding a brief introduction to GPs and regression follows first.

2. Gaussian Processes for Regression

Here we provide a concise summary of GPs for regression, see Chapter 2 of [2] for a detailed account. We can think of a mathematical function $f(x)$ as a collection of random variables defined at locations $x \in X$. A GP can be thought of as defining a probability distribution (prior) over function values $f(x)$. Notationally, we say f is distributed as a GP with mean $m(x)$ and covariance function $k(x, x')$. The

34 covariance function typically depends on some *hyperparameters* that control various properties of the
35 underlying f .

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}')) \quad (1)$$

36 In the context of regression, we are given a training set $\{(\mathbf{x}_i, y_i), i = 1, \dots, N\}$ where $\mathbf{y} = \{y_i\}_{i=1}^N \in \mathbb{R}$
37 are noisy realisations of a latent function f and $X = \{\mathbf{x}_i\}_{i=1}^N \in \mathbb{R}^d$ are d -dimensional inputs,
38 $y_i = f(\mathbf{x}_i) + \mathcal{N}(0, \sigma_n^2)$. The objective is to predict f_* for a new unseen point \mathbf{x}_* .

39
40 Assuming a zero mean GP prior on the latent function $f \sim \mathcal{GP}(0, k(\cdot, \cdot))$ we can write any finite
41 collection of f values as a multivariate Gaussian. In particular, evaluations of f at X , $\mathbf{f} = \{f(\mathbf{x}_i)\}_{i=1}^N$
42 are distributed as a multivariate Gaussian,

$$p(\mathbf{f}|X)^1 = \mathcal{N}(0, K(X, X)) \quad (2)$$

43 where $K(X, X)$ is the covariance matrix constructed by applying the covariance function to the X input
44 points. In our experiments, we use the *squared exponential* (SE) covariance function; section A of the
45 appendix encloses a discussion on covariance functions.

46
47 According to the prior, the targets/outputs \mathbf{y} are a sample from a multivariate Gaussian distribution
48 $\mathbf{y} \sim \mathcal{N}(0, K + \sigma_n^2 \mathbb{I})$ (where we use the shorthand K for $K(X, X)$). Hence, the joint prior distribution of
49 the training targets \mathbf{y} and test output f_* is,

$$\begin{bmatrix} \mathbf{y} \\ f_* \end{bmatrix} \sim \mathcal{N}\left(0, \begin{bmatrix} K + \sigma_n^2 I & K_*^T \\ K_* & K_{**} \end{bmatrix}\right) \quad (3)$$

50 where T denotes matrix transposition and the mean is a $(N + 1)$ zero vector and the covariance matrix
51 has dimension $(N + 1) \times (N + 1)$. For a new input point \mathbf{x}_* , the covariance between the input point
52 and all training data points is summarised by the row vector $K_* = [k(\mathbf{x}_*, \mathbf{x}_1) \ k(\mathbf{x}_*, \mathbf{x}_2) \ \dots \ k(\mathbf{x}_*, \mathbf{x}_N)]$
53 and the covariance with itself (variance) is denoted as $K_{**} = k(\mathbf{x}_*, \mathbf{x}_*)$.

54
55 Once we have specified the joint prior of the training data points along with the test data point, we can
56 derive the posterior distribution by conditioning the joint prior in eq. 3 on the training data points
57 \mathbf{y} . The posterior defined as $p(f_*|\mathbf{y})$ can be derived in closed form using properties of conditional
58 Gaussian distributions (see Appendix A.1 of [2]).

$$\begin{aligned} f_*|\mathbf{y} &\sim \mathcal{N}(\bar{f}_*, \text{var}(f_*)) \\ \bar{f}_* &= K_*(K + \sigma_n^2 I)^{-1} \mathbf{y} \\ \text{var}(f_*) &= K_{**} - K_*(K + \sigma_n^2 I)^{-1} K_*^T \end{aligned} \quad (4)$$

59 The hyperparameters of the covariance function $k(\cdot, \cdot)$ used to generate the covariance matrix K are
60 discovered through optimisation of the marginal likelihood. Gradient based optimizers can be used to
61 minimize the negative log of the marginal likelihood given by,

$$\log p(\mathbf{y}|X) = -\frac{1}{2} \mathbf{y}^T (K + \sigma_n^2 I)^{-1} \mathbf{y} - \frac{1}{2} \log |K + \sigma_n^2 I| - \frac{N}{2} \log(2\pi) \quad (5)$$

62 which follows directly from the fact that, $\mathbf{y} \sim \mathcal{N}(0, K + \sigma_n^2 I)$.

63

¹ We omit the explicit conditioning on X for the sake of brevity

64 Also, the prior $p(\mathbf{f})$ times the likelihood given by $p(\mathbf{y}|\mathbf{f}) = \mathcal{N}(\mathbf{f}, \sigma_n^2)$ gives the marginal likelihood in 5
65 in closed form as the distributions in the integral are both Gaussian.

$$p(\mathbf{y}|X) = \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f})d\mathbf{f} \quad (6)$$

66 3. Summary of sparse approximation approaches

67 The mean and variance of the predictive posterior distribution in eq. 4 involves inversion of matrix
68 $K + \sigma_n^2\mathbb{I}$. The predictive posterior $p(f_*|\mathbf{y})$ is directly derived from the joint prior in eq. 3. Hence, an
69 approximation to the joint prior $p(\mathbf{y}, f_*)$ using a lower rank matrix $Q_{M \times M}$ where $M \ll N$ instead of
70 the full covariance matrix K will reduce the complexity of predictive inference.

