

Scalable and Deep Gaussian Processes

Dmitry A. Kropotov



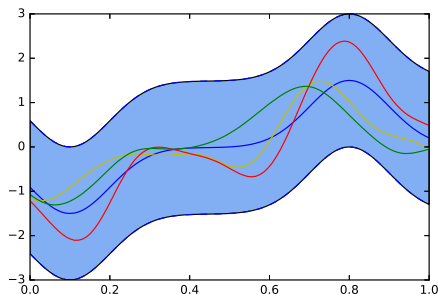
Gaussian Process (GP) is a stochastic process over real-valued functions.

$$f(\mathbf{x}) \sim GP(m(\cdot), k(\cdot, \cdot)) \Leftrightarrow$$

$$\mathbf{f} = [f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)] \sim \mathcal{N}(\mathbf{f} | \mathbf{m}_n, K_{nn}),$$

$$\mathbf{m}_n = [m(\mathbf{x}_1), \dots, m(\mathbf{x}_n)],$$

$$K_{nn} = \{k(\mathbf{x}_i, \mathbf{x}_j)\}_{i,j=1}^{n,n}.$$



Gaussian Process for regression and classification

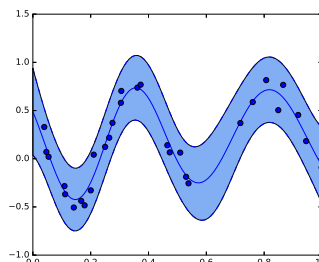
Suppose we have a dataset $(\mathbf{y}, X) = \{y_i, \mathbf{x}_i\}_{i=1}^n$, $\mathbf{x}_i \in \mathbb{R}^d$ – feature vectors, y_i – target values.

$$p(\mathbf{y}, \mathbf{f}|X, \boldsymbol{\theta}) = p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|X, \boldsymbol{\theta}) = \prod_{i=1}^n p(y_i|f_i)p(\mathbf{f}|X, \boldsymbol{\theta}),$$

$$p(\mathbf{f}|X, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{f}|\mathbf{0}, K(X, X; \boldsymbol{\theta})) - \text{GP},$$

$$p(y_i|f_i) = \mathcal{N}(y_i|f_i, \sigma^2) - \text{regression},$$

$$p(y_i|f_i) = \frac{1}{1 + \exp(-y_i f_i)} - \text{classification with 2 classes}.$$



GP probabilistic model:

$$p(\mathbf{y}, \mathbf{f} | X, \boldsymbol{\theta}) = p(\mathbf{y} | \mathbf{f}) p(\mathbf{f} | X, \boldsymbol{\theta}) = \prod_{i=1}^n p(y_i | f_i) p(\mathbf{f} | X, \boldsymbol{\theta}).$$

Training:

$$\log p(\mathbf{y} | X, \boldsymbol{\theta}) \geq \mathbb{E}_{q(\mathbf{f})} [\log p(\mathbf{y}, \mathbf{f} | X, \boldsymbol{\theta}) - \log q(\mathbf{f})] \rightarrow \max_{\boldsymbol{\theta}, q(\cdot)},$$

$$q(\mathbf{f}) = \mathcal{N}(\mathbf{f} | \boldsymbol{\mu}, \Sigma) \approx p(\mathbf{f} | \mathbf{y}, X, \boldsymbol{\theta}).$$

Testing:

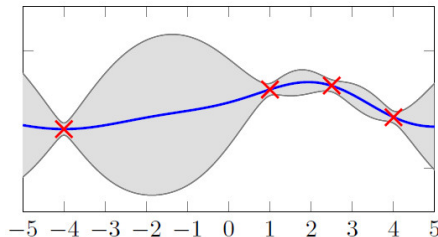
$$\begin{aligned} p(\mathbf{f}_{test} | \mathbf{y}, \boldsymbol{\theta}, X, X_{test}) &= \int p(\mathbf{f}_{test} | \mathbf{f}, \boldsymbol{\theta}, X, X_{test}) p(\mathbf{f} | \mathbf{y}, \boldsymbol{\theta}, X) d\mathbf{f} \approx \\ &\approx \int p(\mathbf{f}_{test} | \mathbf{f}, \boldsymbol{\theta}, X, X_{test}) q(\mathbf{f}) d\mathbf{f}. \end{aligned}$$

Pros:

