

---

# Variational Inference with Censored Gaussian Process Regressors

---

Andrea Karlova<sup>1</sup> Rishabh Kabra<sup>2,1</sup> Daniel Augusto de Souza<sup>1</sup> Brooks Paige<sup>1</sup>

## Abstract

We consider the problem of Bayesian inference when some observations have been censored. In censored data, the dependent variable has been clipped, so we only know that the true value is at least as large (or as small) as the observation. Such data can be modeled using a Tobit likelihood, which can be viewed as a mixture between a normal distribution restricted on the domain without censoring treatment and a point mass at the boundary. This requires careful consideration when evaluating information-theoretic quantities, due to the mixed continuous and discrete probability measures. We introduce a novel approximate inference scheme for Gaussian process models with a Tobit likelihood, derive interpretable analytic expression for the Gaussian process evidence lower bound (ELBO) and demonstrate the resulting model’s efficiency in learning Gaussian process posteriors for censored data relative to uncensored case.

## 1. Introduction

Measurements in the real world come from bounded-scale instruments or finite-time events. Such intervals require careful consideration at the boundaries, where the variable is clipped and the density ceases to be continuous. The Censored Gaussian aka Tobit likelihood is instrumental in modeling such data (Fig 2). While there has been recent work (Gammelli et al., 2022; Basson et al., 2023) advancing the inference of Censored Gaussian Processes, so far, most statistical and information-theoretic quantities of the likelihood distribution were not provided in the closed analytical form.

An appropriate likelihood assumption is essential to any latent modeling task. Choosing the wrong likelihood, such as the standard Gaussian, when it’s impossible to observe

data outside a given range, introduces systematic bias into the model. A misalignment between modeling assumptions and the empirical data distribution not only affects the model’s predictive ability, but also its uncertainty quantification. This presents a two-fold challenge in fields such as medical diagnostics (Rao et al., 2016) or environmental monitoring (Friederichs & Hense, 2007) which not only care more about model uncertainty, but are also likelier to encounter censored or missing data (Chen et al., 2013).

The Tobit likelihood, despite being an appropriate choice when a censoring interval is fixed and given, has found limited use in Bayesian inference (Basson et al., 2023). A key missing piece is an inference procedure that would work in conjunction with expressive models such as Gaussian processes (GPs), without relying on high-variance Monte Carlo estimation. While the posterior is analytically intractable (Ertin, 2007; Groot & Lucas, 2012), we provide feasible analytical approximation to the evidence lower bound, (Jaakkola & Jordan, 2000), which allows for more controllable inference, thus laying a more rigorous foundation for Bayesian inference with censored data. This enables the use of Variational Inference and gradient-based learning for GPs modeling censored data.

The main contribution of our paper is the derivation of the closed form formula for cross entropy of normal and censored normal distribution which we apply for formulating the interpretable closed-form evidence lower bound (ELBO) for Variational Inference when approximating the posterior corresponding to a Tobit likelihood. We demonstrate through GP regression experiments how a Tobit likelihood can be used as a plug-in replacement for the usual uncensored Gaussian.

## 2. Background and Related Work

We set out to infer a (latent) data generating function for censored data. Following a Bayesian approach, there are typically four choices to make: (i) a prior distribution  $p(f)$  over the latent function space, (ii) a likelihood distribution  $p(y_i|f)$  which specifies how the data would be generated given a latent function, (iii) an inference method for the posterior distribution  $p(f|y_1, \dots, y_N) \propto \prod_i p(y_i|f)p(f)$  given observed data, and (iv) a strategy to acquire new data points to aid learning. We go over each of these choices in

---

<sup>1</sup>University College London <sup>2</sup>Google DeepMind. Correspondence to: Andrea Karlova <a.karlova@ucl.ac.uk>.

Accepted by the Structured Probabilistic Inference & Generative Modeling workshop of ICML 2024, Vienna, Austria. Copyright 2024 by the author(s).

the context of our problem, and review the relevant literature below.

## 2.1. A Likelihood for Censored Data

In a traditional regression problem, we observe a set of noise-corrupted output values,  $\mathbf{y} \in \mathbb{R}^N$ , which are dependent variables of a set of input values  $\mathbf{X} \in \mathbb{R}^{N \times D}$ . The goal is to recover the noise-free function  $f(\mathbf{x})$ . A typical choice is to assume a Gaussian likelihood model. However, if we know that not only is the data collection noisy but also clamped to the interval  $[l, u]$ , this gives rise to a different corruption process known as *censoring*. Mathematically, this process is:

$$y(\mathbf{x}) = \begin{cases} l, & \text{if } f(\mathbf{x}) + \varepsilon \leq l \\ f(\mathbf{x}) + \varepsilon, & \text{if } l < f(\mathbf{x}) + \varepsilon < u, \\ u, & \text{if } f(\mathbf{x}) + \varepsilon \geq u \end{cases} \quad (1)$$

where  $\varepsilon \sim \mathcal{N}(0, \sigma_y^2)$ . Since this corruption corresponds to a measurable function, it also defines a distribution, known as the Tobit likelihood in reference to Tobin (Tobin, 1958). This approach has been extensively studied in econometrics (Amemiya, 1984; Robin, 2010) and survival analysis (Klein et al., 2003), with extensions to handle various censoring mechanisms and distributional assumptions, including machine learning (Pearce et al., 2022; Friederichs & Hense, 2007; Moradian et al., 2017).

Unfortunately, the Tobit likelihood is not absolutely continuous with respect to the Lebesgue measure; it lacks a defined probability density function (PDF), in contrast with the Gaussian likelihood. Nonetheless, the cumulative distribution function (CDF) exists as is defined as:

$$p(y \leq \xi | f(\mathbf{x}), l, u) = \begin{cases} 0, & \text{if } \xi < l \\ \Phi(\xi | f(\mathbf{x}), \sigma_y^2), & \text{if } l \leq \xi \leq u, \\ 1, & \text{if } \xi > u \end{cases} \quad (2)$$

where  $\Phi(x | \mu, \sigma^2)$  is the CDF of the Gaussian distribution  $\mathcal{N}(\mu, \sigma^2)$ .

## 2.2. Gaussian Processes (GPs)

After choosing a likelihood, another decision is which prior distribution over functions to use. The most common choice of prior in Bayesian active learning is the Gaussian process (GP) (Li et al., 2024). A GP is a simple, but flexible, distribution over functions built on the assumption that any two values  $f(x)$  and  $f(x')$  are correlated Gaussian random variables (Williams & Rasmussen, 2006). Given training and evaluation sets,  $\mathbf{X}$  and  $\mathbf{X}^*$ , we can represent the output of the latent functions as  $\mathbf{f} = [f(\mathbf{x}_1), \dots, f(\mathbf{x}_N)]$  and

$\mathbf{f}^* = [f(\mathbf{x}_1^*), \dots, f(\mathbf{x}_N^*)]$  with joint distribution:

$$p(\mathbf{f}, \mathbf{f}^*) = \mathcal{N} \left( \begin{bmatrix} \mathbf{f} \\ \mathbf{f}^* \end{bmatrix} \middle| \begin{bmatrix} \boldsymbol{\mu} \\ \boldsymbol{\mu}^* \end{bmatrix}, \begin{bmatrix} \mathbf{K}_f & \mathbf{K}_{f^*} \\ \mathbf{K}_{*f} & \mathbf{K}^* \end{bmatrix} \right), \quad (3)$$

where  $[\boldsymbol{\mu}_f]_i = \mu(x_i)$  and  $[\mathbf{K}_f]_{i,j} = k(x_i, x_j)$ . The mean function  $\mu(x)$  and kernel function  $k(x, x')$  are the parameters of the Gaussian process distribution and specify the mean and covariance of the joint distributions of observations.

The reason for the popularity of GPs is due to their conjugacy with the Gaussian likelihood which allows not only the posterior distribution  $p(\mathbf{f} | \mathbf{y})$  and predictive posterior distribution  $p(\mathbf{f}^* | \mathbf{y})$  to be calculated exactly in closed form, but also the model evidence  $p(\mathbf{y}) = \mathbb{E}_{p(\mathbf{f})} [p(\mathbf{y} | \mathbf{f})]$  enabling gradient-based selection of model hyperparameters by maximizing the model evidence.

## 2.3. Variational Inference for GPs

For likelihoods other than Gaussian, the posterior distribution usually cannot be computed and requires the use of approximations. Possible strategies to deal with this include Markov Chain Monte Carlo (MCMC) methods (Neal, 1997), Laplace approximations (Williams & Barber, 1998; Barrett & Coolen, 2012), and variational inference techniques (Hensman et al., 2015). Due to its deterministic nature compared to MCMC, higher expressivity compared to Laplace approximation, and connections to scalable GPs for big data, variational inference (VI) became the default choice for GP approximations.

Variational inference works by defining an approximate posterior  $p(\mathbf{f}, \mathbf{f}^*) \approx q(\mathbf{f}, \mathbf{f}^*)$  where the posterior of the training data is a Gaussian distribution with learnable parameters  $q(\mathbf{f}) = \mathcal{N}(\mathbf{f} | \hat{\mathbf{f}}, \mathbf{A})$  and the conditional predictive posterior is simply  $q(\mathbf{f} | \mathbf{f}^*) = p(\mathbf{f}^* | \mathbf{f})$ . The parameters of the approximation are trained by maximizing the evidence lower bound (ELBO):

$$\log p(\mathbf{y}) \geq \mathbb{E}_{q(\mathbf{f})} [\log p(\mathbf{y} | \mathbf{f})] - KL[q(\mathbf{f}) || p(\mathbf{f})], \quad (4)$$

with the property that maximizing the ELBO jointly maximizes the evidence, allowing for hyperparameter optimization, and minimizes the divergence between the approximated posterior and the true posterior, reducing the error of the approximation. This property is one of the appealing factors of variational inference.

