

---

# Scalable Gaussian Processes with Billions of Inducing Inputs via Tensor Train Decomposition

---

**Pavel A. Izmailov**  
Lomonosov Moscow State  
University

**Alexander V. Novikov**  
National Research University  
Higher School of Economics

**Dmitry A. Kropotov**  
Lomonosov Moscow State  
University

## Abstract

We propose a method (TT-GP) for approximate inference in Gaussian Process (GP) models. We build on previous results on scalable GPs including stochastic variational inference based on inducing inputs, kernel interpolation, and structure exploiting algebra. The key idea of our method is to use Tensor Train decomposition for variational parameters, which allows us to train GPs with billions of inducing inputs and achieve state-of-the-art results on several benchmarks. Further, our approach allows for training kernels based on deep neural networks without any modifications to the underlying GP model. A neural network learns a multidimensional embedding for the data, which is used by the GP to make the final prediction. We train GP and neural network parameters in end-to-end manner without pretraining through maximization of GP marginal likelihood. We show the efficiency of the proposed approach on several regression and classification benchmark datasets including MNIST, CIFAR-10, and Airline.

## 1 Introduction

Gaussian processes (GPs) provide a prior over functions and allow finding complex regularities in data. The ability of GPs to adjust the complexity of the model to the size of the data makes them appealing to use for big datasets. Unfortunately, standard methods for GP regression and classification scale as  $\mathcal{O}(n^3)$  with the size of the data  $n$  and can not be applied when  $n$  exceeds several thousands.

Numerous approximate inference methods have been proposed in the literature. Many of these methods are

based on the concept of inducing inputs (Quiñonero-Candela and Rasmussen (2005), Snelson and Ghahramani (2006), Williams and Seeger (2000)). These methods build a smaller set  $Z$  of  $m$  points that serve to approximate the true posterior of the process and reduce the complexity to  $\mathcal{O}(nm^2 + m^3)$ . Titsias (2009) proposed to consider the values  $u$  of the Gaussian process at the inducing inputs as latent variables and derived a variational inference procedure to approximate the posterior distribution of these variables. Hensman et al. (2013) and Hensman et al. (2015) extended this framework by using stochastic optimization to scale up the method and generalized it to classification problems.

Inducing input methods allow to use Gaussian processes on datasets containing millions of examples. However, these methods are still limited in the number of inducing points  $m$  they can use (usually up to  $10^4$ ). Small number of inducing inputs limits the flexibility of the models that can be learned with these methods, and does not allow to learn expressive kernel functions (Wilson et al. (2014)). Wilson and Nickisch (2015) proposed KISS-GP framework, which exploits the Kronecker product structure in covariance matrices for inducing points placed on a multidimensional grid in the feature space. KISS-GP has complexity  $\mathcal{O}(n + Dm^{1+1/D})$ , where  $D$  is the dimensionality of the feature space. Note however, that  $m$  is the number of points in a  $D$ -dimensional grid and grows exponentially with  $D$ , which makes the method impractical when the number of features  $D$  is larger than 4.

In this paper, we propose TT-GP method, that can use billions of inducing inputs and is applicable to a much wider range of datasets compared to KISS-GP. We achieve this by combining kernel interpolation and Kronecker algebra of KISS-GP with a scalable variational inference procedure. We restrict the family of variational distributions from Hensman et al. (2013) to have parameters in compact formats. Specifically, we use Kronecker product format for the covariance matrix  $\Sigma$  and Tensor Train format (Oseledets (2011)) for the expectation  $\mu$  of the variational distribution over the values  $u$  of the process at inducing inputs  $Z$ .

Nickson et al. (2015) showed that using Kronecker for-

mat for  $\Sigma$  does not substantially affect the predictive performance of GP regression, while allowing for computational gains. The main contribution of our paper is combining the Kronecker format for  $\Sigma$  with TT-format for  $\mu$ , which, together with efficient inference procedure, allows us to efficiently train GP models with billions of inducing inputs.

Unlike KISS-GP the proposed method has linear complexity with respect to dimensionality  $D$  of the feature space. It means that we can apply TT-GP to datasets that are both large and high-dimensional. Note however, that TT-GP is constructing a grid of inducing inputs in the feature space, and tries to infer the values of the process in all points in the grid. High-dimensional real-world datasets are believed to lie on small-dimensional manifolds in the feature space, and it is impractical to try to recover the complex non-linear transformation that a Gaussian Process defines on the whole feature space. Thus, we use TT-GP on raw features for datasets with dimensionality up to 10. For feature spaces with higher dimensionality we propose to use kernels based on parametric projections, which can be learned from the data.

Wilson et al. (2016a) and Wilson et al. (2016b) demonstrated efficiency of Gaussian processes with kernels based on deep neural networks. They used subsets of outputs of a DNN as inputs for a Gaussian process. As the authors were using KISS-GP, they were limited to using additive kernels, combining multiple low dimensional Gaussian processes. We found that DNN-based kernels are very efficient in combination with TT-GP. These kernels allows us to train TT-GP models on high-dimensional datasets including computer vision tasks. Moreover, unlike the existing deep kernel learning methods, TT-GP does not require any changes in the GP model and allows deep kernels that produce embeddings of dimensionality up to 10.

