

Appendix of Distributed Asynchronous Variational Gaussian Processes

A Derivatives

A.1 Objective

As described in the paper, the objective function to be minimized is $-\mathcal{L} = \sum_{i=1}^n g_i + h$, where

$$\begin{aligned} g_i &= -\ln \mathcal{N}(y_i | \phi_i^T \boldsymbol{\mu}, \beta^{-1}) + \frac{\beta}{2} \phi_i^T \boldsymbol{\Sigma} \phi_i + \frac{\beta}{2} \tilde{k}_{ii} \\ &= \frac{1}{2} \ln 2\pi - \frac{1}{2} \ln \beta + \frac{\beta}{2} \left(y_i^2 - 2y_i \phi_i^T \boldsymbol{\mu} + \boldsymbol{\mu}^T \phi_i \phi_i^T \boldsymbol{\mu} + \phi_i^T \boldsymbol{\Sigma} \phi_i + k_{ii} - \phi_i^T \phi_i \right), \end{aligned} \quad (1)$$

$$\begin{aligned} h &= \text{KL}(q(\mathbf{w}) || p(\mathbf{w})) \\ &= \frac{1}{2} \left(-\ln |\boldsymbol{\Sigma}| - m + \text{tr}(\boldsymbol{\Sigma}) + \boldsymbol{\mu}^T \boldsymbol{\mu} \right), \end{aligned} \quad (2)$$

and we define $\beta = \sigma^{-2}$ and $\phi_i = \phi(\mathbf{x}_i)$.

A.2 Kernel

A common choice for the kernel is the anisotropic squared exponential covariance function:

$$k(\mathbf{x}_i, \mathbf{x}_j) = a_0^2 \exp \left(-\frac{1}{2} (\mathbf{x}_i - \mathbf{x}_j)^T \text{diag}(\boldsymbol{\eta}) (\mathbf{x}_i - \mathbf{x}_j) \right), \quad (3)$$

in which the hyperparameters are the signal variance a_0 and the lengthscales $\boldsymbol{\eta} = \{1/a_k^2\}_{k=1}^d$, controlling how fast the covariance decays with the distance between inputs. Using this covariance function, we can prune input dimensions by shrinking the corresponding lengthscales based on the data (when $\eta_d = 0$, the d -th dimension becomes totally irrelevant to the covariance function value). This pruning is known as Automatic Relevance Determination (ARD) and therefore this covariance is also called the ARD squared exponential.

A.3 Derivative over $\ln \sigma$ ($\ln \beta^{-1/2}$)

The derivative of g_i over $\ln \sigma$ is

$$\frac{\partial g_i}{\partial \ln \sigma} = 1 - \frac{1}{\sigma^2} (y_i^2 - 2y_i \phi_i^T \boldsymbol{\mu} + \phi_i^T (\boldsymbol{\Sigma} + \boldsymbol{\mu} \boldsymbol{\mu}^T) \phi_i + k_{ii} - \phi_i^T \phi_i). \quad (4)$$

A.4 Derivative over $\ln a_0$

The derivative of g_i over $\ln a_0$ is

$$\frac{\partial g_i}{\partial \ln a_0} = \frac{1}{\sigma^2} (-y_i \phi_i^T \boldsymbol{\mu} + \phi_i^T (\boldsymbol{\Sigma} + \boldsymbol{\mu} \boldsymbol{\mu}^T) \phi_i + k_{ii} - \phi_i^T \phi_i). \quad (5)$$

A.5 Derivative over \mathbf{Z}

By defining \mathbf{L} the lower triangular Cholesky factor of \mathbf{K}_{mm}^{-1} , the derivative of g_i over \mathbf{Z} is

$$\begin{aligned} \frac{\partial g_i}{\partial \mathbf{Z}} &= \frac{1}{\sigma^2} \left[((\mathbf{L} \mathbf{p}_i) \circ \mathbf{k}_m(\mathbf{x}_i)) \mathbf{x}_i^T \text{diag}(\boldsymbol{\eta}) - (((\mathbf{L} \mathbf{p}_i) \circ \mathbf{k}_m(\mathbf{x}_i)) \mathbf{1}_d^T) \circ (\mathbf{Z} \text{diag}(\boldsymbol{\eta})) \right. \\ &\quad \left. - (\mathbf{T}_i + \mathbf{T}_i^T) \mathbf{Z} \text{diag}(\boldsymbol{\eta}) + ((\mathbf{T}_i + \mathbf{T}_i^T) \mathbf{1}_m \boldsymbol{\eta}^T) \circ \mathbf{Z} \right], \end{aligned} \quad (6)$$

where

$$\mathbf{p}_i = -\boldsymbol{\mu} y_i + (\boldsymbol{\mu} \boldsymbol{\mu}^T + \boldsymbol{\Sigma}) \phi(\mathbf{x}_i) - \phi(\mathbf{x}_i), \quad (7)$$

$$\mathbf{T}_i = \left[\mathbf{L} ((\phi(\mathbf{x}_i) \mathbf{p}_i^T) \circ \boldsymbol{\Psi}) \mathbf{L}^T \right] \circ \mathbf{K}_{mm}. \quad (8)$$