- Automatically adjust all parameters during training;
- Provides variance during prediction;

Cons:

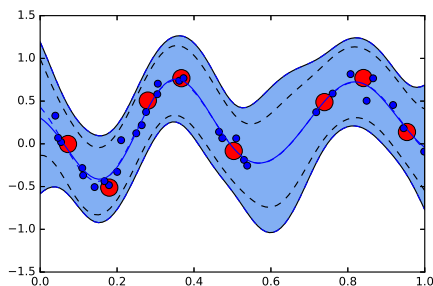
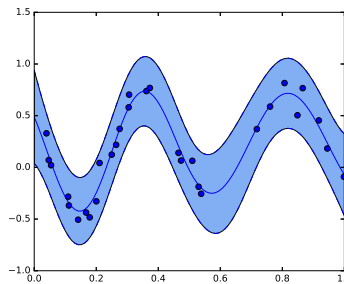
- Training scales as $O(n^3)$, where n – training sample size;
- Do not allow very complex dependencies like deep neural nets.



Gaussian Process with inducing inputs

Let's introduce some new points $Z = \{z_i\}_{i=1}^m$ and suppose we know GP values at these points \mathbf{u} . Key assumption:

$$p(\mathbf{f}|\mathbf{y}, \mathbf{u}, X, Z) \approx p(\mathbf{f}|\mathbf{u}, X, Z).$$



Augmented model:

$$\begin{aligned}p(\mathbf{y}, \mathbf{f}, \mathbf{u} | X, Z, \boldsymbol{\theta}) &= p(\mathbf{y} | \mathbf{f}) p(\mathbf{f}, \mathbf{u} | X, Z, \boldsymbol{\theta}), \\p(\mathbf{f}, \mathbf{u} | X, Z, \boldsymbol{\theta}) &= \mathcal{N}([\mathbf{f}, \mathbf{u}] | [\mathbf{0}_n, \mathbf{0}_m], K_{n+m, n+m}).\end{aligned}$$

Augmented model coincides with previous one:

$$\int p(\mathbf{y}, \mathbf{f}, \mathbf{u} | X, Z, \boldsymbol{\theta}) d\mathbf{u} = p(\mathbf{y}, \mathbf{f} | X, \boldsymbol{\theta}) = p(\mathbf{y} | \mathbf{f}) \mathcal{N}(\mathbf{f} | \mathbf{0}_n, K_{nn}).$$

Augmented model:

$$\begin{aligned}p(\mathbf{y}, \mathbf{f}, \mathbf{u} | X, Z, \boldsymbol{\theta}) &= p(\mathbf{y} | \mathbf{f}) p(\mathbf{f}, \mathbf{u} | X, Z, \boldsymbol{\theta}), \\p(\mathbf{f}, \mathbf{u} | X, Z, \boldsymbol{\theta}) &= \mathcal{N}([\mathbf{f}, \mathbf{u}] | [\mathbf{0}_n, \mathbf{0}_m], K_{n+m, n+m}).\end{aligned}$$

Augmented model coincides with previous one:

$$\int p(\mathbf{y}, \mathbf{f}, \mathbf{u} | X, Z, \boldsymbol{\theta}) d\mathbf{u} = p(\mathbf{y}, \mathbf{f} | X, \boldsymbol{\theta}) = p(\mathbf{y} | \mathbf{f}) \mathcal{N}(\mathbf{f} | \mathbf{0}_n, K_{nn}).$$

Training:

$$\begin{aligned}\log p(\mathbf{y} | X, Z, \boldsymbol{\theta}) &\geq \mathbb{E}_{q(\mathbf{f}, \mathbf{u})} [\log p(\mathbf{y}, \mathbf{f}, \mathbf{u} | X, Z, \boldsymbol{\theta}) - \log q(\mathbf{f}, \mathbf{u})] \rightarrow \max_{\boldsymbol{\theta}, q(\mathbf{f}, \mathbf{u})}, \\q(\mathbf{f}, \mathbf{u}) &\approx p(\mathbf{f}, \mathbf{u} | \mathbf{y}, X, Z, \boldsymbol{\theta}) = p(\mathbf{f} | \mathbf{u}, \mathbf{y}, X, Z, \boldsymbol{\theta}) p(\mathbf{u} | \mathbf{y}, X, Z, \boldsymbol{\theta}).\end{aligned}$$