## 2 Background

### 2.1 Gaussian Processes

A Gaussian process is a collection of random variables, any finite number of which have a joint normal distribution. A GP  $f$  taking place in  $\mathbb{R}^D$  is fully defined by its mean  $m : \mathbb{R}^D \rightarrow \mathbb{R}$  and covariance  $k : \mathbb{R}^D \times \mathbb{R}^D \rightarrow \mathbb{R}$  functions. For every  $x_1, x_2, \dots, x_n \in \mathbb{R}^D$

$$f(x_1), f(x_2), \dots, f(x_n) \sim \mathcal{N}(m, K),$$

where  $m = (m(x_1), m(x_2), \dots, m(x_n))^T \in \mathbb{R}^n$ , and  $K \in \mathbb{R}^{n \times n}$  is the covariance matrix with  $K_{ij} = k(x_i, x_j)$ . Below we will use notation  $K(A, B)$  for the matrix of pairwise values of covariance function  $k$  on points from sets  $A$  and  $B$ .

Consider a regression problem. The dataset consists of  $n$  objects  $X = (x_1, \dots, x_n)^T \in \mathbb{R}^{n \times D}$ , and target

values  $y = (y_1, y_2, \dots, y_n)^T \in \mathbb{R}^n$ . We assume that the data is generated by a latent zero-mean Gaussian process  $f$  plus independent Gaussian noise:

$$\begin{aligned} p(y, f|X) &= p(f|X) \prod_{i=1}^n p(y_i|f_i), \\ p(f|X) &= \mathcal{N}(f|0, K(X, X)), \\ p(y_i|f_i) &= \mathcal{N}(y_i|f_i, \nu^2 I), \end{aligned} \tag{1}$$

where  $f_i = f(x_i)$  is the value of the process at data point  $x_i$  and  $\nu^2$  is the noise variance.

Assume that we want to predict the values of the process  $f_*$  at a set of test points  $X_*$ . As the joint distribution of  $y$  and  $f_*$  is Gaussian, we can analytically compute the conditional distribution  $p(f_*|y, X, X_*) = \mathcal{N}(f_*|\hat{m}, \hat{K})$  with tractable formulas for  $\hat{m}$  and  $\hat{K}$ . The complexity of computing  $\hat{m}$  and  $\hat{K}$  is  $\mathcal{O}(n^3)$  since it involves calculation of the inverse of the covariance matrix  $K(X, X)$ .

Covariance functions usually have a set of hyper-parameters  $\theta$ . For example, the RBF kernel

$$k_{\text{RBF}}(x, x') = \sigma_f^2 \exp(-0.5\|x - x'\|^2/l^2)$$

has two hyper-parameters  $l$  and  $\sigma_f$ . In order to fit the model to the data, we can maximize the marginal likelihood of the process  $p(y|X)$  with respect to these parameters. In case of GP regression this marginal likelihood is tractable and can be computed in  $\mathcal{O}(n^3)$  operations.

For two-class classification problem we use the same model (1) with  $p(y_i|f_i) = 1/(1 + \exp(-y_i f_i))$ , where  $y_i \in \{-1, +1\}$ . In this case both predictive distribution and marginal likelihood are intractable. For detailed description of GP regression and classification see Rasmussen and Williams (2006).

### 2.2 Inducing Inputs

A number of approximate methods were developed to scale up Gaussian processes. Hensman et al. (2013) proposed a variational lower bound that factorizes over observations for Gaussian process marginal likelihood. We rederive this bound here.

Consider a set  $Z \in \mathbb{R}^{m \times D}$  of  $m$  inducing points in the feature space and latent variables  $u \in \mathbb{R}^m$  representing the values of the Gaussian process at these points. Consider the augmented model

$$p(y, f, u) = p(y|f)p(f|u)p(u) = \prod_{i=1}^n p(y_i|f_i)p(f|u)p(u)$$

with

$$\begin{aligned} p(f|u) &= \mathcal{N}(f|K_{nm}K_{mm}^{-1}u, K_{nn} - K_{nm}K_{mm}^{-1}K_{mn}), \\ p(u) &= \mathcal{N}(u|0, K_{mm}), \end{aligned}$$

where  $K_{nn} = K(X, X)$ ,  $K_{nm} = K(X, Z)$ ,  $K_{mn} = K(Z, X) = K_{nm}^T$ ,  $K_{mm} = K(Z, Z)$ .