71
72 Note that,

$$p(f_*|\mathbf{y}) = \int p(\mathbf{f}, f_*|\mathbf{y})d\mathbf{f} \quad (7)$$

73 and

$$p(\mathbf{f}, f_*|\mathbf{y}) = \frac{p(\mathbf{f}, f_*)p(\mathbf{y}|\mathbf{f})}{p(\mathbf{y})} \quad (8)$$

74 Essentially, the joint prior, the likelihood and the posterior are interlinked through the Bayes' rule.
75 Most approaches proposed in literature start with a modified likelihood $p(\mathbf{y}|\mathbf{f}) \approx q(\mathbf{y}|\mathbf{u})$ where
76 $\mathbf{u} = [u_1, \dots, u_m]$ are some arbitrary inducing function values, typically $\mathbf{u} \in \mathbf{y}$. The approaches mainly
77 differ in terms of the theoretical approach in deriving the approximations and how the inducing
78 values are selected.

79
80 [3], [4] and [5] all use the likelihood approximation,

$$p(\mathbf{y}|\mathbf{f}) \approx q(\mathbf{y}|\mathbf{u}) = \mathcal{N}(K_{\mathbf{f},\mathbf{u}}K_{\mathbf{u},\mathbf{u}}^{-1}\mathbf{u}, \sigma_n^2) \quad (9)$$

81 [2] in chapter 8 provides a useful interpretation for this approximation calling it, the *projected process*
82 *approximation* (PPA) as it uses the projection $\mathbf{f} = K_{\mathbf{f},\mathbf{u}}K_{\mathbf{u},\mathbf{u}}^{-1}\mathbf{u}$. The PPA construction uses a likelihood
83 involving all N data points while the posterior predictive distribution involves inverting only a
84 $M \times M$ matrix. The key feature being that information contained in the N points leaves an imprint on
85 the posterior predictive distribution.

86
87 Methods proposed in [6] [5] integrate the idea of online training and sparse representation providing
88 update equations for the posterior mean and covariance. [1] reinterprets all the existing sparse
89 representations and derives the joint prior $p(\mathbf{f}, f_*) \approx q(\mathbf{f}, f_*)$ implied by the different likelihood
90 approximation schemes in literature. Once the approximate joint prior is specified, the predictive
91 distribution $q(f_*|\mathbf{y})$ is derived the same way as for the full GP. The authors point out that this
92 standardisation enables one to compare the different representations.

93
94 The likelihood approximation schemes still require the specification of the choice of inducing variables
95 \mathbf{u} . Let us refer to them as the active set I . While one can select the active set through random sampling,
96 studies have shown that the approach does not work as well as actively selecting the points as
97 optimisers of a selection metric. [7] proposes using a differential entropy score which is a function of
98 the predictive variance to systematically include points in the active set, while [3] proposes using an
99 information gain criteria which can be computed cheaply.

100
101 [8] propose a sparse GP regression scheme using *pseudo inputs* in which active set inputs ($X_{\mathbf{u}}$) and the
102 hyperparameters (θ) for the covariance function are selected through joint optimization of the marginal

103 likelihood. Thus it converts a discrete optimization problem of subset selection to a continuous
 104 paradigm where the active set is not constrained to lie within the training set. The method addresses
 105 the dual selection problems of the active set I and the kernel hyperparameters θ . Some authors refer to
 106 this technique as the *Fully Independent Training Conditional* (FITC), (see section 6 of [1] for a detailed
 107 discussion).

108 4. A Progressive framework for Gaussian Process Regression

109 In this section we describe a framework for training GPs progressively. We assume we have N pairs
 110 of data in the form of inputs X and targets \mathbf{y} , the aim is to select a smaller informative subset $\mathbf{u} \in \mathbf{y}$
 111 (the *active set*) which play an active role in the inference. All other training points $\mathbf{y} \setminus \mathbf{u}$ belong to
 112 the *remainder set*. The active set is constructed incrementally, at each iteration exactly one training
 113 data point is selected. Following the notation from [3] let I denote indices of the active set and
 114 $R = \{1, \dots, N\} \setminus I$ denote the indices of remainder points. Training happens in stages we index by t .
 115 Hence, I_t and R_t denote the index sets at stage t of training.

116
 117 We denote the active set \mathbf{u} by $\mathbf{y}(I_t)$ from here on to clearly incorporate the stage of training. Similarly,
 118 the remainder set is denoted as $\mathbf{y}(R_t)$. The active and remainder input locations are denoted as $X(I_t)$
 119 and $X(R_t)$ respectively. For the purposes of testing we have a hold-out set with N_* test inputs X_* and
 120 targets \mathbf{y}_* .