Augmented model:

$$\begin{aligned}p(\mathbf{y}, \mathbf{f}, \mathbf{u} | X, Z, \boldsymbol{\theta}) &= p(\mathbf{y} | \mathbf{f}) p(\mathbf{f}, \mathbf{u} | X, Z, \boldsymbol{\theta}), \\p(\mathbf{f}, \mathbf{u} | X, Z, \boldsymbol{\theta}) &= \mathcal{N}([\mathbf{f}, \mathbf{u}] | [\mathbf{0}_n, \mathbf{0}_m], K_{n+m, n+m}).\end{aligned}$$

Augmented model coincides with previous one:

$$\int p(\mathbf{y}, \mathbf{f}, \mathbf{u} | X, Z, \boldsymbol{\theta}) d\mathbf{u} = p(\mathbf{y}, \mathbf{f} | X, \boldsymbol{\theta}) = p(\mathbf{y} | \mathbf{f}) \mathcal{N}(\mathbf{f} | \mathbf{0}_n, K_{nn}).$$

Training:

$$\begin{aligned}\log p(\mathbf{y} | X, Z, \boldsymbol{\theta}) &\geq \mathbb{E}_{q(\mathbf{f}, \mathbf{u})} [\log p(\mathbf{y}, \mathbf{f}, \mathbf{u} | X, Z, \boldsymbol{\theta}) - \log q(\mathbf{f}, \mathbf{u})] \rightarrow \max_{\boldsymbol{\theta}, q(\mathbf{f}, \mathbf{u})}, \\q(\mathbf{f}, \mathbf{u}) &\approx p(\mathbf{f}, \mathbf{u} | \mathbf{y}, X, Z, \boldsymbol{\theta}) = p(\mathbf{f} | \mathbf{u}, \mathbf{y}, X, Z, \boldsymbol{\theta}) p(\mathbf{u} | \mathbf{y}, X, Z, \boldsymbol{\theta}).\end{aligned}$$

Using key assumption we can choose $q(\mathbf{f}, \mathbf{u})$ as follows:

$$q(\mathbf{f}, \mathbf{u}) = p(\mathbf{f} | \mathbf{u}, X, Z, \boldsymbol{\theta}) q(\mathbf{u}) = p(\mathbf{f} | \mathbf{u}, X, Z, \boldsymbol{\theta}) \mathcal{N}(\mathbf{u} | \boldsymbol{\mu}, \Sigma).$$

Family for approximate posterior:

$$q(\mathbf{f}, \mathbf{u}) = p(\mathbf{f}|\mathbf{u}, X, Z, \boldsymbol{\theta})q(\mathbf{u}).$$

Training:

$$\begin{aligned}\log p(\mathbf{y}|X, Z, \boldsymbol{\theta}) &\geq \mathbb{E}_{q(\mathbf{f}, \mathbf{u})} \log \frac{p(\mathbf{y}, \mathbf{f}, \mathbf{u}|X, Z, \boldsymbol{\theta})}{q(\mathbf{f}, \mathbf{u})} = \\ &= \mathbb{E}_{q(\mathbf{f}, \mathbf{u})} \log \frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u}, X, Z, \boldsymbol{\theta})p(\mathbf{u}|Z, \boldsymbol{\theta})}{p(\mathbf{f}|\mathbf{u}, X, Z, \boldsymbol{\theta})q(\mathbf{u})} = \\ &= \mathbb{E}_{q(\mathbf{f}, \mathbf{u})} \log p(\mathbf{y}|\mathbf{f}) + \mathbb{E}_{q(\mathbf{u})} \log \frac{p(\mathbf{u}|Z, \boldsymbol{\theta})}{q(\mathbf{u})} = \\ &= \sum_{i=1}^n \mathbb{E}_{q(f_i)} \log p(y_i|f_i) - \text{KL}(q(\mathbf{u})\|p(\mathbf{u}|Z, \boldsymbol{\theta})).\end{aligned}$$

Optimization criterion:

$$\log p(\mathbf{y}|X, Z, \boldsymbol{\theta}) \geq \sum_{i=1}^n \mathbb{E}_{q(f_i)} \log p(y_i|f_i) - \text{KL}(q(\mathbf{u})||p(\mathbf{u}|Z, \boldsymbol{\theta})) \rightarrow \max_{\boldsymbol{\theta}, \boldsymbol{\mu}, \Sigma}.$$