## 3. Our Method

### 3.1. Mixed Continuous and Discrete Measures

Censoring is a data treatment with the following mechanism: events which happen outside our observation window are

assigned the boundary value of this window. This does not change our perspective on how often the events occur (unlike truncation). But it changes the odds of how often events occur at the boundaries. Say we observe data from a normal distribution with values censored to lie in  $[l, u]$ . For an underlying distribution with density  $\mathcal{N}(x; \mu, \sigma)$ , we write the mixed density:

$$p(x; \mu, \sigma, l, u) = \Phi\left(\frac{l-\mu}{\sigma}\right)\delta_{\{l\}}(x) + \mathcal{N}(x; \mu, \sigma)\mathbb{I}_{(l,u)} + [1 - \Phi\left(\frac{u-\mu}{\sigma}\right)]\delta_{\{u\}}(x), \quad (5)$$

where  $\mathbb{I}_{(l,u)}$  is the indicator of belonging into the observation window  $(l, u)$  and  $\delta_{\{z\}}(x)$  is the Dirac delta function, having mass at  $x = z$  and null otherwise. The Dirac delta terms capture the size of the censored tails of the underlying density (see Fig 2).

When working with the censored probability measure, the usual framework which holds either for purely discrete or purely continuous probability distributions cannot be applied directly. We provide an exhaustive introduction to the problem in Appendix B. Essentially, we need to separately treat the areas where the probability density is defined to those with the discontinuities caused by mixing in the point masses of discrete distribution. While (Nair et al., 2006) denotes the mixed distribution by a mixed pair and considers the entropy of a mixed pair of discrete and continuous distribution in connection with entropy rate of MCMC, we carefully treat the problem using basic measure theory (Kallenberg, 2021).

### 3.2. Censored Regression

With the censored targets, the regression problem can be reformulated as:  $y_i = f_i + \epsilon_i$ , where  $i = 1, \dots, n$ , and  $\epsilon_i \sim \mathcal{N}_c(\epsilon|0, \sigma_y^2, l, u)$  comes from a censored normal distribution with the censoring thresholds  $l < u$ . We can see this problem also as a combination of two probit models placed on the lower threshold  $p(y_i^l = 1|\mathbf{x}) = \Phi(f(\mathbf{x})|l, \sigma_y^2)$ , and upper threshold  $p(y_i^u = 1|\mathbf{x}) = 1 - \Phi(f(\mathbf{x})|u, \sigma_y^2)$ , while having a normal regression model  $p(\mathbf{y}|\mathbf{f}, \sigma_y) = \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma_y^2 \mathbf{I}_{nn})$  otherwise.

From this initial motivation, we can write the likelihood  $p(\mathbf{y}|\mathbf{f})$  of the censored model as:

$$\prod_{y_i \leq l} \Phi(l|f_i, \sigma_y^2) \prod_{l < y_i < u} \mathcal{N}(y_i|f_i, \sigma_y^2) \prod_{y_i \geq u} [1 - \Phi(u|f_i, \sigma_y^2)]. \quad (6)$$

We place a GP prior over the latent function  $\mathbf{f}$  and obtain the following posterior distribution over the latent parameters  $\mathbf{f}$ :

$$p(\mathbf{f}|\mathcal{D}, \boldsymbol{\theta}) = \frac{\mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}_{nn})}{p(\mathcal{D}|\boldsymbol{\theta})} p(\mathbf{y}|\mathbf{f}), \quad (7)$$

where we denote the observed data as  $\mathcal{D} = \{(x_i, y_i), i = 1 \dots n\}$  and  $\mathbf{y} = (y_1, \dots, y_n)$  is the vector of labels.

### 3.3. ELBO for Censored Regressor

As the marginal likelihood  $p(\mathcal{D}|\boldsymbol{\theta})$  is not analytically tractable, we assume a Gaussian approximation to the posterior distribution, parametrised by  $\hat{\mathbf{f}}$  and  $\mathbf{A}$ :

$$p(\mathbf{f}|\mathcal{D}, \boldsymbol{\theta}) \approx q(\mathbf{f}|\mathcal{D}, \hat{\boldsymbol{\theta}}) = \mathcal{N}(\mathbf{f}|\hat{\mathbf{f}}, \mathbf{A}), \quad (8)$$

and minimize the KL divergence between the variational approximation of the posterior and the posterior itself. We write the KL divergence in terms of the censored normal likelihood  $p(\mathbf{y}|\mathbf{f})$  and Gaussian prior  $p(\mathbf{f})$  as:

$$\begin{aligned} \ln p(\mathcal{D}|\boldsymbol{\theta}) &\geq \sum_{i=1}^n \mathbb{E}_{q(\mathbf{f}_i|\mathcal{D}, \hat{\boldsymbol{\theta}})} \ln p(y_i|f_i) - \\ &- KL \left[ q(\mathbf{f}|\mathcal{D}, \hat{\boldsymbol{\theta}}) || p(\mathbf{f}|\mathbf{X}, \boldsymbol{\theta}) \right] \equiv ELBO \end{aligned} \quad (9)$$

The derivation of the cross-entropy is presented in Appendix C.1. The mean and scale parameters of the variational posterior  $q$  are denoted  $\hat{\mathbf{f}}$  and  $a_{ii}$ . Taking the final form of the cross-entropy from Eq (22), we arrive at the closed form version of the ELBO:

$$\begin{aligned} \ln p(\mathcal{D}|\boldsymbol{\theta}) &\geq \\ &- \sum_{i=1}^n \log[\sqrt{2\pi}\sigma_y] \left[ \Phi\left(\frac{u-\hat{f}_i}{a_{ii}}\right) - \Phi\left(\frac{l-\hat{f}_i}{a_{ii}}\right) \right] \\ &+ \sum_{i=1}^n \frac{(y_i-\hat{f}_i)^2 + a_{ii}^2}{2\sigma_y^2} \left[ \Phi\left(\frac{u-\hat{f}_i}{a_{ii}}\right) - \Phi\left(\frac{l-\hat{f}_i}{a_{ii}}\right) \right] \\ &+ \sum_{i=1}^n \left[ \log \Phi\left(\frac{y_i-u}{\sigma_y}\right) + \frac{(u+\hat{f}_i-2y_i)a_{ii}}{2\sigma_y^2} \right] \mathcal{N}\left(\frac{u-\hat{f}_i}{a_{ii}}\right) \\ &+ \sum_{i=1}^n \left[ \log \Phi\left(\frac{l-y_i}{\sigma_y}\right) - \frac{(l+\hat{f}_i-2y_i)a_{ii}}{2\sigma_y^2} \right] \mathcal{N}\left(\frac{l-\hat{f}_i}{a_{ii}}\right) \\ &- \frac{1}{2} \left[ \ln |\mathbf{A}| - \ln |\mathbf{K}_{nn}| - n \right] \\ &- \frac{1}{2} \left[ \text{tr}(\mathbf{K}_{nn}^{-1}\mathbf{A}) + \hat{\mathbf{f}}^T \mathbf{K}_{nn}^{-1} \hat{\mathbf{f}} \right]. \end{aligned} \quad (10)$$

The first two terms of the ELBO corresponds to the cross-entropy of the two uncensored Gaussian distributions less one half and scaled by variational approximation of the probability mass of the observation window between  $l$  and  $u$ ; the third term captures the contribution of the upper censoring limit  $u$  while the fourth term corresponds to the contribution at  $l$ . The logarithm of normal CDF times the normal density is the contribution of the atomic mass at the boundary while the second linear term in the equation describes how far we are from the exact fit. It is easy to see that if  $l$  and  $u$  are set to  $-\text{inf}$  and  $\text{inf}$ , then the second and third terms vanish (because  $\mathcal{N}(x)$  becomes 0). The CDF difference in the first line reduces to 1.0, and the whole

Table 1: **Metrics evaluating 1D regression.** The subscript c denotes a metric computed on censored datapoints only.

	Censored Gaussian				Uncensored Gaussian			
	MAE	MAE <sub>c</sub>	NLPD	NLPD <sub>c</sub>	MAE	MAE <sub>c</sub>	NLPD	NLPD <sub>c</sub>
left-censored	0.40	0.14	7.87	15.74	0.67	0.59	2.03	1.93
right-censored	0.61	0.00	549.95	1,075.87	0.62	0.09	56.41	2.77
symmetric-censored	0.40	0.34	37.44	41.56	0.59	0.55	122.10	96.69

 Table 2: **Metrics on real-world regression tasks.**

	Censored Gaussian		Uncensored Gaussian	
	MAE	NLPD	MAE	NLPD
Credit Risk: (40% test data)	9.5663	25.3787	9.5690	11.1385
gbsg cancer: (given test set)	21.1881	17.8265	19.8403	42.9532

ELBO reduces to the ELBO of an uncensored Gaussian. The last two terms corresponds to KL divergence between Gaussian prior and Gaussian variational approximation.