Table 1. Notation used in stagewise training

At stage t				
	Index set	Inputs	Targets	Size
Active points	I_t	$X(I_t)$	$\mathbf{y}(I_t)$	t
Remainder points	R_t	$X(R_t)$	$\mathbf{y}(R_t)$	$N - t$

121 At stage t the active set has exactly t data points as at each stage exactly one point is added to the
 122 active set. We have a fixed set of N training pairs, the active set grows in size as more points are added
 123 to it and the remainder set shrinks in size as points are removed from it. Essentially, we start with all
 124 the training data points in the remainder set and points move from the remainder set to the active set
 125 in each iteration based on a selection criteria.

126
 127 Further,

- 128 • $K = K(X, X)$ denotes the full covariance matrix computed on N inputs X .
- 129 • $K_{I_t} = K(X(I_t), X(I_t))$ denotes the $t \times t$ covariance matrix computed between the active inputs at
 130 the t^{th} stage of training.
- 131 • $K_{\setminus I_t} = K(X(R_t), X(I_t))$ denotes the $(N - t) \times t$ covariance matrix computed between the t active
 132 inputs selected so far and the $(N - t)$ remainder inputs in R_t .
- 133 • $K_{R_t} = K(X(R_t), X(R_t))$ denotes the $(N - t) \times (N - t)$ covariance matrix computed between the
 134 remainder inputs at stage t .
- 135 • $K_{**} = K(X_*, X_*)$ denotes the covariance matrix computed on the test inputs X_* .
- 136 • $K_{*I_t} = K(X_*, X(I_t))$ denotes the $N_* \times t$ covariance matrix computed between the test inputs and
 137 active inputs at stage t .
- 138 • $\boldsymbol{\mu}_t$ and Σ_t denote the predictive posterior mean and covariance computed at stage t for the
 139 remainder inputs $X(R_t)$.
- 140 • $\boldsymbol{\mu}^*$ and Σ^* denote the predictive posterior mean and covariance for the test inputs X_* .

141 The algorithm starts with a single training point $\mathbf{y}(I_1) = [u_1]$ in the active set which is selected
 142 at random from \mathbf{y} , the predictive posterior mean and covariance denoted by $\boldsymbol{\mu}_1$ and Σ_1 at
 143 stage 1 are computed by conditioning on the active set (of 1 point) while the goodness of fit

144 is assessed by predicting on the remainder inputs ($N - 1$ points). In short, we predict at each
 145 stage the mean and covariance of the remainder inputs $X(R_t)$ and compare them to the true
 146 remainder targets $\mathbf{y}(R_t)$, the mean squared error computed between the true remainder targets
 147 and the predicted remainder targets provides the basis for convergence (see section D of the Appendix).

148

149 The main reason for this stagewise iterative training approach is two-fold:

150 1. The selection criteria for the active training target at each stage is tied to the predictive posterior
 151 mean and variance computed on the remainder inputs. We discuss this in more detail in the next
 152 section.

153

154 2. Since the active set is grown one point at a time, we can take advantage of favourable matrix
 155 inverse update identities in order to update the predictive posterior mean and variance at each
 156 stage.

157 In general, at stage t ,

$$p(\mathbf{y}(R_t)|\mathbf{y}(I_t)) = \mathcal{N}(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t)$$

$$\begin{aligned}\boldsymbol{\mu}_t &= K_{\setminus I_t}(K_{I_t} + \sigma_n^2 \mathbb{I})^{-1} \mathbf{y}(I_t) \\ \boldsymbol{\Sigma}_t &= K_{R_t} - K_{\setminus I_t}(K_{I_t} + \sigma_n^2 \mathbb{I})^{-1} K_{\setminus I_t}^T\end{aligned}\quad (10)$$

158 where $\boldsymbol{\mu}_t$ and $\boldsymbol{\Sigma}_t$ denote the predicted posterior mean and covariance.

159

160 The above equations reflect the predictive posterior mean and covariance for a full GP introduced in
 161 eq. 4 where the active set $\mathbf{y}(I_t)$ plays the role of the target vector \mathbf{y} .

162 4.1. The Algorithm

163 Below we give the general algorithm for active set selection using a general selection metric we denote
 164 as Δ .

Algorithm 1 Progressive framework for GPR

Initialisation: Pick a random target $u_1 \in \mathbf{y}$.

Convergence Condition: $(RMSE_{t-1} - RMSE_t) < \delta$ calculated on the remainder set $\mathbf{y}(R_t) = \mathbf{y} \setminus \mathbf{y}(I_t)$.

for each stage t :

GP Train on $(\mathbf{X}(I_t), \mathbf{y}(I_t))$.

GP Predict on $\mathbf{X}(R_t)$.

 Update posterior mean $\boldsymbol{\mu}_t$ and covariance $\boldsymbol{\Sigma}_t$. (see section C for sequential update rules.)