First term:

$$\begin{aligned} q(\mathbf{f}) &= \int q(\mathbf{f}, \mathbf{u}) d\mathbf{u} = \int p(\mathbf{f}|\mathbf{u}, X, Z) q(\mathbf{u}) d\mathbf{u} = \\ &= \int \mathcal{N}(\mathbf{f} | K_{nm} K_{mm}^{-1} \mathbf{u}, K_{nn} - K_{nm} K_{mm}^{-1} K_{mn}) \mathcal{N}(\mathbf{u} | \boldsymbol{\mu}, \Sigma) d\mathbf{u} = \\ &= \mathcal{N}(\mathbf{f} | K_{nm} K_{mm}^{-1} \boldsymbol{\mu}, K_{nn} + K_{nm} K_{mm}^{-1} (\Sigma - K_{mm}) K_{mm}^{-1} K_{mn}). \end{aligned}$$

Optimization criterion:

$$\log p(\mathbf{y}|X, Z, \boldsymbol{\theta}) \geq \sum_{i=1}^n \mathbb{E}_{q(f_i)} \log p(y_i|f_i) - \text{KL}(q(\mathbf{u})||p(\mathbf{u}|Z, \boldsymbol{\theta})) \rightarrow \max_{\boldsymbol{\theta}, \boldsymbol{\mu}, \Sigma}.$$

First term:

$$\begin{aligned} q(\mathbf{f}) &= \int q(\mathbf{f}, \mathbf{u}) d\mathbf{u} = \int p(\mathbf{f}|\mathbf{u}, X, Z) q(\mathbf{u}) d\mathbf{u} = \\ &= \int \mathcal{N}(\mathbf{f}|K_{nm}K_{mm}^{-1}\mathbf{u}, K_{nn} - K_{nm}K_{mm}^{-1}K_{mn}) \mathcal{N}(\mathbf{u}|\boldsymbol{\mu}, \Sigma) d\mathbf{u} = \\ &= \mathcal{N}(\mathbf{f}|K_{nm}K_{mm}^{-1}\boldsymbol{\mu}, K_{nn} + K_{nm}K_{mm}^{-1}(\Sigma - K_{mm})K_{mm}^{-1}K_{mn}). \end{aligned}$$

$$\begin{aligned} q(f_i) &= \mathcal{N}(f_i|\mathbf{k}_i^T K_{mm}^{-1}\boldsymbol{\mu}, k_{ii} + \mathbf{k}_i^T K_{mm}^{-1}(\Sigma - K_{mm})K_{mm}^{-1}\mathbf{k}_i), \\ \mathbf{k}_i &= \{k(\mathbf{x}_i, \mathbf{z}_j)\}_{j=1}^m; \quad k_{ii} = k(\mathbf{x}_i, \mathbf{x}_i). \end{aligned}$$

We don't need to work with K_{nn} , only with its diagonal!

Optimization criterion:

$$\log p(\mathbf{y}|X, Z, \boldsymbol{\theta}) \geq \sum_{i=1}^n \mathbb{E}_{q(f_i)} \log p(y_i|f_i) - \text{KL}(q(\mathbf{u})||p(\mathbf{u}|Z, \boldsymbol{\theta})) \rightarrow \max_{\boldsymbol{\theta}, \boldsymbol{\mu}, \Sigma}.$$

Second term:

$$\begin{aligned} \text{KL}(q(\mathbf{u})||p(\mathbf{u}|Z, \boldsymbol{\theta})) &= \text{KL}(\mathcal{N}(\mathbf{u}|\boldsymbol{\mu}, \Sigma)||\mathcal{N}(\mathbf{u}|\mathbf{0}_m, K_{mm})) = \\ &= -\frac{m}{2} - \frac{1}{2} \log \det K_{mm}^{-1} \Sigma + \frac{1}{2} \text{tr} K_{mm}^{-1} \Sigma + \frac{1}{2} \boldsymbol{\mu}^T K_{mm}^{-1} \boldsymbol{\mu}. \end{aligned}$$

Total costs for all $q(f_i)$ and the second term:

Optimization criterion:

$$\log p(\mathbf{y}|X, Z, \boldsymbol{\theta}) \geq \sum_{i=1}^n \mathbb{E}_{q(f_i)} \log p(y_i|f_i) - \text{KL}(q(\mathbf{u})||p(\mathbf{u}|Z, \boldsymbol{\theta})) \rightarrow \max_{\boldsymbol{\theta}, \boldsymbol{\mu}, \Sigma}.$$

Second term:

$$\begin{aligned} \text{KL}(q(\mathbf{u})||p(\mathbf{u}|Z, \boldsymbol{\theta})) &= \text{KL}(\mathcal{N}(\mathbf{u}|\boldsymbol{\mu}, \Sigma)||\mathcal{N}(\mathbf{u}|\mathbf{0}_m, K_{mm})) = \\ &= -\frac{m}{2} - \frac{1}{2} \log \det K_{mm}^{-1} \Sigma + \frac{1}{2} \text{tr} K_{mm}^{-1} \Sigma + \frac{1}{2} \boldsymbol{\mu}^T K_{mm}^{-1} \boldsymbol{\mu}. \end{aligned}$$

Total costs for all $q(f_i)$ and the second term:

$$O(nm^2 + m^3).$$

We need to estimate $\mathbb{E}_{q(f_i)} \log p(y_i|f_i)$.

In case of regression:

$$p(y_i|f_i) = \mathcal{N}(y_i|f_i, \sigma^2),$$

$$\mathbb{E}_{\mathcal{N}(f_i|m_i, s_i^2)} \log p(y_i|f_i) = \log \mathcal{N}(f_i|m_i, \sigma^2) - \frac{1}{2} s_i^2.$$

In case of classification:

$$p(y_i|f_i) = \frac{1}{1 + \exp(-y_i f_i)},$$

$$\mathbb{E}_{\mathcal{N}(f_i|m_i, s_i^2)} \log p(y_i|f_i) = \mathbb{E}_{\mathcal{N}(\xi|0,1)} \log p(y_i|s_i \xi + m_i).$$

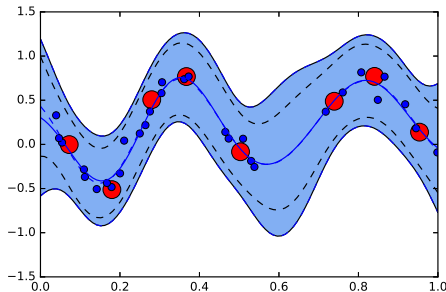
Expectation w.r.t. standard normal distribution can be estimated using Gauss-Hermite quadrature.

Gaussian Process with inducing inputs

- Stochastic optimization with one epoch complexity $O(nm^2 + m^3)$;
- Can optimize w.r.t. positions of inducing points Z , but usually take them fixed, e.g. as cluster centres.

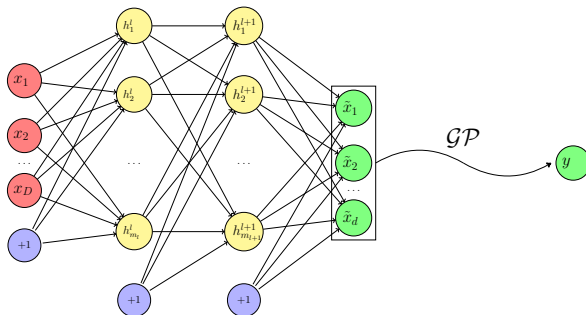
Limitations:

- Can't handle many inducing inputs;
- Current model is not deep.



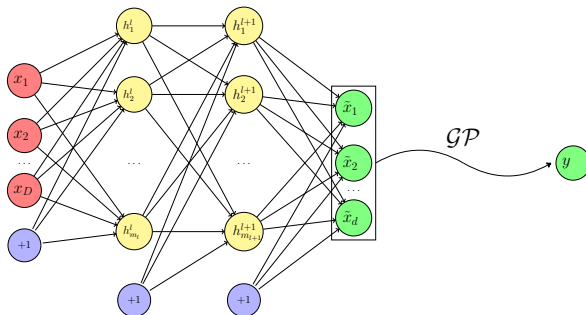
Use new covariance function:

$$k(\mathbf{x}, \mathbf{y}) = k(\text{net}(\mathbf{x}; \boldsymbol{\eta}), \text{net}(\mathbf{y}; \boldsymbol{\eta})).$$



Use new covariance function:

$$k(\mathbf{x}, \mathbf{y}) = k(\text{net}(\mathbf{x}; \boldsymbol{\eta}), \text{net}(\mathbf{y}; \boldsymbol{\eta})).$$



It is not clear where to put inducing inputs!