## 4. Experiments

### 4.1. GP Regression using Variational Inference

In Fig 1, we compare three ways of learning a synthetic function underlying censored observations. We start with uncensored Gaussian likelihood, a default likelihood for GP regression, and assumes no knowledge of the censoring. We use Variational Inference to fit the model. The second model is Censored Gaussian fit via Montecarlo simulation: this method samples from the (Multivariate Normal) variational posterior to parameterize the Censored Gaussian likelihood. We take 10,000 samples to facilitate convergence. The last model we consider is Censored Gaussian fit via Variational Inference: we use the posterior parameters to compute a closed-form expected log-likelihood as part of the ELBO, as derived in (10).

All methods use an RBF kernel with the same initialization, with a full matrix (Cholesky parameterized) covariance. We use L-BFGS (Liu & Nocedal, 1989) for optimization.

Fig 1 shows the benefit of taking censoring into account. The Uncensored Gaussian likelihood matches the empirical data density but completely fails to recover the true structure beyond the censoring bounds. It exhibits very low uncertainty around the censored regions due to the accumulated data density. The Censored Gaussian likelihood shows a significant improvement in capturing the structure of the generative function. Even with MC simulation from the approximate posterior, we achieve an almost unbiased fit. The caveat is the method is high-variance and very inefficient. We achieve the best fits using Variational Inference at significantly lower sample-complexity.

**Metrics.** To compare the uncensored and censored fits quantitatively, we compute the following 2 metrics on test-case data: Mean Absolute Error (MAE) and Negative Log Predictive Density (NLPD). MAE corresponds to assuming a common Laplace predictive likelihood (with scale 1). It is unambiguous and fair to both the uncensored and censored Gaussian likelihood models. For NLPD we evaluate both the uncensored and censored posteriors under the *censored* predictive likelihood—albeit with the respective noise parameters. In simple terms, this corresponds to passing the uncensored likelihood model through a clamping function.

**Real-World Regression Tasks.** We tested our model (10) on two real-life dataset with naturally censored dependent variables: the GBSG Cancer dataset (Katzman et al., 2018) and a Credit Risk dataset<sup>1</sup>. The results in Table 2 show that the Censored Gaussian is a drop-in replacement for the Uncensored likelihood when trained with our ELBO.

## 5. Discussion

We’ve shown that for GP Censored Regression it is in fact possible to derive tractable objective functions and validated our findings empirically by comparing a Censored Gaussian Regressor with the default Uncensored likelihood and comparing Gaussian processes fit using our derived ELBO versus a MC fit. Our derivations provide a general template to handle mixtures of discrete and continuous measures, or measures restricted to specific domains. Our work enables the adoption of more realistic modeling assumptions, due to its interpretability, which reflect how we observe the real world in practice.

<sup>1</sup>[https://github.com/square/pysurvival/raw/master/pysurvival/datasets/credit\\_risk.csv](https://github.com/square/pysurvival/raw/master/pysurvival/datasets/credit_risk.csv)

## References

- Amemiya, T. Tobit models: A survey. *Journal of econometrics*, 24(1-2):3–61, 1984.
- Barrett, J. and Coolen, A. Gaussian process regression for survival analysis with interval censored data. *Biometrika*, 99(1):1–11, 2012.
- Basson, M., Louw, T., and Smith, T. Variational tobit gaussian process regression. *Statistics and Computing*, 33(3), March 2023. ISSN 0960-3174. doi: 10.1007/s11222-023-10225-3.
- Chen, N., Qian, Z., Nabney, I. T., and Meng, X. Wind power forecasts using gaussian processes and numerical weather prediction. *IEEE Transactions on Power Systems*, 29(2): 656–665, 2013.
- Ertin, E. Gaussian process models for censored sensor readings. pp. 665–669, Los Alamitos, CA, USA, aug 2007. IEEE Computer Society. doi: 10.1109/SSP.2007.4301342.
- Friederichs, P. and Hense, A. Statistical downscaling of extreme precipitation events using censored quantile regression. *Monthly Weather Review*, 135(6):2365–2378, June 2007. ISSN 0027-0644. doi: 10.1175/mwr3403.1. URL <http://dx.doi.org/10.1175/MWR3403.1>.
- Gammelli, D., Rolsted, K. P., Pacino, D., and Rodrigues, F. Generalized multi-output gaussian process censored regression. *Pattern Recognition*, 129:108751, 2022.
- Groot, P. and Lucas, P. J. Gaussian process regression with censored data using expectation propagation. 2012.
- Hensman, J., Matthews, A., and Ghahramani, Z. Scalable variational gaussian process classification. In *Artificial Intelligence and Statistics*, pp. 351–360. PMLR, 2015.
- Jaakkola, T. S. and Jordan, M. I. Bayesian parameter estimation via variational methods. *Statistics and Computing*, 10, 2000. doi: 10.1023/A:1008932416310.
- Kallenberg, O. *Foundations of Modern Probability*, volume 1. Springer, 3 edition, 2021.
- Katzman, J. L., Shaham, U., Cloninger, A., Bates, J., Jiang, T., and Kluger, Y. Deepsurv: personalized treatment recommender system using a cox proportional hazards deep neural network. *BMC Medical Research Methodology*, 18(1), February 2018. ISSN 1471-2288. doi: 10.1186/s12874-018-0482-1. URL <http://dx.doi.org/10.1186/s12874-018-0482-1>.
- Klein, J. P., Moeschberger, M. L., et al. *Survival analysis: techniques for censored and truncated data*, volume 1230. Springer, 2003.
- Li, Y. L., Rudner, T. G. J., and Wilson, A. G. A study of bayesian neural network surrogates for bayesian optimization. In *The Twelfth International Conference on Learning Representations, ICLR 2024*, 2024.
- Liu, D. C. and Nocedal, J. On the limited memory bfgs method for large scale optimization. *Mathematical programming*, 45(1):503–528, 1989.
- Moradian, H., Larocque, D., and Bellavance, F. Survival forests for data with dependent censoring. *Statistical Methods in Medical Research*, 28(2):445–461, August 2017. ISSN 1477-0334. doi: 10.1177/0962280217727314. URL <http://dx.doi.org/10.1177/0962280217727314>.
- Nair, C., Prabhakar, B., and Shah, D. On entropy for mixtures of discrete and continuous variables. *arXiv preprint cs/0607075*, 2006.
- Neal, R. M. Monte carlo implementation of gaussian process models for bayesian regression and classification. *arXiv preprint physics/9701026*, 1997.
- Pearce, T., Jeong, J.-H., jia, y., and Zhu, J. Censored quantile regression neural networks for distribution-free survival analysis. In Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., and Oh, A. (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 7450–7461. Curran Associates, Inc., 2022.
- Rao, A., Monteiro, J., and Mourao-Miranda, J. Prediction of clinical scores from neuroimaging data with censored likelihood gaussian processes. In *2016 International Workshop on Pattern Recognition in Neuroimaging (PRNI)*, pp. 1–4. IEEE, 2016.
- Robin, J.-M. *Tobit model*, pp. 323–328. Palgrave Macmillan UK, London, 2010. ISBN 978-0-230-28081-6. doi: 10.1057/9780230280816\_35. URL [https://doi.org/10.1057/9780230280816\\_35](https://doi.org/10.1057/9780230280816_35).
- Tobin, J. Estimation of relationships for limited dependent variables. *Econometrica: journal of the Econometric Society*, pp. 24–36, 1958.
- Williams, C. K. and Barber, D. Bayesian classification with gaussian processes. *IEEE Transactions on pattern analysis and machine intelligence*, 20(12):1342–1351, 1998.
- Williams, C. K. and Rasmussen, C. E. *Gaussian processes for machine learning*, volume 2. MIT press Cambridge, MA, 2006.

## A. Experiments

We present figure with the results<sup>2</sup> from the investigation of the 1D regression. We fit the underlying generative function,  $y = (5x - 10) \sin(10x - 20) + \epsilon$  using  $n = 15$  data points. These are sampled and censored differently. We consider case of left censoring at bound  $-3$ , right censoring at bound  $0$  and symmetric censoring with bounds  $-2$  and  $2$ .