 Compute $\Delta_i \forall i \in R_t$

 Select i where $i = \operatorname{argmax}_{i \in R_t} \Delta_i$

$I_{t+1} \leftarrow I_t \cup \{i\}, R_{t+1} \leftarrow R_t \setminus \{i\}$

If convergence is true:

 break;

end
 return I_T, R_T

165 4.2. Selection of active points

166 In this section we propose a selection criteria Δ that is used to advance the algorithm in the greedy
 167 search. One of the most attractive features of using a GP for regression is the way in which uncertainty

168 is propagated around predictions. The posterior predicted variance contracts surrounding the training
 169 data and expands back to the prior variance away from the training data (see fig. 1).
 170

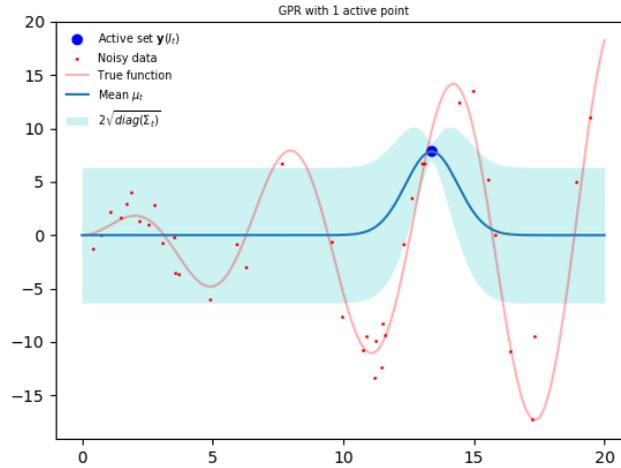


Figure 1. GPR with 1 active point

171 Hence, a simple rule which selects the next active target where the posterior predictive variance is
 172 maximised poses a problem of non-uniqueness. There are usually several remainder targets with the
 173 same posterior predicted variance (captured by the diagonal of the matrix Σ_t). The prior variance σ_f^2
 174 serves as an upper bound for the posterior predicted variance. Let us term the posterior predicted
 175 variance captured in $\text{diag}(\Sigma_t)$ as Δ_1 .
 176

177 An alternative metric that serves as a reasonable choice is the absolute error or deviation from the
 178 posterior predicted mean at each stage t to the true remainder targets. At the end of each stage t
 179 we have, the posterior predicted mean $\mu_t = [\mu_1, \dots, \mu_{N-t}]$ and the true remainder set of targets
 180 $\mathbf{y}(R_t) = \{r_i | i \in R_t\}$. For each remainder target $r_i \in \mathbf{y}(R_t)$ we can compute its absolute deviation from
 181 the mean, let us call this metric Δ_2 . See fig. 2.
 182

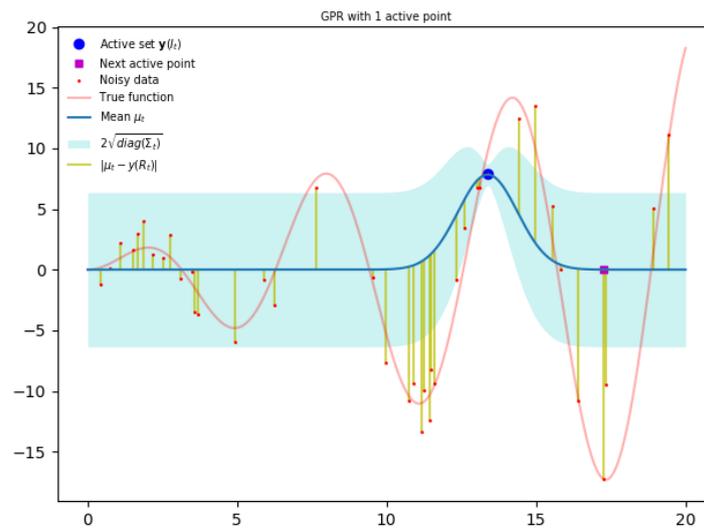


Figure 2. GPR with 1 active point using metric Δ_2 (depicted by the yellow vertical lines). According to this selection rule the next remainder point to be chosen for the active set is the red dot (noisy target) corresponding to the magenta point on the predicted mean.

183 The selection rule can be described as, $\operatorname{argmax}_{u \in y(R_t)} \Delta_2 = |\mu_t - \mathbf{y}(R_t)|$ where $|\cdot|$ denotes absolute
 184 value. This metric is unique for each remainder target value, as seen in fig. 2. However, a drawback
 185 of Δ_2 is that it could end up selecting points in regions where the uncertainty is low (due to a nearby
 186 selected active point) versus points where the uncertainty is high along with a higher deviation from
 187 the mean.

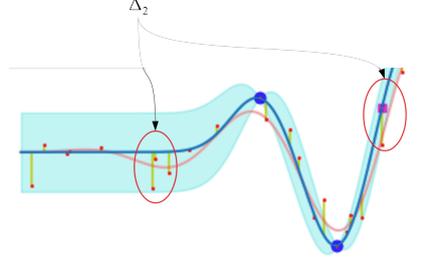


Figure 3. Snapshot of Progressive regression using Δ_2 (shown by the vertical yellow lines), the blue dots depict the active points chosen so far, the magenta square dot is the remainder point that would be chosen next if Δ_2 is used. Selecting a point in the left red circle could be a better choice as there is no active point in that region and could lead to a bigger reduction in overall uncertainty.