The symbol \circ denotes the Hadamard product, and Ψ is an upper triangular matrix with diagonal elements all equal to 0.5 and strictly upper triangular elements all equal to 1, as follows:

$$\Psi = \begin{bmatrix} 0.5 & 1 & \dots & 1 & 1 \\ 0 & 0.5 & \ddots & 1 & 1 \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ 0 & 0 & \ddots & 0.5 & 1 \\ 0 & 0 & \dots & 0 & 0.5 \end{bmatrix}. \quad (9)$$

A.6 Derivative over $\ln \boldsymbol{\eta}$

The derivative of g_i over $\ln \boldsymbol{\eta}$ is

$$\begin{aligned} \frac{\partial g_i}{\partial \ln \boldsymbol{\eta}} = & \frac{1}{2\sigma^2} \left\{ 2\mathbf{1}_m^T [\mathbf{Z} \circ ((\mathbf{L}\mathbf{p}_i) \circ \mathbf{k}_m(\mathbf{x}_i)) \mathbf{x}_i^T] - \mathbf{1}_m^T ((\mathbf{L}\mathbf{p}_i) \circ \mathbf{k}_m(\mathbf{x}_i)) (\mathbf{x}_i \circ \mathbf{x}_i)^T - ((\mathbf{L}\mathbf{p}_i) \circ \mathbf{k}_m(\mathbf{x}_i))^T (\mathbf{Z} \circ \mathbf{Z}) \right. \\ & \left. - \mathbf{1}_m^T [\mathbf{Z} \circ ((\mathbf{T}_i + \mathbf{T}_i^T) \mathbf{Z})] + \mathbf{1}_m^T [(\mathbf{T}_i + \mathbf{T}_i^T)(\mathbf{Z} \circ \mathbf{Z})] \right\} \circ \boldsymbol{\eta}. \end{aligned} \quad (10)$$

B Properties of the ELBO of ADVGP

By defining \mathbf{U} as the upper triangular Cholesky factor of $\boldsymbol{\Sigma}$, i.e., $\boldsymbol{\Sigma} = \mathbf{U}^T \mathbf{U}$, we have

Lemma B.1 *The gradient of g_i in Equation 1, ∇g_i , is Lipschitz continuous with respect to each element in $\boldsymbol{\mu}$ and \mathbf{U} .*

We can prove this by showing the first derivative of ∇g_i with respect to each element of $\boldsymbol{\mu}$ and \mathbf{U} is bounded, which is constant in our case. As shown in our paper, the gradients of g_i with respect to $\boldsymbol{\mu}$ and \mathbf{U} are:

$$\frac{\partial g_i}{\partial \boldsymbol{\mu}} = \frac{1}{\sigma^2} [-y_i \boldsymbol{\phi}_i + \boldsymbol{\phi}_i \boldsymbol{\phi}_i^T \boldsymbol{\mu}], \quad (11)$$

$$\frac{\partial g_i}{\partial \mathbf{U}} = \frac{1}{\sigma^2} \text{triu}[\mathbf{U} \boldsymbol{\phi}_i \boldsymbol{\phi}_i^T], \quad (12)$$

which are affine functions for $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ respectively. Therefore, the first derivative of ∇g_i is constant.