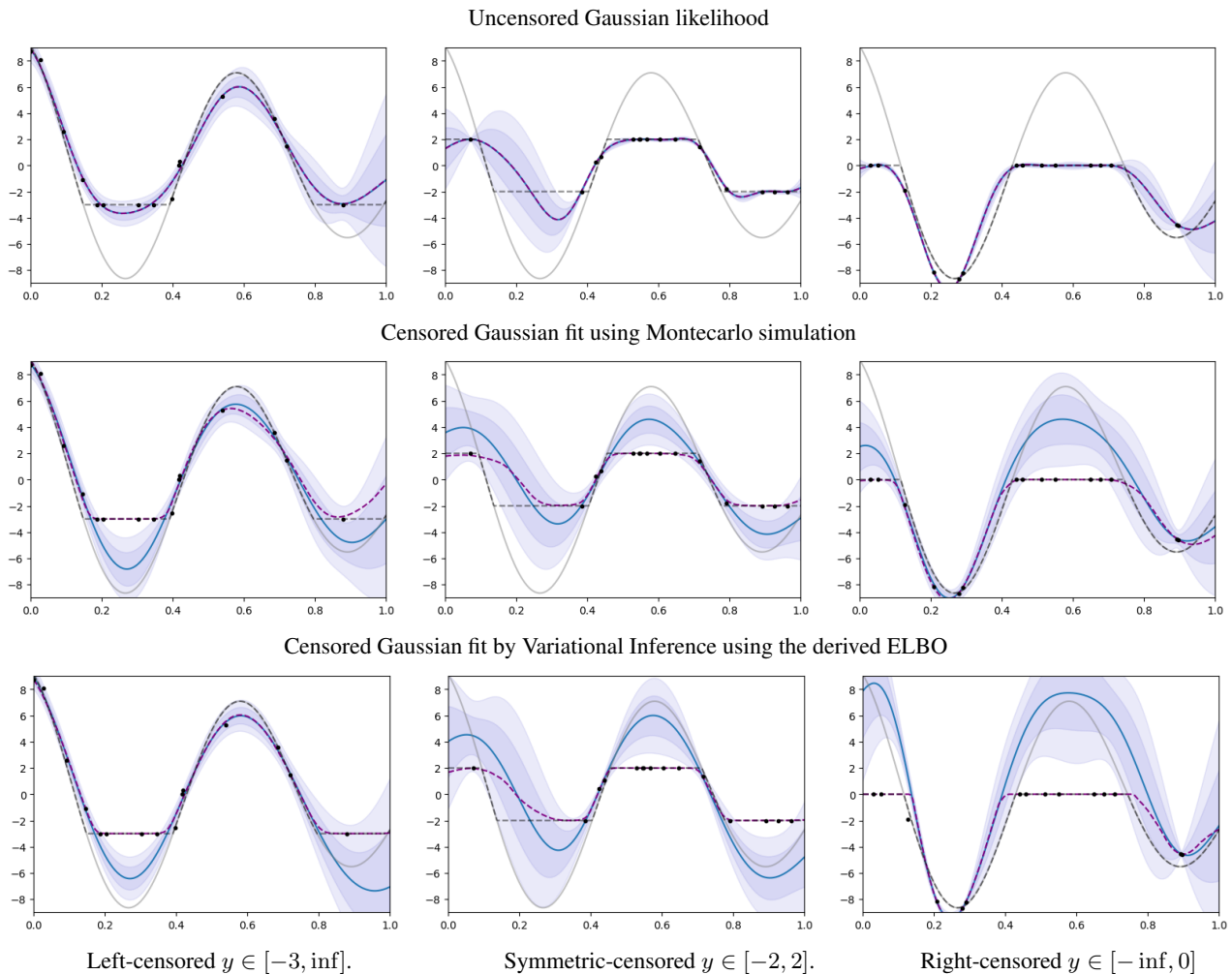


Figure 1: **One-dimensional GP regression.** We fit the underlying generative function,  $y = (5x - 10) \sin(10x - 20) + \epsilon$  using  $n=15$  data points. These are sampled and censored differently per column. The solid blue and gray lines show the uncensored prediction (latent posterior mean) and ground-truth function respectively. The purple and gray dashed lines show the censored predictive posterior mean and censored ground-truth respectively. The regions shaded in dark and light blue show the [15%, 85%] and [2.5%, 97.5%] percentiles of the latent posterior respectively.

## B. Censoring and Truncating

Let us handle the missing data problem by either truncating or censoring the data. In both cases we can assume that the data comes from some latent generative process. The observations are revealed to us only when their values are within the observation window. When truncating the data, we discard the records when the outcome of the measurement falls outside the observation window. This makes us unaware of anything happening outside the observation boundaries including the quantity of the data we drop. On the other hand, during the censoring process, the values outside the observation window are assigned an artificial threshold value. So we keep the additional information about the tails by counting the occurrence

<sup>2</sup><https://sites.google.com/view/ml-with-censored-data>

of possible values outside our observation boundaries. The censoring is often induced by the sensitivity of measuring instruments or the information availability when the value is within a certain bounds. Both data treatments inevitably impacts our modeling choices: the probability distribution of the data which we observe is different to the distribution of the underlying latent data which we could observe if we had a better measuring equipment.

In the following we demonstrate how the operation of data censoring and truncation impacts the underlying probability measures. Consider a diffusion measure  $\nu_{df}$  that is supported on  $\mathbb{R}$ . Say the range of the values we have access to is restricted to some compact interval  $[l, u]$ . This naturally impacts our perception of the odds with which we measure particular events. So the first conclusion we make is that we observe the restriction of measure  $\nu_{df}$  on the interval  $[l, u]$ . Denote the measure of the data distribution observed by us as  $\nu_{emp}$ . From our view point, the values which we cannot observe due to the restriction are considered as impossible and these events have a zero probability. This implies that the supports of the measures  $\nu_{emp}$  and  $\nu_{df}$  differ:  $supp_{\nu_{emp}} \subset supp_{\nu_{df}}$ . We further conclude, that the measures are not equivalent as their domains of impossible events differ.

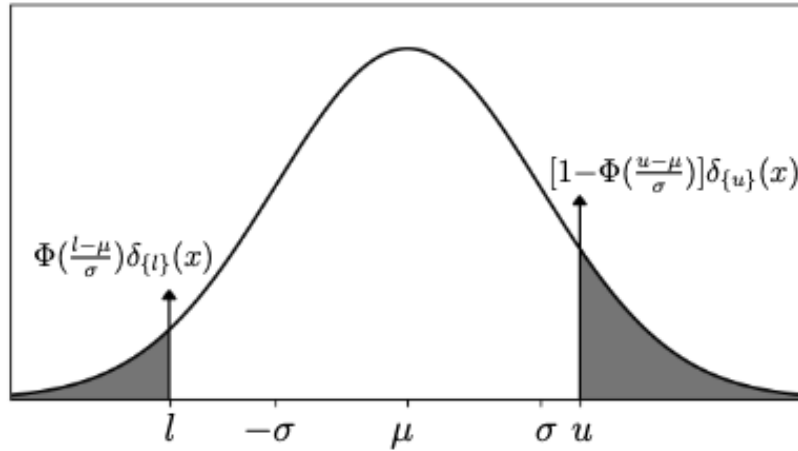


Figure 2: **A Tobit likelihood aka Censored Gaussian.** The support is restricted to the interval  $x \in [l, u]$ . Analytically this is handled by placing point masses at  $x = l$  and  $x = u$ . Despite the non-existence of the pdf at these points, we show how to derive the distribution's entropy and cross-entropy.

Let us associate a normal distribution with density  $\phi(x; \mu, \sigma)$  with our base distribution  $\nu_{df}$ . When we choose to discard the records outside our observation range, the values appear to us to be more frequent. This impacts the probability measure by changing its scale comparing to the original distribution. For example, for the underlying normal distribution with density  $\phi(x; \mu, \sigma)$  we have:

$$c \int_l^u \phi(x; \mu, \sigma) dx = c \left[ \Phi\left(\frac{u-\mu}{\sigma}\right) - \Phi\left(\frac{l-\mu}{\sigma}\right) \right] = 1, \quad (11)$$

and we simply infer the required scaling constant  $c$ . The truncated probability measure  $\nu_{emp}^t$  is given by:

$$\nu_{emp}^t \equiv 0 \cdot \nu_{df} \delta_{(-\infty, l)} + \frac{1}{C} \nu_{df} \delta_{[l, u]} + 0 \cdot \nu_{df} \delta_{(u, \infty)}, \quad \text{with } C \equiv \Phi\left(\frac{u-\mu}{\sigma}\right) - \Phi\left(\frac{l-\mu}{\sigma}\right). \quad (12)$$

When censoring, the events which happens outside our observation domain are assigned the boundary value of the domain. This does not impact our perspective on how often the events occur. Nevertheless, it changes the odds of how often occur the events on the boundaries. We need to modify the latent measure  $\nu_{df}$  restricted on the interval  $[l, u]$  by placing the two atomic measures on its boundaries, where the mass of the atoms corresponds to the masses of the tails of the base latent measure:

$$\nu_{emp}^c \equiv 0 \cdot \nu_{df} \delta_{(-\infty, l)} + \Phi\left(\frac{l-\mu}{\sigma}\right) \cdot \nu_{df} \delta_l + \nu_{df} \delta_{(l, u)} + [1 - \Phi\left(\frac{u-\mu}{\sigma}\right)] \cdot \nu_{df} \delta_u + 0 \cdot \nu_{df} \delta_{(u, \infty)}. \quad (13)$$

These changes to the probability distribution introduce a bias to the standard modeling approaches. For example, if we consider fitting linear regression to the censored or truncated data, the OLS estimator becomes biased and inconsistent. Figure 1 demonstrates the bias introduced to the OLS fit with the true values  $\beta_0 = 1, \beta_1 = 1$  and  $\sigma = 2$  of the parameters

of the underlying generative process  $f(x) = \beta_0 + \beta_1 x + \sigma \varepsilon^2$  and  $\varepsilon \sim N(0, 1)$ . When we consider the Bayesian linear model, choosing the wrong likelihood distribution will result in introducing the bias too. Figure 2 demonstrates the impact of placing the wrong probability distribution over the likelihood of the modified data. When placing the censored normal distribution over the likelihood of the censored data, we recover the correct estimates of the parameters of the underlying generative process. The correction for the linear regression model, which removes the bias of OLS estimators are available in the literature. Correcting the biases for more complex high-dimensional models is yet to be systematically explored.