188 An ideal metric would penalise both high uncertainty captured by $\operatorname{diag}(\Sigma_t)$ and deviation from the
 189 posterior predicted mean $|\mu_t - \mathbf{y}(R_t)|$. The metric Δ_3 is motivated by the desire to capture the two
 190 factors of high uncertainty and underfit. It is given by,

$$\Delta_3 = \sqrt{\operatorname{diag}(\Sigma_t)} + |\mu_t - \mathbf{y}(R_t)|$$

191 *Rule 3:* $\operatorname{argmax}_{u \in y(R_t)} \Delta_3 = \sqrt{\Delta_1} + \Delta_2$

192

193 Note that both components of the addition are vectors of size $(N - t)$ as they are based on the remainder
 194 set. Further, we use $\sqrt{\operatorname{diag}(\Sigma_t)}$ as the standard deviation is comparable in scale to the deviation from
 195 posterior mean. The metric Δ_3 can be evaluated in $\mathcal{O}(1)$ as Σ_t and μ_t are obtained directly from the
 196 training step. A visual depiction of the evolution of progressive training for the function $x \sin x$
 197 is shown in fig. A3.

198 4.3. Complexity of Progressive training

199 The cost of computing an updated inverse for a square matrix of size M grown by 1 row and 1 column is
 200 quadratic $\mathcal{O}(M^2)$ if we use block inversion (Schur complement) and assuming we know the inverse of
 201 the matrix of size M . In the progressive GP framework, in each stage the covariance matrix computed
 202 on active points grows by 1 row and 1 column. If we end up with an active set of M points after M
 203 stages of training we have conducted the matrix update operation M times and the cost each time is
 204 quadratic in the dimension of the matrix we are updating. This give us the following computational
 205 cost in terms of operations.

$$\sum_{i=1}^M i^2 = 1^2 + 2^2 + \dots + M^2 = \frac{M(M+1)(2M+1)}{6} = \mathcal{O}(M^3) \quad (11)$$

206 Hence, while the order of complexity is quadratic per iteration, the overall complexity is $\mathcal{O}(M^3)$ if M
 207 is the size of the final active set. It is important to note that if the inverse is calculated directly in each
 208 iteration, the complexity is $\mathcal{O}(M^4)$; hence, the update with the Schur complement is essential.

209

210 The table below highlights the computational complexity of the full GP and the progressive approach
 211 to training. Note that, we still need to compute and store the full covariance matrix in order to speed
 212 up the matrix updates in the progressive GP approach.

Task	Full GP	Progressive GP
213 Training	$\mathcal{O}(N^3)$	$\mathcal{O}(M^3)$
Prediction	$\mathcal{O}(N^2)$	$\mathcal{O}(M^2N)$
Storage (for K)	$\mathcal{O}(N^2)$	$\mathcal{O}(N^2)$

214 4.4. Selection of hyperparameters

215 The hyperparameters for the SE kernel used in the covariance matrix $\theta = (\sigma_f^2, l, \sigma_n^2)$ are pre-selected
 216 through optimization of the log marginal likelihood (LML) using a random subset of the training data
 217 in a pre-processing step. During the running of the algorithm, the hyperparameters remain fixed. In
 218 this paper, we mainly focus our efforts at providing a framework for selecting active targets and inputs
 219 from the training points while simultaneously training the GP. A framework that weaves together
 220 the hyperparameter selection and active set selection in the context of progressive training of GPs is
 221 being researched. Preliminary experiments where we varied the hyperparameters during stagewise
 222 training using marginal likelihood optimisation lead to instability. The authors of [8] highlight that
 223 active selection causes non-smooth fluctuations in the marginal likelihood making the optimisation
 224 difficult. Hence, a different way for concomitant hyperparameter optimisation needs to be developed.

225 5. Experiments

226 We trained a GP using the progressive training approach by sampling noisy (in order to have a
 227 systematic comparison the noise level for all the experiments was identical, further the random subset
 228 was ensured to be the same size as that for the progressive GP) values from a host of $1d$ functions; we
 229 then predict on a hold out unseen test set X_* . We report the generalisation error (RMSE) under three
 230 training schemes. The full GP, the random subset and the progressive GP algorithm.

Table 2. RMSE on Test data

Function \ Data	Full GP	Random	Progressive GP	% of full dataset
$x^2 \sin(x)$	32.24	91.62	39.29	22%
$x \sin(x)$	2.36	5.95	2.82	18%
$0.5 \sin(x) + 0.5x - 0.02(x - 5)^2$	1.14	2.17	1.96	31%

231 6. Discussion and Further Work

232 Experimentation with some standardised high dimensional datasets is a natural avenue to explore
 233 next. Further, we need to investigate the dynamics of choosing hyperparameters along with the active
 234 points from the data. It is evident that the hyperparameters chosen during initialisation do effect the
 235 choice of active points selected during the greedy search. The mechanics of this interplay between
 236 hyperparameters and active points needs to be investigated. The main criticism of this approach is
 237 that remainder points in $\mathbf{y}(R_t)$ don't play a role in the inference except in a weak sense of guiding the
 238 algorithm to the next stage. This information loss can be justified in situations where there is known
 239 to be redundancy in the data and the prediction based on the subset is accurate enough for the task.
 240 An important caveat is that the progressive approach which prescribes a method to select M points
 241 can be combined with a likelihood approximation scheme, some of which were discussed in section
 242 3 to build a more powerful sparse GP algorithm; and one that makes use of the $N - M$ remainder
 243 points. For instance, the reduced rank approximation of the matrix $K_{NN} \sim \tilde{K} = K_{NM}K_{MM}^{-1}K_{MN}$ can be
 244 used instead of the full covariance matrix K_{NN} . In conclusion, we have proposed a simple criteria to
 245 progressively train GPs along with update equations for the mean and covariance. The progressive
 246 GP approach yields sparse approximations to the full GP with marginal loss in prediction accuracy
 247 relative to training with the full dataset.