For two matrices $A \in \mathbb{R}^{n \times m}$ and $B \in \mathbb{R}^{p \times q}$ their Kronecker product is $np \times mq$ matrix

$$A \otimes B = \begin{bmatrix} a_{11}B & \dots & a_{1m}B \\ \dots & \ddots & \dots \\ a_{n1}B & \dots & a_{nm}B \end{bmatrix}.$$

Properties:

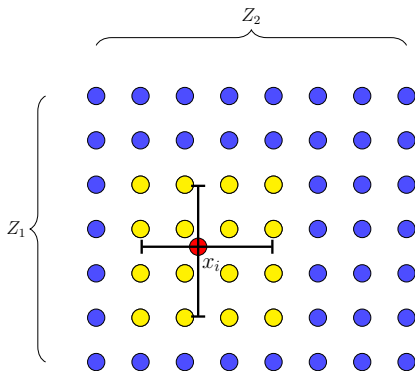
- $(A_1 \otimes A_2 \otimes \dots \otimes A_r)^{-1} = A_1^{-1} \otimes A_2^{-1} \otimes \dots \otimes A_r^{-1}$;
- $\det(A_1 \otimes A_2 \otimes \dots \otimes A_r) = \det(A_1)^{c_1} \det(A_2)^{c_2} \dots \det(A_r)^{c_r}$, where $A_i \in \mathbb{R}^{k_i \times k_i}$, $c_i = \prod_{j \neq i} k_j$.

Inducing points on a regular grid

Let's put inducing inputs Z on a regular grid:

$$Z = Z^1 \times Z^2 \times \cdots \times Z^d,$$

where $Z^i \in \mathbb{R}^{m_i}$, $m = \prod_{i=1}^d m_i$.



Suppose that covariance function can be split over dimensions:

$$k(\mathbf{x}, \mathbf{y}) = k^1(x_1, y_1)k^2(x_2, y_2) \dots k^d(x_d, y_d).$$

E.g. for squared exponential:

$$k(\mathbf{x}, \mathbf{y}) = A \exp(-B\|\mathbf{x} - \mathbf{y}\|^2) = \underbrace{A^{1/d} \exp(-B(x_1 - y_1)^2)}_{k^1(x_1, y_1)} \underbrace{A^{1/d} \exp(-B(x_2 - y_2)^2)}_{k^2(x_2, y_2)} \dots \underbrace{A^{1/d} \exp(-B(x_d - y_d)^2)}_{k^d(x_d, y_d)}$$

Then covariance matrix is given as Kronecker product:

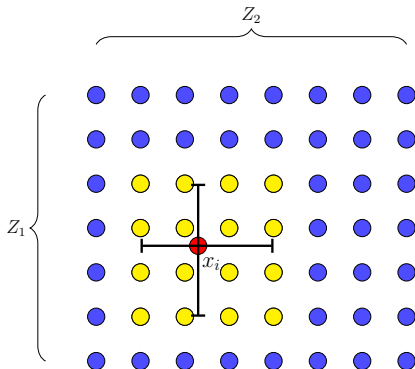
$$K_{mm} = K_{m_1 m_1}^1 \otimes K_{m_2 m_2}^2 \otimes \dots \otimes K_{m_d m_d}^d,$$
$$K_{m_i m_i}^i = K^i(Z^i, Z^i) \in \mathbb{R}^{m_i \times m_i}.$$

Cubic convolution interpolation [Keys, 1981]

We need to estimate covariance between training and inducing inputs K_{mn} . In case of cubic convolution interpolation we have:

$$K_{mn} \approx K_{mm}W, \quad \mathbf{k}_i \approx K_{mm}\mathbf{w}_i,$$

$$\mathbf{w}_i = \mathbf{w}_i^1 \otimes \mathbf{w}_i^2 \otimes \cdots \otimes \mathbf{w}_i^d.$$



Optimization criterion:

$$\log p(\mathbf{y}|X, Z, \boldsymbol{\theta}) \geq \sum_{i=1}^n \mathbb{E}_{q(f_i)} \log p(y_i|f_i) - \text{KL}(q(\mathbf{u})||p(\mathbf{u}|Z, \boldsymbol{\theta})) \rightarrow \max_{\boldsymbol{\theta}, \boldsymbol{\mu}, \Sigma}.$$

