

# ROBUST GAUSSIAN PROCESS REGRESSION WITH HUBER LIKELIHOOD

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Outliers in both covariates and output responses pose significant challenges for Gaussian Process (GP) regression models. We present a novel GP regression approach that effectively integrates the Huber likelihood into the GP framework—with additional parameters that can be set before inference. Specifically, we model the likelihood of observed outputs using the Huber probability distribution: this reduces deviations caused by output outliers. For covariate outliers, we introduce projection pursuit weights—attenuating their influence on the model. To address the analytically intractable, yet unimodal, posterior distribution, we employ Laplace approximation and Gibbs sampling within a Markov Chain Monte Carlo (MCMC) framework. We simplify Gibbs sampling by expressing the likelihood associated with outlying points as normally distributed through the scale mixture representation of the Laplace distribution. This work is particularly important in the field of transmission spectroscopy—where noisy measurements are often neglected in the estimation of planet-to-star radius ratios. We demonstrate the robustness and effectiveness of our method through extensive experiments on synthetic and real-world datasets.

## 1 INTRODUCTION

Bayesian inference which is based on Gaussian likelihood is known to be sensitive to extreme observations and gross errors, called outliers. The estimation of parameters in Gaussian processes (GPs) is affected in non-Gaussian error settings as the predictive uncertainty assigns equal confidence to the measurements, regardless of whether they are outliers or not. We illustrate this problem in a numerical example. Let us consider a 2-d sinc function  $y(x) = \text{sinc}(x) + e$ , where  $x = \sqrt{x_1^2 + x_2^2}$  with an additive error that follows the Student’s t-distribution with 10 degrees of freedom  $e \sim \text{Student’s-t}(2)$ . We add additional large outliers  $y^{(l)}$  with magnitude close to 0.8 and  $x_1^{(l)}$ . Figure 1(a) shows the predicted values at test points  $x = [-10, 10]$ , obtained from standard.

Existing studies addressing the outlier problem in GP regression use various approaches to define the likelihood. Two common strategies are: (1) using a mixture of two normal distributions or (2) employing heavy-tailed distributions. Most of these methods assume the error distribution is known a priori—a condition that is often unrealistic in practical applications. Moreover, their robustness is questionable when faced with extreme observations that do not correspond to the non-normal distribution their heavy tailed likelihood is specified to capture. These models typically struggle to handle both general noise patterns and large errors in covariate and response dimensions, often attempting to fit extreme values. We show this shortcoming in Figure 1(b) with the sinc function data for the GP with the Student’s t-likelihood and employing the MCMC integration approximation method. We notice that the model overfits when large outliers  $x_1^{(l)}$  and  $y^{(l)}$  occur simultaneously, as the Student’s-t likelihood can effectively compensate only for errors in  $y^{(l)}$ .

In this paper, we propose a new way of handling extreme outliers in covariate space and output responses that models the likelihood of the observed data using Huber density function. We significantly enhance downweighting of the outliers compared to the earlier work by Altamirano et al. (2024), which was limited to handling outliers only in the output responses with added hyperparameters  $(\beta, c)$ .

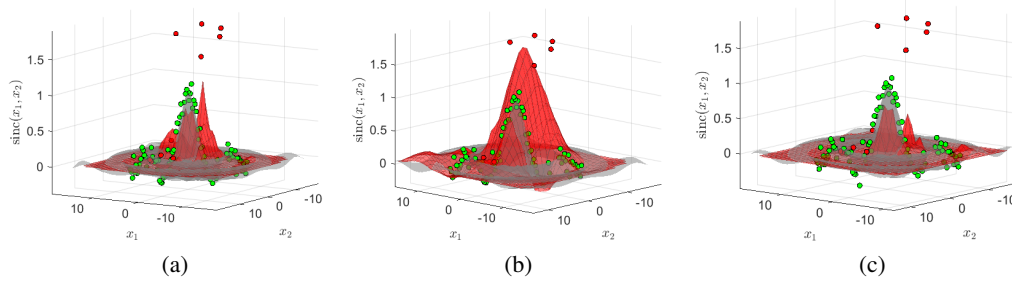


Figure 1: Predictions for  $\text{sinc2d}(\mathbf{x})$ : (a) standard GP, (b) Student’s t-likelihood via MCMC, (c) GP with Huber likelihood. The red surface shows the mean of the model predictions, the grey surface represents the true  $\text{sinc2d}$  function, red dots are outliers, and green dots are training points fitting  $\text{sinc2d}$ . GP with Huber likelihood demonstrates robustness to the outliers  $\{y^{(l)}, \mathbf{x}^{(l)}\}$ .

## 2 RELATED WORK

Goldberg et al. (1997) introduced a dual-model Gaussian process framework to account for covariate-dependent noise. The first Gaussian process model governs the output process  $y$ , while the second Gaussian process governs the noise process. West (1984) investigated heavy-tailed error distributions that are constructed as scale mixtures of normal distributions, which are also used for specifying a priori distribution based on the earlier ideas suggested by De Finetti (1961); Ramsey & Novick (1980). By doing so, the prior distribution discounts any observations highlighting inconsistency between likelihood and prior. Along the same line, Desgagné & Gagnon (2019) assumed a super heavy-tailed error distribution dependent on an explanatory variable to make the estimation of the population mean and ratios robust to outliers. Kuss (2006) extended a mixture of two normal distributions, one to model small errors in regular observations and a second one to model large errors in outlying observations. However, Naish-Guzman & Holden (2007) questioned the adequacy of the two-model approach. They proposed instead twin GP that allow us to choose between the distribution of the regular observations and that of the outliers. Kuss (2006) suggested a GP with a Laplace likelihood model that utilizes a scale mixture representation of Laplace noise distribution where the variance follows an exponential distribution. Vanhatalo et al. (2009) proposed a GP model based on the Student’s t-likelihood function, where the noise is modeled as a scale mixture of Gaussian distributions. Unfortunately with the non-Gaussian likelihood, the Bayesian inference becomes analytically intractable. Consequently, various advanced approximation methods were proposed Kuss (2006); Vanhatalo et al. (2009); Jylänki et al. (2011); Ranjan et al. (2016); Daemi et al. (2019) to overcome the convergence failure of the classical approximation methods such as expectation propagation Minka (2013), Markov Chain Monte Carlo Neal (1997), variational Bayes Ghahramani & Beal (2000), and Laplace approximation Williams (1996). More recently, Li et al. (2021); Andrade & Takeda (2023) presented a robust variants of GPs for datasets with substantial contamination removing the outlier data based on trimming parameters in iterative manner.

In GP regression models with Student’s t-likelihoods Kuss (2006), a scale-mixture representation of the Student’s t-distribution is utilized. A variational approximation is devised presuming the Gaussian likelihood whose individual variances are Gamma distributed. Combined with the Kullback-Leibler divergence,  $\text{KL}(q||p)$ , between the true posterior,  $p$ , and the approximation,  $q$ , an expectation maximization (EM)-type algorithm is implemented. As for the models with Laplace likelihoods, the scale mixture model yields a unimodal posterior enabling the implementation of the EP approximation and the MCMC sampling. Here, a Laplace approximation is inappropriate because the discontinuous derivatives of the Laplace likelihood at zero may cause the Hessian matrix to be undefined.

## 3 CONTRIBUTIONS

Altamirano et al. (2024) proposed a robust Gaussian Process (GP) regression method that leverages generalized Bayesian inference to preserve computational conjugacy. Their

method handles outliers in the output responses through weighting mechanism  $J$  in the noise term:  $\sigma^2 J_{ii} = \sigma^2 (1 + r_i^2/c^2)$ , where  $r_i$  is the residual associated with  $i^{\text{th}}$  data point  $r_i = y_i - m(\mathbf{x}_i)$  and  $c$  is the threshold parameter. However, a potential limitation of this approach is that it may not adequately account for outliers in the output response,  $y_i^{(l)}, y_j^{(c)}$ , when they occur alongside outliers in the covariate dimensions,  $\mathbf{x}_k^{(l)} = [x_1^{(l)}, x_2^{(l)}, \dots, x_d]$ .

Our approach first addresses covariate outliers  $\mathbf{x}_k^{(l)}$  by introducing projection pursuit weights  $w(\mathbf{x}_k)$ . These weights are then applied to scale the residuals  $r$ , ensuring that the influence of an outlier is adjusted based on the presence of extreme covariate outliers  $\mathbf{x}_k^{(l)}$ . This method enables the model to transform contaminated data points  $\{y_i^{(l)}, y_j^{(c)}, \mathbf{x}_k^{(l)}\}$  into a more reliable dataset. **Notably, the projection pursuit weighting operates independently and can be applied to various likelihoods, as shown in our experiments.**

To further handle extreme outliers in output response  $y_i^{(l)}$ , we employ a Huber density function—derived from the exponential of the Huber loss—giving robust  $L_1$  norm treatment for the residuals having over-limit magnitude. **The combination of projection pursuit weighting and Huber likelihood can handle cases where the locations  $i, j$ , and  $k$  coincide. Additionally, when extreme outliers are detected in the covariate dimensions  $\mathbf{x}_k^{(l)}$ , the model selectively retains the corresponding output  $y_k^{(l)}$  if it improves the regression fit.**

## 4 THE MODEL

Let us consider a regression setting  $y_i = f(\mathbf{x}_i) + \epsilon_i$ , where  $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$  is a homoscedastic i.i.d. random variable with constant variance. In GP models, the systematic dependency between the covariates  $\mathbf{x} \in \mathcal{X}$ , where  $\mathcal{X} \subseteq \mathbb{R}^d$ , and the response  $y \in \mathcal{Y}$  is given by a latent function,  $f(\mathbf{x}) : \mathbb{R}^d \rightarrow \mathbb{R}$ . In a truly non-parametric sense, the latent vector function at  $n$  covariates,  $\mathbf{f} = [f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)]^\top$ , is assumed to have a priori probability distribution. This distribution is a joint multivariate normal distribution with zero mean vector and covariance matrix,  $\mathbf{K}$ , that is,

$$\mathbf{f}|\mathbf{X}, \boldsymbol{\theta} \sim \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}). \quad (1)$$

The covariance matrix,  $\mathbf{K}$ , is a positive semi-definite matrix that captures residual spatial association with elements  $K_{i,j} = k(\mathbf{x}_i, \mathbf{x}_j)$ ,  $i, j = 1, \dots, n$ . The function  $k(\cdot, \cdot)$ , chosen from a parametric kernel family such as the Gaussian or the Matérn kernel, is characterized by hyperparameters denoted by  $\boldsymbol{\theta}$ . The likelihood of the data is expressed as  $\mathbf{y}|\mathbf{f}, \sigma \sim \mathcal{N}(\mathbf{y}|\mathbf{f}, \boldsymbol{\Sigma})$ , and the resulting posterior distribution on  $\mathbf{f}$  as where  $\boldsymbol{\Sigma} = \text{diag}(\sigma_1^2, \dots, \sigma_n^2)$ .

Next, we develop three aspects of the proposed GP-Huber model: Huber likelihood, projection pursuit weights, and the resulting unimodal posterior distribution. Following that, we discuss the hyperparametric settings of the GP-Huber.