The standard variational lower bound is given by

$$\begin{aligned} \log p(y) &\geq \mathbb{E}_{q(u, f)} \log \frac{p(y, f, u)}{q(u, f)} = \\ &= \mathbb{E}_{q(f)} \log \prod_{i=1}^n p(y_i | f_i) - \text{KL}(q(u, f) || p(u, f)), \end{aligned} \quad (2)$$

where  $q(u, f)$  is the variational distribution over latent variables. Consider the following family of variational distributions

$$q(u, f) = p(f|u) \mathcal{N}(u | \mu, \Sigma), \quad (3)$$

where  $\mu \in \mathbb{R}^m$  and  $\Sigma \in \mathbb{R}^{m \times m}$  are variational parameters. Then the marginal distribution over  $f$  can be computed analytically

$$\begin{aligned} q(f) &= \mathcal{N}(f | K_{nm} K_{mm}^{-1} \mu, \\ &\quad K_{nn} + K_{nm} K_{mm}^{-1} (\Sigma - K_{mm}) K_{mm}^{-1} K_{mn}). \end{aligned} \quad (4)$$

We can then rewrite (2) as

$$\begin{aligned} \log p(y) &\geq \sum_{i=1}^n \mathbb{E}_{q(f_i)} \log p(y_i | f_i) - \\ &\quad \text{KL}(q(u) || p(u)). \end{aligned} \quad (5)$$

Note, that the lower bound (5) factorizes over observations and thus stochastic optimization can be applied to maximize this bound with respect to both kernel hyper-parameters  $\theta$  and variational parameters  $\mu$  and  $\Sigma$ . In case of regression we can rewrite (5) in the closed form

$$\begin{aligned} \log p(y) &\geq \sum_{i=1}^n \left( \log \mathcal{N}(y_i | k_i^T K_{mm}^{-1} \mu, \nu^2) - \right. \\ &\quad \left. - \frac{1}{2\nu^2} \tilde{K}_{ii} - \frac{1}{2\nu^2} \text{tr}(k_i^T K_{mm}^{-1} \Sigma K_{mm}^{-1} k_i) \right) - \\ &\quad - \frac{1}{2} \left( \log \frac{|K_{mm}|}{|\Sigma|} - m + \text{tr}(K_{mm}^{-1} \Sigma) + \mu^T K_{mm}^{-1} \mu \right), \end{aligned} \quad (6)$$

where  $k_i \in \mathbb{R}^m$  is the  $i$ -th column of  $K_{mn}$  matrix and  $\tilde{K} = K_{nn} - K_{nm} K_{mm}^{-1} K_{mn}$ .

At prediction time we can use the variational distribution as a substitute for posterior

$$\begin{aligned} p(f_* | y) &= \int p(f_* | f, u) p(f, u | y) df du \approx \\ &\approx \int p(f_* | f, u) q(f, u) df du = \int p(f_* | u) q(u) du. \end{aligned}$$

The complexity of computing the bound (6) is  $\mathcal{O}(nm^2 + m^3)$ . Hensman et al. (2015) proposes to use Gauss-Hermite quadratures to approximate the expectation

term in (5) for binary classification problem to obtain the same computational complexity  $\mathcal{O}(nm^2 + m^3)$ . This complexity allows to use Gaussian processes in tasks with millions of training samples, but these methods are limited to use small numbers of inducing points  $m$ , which hurts the predictive performance and doesn't allow to learn expressive kernel functions.

### 2.3 KISS-GP

Saatçi (2012) noted that the covariance matrices computed at points on a multidimensional grid in the feature space can be represented as a Kronecker product if the kernel function factorizes over dimensions

$$k(x, x') = k_1(x^1, x'^1) \cdot k_2(x^2, x'^2) \dots k_D(x^D, x'^D). \quad (7)$$

Note, that many popular covariance functions, including RBF, belong to this class. Kronecker structure of covariance matrices allows to perform efficient inference for full Gaussian processes with inputs  $X$  on a grid.

Wilson and Nickisch (2015) proposed to set inducing inputs  $Z$  on a grid:

$$Z = Z^1 \times Z^2 \times \dots \times Z^D, \quad Z^i \in \mathbb{R}^{m_i} \quad \forall i = 1, 2, \dots, D.$$

The number  $m$  of inducing points is then given by  $m = \prod_{i=1}^D m_i$ .

Let the covariance function satisfy (7). Then the covariance matrix  $K_{mm}$  can be represented as a Kronecker product over dimensions

$$K_{mm} = K_{m_1 m_1}^1 \otimes K_{m_2 m_2}^2 \otimes \dots \otimes K_{m_D m_D}^D,$$

where

$$K_{m_i m_i}^i = K_i(Z_i, Z_i) \in \mathbb{R}^{m_i \times m_i} \quad \forall i = 1, 2, \dots, D.$$

Kronecker products allow efficient computation of matrix inverse and determinant:

$$\begin{aligned} (A_1 \otimes A_2 \otimes \dots \otimes A_D)^{-1} &= A_1^{-1} \otimes A_2^{-1} \otimes \dots \otimes A_D^{-1}, \\ |A_1 \otimes A_2 \otimes \dots \otimes A_D| &= |A_1|^{c_1} \cdot |A_2|^{c_2} \cdot \dots \cdot |A_D|^{c_D}, \end{aligned}$$

where  $A_i \in \mathbb{R}^{k_i \times k_i}$ ,  $c_i = \prod_{j \neq i} k_j$ ,  $\forall i = 1, 2, \dots, D$ .