248 **Acknowledgments:** This work was supported by the Alan Turing Institute through the Doctoral Studentship for
 249 International Students.

250 Abbreviations

251 The following abbreviations are used in this manuscript:

252 GP Gaussian processes
 253 GPR Gaussian process regression

254 Appendix A. The covariance function

255 The covariance function controls the properties of functions a GP can model, for example, smoothness,
 256 stationarity and shape; and can typically have multiple parameters called *hyperparameters*. The
 257 covariance function k takes as input, two input locations and hence can be viewed as representing
 258 some form of a similarity between the two data points. The process of learning the hyperparameters of
 259 the covariance function is an important part of GP training.

260
 261 The *squared exponential*(SE) kernel is defined as,

$$k(x_i, x_j) = \sigma_f^2 \exp\left(-\frac{(x_i - x_j)^2}{2l^2}\right) + \delta_{ij}\sigma_n^2 \quad (\text{A1})$$

262 where $\theta = (\sigma_f^2, l, \sigma_n^2)$ is the set of hyperparameters, comprising the signal variance $\sigma_f^2 > 0$ which
 263 controls the variation of function values from their mean, the lengthscale $l > 0$ which controls how
 264 smooth a function is and $\sigma_n^2 \geq 0$ is the noise variance which allows for noise to be present in the
 265 data. The noise variance applies only when $i = j$. In the noiseless case, we just drop the additive noise
 266 variance term.

267

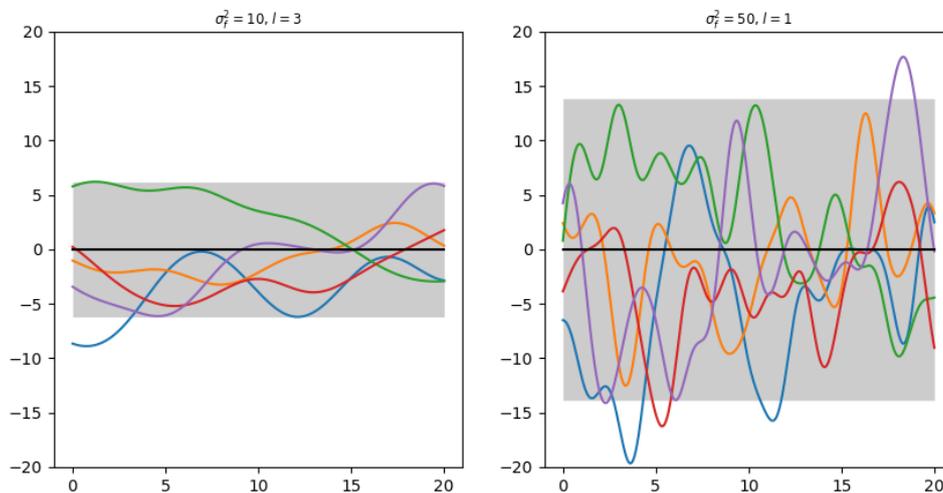


Figure A1. Samples from GP prior with SE covariance function and 95% confidence intervals ($\pm 1.96\sigma_f$)

268 Before performing GPR we calculate the covariance matrix by applying the covariance function to all
 269 possible combinations of the input points,

$$K(X, X) = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \dots & k(\mathbf{x}_1, \mathbf{x}_N) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \dots & k(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & k(\mathbf{x}_N, \mathbf{x}_2) & \dots & k(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix} \quad (\text{A2})$$

270 For the sake of notational convenience we denote $K(X, X)$ as K .

271 Appendix B. GP Regression (GPR)

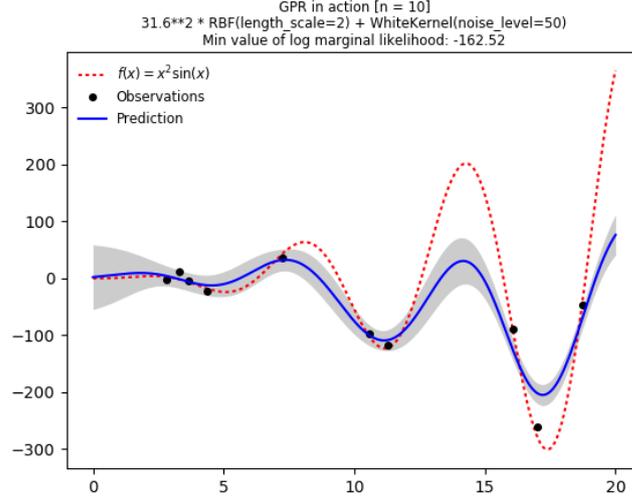


Figure A2. GPR in action where the kernel hyperparameters have been optimized through minimization of the negative log marginal likelihood

272 An important aspect of GPR is the selection of hyperparameters for the covariance function. The SE
 273 kernel shown in A1 has three free parameters. There are several approaches to setting the values of
 274 these hyperparameters based on training data, the joint problem of selecting an appropriate kernel
 275 function and its hyperparameters is called *model selection* [2].