### 4.1 HUBER LIKELIHOOD

We propose to use the Huber density function based on the Huber loss proposed by Huber (1992) to model the likelihood of the observed data. The Huber loss function  $\rho(\cdot)$  is a truncated mixture of two commonly used loss functions: squared loss,  $l(r) = r^2$  for residuals below threshold  $b$ , and absolute loss,  $l(r) = |r|$  for residuals  $r_i = y_i - f(\mathbf{x}_i)$  below threshold  $b$ , given by

$$\rho(r) = \begin{cases} \frac{1}{2}r^2, & \text{if } |r| \leq b \\ b|r| - \frac{1}{2}b^2, & \text{otherwise} \end{cases} \quad (2)$$

Huber (1992) considered the contamination model  $(1 - \varepsilon)G(r) + \varepsilon H(r)$ , where  $G(r)$  is the Gaussian cumulative density function and  $H(r)$  is the unknown cumulative density function. The associated least favorable Huber density function with a fraction of contamination  $\varepsilon$  is defined as

$$p_H(\mathbf{y}|\mathbf{f}, \phi) = \prod_{i=1}^n \frac{1 - \varepsilon}{\sqrt{2\pi}\sigma} \exp(-\rho(r_i)). \quad (3)$$

The parameter  $\varepsilon$ , symbolizing the fraction of the dataset presumed to deviate from the underlying model, can be computed utilizing the minimum covariance determinant estimator Hubert & Debruyne

(2010). The threshold  $b$  is selected to protect estimation of the model parameters and hyperparameters against the fraction of contamination  $\varepsilon$ . The Huber likelihood provides a balance between sensitivity to inliers and robustness to outliers, controlled by the threshold  $b$ , which has a theoretical interpretation and can be set based on domain knowledge or easily tuned empirically. Student's-t likelihood, while also robust to outliers, may give undue influence to extreme observations because of its heavy tails. The Laplace likelihood's uniform linear loss may underweight small residuals—potentially leading to less efficient estimates when the data contains mostly inliers.

## 4.2 PROJECTION PURSUIT WEIGHTING

The idea is to scale the residual  $r_i$  associated with the  $i^{\text{th}}$  data point with projection pursuit weight  $w(\mathbf{x}_i)$  based on robust variant of Mahalanobis distances, called projection statistics  $\text{PS}(\mathbf{x}_i) : \mathbb{R}^d \rightarrow \mathbb{R}^d$ . This scaling highlights the impact of outliers in single or multiple dimensions masking each other in the covariate space. Residual larger than the threshold  $b$  gets robust  $L1$  norm treatment, while those smaller than  $b$  are treated with an efficient  $L2$  norm within the Huber loss  $\rho(r)$ .

We obtain standardized the residual  $r_{S_i} = r_i/(w_i \sigma s)$  by scaling  $r_i$  by its corresponding projection pursuit weight  $w_i$  and using a scaling factor  $s = b_d \text{ med}|\mathbf{r}|$ , where  $b_d = 1 + 5/(n - d)$  is the dimensionality correction factor. When the error distribution is unknown,  $s$  accounts for its spread parameter. The projection pursuit weights  $w$  limit the influence of outliers simultaneously arising in multiple covariate dimensions at multiple locations on the loss function, are based on projection statistics  $\text{PS}_i$ , calculated as

$$w_i = \begin{cases} 1, & \text{for } \text{PS}_i^2 \leq c_i, \\ \frac{c_i}{\text{PS}_i^2}, & \text{for } \text{PS}_i^2 > c_i. \end{cases} \quad (4)$$

The projection statistics (Stahel, 1981; Donoho, 1982) are a robust version of Mahalanobis distances based on the median absolute distance from the median. Formally defined as the maxima of the standardized projection distances obtained by projecting the point cloud in the directions that originate from the co-ordinate wise median and that pass through each of the data points,  $\mathbf{x}_i$  (Mili et al., 1996). They're easy to calculate:

$$\text{PS}_i = \max_{\|\mathbf{u}_j\|=1} \frac{|\mathbf{x}_i^T \mathbf{u}_j - \text{median}_k(\mathbf{x}_k^T \mathbf{u}_j)|}{1.4826 \text{ median}_i |\mathbf{x}_i^T \mathbf{u}_j - \text{median}_k(\mathbf{x}_k^T \mathbf{u}_j)|}, \quad (5)$$

where  $\mathbf{u}_j = \frac{\mathbf{x}_j - \mathbf{M}}{\|\mathbf{x}_j - \mathbf{M}\|}$ ;  $j, k = 1, \dots, n$ . The co-ordinate wise median  $\mathbf{M}$  is given by  $\mathbf{M} = \{\text{med}_{j=1, \dots, n} \mathbf{x}_{j1}, \dots, \text{med}_{j=1, \dots, n} \mathbf{x}_{jd}\}$ . The projection statistics attain the maximum breakdown point given by  $[(n - d - 1)/2]/n$  (Maronna & Yohai, 1995).

Stahel et al. (1991) and Mili et al. (1996) showed that, when  $n > 5d$ , the squared projection statistics  $\text{PS}_i^2$  roughly follow a  $\chi^2$  distribution with a degree of freedom equal to the number of non-zero elements  $\nu_i$  in the row vector of the associated regressor,  $\mathbf{x}_i$ , i.e.,  $\text{PS}_i^2 \sim \chi_{\nu_i}^2$ . However, when  $n \leq 5d$ , it is the PS that roughly follow a  $\chi^2$  distribution, that is,  $\text{PS}_i \sim \chi_{\nu_i}^2$ . Consequently, the threshold  $c_i$  is chosen as the 97.5 percentile of the chi-square distribution with  $\nu_i$  degrees of freedom while defining weights in equation 4.

Throughout the inference process (as detailed in Section 5), we use standardized residuals  $r_{S_i}$  within the Huber likelihood.

$$p_H(\mathbf{y}|\mathbf{f}, \phi) = \prod_{i=1}^n \frac{1 - \varepsilon}{\sqrt{2\pi}\sigma} \exp(-\rho(r_{S_i})). \quad (6)$$

## 4.3 GP-HUBER POSTERIOR

The posterior distribution resulting from our model, which incorporates a non-conjugate prior, is given as:

$$p(\mathbf{f}|\mathcal{D}, \boldsymbol{\theta}, \sigma) = \frac{p_G(\mathbf{f}|\mathbf{0}, \mathbf{K})}{p(\mathcal{D}|\boldsymbol{\theta}, \sigma)} p_H(\mathbf{y}|\mathbf{f}, \sigma), \quad (7)$$

where where  $p_G(\mathbf{f}|\mathbf{0}, \mathbf{K})$  is the Gaussian prior  $\mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K})$  and  $p_H(\mathbf{y}|\mathbf{f}, \sigma)$  is the likelihood modeled using the Huber density. This formulation leads to a posterior that does not have a closed-form expression due to the non-conjugate nature of the Huber likelihood.

The marginal likelihood (or evidence) of the data, which plays a crucial role in model selection and hyperparameter optimization, is expressed as:

$$p(\mathcal{D}|\sigma, \theta) = \int p_G(\mathbf{f}|\mathbf{0}, \mathbf{K}) p_H(\mathbf{y}|\mathbf{f}, \sigma) d\mathbf{f}. \quad (8)$$

**Theorem 1.** *Let  $\mathcal{D} = (\mathbf{x}_i, y_i)_{i=1}^n$  be a dataset with distinct covariates  $\mathbf{x}_i \in \mathcal{X}$  and response  $y_i \in \mathcal{Y}$ , where  $n < \infty$ . The kernel matrix  $\mathbf{K} \in \mathbb{R}^{n \times n}$  is positive definite, with elements  $K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$  defined by a continuous kernel function  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ . Assume the Huber likelihood function  $p_H(\mathbf{y}|\mathbf{f}, \sigma)$  based on strictly convex and continuous Huber loss  $\rho(r_i) : \mathbb{R} \rightarrow \mathbb{R}$ . Then the posterior distribution  $p(\mathbf{f}|\mathcal{D}, \theta, \sigma)$  is unimodal.*

The proof is presented in appendix B.1. This theorem shows that despite the non-Gaussian and potentially complex nature of the Huber likelihood, the posterior retains a single peak. This simplifies both parameter inference and hyperparameter optimization.

We can set the threshold  $b = 1.5$  to achieve high efficiency at the Gaussian distribution (see appendix B.3). This would make our model robust to 10% outliers (since fraction of contamination is  $\varepsilon = 0.1$ ). Note that, in the context of our work, "efficiency" refers to the estimator's ability to achieve low variance when the noise follows a Gaussian distribution. Specifically, a highly efficient estimator can make the best use of data that is predominantly Gaussian, leading to more accurate parameter estimation. The contamination fraction  $\varepsilon$  defines the model's tolerance to deviations from the Gaussian assumption, allowing it to handle a proportion of outlier points without being overly influenced by them. The parameter  $b$  controls the threshold for identifying outliers and thus influences the transition between  $L2$  and  $L1$  norm treatment. By setting  $b = 0.45$ , we get  $\varepsilon = 0.45$  for heavy-tailed and Gaussian error distributions, we aim to accommodate up to 45% outliers while maintaining reasonable efficiency. The only hyperparameter of the likelihood function requiring estimation is  $\phi = \sigma^2$ .

## 5 APPROXIMATE BAYESIAN INFERENCE

By retaining the optimization-friendly properties of convex problems ensured by to unimodality (see Theorem 1), our method enables the use of the Laplace approximation (Tierney & Kadane, 1986) for the posterior. To facilitate predictions  $f^*$ , we develop Gibbs sampling and Laplace's method. The key requirement for the latter is the continuity of the Huber density function, ensuring that its derivatives exist for all  $r_S$  in the interval  $(-\infty, \infty)$ . In Gibbs sampling, the joint posterior distribution  $p(\mathbf{f}, \theta, \sigma^2)$  can be simplified using the scale mixture model of the Laplace distribution for data points with residuals  $r \geq b$ : this representation expresses the likelihood of these points as a normal distribution—making the sampling process more efficient.

### 5.1 GIBBS SAMPLING

The Huber density function is a mixture of a truncated normal and a Laplace density function for an absolute standardized residual respectively lying within and outside the threshold  $b$ . This yields

$$p_H(y|f, \sigma) = \begin{cases} \frac{C_1}{\sqrt{2\pi}w_i\sigma_g s} \exp\left(-\frac{r_i^2}{2w_i^2\sigma_g^2 s^2}\right) & |r_{S_i}| \leq b, \\ \frac{C_2}{2w_i a s} \exp\left(-\frac{b|r_i|}{w_i a s}\right) & |r_{S_i}| > b, \end{cases} \quad (9)$$

where  $C_1$  and  $C_2$  are the constants respectively, defined as  $C_1 = 1 - \varepsilon$  and  $C_2 = \sqrt{\frac{\pi}{2}} \exp(b^2/2)$ . The Laplace distribution  $p_L(y_i|f(\mathbf{x}_i), a)$  with location parameter  $a$  can be represented as a scale mixture of normal distributions  $\mathcal{N}(y_i|f(\mathbf{x}_i), \sigma_i^2)$  where  $\sigma_i^2$  follows an exponential distribution  $p_E(\sigma_i^2|\beta)$  Andrews & Mallows (1974) and  $i = 1, \dots, n_l$  are the indices of the points associated with the standardized residuals larger than the threshold  $b$  hereafter referred to as outlying points. Formally, we have

$$p_L(y_i|f(\mathbf{x}_i), a) = \int p_G(y_i|f(\mathbf{x}_i), \sigma_i^2) p_E(\sigma_i^2|\beta) d\sigma_i^2. \quad (10)$$

Using this property, we represent the individual standard deviations corresponding to  $n_l$  outlying training points as  $\{\sigma_{l_1}, \dots, \sigma_{l_{n_l}}\}$ , which are elements of the vector  $\sigma_l$ . The variance associated with