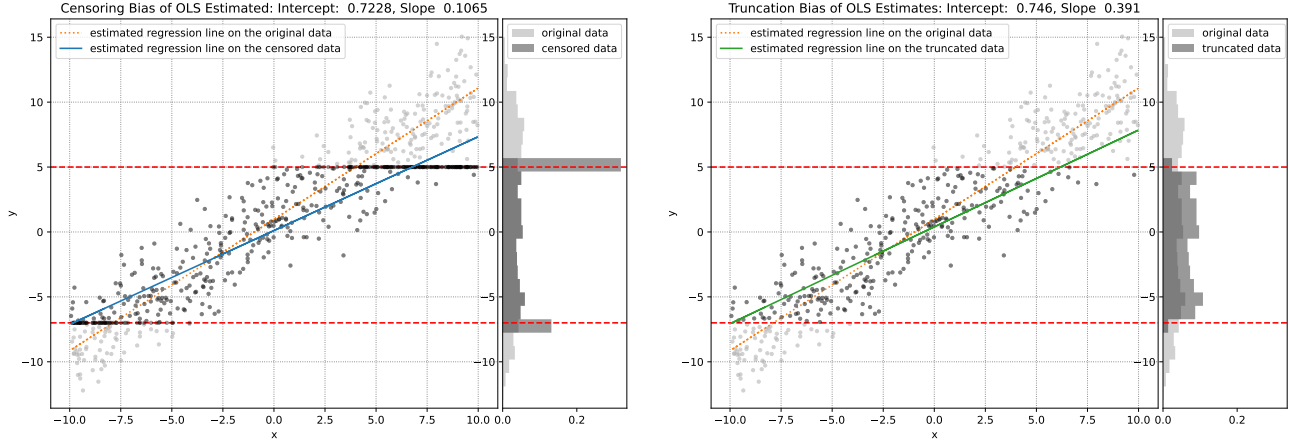


Figure 3: The bias of the OLS estimator fitted to the censored and truncated data. The original data are sampled from the regression line  $f(x) = 1 + x + 2\varepsilon^2$  where  $\varepsilon \sim N(0, 1)$ .

### C. Cross Entropy of Censored Normal

Using the measure defined in (13) we derive the cross entropy of the censored normal distribution. The cross entropy of  $\nu$  wrt to reference measure  $\rho$  defined on some measurable space  $(\Omega, \mathcal{A})$ , with  $\nu$  being absolutely continuous wrt to a reference measure  $\rho$ , is defined as:

$$H\left(\frac{d\nu}{d\rho}\right) = \int_{\Omega} \log \frac{d\nu}{d\rho} \rho(\omega), \quad (14)$$

where  $d\nu/d\rho$  corresponds to the Radon-Nikodym derivative of the two considered measures.

Recall from section B that the latent diffusive measure  $\nu_{df}$  is associated with the Lebesgue measure with mixed-in Dirac measures which guarantees the existence of the density. The Radon-Nikodym decomposition of the measures characterises the decomposition of the measure into its diffusive part and the discrete part, i.e. allows to mix in the atoms which corresponds to the singular measures wrt Lebesgue measure and atomic measures without jeopardizing the existence of the integral (14). This does not pose any issue for the existence of the cross-entropy as long as the absolute continuity of measure  $\nu$  wrt  $\rho$  is satisfied. We define the reference measure as  $\rho = \lambda + \delta_l + \delta_u$ , where  $\lambda$  denotes the Lebesgue measure. The measure  $\nu_{emp}^c$  is absolutely continuous wrt  $\rho$  as anytime  $\rho(A) = 0$ , then also  $\nu_{emp}^c(A) = 0$  for any  $A \in \mathcal{A}$ . The opposite does not hold and so the measures  $\nu_{emp}^c$  and  $\nu$  are not equivalent. The Radon-Nykodym derivative of  $\frac{d\nu_{emp}^c}{d\rho}$  defines the density  $p(\mu, \sigma, l, u)$ , where:

$$p(x; \mu, \sigma, l, u) = \Phi\left(\frac{l-\mu}{\sigma}\right)\delta_l(x) + \phi(x; \mu, \sigma)\mathbb{I}_{(l,u)}(x) + [1 - \Phi\left(\frac{u-\mu}{\sigma}\right)]\delta_u(x). \quad (15)$$

KL-divergence for the probability measures  $\nu_1, \nu_2$ , where  $\nu_2$  is absolutely continuous wrt to  $\nu_1$  and both measures are absolutely continuous wrt to  $\rho$ , is defined as follows:

$$D\left(\frac{d\nu_1}{d\rho} \parallel \frac{d\nu_2}{d\rho}\right) = \int_{\Omega} \left[ \log \frac{d\nu_2}{d\nu_1} \right] \frac{d\nu_1}{d\rho} \rho(\omega) \quad (16)$$



## Variational Inference with Censored Gaussian Process Regressors

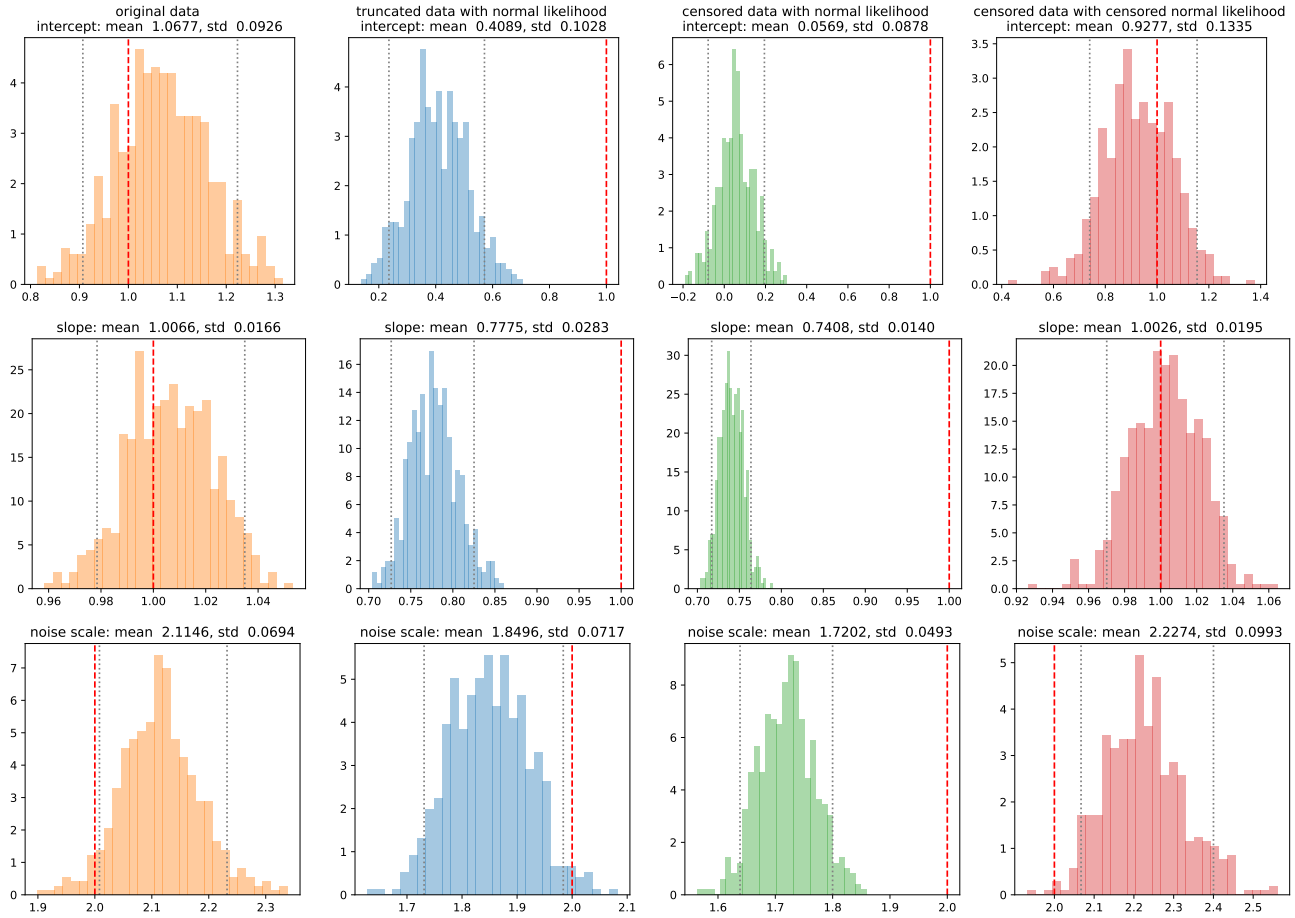


Figure 4: The posterior distribution of the parameters of the estimated linear model fitted to the dataset 1) without any data treatment, 2) with applying the truncation, 3) applying the the censoring. The red vertical line indicates the true value of the parameter used to generate the data. The dotted lines corresponds to 5% and 95% quantile.

This implies that we can only compute one-sided KL divergence between normal and censored normal probability distribution. We cannot swap the sides and the symmetrisation of the KL-divergence between normal and censored normal probability distribution is not possible. This is because the censored normal distribution is absolutely continuous wrt to the normal distribution. The normal distribution is not absolutely continuous wrt to the censored normal as it assigns non-zero mass to the tails, which the censored normal measures as a null set. This also provides guidance on how to compute KL-divergence between two censored normal distributions with different censoring bounds. The cross-entropy is the KL-divergence less the entropy of left-hand side probability measure, so the same rules applies to satisfy the existence of the integral.