Another major idea of KISS-GP is to use interpolation to approximate  $K_{mn}$ . Considering inducing inputs as interpolation points for the function  $k(\cdot, z_i)$  we can write

$$K_{mn} \approx K_{mm} W, \quad k_i \approx K_{mm} w_i, \quad (8)$$

where  $W \in \mathbb{R}^{m \times n}$  contains the coefficients of interpolation, and  $w_i$  is its  $i$ -th column. Authors of KISS-GP suggest using cubic convolutional interpolation (Keys (1981)), in which case the interpolation weights  $w_i$  can be represented as a Kronecker product over dimensions

$$w_i = w_i^1 \otimes w_i^2 \otimes \dots \otimes w_i^D, \quad w_i \in \mathbb{R}^{m_i} \quad \forall i = 1, 2, \dots, D.$$

Wilson and Nickisch (2015) combine these ideas with SOR (Silverman (1985)) in the KISS-GP method yielding  $\mathcal{O}(n + Dm^{1+1/D})$  computational complexity. This complexity allows to use KISS-GP with a large number (possibly greater than  $n$ ) of inducing points. Note, however, that  $m$  grows exponentially with the dimensionality  $D$  of the feature space and the method becomes impractical when  $D > 4$ .

## 2.4 Tensor Train Decomposition

Tensor Train (TT) decomposition, proposed in Oseledets (2011), allows to efficiently store tensors (multidimensional arrays of data), large matrices, and vectors. For matrices and vectors in the TT-format linear algebra operations can be implemented efficiently. The TT format was successfully applied for different machine learning tasks (see Novikov et al. (2014), Novikov et al. (2015)).

Consider a  $D$ -dimensional tensor  $\mathcal{A} \in \mathbb{R}^{k_1 \times k_2 \times \dots \times k_D}$ .  $\mathcal{A}$  is said to be in the Tensor Train format if

$$\mathcal{A}(i_1, i_2, \dots, i_D) = G_1[i_1] \cdot G_2[i_2] \cdot \dots \cdot G_D[i_D], \quad (9)$$

$$i_k \in \{1, 2, \dots, n_k\} \quad \forall k,$$

where  $G_k[i_k] \in \mathbb{R}^{r_{k-1} \times r_k} \quad \forall k, i_k, r_0 = r_D = 1$ . Matrices  $G_k$  are called TT-cores, and numbers  $r_k$  are called TT-ranks of tensor  $\mathcal{A}$ .

In order to represent a vector in TT-format, it is reshaped to a multidimensional tensor (possibly with zero padding) and then format (9) is used. We will use TT-format for the vector  $\mu$  of expectations of the values  $u$  of the Gaussian process in points  $Z$  placed on a multidimensional grid. In this case,  $\mu$  is naturally represented as a  $D$ -dimensional tensor.

For matrices TT format is given by

$$M(i_1, i_2, \dots, i_D; j_1, j_2, \dots, j_D) =$$

$$G_1[i_1, j_1] \cdot G_2[i_2, j_2] \cdot \dots \cdot G_D[i_D, j_D],$$

where  $G_k[i_k, j_k] \in \mathbb{R}^{r_{k-1} \times r_k} \quad \forall k, i_k, j_k, r_0 = r_D = 1$ . Note, that Kronecker product format is a special case of the TT-matrix with TT-ranks  $r_1 = r_2 = \dots = r_D = 1$ .

Let  $u, v \in \mathbb{R}^{n_1 \cdot n_2 \cdot \dots \cdot n_D}$  be vectors in TT-format with TT-ranks not greater than  $r$ . Let  $A$  and  $B$  be represented as a Kronecker product

$$A = A_1 \otimes A_2 \otimes \dots \otimes A_D, \quad A_k \in \mathbb{R}^{n_k \times n_k} \quad \forall k,$$

and the same for  $B$ . Let  $n = \max_k n_k$ . Then the computational complexity of computing the quadratic form  $u^T A v$  and  $\text{tr}(AB)$  is  $\mathcal{O}(Dnr^3)$  and  $\mathcal{O}(Dn^2)$  correspondingly. We will need these two operations below. See Oseledets (2011) for a detailed description of TT format and efficient algorithms implementing linear algebraic operations with it.

## 3 TT-GP

In the previous section we described several methods for GP regression and classification. All these methods have different limitations: standard methods are limited to small-scale datasets, KISS-GP requires small dimensionality of the feature space, and other methods based on inducing points are limited to use a small number  $m$  of these points. In this section, we propose the TT-GP method that can be used with big datasets and can incorporate billions of inducing inputs. Additionally, TT-GP allows for training expressive deep kernels to work with structured data (e.g. images).