$n_g$  inlying points is denoted as  $\sigma_g^2$ . Conclusively, the Huber probability density function takes the form

$$\mathbf{y}|\mathbf{f}, \sigma_g^2, \sigma_l^2, \beta \sim \begin{cases} \prod_{i=1}^{n_g} C_1 \mathcal{N}(y_i | f(\mathbf{x}_i), \sigma_g^2) & |r_{S_i}| \leq b, \\ \prod_{i=1}^{n_l} C_2 \mathcal{N}(y_i | f(\mathbf{x}_i), \sigma_l^2) \text{Exponential}(\sigma_l^2, \beta) & |r_{S_i}| > b, \end{cases} \quad (11)$$

where  $n_g + n_l = n$  is the total number of points in the training dataset. An alternative representation of the likelihood function is given by

$$\mathbf{y}_g, \mathbf{y}_l | \mathbf{f}_g, \mathbf{f}_l, \sigma_g^2, \sigma_l^2 \sim \mathcal{N} \left( \begin{bmatrix} \mathbf{y}_g | \mathbf{f}_g \\ \mathbf{y}_l | \mathbf{f}_l \end{bmatrix}, \begin{bmatrix} \Sigma_{gg} & \mathbf{0} \\ \mathbf{0} & \Sigma_{ll} \end{bmatrix} \right), \quad (12)$$

where  $\Sigma_{gg}$  and  $\Sigma_{ll}$  both are diagonal matrices, the former with constant diagonal elements equal to  $\sigma_g^2$  and the latter with diagonal entries  $\{\sigma_{l_1}^2, \dots, \sigma_{l_{n_l}}^2\}$ . Let the hyperparameter vector  $\sigma^2$  consist of the diagonal entries of the matrix  $\Sigma_{gg}$ , which are  $\sigma_g^2$  and  $\sigma_l^2$ . The joint posterior probability density function of  $\mathbf{f}$ ,  $\sigma^2$ , and  $\theta$  is given by

$$p(\mathbf{f}, \sigma^2, \theta) \propto p(\mathbf{y} | \mathbf{f}, \sigma^2) p_G(\mathbf{f} | \mathbf{0}, \mathbf{K}) p(\sigma^2 | \beta) p(\beta | \zeta) p(\theta | \zeta). \quad (13)$$

We assume that the hyper-hyperparameter vector  $\beta$  and the hyperparameter vector  $\theta$  follow the log-uniform distribution with parameters contained in  $\zeta$ . Since the distribution of the variance parameter  $\sigma_g^2$  of  $n_g$  inlying training points is degenerate, the hyper-hyperparameter vector  $\beta = [\beta_g, \beta_l]^T$  corresponding to the  $n_g$  points follows a degenerate distribution as well. Therefore,  $p(\sigma_g^2 | \beta_g)$  is a Dirac impulse while  $\sigma_l^2 | \beta_l \sim \text{Exponential}(\sigma_l^2 | \beta_l)$ . The samples generated from this distribution are highly correlated. Therefore, in order to better mix the Monte Carlo chains, we follow the trick used by Kuss (2006) as follows:

$$p(\sigma^2, \beta, \theta) \propto \left[ \int p_G(\mathbf{y} | \mathbf{f}, \Sigma) p_G(\mathbf{f} | \mathbf{0}, \mathbf{K}) d\mathbf{f} \right] p(\sigma^2 | \beta) p(\beta | \zeta) p(\theta | \zeta), \quad (14)$$

where the covariance matrix of the  $n_g$  inlying samples and the  $n_l$  outlying samples is given by  $\Sigma = \begin{bmatrix} \Sigma_{gg} & \mathbf{0} \\ \mathbf{0} & \Sigma_{ll} \end{bmatrix}$ . The samples can be used to obtain the approximated probability density functions of the latent vector function,  $p(\mathbf{f}^* | \mathcal{D}, \mathbf{X}^*)$ , at the new test covariates contained in  $\mathbf{X}^*$  by averaging over all unknowns. Formally, we have

$$p(\mathbf{f}^* | \mathcal{D}, \mathbf{X}^*) = \int p(\mathbf{f}^* | \mathbf{f}, \sigma^2, \theta, \mathbf{X}^*, \mathcal{D}) p(\mathbf{f}, \sigma^2, \theta | \mathcal{D}) d\mathbf{f} d\sigma^2 d\theta. \quad (15)$$

For  $T$  samples, it can be evaluated as

$$p(\mathbf{f}^* | \mathcal{D}, \mathbf{X}^*, \zeta) = \frac{1}{T} \sum_{t=1}^T \int p(\mathbf{f}^* | \mathbf{f}, \mathbf{X}, \mathbf{X}^*, \theta_t) p(\mathbf{f} | \mathcal{D}, \sigma_t^2, \theta_t) d\mathbf{f}. \quad (16)$$

## 5.2 LAPLACE APPROXIMATION

To ensure the continuity of the derivative of the Huber density function with respect to the latent vector function  $\mathbf{f}$ , we utilize the pseudo-Huber loss function Charbonnier et al. (1997), which is defined as

$$\rho(r_S) = b^2 \left( \sqrt{1 + \left( \frac{r_S}{b} \right)^2} - 1 \right). \quad (17)$$

Laplace approximation of the posterior requires the likelihood to be log-concave in order for it to be represented by a unimodal multivariate normal distribution. It is executed by approximating the posterior distribution of  $\mathbf{f}$  with a normal distribution Rue et al. (2009), that is,

$$\mathbf{f} | \mathcal{D}, \sigma, \theta \sim \mathcal{N}(\hat{\mathbf{f}} | \mathbf{f}, \mathbf{A}). \quad (18)$$

The remainder of the method is detailed in appendix A.

Finally, we present the following theorem, which guarantees the robustness of GP-Huber to outliers.

**Theorem 2.** *Under the same assumptions as Theorem 1, the influence of an individual observation  $y$  on the posterior mean  $\mathbb{E}[f | y]$  is bounded:*

$$\left| \frac{\partial}{\partial y} \mathbb{E}[f | y] \right| \leq \frac{b}{\sigma}.$$

Proof is provided in appendix B.2



	SCtMCMC	tLA	HuberMCMC <sup>+pw</sup>	HuberLA <sup>+pw</sup>	RCGP	GP	LaplaceMCMC
	$\varepsilon \sim \mathcal{N}(0.01, 0.08)$						
RMSE	0.74 (0.52)	0.75 (1.31)	0.37 (0.42)	<b>0.25 (0.25)</b>	1.84 (0.82)	1.44 (0.90)	0.43 (0.46)
MAE	0.47 (0.25)	0.48 (0.61)	0.31 (0.25)	<b>0.14 (0.14)</b>	1.28 (0.54)	1.24 (0.68)	0.33 (0.26)
	$\varepsilon \sim \text{Student-}t(10)$						
RMSE	4.86 (11.56)	1.22 (1.31)	<b>0.50 (0.81)</b>	1.17 ( <b>0.37</b> )	1.89 (0.88)	1.52 (0.98)	0.59 (0.93)
MAE	1.67 (1.25)	0.77 (0.65)	<b>0.41 (0.39)</b>	0.79 ( <b>0.18</b> )	1.71 (0.85)	1.34 (0.22)	0.43 (0.35)
	$\varepsilon \sim \text{Laplace}(0, 0.1)$						
RMSE	4.76 (0.48)	1.23 (1.31)	<b>0.58 (0.42)</b>	1.17 ( <b>0.35</b> )	1.95 (0.86)	1.51 (0.89)	1.06 (0.82)
MAE	1.64 (0.23)	0.76 (0.61)	<b>0.41 (0.24)</b>	0.68 ( <b>0.18</b> )	1.27 (0.46)	1.23 (0.41)	0.75 (0.34)
	$\varepsilon \sim \text{Student-}t(1) \text{ (Cauchy)}$						
RMSE	4.75 (0.57)	1.25 (1.32)	<b>0.61 (0.49)</b>	1.20 ( <b>0.17</b> )	1.97 (0.62)	1.50 (0.89)	0.42 (0.75)
MAE	1.65 (0.27)	0.78 (0.67)	<b>0.47 (0.27)</b>	0.81 ( <b>0.11</b> )	1.78 (0.42)	1.32 (0.66)	0.66 (0.38)

Table 1: RMSE and MAE values on the Neal dataset for the Case 1. Values in parentheses represent the performance for Case 3. Bold values highlight the best performance with the lowest RMSE and MAE.

## 6 EXPERIMENTS

Through our experiments, we aim to address the following questions:

(Q1) When is HuberLA<sup>+pw</sup> (GP-Huber with Laplace’s method with pursuit weights) preferable, and under which outlier scenarios is HuberMCMC<sup>+pw</sup> (GP-Huber with Gibbs sampling with pursuit weights) more suitable?

(Q2) Does GP-Huber show a significant performance improvement over standard GP regression and the RCGP method proposed by Altamirano et al. (2024) under their experimental settings?

(Q3) Does projection pursuit weighting give GP-Huber an edge over baselines with the same weighting?

(Q4) Does GP-Huber provide more accurate estimates of the planet-to-star radius ratio compared to the standard GP method used by Gibson et al. (2012) in the transmission spectroscopy experiment?

We conducted experiments on benchmark datasets with extreme outliers in location, magnitude, and error distribution. The threshold  $b$  was 1.5 for Gaussian errors and 0.45 for Student’s-t, Laplace, and Cauchy errors. An anisotropic squared exponential kernel was used, with a zero mean function except in the spectroscopy experiment. Performance was measured using RMSE and MAE. .

### 6.1 NEAL DATASET

We evaluate the proposed GP-Huber on the Neal dataset (Neal, 1997) for the following cases of extreme outliers:

**Case 1:** Extreme outliers  $y_i^{(l)}, x_k^{(l)}$  in added in output and covariate dimensions, respectively.

**Case 2:** Only output dimensions  $y_i^{(l)}$  were contaminated with extreme data points.

**Case 3:** Bad data points  $y_j^{(c)}, x_k^{(l)}$  in added to both output and covariate dimensions, respectively, with the former being relatively close to the main data cluster compared to Case 1.

**Case 4:** Only output dimensions were contaminated with data points  $y_j^{(c)}$  relatively close to the data cloud compared to Case 1.

In all the cases above, the locations  $i, j$  and  $k$  may differ or coincide (refer to appendix C.1 for the location and magnitude details on outliers). For each case, we considered four different error distributions:  $\mathcal{N}(0.01, 0.08)$ , Student- $t(10)$ , Laplace(0, 0.1), Student’s- $t(1)$ .

The baseline models considered for comparison on the Neal dataset, along with RCGP, include: GP with a Student’s t error model solved using MCMC integration (SCtMCMC), GP with a Student’s t error model using Laplace approximation (tLA), and GP with a Laplace likelihood solved via MCMC integration (LaplaceMCMC). Table 1 presents the RMSE and MAE values comparing GP-Huber against these baselines for Cases 1 and 3. Refer to appendix C.1 for the Tables 5, 6 for the Cases 2, 4 and appendix B.4 for the implementation details of the baselines. Furthermore, pursuit weighting is incorporated into all baseline models, and their performances are compared in Tables 7 to 10 (provided in appendix C.2). Now, we are in position to answer Q1.

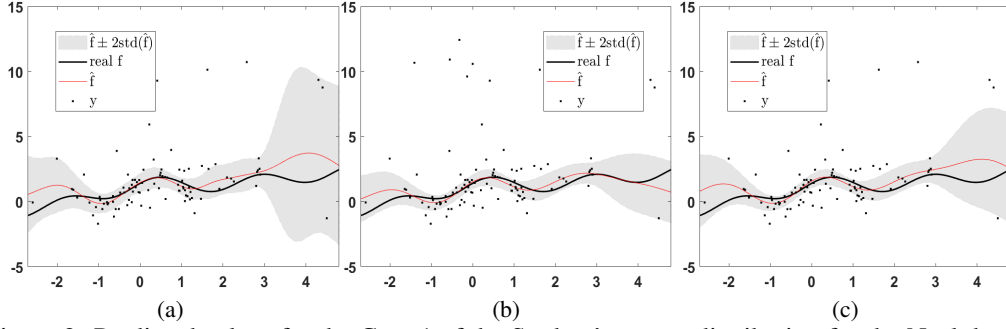


Figure 2: Predicted values for the Case 1 of the Student’s t-error distribution for the Neal dataset obtained from the eight considered GP regression models: (a) HuberMCMC ; (b) HuberLA; (c) RCGP.