276

277 Appendix C. Progressive updates

278 In this section, we discuss in detail the stagewise updates to the posterior predicted mean and
 279 covariance given in eq. 10.

280 Appendix C.0.1. Mean

$$\text{At stage } t: \mu_t = K_{\setminus I_t} (K_{I_t} + \sigma^2 \mathbb{I})^{-1} \mathbf{y}(I_t)$$

↓

$$\text{At stage } t+1: \mu_{t+1} = K_{\setminus I_{t+1}} (K_{I_{t+1}} + \sigma^2 \mathbb{I})^{-1} \mathbf{y}(I_{t+1})$$

281 Notice that the mean update from stage $t \rightarrow t + 1$ involves updating two matrices.

282

First,

$$K_{\setminus I_t} \rightarrow K_{\setminus I_{t+1}}$$

283 In this update we evolve a $(N - t) \times t$ matrix into a $(N - (t + 1)) \times (t + 1)$. We are dropping a row
 284 and adding a column. The newly added column contains the covariances computed between the newly
 285 selected active point, say s and all the $(N - (t + 1))$ points in the remainder set $(k(s, r_i) | i \in R_{t+1})$.
 286 If we assume that the full covariance matrix $K(X, X)$ is computed at the start then this update just
 287 requires selecting the corresponding entries from the matrix $K(X, X)$. The complexity of this operation
 288 is just $\mathcal{O}(1)$.

289

290 Second,

$$K_{I_t} \rightarrow K_{I_{t+1}}$$

291 This is the covariance matrix computed on the active points, it grows by 1 row and 1 column in each
 292 stage.



293 Further, its inverse is required in each stage. We make use of block inversion to update the inverse.
 294 The complexity of this operation is quadratic per iteration, this is shown in section 4.3.

295 Appendix C.0.2. Covariance

$$\begin{aligned} \text{At stage } t: \Sigma_t &= K_{R_t} - K_{\setminus I_t} (K_{I_t} + \sigma^2 \mathbb{I}_t)^{-1} K_{\setminus I_t}^T \\ &\downarrow \\ \text{At stage } t+1: \Sigma_{t+1} &= K_{R_{t+1}} - K_{\setminus I_{t+1}} (K_{I_{t+1}} + \sigma^2 \mathbb{I}_t)^{-1} K_{\setminus I_{t+1}}^T \end{aligned}$$

296 The updates of $K_{\setminus I_t} \rightarrow K_{\setminus I_{t+1}}$ and $K_{I_t} \rightarrow K_{I_{t+1}}$ were already discussed in the previous section. The only
 297 new matrix update involved here is,

$$K_{R_t} \rightarrow K_{R_{t+1}}$$

298 Since, K_{R_t} is the covariance matrix computed on the remainder inputs and the size of the remainder
 299 inputs shrinks as the training progresses, $K_{R_t} \rightarrow K_{R_{t+1}}$ involves dropping a row and a column as the
 300 matrix shrinks from size $(N - t) \times (N - t)$ to $(N - (t + 1)) \times (N - (t + 1))$.

301 Appendix D. Convergence of progressive training

302 At each stage of progressive training, the algorithm tracks the generalisation error by computing the
 303 posterior predictive mean and covariance on the remainder set $X(R_t)$. We define generalisation error
 304 as the square root of the sum of squared errors (abbreviated as RMSE) between the mean μ_t and the
 305 true remainder targets $\mathbf{y}(R_t) = \{r_i | i \in R_t\}$.

306

The RMSE is given by the l^2 norm,

$$\|\boldsymbol{\mu}_t - \mathbf{y}(R_t)\|_2 = \sum_{i=1}^{N-t} (\mu_i - r_i)^2 \quad (\text{A3})$$

307 where the mean prediction for test inputs at stage t is given by eq. 10.

308

309 If we denote the RMSE at stage t as ζ_t the algorithm terminates when,

$$\zeta_{t-1} - \zeta_t < \delta \quad (\text{A4})$$

310 where δ denotes a threshold value usually tied to the range of the target values \mathbf{y} .