If $\mathbf{k}_i = K_{mm}\mathbf{w}_i$, then

$$\begin{aligned} q(f_i) &= \mathcal{N}(f_i | \mathbf{k}_i^T K_{mm}^{-1} \boldsymbol{\mu}, k_{ii} + \mathbf{k}_i^T K_{mm}^{-1} (\Sigma - K_{mm}) K_{mm}^{-1} \mathbf{k}_i) = \\ &= \mathcal{N}(f_i | \mathbf{w}_i^T \boldsymbol{\mu}, k_{ii} + \mathbf{w}_i^T (\Sigma - K_{mm}) \mathbf{w}_i). \end{aligned}$$

Second term:

$$\begin{aligned} \text{KL}(q(\mathbf{u})||p(\mathbf{u}|Z, \boldsymbol{\theta})) &= \\ &= -\frac{m}{2} - \frac{1}{2} \log \det K_{mm}^{-1} \Sigma + \frac{1}{2} \text{tr} K_{mm}^{-1} \Sigma + \frac{1}{2} \boldsymbol{\mu}^T K_{mm}^{-1} \boldsymbol{\mu}. \end{aligned}$$

Tensor Train format [Oseledets, 2011]

Tensor Train format (TT format) gives a compact representation of multidimensional tensors. If $\mathcal{A} \in \mathbb{R}^{m_1 \times m_2 \times \dots \times m_d}$, then

$$A(i_1, i_2, \dots, i_d) = G[i_1]G[i_2] \dots G[i_d],$$
$$G[i_k] \in \mathbb{R}^{r_k \times r_{k+1}}, \quad r_1 = r_{d+1} = 1.$$

Here $G[i_k]$ are TT cores and r_k – TT ranks.

TT format allows many linear algebra operations to perform efficiently.

$$\mathcal{A}(2, 4, 2, 3) = G_1 \times G_2 \times G_3 \times G_4$$

$i_1 = 2 \quad i_2 = 4 \quad i_3 = 2 \quad i_4 = 3$

A family for variational distribution $q(\mathbf{u}) = \mathcal{N}(\mathbf{u}|\boldsymbol{\mu}, \Sigma)$:

- $\boldsymbol{\mu}$ in TT format with fixed rank r ;
- $\Sigma = \Sigma^1 \otimes \Sigma^2 \otimes \dots \otimes \Sigma^d$.

Optimization criterion:

$$\log p(\mathbf{y}|X, Z, \boldsymbol{\theta}) \geq \sum_{i=1}^n \mathbb{E}_{q(f_i)} \log p(y_i|f_i) - \frac{m}{2} + \frac{1}{2} \log \det K_{mm}^{-1} \Sigma -$$

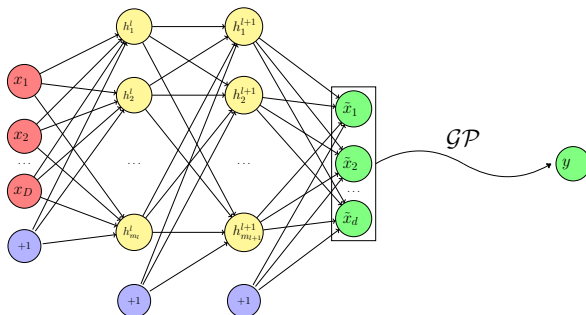
$$- \frac{1}{2} \text{tr} K_{mm}^{-1} \Sigma - \frac{1}{2} \boldsymbol{\mu}^T K_{mm}^{-1} \boldsymbol{\mu},$$

$$q(f_i) = \mathcal{N}(f_i | \mathbf{w}_i^T \boldsymbol{\mu}, k_{ii} + \mathbf{w}_i^T (\Sigma - K_{mm}) \mathbf{w}_i).$$

If m_0 is number of inducing inputs per dimension, then all calculations can be performed in $O(ndm_0r^2 + dm_0r^3 + dm_0^3)$.

TT-GP can be easily combined with deep net covariance function

$$k(\mathbf{x}, \mathbf{y}) = k(\text{net}(\mathbf{x}; \boldsymbol{\eta}), \text{net}(\mathbf{y}; \boldsymbol{\eta})).$$



Representation for Digits