Lemma B.2 *h in Equation 2 is a convex function with respect to $\boldsymbol{\mu}$ and \mathbf{U} .*

This can be proved by verifying that the Hessian matrices of h with respect to $\boldsymbol{\mu}$ and $\text{vec}(\mathbf{U})$ are both positive semidefinite, where we denote $\text{vec}(\cdot)$ as the operator that stacks the columns of a matrix as a vector. To show this, we first compute the partial derivatives of h with respect to $\boldsymbol{\mu}$ and \mathbf{U} as

$$\frac{\partial h}{\partial \boldsymbol{\mu}} = \boldsymbol{\mu}, \quad (13)$$

$$\frac{\partial h}{\partial \mathbf{U}} = -\text{diag}(\mathbf{U}^{-1}) + \mathbf{U}. \quad (14)$$

The Hessian matrix of h with respect to $\boldsymbol{\mu}$ is

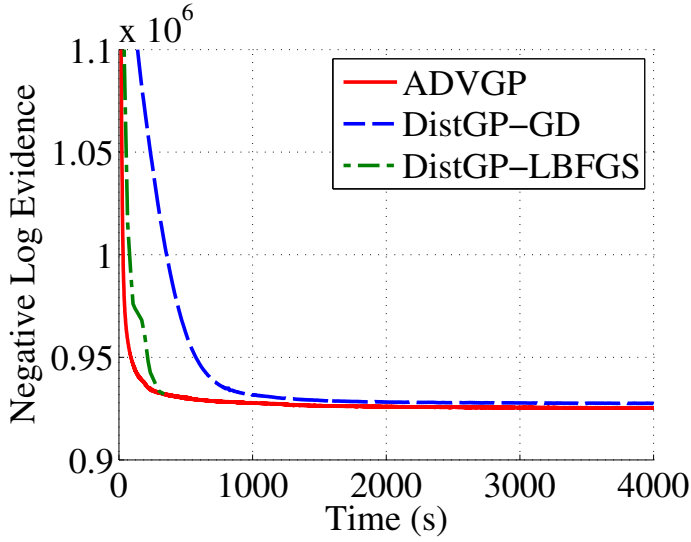
$$\mathbf{H}(\boldsymbol{\mu}) = \mathbf{I}_{m \times m} \succeq 0. \quad (15)$$

The Hessian matrix of h with respect to $\text{vec}(\mathbf{U})$ is

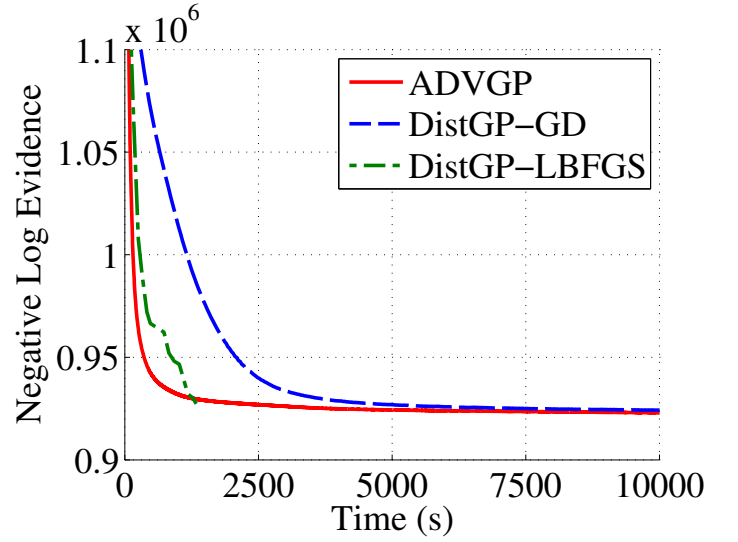
$$\mathbf{H}(\text{vec}(\mathbf{U})) = \text{diag}(\mathbf{h}) \succeq 0, \quad (16)$$

where $\mathbf{h} = [\frac{\partial h}{\partial U_{11}^2}, \dots, \frac{\partial h}{\partial U_{1m}^2}, \dots, \frac{\partial h}{\partial U_{m1}^2}, \dots, \frac{\partial h}{\partial U_{mm}^2}]$, and $\frac{\partial h}{\partial U_{ij}^2} = 1 + \delta(i, j) \frac{1}{U_{i,i}^2}$.