### 3.1 Variational Parameters Approximation

In section 2.2 we derived the variational lower bound of Hensman et al. (2013). We will place the inducing inputs  $Z$  on a multidimensional grid in the feature space and we will assume the covariance function satisfies (7). Let the number of inducing points in each dimension be  $m_0$ . Then the total number of induced points is  $m = m_0^D$ . As shown in Section 2.3, in this case  $K_{mm}$  matrix can be expressed as a Kronecker product over dimensions. Substituting the approximation (8) into the lower bound (6), we obtain

$$\log p(y) \geq \sum_{i=1}^n \left( \log \mathcal{N}(y_i | w_i^T \mu, \nu^2) - \frac{1}{2\nu^2} \tilde{K}_{ii} - \frac{1}{2\nu^2} \text{tr}(w_i^T \Sigma w_i) \right) - \frac{1}{2} \left( \log \frac{|K_{mm}|}{|\Sigma|} - m + \text{tr}(K_{mm}^{-1} \Sigma) + \mu^T K_{mm}^{-1} \mu \right), \quad (10)$$

where  $\tilde{K}_{ii} = k(x_i, x_i) - w_i^T K_{mm} w_i$ .

Note that  $K_{mm}^{-1}$  and  $|K_{mm}|$  can be computed with  $\mathcal{O}(Dm_0^3) = \mathcal{O}(Dm^{3/D})$  operations due to the Kronecker product structure. Now the most computationally demanding terms are those containing variational parameters  $\mu$  and  $\Sigma$ .

Let us restrict the family of variational distributions (3). Let  $\Sigma$  be a Kronecker product over dimensions, and  $\mu$  be in the TT-format with TT-rank  $r$  ( $r$  is a hyper-parameter of our method). Then, according to section 2.4, we can compute the lower bound (10) with  $\mathcal{O}(nDm_0r^2 + Dm_0r^3 + Dm_0^3) = \mathcal{O}(nDm^{1/D}r^2 + Dm^{1/D}r^3 + Dm^{3/D})$  complexity.

The proposed TT-GP method has linear complexity with respect to dimensionality  $D$  of the feature space, despite the exponential growth of the number of inducing inputs. Lower bound (10) can be maximized with respect to kernel hyper-parameters  $\theta$ , TT-cores of  $\mu$ , and Kronecker multipliers of  $\Sigma$ . Note that stochastic optimization can be applied, as the bound (10) factorizes over data points.

### 3.2 Classification

In this section we describe a generalization of the proposed method for multiclass classification. In this case the dataset consists of features  $X = (x_1, x_2, \dots, x_n)^T \in \mathbb{R}^{n \times D}$  and target values  $y = (y_1, y_2, \dots, y_n)^T \in \{1, 2, \dots, C\}^n$ , where  $C$  is the number of classes.

Consider  $C$  Gaussian processes taking place in  $\mathbb{R}^D$ . Each process corresponds to its own class. We will place  $m = m_0^D$  inducing points  $Z$  on a grid in the feature space, and they will be shared between all processes. Each process has its own set of latent variables representing the values of the process at data points  $f^c \in \mathbb{R}^n$ , and inducing inputs  $u^c \in \mathbb{R}^m$ . We will use the following model

$$p(y, f, u) = \prod_{i=1}^n p(y_i | f_i^{1,2,\dots,C}) \prod_{c=1}^C p(f^c | u^c) p(u^c),$$

where  $f_i^{1,2,\dots,C}$  is the vector consisting of the values of all processes  $1, 2, \dots, C$  at data point  $i$ ,  $p(f^c | u^c)$  and  $p(u^c)$  are defined as in (2) and

$$p(y_i | f_i^{1,2,\dots,C}) = \frac{\exp(f_i^{y_i})}{\sum_{j=1}^C \exp(f_i^j)}.$$

We will use variational distributions of the form

$$\begin{aligned} q(f^1, f^2, \dots, f^C, u^1, u^2, \dots, u^C) &= \\ &= q(f^1, u^1) \cdot q(f^2, u^2) \cdot \dots \cdot q(f^C, u^C), \end{aligned}$$

where  $q(f^c, u^c) = p(f^c | u^c) \mathcal{N}(u^c | \mu^c, \Sigma^c)$   $c = 1, 2, \dots, C$ , all  $\mu^c$  are represented in TT-format with TT-ranks not greater than  $r$  and all  $\Sigma^c$  are represented as Kronecker products over dimensions. Similarly to (5), we obtain

$$\begin{aligned} \log p(y) &\geq \sum_{i=1}^n \mathbb{E}_{q(f_i^{1,2,\dots,C})} \log p(y_i | f_i^{1,2,\dots,C}) \\ &\quad - \sum_{c=1}^C \text{KL}(q(u_c) || p(u_c)) \end{aligned} \quad (11)$$

The second term in (11) can be computed analytically as a sum of KL-divergences between normal distributions. The first term is intractable. In order to approximate the first term we will use a lower bound. Let  $y_i$  belong to class  $c$ . Then, we can rewrite

$$\begin{aligned} \mathbb{E}_{q(f_i^{1,2,\dots,C})} \log p(y_i | f_i^{1,2,\dots,C}) &= \\ \mathbb{E}_{q(f_i^c) | f_i^c} \log \left( \sum_{j=1}^C \exp(f_i^j) \right), \end{aligned} \quad (12)$$

where  $q(f_i^{1,2,\dots,C}) = \mathcal{N}(f^1 | m_i^1, s_i^1) \cdot \mathcal{N}(f^2 | m_i^2, s_i^2) \cdot \dots \cdot \mathcal{N}(f^C | m_i^C, s_i^C)$ .