### When is HuberMCMC better?

In scenarios with  $y^{(l)}$ ,  $x^{(l)}$  (Case 1), HuberMCMC performed better than HuberLA (see, Tables 1 and 5). HuberMCMC also outperformed tLA in predictive accuracy, demonstrating a more robust fit that is less influenced by  $x^{(l)}$  (Figure 2). HuberLA generally provided better uncertainty quantification compared to HuberMCMC (see Figures 2 and 6), while maintaining competitive predictive performance. In outlier scenarios with  $y^{(l)}$  (Case 2), HuberMCMC exhibited superior performance across Student’s-t, Laplace, and Cauchy error distributions (see, Table 5). This suggests that HuberMCMC is a robust choice for datasets containing extreme output outliers i.e. outlier scenarios similar to Cases 1 and 2.

### When is HuberLA better?

HuberLA exhibited superior performance in handling closer output outliers  $y^{(c)}$  compared to HuberMCMC (values in parenthesis in the Table 1 and Table 6). Figure 7 highlights HuberLA’s robustness to  $x^{(l)}$ , in contrast to tLA which is influenced by such points. While HuberLA generally provided more accurate predictions and reliable uncertainty quantification than both HuberMCMC and tLA, HuberMCMC performed competitively for the Cases 3 and 4.

From Tables 7 and 8 in appendix C.2 (Cases 1 and 3), where projection pursuit weights were added to other baselines, we observe that HuberLA and HuberMCMC benefit the most from these weights. While Student’s-t likelihood also scales residuals by pursuit weights, its logarithmic penalty  $\propto \log(1 + r^2/\nu)$  (with  $\nu$  controlling tail heaviness) is less sensitive to large residuals than Huber likelihood’s linearized penalty  $\propto |r|$ . Laplace likelihood similarly penalizes  $|r|$  but lacks a quadratic center, making Huber likelihood the optimal balance of robustness and efficiency. Tables 9 and 10 show the results for outlier Cases 2 and 4, where projection pursuit weighting is added to other baselines. The weights equal 1 for all data points due to the absence of covariate outliers  $x^{(l)}$ . GP-Huber performs comparably to other baselines in these cases. This demonstrates that the weighting mechanism enhances GP-Huber’s accuracy, addressing Q3.

## 6.2 UCI DATASETS

In this set of experiments, we compared the performance of GP-Huber on the UCI datasets, Energy and Yacht, against RCGP and other baselines: t-GP, m-GP, and standard GP, using the outlier settings from Altamirano et al. (2024). We specifically focused on the "focused outlier" and "asymmetrical outlier" scenarios, as they closely resemble our extreme and close outlier cases.

MAE values of the comparison are presented in Table 2. As expected, HuberLA demonstrates to be more robust than HuberMCMC since the asymmetrical and focused outliers cases considered in the study of Altamirano et al. (2024) broadly fall under the Cases 3 and 4 in our study. On the Energy dataset, HuberLA outperformed both tLA and RCGP.

In our experiments, HuberLA outperformed RCGP and other baselines significantly in asymmetric outlier case and also showing the good computational efficiency, thus answering Q2.

*Remark.* Note that, as the outliers are present only in the response and not in the covariate dimensions, the projection pursuit weights are equal to 1 for these datasets.



	GP	RCGP	t-GP	m-GP	HuberMCMC	HuberLA
	Focused Outliers					
Energy	0.03 (0.04)	<b>0.02</b> (0.00)	0.03 (0.05)	0.24 (0.00)	0.12 (0.01)	0.04 (0.01)
Yacht	0.26 (0.15)	<b>0.10</b> (0.14)	0.20 (0.04)	0.24 (0.00)	0.24 (0.02)	0.28 (0.00)
	Asymmetric Outliers					
Energy	0.54 (0.02)	0.44 (0.04)	0.42 (0.02)	0.41 (0.00)	0.47 (0.02)	<b>0.11 (0.00)</b>
Yacht	0.54 (0.06)	0.35 (0.02)	0.41 (0.00)	0.40 (0.00)	0.51 (0.01)	<b>0.12 (0.00)</b>

Table 2: MAE values for energy and yacht. Bold values indicate the best performance for each row.

Computational costs for the experiments on Neal and UCI datasets are presented in Table 12 and 13 in appendix D. HuberLA—similar to RCGP and tLA—requires less computational time than HuberMCMC, as expected. The models converge faster for unidimensional data: HuberMCMC performs comparably to MCMC techniques with Student-t likelihood. For multidimensional cases, HuberMCMC, as expected for sampling-based methods, requires more time to converge, while HuberLA achieves faster convergence (between 5 to 10 s).

### 6.3 TRANSMISSION SPECTROSCOPY

Transmission spectroscopy records the relative change in the stellar flux, which is the incident photons per unit area, as a planet travels in front of the star. The sources of error, such as photon noise and instrumental and astrophysical systematics, raise many potential challenges for precise planet’s atmosphere characterization. The goal is to infer the planet to star radius ratio  $\rho_{radius}$  from the observed flux as the planet passes in front of the star. The optical state parameters are metered via auxiliary measurements of the spectral trace such as position, width, angle, or other parameters, indicating the state of the detector and optics, which are thought to be the cause of instrumental systematics. Instead of modeling the latter as a linear function of the optical state parameters, Gibson et al. (2012) proposed a non-parametric model by leveraging GPs.

The observation set obtained from HST-NICMOS includes the light curves for 18 wavelength channels extracted from  $n = 638$  spectra of the planetary system HD-189733. The flux measurements contained in the vector,  $\mathbf{f} = [f_1, f_2, \dots, f_n]^T$ , are recorded at  $n$  time instants,  $\{t_1, t_2, \dots, t_n\}$  and the optical state parameters  $\mathbf{x}_{t_i}$  collected in the matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$  constitute the training dataset. We extend the work of Gibson et al. (2012) by using the GP-Huber model to estimate the planet-to-star radius ratio  $\rho_{radius}$ . As demonstrated earlier, the robustness to outliers of GP-Huber allows us to utilize 517 measurements associated with four out-of-transit orbits, namely orbit numbers,  $\{2, 3, 4, 5\}$ , and 137 measurements associated with one in-transit orbit, namely orbit number 1. The latter was excluded from the analysis performed by Gibson et al. (2012) as it constitutes much larger systematics effects attributed to the spacecraft settling. The observed transit flux modeled in the GP framework follows a normal distribution, that is,

$$\mathbf{f}(t, \mathbf{X}) \sim \mathcal{N}(\mathbf{T}(t, \phi), \mathbf{K}), \quad (19)$$

where the parameter vector,  $\phi$ , include the parameter of interest,  $\rho_{radius}$ , and other parameters. We consider the analytical quadratic limb darkening transit function proposed by Mandel & Agol (2002). Analogous with equation 11, we assume that the observed transit flux vector,  $\mathbf{f} = \mathbf{f}(t, \mathbf{X})$ , in the GP-Huber framework follows a normal distribution, that is,

$$\mathbf{f} | \mathbf{T}(t, \phi), \mathbf{X}, \phi, \theta, \sigma^2 \sim \mathcal{N}(\mathbf{T}(t, \mathbf{X}), \Sigma + \mathbf{K}). \quad (20)$$

The joint un-normalized log-posterior function of  $\phi$ ,  $\beta$ , and  $\theta$  with the gamma aprior probability density function,  $p(\theta) = \frac{1}{t} \exp(-\frac{\theta}{t})$ , over the covariance function hyperparameters is given by

$$\log P(\phi, \theta, \sigma^2, \beta | \mathbf{f}, \mathbf{X}, \zeta) = \log(\mathcal{L}(\mathbf{r}_S | \mathbf{X}, \phi, \theta, \sigma^2)) - \frac{\tau}{t_\tau} - \sum_{i=1}^d \left( \frac{1}{s_i t_i} \right) + \log(\beta) - \beta^T \sigma^2 + \log(p(\beta | \zeta)) + C. \quad (21)$$

The challenging task now is to infer the parameter  $\rho_{radius}$  from the joint posterior distribution of  $(\phi, \theta, \sigma^2, \beta)$ . The log-likelihood  $\mathcal{L}$  term is expressed as

$$\log \mathcal{L}(\mathbf{r}_S | \mathbf{X}, \phi, \theta, \sigma^2) = -\frac{1}{2} \mathbf{r}_S^T (\Sigma + \mathbf{K})^{-1} \mathbf{r}_S - \frac{1}{2} \log |\Sigma + \mathbf{K}| - \frac{n}{2} \log(2\pi) + \log(1 - \varepsilon), \quad (22)$$

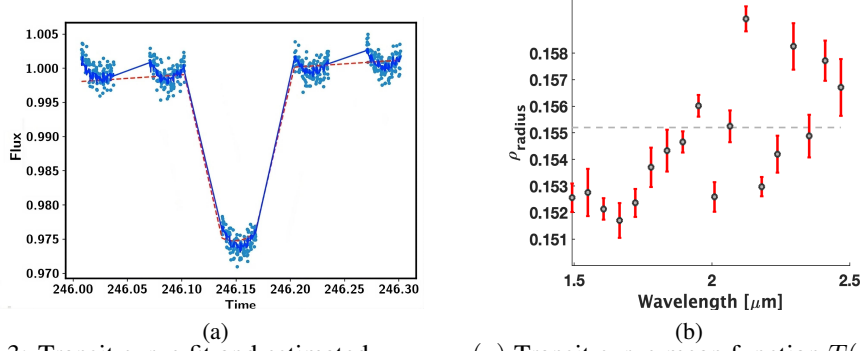


Figure 3: Transit curve fit and estimated  $\rho_{radius}$ . (a) Transit curve mean function  $T(t, \theta)$  (in red dotted line) and GP-Huber model fit (in blue solid line); (b) results of planet-to-star radius ratios ( $\rho_{radius}$ ) obtained from GP-Huber with error-bars. The dashed grey line represents the sanity check values.

where  $\mathbf{r} = \mathbf{f} - \mathbf{T}(t, \mathbf{X})$ . One of the approaches is to use the Bayesian method that seeks the posterior distribution of  $\rho_{radius}$  by marginalizing over the other parameters of the mean function parameters  $\phi$  and the covariance function hyperparameters,  $\theta$  using MCMC methods. The other method proposed as the type-II maximum likelihood method by Gibson et al. (2012), where the hyperparameters,  $\theta$  and  $\sigma^2$ . Formally, we have

$$(\hat{\phi}, \hat{\theta}, \hat{\sigma}^2, \hat{\beta}) = \arg \max_{\phi, \theta, \sigma^2, \beta} \log P(\phi, \theta, \sigma^2, \beta | \mathbf{f}, \mathbf{X}, \zeta). \quad (23)$$

And the posterior distribution of the parameter of interest  $\rho_{radius}$  is obtained by marginalizing the joint posterior distribution  $p(\phi, \theta, \sigma^2, \beta)$  over the hyperparameters and the rest of the mean function parameters. In the standard type II maximum likelihood method, the hyperparameters are fixed to their maximum likelihood estimates i.e. by maximizing the evidence  $p(\mathcal{D} | \phi, \theta, \sigma^2)$ .

Figure 3(a) shows the transit fit obtained for one wavelength channel. Figure 3(b) shows the estimated  $\rho_{radius}$  obtained using MCMC integration over the rest of the mean function parameters  $\phi$  and hyperparameters  $\theta$  along with the values estimated from the white light curve represented as the white dashed line. Note that the estimated  $\rho_{radius}$  values are very close to the white light curve value of 0.155. Most of our results agree with the results obtained from the Gibson model except for wavelength channels  $1.665\mu\text{m}$  and  $2.124\mu\text{m}$  (see, appendix C.3), which effectively answers Q4.

Our code<sup>1</sup> was implemented in Matlab R2023a with the help of package gpstuff on Intel i7 processor.