C Negative Log Evidences on US Flight Data



(A) $m = 100$

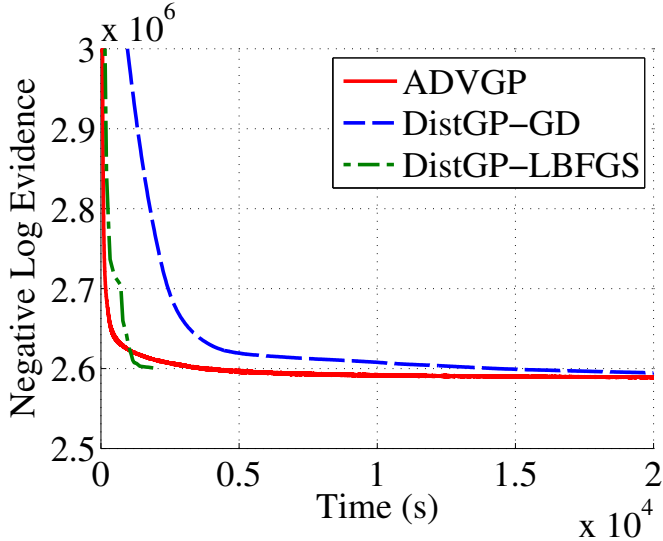


(B) $m = 200$

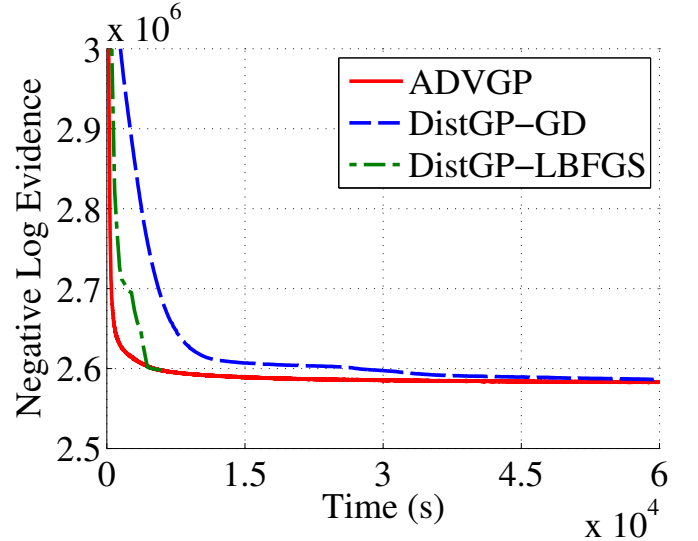
Figure C.1: Negative log evidences for 700K/100K US Flight data as a function of training time.

Method	$m = 100$	$m = 200$
ADVGP	925236	922907
DistGP-GD	927414	924208
DistGP-LBFGS	932179	927331

Table C.1: Negative log evidences for 700K/100K US Flight data.



(A) $m = 100$



(B) $m = 200$

Figure C.2: Negative log evidences for 2M/100K US Flight data as a function of training time.

Method	$m = 100$	$m = 200$
ADVGP	2.58921×10^6	2.58267×10^6
DistGP-GD	2.59471×10^6	2.58601×10^6
DistGP-LBFGS	2.59971×10^6	2.59817×10^6

Table C.2: Negative log evidences for 2M/100K US Flight data.