The first term in (12) is obviously tractable, while the second term has to be approximated. Bouchard (2007) discusses several lower bounds for expectations of this type. Below we derive one of these bounds, which we use in TT-GP.

Concavity of logarithm implies  $\log \left( \sum_{j=1}^C \exp(f_i^j) \right) \leq \varphi \sum_{j=1}^C \exp(f_i^j) - \log \varphi - 1$ ,  $\forall \varphi > 0$ . Taking expectation of both sides of the inequality and minimizing with respect to  $\varphi$ , we obtain

$$\begin{aligned} \mathbb{E}_{q(f_i^{1,2,\dots,C})} \log \left( \sum_{j=1}^C \exp(f_i^j) \right) &\leq \\ \log \left( \sum_{j=1}^C \exp \left( m_i^j + \frac{1}{2} s_i^j \right) \right). \end{aligned} \quad (13)$$

Substituting (13) back into (11) we obtain a tractable lower bound for multiclass classification task, that can be maximized with respect to kernel hyper-parameters  $\theta^c$ , TT-cores of  $\mu^c$  and Kronecker factors of  $\Sigma^c$ . The complexity of the method is  $C$  times higher, than in regression case.

### 3.3 Deep kernels

Wilson et al. (2016b) and Wilson et al. (2016a) showed the efficiency of using expressive kernel functions based on deep neural networks with Gaussian processes on a variety of tasks. The proposed TT-GP method is naturally compatible with this idea.

Consider a covariance function  $k$  satisfying (7) and a neural network (or in fact any parametric transform)  $net$ . We can define a new kernel as follows

$$k_{net}(x, x') = k(net(x), net(x')).$$

We can train the neural network weights through maximization of GP marginal likelihood, the same way, as we normally train kernel hyper-parameters  $\theta$ . This way, the network learns a multidimensional embedding for the data, and GP is making the prediction working with this embedding. Wilson et al. (2016b) trained additive deep kernels combining one-dimensional GPs on different outputs of a neural network. TT-GP allows us to train Gaussian processes on all network outputs, and train the whole model end-to-end without any pretraining.

## 4 Experiments

In this section we first compare the proposed TT-GP<sup>1</sup> method with SVI-GP (Hensman et al. (2013)) on regression tasks and KLSP-GP (Hensman et al. (2015))

<sup>1</sup>for TT-GP we use our implementation available at <https://anonymized-link>

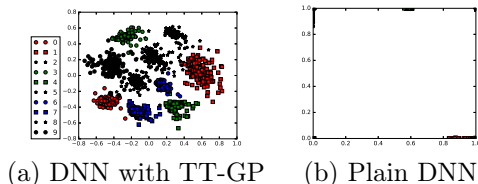


Figure 1: Learned representation for Digits dataset.

on binary classification tasks using standard RBF kernel functions. Then, we test the ability of our method to learn expressive deep kernel functions and compare it with SV-DKL (Wilson et al. (2016b)).

#### 4.1 Expectation approximation

In this section we provide some numerical justification of using Tensor Train format for the mean  $\mu$  of the variational distribution.

- Visualization of the  $\mu$  for standard inducing input GP and TT-GP.
- Plot of quality of TT approximation against TT-rank.

#### 4.2 Standard Kernels

For testing our method with standard covariance functions we used a range of classification and regression tasks from UCI and LIBSVM archives and the Airline dataset, that is popular for testing scalable GP models (Hensman et al. (2013), Hensman et al. (2015), Wilson et al. (2016b), Cutajar et al. (2016)).

For SVI-GP and KLSP-GP we used the implementations provided in GPfLow (Matthews et al. (2016)). For Airline dataset we provide results reported in the original paper (Hensman et al. (2015)). For our experiments, we use a cluster of Intel Xeon E5-2698B v3 CPUs having 16 cores and 230 GB of RAM.

For YearPred, EEG and covtype datasets we used a  $d$ -dimensional linear embedding inside the RBF kernel for TT-GP, as the number  $D$  of features makes it impractical to set inducing inputs on a grid in a  $D$ -dimensional space in this case.

Table 1 shows the results on different regression and classification tasks. We can see, that TT-GP is able to achieve better predictive quality on all datasets except EEG. We also note that the method is able to achieve good predictive performance with linear embedding, which makes it practical for a wide range of datasets.

Table 1: Experimental results for standard kernels. In the table acc. stands for  $r^2$  for regression and accuracy for classification tasks.  $n$  is the size of the training set,  $D$  is the dimensionality of the feature space,  $m$  is the number of inducing inputs,  $r$  is TT-ranks of  $\mu$  for TT-GP;  $t$  is the time per one pass over the data (epoch) in seconds; where provided,  $d$  is the dimensionality of linear embedding.

\* for KLSP-GP on Airline we provide results from the original paper where the accuracy is given as a plot, and detailed information about experiment setup is not available.