## 7 CONCLUSIONS

The proposed GP-Huber model shows promise for handling a variety of heavy-tailed and Gaussian error distributions with extreme outliers in both covariate and output dimensions. **Notably, it introduces additional parameters,  $b$  and  $\varepsilon$ , which can be heuristically set prior to parameter inference.** The model’s unimodal posterior simplifies Gibbs sampling and allows for an efficient Laplace approximation. **We prove the bounded influence of observations on the posterior mean.** From our experiments on the Neal and UCI datasets, we observe that HuberMCMC and HuberLA offer superior robustness compared to RCGP and other baselines. Additionally, the transmission spectroscopy experiment demonstrates their potential in real-world applications.

In future work, we will examine GP-Huber’s performance with skewed error distributions and investigate the use of high breakdown estimators for highly corrupted real-world datasets. Another direction for future work involves extending the scalability of GP-Huber to handle large datasets by implementing sparse inference techniques.

<sup>1</sup><https://anonymous.4open.science/r/GpHuber-6A2D>

## REFERENCES

- Matias Altamirano, Francois-Xavier Briol, and Jeremias Knoblauch. Robust and conjugate gaussian process regression. In *Forty-first International Conference on Machine Learning*, 2024.
- Daniel Andrade and Akiko Takeda. Robust gaussian process regression with the trimmed marginal likelihood. In *Uncertainty in Artificial Intelligence*, pp. 67–76. PMLR, 2023.
- David F Andrews and Colin L Mallows. Scale mixtures of normal distributions. *Journal of the Royal Statistical Society: Series B (Methodological)*, 36(1):99–102, 1974.
- Pierre Charbonnier, Laure Blanc-Féraud, Gilles Aubert, and Michel Barlaud. Deterministic edge-preserving regularization in computed imaging. *IEEE Transactions on image processing*, 6(2): 298–311, 1997.
- Atefeh Daemi, Yousef Alipouri, and Biao Huang. Identification of robust gaussian process regression with noisy input using em algorithm. *Chemometrics and Intelligent Laboratory Systems*, 191:1–11, 2019.
- Bruno De Finetti. The bayesian approach to the rejection of outliers. In *Proceedings of the fourth Berkeley Symposium on Probability and Statistics*, volume 1, pp. 199–210. University Press Berkeley, Calif., 1961.
- Alain Desgagné and Philippe Gagnon. Bayesian robustness to outliers in linear regression and ratio estimation. *Brazilian Journal of Probability and Statistics*, 33(2):205–221, 2019.
- David L Donoho. Breakdown properties of multivariate location estimators. Technical report, Technical report, Harvard University, Boston. URL <http://www-stat.stanford.edu/~donoho/>, 1982.
- Zoubin Ghahramani and Matthew Beal. Propagation algorithms for variational bayesian learning. *Advances in neural information processing systems*, 13, 2000.
- NP Gibson, Suzanne Aigrain, S Roberts, TM Evans, Michael Osborne, and F Pont. A gaussian process framework for modelling instrumental systematics: application to transmission spectroscopy. *Monthly notices of the royal astronomical society*, 419(3):2683–2694, 2012.
- Paul Goldberg, Christopher Williams, and Christopher Bishop. Regression with input-dependent noise: A gaussian process treatment. *Advances in neural information processing systems*, 10, 1997.
- Peter J Huber. Robust estimation of a location parameter. In *Breakthroughs in statistics: Methodology and distribution*, pp. 492–518. Springer, 1992.
- Mia Hubert and Michiel Debruyne. Minimum covariance determinant. *Wiley interdisciplinary reviews: Computational statistics*, 2(1):36–43, 2010.
- Pasi Jylänki, Jarno Vanhatalo, and Aki Vehtari. Robust gaussian process regression with a student-t likelihood. *Journal of Machine Learning Research*, 12(11), 2011.
- Laura Kreidberg. Exoplanet atmosphere measurements from transmission spectroscopy and other planet-star combined light observations. *arXiv preprint arXiv:1709.05941*, 2017.
- Malte Kuss. *Gaussian process models for robust regression, classification, and reinforcement learning*. PhD thesis, Technische Universität Darmstadt Darmstadt, Germany, 2006.
- Zhao-Zhou Li, Lu Li, and Zhengyi Shao. Robust gaussian process regression based on iterative trimming. *Astronomy and Computing*, 36:100483, 2021.
- Kaisey Mandel and Eric Agol. Analytic light curves for planetary transit searches. *The Astrophysical Journal*, 580(2):L171, 2002.
- Ricardo A Maronna and Victor J Yohai. The behavior of the stahel-donoho robust multivariate estimator. *Journal of the American Statistical Association*, 90(429):330–341, 1995.
- Lamine Mili, MG Cheniae, NS Vichare, and Peter J Rousseeuw. Robust state estimation based on projection statistics [of power systems]. *IEEE Transactions on Power Systems*, 11(2):1118–1127, 1996.

- Thomas P Minka. Expectation propagation for approximate bayesian inference. *arXiv preprint arXiv:1301.2294*, 2013.
- Andrew Naish-Guzman and Sean Holden. Robust regression with twinned gaussian processes. *Advances in neural information processing systems*, 20, 2007.
- Radford M Neal. Monte carlo implementation of gaussian process models for bayesian regression and classification. *arXiv preprint physics/9701026*, 1997.
- JO Ramsey and MR Novick. Plu robust bayesian decision theory: point estimation. *J. Amer. Statist. Assoc.*, 75:401–407, 1980.
- Rishik Ranjan, Biao Huang, and Alireza Fatehi. Robust gaussian process modeling using em algorithm. *Journal of Process Control*, 42:125–136, 2016.
- Håvard Rue, Sara Martino, and Nicolas Chopin. Approximate bayesian inference for latent gaussian models by using integrated nested laplace approximations. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 71(2):319–392, 2009.
- Werner Stahel, Sanford Weisberg, Peter J Rousseeuw, and Bert C van Zomeren. Robust distances: simulations and cutoff values. In *Directions in Robust Statistics and Diagnostics: Part II*, pp. 195–203. Springer, 1991.
- Werner A Stahel. *Robuste schätzungen: infinitesimale optimalität und schätzungen von kovarianzmatrizen*. PhD thesis, ETH Zurich, 1981.
- Luke Tierney and Joseph B Kadane. Accurate approximations for posterior moments and marginal densities. *Journal of the american statistical association*, 81(393):82–86, 1986.
- Jarno Vanhatalo, Pasi Jylänki, and Aki Vehtari. Gaussian process regression with student-t likelihood. *Advances in neural information processing systems*, 22, 2009.
- Mike West. Outlier models and prior distributions in bayesian linear regression. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 46(3):431–439, 1984.
- Christopher Williams. Computing with infinite networks. *Advances in neural information processing systems*, 9, 1996.

## A LAPLACE APPROXIMATION

A Taylor series expansion about the largest mode of the un-normalized posterior density function of  $\mathbf{f}$  yields  $q(\mathbf{f}|\mathcal{D}, \sigma, \boldsymbol{\theta}) \approx p_H(\mathbf{y}|\mathbf{f}, \sigma)p_G(\mathbf{f}|\mathbf{0}, \mathbf{K})$ . The latter is used to define the MAP estimate  $\hat{\mathbf{f}}$ , given by

$$\hat{\mathbf{f}} = \arg \max_{\mathbf{f}} \ln q(\mathbf{f}|\mathcal{D}, \sigma, \boldsymbol{\theta}), \quad (24)$$

which may converge to a local mode in case of multimodal likelihood. As for the posterior covariance matrix,  $\mathbf{A}$ , it is given by

$$\mathbf{A} = (\mathbf{K}^{-1} + \mathbf{W})^{-1}, \quad (25)$$

where  $\mathbf{W} = -\nabla \nabla_{\mathbf{f}} \ln (p_H(\mathbf{y}|\hat{\mathbf{f}}, \sigma))$ . The hyperparameter vector  $(\sigma, \boldsymbol{\theta})$  is estimated by maximizing the log of the approximate evidence given by equation 8 using the gradient descent or the conjugate gradient method since the gradient can be analytically derived. Formally, we have

$$(\hat{\sigma}, \hat{\boldsymbol{\theta}}) = \arg \max_{(\sigma, \boldsymbol{\theta})} \ln q(\mathcal{D}|\sigma, \boldsymbol{\theta}), \quad (26)$$

where  $q(\mathcal{D}|\sigma, \boldsymbol{\theta}) \approx p(\mathcal{D}|\sigma, \boldsymbol{\theta})$  is the approximate log evidence given by

$$\ln q(\mathcal{D}|\sigma, \boldsymbol{\theta}) = \ln p_H(\hat{\mathbf{f}}|\mathbf{f}) - \frac{1}{2} \ln |\mathbf{K}| - \frac{1}{2} \mathbf{f}^T \mathbf{K}^{-1} \mathbf{f} + \frac{1}{2} \ln |\mathbf{A}|. \quad (27)$$

## B PROOFS AND THEORETICAL EXPLANATIONS

### B.1 PROOF OF THEOREM 1

*Proof.* The GP-Huber posterior distribution is proportional to the expression:

$$p(\mathbf{f}|\mathcal{D}, \boldsymbol{\theta}, \sigma) \propto \exp \left( -\frac{1}{2} \mathbf{f}^\top \mathbf{K}^{-1} \mathbf{f} - \sum_{i=1}^n \rho(y_i - f_i) \right),$$

where  $\rho$  denotes the Huber loss function, which is continuous and strictly convex. The derivative of the log-posterior with respect to  $\mathbf{f}$  is:

$$\nabla_{\mathbf{f}} \log p(\mathbf{f}|\mathcal{D}, \boldsymbol{\theta}, \sigma) \propto -\mathbf{K}^{-1} \mathbf{f} - \nabla_{\mathbf{f}} \rho(y_i - f_i),$$

and for each component  $f_i$ , the derivative becomes:

$$h_i(f_i) = \frac{-(y_i - f_i)}{\sqrt{1 + (y_i - f_i)^2}} - v_i,$$

where  $v_i$  represents the  $i^{\text{th}}$  component of  $\mathbf{v} = \mathbf{K}^{-1} \mathbf{f}$ . The term  $\frac{-(y_i - f_i)}{\sqrt{1 + (y_i - f_i)^2}}$  is strictly monotonic in  $f_i$ , as its derivative is positive for all  $f_i$ . Its behavior at the limits is given by:

$$\frac{-(y_i - f_i)}{\sqrt{1 + (y_i - f_i)^2}} = \begin{cases} 0 & \text{if } f_i \rightarrow y_i, \\ -1 & \text{if } f_i \rightarrow \infty. \end{cases}$$

The second term,  $v_i$ , arises from the precision matrix  $\mathbf{K}^{-1}$ , which is symmetric and positive definite. By the spectral theorem,  $\mathbf{K}^{-1}$  can be diagonalized as  $\mathbf{K}^{-1} = \mathbf{Q} \boldsymbol{\Lambda} \mathbf{Q}^\top$ , where  $\boldsymbol{\Lambda}$  is the diagonal matrix of positive eigenvalues, and  $\mathbf{Q}$  is an orthogonal matrix. This ensures that  $v_i$  is a linear combination of the entries of  $\mathbf{f}$  and is therefore continuous and differentiable in  $f_i$ . Combining these terms,  $h_i(f_i)$  is strictly monotonic because the first term is monotonic and the second term is linear. By the intermediate value theorem,  $h_i(f_i)$  crosses zero exactly once because:

$$\lim_{f_i \rightarrow y_i} h_i(f_i) = 0 \quad \text{and} \quad \lim_{f_i \rightarrow \infty} h_i(f_i) = -1.$$

Thus,  $h_i(f_i)$  has a unique root, ensuring that the log-posterior  $\log p(\mathbf{f}|\mathcal{D}, \boldsymbol{\theta}, \sigma)$  is strictly concave. The strict concavity of the log-posterior implies that the posterior distribution  $p(\mathbf{f}|\mathcal{D}, \boldsymbol{\theta}, \sigma)$  has a unique mode, proving that it is unimodal.  $\square$