### C.1. Derivation of the Cross Entropy

Let us derive the (negative) cross entropy of the normal distribution with the density  $\phi(m, s)$  and censored normal distribution  $p(\mu, \sigma; l, u)$  with the underlying normal density  $\phi(\mu, \sigma)$ :

$$\begin{aligned}
 & \mathbb{E}_{\phi(m, s)} \log p(l, u) \\
 &= \int_{\mathbb{R}} \phi(x; m, s) \log \left[ \Phi\left(\frac{l-\mu}{\sigma}\right) \delta_l(x) + \phi(x; \mu, \sigma) \delta_{(l, u)}(x) + [1 - \Phi\left(\frac{u-\mu}{\sigma}\right)] \delta_u(x) \right] dx \\
 &= \left[ \phi(x; m, s) \log \left[ \Phi\left(\frac{l-\mu}{\sigma}\right) \delta_l(x) + \phi(x; \mu, \sigma) \delta_{(l, u)}(x) + [1 - \Phi\left(\frac{u-\mu}{\sigma}\right)] \delta_u(x) \right] \right]_{x=l} \\
 &\quad + \int_l^u \phi(x; m, s) \log \left[ \Phi\left(\frac{l-\mu}{\sigma}\right) \delta_l(x) + \phi(x; \mu, \sigma) \delta_{(l, u)}(x) + [1 - \Phi\left(\frac{u-\mu}{\sigma}\right)] \delta_u(x) \right] dx \\
 &\quad + \left[ \phi(x; m, s) \log \left[ \Phi\left(\frac{l-\mu}{\sigma}\right) \delta_l(x) + \phi(x; \mu, \sigma) \delta_{(l, u)}(x) + [1 - \Phi\left(\frac{u-\mu}{\sigma}\right)] \delta_u(x) \right] \right]_{x=u} \\
 &= \phi\left(\frac{l-m}{s}\right) \log[\Phi\left(\frac{l-\mu}{\sigma}\right)] + \int_l^u \phi(x; m, s) \log \phi(x; \mu, \sigma) dx + \phi\left(\frac{u-m}{s}\right) \log[\Phi\left(\frac{\mu-u}{\sigma}\right)].
 \end{aligned} \tag{17}$$

Let us compute the middle term which corresponds to the cross-entropy of two unscaled normal distributions restricted to the compact interval  $[l, u]$ :

$$\begin{aligned}
 I &\equiv \int_l^u \phi(x; m, s) \log \phi(x; \mu, \sigma) dx \\
 &= -\frac{1}{2} \log[2\pi\sigma^2] \int_l^u \phi(x; m, s) dx - \frac{1}{\sqrt{2\pi s^2}} \int_l^u \frac{(x-\mu)^2}{2\sigma^2} e^{-\frac{(x-m)^2}{2s^2}} dx.
 \end{aligned} \tag{18}$$

The first term is the difference of the normal CDF at the transformed boundaries. For computing the second term we use the following trick:  $[(x-m) + (m-\mu)]^2 = (x-m)^2 + 2x(m-\mu) + \mu^2 - m^2$ :

$$\frac{1}{\sqrt{2\pi s^2}} \int_l^u \frac{(x-\mu)^2}{2\sigma^2} e^{-\frac{(x-m)^2}{2s^2}} dx = \frac{s^2}{\sigma^2 \sqrt{2\pi s^2}} \int_l^u \frac{(x-m)^2 + 2x(m-\mu) + \mu^2 - m^2}{2s^2} e^{-\frac{(x-m)^2}{2s^2}} dx. \tag{19}$$

To evaluate the first term in (19) we use substitution  $z = \frac{x-m}{s}$ . Because  $\left[-ze^{-\frac{z^2}{2}}\right] dz = z^2 e^{-\frac{z^2}{2}} - e^{-\frac{z^2}{2}}$ , we can compute the second term of the above integral in the closed form:

$$\begin{aligned}
 \frac{1}{\sqrt{2\pi}} \int_{\frac{l-m}{s}}^{\frac{u-m}{s}} z^2 e^{-\frac{z^2}{2}} dz &= \frac{1}{\sqrt{2\pi}} \int_{\frac{l-m}{s}}^{\frac{u-m}{s}} e^{-\frac{z^2}{2}} dz - \frac{1}{\sqrt{2\pi}} \left[ ze^{-\frac{z^2}{2}} \right]_{\frac{l-m}{s}}^{\frac{u-m}{s}} \\
 &= \Phi\left(\frac{u-m}{s}\right) - \Phi\left(\frac{l-m}{s}\right) - \frac{u-m}{s} \phi\left(\frac{u-m}{s}\right) + \frac{l-m}{s} \phi\left(\frac{l-m}{s}\right).
 \end{aligned} \tag{20}$$

Let us evaluate the middle term in (19). Recall  $\left[-e^{-\frac{z^2}{2}}\right] dz = ze^{-\frac{z^2}{2}}$  and substitute  $z = \frac{x-m}{s}$ :

$$\begin{aligned}
 \frac{1}{\sqrt{2\pi s^2}} \int_l^u x e^{-\frac{(x-m)^2}{2s^2}} dx &= \frac{s}{\sqrt{2\pi}} \int_{\frac{l-m}{s}}^{\frac{u-m}{s}} z e^{-\frac{z^2}{2}} dz + \frac{m}{\sqrt{2\pi s^2}} \int_{\frac{l-m}{s}}^{\frac{u-m}{s}} e^{-\frac{z^2}{2}} dz \\
 &= m[\Phi\left(\frac{u-m}{s}\right) - \Phi\left(\frac{l-m}{s}\right)] - s[\phi\left(\frac{u-m}{s}\right) - \phi\left(\frac{l-m}{s}\right)].
 \end{aligned} \tag{21}$$

Putting it all together we have the expression for the integral  $I$ :

$$\begin{aligned}
 I &= -\frac{1}{2} \log[2\pi\sigma^2] [\Phi\left(\frac{u-m}{s}\right) - \Phi\left(\frac{l-m}{s}\right)] - \frac{s^2}{2\sigma^2} [\Phi\left(\frac{u-m}{s}\right) - \Phi\left(\frac{l-m}{s}\right) - \frac{u-m}{s} \phi\left(\frac{u-m}{s}\right) + \frac{l-m}{s} \phi\left(\frac{l-m}{s}\right)] \\
 &\quad - \frac{2(m-\mu)}{2\sigma^2} \{m[\Phi\left(\frac{u-m}{s}\right) - \Phi\left(\frac{l-m}{s}\right)] - s[\phi\left(\frac{u-m}{s}\right) - \phi\left(\frac{l-m}{s}\right)]\} - \frac{\mu^2 - m^2}{2\sigma^2} [\Phi\left(\frac{u-m}{s}\right) - \Phi\left(\frac{l-m}{s}\right)] \\
 &= \left\{ \frac{m^2 - \mu^2 - 2m(m-\mu) - s^2}{2\sigma^2} - \frac{1}{2} \log[2\pi\sigma^2] \right\} [\Phi\left(\frac{u-m}{s}\right) - \Phi\left(\frac{l-m}{s}\right)] + \frac{(u+m-2\mu)s}{2\sigma^2} \phi\left(\frac{u-m}{s}\right) - \frac{(l+m-2\mu)s}{2\sigma^2} \phi\left(\frac{l-m}{s}\right).
 \end{aligned}$$

The formula for the cross entropy is:

$$\begin{aligned}
 & -\mathbb{E}_{\phi(m,s)} \log p(\mu, \sigma, l, u) \\
 &= -\phi\left(\frac{l-m}{s}\right) \log[\Phi\left(\frac{l-\mu}{\sigma}\right)] - I - \phi\left(\frac{u-m}{s}\right) \log[\Phi\left(\frac{\mu-u}{\sigma}\right)] \\
 &= \left\{\frac{1}{2} \log[2\pi\sigma^2] + \frac{(\mu-m)^2 + s^2}{2\sigma^2}\right\} [\Phi\left(\frac{u-m}{s}\right) - \Phi\left(\frac{l-m}{s}\right)] \\
 &\quad - \left\{\log \Phi\left(\frac{\mu-u}{\sigma}\right) + \frac{(u+m-2\mu)s}{2\sigma^2}\right\} \phi\left(\frac{u-m}{s}\right) \\
 &\quad - \left\{\log \Phi\left(\frac{l-\mu}{\sigma}\right) - \frac{(l+m-2\mu)s}{2\sigma^2}\right\} \phi\left(\frac{l-m}{s}\right)
 \end{aligned} \tag{22}$$

## D. Approximate Inference for Tobit Gaussian Process Regressors

Gaussian processes (GPs) are a popular choice of a prior for Bayesian non-parametric regression. Gaussian process is defined as a stochastic process indexed by a set  $\mathcal{X} \in \mathbb{R}^d$ : such that  $\{f(\mathbf{x}) : \mathbf{x} \in \mathcal{X}\}$ . Every finite combination of random variables of the process has a joint Gaussian distribution. Formally, GPs are distributions over functions. A GP is fully specified by its mean function  $\mu$  and covariance structure over some finite set  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ , GP is uniquely defined by  $p(\mathbf{f}|\mathcal{D}) = \mathcal{N}(\mathbf{f}|\boldsymbol{\mu}, \mathbf{K})$  with latent vector  $\mathbf{f} = (f(\mathbf{x}_1), \dots, f(\mathbf{x}_n))$ , mean vector  $\boldsymbol{\mu} = (\mu(\mathbf{x}_1), \dots, \mu(\mathbf{x}_n))$  and covariance matrix  $\mathbf{K} = k(\mathbf{x}_i, \mathbf{x}_j)_{i,j=1,\dots,n}$ .