Dataset			SVI-GP / KLSP-GP			
Name	$n$	$D$	acc.	$m$	$t$ (s)	acc.
Powerplant	7654	4	0.94	200	10	0.95
Protein	36584	9	0.50	200	45	0.56
YearPred	463K	90	0.30	1000	597	0.32
Airline	6M	8	0.665*	-	-	0.694
svmguidel	3089	4	0.967	200	4	0.969
EEG	11984	14	0.915	1000	18	0.908
covtype bin	465K	54	0.817	1000	320	0.852

#### 4.3 Deep Kernels

##### 4.3.1 Representation learning

We first explore the representation our model learns for data on the small Digits<sup>2</sup> dataset containing  $n = 1797$   $8 \times 8$  images of handwritten digits. We used a TT-GP with a kernel based on a small fully-connected neural network with two hidden layers with 50 neurons each and  $d = 2$  neurons in the output layer to obtain a 2-dimensional embedding. We trained the model to classify the digits to 10 classes corresponding to different digits. Fig. 1,a shows the learned embedding. We also trained the same network standalone, adding another layer with 10 outputs and softmax activations. The embedding for this network is shown in fig. 1,b.

We can see that the stand-alone DNN with linear classifiers is unable to learn a good 2-dimensional embedding. On the other hand, using a flexible GP classifier our model learns to group objects of the same class into compact regions.

##### 4.3.2 Classification tasks

To test our model with deep kernels we used Airline, CIFAR-10 (Krizhevsky (2009)) and MNIST (LeCun et al. (1998)) datasets. The corresponding DNN architecture is shown in Table 2. For CIFAR-10 dataset we also use standard data augmentation techniques with random cropping of  $24 \times 24$  parts of the image, horizontal flipping, randomly adjusting brightness and

<sup>2</sup>[http://scikit-learn.org/stable/auto\\_examples/datasets/plot\\_digits\\_last\\_image.html](http://scikit-learn.org/stable/auto_examples/datasets/plot_digits_last_image.html)

Table 2: DNN architecture used in experiments with deep kernels. Here  $F(h)$  means a fully-connected layer with  $h$  neurons;  $C(h \times w, f)$  means a convolutional layer with  $f$   $h \times w$  filters;  $P(h \times w)$  means max-pooling with  $h \times w$  kernel; ReLU stands for rectified linear unit and BN means batch normalization (Ioffe and Szegedy (2015)).

Dataset	Architecture
Airline	$F(1000)$ -ReLU- $F(1000)$ -ReLU- $F(500)$ -ReLU- $F(50)$ -ReLU- $F(2)$
CIFAR-10	$C(3 \times 3, 128)$ -BN-ReLU- $C(3 \times 3, 128)$ -BN-ReLU- $C(3 \times 3, 256)$ -BN-ReLU- $P(3 \times 3)$ - $C(3 \times 3, 256)$ -BN-ReLU- $C(3 \times 3, 256)$ -BN-ReLU- $P(3 \times 3)$ - $F(1536)$ -BN-ReLU- $P(5 \times 5)$ -BN-ReLU- $F(10)$
MNIST	$C(5 \times 5, 32)$ -ReLU- $P(2 \times 2)$ - $C(5 \times 5, 64)$ -ReLU- $P(2 \times 2)$ - $F(1024)$ -ReLU- $P(4)$

Table 3: Results of experiments with deep kernels. Here acc. is classification accuracy;  $C$  is the number of classes;  $d$  is the dimensionality of embedding learned by the model;  $t$  is the time per one pass over data (epoch) in seconds.

	Dataset			SV-DKL
Name	$n$	$D$	$C$	acc.
Airline	$6M$	8	2	0.781
CIFAR-10	$50K$	$32 \times 32 \times 3$	10	0.770
MNIST	$60K$	$28 \times 28$	10	0.992

contrast. In all experiments we also add a BN without trainable mean and variance after the DNN output layer to project the outputs into the region where inducing inputs are placed. We use  $m_0 = 10$  inducing inputs per dimension and set TT-ranks of  $\mu$  to  $r = 10$  for all three datasets. For experiments with convolutional neural networks, we used Nvidia Tesla K80 GPUs to train the model.

Table 3 shows the results of the experiments for our TT-GP with DNN kernel and SV-DKL. Note, that the comparison is not absolutely fair on CIFAR-10 and MNIST datasets, as we didn’t use the same exact architecture and preprocessing as Wilson et al. (2016b). On Airline dataset we used the same exact architecture and preprocessing as SV-DKL and TT-GP achieves a higher accuracy on this dataset.

We also provide results of stand-alone DNNs for comparison. We used the same networks that were used in TT-GP kernels with the last linear layers replaced by layers with  $C$  outputs and softmax activations. Overall, we can see, that our model is able to achieve good predictive performance, improving the results of stand-alone DNN on Airline and MNIST.

We train all the models from random initialization without pretraining. We also tried using pretrained DNNs as initialization for the kernel of our TT-GP model, which sometimes leads to faster convergence, but does not improve the final accuracy.