### B.2 PROOF OF THEOREM 2

*Proof.* The influence function  $I(r_i) = -\frac{\partial}{\partial r_i} \log p(y_i | f_i)$  of the likelihood  $p(y_i | f_i, \sigma^2)$  for inliers ( $|r_i| \leq b$ ) is:

$$I(r_i) = \frac{r_i}{\sigma^2},$$

and for outliers ( $|r_i| > b$ ) becomes:

$$I(r_i) = \frac{b \operatorname{sign}(r_i)}{\sigma}.$$

We observe that the  $I(r_i)$  is bounded i.e. the contribution of any individual observation  $y_i$  to  $\mathbb{E}[f(x) | \mathbf{y}]$  is bounded, ensuring robustness to outliers.  $\square$

### B.3 SELECTION OF THE THRESHOLD $b$ AND $\varepsilon$

The Huber estimator is a maximum likelihood estimator associated with the least favorable density function given by

$$\tilde{g}(r) = \frac{1 - \varepsilon}{\sqrt{2\pi}\sigma} e^{-\rho(\frac{r}{\sigma})}, \quad (28)$$

which can be further elaborated as

$$\tilde{g}(r) = \begin{cases} \frac{1 - \varepsilon}{\sqrt{2\pi}} e^{-\frac{r^2}{2}} & \text{for } |r| \leq b \\ \frac{1 - \varepsilon}{\sqrt{2\pi}} e^{-|b||r| - \frac{b^2}{2}} & \text{for } |r| > b \end{cases} \quad (29)$$

This distribution is Gaussian in the center and Laplacian in the tails. The threshold  $b$  is related to the fraction of contamination  $\varepsilon$  against which we want to be protected. This relation is obtained by setting

$$\int_{-\infty}^{\infty} \tilde{g}(r) dr = 1 \quad (30)$$

yielding

$$\begin{aligned} \int_{-b}^b \frac{(1-\varepsilon)}{\sqrt{2\pi}} e^{-\frac{r^2}{2}} dr + 2(1-\varepsilon) \int_b^{\infty} \frac{1}{\sqrt{2\pi}} e^{(-br+\frac{b^2}{2})} dr &= 1 \\ (1-\varepsilon) \int_{-b}^b \frac{1}{\sqrt{2\pi}} e^{-\frac{r^2}{2}} dr &= (1-\varepsilon)[1 - 2(1 - \Phi(b))] = (1-\varepsilon)(2\Phi(b) - 1); \end{aligned}$$

and

$$\begin{aligned} 2(1-\varepsilon) \int_b^{\infty} \frac{1}{\sqrt{2\pi}} e^{(-br+\frac{b^2}{2})} dr &= 2(1-\varepsilon) \frac{1}{\sqrt{2\pi}} \left[ -\frac{1}{b} e^{(-br+\frac{b^2}{2})} \right]_b^{\infty} \\ &= \frac{2(1-\varepsilon)}{\sqrt{2\pi}} \frac{1}{b} e^{-\frac{b^2}{2}} = \frac{2(1-\varepsilon)}{b} \phi(b); \end{aligned}$$

Solving further, we get

$$2\Phi(b) - 1 + \frac{2}{b}\Phi(b) = \frac{1}{1-\varepsilon} \quad (31)$$

We observe that  $b$  decreases to 0 as  $\varepsilon$  increases to 1. At  $b = 1.5$ , the Huber loss can handle roughly  $\varepsilon = 0.1$  i.e. 10% of outliers. The values are listed in the table below:

$\varepsilon$	$b$
0.10	3.16
0.20	2.24
0.30	1.83
0.40	1.58
0.50	1.41
0.60	1.29
0.70	1.20
0.80	1.12
0.90	1.05
1.00	1.00

Table 3: Values of  $\varepsilon$  and  $b$  obtained by solving Equation 36.

Figure 4 presents the RMSE and MAE values on the Neal dataset for different choices of  $b$  and  $\varepsilon$ .

#### B.4 IMPLEMENTATION DETAILS OF THE BASELINES

Model #	Description
SCTMCMC	GP with a Student’s t error model using a scale mixture representation, solved via MCMC integration of the latent vector $\mathbf{f}$ , likelihood hyperparameters $\phi = (\nu, \sigma^2)$ , and kernel hyperparameters $\theta$ .
tLA	Student’s t likelihood model with the Laplace approximation, solved via Laplace integration over $\mathbf{f}$ and MAP estimates of $\phi = (\nu, \sigma^2)$ and $\theta$ .
GP	Conjugate model: GP with normal error, where the hyperparameters include $\mathbf{f}$ , $\phi = \sigma^2$ , and $\theta$ .
LaplaceMCMC	GP with Laplace likelihood, solved via MCMC integration over $\mathbf{f}$ , $\phi = \sigma$ , and $\theta$ .

Table 4: Descriptions of the baseline models and their corresponding inference techniques.



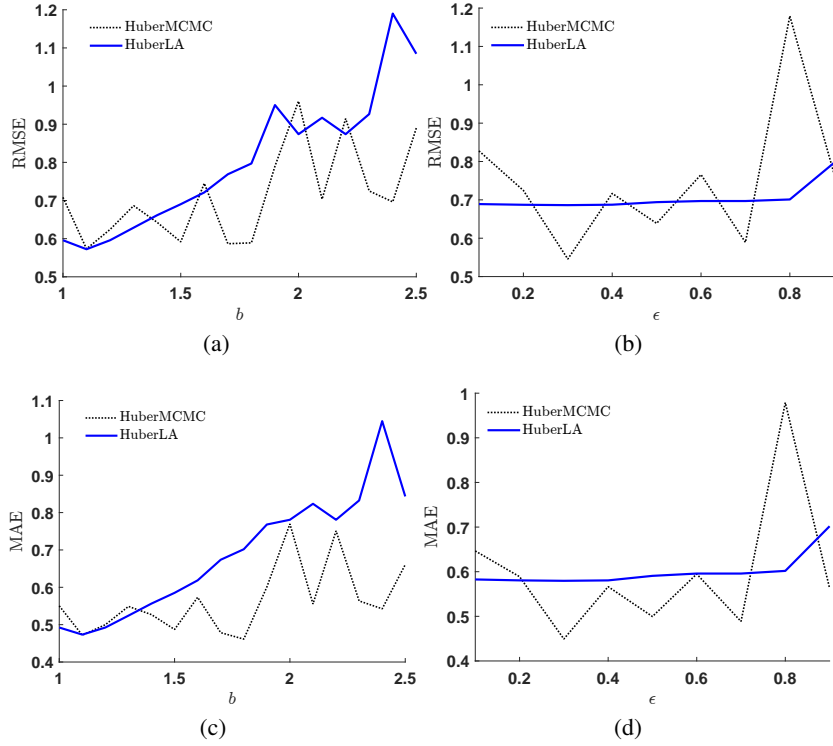


Figure 4: Impact of  $b$  and  $\epsilon$  on predictions for Neal dataset (case 1 with Student-t noise): (a) RMSE vs  $b$  ( $\epsilon = 0.45$ ); (b) RMSE vs  $\epsilon$  ( $b = 1.5$ ); (c) MAE vs  $b$ ; (d) MAE vs  $\epsilon$ .

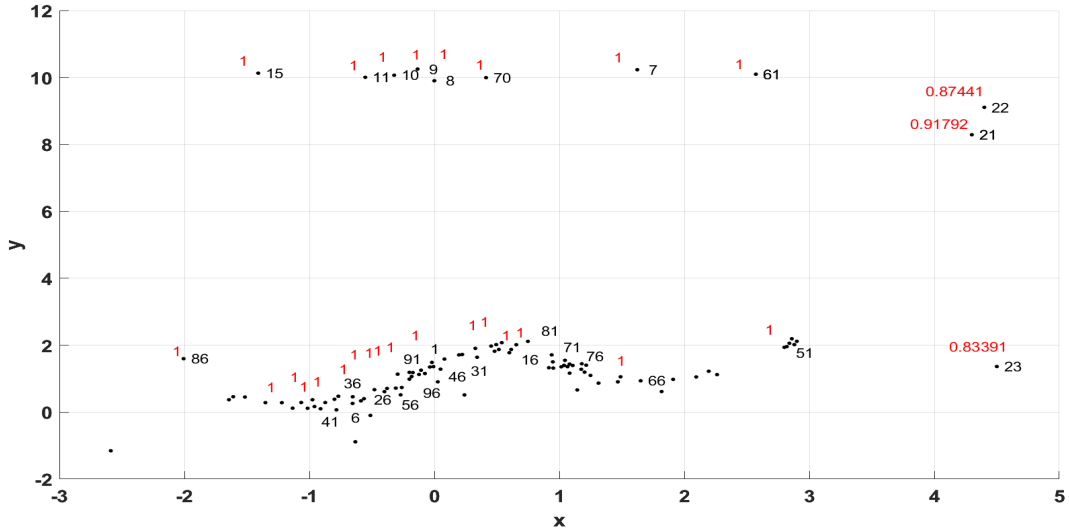


Figure 5: Weights based on PS for the Neal data. The numbers right to the data points indicate index numbers and the ones to the left in red color indicate the weights associated with that data point.

## C ADDITIONAL EXPERIMENTS

### C.1 NEAL DATASET

(Neal, 1997) proposed the following artificial model:

$$g(\mathbf{x}_i) = 0.3 + 0.4x + 0.5\sin(2.7x) + 1.1/(1 + x^2). \quad (32)$$

A sample of  $n = 100$  points constitutes the training data set,  $(\mathbf{X}, \mathbf{y})$ . The predictions of the vector function,  $\mathbf{f}^*$ , are made at  $n^* = 541$  test covariates contained in  $\mathbf{x}^*$  over the interval  $[-2.7, 5]$ . Since the projection statistics require at least a two-dimensional covariate space, they are calculated on the regressors' vector,  $\mathbf{x}$  combined with the column of ones, i.e., on the matrix  $\mathbf{H} = [\mathbf{1}, \mathbf{x}]$ . Specifically, for a test point  $\mathbf{h}_i = [1, x_i]$ ,  $\text{PS}(\mathbf{h}_i)$  is calculated using equation 4 in the paper. The training covariate,  $x_i$ , is flagged as an outlier if the associated projection pursuit weights,  $w_i = \min\left(1, \frac{c}{\text{PS}(\mathbf{h}_i)^2}\right)$ , has a value less than one.

We demonstrate the proposed GP-Huber in four cases of error probability distribution: (i)  $\mathcal{N}(0.01, 0.08)$ ; (ii) the Student's t-distribution with 10 degrees of freedom; (iii) Laplace(0, 0.1); and (iv) the Cauchy distribution. For each of these error distributions, we introduce extreme output outliers  $y^l = \{90.5, 8.6, 98.1, 5.3, 5.2, 6.1, 1, 8\}$  at locations  $j = \{7, 8, 9, 10, 11, 15, 61, 70\}$ , extreme covariate data points  $x^{(l)} = \{4.3, 4.4, 4.5\}$  at locations  $i = \{21, 22, 23\}$ . We also add large magnitudes to introduce group of good data points to the covariates  $\{x_{50}, x_{51}, x_{52}, x_{53}, x_{54}, x_{55}\}$  for which  $y_i = g(x_i)$ .

We observe that the projection pursuit weights based on the PS corresponding to the bad leverage points are  $\{0.9179, 0.8744, 0.8339\}$  while those corresponding to the good leverage points are equal to 1 (see Figure 5).