We denote the observed data as  $\mathcal{D} = \{(\mathbf{x}_i, y_i), i = 1 \dots n\}$ , where  $\mathbf{y} = (y_1, \dots, y_n)$  is the vector of labels. To simplify notation, we also denote by  $\boldsymbol{\theta}$  the set of the Gaussian process parameters in general, e.g.  $\boldsymbol{\theta} = (\boldsymbol{\mu}, \mathbf{K})$ .

### D.1. Gaussian Process Model for Regression

To simplify the notation, we denote  $f_i = f(\mathbf{x}_i), i = 1 \dots n$ . The standard GP regression task is focused on estimating the latent function  $\mathbf{f}$  out of noisy observations  $\mathbf{y} \in \mathbb{R}^n$ :  $y_i = f_i + \varepsilon_i, i = 1, \dots, n$  and where  $\varepsilon_i \sim \mathcal{N}(\varepsilon|0, \sigma_y^2)$ .

To estimate the parameters, we can write the likelihood function, which is Gaussian due to the assumption on the noise term:

$$p(\mathbf{y}|\mathbf{f}, \sigma_y) = \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma_y^2 \mathbf{I}_{nn}), \tag{23}$$

where  $\mathbf{I}_{nn}$  denotes  $n \times n$  identity matrix.

To introduce Bayesian inference we assume zero mean prior GP with kernel function  $k(\mathbf{x}_i, \mathbf{x}_j): f \sim GP(0, k(\mathbf{x}_i, \mathbf{x}_j))$ . So the joint distribution of the latent function values corresponding to any set  $\mathbf{X}$  of input data is a multivariate Gaussian:

$$p(\mathbf{f}|\boldsymbol{\theta}) = \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}_{nn}), \tag{24}$$

where  $\boldsymbol{\theta}$  denotes the parameters of the kernel function and the noise parameter  $\sigma_y$ . Using Bayes rule we obtain for of posterior:

$$p(\mathbf{f}|\mathcal{D}, \boldsymbol{\theta}) = \frac{p(\mathbf{f}|\mathbf{X}, \boldsymbol{\theta})p(\mathbf{y}|\mathbf{f})}{p(\mathcal{D}|\boldsymbol{\theta})} = \frac{\mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}_{nn})}{p(\mathcal{D}|\boldsymbol{\theta})} \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma_y^2 \mathbf{I}_{nn}). \tag{25}$$

The marginal likelihood can be evaluated as:

$$p(\mathcal{D}|\boldsymbol{\theta}) = \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{X}, \boldsymbol{\theta})d\mathbf{f} = \int \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma_y^2 \mathbf{I}_{nn})\mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}_{nn})d\mathbf{f}. \tag{26}$$

Because we marginalise out from the two multivariate normals, the resulting distribution is also multivariate normal:  $p(\mathbf{y}|\boldsymbol{\theta}) = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{K}_{nn} + \sigma_y^2 \mathbf{I}_{nn})$ . To estimate the parameters  $\boldsymbol{\theta}$  of the marginal likelihood we can use gradient based optimisation to maximise the marginal log-likelihood of the model.

### D.2. Gaussian Process for Binary Classification

While the inference for regression task has many analytic advantages, even the binary classification loses analytical tractability. For the binary classification task:  $y_i$  takes one of the two values  $\{-1, 1\}$  only. The probability of the positive outcome  $p(y_i = 1|\mathbf{x}) = \tilde{\sigma}(f(\mathbf{x}))$ , where  $\tilde{\sigma}(\cdot)$  denotes the sigmoid transformation. For example, lets consider a probit model:  $p(y_i = 1|\mathbf{x}) = \Phi(f(\mathbf{x}))$ , where  $\Phi(\cdot)$  is the CDF of standard normal distribution. The class labels are Bernoulli variables

with parameter  $\Phi(f(x))$ . The likelihood for the binary classification case then corresponds to the product distribution:

$$p(\mathbf{y}|\mathbf{f}) = \prod_{i=1}^n p(y_i|f_i) = \prod_{i=1}^n \Phi(y_i f_i) \quad (27)$$

Placing the GP prior over the latent function  $\mathbf{f}$ , as in (28) leads to posterior:

$$p(\mathbf{f}|\mathcal{D}, \boldsymbol{\theta}) = \frac{p(\mathbf{f}|\mathbf{X}, \boldsymbol{\theta})}{p(\mathcal{D}|\boldsymbol{\theta})} \prod_{i=1}^n p(y_i|f_i) = \frac{\mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}_{nn})}{p(\mathcal{D}|\boldsymbol{\theta})} \prod_{i=1}^n \Phi(y_i f_i) \quad (28)$$

Following similar methodology as for the regression task, we would like to evaluate the marginal likelihood:

$$p(\mathcal{D}|\boldsymbol{\theta}) = \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{X}, \boldsymbol{\theta})d\mathbf{f} = \int \prod_{i=1}^n \Phi(y_i f_i) \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}_{nn})d\mathbf{f}, \quad (29)$$

however due to the Bernoulli likelihood, the marginal likelihood is analytically intractable.

### D.3. Gaussian Process Censored Regression

The censored regression problem corresponds to the case where the true target value  $\eta \in \mathbb{R}$  are partially unobservable and instead we observe the censored version of the target  $y$ . For the censoring lower and upper thresholds:  $l < u \in \mathbb{R}$ , we describe the censoring process is described as:

$$y = \begin{cases} l, & \text{if } \eta \leq l, \\ \eta, & \text{if } l < \eta < u, \\ u, & \text{if } \eta \geq u, \end{cases} \quad (30)$$

i.e. the values inside the censoring range are shifted to the boundaries of the range and assigned a single value.

Consider the standard GP regression from D.1, where we estimate the latent function  $\mathbf{f}$  out of noisy observations  $\boldsymbol{\eta} \in \mathbb{R}^n$ :  $\eta_i = f_i + \varepsilon_i, i = 1, \dots, n$  and where  $\varepsilon_i \sim \mathcal{N}(\varepsilon|0, \sigma_y^2)$ . Because the latent noisy targets are censored, our problem reformulates as:  $y_i = f_i + \varepsilon_i, i = 1, \dots, n$ , where  $\varepsilon_i \sim \mathcal{N}_c(\varepsilon|0, \sigma_y^2, l, u)$  comes from censored normal distribution with the censoring thresholds  $l < u$ . We can see this problem also as a combination of two probit models placed on the lower threshold  $p(y_i^l = 1|\mathbf{x}) = \Phi(f(\mathbf{x})|l, \sigma_y^2)$ , upper threshold  $p(y_i^u = 1|\mathbf{x}) = 1 - \Phi(f(\mathbf{x})|u, \sigma_y^2)$ , while having a normal regression model  $p(\mathbf{y}|\mathbf{f}, \sigma_y) = \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma_y^2 \mathbf{I}_{nn})$  otherwise.

From this initial motivation, we can write the likelihood of the censored model as:

$$p(\mathbf{y}|\mathbf{f}) = \prod_{y_i \leq l} \Phi(l|f_i, \sigma_y^2) \prod_{l < y_i < u} \mathcal{N}(y_i|f_i, \sigma_y^2) \prod_{y_i \geq u} [1 - \Phi(u|f_i, \sigma_y^2)]. \quad (31)$$

We place the GP prior over the latent function  $\mathbf{f}$  as in (24) and obtain the posterior distribution over these latent parameters  $\mathbf{f}$ :

$$\begin{aligned} p(\mathbf{f}|\mathcal{D}, \boldsymbol{\theta}) &= \frac{p(\mathbf{f}|\mathbf{X}, \boldsymbol{\theta})}{p(\mathcal{D}|\boldsymbol{\theta})} \prod_{i=1}^n p(y_i|f_i) \\ &= \frac{\mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}_{nn})}{p(\mathcal{D}|\boldsymbol{\theta})} \prod_{y_i \leq l} \Phi(l|f_i, \sigma_y^2) \prod_{l < y_i < u} \mathcal{N}(y_i|f_i, \sigma_y^2) \prod_{y_i \geq u} [1 - \Phi(u|f_i, \sigma_y^2)]. \end{aligned} \quad (32)$$

The marginal likelihood is analytically intractable due to the atoms on the censoring boundaries. Also note, that  $\Phi(x|m, s^2) = 1 - \Phi(m|x, s^2)$ . We will use this relationship to split the posterior into the truncated Gaussian part and the analytically

intractable mixture.

$$\begin{aligned}
 p(\mathcal{D}|\boldsymbol{\theta}) &= \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{X}, \boldsymbol{\theta})d\mathbf{f} \\
 &= \int \prod_{y_i \leq l} \Phi(l|f_i, \sigma_y^2) \prod_{l < y_i < u} \mathcal{N}(y_i|f_i, \sigma_y^2) \prod_{y_i \geq u} [1 - \Phi(u|f_i, \sigma_y^2)] \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}_{nn})d\mathbf{f} \\
 &= \int \mathcal{N}(\mathbf{y}\delta_{\mathbf{y} \in (l,u)}|\mathbf{f}, \sigma_y^2) \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}_{nn})d\mathbf{f} + \int \prod_{y_i \leq l} [1 - \Phi(f_i|l, \sigma_y^2)] \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}_{nn})d\mathbf{f} + \\
 &+ \int \prod_{y_i \geq u} \Phi(f_i|u, \sigma_y^2) \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}_{nn})d\mathbf{f} \\
 &= \mathcal{N}(\mathbf{y}\delta_{\mathbf{y} \in (l,u)}|\mathbf{0}, \mathbf{K}_{nn} + \sigma_y^2 \mathbf{I}_{nn}) + \int \prod_{y_i \leq l} [1 - \Phi(f_i|l, \sigma_y^2)] \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}_{nn})d\mathbf{f} + \\
 &+ \int \prod_{y_i \geq u} \Phi(f_i|u, \sigma_y^2) \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}_{nn})d\mathbf{f}.
 \end{aligned} \tag{33}$$