## 5 Discussion

We proposed TT-GP method for scalable inference in Gaussian process models for regression and classification. The proposed method is capable of using billions of inducing inputs, which is impossible for existing methods. This allows us to improve the performance over state-of-the-art both with standard and deep kernels on several benchmark datasets. Further, we believe that our model provides a more natural way of learning deep kernel functions than the existing approaches since it makes several hacky tricks unnecessary, such as pretraining or using architectures bottlenecked with 1-dimensional embedding.

Our preliminary experiments showed that TT-GP is inferior in terms of uncertainty quantification compared to existing methods. We suspect that the reason for this is restricting Kronecker structure for the covariance matrix  $\Sigma$ . We hope to alleviate this limitation by using Tensor Train format for  $\Sigma$  and corresponding approximations to its determinant.

As a promising direction for future work we consider training TT-GP with deep kernels incrementally, using the variational approximation of posterior distribution as a prior for new data. We also find it interesting to try using the low-dimensional embeddings learned by our model for transfer learning. Finally, we want to explore the performance of our model in different practical applications including hyper-parameter optimization.

## References

- G. Bouchard. Efficient bounds for the softmax function and applications to approximate inference in hybrid models. In *NIPS 2007 workshop for approximate Bayesian inference in continuous/hybrid systems*, 2007.
- K. Cutajar, E. V. Bonilla, P. Michiardi, and M. Filippone. Practical learning of deep gaussian processes via random fourier features. *arXiv preprint arXiv:1610.04386*, 2016.
- J. Hensman, N. Fusi, and N. D. Lawrence. Gaussian processes for big data. In *Proceedings of the Twenty-Ninth Conference on Uncertainty in Artificial Intelligence*, pages 282–290. AUAI Press, 2013.
- J. Hensman, A. G. de G. Matthews, and Z. Ghahramani. Scalable variational gaussian process classification. In *AISTATS*, 2015.
- S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *Proceedings of the 32nd International Conference on Machine Learning (ICML-15)*, pages 448–456, 2015.
- R. Keys. Cubic convolution interpolation for digital image processing. *IEEE transactions on acoustics, speech, and signal processing*, 29(6):1153–1160, 1981.

- A. Krizhevsky. Learning multiple layers of features from tiny images. 2009.
- Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- A. G. de G. Matthews, M. van der Wilk, T. Nickson, K. Fujii, A. Boukouvalas, P. León-Villagrà, Z. Ghahramani, and J. Hensman. GPflow: A Gaussian process library using TensorFlow. *arXiv preprint 1610.08733*, October 2016.
- T. Nickson, T. Gunter, C. Lloyd, M. A. Osborne, and S. Roberts. Blitzkriging: Kronecker-structured stochastic gaussian processes. *arXiv preprint arXiv:1510.07965*, 2015.
- A. Novikov, A. Rodomanov, A. Osokin, and D. Vetrov. Putting MRFs on a tensor train. In *International Conference on Machine Learning*, pages 811–819, 2014.
- A. Novikov, D. Podoprikin, A. Osokin, and D. Vetrov. Tensorizing neural networks. In *Advances in Neural Information Processing Systems*, pages 442–450, 2015.
- I. V. Oseledets. Tensor-train decomposition. *SIAM Journal on Scientific Computing*, 33(5):2295–2317, 2011.
- J. Quiñonero-Candela and C. E. Rasmussen. A unifying view of sparse approximate gaussian process regression. *Journal of Machine Learning Research*, 6 (Dec):1939–1959, 2005.
- C. E. Rasmussen and C. K. I. Williams. Gaussian processes for machine learning, 2006.
- Y. Saatçi. *Scalable inference for structured Gaussian process models*. PhD thesis, University of Cambridge, 2012.
- B. W. Silverman. Some aspects of the spline smoothing approach to non-parametric regression curve fitting. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 1–52, 1985.
- E. Snelson and Z. Ghahramani. Sparse gaussian processes using pseudo-inputs. *Advances in neural information processing systems*, 18:1257, 2006.
- M. K. Titsias. Variational learning of inducing variables in sparse gaussian processes. In *AISTATS*, volume 5, pages 567–574, 2009.
- C. K. I. Williams and M. Seeger. Using the nyström method to speed up kernel machines. In *Proceedings of the 13th International Conference on Neural Information Processing Systems*, pages 661–667. MIT press, 2000.
- A. G. Wilson and H. Nickisch. Kernel interpolation for scalable structured gaussian processes (kiss-gp). In *International Conference on Machine Learning*, pages 1775–1784, 2015.
- A. G. Wilson, Z. Hu, R. Salakhutdinov, and E. P. Xing. Deep kernel learning. In *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics*, pages 370–378, 2016a.
- A. G. Wilson, Z. Hu, R. Salakhutdinov, and E. P. Xing. Stochastic variational deep kernel learning. In *Advances in Neural Information Processing Systems*, pages 2586–2594, 2016b.
- Andrew Wilson, Elad Gilboa, John P Cunningham, and Arye Nehorai. Fast kernel learning for multidimensional pattern extrapolation. In *Advances in Neural Information Processing Systems*, pages 3626–3634, 2014.