D Mean Negative Log Predictive Likelihoods (MNLPs) on US Flight Data

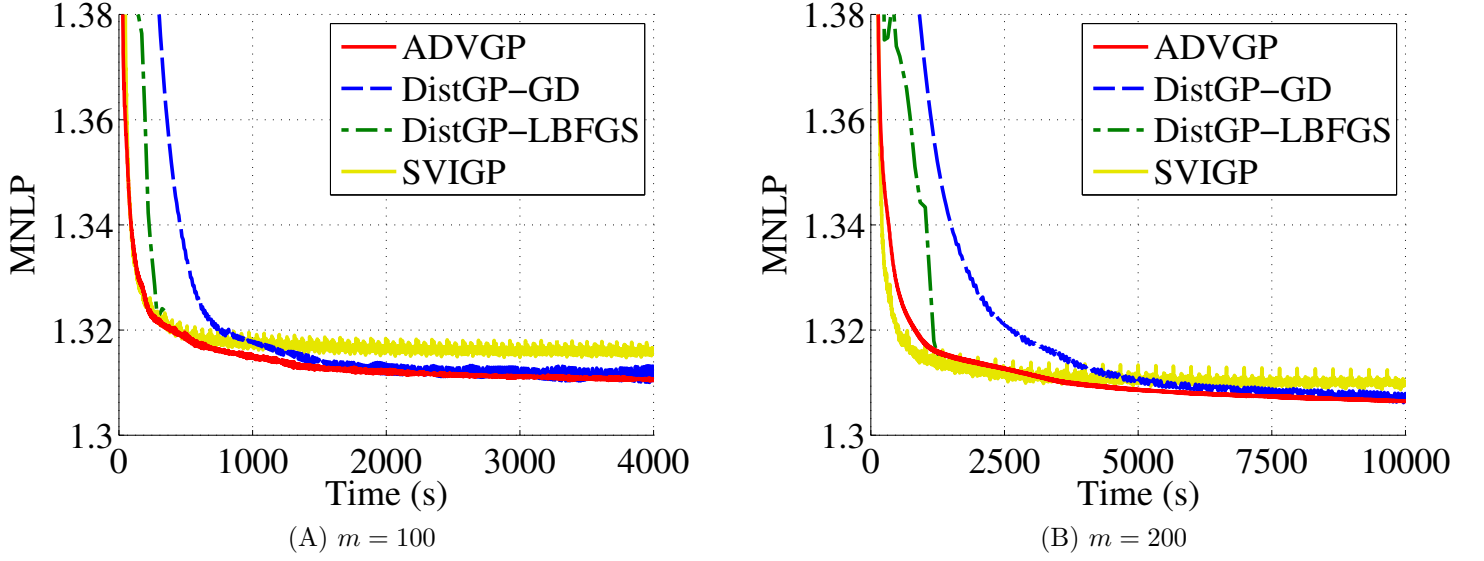


Figure D.1: Mean negative log predictive likelihoods for 700K/100K US Flight data as a function of training time.

Method	$m = 100$	$m = 200$
ADVGP	1.3106	1.3066
DistGP-GD	1.3099	1.3062
DistGP-LBFGS	1.3237	1.3136
SVIGP	1.3157	1.3096

Table D.1: Mean negative log predictive likelihoods for 700K/100K US Flight data.

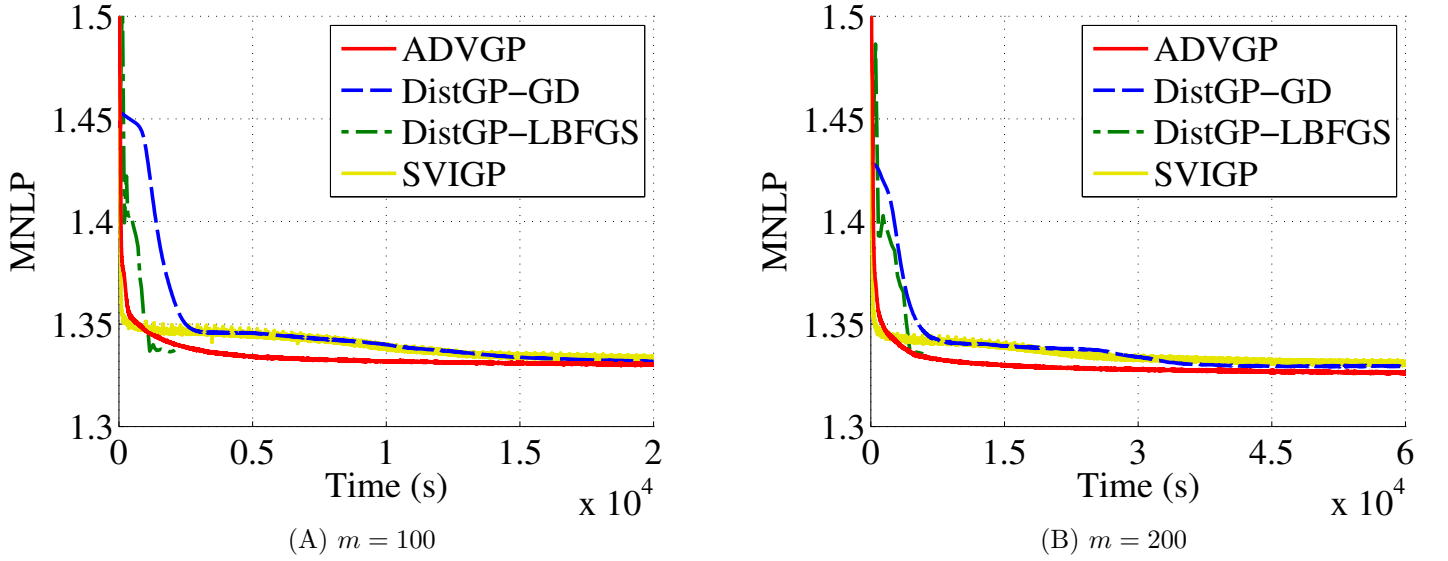


Figure D.2: Mean negative log predictive likelihoods for 2M/100K US Flight data as a function of training time.

Method	$m = 100$	$m = 200$
ADVGP	1.3301	1.3258
DistGP-GD	1.3317	1.3297
DistGP-LBFGS	1.3380	1.3355
SVIGP	1.3335	1.3306

Table D.2: Mean negative log predictive likelihoods for 2M/100K US Flight data.