### D.3.1. VARIATIONAL INFERENCE FOR CENSORED REGRESSOR

Let us assume the Gaussian approximation of the posterior distribution:

$$p(\mathbf{f}|\mathcal{D}, \boldsymbol{\theta}) \approx q(\mathbf{f}|\mathcal{D}, \hat{\boldsymbol{\theta}}) = \mathcal{N}(\mathbf{f}|\hat{\mathbf{f}}, \mathbf{A}), \tag{34}$$

The goal is to minimize the KL divergence between the Gaussian approximation of the posterior and the posterior itself. In the following we provide the derivation of the evidence lower bound (ELBO) with the censored normal likelihood  $p(\mathbf{y}|\mathbf{f})$  and Gaussian prior  $p(\mathbf{f})$ :

$$\begin{aligned}
 &KL [q(\mathbf{f}|\mathcal{D}, \hat{\boldsymbol{\theta}})||p(\mathbf{f}|\mathcal{D}, \boldsymbol{\theta})] \\
 &= \int q(\mathbf{f}|\mathcal{D}, \hat{\boldsymbol{\theta}}) [\ln q(\mathbf{f}|\mathcal{D}, \hat{\boldsymbol{\theta}}) - \ln p(\mathbf{f}|\mathcal{D}, \boldsymbol{\theta})] d\mathbf{f} \\
 &= \int q(\mathbf{f}|\mathcal{D}, \hat{\boldsymbol{\theta}}) [\ln q(\mathbf{f}|\mathcal{D}, \hat{\boldsymbol{\theta}}) - \ln p(\mathbf{f}|\mathbf{X}, \boldsymbol{\theta}) - \ln p(\mathbf{y}|\mathbf{f}) + \ln p(\mathcal{D}|\boldsymbol{\theta})] d\mathbf{f} \\
 &= KL [q(\mathbf{f}|\mathcal{D}, \hat{\boldsymbol{\theta}})||p(\mathbf{f}|\mathbf{X}, \boldsymbol{\theta})] - \int q(\mathbf{f}|\mathcal{D}, \hat{\boldsymbol{\theta}}) \ln p(\mathbf{y}|\mathbf{f})d\mathbf{f} + \ln p(\mathcal{D}|\boldsymbol{\theta}).
 \end{aligned} \tag{35}$$

Because left handside term  $KL [q(\mathbf{f})||p(\mathbf{f}|\mathcal{D}, \boldsymbol{\theta})]$  is non-negative and we aim to minimize this distance, we write the ELBO as:

$$\begin{aligned}
 \ln p(\mathcal{D}|\boldsymbol{\theta}) &\geq \mathbb{E}_{q(\mathbf{f}|\mathcal{D}, \hat{\boldsymbol{\theta}})} [\ln p(\mathbf{y}|\mathbf{f})] - KL [q(\mathbf{f}|\mathcal{D}, \hat{\boldsymbol{\theta}})||p(\mathbf{f}|\mathbf{X}, \boldsymbol{\theta})] \\
 &= \sum_{i=1}^n \mathbb{E}_{q(f_i|\mathcal{D}, \hat{\boldsymbol{\theta}})} \ln p(y_i|f_i) - KL [q(\mathbf{f}|\mathcal{D}, \hat{\boldsymbol{\theta}})||p(\mathbf{f}|\mathbf{X}, \boldsymbol{\theta})].
 \end{aligned} \tag{36}$$

The first term comprises the (negative) cross entropies of the approximate posterior with respect to the likelihood. The cross entropy can be split into two terms: a factorised likelihood and a marginalisation term of the remaining variables. We can

then apply the derived crossed entropy between Gaussian and censored distribution. Let us derive the analytic form:

$$\begin{aligned}
 & \sum_{i=1}^n \mathbb{E}_{q(f_i|\mathcal{D}, \hat{\theta})} \log p(y_i|f_i) \\
 &= - \sum_{i=1}^n \left( \log[\sqrt{2\pi}\sigma_y] + \frac{(f_i - \hat{f}_i)^2 + a_{ii}^2}{2\sigma_y^2} \right) \left[ \Phi\left(\frac{u - \hat{f}_i}{a_{ii}}\right) - \Phi\left(\frac{l - \hat{f}_i}{a_{ii}}\right) \right] \\
 &+ \sum_{i=1}^n \left[ \log \Phi\left(\frac{f_i - u}{\sigma_y}\right) + \frac{(u + \hat{f}_i - 2f_i)a_{ii}}{2\sigma_y^2} \right] \phi\left(\frac{u - \hat{f}_i}{a_{ii}}\right) \\
 &+ \sum_{i=1}^n \left[ \log \Phi\left(\frac{l - f_i}{\sigma_y}\right) - \frac{(l + \hat{f}_i - 2f_i)a_{ii}}{2\sigma_y^2} \right] \phi\left(\frac{l - \hat{f}_i}{a_{ii}}\right).
 \end{aligned} \tag{37}$$

The second term is the KL divergence between two multivariate Gaussian distributions. It can be written in matrix calculus as follows:

$$\begin{aligned}
 KL[q(\mathbf{f}|\mathcal{D}, \hat{\theta})||p(\mathbf{f}|\mathbf{X}, \theta)] &= KL[\mathcal{N}(\mathbf{f}|\hat{\mathbf{f}}, \mathbf{A})||\mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}_{nn})] \\
 &= \mathbb{E}_{\mathcal{N}(\mathbf{f}|\hat{\mathbf{f}}, \mathbf{A})} \left[ \ln \mathcal{N}(\mathbf{f}|\hat{\mathbf{f}}, \mathbf{A}) - \ln \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}_{nn}) \right] \\
 &= \frac{1}{2} \mathbb{E}_{\mathcal{N}(\mathbf{f}|\hat{\mathbf{f}}, \mathbf{A})} \left[ -\ln |\mathbf{A}| - (\mathbf{f} - \hat{\mathbf{f}})^T \mathbf{A}^{-1} (\mathbf{f} - \hat{\mathbf{f}}) + \ln |\mathbf{K}_{nn}| + \mathbf{f}^T \mathbf{K}_{nn}^{-1} \mathbf{f} \right] \\
 &= \frac{1}{2} \ln \frac{|\mathbf{A}|}{|\mathbf{K}_{nn}|} + \frac{1}{2} \mathbb{E}_{\mathcal{N}(\mathbf{f}|\hat{\mathbf{f}}, \mathbf{A})} \left[ -\text{tr}[\mathbf{A}^{-1} (\mathbf{f} - \hat{\mathbf{f}}) (\mathbf{f} - \hat{\mathbf{f}})^T] + \text{tr}[\mathbf{K}_{nn}^{-1} \mathbf{f} \mathbf{f}^T] \right] \\
 &= \frac{1}{2} \ln \frac{|\mathbf{A}|}{|\mathbf{K}_{nn}|} - \frac{1}{2} \text{tr}(\mathbf{A} \mathbf{A}^{-1}) + \int \text{tr}[\mathbf{K}_{nn}^{-1} \mathbf{f} \mathbf{f}^T] \mathcal{N}(\mathbf{f}|\hat{\mathbf{f}}, \mathbf{A}) d\mathbf{f} \\
 &= \frac{1}{2} \left[ \ln |\mathbf{A}| - \ln |\mathbf{K}_{nn}| - n + \text{tr}(\mathbf{K}_{nn}^{-1} \mathbf{A}) + \hat{\mathbf{f}}^T \mathbf{K}_{nn}^{-1} \hat{\mathbf{f}} \right].
 \end{aligned} \tag{38}$$

We arrive to the final version of the ELBO:

$$\ln p(\mathcal{D}|\theta) \geq - \sum_{i=1}^n \left( \log[\sqrt{2\pi}\sigma_y] + \frac{(f_i - \hat{f}_i)^2 + a_{ii}^2}{2\sigma_y^2} \right) \left[ \Phi\left(\frac{u - \hat{f}_i}{a_{ii}}\right) - \Phi\left(\frac{l - \hat{f}_i}{a_{ii}}\right) \right] \tag{39}$$

$$+ \sum_{i=1}^n \left[ \log \Phi\left(\frac{f_i - u}{\sigma_y}\right) + \frac{(u + \hat{f}_i - 2f_i)a_{ii}}{2\sigma_y^2} \right] \phi\left(\frac{u - \hat{f}_i}{a_{ii}}\right) \tag{40}$$

$$+ \sum_{i=1}^n \left[ \log \Phi\left(\frac{l - f_i}{\sigma_y}\right) - \frac{(l + \hat{f}_i - 2f_i)a_{ii}}{2\sigma_y^2} \right] \phi\left(\frac{l - \hat{f}_i}{a_{ii}}\right) \tag{41}$$

$$- \frac{1}{2} \left[ \ln |\mathbf{A}| - \ln |\mathbf{K}_{nn}| - n + \text{tr}(\mathbf{K}_{nn}^{-1} \mathbf{A}) + \hat{\mathbf{f}}^T \mathbf{K}_{nn}^{-1} \hat{\mathbf{f}} \right]. \tag{42}$$