	SCtMCMC	tLA	HuberMCMC <sup>+pw</sup>	HuberLA <sup>+pw</sup>	RCGP	GP	LaplaceMCMC
$\varepsilon \sim \mathcal{N}(0.01, 0.08)$							
RMSE	1.41	1.30	1.40	1.36	2.04	1.74	1.48
MAE	0.90	0.81	0.99	0.95	1.93	1.51	0.98
$\varepsilon \sim \text{Student-t}(10)$							
RMSE	1.22	1.14	0.91	1.12	2.04	1.66	1.01
MAE	0.63	0.56	0.62	0.67	1.85	1.34	0.92
$\varepsilon \sim \text{Laplace}(0, 0.1)$							
RMSE	1.38	2.73	1.33	1.37	2.06	1.73	1.33
MAE	0.88	1.82	0.97	0.96	1.95	1.51	0.95
$\varepsilon \sim \text{Student-t}(1) \text{ (Cauchy)}$							
RMSE	4.74	2.11	1.33	1.38	2.11	1.75	1.33
MAE	1.67	1.36	0.96	0.98	1.84	1.50	0.95

Table 5: Neal results for Case 2

	SCtMCMC	tLA	HuberMCMC <sup>+pw</sup>	HuberLA <sup>+pw</sup>	RCGP	GP	LaplaceMCMC
$\varepsilon \sim \mathcal{N}(0.01, 0.08)$							
RMSE	1.02	1.01	1.48	1.50	1.10	1.17	0.98
MAE	0.51	0.52	0.79	0.54	0.76	0.78	0.53
$\varepsilon \sim \text{Student-t}(10)$							
RMSE	1.58	1.02	1.17	1.13	1.11	1.17	0.61
MAE	1.28	0.52	0.53	0.78	0.76	0.85	0.35
$\varepsilon \sim \text{Laplace}(0, 0.1)$							
RMSE	1.04	1.01	1.06	1.18	1.16	1.08	1.16
MAE	0.51	0.52	0.53	0.78	0.66	0.66	0.58
$\varepsilon \sim \text{Student-t}(1) \text{ (Cauchy)}$							
RMSE	1.58	1.02	1.18	1.02	1.10	1.07	1.04
MAE	1.28	0.52	0.63	0.78	0.56	0.56	0.52

Table 6: Neal results for the Case 4.

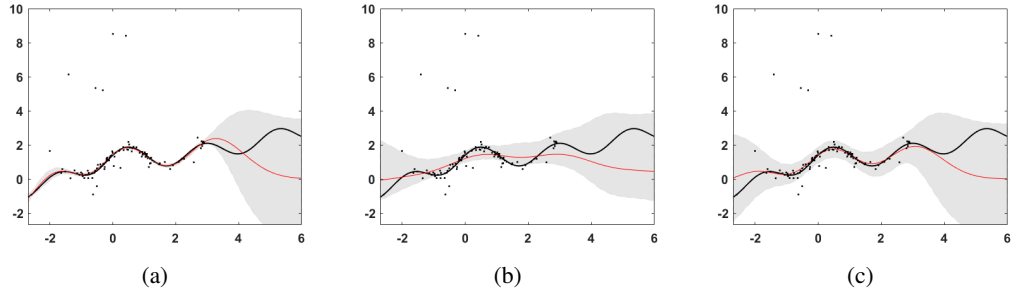


Figure 6: Predicted values for (a) tLA; (b) HuberMCMC; (c) HuberLA with standard deviations for the Case 2 with error following Student’s t distribution on Neal dataset.

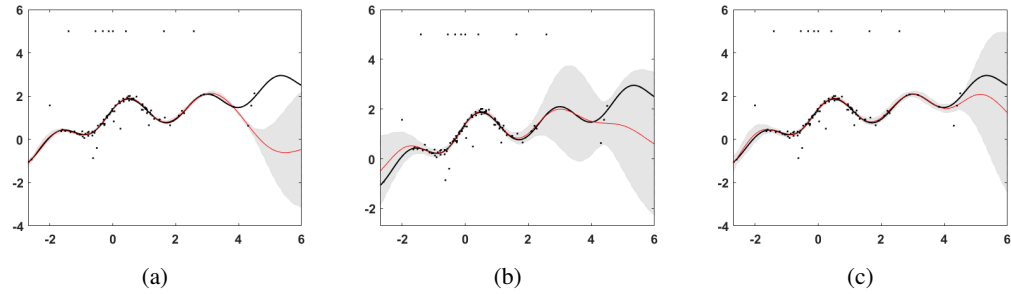


Figure 7: Predicted values for (a) tLA; (b) HuberMCMC; (c) HuberLA with standard deviations for the Case 3 with error following Student’s t distribution on Neal dataset.

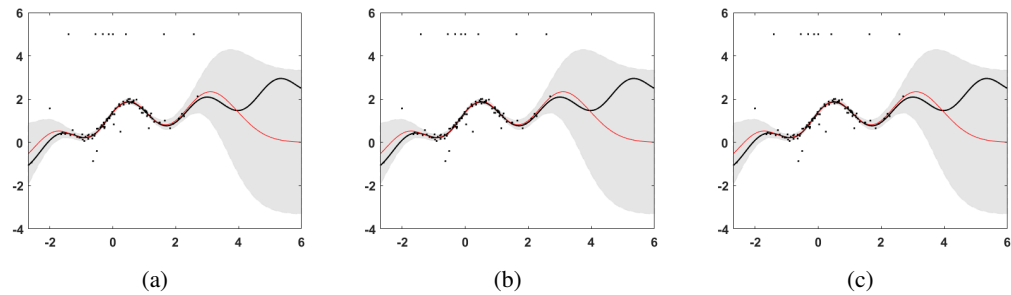


Figure 8: Predicted values for (a) tLA; (b) HuberMCMC; (c) HuberLA with standard deviations for the Case 4 with error following Student’s t distribution on Neal dataset.

## C.2 RESULTS WHEN PURSUIT WEIGHTING ADDED TO BASELINES

We ran 10 simulations with random hyperparameter initialization, reporting standard deviations in parentheses in Tables 7, 8 (Cases 1 and 3) and Tables 9, 10 (Cases 2 and 4).

In Cases 2 and 4, pursuit weights equal 1 since no covariate outliers ( $x^{(l)}$ ) are added. Thus, we denote these as ( $-pw$ ) as they are redundant.

	SCtMCMC <sup>+pw</sup>	tLA <sup>+pw</sup>	HuberMCMC <sup>+pw</sup>	HuberLA <sup>+pw</sup>	RCGP <sup>+pw</sup>	GP <sup>+pw</sup>	LaplaceMCMC <sup>+pw</sup>
	$\varepsilon \sim \mathcal{N}(0.01, 0.08)$						
RMSE	2.14 (2.32)	1.26 (0.05)	<b>0.63</b> (0.14)	1.18 (0.01)	2.04 (0.005)	1.56 (0.01)	0.71 (0.11)
MAE	1.07 (0.62)	0.82 (0.04)	<b>0.47</b> (0.08)	0.81 (0.01)	1.78 (0.003)	1.38 (0.02)	0.50 (0.07)
	$\varepsilon \sim \text{Student-}t(10)$						
RMSE	2.25 (1.56)	1.31 (0.00)	<b>0.61</b> (0.18)	1.19 (0.00)	2.06 (0.00)	1.58 (0.00)	0.74 (0.16)
MAE	1.11 (0.35)	0.86 (0.00)	<b>0.47</b> (0.11)	0.83 (0.00)	1.79 (0.00)	1.40 (0.00)	0.53 (0.10)
	$\varepsilon \sim \text{Laplace}(0, 0.1)$						
RMSE	1.75 (1.05)	1.31 (0.00)	<b>0.65</b> (0.15)	1.19 (0.00)	2.06 (0.002)	1.58 (0.00)	0.69 (0.16)
MAE	1.00 (0.24)	0.86 (0.00)	<b>0.49</b> (0.10)	0.83 (0.00)	1.08 (0.00)	1.40 (0.00)	0.50 (0.09)
	$\varepsilon \sim \text{Student-}t(1) \text{ (Cauchy)}$						
RMSE	1.33 (0.12)	1.31 (0.00)	<b>0.67</b> (0.23)	1.19 (0.00)	2.06 (0.00)	1.57 (0.00)	0.71 (0.16)
MAE	0.90 (0.05)	0.86 (0.00)	<b>0.50</b> (0.15)	0.83 (0.00)	1.05 (0.00)	1.39 (0.00)	0.52 (0.09)

Table 7: RMSE and MAE values on the Neal dataset for the Case 1 (with standard errors in parenthesis). Bold values highlight the best performance with the lowest RMSE and MAE.

	SCtMCMC <sup>+pw</sup>	tLA <sup>+pw</sup>	HuberMCMC <sup>+pw</sup>	HuberLA <sup>+pw</sup>	RCGP <sup>+pw</sup>	GP <sup>+pw</sup>	LaplaceMCMC <sup>+pw</sup>
	$\varepsilon \sim \mathcal{N}(0.01, 0.08)$						
RMSE	0.91 (0.04)	0.97 (0.14)	1.14 (0.23)	<b>0.84</b> (0.09)	2.04 (0.004)	0.90 (0.02)	0.90 (0.09)
MAE	0.62 (0.03)	0.59 (0.05)	0.76 (0.16)	<b>0.57</b> (0.06)	1.68 (0.002)	0.67 (0.02)	0.61 (0.04)
	$\varepsilon \sim \text{Student-}t(10)$						
RMSE	0.85 (0.02)	0.81 (0.00)	0.98 (0.12)	<b>0.80</b> (0.08)	2.05 (0.00)	0.88 (0.00)	0.87 (0.12)
MAE	0.61 (0.01)	<b>0.56</b> (0.00)	0.67 (0.08)	<b>0.56</b> (0.06)	1.70 (0.001)	0.66 (0.00)	0.62 (0.08)
	$\varepsilon \sim \text{Laplace}(0, 0.1)$						
RMSE	0.85 (0.02)	<b>0.81</b> (0.00)	0.99 (0.13)	0.83 (0.07)	2.05 (0.002)	0.88 (0.00)	0.84 (0.07)
MAE	0.61 (0.01)	<b>0.55</b> (0.00)	0.68 (0.10)	0.58 (0.05)	1.69 (0.00)	0.66 (0.00)	0.59 (0.04)
	$\varepsilon \sim \text{Student-}t(1) \text{ (Cauchy)}$						
RMSE	0.86 (0.04)	0.82 (0.01)	0.91 (0.10)	<b>0.79</b> (0.06)	2.05 (0.00)	0.88 (0.01)	0.83 (0.07)
MAE	0.62 (0.03)	0.56 (0.01)	0.61 (0.08)	<b>0.55</b> (0.05)	1.69 (0.00)	0.66 (0.01)	0.58 (0.05)

Table 8: RMSE and MAE values on the Neal dataset for the Case 3 (with standard errors in parenthesis). Bold values highlight the best performance with the lowest RMSE and MAE.