311 **References**

- 312 1. Quiñero-Candela, J.; Rasmussen, C.E. A unifying view of sparse approximate Gaussian process
313 regression. *Journal of Machine Learning Research* **2005**, *6*, 1939–1959.
- 314 2. Rasmussen, C.E. Gaussian processes in machine learning. In *Advanced lectures on machine learning*; Springer,
315 2004; pp. 63–71.
- 316 3. Seeger, M.; Williams, C.; Lawrence, N. Fast forward selection to speed up sparse Gaussian process
317 regression. *Artificial Intelligence and Statistics 9*, 2003, number EPFL-CONF-161318.
- 318 4. Smola, A.J.; Bartlett, P.L. Sparse greedy Gaussian process regression. *Advances in neural information
319 processing systems*, 2001, pp. 619–625.
- 320 5. Csató, L.; Opper, M. Sparse on-line Gaussian processes. *Neural computation* **2002**, *14*, 641–668.
- 321 6. Huber, M.F. Recursive Gaussian process: On-line regression and learning. *Pattern Recognition Letters* **2014**,
322 *45*, 85–91.
- 323 7. Lawrence, N.D.; Seeger, M.; Herbrich, R. Fast sparse gaussian process methods: The informative vector
324 machine. *Advances in Neural Information Processing Systems 15*. Citeseer, 2003.
- 325 8. Snelson, E.; Ghahramani, Z. Sparse Gaussian processes using pseudo-inputs. *Advances in neural
326 information processing systems*, 2006, pp. 1257–1264.

327 © 2018 by the authors. Submitted to *Proceedings* for possible open access publication
328 under the terms and conditions of the Creative Commons Attribution (CC BY) license
329 (<http://creativecommons.org/licenses/by/4.0/>).

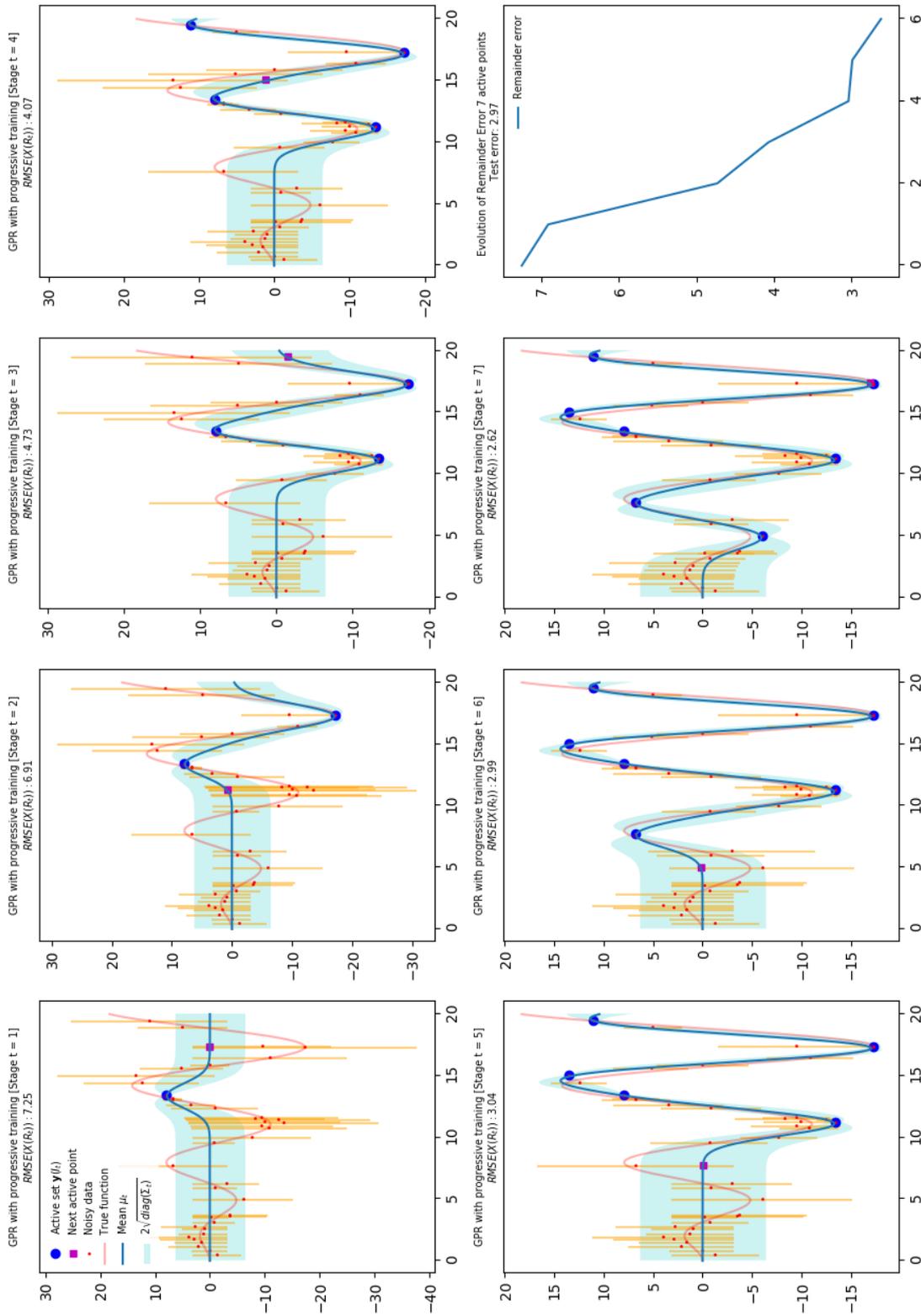


Figure A3. The progressive GP framework in action. The orange vertical bars show the Δ_3 metric on the basis of which the next active point is selected. The algorithm terminates after 7 iterations when the reduction in the RMSE on the remainder set is below a preset threshold.