	SCtMCMC <sup>-pw</sup>	tLA <sup>-pw</sup>	HuberMCMC <sup>-pw</sup>	HuberLA <sup>-pw</sup>	RCGP <sup>-pw</sup>	GP <sup>-pw</sup>	LaplaceMCMC <sup>-pw</sup>
	$\varepsilon \sim \mathcal{N}(0.01, 0.08)$						
RMSE	2.74 (1.84)	<b>1.23</b> (0.10)	1.24 (0.04)	1.30 (0.01)	2.03 (0.01)	1.68 (0.02)	1.29 (0.12)
MAE	1.10 (0.52)	<b>0.75</b> (0.09)	0.88 (0.03)	0.89 (0.01)	1.69 (0.008)	1.46 (0.02)	0.88 (0.06)
	$\varepsilon \sim \text{Student-}t(10)$						
RMSE	1.92 (1.54)	1.18 (0.02)	<b>1.17</b> (0.04)	1.29 (0.00)	2.07 (0.00)	1.70 (0.00)	1.26 (0.10)
MAE	0.92 (0.39)	<b>0.71</b> (0.03)	0.84 (0.04)	0.88 (0.00)	1.70 (0.00)	1.48 (0.00)	0.89 (0.06)
	$\varepsilon \sim \text{Laplace}(0, 0.1)$						
RMSE	3.26 (1.83)	1.20 (0.04)	<b>1.16</b> (0.10)	1.29 (0.00)	2.08 (0.00)	1.70 (0.00)	1.17 (0.08)
MAE	1.26 (0.47)	<b>0.74</b> (0.05)	0.84 (0.06)	0.88 (0.00)	1.70 (0.00)	1.48 (0.00)	0.82 (0.06)
	$\varepsilon \sim \text{Student-}t(1) \text{ (Cauchy)}$						
RMSE	3.40 (1.82)	1.16 (0.06)	<b>1.12</b> (0.08)	1.29 (0.00)	2.08 (0.00)	1.70 (0.00)	1.19 (0.10)
MAE	1.31 (0.47)	<b>0.69</b> (0.06)	0.81 (0.05)	0.89 (0.00)	1.70 (0.00)	1.48 (0.00)	0.84 (0.07)

Table 9: Results for Case 2 (standard errors in parentheses)

## C.3 TRANSMISSION SPECTROSCOPY

Transmission spectroscopy records the relative change in the stellar flux, which is the incident photons per unit area, as a planet travels in front of the star around which it revolves. When the planet faces the star directly, known as a transit, it occludes a fraction of the stellar flux emitted by the star equal to the sky-projected area of the planet as compared to the area of the star, which is referred to as

	SctMCMC <sup>-pw</sup>	tLA <sup>-pw</sup>	HuberMCMC <sup>-pw</sup>	HuberLA <sup>-pw</sup>	RCGP <sup>-pw</sup>	GP <sup>-pw</sup>	LaplaceMCMC <sup>-pw</sup>
	$\varepsilon \sim \mathcal{N}(0.01, 0.08)$						
RMSE	3.45 (1.78)	1.17 (0.06)	1.28 (0.02)	1.18 (0.05)	1.99 (0.003)	1.68 (0.02)	<b>1.13</b> (0.11)
MAE	1.28 (0.50)	<b>0.67</b> (0.06)	0.87 (0.02)	0.85 (0.04)	1.74 (0.002)	1.45 (0.02)	0.79 (0.07)
	$\varepsilon \sim \text{Student-}t(10)$						
RMSE	2.43 (1.49)	1.21 (0.15)	1.24 (0.00)	<b>1.15</b> (0.09)	1.99 (0.00)	1.72 (0.00)	1.17 (0.10)
MAE	1.06 (0.45)	<b>0.76</b> (0.14)	0.91 (0.00)	0.83 (0.06)	1.75 (0.00)	1.47 (0.00)	0.82 (0.07)
	$\varepsilon \sim \text{Laplace}(0, 0.1)$						
RMSE	3.53 (1.69)	1.19 (0.10)	1.24 (0.00)	<b>1.14</b> (0.10)	1.99 (0.00)	1.72 (0.00)	1.17 (0.07)
MAE	1.33 (0.49)	<b>0.74</b> (0.10)	0.90 (0.00)	0.82 (0.06)	1.75 (0.00)	1.47 (0.00)	0.82 (0.06)
	$\varepsilon \sim \text{Student-}t(1) \text{ (Cauchy)}$						
RMSE	2.53 (1.82)	1.13 (0.04)	1.21 (0.03)	<b>1.08</b> (0.09)	2.00 (0.00)	1.71 (0.00)	1.15 (0.08)
MAE	1.04 (0.51)	<b>0.66</b> (0.05)	0.88 (0.02)	0.78 (0.07)	1.76 (0.00)	1.47 (0.00)	0.81 (0.06)

Table 10: Results for Case 4 (standard errors in parentheses)

transit depth. The measurement of the total flux over time is known as the light curve. The property on which the transmission spectroscopy relies to estimate the transit curve parameters is the planet’s transit depth, which depends on the wavelengths of the transmitted flux. For the wavelengths where the planet’s atmosphere is opaque due to the absorption of the emitted electromagnetic waves by constituent atoms or molecules, the planet blocks slightly more stellar flux. The variations are measured by binning the light curve into spectrophotometric channels of different wavelengths and by fitting the light curve from each channel separately with a transit model Kreidberg (2017).

The sources of error, such as photon noise and instrumental and astrophysical systematics, raise many potential challenges for precise atmosphere characterization. Pointing drift or modifications in the telescope focus influence the spectrum position on the detector to a small degree during transit due to instrumental systematics. Note that instrumental systematics are nothing but what is popularly known as systematic errors in statistics, which are here attributed to the atmospheric effects on the physical properties of an instrument. The optical state parameters are metered via auxiliary measurements of the spectral trace such as position, width, angle, or other parameters, indicating the state of the detector and optics, which are thought to be the cause of instrumental systematics. Instead of modeling the latter as a linear function of the optical state parameters, Gibson et al. (2012) proposed a non-parametric model by leveraging GPs.

The observation set obtained from HST- NICMOS includes the light curves for 18 wavelength channels extracted from  $n=638$  spectra along with six optical state parameters, namely the position of the spectral trace along the dispersion axis,  $\Delta X$ , the average position of the spectral trace along the cross-dispersion axis,  $\Delta Y$ , the angle of the spectral trace with the x-axis,  $W$ , the average width of the spectral trace,  $\psi^s$ , the temperature,  $T$ , and the orbital phase,  $\psi^H$ . The flux measurements contained in the vector,  $\mathbf{f} = [f_1, f_2, \dots, f_n]^T$ , are recorded at  $n$  time instants,  $\{t_1, t_2, \dots, t_n\}$ , contained in the time vector,  $\mathbf{t}$ , and the optical state parameters are given by  $\mathbf{x}_i = [\Delta X_i, \Delta Y_i, W_i, \psi_i^H, T_i, \psi_i^s]^T$  at time instant,  $t_i$ , collected in the matrix  $\mathbf{X} \in \mathbb{R}^{6 \times N}$  given by  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ .

The observed transit flux modeled in the GP framework follows a normal distribution, that is,

$$\mathbf{f}(\mathbf{t}, \mathbf{X}) \sim \mathcal{N}(\mathbf{T}(\mathbf{t}, \phi), \mathbf{K}(\mathbf{X}, \mathbf{X}|\theta)). \quad (33)$$

where the parameter vector,  $\phi$ , include the parameter of interest,  $\rho_{radius}$ , and other parameters, namely out-of-transit flux,  $f_{oot}$ , time gradient,  $T_{grad}$ , fixed central transit time,  $T_0$ , orbital period,  $P$ , limb darkening coefficient,  $c_1$ , limb darkening coefficient,  $c_2$ . The transit vector function,  $\mathbf{T}(\mathbf{t}, \phi)$ , is hereafter referred to as mean function parameter vector. The non-variable mean function parameters are fixed or calculated as stated in Gibson et al. (2012). Along with the planet-to-star radius ratio, the other mean function parameters are the parameters of a linear baseline model,  $f_{oot}$  and  $T_{grad}$ . The covariance matrix,  $\Sigma(\mathbf{x}_i, \mathbf{x}_j|\theta)$ , is the covariance between two output flux measurements defined as a function of the distance between optical state parameters,  $(\mathbf{x}_i, \mathbf{x}_j)$ , given by

$$K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j) + \delta_{ij}\sigma^2, \quad (34)$$

where  $k(\cdot, \cdot)$  is a Gaussian kernel. The threshold parameter,  $b$ , is set to 1.5 to achieve good robustness and efficiency at data distributed normally.

Table 11: Results of the planet-to-star radius ratio obtained from Gibson (2012) and GP-Huber.

Wavelength ( $\mu\text{m}$ )	Results from model in Gibson2012		Results obtained from GP-Huber	
	$\rho_{radius}$	$\Delta\rho_{radius}$	$\rho_{radius}$	$\Delta\rho_{radius}$
2.468	0.15545	0.00077	0.15525	0.00071
2.411	0.15520	0.00052	0.15771	0.0008911
2.353	0.15455	0.00044	0.15488	0.0004021
2.296	0.15513	0.00057	0.15825	0.0006526
2.238	0.15512	0.00041	0.1542	0.0005276
2.181	0.15504	0.00051	0.15297	0.0007462
2.124	0.15417	0.00066	0.15928	0.0007869
2.066	0.15508	0.00066	0.15525	0.000399
2.009	0.15393	0.00036	0.15259	0.0004077
1.951	0.15595	0.00051	0.15602	0.0005586
1.894	0.15549	0.0006	0.15466	0.0005988
1.837	0.15513	0.00053	0.15433	0.0004704
1.779	0.15534	0.00051	0.1537	0.0003601
1.722	0.15447	0.00087	0.14937	0.0006938
1.665	0.15429	0.00064	0.1517	0.000871
1.607	0.15266	0.00062	0.15213	0.0008045
1.55	0.15359	0.00073	0.15276	0.0007583
1.492	0.15367	0.00118	0.15256	0.0010653

The joint un-normalized log-posterior function of  $\phi$ ,  $\beta$ , and  $\theta$  is given by

$$\log P(\phi, \theta, \sigma^2, \beta | \mathbf{f}, \mathbf{X}, \zeta) = \log (\mathcal{L}(\mathbf{r} | \mathbf{X}, \phi, \theta, \sigma^2)) - \frac{\tau}{l_\tau} - \sum_{i=1}^d \left( \frac{1}{s_i l_i} \right) + \log(\beta) - \beta^T \sigma^2 + \log(p(\beta | \zeta)) + C. \quad (35)$$

Here, we lay the gamma a priori probability density function,  $p(\theta) = \frac{1}{l} \exp(-\frac{\theta}{l})$  over the covariance function hyperparameters  $\theta$ . The parameter  $l_\tau$  is of the gamma a priori associated with hyperparameter  $\tau$  and C represents additional constant terms. The samples of  $\beta_i$  are generated from log uniform distribution to lay a non-informative prior with parameter vector,  $\zeta$ , whereas  $p(\beta_g)$  is a degenerate probability density function.

The values of the planet-to-star radius ratio  $\rho_{radius}$  for each wavelength obtained from the GP-Huber model are shown in Table 11 along with those obtained from the model described in Gibson et al. (2012) referred to as Gibson2012, where  $\Delta\rho_{radius}$  represents the estimated uncertainty.

## D COMPUTATIONAL COSTS

Table 12 lists the computational costs of 10 random simulations for the Neal dataset under the Student’s-t noise setting. Similar costs were observed for other noise settings.

	SCtMCMC	tLA	HuberMCMC	HuberLA	RCGP	GP	LaplaceMCMC
Case 1	7.14 (0.31)	0.50 (0.14)	8.14 (0.38)	1.84 (0.47)	1.12 (0.01)	0.17 (0.04)	8.12 (0.37)
Case 2	7.32 (0.41)	0.69 (0.09)	8.37 (0.77)	1.64 (0.16)	1.23 (0.00)	0.19 (0.03)	8.34 (1.01)
Case 3	10.86 (1.20)	0.76 (0.27)	12.09 (0.41)	2.75 (0.25)	1.18 (0.00)	0.28 (0.01)	11.41 (0.45)
Case 4	10.14 (1.13)	0.56 (0.22)	11.62 (0.31)	2.45 (0.26)	1.27 (0.00)	0.28 (0.02)	11.20 (0.33)

Table 12: Computational costs (in seconds) on Neal dataset for Student’s-t noise setting. Standard deviations are given in parenthesis.



	<b>Outlier Case</b>	<b>RCGP</b>	<b>HuberMCMC</b>	<b>HuberLA</b>
Energy	Focused	3.01 (1.20)	155.71 (20.16)	8.43 (1.75)
	Asymmetric	3.26 (2.00)	131.23 (12.02)	9.79 (1.55)
Yacht	Focused	2.63 (0.4)	69.48 (18.37)	4.41 (0.98)
	Asymmetric	2.72 (0.00)	57.85 (8.12)	5.89 (1.75)

Table 13: Execution times (in seconds) for Energy and Yacht datasets.